

---

# MODELING OF SEMICONDUCTOR DEVICES

---

TIBOR GRASSER

---

360.033

SS 2010

GHOSTWRITTEN BY:  
KARL RUPP

BASED ON INPUT FROM:  
ANDREAS GEHRING  
WOLFGANG GÖS  
TIBOR GRASSER  
PHILIPP HEHENBERGER  
FRANZ SCHANOVSKY  
PAUL-JÜRGEN WAGNER

INSTITUTE FOR MICROELECTRONICS

TU VIENNA

# Contents

<b>1</b>	<b>Semiconductor Equations</b>	<b>1</b>
1.1	The (Electrostatic) Field Side of the Problem . . . . .	1
1.2	Continuity Equations . . . . .	2
1.3	Charge Transport – The Drift-Diffusion-Model . . . . .	3
1.4	Summary – The Basic Semiconductor Equations . . . . .	4
1.5	Outlook – Lattice Heat-Flow . . . . .	5
<b>2</b>	<b>Carrier Transport</b>	<b>7</b>
2.1	The Distribution Function . . . . .	7
2.2	The Equilibrium Case . . . . .	10
2.3	Boundary Conditions . . . . .	13
2.3.1	Quasi-Fermi Levels . . . . .	13
2.3.2	Contact Potentials . . . . .	17
2.3.3	The Built-in Potential . . . . .	20
2.4	Ohmic Contacts . . . . .	21
<b>3</b>	<b>Basics of Numerical Analysis</b>	<b>25</b>
3.1	Introduction to Finite Differences . . . . .	25
3.2	Numerical Solution of Differential Equations . . . . .	28
3.3	The Second Order Derivative . . . . .	30
3.4	Accuracy of Different Discretization Schemes . . . . .	32
3.5	Solution of Linear Systems of Equations . . . . .	34
3.6	Solution of Nonlinear Systems of Equations . . . . .	35
<b>4</b>	<b>Two-Dimensional Simulation and Grids</b>	<b>41</b>
4.1	Two-Dimensional Laplace Equation . . . . .	41
4.2	Box Integration Method . . . . .	44
4.2.1	Example: The Poisson Equation . . . . .	45
4.2.2	Example: Extraction of Capacitances . . . . .	46

---

<b>5</b>	<b>Tessellation of Unstructured Meshes</b>	<b>51</b>
5.1	Voronoi Tessellation . . . . .	52
5.2	Triangular Delaunay Meshes . . . . .	52
5.3	Skeleton Mesh . . . . .	54
5.4	Auxiliary Functions . . . . .	57
5.5	Skeleton Mesh – Boundaries . . . . .	59
5.6	Skeleton Mesh – Interfaces . . . . .	62
5.7	Mesh Refinement . . . . .	64
<b>6</b>	<b>Transport Phenomena and their Numerical Analysis</b>	<b>69</b>
6.1	Discretization in the Time Domain . . . . .	70
6.2	Stability of Discretization Schemes . . . . .	73
6.3	Diffusive Problems . . . . .	75
6.4	Convective Problems . . . . .	80
6.5	Diffusive and Convective Problems . . . . .	83
6.5.1	Scharfetter-Gummel Discretization . . . . .	87
<b>7</b>	<b>Parameter Modeling</b>	<b>91</b>
7.1	Carrier Mobilities . . . . .	91
7.1.1	Lattice Scattering . . . . .	92
7.1.2	Ionized Impurity Scattering . . . . .	93
7.1.3	Surface/Interface Scattering . . . . .	93
7.1.4	Carrier Heating . . . . .	95
7.2	Carrier Generation and Recombination . . . . .	96
7.2.1	Phonon Assisted Recombination and Generation . . . . .	97
7.2.2	Photon Transition . . . . .	100
7.2.3	Auger Generation-Recombination . . . . .	101
7.2.4	Impact Ionization . . . . .	102
<b>8</b>	<b>Devices in Detail</b>	<b>105</b>
8.1	Analytical Diode Model . . . . .	105
8.2	Small Signal Analysis of a Diode . . . . .	110
8.3	Analytical Bipolar Junction Transistor Model . . . . .	112
8.4	The Metal-Oxide-Semiconductor Capacitor . . . . .	116
8.5	The Metal-Oxide-Semiconductor Field-Effect-Transistor (MOSFET) . . . . .	122
8.6	CMOS Design Issues . . . . .	125
<b>A</b>	<b>Partial Differential Equations</b>	<b>131</b>

---

## CONTENTS

---

A.1	Boundary and Initial Conditions . . . . .	131
A.2	Classification . . . . .	133
<b>B</b>	<b>Vector Analysis and Its Implementation in SGFramework</b>	<b>135</b>
B.1	Divergence . . . . .	135
B.2	Curl (Rotation) . . . . .	137
B.3	Gradient . . . . .	138
B.4	Nabla . . . . .	139
B.5	Manipulating Expressions . . . . .	140
B.6	Identities . . . . .	141
B.7	Integral Theorems of Stokes and Gauss . . . . .	141
<b>C</b>	<b>Basics of Electromagnetism – Maxwell’s Equations</b>	<b>143</b>
C.1	Interface Conditions . . . . .	144
C.2	Continuity Equation . . . . .	145
<b>D</b>	<b>Vocabulary</b>	<b>147</b>
	<b>Bibliography</b>	<b>150</b>
	<b>Index</b>	<b>150</b>



# Preface<sup>1</sup>

Attempts to model semiconductor devices are nearly as old as the devices themselves. Any good model should be based on a **sound**<sup>2</sup> physical understanding of the particular device and the use of **dumb**<sup>3</sup> empirical fits should be avoided. The theoretical basis we require in this lecture has been developed in the various lectures dealing with the fundamentals of semiconductor physics. The response of carriers to external stimulus and internal forces and fields essentially **determines**<sup>4</sup> the device behavior and can be expressed by a set of differential equations. Various degrees of **sophistication**<sup>5</sup> exist but we will limit ourselves to the most **rudimentary**<sup>6</sup> model, the drift-diffusion model.

Using the drift-diffusion model, the behavior of most (conventional) semiconductor devices can be understood and modeled. Unfortunately, no closed form analytic solutions of the drift-diffusion model exists and one has to **resort**<sup>7</sup> to alternative techniques. Roughly speaking, modeling attempts fall into the two categories compact modeling and numerical modeling. By considering the **peculiarities**<sup>8</sup> of each device **subjected**<sup>9</sup> to particular bias conditions, the art of compact modeling tries to simplify the underlying equations in such a way that an approximate analytic expression is obtained. Conventionally, a number of assumptions enters the derivation and sometimes these assumptions are difficult to justify. Also the analytical solutions are often used without detailed knowledge of these simplifying assumptions, which can dramatically limit the validity of compact models. Thus, it is often of fundamental importance to solve the semiconductor device equations with as few approximations as possible using numerical techniques. Unfortunately, a number of difficulties have to be **surmounted**<sup>10</sup> in order to achieve this goal. As such, numerical device modeling is a highly interdisciplinary subject, requiring a good knowledge of semiconductor device physics, electrical engineering, numerical mathematics, and computer science. In this lecture we try to provide you with the necessary basics, the **big picture**<sup>11</sup>, so to speak. It has to be understood, however, that it is impossible to cover the rich spectrum of knowledge available in all these fields and the interested student is **encouraged**<sup>12</sup> to consult appropriate scientific literature. Access to numerous scientific journals is provided for free to our students by the Technology Library ([www.ub.tuwien.ac.at/eng/index.html](http://www.ub.tuwien.ac.at/eng/index.html)). A good starting point would be search engines such as Scopus ([www.scopus.com](http://www.scopus.com)) and IEEE Explore ([ieeexplore.ieee.org](http://ieeexplore.ieee.org)).

Although a good theoretical knowledge is mandatory, the real understanding of a topic is best gained in a hands-on approach. As such, a strong focus is put on the practical (home-work) part of this lecture where (simple) realistic device structures have to be numerically modeled.

---

<sup>1</sup> **preface** [ˈpreɪfɪs], **NOT** [ˈprɪːfeɪs]: Vorwort    <sup>2</sup> **sound** [saʊnd]: auch: vernünftig, sinnvoll    <sup>3</sup> **dumb** [dʌm]: einfältig, primitiv    <sup>4</sup> **to determine** [dɪˈtɛː.mɪn], **NOT** [determɪn]: bestimmen, festlegen    <sup>5</sup> **sophistication** [səˌfɪs.tɪˈkeɪ.ʃən]: Raffinesse    <sup>6</sup> **rudimentary** [ruːdɪˈmen.tɪr.i]: rudimentär, elementar    <sup>7</sup> **to resort to sth.** [rɪˈzɔːrt]: auf etw. zurückgreifen    <sup>8</sup> **peculiarity** [pɪˌkjuː.liˈer.ə.ti]: Eigenheit, Ausprägung    <sup>9</sup> **to be subjected to sth.** [sʌb.dʒekt]: etwas ausgesetzt werden    <sup>10</sup> **to surmount** [səˈmaʊnt]: bewältigen, überwinden    <sup>11</sup> **big picture** [bɪɡ pɪk.tʃəː]: Das große Ganze, ein erster Einstieg    <sup>12</sup> **to encourage** [ɪnˈkʌr.ɪdʒ]: animieren, ermuntern

In order to avoid the typical **pitfalls**<sup>1</sup> associated with C++ compilers, operating systems, file formats, etc., we provide you with a web interface to the simulation engine SGFRAMEWORK (`meb.iue.tuwien.ac.at`).

Being the dominant language in technical sciences, these lecture notes are written in English, to allow you to develop the necessary vocabulary required in this field. We have taken the liberty to remind you of the **pronunciation**<sup>2</sup> and translation of frequently used words and terms using footnotes. In particular, we have tried to highlight words which are often subjected to creative pronunciation solutions by non-native speakers, such as **determine**<sup>3</sup>. Pronunciation notes are taken from the Cambridge online dictionaries (`dictionary.cambridge.org`) and follow the conventions of American English.

These lecture notes are based on my slides compiled for this lecture and additional input contributed by various people. I am particularly **indebted**<sup>4</sup> to Karl Rupp for **painstakingly**<sup>5</sup> collecting, unifying, extending, and in very large parts completely rewriting the existing material to form this hopefully consistent first version of our lecture notes.

Tibor Grasser  
Wien, March 2009

---

<sup>1</sup> **pitfall** [pɪt.fɑ:l]: Fallgrube, Fallstrick, Fehler    <sup>2</sup> **pronunciation** [prəˌnʌnʃi'eɪʃən]: Aussprache    <sup>3</sup> **to determine** [drɪ'teɪ.nəm], **NOT** [determən]: bestimmen, festlegen    <sup>4</sup> **indebted** [ɪn'det.ɪd]: verpflichtet, verschuldet  
<sup>5</sup> **painstakingly** [peɪnz,tet.kiŋ.li]: sorgfältig [Man beachte den Wortstamm *pain* ;-), Anm. K.R.]

# Symbols

## Mathematical Symbols

$\Omega$	bounded domain
$\partial\Omega$	boundary of a domain $\Omega$
$\mathbb{R}$	real numbers
$A = (a_{ij})_{i,j=1}^n$	$n \times n$ matrix
$\mathbf{b} = (b_i)_{i=1}^n, \mathbf{x} = (x_i)_{i=1}^n$	vectors with $n$ entries
$\frac{d}{dx}, \frac{d}{dt}, \dots$	derivative with respect to $x, t, \dots$
$\frac{d^2}{dx^2}, \frac{d^2}{dt^2}, \dots$	second order derivative with respect to $x, t, \dots$
$\frac{\partial}{\partial x}, \frac{\partial}{\partial t}, \dots$	partial derivative with respect to $x, t, \dots$
$\frac{\partial^2}{\partial x^2}, \frac{\partial^2}{\partial t^2}, \dots$	second order partial derivative with respect to $x, t, \dots$
$\exp(x), e^x$	exponential function
$\nabla$	Nabla operator
$\nabla \cdot$	divergence
$\nabla \times$	rotation
$\Delta = \nabla^2$	Laplace operator

## Physical Quantities and Constants

$E$	electric field
$D$	electric flux density
$B$	magnetic field
$H$	magnetizing field
$\psi$	(electrostatic) potential
$\hat{\epsilon}$ or $\epsilon$	permittivity (first case: a tensor)
$\rho$	charge density
$n$	electron concentration
$p$	hole concentration
$N_A, N_D$	acceptor and donor concentration
$n_i^2$	intrinsic carrier concentration
$q$	elementary charge ( $1.602176487(40) \times 10^{-19}$ C)
$J_n, J_p$	electron/hole current density
$R$	recombination rate



---

$\mu_n, \mu_p$	electron/hole mobility
$v_n, v_p$	electron/hole velocity
$\sigma_n, \sigma_p$	electron/hole conductivity
$D_n, D_p$	electron/hole diffusion coefficient
$k_B$	Boltzmann constant ( $1.3806503 \times 10^{-23} \text{ m}^2 \text{ kg s}^{-2} \text{ K}^{-1}$ )
$T_L$	lattice temperature
$V_T$	thermal voltage
$E_g$	band gap energy
$E_F$	Fermi energy
$V_T$	thermal voltage

## Chapter-specific Symbols and Quantities

### Chapter 1

$\kappa$	thermal conductivity
$\rho$	mass density
$c$	specific heat
$H$	heat source term

### Chapter 2

$q$	momentum
$F$	force field
$R$	random force field (scattering)
$u$	(particle) velocity
$f(\mathbf{p}, \mathbf{r}, t)$	(carrier) distribution function
$Q(f)$	scattering operator
$S(\cdot, \cdot)$	scattering rate
$\mathcal{E}_{\text{pot}}$	potential energy
$\mathcal{E}_{\text{kin}}$	kinetic energy
$E_{\text{tot}}$	total energy
$E_c$	conduction band edge energy
$E_{c,0}$	conduction band edge energy without external bias
$m^*$	effective mass
$v_{\text{th}}$	thermal velocity
$N_c, N_v$	effective densities of states (aka. band weights) of conduction and valence band
$E_{Fn}, E_{Fp}$	quasi Fermi levels for electrons and holes
$\mathcal{E}_v, v \in \{c, v, i\}$	Energy of conduction band, valence band and intrinsic energy level.
$V_{12}$	contact potential between material 1 and material 2
$q\Phi$	work function
$E_{\text{vac}}$	vacuum energy
$V_D$	diffusion voltage (aka. built-in potential)

### Chapter 3

$h_i, \Delta x_i$	grid size, distance between adjacent grid points
$x_i$	$i$ -th grid point (1 dimensional grid)
$u_i$	value of a quantity $u$ evaluated at the $i$ -th grid point $x_i$
$J_F(x)$	Jacobian matrix of a vector-valued function $F$ evaluated at $x$

### Chapter 4

$\mathcal{V}$	a small box
$\partial\mathcal{V}$	surface of a box $\mathcal{V}$

### Chapter 5

$\Omega_i$	$i$ -th box of a tessellation of $\Omega$
$V_i$	volume of the box $\Omega_i$
$\cup$	union
$\mathcal{N}_i$	set of all nodes that are neighbors of the $i$ -th node
$\Delta x, \Delta y$	grid size in $x$ and $y$ direction, distance between adjacent grid points
$x_{i,j}$	grid point
$u_{i,j}$	value of a quantity $u$ evaluated at $x_{i,j}$
$d_{i,j}$	length of the edge that connects $i$ and its neighbor-point $j$
$A_{i,j}$	area of the surface element of $\Omega_i$ that interfaces the box $\Omega_j$
$D_{i,j}$	electric flux density in the outward direction at $A_{i,j}$
$M_{i,j}$	refinement measure

### Chapter 6

$\frac{\partial}{\partial n}$	normal derivative
$\mathcal{F}$	Fourier transformation
$\eta = \frac{\Delta t}{(\Delta x)^2}$	discretization parameter
$\lambda = \frac{\Delta t}{\Delta x}$	discretization parameter
$\mathcal{B}(x) = x/(e^x - 1)$	Bernoulli function

### Chapter 7

$\mu^L$	mobility considering scattering at lattice atoms or defects
$\mu^I$	mobility considering scattering at charges or neutral impurities
$\mu^S$	mobility considering surface roughness scattering
$\mu^F$	mobility considering increased scattering due to heating
$\mu^{\text{LISF}}$	effective mobility
$\tau_n, \tau_p$	carrier lifetimes for electrons and holes respectively
$f_{t,0}$	Fermi-Dirac distribution function
$e^-$	electron
$h^+$	hole

---

$R^{\text{RSH}}$	Read-Shockley-Hall recombination rate
$R^{\text{opt}}$	recombination due to optical generation and radiative recombination
$R^{\text{AU}}$	Auger recombination rate
$G^{\text{II}}$	impact ionization generation rate

## Chapter 8

$l_p$	width of the space-charge region in $p$ -Si
$l_n$	width of the space-charge region in $n$ -Si
$V_c$	contact voltage
$n_{n0}, p_{p0}$	equilibrium majority carrier concentration in $n$ - and $p$ -region respectively
$p_{n0}, n_{p0}$	equilibrium minority carrier concentration in $n$ - and $p$ -region respectively
$C'_J$	junction capacitance (or depletion capacitance) per unit area
$C'_D$	diffusion capacitance (or storage capacitance) per unit area
$\Delta n := n - n_{p0}$	deviation from equilibrium minority carrier concentration in the $npn$ -base
$L_B = \sqrt{D_n \tau_n}$	diffusion length
$W$	base width
$V_{\text{FB}}$	flat band voltage
$\chi_s$	electron affinity
$\phi_s$	surface potential
$V_t$	threshold voltage
$I_D$	drain current
$\gamma$	body factor
$\lambda$	channel length parameter
$S$	subthreshold slope

# Pronunciation

Since the **pronunciation**<sup>1</sup> of English words is often a source of confusion among non-native speakers, a phonetic transcription of critical words is given. This phonetic transcription uses the standard International Phonetic Alphabet (IPA) and is based on the Cambridge Dictionary of American English.

Symbol	Description	Example
ɑ:		calm [kɑ:m], heart [hɑ:rt], far [fɑ:r]
æ		act [ækt]
aɪ		dive [daɪv], cry [kraɪ]
e		met [met], lend [lend], pen [pen]
eɪ		say [seɪ], weight [weɪt]
ɚ	r-colored schwa	mother [mʌðɚ]
ɪ		fit [fɪt], win [wɪn]
i:		feed [fi:d], me [mi:]
oʊ		note [noʊt], coat [koʊt]
ɔ:		more [mɔ:r], cord [kɔ:rd]
ɔɪ		boy [bɔɪ], joint [dʒɔɪnt]
ʊ		could [kʊd], stood [stʊd]
u:		you [ju:], use [ju:z]
ɛ:		turn [tɛ:n], third [θɛ:rd]
ʌ		fund [fʌnd], must [mʌst]
ə	schwa	about [ə'baʊt]
ə	optional schwa	label ['leɪbəl] or ['leɪb]
i		very ['veri]
ɫ		handle ['hændɫ]
ŋ	velar nasal n (no g!)	bring [brɪŋ]
ʃ		ship [ʃɪp]
ʒ		measure [meʒɚ]
ɾ	american flapped t	butter ['bʌɾɚ]
θ	voiceless dental fricative	thing [θɪŋ]
ð	voiced dental fricative	then [ðen]
dʒ		joy [dʒɔɪ]

<sup>1</sup> **pronunciation** [prəˌnʌntsi'eɪʃən]: Aussprache



# Chapter 1

## Semiconductor Equations

In this section we will derive the set of equations needed to describe microelectronic devices. We reduce Maxwell's equations to the absolute minimum necessary and add equations that describe the behavior of the semiconductor material. Later in the lecture, this equation system will be solved numerically and approximated analytically.

### 1.1 The (Electrostatic) Field Side of the Problem

As a first step we show that the use of quasi-electrostatic approximations is **justified**<sup>1</sup>. Our condition for that is that the characteristic length of our system is considerably (by a factor of, say 10) smaller than the shortest electromagnetic wavelength present in the device. Given an upper limit of 100 GHz for the frequency of our electromagnetic field results in a wavelength of  $\lambda = c/f = 877\mu\text{m}$ ; the characteristic device dimension in microelectronics is about 1  $\mu\text{m}$ , confirming our **assumption**<sup>2</sup> of a quasi-static situation.

But what exactly does 'quasi-electrostatic approximation' mean? First of all, both the displacement current  $\partial\mathbf{D}/\partial t$  and the induction  $\partial\mathbf{B}/\partial t$  are **neglected**<sup>3</sup>. This simplifies Maxwell's equations to a much higher extent than one may assume at a first sight: Not only do two terms **vanish**<sup>4</sup>, but the former coupled system of differential equations **decouples**<sup>5</sup>, i.e. the direct coupling between the electric and the magnetic part of the field and their corresponding equations disappears. The only remaining interaction between the two field components is through the relation between the electric field  $\mathbf{E}$  and the electric current density  $\mathbf{J}$ , which causes a magnetic field  $\mathbf{H}$ . To simplify things further, even this magnetic field is usually neglected, making the two Maxwell equations for the magnetic field completely **obsolete**<sup>6</sup> in our applications. The fact that the right hand side of  $\nabla \times \mathbf{E} = -\partial\mathbf{B}/\partial t$  is zero permits the introduction of a *scalar potential*  $\psi$ , with the electric field being  $\mathbf{E} = -\nabla\psi$  (the minus sign is here for historical reasons; it is in no way physically justified, nor mandatory). Assuming linear, but possibly anisotropic and inhomogeneous material, the electric elasticity equation  $\mathbf{D} = \hat{\epsilon} \cdot \mathbf{E}$  holds, where the permittivity tensor  $\hat{\epsilon}$  in general depends on the spatial coordinates. Together with Gauss' law we get

$$\nabla \cdot (\hat{\epsilon} \cdot \nabla\psi) = -\rho \quad . \quad (1.1)$$

---

<sup>1</sup> **to justify** [dʒʌs.tɪ.fʌɪ]: rechtfertigen    <sup>2</sup> **assumption** [ə'sʌmp.tʃən]: Annahme    <sup>3</sup> **to neglect** [nɪ'gʌlekt]: vernachlässigen    <sup>4</sup> **to vanish** [væn.ɪʃ]: verschwinden    <sup>5</sup> **to decouple** [dɪkʌp.l]: entkoppeln    <sup>6</sup> **obsolete** [ɒb.sə'li:t]: hinfällig

The charge density  $\rho$  is composed of three components. Electrons possessing enough thermal energy to detach from the dopants and therefore obtain the ability to move around constitute the electron density  $n$ . Clearly, they leave behind a positively charged atom unable to move. But the electron's place can be taken by another electron, which itself leaves an empty electron position somewhere else. As more electrons move 'downstream' to fill the electron **vacancies**<sup>1</sup>, the vacancy itself moves 'upstream'. This constitutes a **fictitious**<sup>2</sup>, positively charged carrier type called a *hole*. The hole density, which is actually the density of 'missing' electrons, is denoted by  $p$ . Electrons and holes are present even in a perfectly pure semiconductor. In real semiconductors, impurities are always present. Moreover, they are often introduced **deliberately**<sup>3</sup> to control the electric conductivity of the material, in which case they are called *dopants*. The concentration of ionized **impurities**<sup>4</sup> and dopants is summed up in the concentration  $C$ . Putting all together **yields**<sup>5</sup> the space charge density

$$\rho = q(p - n + C) , \quad (1.2)$$

where  $q$  is the elementary charge. Finally, assuming the permittivity to be scalar and spatially independent (so it can be moved out of the nabla operator), we arrive at the equation

$$\nabla^2 \psi = q(n - p - C)/\epsilon , \quad (1.3)$$

which is known as *Poisson's equation*. In case  $\rho \equiv 0$ , one obtains *Laplace's equation*.

## 1.2 Continuity Equations

In pretty much the same way as the charge density consists of various contributions, the current density is **decomposed**<sup>6</sup> into an electron current density  $J_n$  and a hole current density  $J_p$ —the impurities and dopants are fixed in the crystal and therefore do not contribute to the current density. Hence, the charge continuity equation reads

$$\nabla \cdot (J_n + J_p) + q \frac{\partial}{\partial t} (p - n) = 0 . \quad (1.4)$$

(Note that since  $\partial C / \partial t \equiv 0$ , the respective term in  $\rho$  was dropped in the equation above.) This equation can formally be split into two equations; at the right hand side, a new term  $R$  is introduced:

$$\nabla \cdot J_n - q \frac{\partial n}{\partial t} = qR , \quad (1.5)$$

$$\nabla \cdot J_p + q \frac{\partial p}{\partial t} = -qR . \quad (1.6)$$

The interpretation of these two equations is as follows: Since charge particles actually can not be 'generated' or 'destroyed' (the right hand side of the continuity equation is zero), every additional electron that shows up in a left-alone semiconductor leaves an additional hole. Since these two have opposite charges, they appear in their respective continuity equation with opposite signs. The quantity  $R$  is the rate at which electron-hole-pairs are generated minus the rate at which they *recombine*. In equilibrium, the *net* recombination rate is zero (detailed balance), thus  $R \equiv 0$ ; but also out of equilibrium  $R$  is often neglected, because it considerably simplifies the problem. Depending on the type of the device, the inclusion of carrier generation and recombination models is mandatory for a realistic description.

---

<sup>1</sup> **vacancy** [veɪ.kəˈnɪ.sɪ]: freie Stelle, insbes. Gitterfreistelle    <sup>2</sup> **fictitious** [fɪkˈtɪʃ.əs]: fiktiv    <sup>3</sup> **deliberately** [dɪˈlɪb.ə.rət.li]: absichtlich    <sup>4</sup> **impurity** [ɪmˈpʊə.rɪ.ti]: Störstelle, Störatom    <sup>5</sup> **to yield sth.** [jɪːld]: etwas ergeben, etwas hervorbringen    <sup>6</sup> **to decompose** [diː.kəmˈpəʊz]: aufteilen, spalten

### 1.3 Charge Transport – The Drift-Diffusion-Model

The *structure* of our field problem is described by Poisson's equation and the two continuity equations. Neglecting  $R$ , five quantities are involved ( $\psi$ ,  $n$ ,  $p$ ,  $J_n$ , and  $J_p$ ), so we are short of two equations. They are provided by the microscopic model of the material considered, which describes how the field (mainly governed by Poisson's equation) acts on the charge particles (mainly governed by the continuity equations). The simplest model available is the so-called *drift-diffusion model*, which considers two distinct charge carrier transport mechanisms: Charge carrier *drift* due to the presence of an electric field, which is the customary transport mechanism in ordinary conductors, e.g. metals; and charge carrier *diffusion*. Diffusion is a fundamental process, which tries to establish a thermodynamic equilibrium in an initially imbalanced physical system. In our case, the physical system is the spatially non-constant **distribution**<sup>1</sup> of charge particles ( $n$  and  $p$ , where  $\nabla n \neq 0$  and  $\nabla p \neq 0$  in general). The thermodynamic equilibrium would be a situation in which  $\nabla n = \nabla p = 0$ , i.e. the carriers are **evenly**<sup>2</sup> distributed in the crystal. This situation is established through carrier migration from areas with high concentration to areas where the concentration of particles is lower. Mathematically speaking, the particles move in the opposite direction of the concentration gradient.

The drift component is expressed using the concept of *carrier mobility*, which is the proportionality factor between field strength and (average) carrier velocity. Denoting the individual mobilities for electrons and holes by  $\mu_n$  and  $\mu_p$ , we have

$$v_n = -\mu_n E \quad \text{and} \quad v_p = \mu_p E . \quad (1.7)$$

(Note the minus sign, because electrons with their negative charge travel – for historical reasons – *against* the field direction!). Charge carriers moving with some average velocity  $v$  constitute an electric current density  $J$ , whose magnitude is proportional not only to the velocity itself, but also to the *absolute number* of charge carriers per unit area that are on their way. But this number of carriers per unit area is directly related to the carrier density, which is  $n$  for electrons, each carrying a charge  $-q$  (mind the sign!) and  $p$  for holes (with charge  $q$  each). All in all we have

$$J_n^{\text{Drift}} = -qn v_n = qn \mu_n E \quad \text{and} \quad J_p^{\text{Drift}} = qp v_p = qp \mu_p E ; \quad (1.8)$$

again note the signs: Electrons move against the field direction, because of their negative charge. But for the same reason, the current they constitute points in the opposite direction of their velocity vector. In the end, both minus signs cancel. By introducing the conductivities  $\sigma_n = qn \mu_n$  and  $\sigma_p = qp \mu_p$  the relations above take the form of Ohm's law,

$$J_n^{\text{Drift}} = \sigma_n E \quad \text{and} \quad J_p^{\text{Drift}} = \sigma_p E . \quad (1.9)$$

The diffusion component is described using the notion of a *particle flux density*  $F$ , which is proportional to the negative gradient of the particle density. The proportionality factor is called the *diffusion coefficient*—since electrons and holes diffuse separately, two distinct diffusion coefficients  $D_n$  and  $D_p$  are involved:

$$F_n = -D_n \nabla n , \quad F_p = -D_p \nabla p . \quad (1.10)$$

The current densities are simply the flux densities multiplied with the individual charge of the carriers,

$$J_n^{\text{Diffusion}} = -q F_n = q D_n \nabla n , \quad J_p^{\text{Diffusion}} = q F_p = -q D_p \nabla p . \quad (1.11)$$

<sup>1</sup> **distribution** [dɪ'strɪb.juː.ʃən]: Verteilung    <sup>2</sup> **evenly** [iː.vən.li]: gleichmäßig



Close to equilibrium the carrier mobility and the diffusion coefficient are linked by the *Einstein relation*:

$$D_{n,p} = \frac{k_B T}{q} \mu_{n,p} = V_T \mu_{n,p} ; \quad (1.12)$$

the quantity  $V_T$  is referred to as the *thermal voltage* which is around 26mV at room temperature. The Einstein relation is only approximately valid for the non-equilibrium case we are interested in. The temperature  $T$  in (1.12) is the temperature of the electrons<sup>1</sup>. Although in many situations the **lattice**<sup>2</sup> temperature matches the electron temperature, there are cases (like in a MOSFET-channel) where electrons have a much higher “temperature” (i.e.<sup>3</sup> energy) than the lattice (*hot electrons*).

Note that carrier diffusion is not part of classical electrodynamics, where  $J$  and  $E$  are linked by Ohm’s law. The reason for this is simple: In classical electrodynamics, matter is either a conductor or an insulator. In the latter case, there’s no carrier transport at all, while in the former case the number of free charge carriers is assumed to be very high. This is the case for metals, where each atom at least contributes one electron to the electron gas, and for **ionic conductors**<sup>4</sup>, where each atom is a free charge carrier itself. In a system with that many free charged particles, a local disturbance in the carrier density is equilibrated almost instantaneously by the surrounding charges, causing the carrier densities to be virtually constant **throughout**<sup>5</sup> the whole body. The fact that in a semiconductor there is only a single electron for, say, every 100.000th atom makes the diffusion of carriers visible to a macroscopic observer.

## 1.4 Summary – The Basic Semiconductor Equations

Finally, we are able to put together a set of equations of first and second order that describes the basic behavior of a semiconductor:

$$\nabla^2 \psi = q(n - p - C)/\varepsilon , \quad (1.13)$$

$$\nabla \cdot J_n - q \frac{\partial n}{\partial t} = qR , \quad (1.14)$$

$$\nabla \cdot J_p + q \frac{\partial p}{\partial t} = -qR , \quad (1.15)$$

$$J_n = qn\mu_n E + qD_n \nabla n , \quad (1.16)$$

$$J_p = qp\mu_p E - qD_p \nabla p . \quad (1.17)$$

Substituting for the electric field  $E = -\nabla\psi$  and using the Einstein relation,

$$\nabla^2 \psi = q(n - p - C)/\varepsilon , \quad (1.18)$$

$$\nabla \cdot J_n - q \frac{\partial n}{\partial t} = qR , \quad (1.19)$$

$$\nabla \cdot J_p + q \frac{\partial p}{\partial t} = -qR , \quad (1.20)$$

$$J_n = -q\mu_n (n \nabla \psi - V_T \nabla n) , \quad (1.21)$$

$$J_p = -q\mu_p (p \nabla \psi + V_T \nabla p) . \quad (1.22)$$

<sup>1</sup> The temperature of electrons will be discussed in Chapter 2    <sup>2</sup> **lattice** [læt.ɪs]: Kristallgitter    <sup>3</sup> i.e. (lat. id est): “that is”    <sup>4</sup> **ionic conductor** [aɪ'ɒn.ɪk kən'dʌk.təʃ]: Ionenleiter    <sup>5</sup> **throughout** [θru:'aʊt]: durchweg, hindurch

The current relations can be inserted into the continuity equations, yielding the system of second order partial differential equations

$$\nabla^2 \psi = q(n - p - C) / \epsilon , \quad (1.23)$$

$$\nabla \cdot (\mu_n n \nabla \psi - \mu_n V_T \nabla n) + \frac{\partial n}{\partial t} = -R , \quad (1.24)$$

$$\nabla \cdot (\mu_p p \nabla \psi + \mu_p V_T \nabla p) - \frac{\partial p}{\partial t} = R . \quad (1.25)$$

These equations constitute the *drift-diffusion-model (DD-model)*, which, despite its limitations, is still the most widely used semiconductor model today. It was first derived by Van Roosbroeck back in 1950. [?]

## 1.5 Outlook – Lattice Heat-Flow

Basically all microphysic phenomena in solids are temperature dependent. Temperature enters the semiconductor equations directly via the thermal voltage  $V_T$  and indirectly via the recombination rate  $R$  and the temperature dependence of the mobilities  $\mu_n$  and  $\mu_p$ . Since the local temperature is an additional quantity in our semiconductor model, an additional equation must be provided. Heat enters the balance equations in the same way as charge does; the equation for the temperature distribution therefore pretty much looks like the charge continuity equation:

$$\nabla \cdot (\kappa \nabla T_L) - \rho c \frac{\partial T_L}{\partial t} = -H \quad (1.26)$$

and is of a structure similar as (1.25). Heat is redistributed via phonons, but there is no drift component since atoms are fixed within the lattice.  $T_L$  is the temperature, more precisely, the temperature of the lattice—in some cases, e.g. in the channels of field effect transistors where the electrons reach saturation velocity, they are attributed a higher temperature than that of the lattice (*hot electrons*), requiring a **distinction**<sup>1</sup> between lattice temperature and carrier temperature. The stationary temperature distribution is governed by the thermal conductivity  $\kappa$ , while the initial transient response to a change in the heat sources  $H$  is **determined**<sup>2</sup> by the mass density  $\rho$  ( $2328 \text{ VAs}^3 \text{m}^{-5}$  in silicon) and the specific heat  $c$  ( $703 \text{ m}^2 \text{s}^{-2} \text{K}^{-1}$  in silicon). Lastly, the heat generation term  $H$  provides the ‘back end’ of the coupling between the heat-flow and the drift-diffusion equations, since heat in the semiconductor is either generated by a first order Joule-term  $E \cdot J$  or by carrier recombination: Every recombination process sets free an energy amount at least equal to the semiconductor’s band gap energy  $E_g$ ; every generation of an electron-hole-pair withdraws  $E_g$  from the crystal. Therefore, we have

$$H \approx E \cdot J + R E_g . \quad (1.27)$$

Interestingly, a temperature gradient inside the crystal also **provokes**<sup>3</sup> a carrier diffusion process, yielding additional terms in the drift-diffusion current relations:

$$J_{n,\text{th}} = q D_{n,\text{th}} \nabla T_L , \quad J_{p,\text{th}} = -q D_{p,\text{th}} \nabla T_L ; \quad D_{n,p,\text{th}} = D_{n,p} / (2T_L) . \quad (1.28)$$

---

<sup>1</sup> **distinction** [dɪˈstɪŋk.ʃən]: Unterscheidung    <sup>2</sup> **to determine** [dɪˈteɪ.mɪn], **NOT** [dɪˈtɛr.mɪn]: bestimmen, festlegen

<sup>3</sup> **to provoke** [prəˈvʊk]: auslösen, bewirken

These current contributions are required for the modeling of the *thermoelectric effect*, aka<sup>1</sup> *Seebeck effect*. Also the reverse process can be observed, the transfer of heat by an electric current is known under the name *Peltier effect*.

---

<sup>1</sup> aka, also known as [ɔːlsou noʊ] æz]: auch bekannt unter, so genannt

## Chapter 2

# Carrier Transport

In this chapter we will look at carrier transport in greater detail, that is, how electrons and holes in a semiconductor respond to local and external forces. In the first chapter we have briefly introduced the drift-diffusion model, where we solve for the electron and hole concentrations. This means that electrons and holes in this model are **solely**<sup>1</sup> characterized by their locations in space, just as impurities have fixed positions in the lattice. As we know from semiconductor physics, an energy (or momentum) can be assigned to each electron, and the allowed energy levels are obtained from the band structure of a material. Instead of looking at the spatial distribution of electrons only, we will also briefly discuss electron energies, i.e. the population density of the energy bands. As we will see, the additional effort of considering the electron momentum pays off in the form of increased insight in the semiconductor device internals.

### 2.1 The Distribution Function

When no external electric field is applied and all transient processes have relaxed, we say that the semiconductor is in *thermal equilibrium*. In thermal equilibrium the electron gas is in equilibrium with the lattice. This does not mean, however, that the electrons rest in their equilibrium position: They rather fly around in a random fashion with their average energy equal to the thermal energy. Their average velocity, however, is zero because there is no preferred direction of movement and thus on average no electrical current flows. When an electric field is applied, the electrons are accelerated in **accordance**<sup>2</sup> with Newton's law of motion

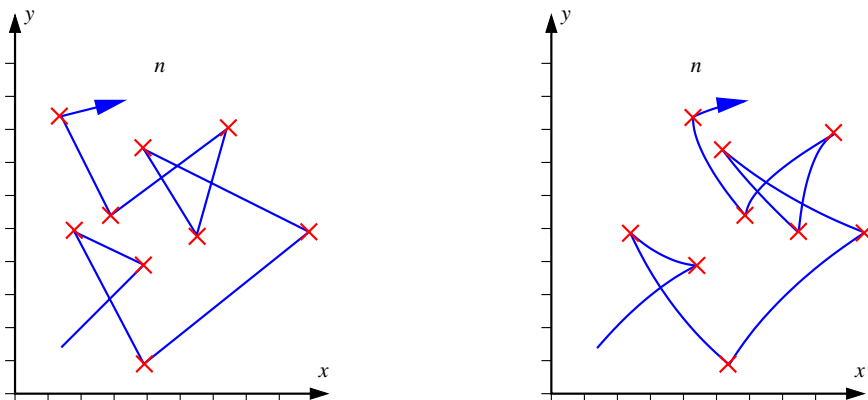
$$\frac{\partial p}{\partial t} = -qE . \quad (2.1)$$

Assuming that the electric field is kept constant and never turned off, it follows from (2.1) that the electrons reach an infinite velocity, or, taking relativistic effects into account, the speed of light. In a real semiconductor, however, an important process sets a limit to the maximum velocity: **scattering**<sup>3</sup>. In the models we will consider here, scattering processes are traced back to changes of the local potential seen by electrons. These changes can be caused (besides many other effects) by lattice vibrations which change the local band edge energy or by the electric field emerging from ionized impurities.

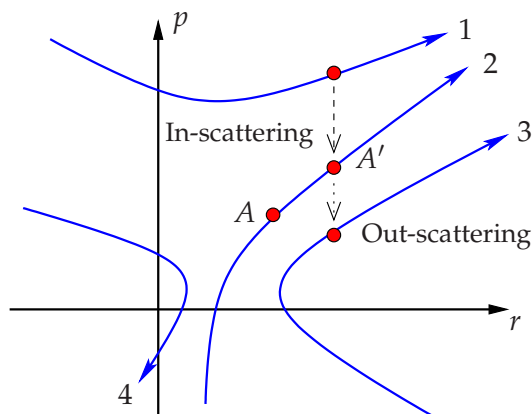
At any temperature above zero Kelvin, lattice atoms vibrate, **disturbing**<sup>4</sup> the perfect periodicity of the otherwise stationary lattice. This breaking of the periodicity leads to an energy exchange

---

<sup>1</sup> **solely** [soʊl.i]: lediglich    <sup>2</sup> **accordance** [ə'kɔ:r.dʌnts]: Übereinstimmung    <sup>3</sup> **to scatter** [skæɪ.ə]: streuen, zerstreuen    <sup>4</sup> **to disturb** [dɪ'stɜ:b]: stören, durcheinanderbringen



**Figure 2.1:** Thermal motion of a carrier without an externally applied field (left) and under the influence of a very strong electric field (right).



**Figure 2.2:** Illustration of trajectories in position-momentum space. Carriers move along a trajectory according to Newton's Laws. Occasionally they scatter to another trajectory. Scattering instantly changes the carrier's momentum, but does not affect its position.

of the electrons with the lattice. Lattice vibrations are characterized by *phonons*, that is why the resulting scattering process is termed *phonon scattering*. Where the electron movement is impeded by phonon scattering, the mobility decreases with increasing temperature. Consequently, a reduction of phonon scattering can be achieved by lowering the temperature, but since the ambient temperature is around 300 Kelvin in almost all cases, we have to live with this scattering mechanism.

The second possibility for a **deviation**<sup>1</sup> from perfect periodicity of the lattice potential is due to impurities: When there is a significant number of impurity atoms, the local lattice potential is distorted and scatters passing electrons. This scattering mechanism is called *impurity scattering* and can be controlled — unlike phonon scattering — rather easily: It can be eliminated by removing impurities from the material. However, a removal of impurities has other side-effects like increased resistance, which is often not desired. Since impurities are the essential building blocks of semiconductor devices and are used to form *pn* junctions, the scattering associated with impurities cannot be completely avoided. Nonetheless, under certain circumstances it is

<sup>1</sup> **deviation** [di:vi'eɪ.ʃən]: Abweichung

indeed required to find a suitable tradeoff between shorter space charge regions resulting due to higher doping and increased mobility resulting from lower doping concentrations.

Thus, an application-dependent tradeoff between increased mobility and reduced resistance is necessary. Most importantly, impurities are the fundamental basis for the functionality of all semiconductor devices.

Under the influence of an electric field the electrons accelerate and change their momentum according to (2.1). It is important to realize, however, that this change in momentum is normally small compared to the thermal energy. Therefore the electrons still move in a chaotic thermal way and only slightly change their momentum. Assuming that these particles behave like classical particles we can specify both their position and momentum at the same time. It is important to understand that in a more rigorous approach, where these particles are treated as quantum-mechanical wave packets, this is not possible because it violates Heisenberg's uncertainty principle

$$\Delta p \Delta r \geq \frac{\hbar}{2} .$$

In theory, we could solve Newton's equations of motion

$$\frac{d\mathbf{p}_i}{dt} = \mathbf{F}(\mathbf{p}, \mathbf{r}, t) + \mathbf{R}(\mathbf{p}, \mathbf{r}, t) \quad (2.2)$$

$$\frac{d\mathbf{r}_i}{dt} = \mathbf{u}_i(t) \quad i = 1, \dots, N \quad (2.3)$$

for each of the  $N$  carriers. The position of the carrier  $i$  is given as  $\mathbf{r}_i(t)$  and its momentum as  $\mathbf{p}_i(t)$ .  $\mathbf{R}(\mathbf{p}, \mathbf{r}, t)$  is a random force which introduces the effect of lattice vibrations and impurities into the model,  $\mathbf{F}(\mathbf{p}, \mathbf{r}, t)$  is the externally applied force, which is given for electrons as  $\mathbf{F}(\mathbf{p}, \mathbf{r}, t) = -q\mathbf{E}(\mathbf{r}, t)$  (for negligible magnetic fields).

The current flowing out of a device can be obtained by counting the carriers moving through the contact area. In a realistic device the number of free carriers  $N$  can be very large ( $N \gg 10^{20}$ ) and therefore only a representative sample of carriers can be considered in a practical implementation. However, even such a representative sample requires a large number of particles ( $N > 10^5$ ) and a direct solution of (2.2) and (2.3) is thus very time consuming. Such an approach is, despite its pitfalls, used in so-called *Monte Carlo simulations* to provide very accurate solutions of the problem.

Instead of considering a large number of carriers represented by their momentum and position ( $\mathbf{p}_i(t), \mathbf{r}_i(t)$ ) we just consider their statistical properties. We do this by defining a distribution function  $f(\mathbf{p}, \mathbf{r}, t)$  in such a way that  $f(\mathbf{p}, \mathbf{r}, t)d\mathbf{p}d\mathbf{r}$  gives the probability of finding a carrier with a momentum in the range  $[\mathbf{p}, \mathbf{p}+d\mathbf{p}]$  and a position inside the volume  $[\mathbf{r}, \mathbf{r}+d\mathbf{r}]$ . The distribution function considered here is therefore a *classical* concept as it defines both momentum and position of the particles and thus is in contradiction to Heisenberg's uncertainty principle.

The distribution function  $f(\mathbf{p}, \mathbf{r}, t)$  satisfies the *Boltzmann Transport Equation*

$$\frac{\partial f}{\partial t} + \mathbf{u} \cdot \nabla_{\mathbf{r}} f + \mathbf{F} \cdot \nabla_{\mathbf{p}} f = Q(f) , \quad (2.4)$$

where the *scattering operator*  $Q$  is given as

$$Q(f) = \sum_{\mathbf{p}'} f(\mathbf{p}') [1 - f(\mathbf{p})] S(\mathbf{p}', \mathbf{p}) - \sum_{\mathbf{p}'} f(\mathbf{p}) [1 - f(\mathbf{p}')] S(\mathbf{p}, \mathbf{p}') . \quad (2.5)$$

The equation describes the kinetic behavior of gases, in our case the electron gas.

Scattering is modeled as a number of (elastic or inelastic) collisions between particles and the lattice. The function  $S(\mathbf{p}', \mathbf{p})$  gives the transition rates for particles from the state denoted by the momentum  $\mathbf{p}'$  to the state of momentum  $\mathbf{p}$  after the scattering event and vice versa for  $S(\mathbf{p}, \mathbf{p}')$ . It has to be **emphasized**<sup>1</sup> that the scattering operator  $Q(f)$  in (2.5) takes the *Pauli Principle* into account: The terms  $f(\mathbf{p}')$  and  $f(\mathbf{p})$  require the present states  $\mathbf{p}'$  and  $\mathbf{p}$  to be occupied and the new states  $\mathbf{p}$  and  $\mathbf{p}'$  to be empty, hence the terms  $[1 - f(\mathbf{p})]$  and  $[1 - f(\mathbf{p}')$ . Note that the Pauli principle makes the Boltzmann transport equation *nonlinear* and thus much more complicated to solve. For non-degenerate semiconductors, where the carrier concentration is relatively small ( $f \ll 1$ ), the Pauli principle can be neglected and we obtain with

$$Q(f) \approx \sum_{\mathbf{p}'} f(\mathbf{p}')S(\mathbf{p}', \mathbf{p}) - \sum_{\mathbf{p}'} f(\mathbf{p})S(\mathbf{p}, \mathbf{p}') \quad (2.6)$$

a linear equation.

The solution of Boltzmann's transport equation provides excellent results, but it is very difficult to solve due to its higher dimensionality compared to the drift-diffusion model. For a full three-dimensional simulation, we have three spatial coordinates, three momentum coordinates and one time coordinate, thus a seven-dimensional simulation would be necessary. This means that using a **discretization**<sup>2</sup> with 100 unknowns in each coordinate direction, we would have to deal with a total of  $10^{14}$  points! If we assume  $7 \times 8$  bytes (i.e. eight byte for each coordinate) of memory for storing the location of each point, 5.600 Terabytes of memory would be needed for storing the locations of the points in the seven-dimensional simulation space only!

## 2.2 The Equilibrium Case

In the equilibrium case, the time **derivative**<sup>3</sup> in the Boltzmann transport equation (2.4) vanishes, hence  $f(\mathbf{p}, \mathbf{r}, t) = f(\mathbf{p}, \mathbf{r})$ . Furthermore, in equilibrium the scattering operator  $Q$  vanishes, because transitions from momentum  $\mathbf{p}'$  to momentum  $\mathbf{p}$  have the same probability as transitions from momentum  $\mathbf{p}$  to momentum  $\mathbf{p}'$  (principle of detailed balance). It can be shown that the equilibrium solution is the *Fermi-Dirac distribution*

$$f(\mathbf{p}, \mathbf{r}) = \frac{1}{1 + \exp\left(\frac{E_{\text{tot}}(\mathbf{p}, \mathbf{r}) - E_F}{k_B T_L}\right)}, \quad (2.7)$$

where  $E_F$  is the Fermi level and  $E_{\text{tot}}(\mathbf{p})$  the total carrier energy given as the sum of potential and kinetic energy of the carriers:

$$E_{\text{tot}}(\mathbf{p}) = E_{\text{pot}}(\mathbf{r}) + \mathcal{E}_{\text{kin}}(\mathbf{p}, \mathbf{r}). \quad (2.8)$$

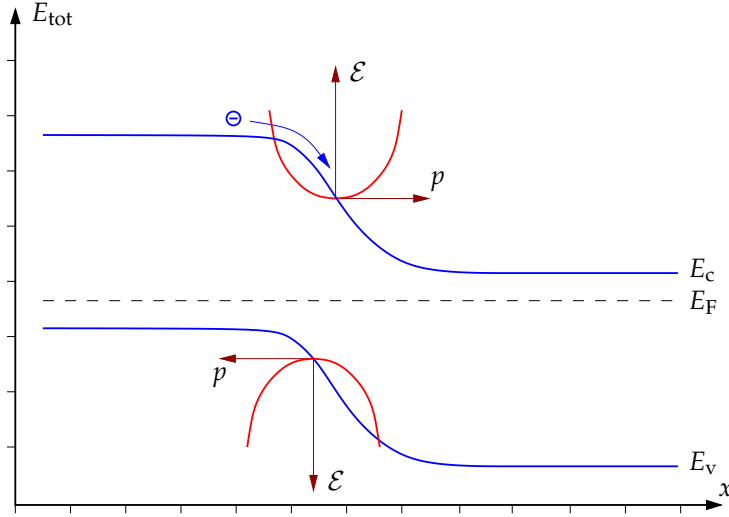
The potential energy of an electron is given by the band edge of the conduction band subject to an externally applied potential  $\psi$

$$E_{\text{pot}}(\mathbf{p}) = E_c = E_{c,0}(\mathbf{r}) - q\psi(\mathbf{r}) \quad (2.9)$$

for electrons, so that we find

$$E_{\text{tot}}(\mathbf{p}) = E_c(\mathbf{r}) + \mathcal{E}(\mathbf{p}, \mathbf{r}). \quad (2.10)$$

<sup>1</sup> **to emphasize** [emp.fə.saɪz]: betonen, hervorheben <sup>2</sup> **discretization** [dɪ'skri:tɪ'seɪ.ʃən]: Diskretisierung, vgl. Kapitel 3 <sup>3</sup> **derivative** [dɪ'rɪv.ə.tɪv]: Ableitung



**Figure 2.3:** Close to the band edges, which determine the potential energy (blue), the dispersion relation can be approximated by parabolas. The effective masses of electrons and holes give the curvature (which may also depend on the location in space) and determine the kinetic energy (red).

$\mathcal{E}(\mathbf{p}, \mathbf{r})$  is the kinetic energy given by the band structure. In many cases, this *dispersion relation* (or  $E - k$ -relation) is assumed to be spheric<sup>1</sup> (isotropic) parabolic<sup>2</sup> near the band edge, hence

$$\mathcal{E}(\mathbf{p}, \mathbf{r}) = \frac{|\mathbf{p}|^2}{2m^*(\mathbf{r})} = \frac{m^*(\mathbf{r})|\mathbf{u}|^2}{2}, \quad (2.11)$$

where  $m^*$  is the *effective mass*. Please keep in mind the analogy with the classical formula  $E_{\text{kin}} = p^2/2m = mv^2/2$ : The effective mass can thus be interpreted as a (fictitious) mass of the electron such that the classical energy formula holds. In more mathematical terms,  $m^*$  arises from the Taylor expansion of any arbitrarily complicated band structure and is defined in such a way that the parabolic relation for  $\mathcal{E}(\mathbf{p}, \mathbf{r})$  has the “correct” curvature near the band edge with respect to  $\mathbf{p}$ .

As already mentioned above, the Fermi-Dirac distribution honors the *Pauli Principle*: Each state can only be occupied by two electrons. Thus, by means of the Fermi-Dirac distribution it is possible to simulate degenerate semiconductors<sup>3</sup>.

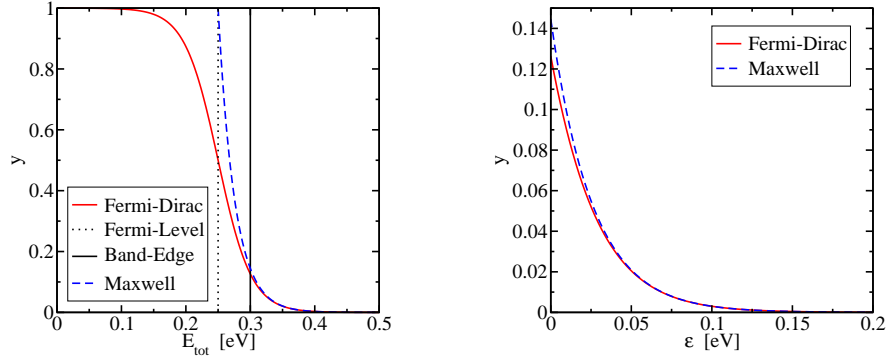
For non-degenerate semiconductors there holds  $E_c - E_F \gg k_B T_L$ , so (2.7) simplifies to the *Maxwell-Boltzmann distribution*

$$f(\mathbf{p}, \mathbf{r}) = \exp\left(\frac{E_F - E_{\text{tot}}(\mathbf{p})}{k_B T_L}\right) = \exp\left(\frac{E_F - E_c}{k_B T_L}\right) \exp\left(-\frac{\mathcal{E}(\mathbf{p})}{k_B T_L}\right) = A \exp\left(-\frac{\mathcal{E}(\mathbf{p})}{k_B T_L}\right). \quad (2.12)$$

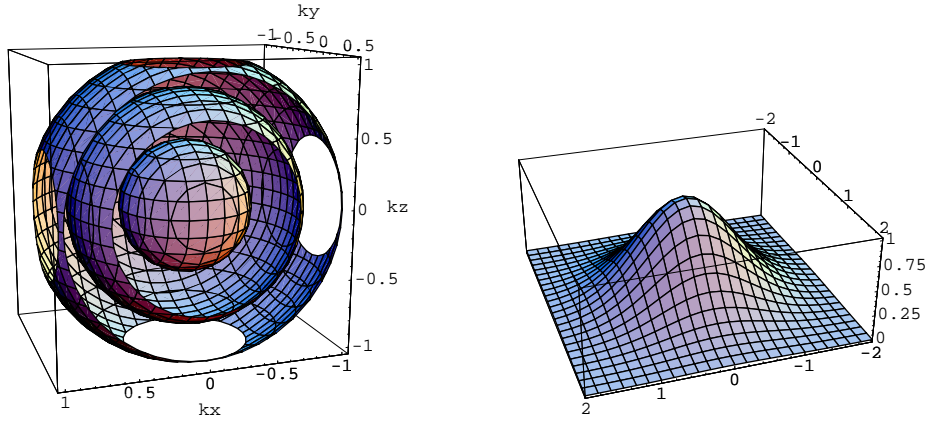
Unlike the Fermi-Dirac distribution, the Maxwell-Boltzmann distribution neglects the Pauli-Principle, thus it is not suitable for the simulation and modeling of degenerate semiconductors. From a formal point of view, the Maxwell-Boltzmann distribution is the solution of the Boltzmann equation if the scattering operator  $Q(f)$  in (2.6) is used instead of (2.5).

<sup>1</sup> i.e. independent of the angle, a function of  $|\mathbf{p}|$  only    <sup>2</sup> i.e. a term of power two    <sup>3</sup> A semiconductor is called *degenerate*, if its doping is so high that the individual dopant states in the band diagram merge to a so-called impurity band (located in the band gap of the intrinsic material). The material does not show the typical characteristics of a semiconductor anymore.





**Figure 2.4:** A comparison of the Fermi-Dirac and the Maxwell-Boltzmann distributions. It can be seen that the approximation in the conduction band (right) is quite good, but as soon as the Fermi level  $E_F$  approaches  $E_c$  (or  $E_v$ ), the Fermi-Dirac distribution has to be used.



**Figure 2.5:** Maxwell-Boltzmann distribution for parabolic bands (2.12). Isosurfaces are plotted in three dimensions (left). The two-dimensional plot better illustrates the underlying Gaussian normal distribution (right).

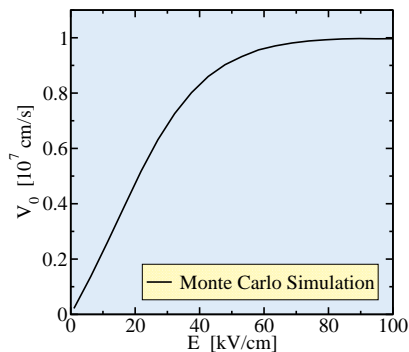
For parabolic bands with  $\mathcal{E}(\mathbf{p}) = |\mathbf{p}|^2/(2m^*)$  we obtain:

$$f(\mathbf{p}, \mathbf{r}) = A \exp\left(-\frac{|\mathbf{p}|^2}{2m^*k_B T_L}\right), \quad (2.13)$$

which corresponds to a normal distribution with mean zero and variance  $\sigma^2 = m^*k_B T_L$ . A comparison of the Fermi-Dirac and the Maxwell-Boltzmann distributions is given in Fig. 2.4, while from Fig. 2.5 we see that the average momentum is zero. However, a significant number of electrons occupy higher momentum states.

From the distribution of momentum we are now able to find the associated electron velocities: Let us consider the momentum  $\mathbf{p}_{\text{th}}$  such that  $f(\mathbf{p}_{\text{th}}) = A \exp(-1)$ , such that 85 percent (which corresponds to a deviation of  $\pm\sigma\sqrt{2}$  from the mean value of the normal distribution) of the particles are covered. From (2.12) we find the *thermal velocity*  $v_{\text{th}}$  as

$$|\mathbf{p}_{\text{th}}| = \sqrt{2m^*k_B T_L} \quad \Rightarrow \quad v_{\text{th}} = \frac{|\mathbf{p}_{\text{th}}|}{m^*} = \sqrt{\frac{2k_B T_L}{m^*}} \approx 10^7 \text{ cm/s}. \quad (2.14)$$



**Figure 2.6:** Results obtained by a Monte Carlo simulation show that the velocity of electrons due to an externally applied field saturates.

Let us compare the thermal velocity with typical velocities due to externally applied fields: With an electron mobility of  $\mu_n = 200 \text{ cm}^2/(\text{Vs})$  and an electric field of  $E = 50 \text{ kV/cm}$ , the relation  $v = \mu_n E$  yields  $v \approx 10^7 \text{ cm/s}$ . We will see in Chapter 7 that higher electric fields do not lead to higher velocities — the velocity saturates. This so-called *saturation velocity* is approximately equal to the thermal velocity in silicon, but in other semiconductors these two velocities differ. Results obtained from Monte Carlo simulations are shown in Fig. 2.6.

## 2.3 Boundary Conditions

The basic semiconductor equations are posed in a bounded domain  $\Omega \in \mathbb{R}^n$  ( $n = 1, 2, 3$ ). At the boundary  $\partial\Omega$  of  $\Omega$  appropriate boundary conditions need to be specified for the variables  $\psi$ ,  $n$ , and  $p$ . Two simple cases will be discussed in the following: Ideal Ohmic contacts and artificial boundary conditions. Models for another important case, Schottky contacts, can be found for instance in [?, ?]. Before we can discuss these boundary conditions, the issue of contact potentials, Fermi- and quasi-Fermi levels is briefly reviewed.

### 2.3.1 Quasi-Fermi Levels

Let us consider the following  $pn$ -diode consisting of three consecutive regions: The width of each region is  $L = 1 \mu\text{m}$  and the dopings are  $N_A = N_D = 10^{16} \text{ cm}^{-3}$ . In the stationary case all time derivatives vanish, hence we have to solve the following equations:

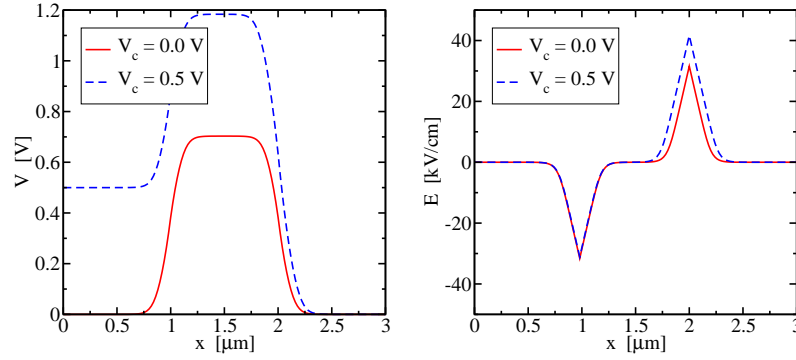
$$\nabla^2 \psi = \frac{q}{\epsilon} (n - p + N_A - N_D) \quad (2.15)$$

$$\nabla \cdot (+\mu_n k_B T_L \nabla n - q \mu_n n \nabla \psi) = +qR \quad (2.16)$$

$$\nabla \cdot (-\mu_p k_B T_L \nabla p - q \mu_p p \nabla \psi) = -qR \quad (2.17)$$

The first equation is the Poisson equation for the potential, while the second and third equations refer to electron and hole conservation and are obtained from (1.23)-(1.25) by setting the time derivative to zero. We have to supply boundary conditions for each second order partial differential equation<sup>1</sup>. Intuitively, we may use as boundary conditions for Poisson's equations

<sup>1</sup> In the present one-dimensional stationary setting, each unknown quantity depends on one variable only, thus we are dealing with a system of ordinary differential equations.



**Figure 2.7:** Simulation results for a  $pnp$ -diode for two different bias conditions. For opposite polarity of  $V_c$ , the results are reversed due to symmetry considerations.

the potentials applied to the contacts on either side of the device:

$$\psi(x = 0) = V_c, \quad \psi(x = 3L) = 0.$$

As boundary conditions for the continuity equations it appears natural to use the thermal equilibrium concentrations at the contacts (Ohmic contacts, that will be explained later):

$$p(x = 0) = N_A, \quad p(x = 3L) = N_A, \quad (2.18)$$

$$n(x = 0) = n_i^2 / N_A, \quad n(x = 3L) = n_i^2 / N_A. \quad (2.19)$$

Simulation results for this setting are given in Fig. 2.7. It is interesting to see that there is a **overshoot**<sup>1</sup> of the electrostatic potential in the center ( $n$ -region). This contradicts at first sight the maximum principle for elliptic PDEs, but on closer inspection we see that the maximum principle cannot be applied for the stationary semiconductor device equation. We will come back to this later.

In thermal equilibrium, the number of electrons and holes is given by Boltzmann statistics:

$$n_0 = N_c \exp\left(\frac{E_F - E_c}{k_B T_L}\right) = n_i \exp\left(\frac{E_F - E_i}{k_B T_L}\right), \quad (2.20)$$

$$p_0 = N_v \exp\left(\frac{E_v - E_F}{k_B T_L}\right) = p_i \exp\left(\frac{E_i - E_F}{k_B T_L}\right). \quad (2.21)$$

Here,  $N_c$  and  $N_v$  are the *effective density of states* or *band weights* of the conduction and valence band respectively. Multiplying these two equations with each other we obtain  $n_0 p_0 = n_i p_i = n_i^2$ . One has to keep in mind that the thermal equilibrium does *not* require the potential to be position-independent: Typically, we have

$$E_c = E_{c,0}(\mathbf{r}) - q\psi(\mathbf{r}), \quad (2.22)$$

$$E_v = E_{v,0}(\mathbf{r}) - q\psi(\mathbf{r}), \quad (2.23)$$

$$E_i = E_{i,0}(\mathbf{r}) - q\psi(\mathbf{r}). \quad (2.24)$$

<sup>1</sup> **overshoot** [əʊ.vəʃu:t]: die Überhöhung

Away from the thermal equilibrium, the situation is more complicated. Let us rewrite the current relation (1.21) in such a way that the current is the gradient of a some quantity:

$$\begin{aligned}
 J_n &= q\mu_n V_T \nabla n - q\mu_n n \nabla \psi \\
 &= q\mu_n n \left( V_T \frac{1}{n} \nabla n - \nabla \psi \right) \\
 &= q\mu_n n \left( V_T \frac{n_i}{n} \nabla \frac{n}{n_i} - \nabla \psi \right) \\
 &= q\mu_n n \left( V_T \nabla \ln \left( \frac{n}{n_i} \right) - \nabla \psi \right) \\
 &= q\mu_n n \nabla \underbrace{\left( V_T \ln \left( \frac{n}{n_i} \right) - \psi \right)}_{=: -\phi_n} .
 \end{aligned}$$

Writing the electron concentration  $n$  as a function of  $\phi_n$ , we obtain

$$n = n_i \exp\left(-\frac{\phi_n}{V_T}\right) \exp\left(\frac{\psi}{V_T}\right). \quad (2.25)$$

Comparing with (2.20), the introduction of *quasi-Fermi levels* is apparent. They are also called *imrefs*, which stands for *imaginary reference* and quite conveniently corresponds to the word Fermi spelled backwards. In contrast to the Fermi level, which is a unique global energy level, quasi-Fermi levels refer to a subsystem such as electrons and holes only, may depend on the location, and are introduced via the relationships

$$n = n_i \exp\left(\frac{E_{Fn} - E_i}{k_B T_L}\right) = n_i \exp\left(\frac{E_{Fn} - E_{i,0}}{k_B T_L}\right) \exp\left(\frac{\psi}{V_T}\right), \quad (2.26)$$

$$p = p_i \exp\left(\frac{E_i - E_{Fp}}{k_B T_L}\right) = p_i \exp\left(\frac{E_{i,0} - E_{Fp}}{k_B T_L}\right) \exp\left(-\frac{\psi}{V_T}\right). \quad (2.27)$$

While in thermal equilibrium  $n_0 p_0 = n_i^2$  holds, we find after multiplication of (2.26) with (2.27) that

$$pn = n_i^2 \exp\left(\frac{E_{Fn} - E_{Fp}}{kT}\right) \quad (2.28)$$

and we see that in equilibrium indeed  $E_{Fn} = E_{Fp} = E_F$  holds, so that the definition of quasi Fermi levels is consistent. The auxiliary quantity  $\phi_n$  relates to the quasi-Fermi levels as

$$-q\phi_n = E_{Fn} - E_{i,0}. \quad (2.29)$$

Since the intrinsic energy  $E_{i,0}$  is a globally constant level, we therefore arrive at

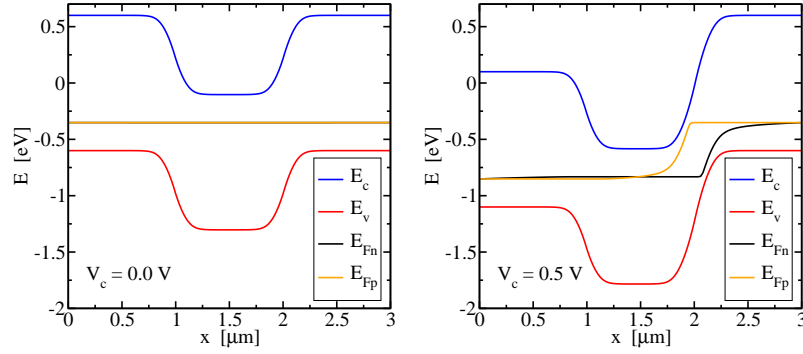
$$J_n = n\mu_n \nabla E_{Fn}, \quad (2.30)$$

while a similar **calculation**<sup>1</sup> for holes gives

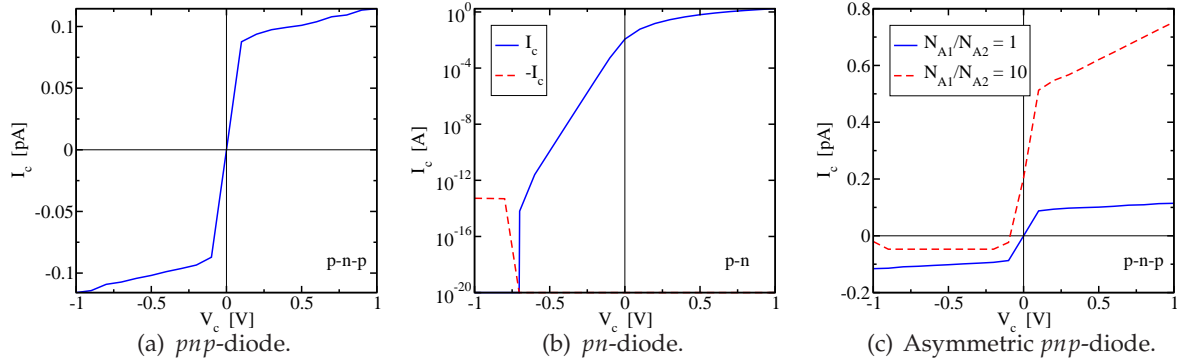
$$J_p = p\mu_p \nabla E_{Fp}. \quad (2.31)$$

Thus, in the most general case, *the current depends on the gradient of the quasi-Fermi-levels, not on the gradient of the potential!*

<sup>1</sup> **calculation** [kæɪ.lkjʊˈleɪ.ʃən]: Berechnung



**Figure 2.8:** Quasi-Fermi levels of the  $pnp$ -diode for the equilibrium and non-equilibrium case.



**Figure 2.9:** Current-voltage relations for simple devices.

Let us come back to the  $pnp$ -diode. If we assume that a solution of (2.15) - (2.17) is available, we can compute the quasi-Fermi levels in a post-processing step by rearranging (2.26) and (2.27):

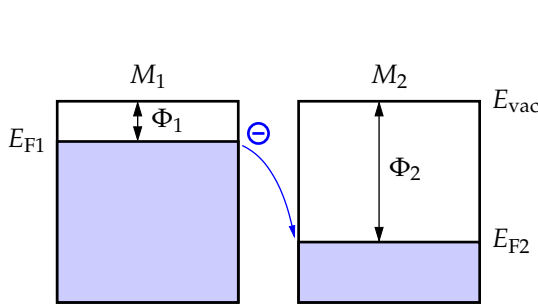
$$E_{Fn} = E_{i,0} - q\psi + qV_T \ln\left(\frac{n}{n_i}\right) = E_{i,0} - q\psi + k_B T_L \ln\left(\frac{n}{n_i}\right), \quad (2.32)$$

$$E_{Fp} = E_{i,0} - q\psi - qV_T \ln\left(\frac{p}{n_i}\right) = E_{i,0} - q\psi - k_B T_L \ln\left(\frac{p}{n_i}\right). \quad (2.33)$$

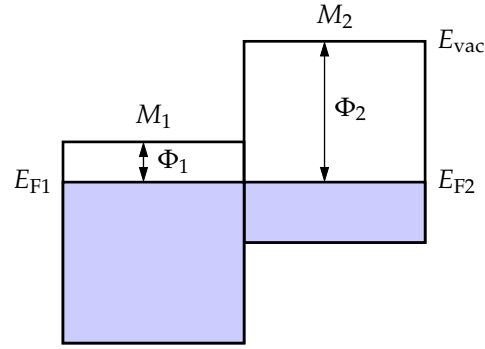
Since we are free to select a reference energy level, we select  $E_{i,0} = 0$  and find with  $E_g \approx 1.2\text{eV}$ ,  $E_{c,0} \approx E_{i,0} + E_g/2$  and  $E_{v,0} \approx E_{i,0} - E_g/2$  the quasi-Fermi levels shown in Fig. 2.8, where one should note the voltage drop on the second  $np$ -junction because of the reverse-bias there.

Since we can see the  $pnp$ -diode as a sequence of two  $pn$ -diodes with opposite conducting direction, we expect a current blockade for both forward- and reverse-bias. Looking at the current-voltage relation of the  $pnp$ -diode (Fig. 2.9) for applied voltages between  $-1\text{V}$  and  $+1\text{V}$  shows — despite the rather limited range of applicability of such a device — indeed good agreement with the expected behavior. However, if we remove the second  $p$ -region, the obtained  $pn$ -diode behaves like a battery (i.e. it sources current without any external voltage applied). The situation is similar if we choose different doping levels for the two  $p$ -regions in our  $pnp$ -diode (Fig. 2.9).

Such passive devices cannot behave like batteries, therefore a quite natural idea is to “shift” the current-voltage relation such that it crosses the origin. In the next section we will investigate the reasons for the observed “offset” and find that such a “shift” of the current-voltage relation



**Figure 2.10:** Materials  $M_1$  and  $M_2$  before joining, where each material has its own Fermi-level.



**Figure 2.11:** After joining materials  $M_1$  and  $M_2$ , band energies are shifted so that a single global Fermi-level is obtained.

is not deliberate but in fact required and has a sound physical background.

### 2.3.2 Contact Potentials

The Fermi level plays an important role in formulating equilibrium conditions when two materials are brought into contact. The combined system will be in thermal equilibrium only when  $E_F$  is the same in both parts, because all quantum levels at a given energy must have equal occupation probability at thermal equilibrium. If  $E_F$  in each material – relative to a common reference – is not initially equal before contact, then on contact there will be a flow of electrons from the material with the higher initial  $E_F$  to the material with the lower initial  $E_F$ . This electron flow will continue until equality of the Fermi energies of the two systems is achieved [?].

No electric field exists in the initially neutral materials. However, as each **electron**<sup>1</sup> crosses the **junction**<sup>2</sup>, it leaves behind a net charge of opposite polarity, and an electric field is thus established in the **vicinity**<sup>3</sup> of the junction, which tends to inhibit the movement of the electrons. This movement results in space-charge regions of different sign on either side of the junction, where a **dipole**<sup>4</sup> layer of finite thickness is formed which we are going to model with an idealized abrupt profile. An electrostatic potential change is then encountered when going from one material, through the junction, to the other material. In thermal equilibrium, the total potential drop in going from  $M_1$  to  $M_2$  is called the *contact potential*  $\psi_{12}$  of material  $M_1$  to  $M_2$  [?].

$$q\psi_{12} = E_{F1} - E_{F2} \tag{2.34}$$

Consider the case depicted in Fig. 2.10, where  $E_{F1} > E_{F2}$ . Material  $M_1$  will lose electrons and will thus become positively charged, while  $M_2$  will receive electrons and will thus become negatively charged. Therefore, the potential difference between  $M_1$  and  $M_2$  will be positive, as stated by (2.34). From this model it appears that there is a discontinuity of the potential at the metal-metal interface. This is, however, not the case: The potential changes continuously in a very thin (dipole) layer located at the interface.

Frequently one uses the *work-functions*  $q\Phi_1$  and  $q\Phi_2$ , which are positive quantities and give the distance between the Fermi energies and the *vacuum energy*  $E_{vac}$ , i.e. they correspond to the

<sup>1</sup> **electron** [ɪˈlektɹɒn], **NOT** [ˈɛlektɹɒn]: Elektron    <sup>2</sup> **junction** [ˈdʒʌŋkʃən]: die Sperrschicht    <sup>3</sup> **vicinity** [vəˈsɪnəti]: Umgebung    <sup>4</sup> **dipole** [ˈdaɪpəl]: Dipol

work needed to remove an electron from the lattice,

$$\begin{aligned} q\Phi_1 &= E_{\text{vac}} - E_{\text{F1}} , \\ q\Phi_2 &= E_{\text{vac}} - E_{\text{F2}} . \end{aligned}$$

With this definition the contact potential  $\psi_{12}$  can be expressed in terms of the work functions of the materials as (note the change in sign!)

$$\psi_{12} = \Phi_2 - \Phi_1 . \quad (2.35)$$

If  $N$  materials are considered in series, the potential  $\psi_{1N}$  between the first and the last material is expressed in terms of the work functions in the loop. Using (2.35) yields

$$\psi_{1N} = (\Phi_2 - \Phi_1) + (\Phi_3 - \Phi_2) + \dots + (\Phi_N - \Phi_{N-1}) .$$

It is clear that, with the exception of  $\Phi_1$  and  $\Phi_N$ , each contact potentials appears twice in the sum: once with a plus sign and once with a minus sign. Therefore

$$\psi_{1N} = \Phi_N - \Phi_1 . \quad (2.36)$$

Thus **no matter**<sup>1</sup> how many materials are in the chain, the electrostatic potential difference between its two ends depends *only* on the *first* and the *last* material (cf. electrochemical series).

When neutral  $n$ -type and  $p$ -type semiconductors are brought together to form a junction, at the interface a space-charge region forms, and a contact potential  $\psi_{\text{D}}$  appears.  $\psi_{\text{D}}$  is also referred to as *diffusion voltage* or *built-in potential* of the  $pn$ -junction, and is given by (2.34) as

$$q\psi_{\text{D}} = E_{\text{Fn}} - E_{\text{Fp}} .$$

Note that the built-in potential is defined as the contact potential between the  $n$ - and the  $p$ -layer, which is thus a positive quantity, but one commonly speaks of  $pn$ -junctions rather than of  $np$ -junctions.

According to (2.32) and (2.33) the Fermi energies of neutral semiconductors are

$$\begin{aligned} E_{\text{Fn}} &= E_{\text{i}} + k_{\text{B}}T_{\text{L}} \ln\left(\frac{n}{n_{\text{i}}}\right) , \\ E_{\text{Fp}} &= E_{\text{i}} - k_{\text{B}}T_{\text{L}} \ln\left(\frac{p}{n_{\text{i}}}\right) . \end{aligned}$$

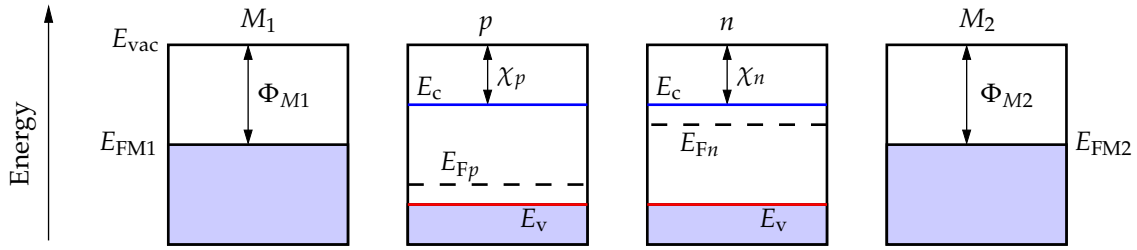
Setting  $n = N_{\text{D}}$  and  $p = N_{\text{A}}$  the contact potential of the  $pn$ -junction is obtained as

$$\psi_{\text{D}} = \frac{k_{\text{B}}T_{\text{L}}}{q} \ln\left(\frac{np}{n_{\text{i}}^2}\right) = V_{\text{T}} \ln\left(\frac{N_{\text{A}}N_{\text{D}}}{n_{\text{i}}^2}\right) . \quad (2.37)$$

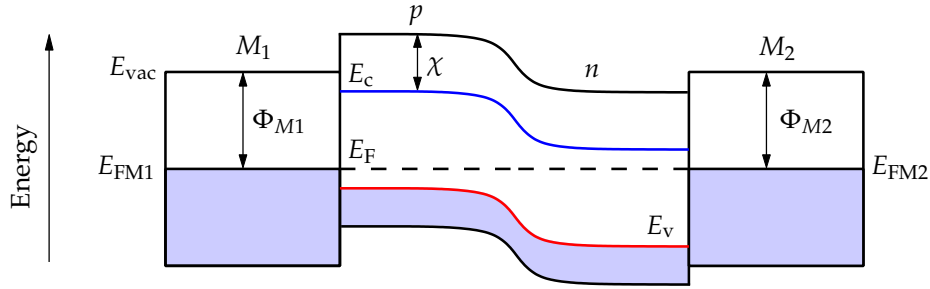
As a next step we consider such a  $pn$ -diode sandwiched between two metals  $M_1$  and  $M_2$  as shown in Fig. 2.12 and Fig. 2.13. The contact potential to the metal contacts is idealized as a discontinuity in the potential, even though there cannot be discontinuities of the potential. We can, however, think of a very thin layer at the interface in which the potential changes continuously, so that the assumption of a discontinuous potential can be justified from a modeling point of view. The contact potentials between the individual material layers are then obtained as

$$\begin{aligned} q\psi_{\text{M1p}} &= E_{\text{FM1}} - E_{\text{Fp}} , \\ q\psi_{\text{pn}} &= E_{\text{Fp}} - E_{\text{Fn}} = -q\psi_{\text{D}} , \\ q\psi_{\text{nM2}} &= E_{\text{Fn}} - E_{\text{FM2}} . \end{aligned}$$

<sup>1</sup> **no matter** [noʊ 'mæɪ.tər]: ganz egal



**Figure 2.12:** A  $pn$ -junction sandwiched between two metals  $M_1$  and  $M_2$  before contact.



**Figure 2.13:** A  $pn$ -junction sandwiched between two metals  $M_1$  and  $M_2$  after contact. Shown are the band edges inside the semiconductor. The contact potential to the metal contacts is idealized as a discontinuity in the potential and models an infinitesimal thin dipole layer.

By recalling that the difference between the two metal Fermi levels equals the contact potential

$$E_{FM2} - E_{FM1} = q\psi_{12} ,$$

we rewrite the contact potential  $q\psi_{nM2}$  as

$$q\psi_{nM2} = E_{Fn} - E_{FM2} = E_{Fn} - E_{FM2} + E_{FM1} - E_{FM1} = E_{Fn} - q\psi_{12} - E_{M1} .$$

As the reference energy is **arbitrary**<sup>1</sup> we choose  $E_{M1} = 0$  and obtain

$$\begin{aligned} \psi_{M1p} &= -E_{Fp}/q , \\ \psi_{nM2} &= +E_{Fn}/q - \psi_{12} . \end{aligned}$$

Since the metals used to contact a semiconductor are normally the same, we continue with the special case where  $M_1 = M_2 = M$  and we obtain

$$\begin{aligned} \psi_{M1p} &= -E_{Fp}/q , \\ \psi_{nM2} &= +E_{Fn}/q . \end{aligned}$$

The potential distribution inside a  $pn$ -junction is shown in Fig. 2.14. The potential  $\psi_L$  on the left side of the semiconductor is given by

$$\psi_{M1} - \psi_L = \psi_{M1p} ,$$

<sup>1</sup> **arbitrary** [ˈɑːrbətəri]: willkürlich, beliebig



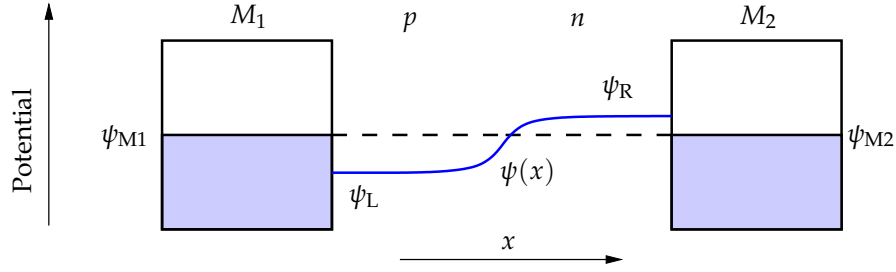


Figure 2.14: The potential distribution of a  $pn$ -junction.

and therefore we obtain

$$\psi_L = \psi_{M1} - \psi_{M1p} = \psi_{M1} + E_{Fp}/q.$$

Similarly we obtain for  $\psi_R$

$$\psi_R - \psi_{M2} = \psi_{nM2},$$

and thus

$$\psi_R = \psi_{M2} + \psi_{nM2} = \psi_{M2} + E_{Fn}/q.$$

Note that  $E_{Fp} < 0$  and  $E_{Fn} > 0$ .

In summary, the potentials at the left and right sides of the semiconductor are given by

$$\begin{aligned} \psi_L &= \psi_{M1} + E_{Fp}/q \\ \psi_R &= \psi_{M2} + E_{Fn}/q \end{aligned}$$

### 2.3.3 The Built-in Potential

Assume we have given an inhomogeneous electron concentration  $n(\mathbf{r})$  inside a semiconductor, which is in thermal equilibrium. We now ask for the potential distribution  $\psi(\mathbf{r})$  which is necessary to keep the electron distribution  $n(\mathbf{r})$  in this steady state.

In thermal equilibrium the current densities of either carrier type vanish,

$$\mathbf{J}_n = \mathbf{J}_p = \mathbf{0},$$

which means that an existing diffusion component has to be compensated exactly by a drift component. To guarantee  $\mathbf{J}_n = \mathbf{0}$ , it follows from (2.30), that the quasi-Fermi potential has to be constant across the semiconductor and we obtain

$$\psi_n = \psi - V_T \ln\left(\frac{n}{n_i}\right) = \text{const}.$$

In this equation one has again the freedom to choose a reference point. It is sound, although arbitrary, to define the quasi-Fermi potential to be zero if the distribution function represents thermal equilibrium, which is the case for a structure to which no external forces are applied.

With this convention, the potential required to retain a given equilibrium carrier concentration  $n(\mathbf{r})$  is found to be

$$\psi_{\text{bi}}(\mathbf{r}) = V_T \ln\left(\frac{n(\mathbf{r})}{n_i}\right). \quad (2.38)$$

$\psi_{\text{bi}}$  is termed *built-in potential*. Note that  $\psi_{\text{bi}}$  depends on the position  $\mathbf{r}$ . Equivalently, one obtains a built-in potential of the form

$$\psi_{\text{bi}}(\mathbf{r}) = -V_T \ln\left(\frac{p(\mathbf{r})}{n_i}\right) \quad (2.39)$$

by setting  $\psi_p = \text{const}$ . This equation can also be obtained from (2.38) by setting  $n(\mathbf{r}) = n_i^2/p(\mathbf{r})$ . Although it is already implicitly included in the calculation of this subsection, we have to emphasize again that the built-in potential is the solution for the potential (up to a constant shift) in the drift-diffusion model for the case that the current through the device is zero (i.e.  $J = 0$ ).

The exact built-in potential is a-priori unknown, but for a first guess one can derive an approximation by assuming charge neutrality:

$$p - n + N_D^+ - N_A^- = 0. \quad (2.40)$$

From (2.38) and (2.39), both  $n$  and  $p$  can be expressed as functions of the built-in potential so that (2.40) becomes

$$n_i \exp\left(-\frac{\psi}{V_T}\right) - n_i \exp\left(\frac{\psi}{V_T}\right) + N_D^+ - N_A^- = 0.$$

The potential  $\psi$  occurring in the above equation is the built-in potential we have been looking for. This equation can easily be solved and gives

$$\psi_{\text{bi}} = V_T \operatorname{arsinh}\left(\frac{N_D^+ - N_A^-}{2n_i}\right). \quad (2.41)$$

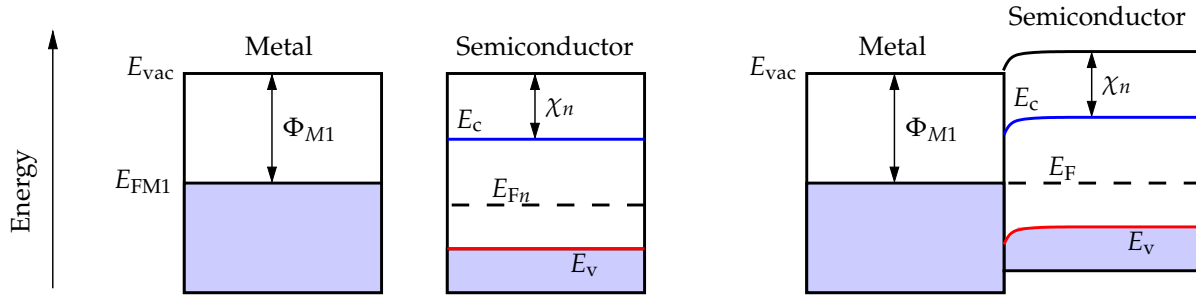
One should keep in mind that this derivation is strictly valid only for homogeneously doped semiconductors when  $\rho = 0$ . For non-uniform doping the potential will be a function of position and hence, an electric field will exist. This is possible only if there is a space-charge density, a fact which is contrary to (2.40). Nevertheless, the built-in potential for the inhomogeneous case is frequently approximated by (2.41) in a first step. The obtained solution is then used as an initial guess for the solution of the full non-linear drift-diffusion model by means of Newton's method.

## 2.4 Ohmic Contacts

There are two possible scenarios of a metal and a semiconductor being brought into contact: The first case is obtained when the Fermi level of the metal is lower than the Fermi level of the semiconductor. One then obtains the so-called *Schottky contact* with a non-linear, rectifying transfer characteristic (Schottky diode). On the other hand, contacting a metal with a semiconductor such that the Fermi level of the metal is higher than that of the semiconductor, a non-rectifying *Ohmic contact* is found in general<sup>1</sup>. The latter is shown in Fig. 2.15.

---

<sup>1</sup> It is possible to get an Ohmic contact even if the Fermi-level of the metal is lower than that of the semiconductor: With high doping concentrations, the resulting potential barrier is so thin that electrons tunnel through. This is the case for doping concentrations larger than  $10^{17} \text{cm}^{-3}$  (concentrations around  $10^{20} \text{cm}^{-3}$  are normally used).



**Figure 2.15:** Band diagram before (left) and after (right) contacting a metal and a semiconductor in the case of an Ohmic contact.

Ideal Ohmic contacts are the most commonly used boundary conditions for the simulation of semiconductor devices. The primary reason is that they are simple to implement in the form of Dirichlet boundary conditions for the potential and carrier concentrations, leading to very stable computations. From a physical point of view ohmic contacts are rather **crude**<sup>1</sup> and unlikely to be correct. Nevertheless, the hope is that there is minor impact of the contact region on the device behavior, so that ohmic contacts can be justified.

In the contact region any voltage drops are neglected in the simple model for an Ohmic contact. In the band diagram the Fermi level of the metal is continuously connected to the quasi-Fermi level of the semiconductor. Near the contact thermal equilibrium and vanishing space charge are assumed, so that we have (with  $C = N_D^+ - N_A^-$ )

$$np = n_i^2 \quad p - n + C = 0.$$

From the resulting quadratic equation we get

$$n = \frac{1}{2} \left( \sqrt{C^2 + 4n_i^2} + C \right),$$

$$p = \frac{1}{2} \left( \sqrt{C^2 + 4n_i^2} - C \right).$$

The important limiting cases for high and low doping are

$$\begin{array}{lll} C \gg n_i : & n = C, & p = n_i^2 / C, \\ C \ll n_i : & n = n_i, & p = n_i \end{array}$$

and vice versa for holes. The case  $C \gg n_i$ , which is usually fulfilled, leads to the boundary conditions we have used so far, thus we can keep them in most cases.

When we apply a bias  $V_c$ , we have for the boundary conditions of the potential at an ideal Ohmic contact

$$\psi = V_c - \psi_{MS},$$

<sup>1</sup> **crude** [kru:d]: grob, ungehobelt

where  $\psi_{\text{MS}}$  is the contact potential of the metal-semiconductor contact obtained from

$$q\psi_{\text{MS}} = E_{\text{FM}} - E_{\text{FS}} .$$

The Fermi level  $E_{\text{FS}}$  in the semiconductor can be found from the charge neutrality condition  $\rho = 0$  ( $np = n_i^2$  and thus  $E_{\text{Fn}} = E_{\text{Fp}}$ ) at the contact:

$$\begin{aligned} E_{\text{FS}} = E_{\text{Fn}} = E_{\text{Fp}} &= E_i + k_B T_L \ln\left(\frac{n}{n_i}\right) \\ &= E_i + k_B T_L \ln\left(\frac{1}{2n_i} \left(\sqrt{C^2 + 4n_i^2} + C\right)\right) \\ &= E_i + k_B T_L \operatorname{arcsinh}\left(\frac{C}{2n_i}\right) \\ &= E_i + q\psi_{\text{bi}} . \end{aligned} \tag{2.42}$$

This finally gives the doping concentration dependent contact potential

$$\psi_{\text{MS}} = \psi_{\text{MS}}(C) = \Phi'_{\text{MS}} - V_T \operatorname{arcsinh}\left(\frac{C}{2n_i}\right) , \tag{2.43}$$

with the zero-doping work function difference  $\Phi'_{\text{MS}} = (E_{\text{FM}} - E_i)/q$ .

Let us consider the two special cases of highly doped regions as will be used in all simulation examples in this lecture:

- **Highly doped  $n$ -region:** For  $N_D/n_i > 20$ , and with  $\sqrt{C^2 + 4n_i^2} \approx \sqrt{N_D^2 + 4n_i^2} \approx N_D$  we obtain

$$\begin{aligned} n &\approx N_D , \\ p &\approx n_i^2/n , \\ \psi &= V_c - \Phi'_{\text{MS}} + V_T \ln(N_D/n_i) . \end{aligned} \tag{2.44}$$

- **Highly doped  $p$ -region:** For  $N_A/n_i > 20$  and with  $\sqrt{C^2 + 4n_i^2} \approx \sqrt{N_A^2 + 4n_i^2} \approx N_A$  we obtain

$$\begin{aligned} p &\approx N_A , \\ n &\approx n_i^2/p , \\ \psi &= V_c - \Phi'_{\text{MS}} - V_T \ln(N_A/n_i) . \end{aligned} \tag{2.45}$$

As closing example we come back to the  $pn$ -diode where we have found shifted current-voltage characteristics in Sec. 2.3.1. Our aim is to find the device characteristics for an externally applied bias voltage  $V_c$ . We already know that applying  $V_c$  at the  $p$ -region and grounding the  $n$ -region does not lead to the expected results, because the built-in potential is not taken into account. So, let us denote with  $V_L$  the voltage applied to the  $p$ -region (on the left) and with  $V_R$  the voltage applied to the  $n$ -region, taking built-in potentials into account. From (2.44) and (2.45) we deduce

$$\begin{aligned} V_L &= V_c - \Phi'_{\text{MS}} - V_T \ln(N_A/n_i) , \\ V_R &= 0 - \Phi'_{\text{MS}} + V_T \ln(N_D/n_i) . \end{aligned}$$

This way, we have to modify the voltages applied to both electrodes. However, since the potential reference can be set arbitrarily, we may add  $\Phi'_{\text{MS}} + V_T \ln(N_A/n_i)$  to both  $V_L$  and  $V_R$  and get

$$V_L = V_c,$$

$$V_R = V_T \ln(N_D/n_i) + V_T \ln(N_A/n_i) = V_T \ln(N_A N_D/n_i^2) = V_D.$$

Thus, the simulation results in Fig. 2.9 for the  $pn$ -diode are indeed shifted by  $\psi_D$  because contact potentials have not been taken into account.

## Chapter 3

# Basics of Numerical Analysis

So far we have considered the continuous mathematical description of physical quantities in a semiconductor. However, for numerical simulation these continuous quantities have to be discretized. This chapter is **devoted**<sup>1</sup> to the basic numerical principles such as the discretization of derivatives and the solution of linear and non-linear systems of equations. We will **readdress**<sup>2</sup> these topics many times throughout this course. A direct translation of vector operators (grad =  $\nabla$ , div =  $\nabla \cdot$  and curl =  $\nabla \times$ ) from their continuous to their numerical representations can be found in Appendix B. In the following chapters we are going to extend the methods shown here with more involved techniques, until we can finally handle the complete set of the semiconductor equations on complex geometries.

### 3.1 Introduction to Finite Differences

The goal of *finite differences methods* is to approximate differential equations by a system of algebraic equations. This process, known as discretization, involves replacing the derivatives of quantities with differences of the same quantities, **evaluated**<sup>3</sup> at discrete locations.

To introduce numerical differentiation we will make use of variables that are set up as *arrays*. A continuous function is sampled at discrete points such as in Fig. 3.1. Recalling the Nyquist-Shannon sampling theorem in signal processing, the distance between two **adjacent**<sup>4</sup> discrete points is an important parameter for the quality of the numerical approximation.

For the hands-on part we will use SGFRAMEWORK, a **handy**<sup>5</sup> simulation environment developed for the simulation of semiconductor devices. The underlying equations are specified directly as code and a built-in solver computes a solution of the specified system of equations. The equations have to be supplied in a discretized form, which means that we have to apply a discretization method to the continuous formulation. One of the features of SGFRAMEWORK that makes it particularly useful is that it can solve simultaneous equations when the quantity that is being calculated is represented by an array. For example, if we want to find the potential  $\psi$  inside a rectangular simulation domain, we would have to settle for finding it at a set of points in the interior. The values of  $\psi$  at these points may be represented by an array. The points will form a grid in the simulation domain.

If the independent coordinate **employed**<sup>6</sup> is  $x$  and each grid point is at  $x_i = i\Delta x$ , where  $i$  is

---

<sup>1</sup> to devote [dr'vov.tɪd]: widmen    <sup>2</sup> to readdress [rɪ:ə'dres]: sich nochmal zuwenden    <sup>3</sup> to evaluate sth. [r'væl.ju.ert]: etwas auswerten    <sup>4</sup> adjacent [ə'dʒeɪ.sənt]: benachbart, angrenzend    <sup>5</sup> handy [hæ.n.dɪ]: praktisch, geschickt    <sup>6</sup> to employ [ɪm'plɔɪ]: einführen, einsetzen

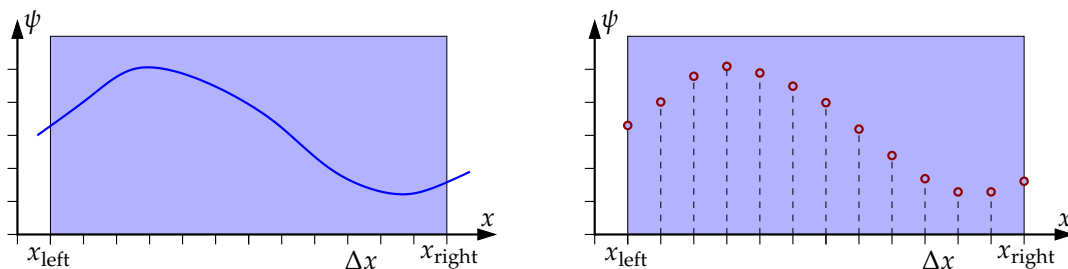


Figure 3.1: A continuous function (left) is represented by a discretized version (right).

an integer and  $\Delta x$  is the distance between two adjacent grid points, we can *label (address)* each  $x$ -value by its  $i$ -value (Fig. 3.1). If the problem is one-dimensional, we only use one coordinate  $x = i\Delta x$ , and the array for the quantity is simply  $\psi[i]$ . From the definition of the derivative

$$\frac{d\psi(x)}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta\psi}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{\psi(x + \Delta x) - \psi(x)}{\Delta x} \quad (3.1)$$

we can find an approximate expression for  $d\psi/dx$  at the point  $x_i = i\Delta x$ . This leads in a natural way to the one-sided differences

$$\left. \frac{d\psi(x)}{dx} \right|_{x=i\Delta x} \approx \frac{\psi[i+1] - \psi[i]}{\Delta x} \quad (\text{right-sided difference}) \quad (3.2)$$

and

$$\left. \frac{d\psi(x)}{dx} \right|_{x=i\Delta x} \approx \frac{\psi[i] - \psi[i-1]}{\Delta x} \quad (\text{left-sided difference}) \quad (3.3)$$

Taking the mean value of these two schemes (see Fig. 3.2),

$$\left. \frac{d\psi(x)}{dx} \right|_{x=i\Delta x} \approx \frac{\psi[i+1] - \psi[i-1]}{2\Delta x} \quad (\text{central difference}) \quad (3.4)$$

is obtained. Its name is derived from the fact that it uses values of  $\psi$  at points on either side of the point  $i$ .

One-sided approximations are usually less **accurate**<sup>1</sup> than the central difference scheme, since they give weight to one side of the point  $i$  and ignore the other side, but sometimes they are useful, especially at the boundary of a domain. The replacement of the derivative by the difference's quotient leads to *approximation errors*, which will become larger the further apart the

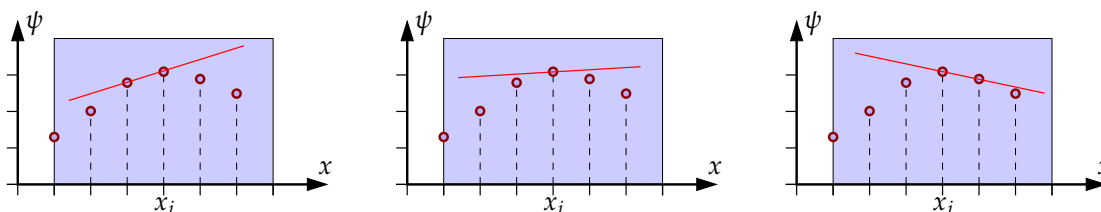


Figure 3.2: Comparison of various numerical difference schemes for  $d\psi/dx$  at the point  $x_i$ : left-sided differences (left), central differences (middle), right-sided differences (right).

<sup>1</sup> **accurate** [æk.jʊ.rət]: genau

discrete points are. On the other hand,  $\Delta x$  cannot be made too small due to the finite representation of floating point numbers in a computer.

As an example let us consider the calculation of the electric field from a given potential. In the one-dimensional case the electric field is defined as

$$E(x) = -\nabla\psi(x) = -\frac{d\psi(x)}{dx} . \quad (3.5)$$

Discretizing (3.5) using right-sided differences, we obtain

$$\begin{aligned} E[0] &= -(\psi[1] - \psi[0])/\Delta x \\ E[1] &= -(\psi[2] - \psi[1])/\Delta x \\ &\vdots \\ E[N-1] &= -(\psi[N] - \psi[N-1])/\Delta x . \end{aligned}$$

This can be written in matrix form

$$\begin{pmatrix} E[0] \\ E[1] \\ \vdots \\ E[N-1] \end{pmatrix} = \frac{1}{\Delta x} \begin{pmatrix} 1 & -1 & 0 & 0 & \cdots \\ 0 & 1 & -1 & 0 & \cdots \\ \vdots & \ddots & \ddots & \ddots & \ddots \\ 0 & \cdots & 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} \psi[0] \\ \psi[1] \\ \vdots \\ \psi[N-1] \\ \psi[N] \end{pmatrix} . \quad (3.6)$$

or in a more compact form as  $E = A\psi$ . Note that the dimension of the matrix  $A$  is  $N \times (N + 1)$  and that we cannot evaluate  $E[N]$  with right-sided differences because there is no right neighbor at  $N + 1$ . We could evaluate  $E[N]$  by using left-sided differences as well. This is done in the next example, where left- and right-sided differences are used at the boundary points  $i = 0$  and  $i = N$  and central differences elsewhere. We obtain

$$\begin{aligned} E[0] &= -(\psi[1] - \psi[0])/\Delta x \\ E[1] &= -(\psi[2] - \psi[0])/(2\Delta x) \\ &\vdots \\ E[N] &= -(\psi[N] - \psi[N-1])/\Delta x \end{aligned}$$

or, in compact form

$$\begin{pmatrix} E[0] \\ E[1] \\ \vdots \\ E[N-1] \\ E[N] \end{pmatrix} = \frac{1}{2\Delta x} \begin{pmatrix} 2 & -2 & 0 & 0 & 0 & \cdots \\ 1 & 0 & -1 & 0 & 0 & \cdots \\ 0 & 1 & 0 & -1 & 0 & \cdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \\ 0 & \cdots & 0 & 1 & 0 & -1 \\ 0 & \cdots & 0 & 0 & 2 & -2 \end{pmatrix} \begin{pmatrix} \psi[0] \\ \psi[1] \\ \vdots \\ \psi[N-1] \\ \psi[N] \end{pmatrix} . \quad (3.7)$$

In matrix form we again have  $E = A\psi$ , but this time  $E$  is of dimension  $N + 1$  as it should be and  $A$  is a square  $(N + 1) \times (N + 1)$  matrix.

File `basics1.sg` shows an SGFRAMEWORK program that allows us to differentiate a one-dimensional function. To do this we use two arrays. We start with  $\psi[i]$ , the electrostatic potential. The electric field is given by (3.5), so if we know the values in the array  $\psi[i]$  we can find the values of the electric field array  $E[i]$  from it.



```

1  const PI = 3.14159265358979323846; // pi
2  const DIM = 101; // number of mesh points
3  const DX = 2.0*PI / (DIM-1); // mesh spacing (cm) into (DIM-1) parts
4
5  var x[DIM], psi[DIM], E[DIM];
6
7  // these equations could be assignment statements
8  equ E[i=0] -> E[i] = -(psi[i+1] - psi[i]) / (1.0*DX);
9  equ E[i=1..DIM-2] -> E[i] = -(psi[i+1] - psi[i-1]) / (2.0*DX);
10 equ E[i=DIM-1] -> E[i] = -(psi[i] - psi[i-1]) / (1.0*DX);
11
12 begin main
13   assign x [i=0..DIM-1] = i*DX;
14   assign psi[i=0..DIM-1] = cos(x[i]);
15   solve;
16   write;
17 end

```

source\_code/basics1.sg

SGFRAMEWORK can solve equation systems of the form  $x = Ay$ . In our special case  $y$  is given and  $x$  can be obtained via a simple matrix multiplication. In real applications, however,  $x$  is given and the solution  $y$  has to be obtained by formally inverting  $A$ , that is, solving the equation system  $x = Ay$ .

The `assign` command was used to initialize  $\psi[i]$ . An `assign` statement could have been used as well to find  $E[i]$  from  $\psi[i]$  since in this case all the **quantities**<sup>1</sup> on the right hand side of the equation are known, but for demonstration purposes this was not done here.

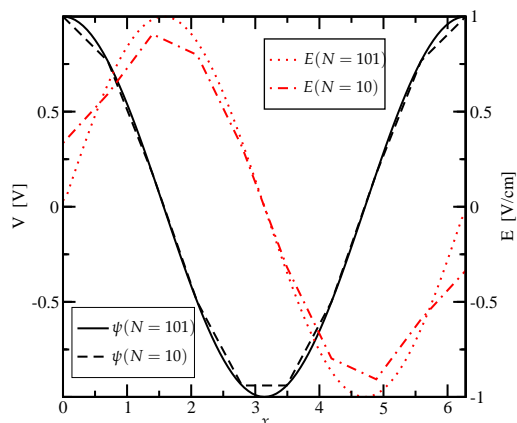
The `equ` statement implies that those values of  $E$  indicated are to be found from that equation. The `solve` statement tells the software to solve the equations appearing in the `equ` statement(s). The `equ` statement is a more powerful way of writing an equation than the `assign` statement. In an `assign` statement the right hand side of the equation must be an explicit expression in terms of known quantities, which will be evaluated to give the value of the left-hand side. The expression to the right of the arrow in an `equ` statement is evaluated to give the value of the quantity on the left of the arrow. The quantity to the left of the arrow must appear in the expression to the right. However, in an `equ` statement the expression to the right does not need to be an explicit expression for the quantity on the left. It could (for instance) be a non-linear equation for the quantity on the left which cannot be solved analytically. The left and right hand sides of the expression could both contain the quantity being solved for, or just one of them may contain it. That is, the expression to the right of the arrow can read  $f(x) = g(x)$  where  $x$  is the unknown. Finally, the expression may contain other unknown quantities which will only become known when a matrix is inverted, as is the case in the next example we will be looking at.

If  $E$  and  $\psi$  are plotted (Figs. 3.3, 3.4), it should be clear that  $E$  does look like  $-d\psi/dx$  with some small errors due to the way  $\psi$  and  $d\psi/dx$  are discretized using arrays at points  $x = i\Delta x$ .

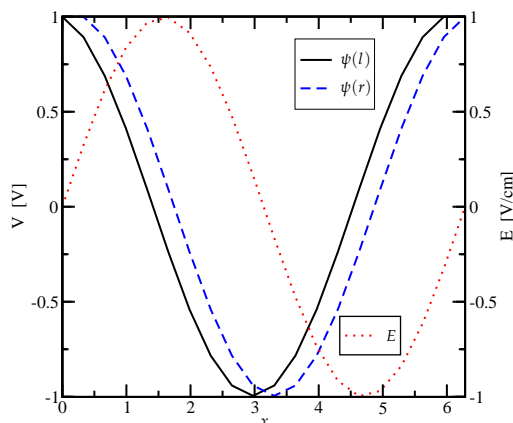
## 3.2 Numerical Solution of Differential Equations

We can immediately modify the equations we just solved and use the same expressions for the derivatives in a different way. Up to now we have used the equations to find the derivative,

<sup>1</sup> **quantities** ['kwɑ:n.tə.ti]: Größe, hier speziell: Matricelemente



**Figure 3.3:** Influence of the number of grid points on the final result (10 versus 101 grid points).



**Figure 3.4:** Influence of the discretization scheme on the final result (left-sided versus right-sided).

that is, we computed  $d\psi/dx$  from  $\psi$ . Instead, we can easily use the same approach to find the quantity being differentiated, for instance, we can compute the solution  $\psi$  of a differential equation from  $d\psi/dx$ . In the program `basics2.sg`, we solve a differential equation to illustrate the procedure. (In many of the numerical simulations given here, even small changes in the input file can cause a failure of convergence of the solution process. This in turn can lead to overflows or to domain errors. These will be reported as such, in most cases; but sometimes they will cause the simulation to crash.)

```

1  const PI = 3.14159265358979323846; // pi
2  const DIM = 101; // number of mesh points
3  const DX = 2.0 * PI / (DIM-1); // mesh spacing (cm)
4
5  var x[DIM], psi[DIM], E[DIM];
6
7  equ psi[i=0] -> psi[i] = 1.0; // boundary condition at the left side
8  equ psi[i=1..DIM-1] -> E[i] = -(psi[i]-psi[i-1]) / DX;
9
10 begin main
11  assign x[i=0..DIM-1] = i * DX;
12  assign E[i=0..DIM-1] = sin(x[i]);
13  solve;
14  write;
15 end

```

source\_code/basics2.sg

In this file, instead of initially knowing  $\psi$  and calculating  $E$  from it, suppose we know  $E$  and we **seek**<sup>1</sup>  $\psi$ . As a result we have to solve a differential equation of  $\psi$ . There is only one vector representing an ‘independent variable’ in this problem, namely  $x$ .

One thing to notice in this example is the boundary condition in this simulation, namely  $\psi[0] = 1$ . If the statement were **omitted**<sup>2</sup>,  $\psi[0]$  would be automatically set to zero, since all variables and array elements are initialized to zero until we explicitly change them. Then all the other  $\psi[i]$  can be computed from this. The equation in this input file finds  $\psi[i]$  from  $E[i]$

<sup>1</sup> to seek sth. [si:k]: etwas suchen    <sup>2</sup> to omit [ou'mit]: auslassen

and from  $\psi[i - 1]$ , and it starts at  $i = 1$ , so in case the system is solved in this way,  $\psi[0]$  has to be specified. Internally, SGFRAMEWORK automatically handles boundary conditions and then solves the equations appropriately.

The other thing **worth noticing**<sup>1</sup> here is that one of the less accurate one-sided differences was used to represent  $d\psi/dx$ . This type of difference can be easier to handle than the central difference used in the previous example since it allows us to start calculating values at one side of the range of  $x$  and work across the domain towards the other side. A central difference scheme would have complicated this process, since  $\psi[i + 1]$  would appear in the equation for  $\psi[i]$ , and  $\psi[i + 1]$  is not found until after  $\psi[i]$  is found, if we are working from one side to the other. The one-sided difference allows an explicit expression for  $\psi[i]$  to be written in terms of known quantities.

Since we use the `equ` statement to find  $\psi[i]$  the use of the one-sided difference was not actually necessary; the `equ` statement calls matrix routines which could have handled the central difference. On the other hand, a little algebra would have been required before we could use an assignment statement appropriate for central differences. The algebra would have allowed us to rewrite our equation so that the right hand side was explicitly known, with  $\psi[i]$  being on the left hand side.

### 3.3 The Second Order Derivative

The equation solved in the previous example was a first order ordinary differential equation (ODE), because the ‘highest’ derivative that appeared was a first order derivative. The electrostatic potential can also be determined from a given charge density  $\rho$  instead of an electric field  $E$ , using Poisson’s equation. In one dimension Poisson’s equation is

$$\frac{d^2\psi}{dx^2} = -\frac{\rho}{\epsilon} . \tag{3.8}$$

This is a second order elliptic ODE. To solve it we need to rewrite  $d^2\psi/dx^2$  using the differences between adjacent values of  $\psi$  the same way as we have done for  $d\psi/dx$ .

If we use the same approach as before, the derivative evaluated at the point  $x = i\Delta x$  is

$$\left(\frac{d\psi}{dx}\right)_{x=i\Delta x} \approx \frac{\psi[i + 1] - \psi[i - 1]}{2\Delta x} .$$

But if we use the same **reasoning**<sup>2</sup> to find the derivative of  $d\psi/dx$  then

$$\frac{d}{dx} \left(\frac{d\psi}{dx}\right)_{x=i\Delta x} \approx \left( \left(\frac{d\psi}{dx}\right)_{i+1} - \left(\frac{d\psi}{dx}\right)_{i-1} \right) \frac{1}{2\Delta x} .$$

To find the derivative at  $x = (i + 1)\Delta x$  we take the expression for the derivative at  $x = i\Delta x$  and add one to every occurrence of  $i$ . Similarly, at  $x = (i - 1)\Delta x$  we subtract one from  $i$  at each occurrence in the approximated derivative. Inserting these differences in the expression above gives

$$\begin{aligned} \left(\frac{d^2\psi}{dx^2}\right)_{x=i\Delta x} &\approx \left( \frac{\psi[i + 2] - \psi[i]}{2\Delta x} - \frac{\psi[i] - \psi[i - 2]}{2\Delta x} \right) \frac{1}{2\Delta x} \\ &= \frac{\psi[i + 2] - 2\psi[i] + \psi[i - 2]}{(2\Delta x)^2} . \end{aligned}$$

<sup>1</sup> **worth noticing** [wɜ:θ nɒv.tʃɪnɪŋ]: erwähnenswert    <sup>2</sup> **to reason** [ri:zən]: begründen, überlegen

This second order derivative is evaluated using points  $2\Delta x$  away from the central point  $x = i\Delta x$  and therefore its **denominator**<sup>1</sup> is  $(2\Delta x)^2$ . We can think of using an imaginary grid with half the spacing of the real grid. If we had evaluated  $d^2\psi/dx^2$  on the finer grid, we would have had taken points located  $\pm\frac{1}{2}\Delta x$  away from the center point  $x = i\Delta x$  and would have had divided by  $\Delta x^2$  instead of  $(2\Delta x)^2$ . This means that we can also use

$$\left(\frac{d^2\psi}{dx^2}\right)_{x=i\Delta x} \approx \frac{\psi[i+1] - 2\psi[i] + \psi[i-1]}{(\Delta x)^2}. \quad (3.9)$$

An alternative way to obtain the second derivative involves finding the first central derivative at each of two fictional points,  $i + 1/2$  and  $i - 1/2$ . These are found from the differences  $\psi[i + 1] - \psi[i]$  and  $\psi[i] - \psi[i - 1]$  respectively, each divided by  $\Delta x$ . To find the second derivative one then takes the difference between these first derivatives, and divides again by  $\Delta x$ , which again results in the above equation.

The only thing we still have to care about before being able to numerically solve a differential equation are the necessary boundary conditions. A **straightforward**<sup>2</sup> way to see this is to look at the finite difference equation. Whichever value of  $\psi$  we are solving for,  $\psi[i - 1]$ ,  $\psi[i]$  or  $\psi[i + 1]$ , we need to know the other two to find the value we want. For general  $i$  we for example find  $\psi[i + 1]$  from  $\psi[i]$  and  $\psi[i - 1]$ . To be able to start this iterative process, we need to know any two  $\psi[k]$  and  $\psi[j]$  for  $k \neq j$ . Clearly, the scheme can be applied directly if  $k = j \pm 1$  (allows the computation of  $\psi[k - 1]$  and  $\psi[k + 2]$ , which allows the computation of  $\psi[k - 2]$  and  $\psi[k + 3]$  and so on), but even the general case allows the computation of all  $\psi[i]$  by solving a system of equations. Our difference equation of the order two therefore requires indeed two boundary conditions<sup>3</sup>. Note that, in contrast to analytical solutions of linear differential equations or difference equations where a *general solution* can be formulated as a linear combination of elements from a solution basis with unknown coefficients, numerical solution by nature can only yield a single solution of a particular problem, and therefore *always requires boundary conditions*.

```

1  const DIM = 100;           // number of mesh points
2  const W   = 1.0;           // width of mesh (cm)
3  const DX  = W/(DIM-1);    // mesh spacing (cm)
4  const EPS = 8.854e-14;    // permittivity of free space (F/cm)
5
6  var x[DIM], psi[DIM], rho[DIM];
7
8  // this equation cannot be an assignment statement
9  // the boundary conditions are implicit (psi[0] = 0.0, psi[DIM-1] = 0.0)
10 equ psi[i=1..DIM-2] -> (psi[i+1]-2.0*psi[i]+psi[i-1])/sq(DX) = -rho[i]/EPS;
11
12 begin main
13   assign x[i=all] = i*DX;
14   assign rho[i=1*DIM/4..3*DIM/4] = 1.0e-10*sign(W/2.0-x[i]);
15   solve;
16   write;
17 end

```

source\_code/basics3.sg

In program `basics3.sg` the Poisson equation is discretized using the scheme (3.9). In this case, opposed to the examples involving first order differences above, one cannot write an explicit expression for  $\psi[i]$  in terms of known quantities. Hence an `equ` statement must be used. SGFRAMEWORK solves the system of equations using matrix techniques when instructed

<sup>1</sup> **denominator** [di'na: mæ.ni:tə]: Nenner    <sup>2</sup> **straightforward** [streit'fɔ: wəd]: unkompliziert    <sup>3</sup> to be more precise, on each characteristic line, two boundary or initial conditions must be given

to do so by the `solve` statement. The boundary conditions are formulated implicitly: Instead of explicitly specifying values for  $\psi$ , the initial setting of  $\psi[i] = 0$  provided by `SGFRAMEWORK` is used. Since the second order finite difference scheme for point  $i$  uses the neighboring points  $i - 1$  and  $i + 1$ , no equations must be given for points  $i = 0$  and  $i = DIM$ .

The ‘boxed’ shape of  $\rho$  in the example leads to a parabolic shape of the potential where  $\rho \neq 0$  (cf<sup>1</sup>. Fig. 3.5. Remember: A constant charge density gives a linearly increasing magnitude of the electric field by integrating once, and a quadratic dependence of the potential after the second integration). In the regions where  $\rho = 0$ , the potential increases linearly to satisfy the boundary conditions  $\psi[0] = \psi[DIM - 1] = 0$ .

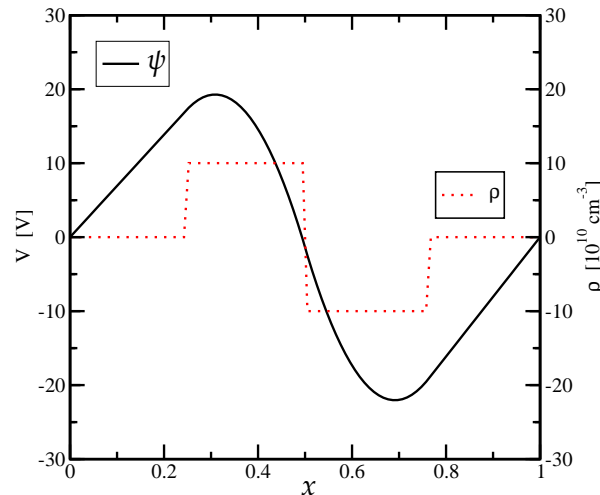


Figure 3.5: The solution of Poisson’s equation for a given density of charges.

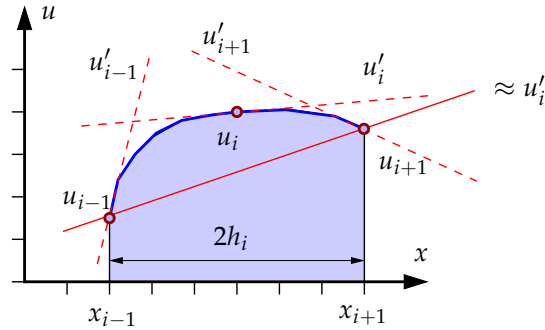
### 3.4 Accuracy of Different Discretization Schemes

In numerical analysis a continuous function is replaced by a sampled (*discretized*) one. This discretization happens on a *grid*. The quality of approximation of the continuous function by its sampled counterpart clearly depends on how fine the underlying grid is. To quantify how fine a grid is, the increment (*grid size*)  $h_i = x_{i+1} - x_i$  can be defined. Using the increment, grids can be categorized as follows:

- Uniform grid: All  $h_i$  are equal, so that the grid consists of equidistant points.
- Quasi-uniform grid: Consecutive increments are of almost the same magnitude. More precisely:  $h_{i+1} = h_i(1 + O(h_i))$ , where  $O(h_i)$  has the meaning of a small local perturbation—the grid distance  $h_i$  only slightly differs from its neighbors.
- Non-uniform grid: The  $h_i$  are arbitrary.

Even though the discretized version of a function contains less information than the continuous one, it is still possible to find numerical approximations of its derivatives. For a twice differentiable function  $u$ , and using the **abbreviations**<sup>2</sup>  $u_i = u(x_i)$  (and hence  $u_{i+1} = u(x_{i+1}) =$

<sup>1</sup> cf, to confer [kən'fɜːr]: vergleichen, konsultieren    <sup>2</sup> abbreviation [ə.bri.vi'eɪ.ʃən]: Abkürzung, Kurzwort



**Figure 3.6:** Derivative of  $u$  at different locations, and the approximation by finite differences.

$u(x_i + h_i)$ ) for the sake of **brevity**<sup>1</sup>, a Taylor **expansion**<sup>2</sup> of  $u$  at  $u_i$  gives

$$u_{i+1} = u_i + u'_i h_i + O(h_i^2) \quad \Rightarrow \quad u'_i = \frac{u_{i+1} - u_i}{h_i} + O(h_i) .$$

The first term on the right hand side is nothing more than the right-sided difference encountered in the previous section. The second term tells us that the *error* introduced by discretization **diminishes**<sup>3</sup> at least linearly with the ‘fineness’ of the underlying grid, i.e. doubling the number of grid points reduces the discretization error bound by a factor of two. We call a discretization scheme with such a behavior as *first order accurate*. Using Taylor expansions of  $u_{i+1}$  and  $u_{i-1}$  and assuming a uniform grid with grid size  $h$ ,

$$\left. \begin{aligned} u_{i+1} &= u_i + u'_i h + O(h^2) \\ u_{i-1} &= u_i - u'_i h + O(h^2) \end{aligned} \right\} \Rightarrow u'_i = \frac{u_{i+1} - u_{i-1}}{2h} + O(h) ,$$

it can easily be seen that the discretization error of a central difference scheme also depends linearly on the grid increment. Fig. 3.6 shows that the secant representing the central finite difference is an approximation to the ‘real’ derivative  $u'_i$ . The approximation error clearly depends on the grid size  $h_i$  and on the *curvature*, i.e. on the higher order derivatives, of  $u$ .

The general error introduced by the finite difference approximation can be analyzed by expanding  $u$  into a Taylor series at  $x_i$

$$\begin{aligned} u_{i+1} &= u_i + u'_i h_i + \frac{1}{2!} u''_i h_i^2 + \frac{1}{3!} u_i^{(3)} h_i^3 + O(h_i^4) , \\ u_{i-1} &= u_i - u'_i h_{i-1} + \frac{1}{2!} u''_i h_{i-1}^2 - \frac{1}{3!} u_i^{(3)} h_{i-1}^3 + O(h_{i-1}^4) . \end{aligned}$$

Subtracting both equations gives for a central difference approximation

$$u'_i = \frac{u_{i+1} - u_{i-1}}{h_i + h_{i-1}} + \frac{h_i - h_{i-1}}{2} u''_i + O\left(\frac{h_i^3 + h_{i-1}^3}{h_i + h_{i-1}}\right) .$$

The first order term in  $h_i$  disappears for  $h_i = h_{i-1}$  (uniform grid), so that the local **truncation error**<sup>4</sup>  $T$  is bounded by  $O(h^2)|u''(x)|$ , i.e. quadratically for a uniform grid with grid size  $h$ .

<sup>1</sup> **brevity** [brev.i.ti]: Kürze    <sup>2</sup> **expansion** [ik'spæn.tʃən]: hier: Entwicklung    <sup>3</sup> **to diminish** [dr'mi.niʃ]: abnehmen, abklingen    <sup>4</sup> **truncation error** [trʌŋ'kei.ʃən 'er.ə]: Abschneidefehler

Thus, if the grid size is reduced by a factor of two,  $T$  is usually reduced by a factor of four, which we will refer to as *second order accurate*. For quasi-uniform grids, the truncation error is found to be bounded by  $O(h^2)|u''| + O(h^2)|u^{(3)}|$ , while for non-uniform grids only a linear dependence ( $T < O(h)|u''(x)|$ ) between truncation error and grid size is obtained. A similar analysis can be done for the second order derivative. Starting from a Taylor series expansion at  $x_i$  and eliminating  $u'_i$ , the second order derivative is found as

$$u''_i = \frac{\left(\frac{u_{i+1} - u_i}{h_i}\right) - \left(\frac{u_i - u_{i-1}}{h_{i-1}}\right)}{\left(\frac{h_i + h_{i-1}}{2}\right)} + \left(\frac{h_i - h_{i-1}}{3}\right) u_i^{(3)} + O(h^2). \quad (3.10)$$

Consequently, for the truncation error  $T$  we find similar results: Quadratic dependence of the truncation error with the grid size for uniform and quasi-uniform grids and linear dependence for non-uniform grids.

### 3.5 Solution of Linear Systems of Equations

Although a call to `solve` in `SGFramework` is **sufficient**<sup>1</sup> to solve a system of equations, some deeper knowledge about the solution process is of advantage. Especially three-dimensional device simulations lead to very high numbers of unknowns, which may slow down the solution process considerably if an improper solution algorithm is chosen.

For linear systems, a vector of unknowns  $\mathbf{x} = (x_i)_{i=1}^n$  is sought such that

$$\mathbf{Ax} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} = \mathbf{b},$$

with system matrix  $\mathbf{A} = (a_{ij})_{i,j=1}^n$  and right hand side vector  $\mathbf{b} = (b_i)_{i=1}^n$ . Formally, the solution is given as  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ , but computing the matrix inverse explicitly is computationally too expensive and usually requires a lot of computer memory.

We will start with *direct solvers*, where the solution is found by a rather well known number of manipulations:

- *Gauss' method*. This method is one of the oldest methods. Matrix entries  $a_{ij}$  with  $i > j$  in the lower left triangle of  $\mathbf{A}$  are eliminated by means of linear combinations with rows  $i - 1, i - 2, \dots, 1$ . Finally, the solution vector is found using bottom-up insertion. The computational effort is  $O(2n^3/3)$ .
- *LU decomposition*. An invertible matrix  $\mathbf{A}$  whose **principal minors**<sup>2</sup> are all different from zero can be decomposed into

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ l_{21} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & \dots & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ 0 & u_{22} & \dots & u_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & u_{nn} \end{pmatrix} = \mathbf{LU},$$

<sup>1</sup> **sufficient** [sɒf.ɪ.ʃ.ənt]: ausreichend    <sup>2</sup> **principal minors** [prɪn.t.sɪ.pəl maɪ.nəʳ]: Hauptminoren

where  $L$  is a lower triangular matrix with ones in the main diagonal and  $U$  is an upper right triangle matrix. The solution  $x$  is found in two steps: First, the (triangular) system  $Ly = b$  is solved for  $y$ , then the triangular system  $Ux = y$  is solved for  $x$ .

The overall computational effort is again  $O(2n^3/3)$  and is due to the decomposition process. Since the solution of  $Ly = b$  and  $Ux = y$  requires only  $O(n^2)$  operations (triangular matrices!), the  $LU$ -decomposition is attractive in case several systems  $Ax_1 = b_1, Ax_2 = b_2, \dots, Ax_k = b_k$  with the same system matrix  $A$  need to be solved one after another, but cannot be solved simultaneously.

- *Cholesky decomposition.* A symmetric, positive definite matrix  $A$  can be decomposed into  $A = LL^T$ , where  $L$  is a lower left triangular matrix. The solution vector  $x$  is then found via a two-step process just as for the  $LU$ -decomposition.

The computational effort is only  $O(n^3/3)$ , so it is asymptotically twice as fast as the Gauss method.

The above methods work well with dense matrices, but show poor performance for *sparse* matrices. In a sparse matrix, most entries are equal to zero and thus need not be stored, reducing memory consumption tremendously. We already encountered sparse matrices in (3.6) and (3.7) with entries unequal to zero around the main diagonal only. Even though direct methods for such matrices exist, *iterative solvers* are preferred when some properties like diagonal dominance or positive definiteness of  $A$  can be assured. These methods include Jacobi iteration, Gauss-Seidel iteration, Steepest Descent Method and the Conjugate Gradient algorithm.

We will not go into further detail here, however, the reader should be aware of the existence of iterative solvers, whose asymptotic run-time behavior is usually superior to that of direct solvers.

### 3.6 Solution of Nonlinear Systems of Equations

For the solution of non-linear equations the most widespread method is *Newton's method*. It is best explained in one dimension (Fig. 3.8): For a point  $x_k$  where a function  $f$  is evaluated to  $f(x_k)$ , a reasonable guess  $x_{k+1}$  for the location of a **root**<sup>1</sup>  $y$  (such that  $f(y) = 0$ ) is the intersection of the tangent with the  $x$ -axis. The tangent can be obtained from a first order Taylor expansion of  $f$  in  $x_k$ :

$$f(x_{k+1}) \approx f(x_k) + f'(x_k)(x_{k+1} - x_k).$$

Since  $x_{k+1}$  should be a root of  $f$ , we set  $f(x_{k+1}) = 0$ , so that

$$0 \stackrel{!}{=} f(x_k) + f'(x_k)(x_{k+1} - x_k).$$

remains. After a rearrangement we find

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}.$$

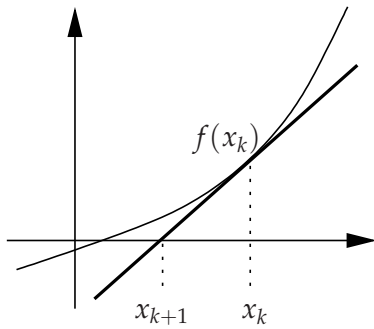
Newton's method in algorithmic form is thus:

1. Start with an initial guess  $x_0$  and  $k = 0$ .

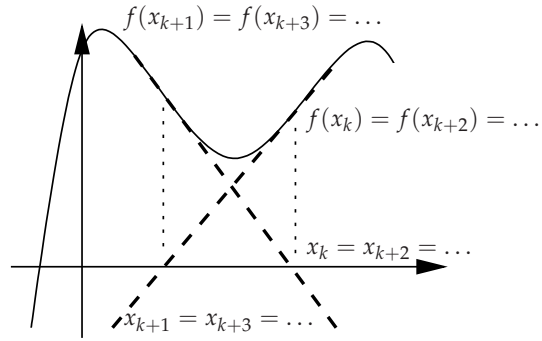
---

<sup>1</sup> **root** [ru:t]: Wurzel, Ursprung, hier: Nullstelle





**Figure 3.7:** One-dimensional motivation for Newton's method.



**Figure 3.8:** Newton's method may fail if the initial guess is too far away from a root.

$k$	$x_k$	$f(x_k)$	$\Delta x = -f(x_k)/f'(x_k)$
0	1.65	0.17186502845391	0.2967688297482
1	1.94676882974821	-0.04323343404123	-0.0498553585067
2	1.89691347124147	-0.00116331450812	-0.0014180411387
3	1.89549543010275	-9.52580161861505E-07	-1.1630679886549E-06
4	1.89549426703476	-6.41042774418565E-13	-7.8269248466524E-13
5	1.89549426703398	0.0	0.0

**Table 3.1:** Finding the root of  $f(x) = \sin(x) - 0.5x$  using Newton's method.

- Evaluate  $f(x_k)$  and stop if  $|f(x_k)| < \varepsilon$ .
- Evaluate  $f'(x_k)$  to set up the tangent at  $f(x_k)$  and find the intersection  $x_{k+1}$  with the abscissa. It is given as  $x_{k+1} = x_k - f(x_k)/f'(x_k)$ .
- Increase  $k$  (i.e. take the new approximation as the new guess) and go to 2.

Instead of checking the absolute value of  $f(x_k)$  in step 2 only, one can also check the magnitude of  $\Delta x_k = -f(x_k)/f'(x_k)$  compared to  $x_k$ .

Unfortunately, Newton's method fails if  $f'(x_k) = 0$ , which means that the tangent is parallel to the  $x$ -axis and no intersection can be found. If the function  $f(x)$  is continuously differentiable, its derivative does not vanish at  $\alpha$  and it has a second derivative at  $\alpha$  then the convergence is quadratic or faster, provided that the initial guess is in a neighborhood of  $\alpha$ . This means that the number of valid digits doubles with each step, as long as we start with a good initial guess. What a "good" initial guess actually is, strongly depends on the given problem.

Although Newton's method appears to be a root-finding algorithm only, it can be used to solve arbitrary non-linear equations: In order to solve  $f(x) = g(x)$  for  $x$ , rewrite the equation in implicit form as  $f(x) - g(x) = 0$ , set  $\tilde{f}(x) := f(x) - g(x)$  and apply Newton's method to  $\tilde{f}(x)$ .

For example, consider the equation  $\sin(x) = x/2$  and solve for  $x$ . With  $f(x) = \sin(x) - 0.5x$  we get  $f'(x) = \cos(x) - 0.5$ . The numbers finally resulting from Newton's method with initial guess  $x_0 = 1.65$  can be found in Tab. 3.1.

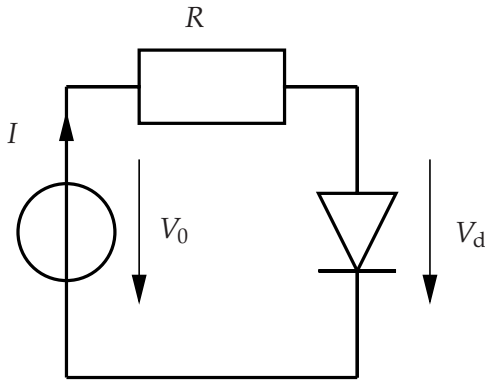


Figure 3.9: Example circuit for demonstration of Newton's method.

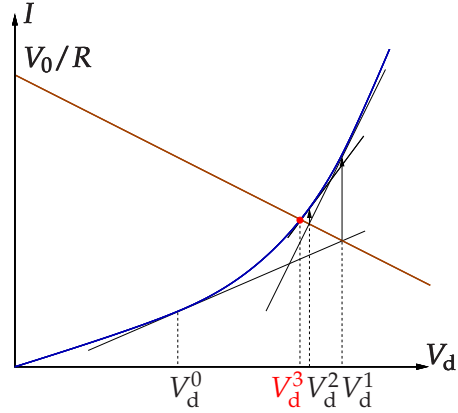


Figure 3.10: Newton's method for the solution of (3.11).

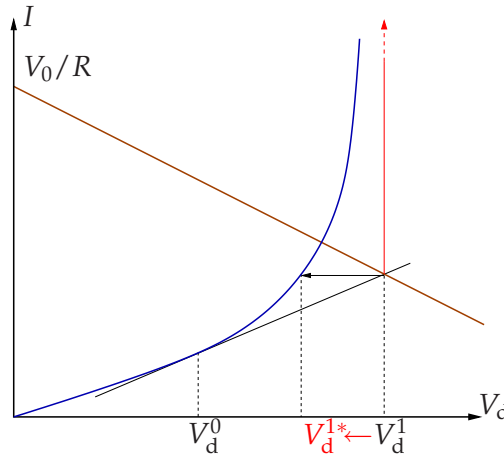


Figure 3.11: Newton's method may fail if the starting point is not chosen properly. The use of damping improves stability, but a good initial guess is still needed.

A more practical example is the determination of a current through a diode. Consider the circuit given in Fig. 3.9. The governing equation is

$$\frac{V_0 - V_d}{R} = I_d(V_d) = I_s(\exp(V/V_T) - 1), \quad (3.11)$$

which is non-linear in  $V_d$  with known  $I_s$ ,  $V_0$  and  $V_T$  because of the exponential current characteristic of the diode. In Fig. 3.10 the solution using Newton's method is demonstrated: Starting with the initial guess  $V_d^0$ , the tangent to the diode characteristic is drawn. Then, the **intersection**<sup>1</sup> with the load line (i.e. the right hand side of (3.11)) is determined, resulting in a new approximation  $V_d^1$ . This procedure is repeated until one is sufficiently close to the solution.

Nevertheless, the above example can easily be modified so that Newton's method fails: Assume that the diode characteristic has a pole instead of exponential behavior. Then, a starting point that is too far away from the true solution does not converge, as is illustrated in Fig. 3.11. The numerical stability of Newton's method can be improved by *damping*: With a damping factor  $d \in [0, 1]$  and notation as in the previous example, the next approximation  $V_d^{i+1}$  is now

<sup>1</sup> **intersection** [ɪn.tə'sek.ʃən]: Schnittmenge, Schnittpunkt

computed as

$$V_d^{i+1} = V_d^i - d \frac{f(V_d^i)}{f'(V_d^i)}. \quad (3.12)$$

If  $d$  is chosen too small, the number of iterations increases unnecessarily, while  $d$  chosen too large may cause the solution process to fail. A clever trade-off is thus necessary in practice.

The issue of initial solutions that are sufficiently close to the true solution is especially a considerable problem for complicated systems. As for semiconductor device simulation, an initial guess uses built-in potentials (Chapter 2) and assumes charge neutrality (like for Ohmic contacts, Chapter 2).

A generalization to systems of  $n$  non-linear equations is rather straightforward, but mathematical analysis becomes much harder. In order to solve the coupled equation system

$$\begin{aligned} F_1(\mathbf{x}) &= 0 \\ F_2(\mathbf{x}) &= 0 \\ &\vdots \\ F_n(\mathbf{x}) &= 0 \end{aligned}$$

for a solution vector  $\mathbf{x}$  of dimension  $n$ , introduce the vector-valued function

$$\mathbf{F}: \mathbf{x} \mapsto \begin{pmatrix} F_1(\mathbf{x}) \\ F_2(\mathbf{x}) \\ \vdots \\ F_n(\mathbf{x}) \end{pmatrix}.$$

In the algorithm given above, one then has to multiply with the inverse of the *Jacobian matrix*  $\mathbf{J}_F(\mathbf{x}_k)$  of dimension  $n \times n$  instead of dividing by  $f'(x_k)$ :

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{J}_F^{-1}(\mathbf{x}_k) \mathbf{F}(\mathbf{x}_k).$$

However, the inverse of the Jacobian is not computed explicitly, but the system

$$\mathbf{J}_F(\mathbf{x}_k) \mathbf{y} = \begin{pmatrix} \frac{\partial F_1}{\partial x_1}(\mathbf{x}_k) & \frac{\partial F_1}{\partial x_2}(\mathbf{x}_k) & \cdots & \frac{\partial F_1}{\partial x_n}(\mathbf{x}_k) \\ \frac{\partial F_2}{\partial x_1}(\mathbf{x}_k) & \frac{\partial F_2}{\partial x_2}(\mathbf{x}_k) & \cdots & \frac{\partial F_2}{\partial x_n}(\mathbf{x}_k) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_n}{\partial x_1}(\mathbf{x}_k) & \frac{\partial F_n}{\partial x_2}(\mathbf{x}_k) & \cdots & \frac{\partial F_n}{\partial x_n}(\mathbf{x}_k) \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} F_1(\mathbf{x}_k) \\ F_2(\mathbf{x}_k) \\ \vdots \\ F_n(\mathbf{x}_k) \end{pmatrix} = \mathbf{F}(\mathbf{x}_k)$$

is solved using a linear equation solver. To sum up, the vector-valued Newton's method reads as follows:

1. Start with an initial guess  $\mathbf{x}_0$  and  $k = 0$ .
2. Evaluate  $\mathbf{F}(\mathbf{x}_k)$  and stop if  $\|\mathbf{F}(\mathbf{x}_k)\| < \varepsilon$  (i.e. a sufficiently good approximation to the true solution is found).
3. Solve  $\mathbf{J}_F(\mathbf{x}_k) \mathbf{y} = \mathbf{F}(\mathbf{x}_k)$  for  $\mathbf{y}$ .
4. Set  $\mathbf{x}_{k+1} := \mathbf{x}_k - \mathbf{y}$ , increase  $k$  and go to 2.

To improve the convergence behavior, a damping factor  $0 < d < 1$  is commonly used. The update is then  $x_k := x_k - d\mathbf{y}$ .

As a final example, consider the system

$$\begin{aligned} f_1(x, y) &= 3x + 4y = 4, \\ f_2(x, y) &= 7x + 3y = 9. \end{aligned}$$

Define the implicit form

$$\begin{aligned} F_1(x, y) &= 3x + 4y - 4 = 0, \\ F_2(x, y) &= 7x + 3y - 9 = 0, \end{aligned}$$

whose Jacobian is

$$J_F = \begin{pmatrix} 3 & 4 \\ 7 & 3 \end{pmatrix}.$$

Starting with  $\begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  results in

$$\begin{pmatrix} x_1 \\ y_1 \end{pmatrix} = -J_F^{-1} \begin{pmatrix} F_1(0,0) \\ F_2(0,0) \end{pmatrix} = - \begin{pmatrix} 3 & 4 \\ 7 & 3 \end{pmatrix}^{-1} \begin{pmatrix} -4 \\ -9 \end{pmatrix} = \begin{pmatrix} 3 & 4 \\ 7 & 3 \end{pmatrix}^{-1} \begin{pmatrix} 4 \\ 9 \end{pmatrix}.$$

For this system of linear equations the Jacobian is nothing but the initial system in matrix-vector form! Please note that an additional damping factor requires several solutions for the same system matrix, although the solution could be obtained in one step only.



## Chapter 4

# Two-Dimensional Simulation and Grids

In the previous chapter the numerical approximation of derivatives in one dimension was shown. This chapter goes one step further and applies these approximation techniques in two dimensions. First, the discretization and solution of the two-dimensional Laplace equation using finite differences will be discussed, after that, the flux-conserving box-integration method will be presented. The latter is also applicable to unstructured meshes, which will be the topic of Chapter 5.

### 4.1 Two-Dimensional Laplace Equation

The Laplace operator  $\Delta = \nabla^2$  is a second order differential operator in the  $n$ -dimensional Euclidean space. It is defined as the divergence of the gradient  $\nabla$ . Laplace's equation is the homogeneous form of the Poisson equation, i.e. without charge density  $\rho$  on the right hand side. In two dimensions it reads

$$\nabla^2\psi(x,y) = \frac{\partial^2\psi(x,y)}{\partial x^2} + \frac{\partial^2\psi(x,y)}{\partial y^2} = 0 \quad (4.1)$$

in **Cartesian**<sup>1</sup> coordinates. The solutions of the Laplace equation are called *harmonic functions*. To determine the solution of (4.1) completely it is necessary to specify appropriate boundary conditions.

We begin our discussion of differencing schemes for elliptic equations by looking at the Poisson equation (4.1) in the unit square with equidistant grid spacing in the  $x$ - and  $y$ - directions. The standard central difference approximations for the second order derivatives with

$$\left(\frac{d^2\psi}{dx^2}\right)_{x=i\Delta x} \approx \frac{\psi[i+1] - 2\psi[i] + \psi[i-1]}{(\Delta x)^2},$$

are now applied in two dimensions thus giving

$$\nabla^2\psi(x,y) \approx \frac{\psi[i+1,j] - 2\psi[i,j] + \psi[i-1,j]}{(\Delta x)^2} + \frac{\psi[i,j+1] - 2\psi[i,j] + \psi[i,j-1]}{(\Delta y)^2}. \quad (4.2)$$

---

<sup>1</sup> **Cartesian** [ka:'ti.zi.ən]: kartesisch

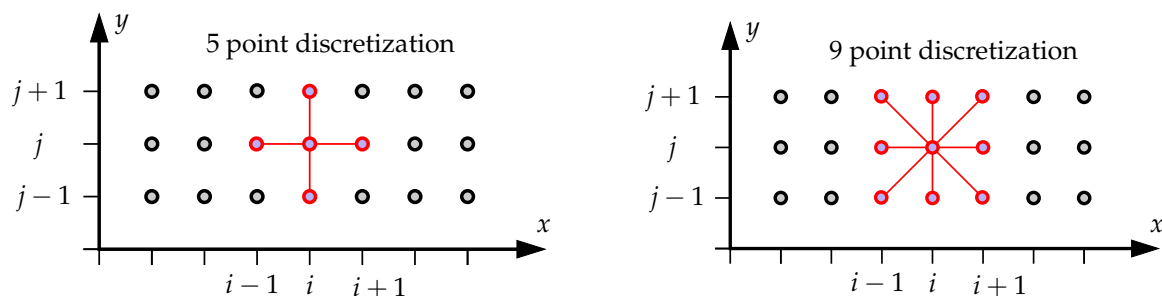


Figure 4.1: Five-point Laplacian (left) and nine-point Laplacian (right).

On an equidistant square grid ( $\Delta y = \Delta x$ ) this becomes

$$\nabla^2 \psi(x, y) \approx \frac{\psi[i+1, j] + \psi[i, j+1] - 4\psi[i, j] + \psi[i-1, j] + \psi[i, j-1]}{(\Delta x)^2}.$$

The difference operator on the right hand side is called *five-point* (discrete) **Laplacian**<sup>1</sup> (aka *five-point stencil*<sup>2</sup>, because the point operation pattern is applied to all points of the grid) with second order accuracy. Another possible approximation of the Laplace equation is the fourth order accurate *nine-point* Laplacian. For further information the reader is referred to the book of Strikwerda [?].

```

1 const NX = 21,    NY = 21;    // number of mesh points
2 const LX = NX - 1, LY = NY - 1; // index of last mesh point
3 const DX = 1.0;    // mesh spacing
4 var psi[NX, NY];
5
6 // Laplace's equation (only inner points)
7 equ psi[i=1..LX-1, j=1..LY-1] ->
8     {psi[i-1, j] + psi[i, j-1] - 4*psi[i, j] +
9     psi[i+1, j] + psi[i, j+1]} / sq(DX) = 0.0;
10 begin main
11   assign psi[i=all, j=all] = 0.0; // Dirichlet boundary condition
12   assign psi[i=0, j=all] = 10.0; // Dirichlet boundary condition
13   solve;
14   write;
15 end

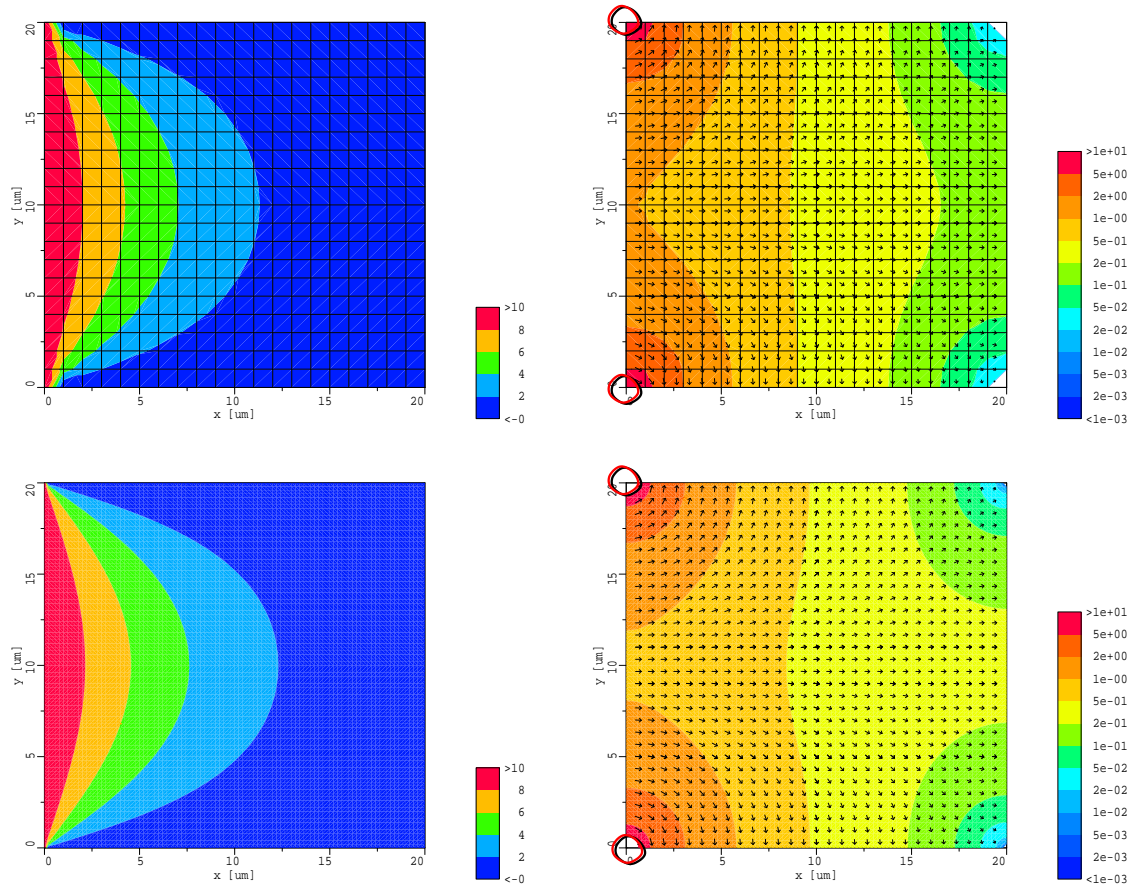
```

source\_code/laplace.sg

In `laplace.sg` the potential of `psi[i=0, j=all]` is set equal to `10.0`, which constitutes a Dirichlet boundary condition and overwrites the previous assignment `psi[i=all, j=all] = 0.0`. The command `solve` lets SGFRAMEWORK solve the system of equations given by the `equ` command. The resulting potential can be seen in Fig. 4.2. Due to the discontinuous boundary conditions given, singularities of the electrical field (which is the negative gradient of the potential) occur. Such singularities cannot simply be swept under the carpet: They may influence the global quality of approximation through a reduction of the convergence rate.

Motivated from the singularities that show up with Dirichlet boundary conditions, one may think that large gradients are in general a consequence of the source/sink term  $g$  and of discontinuities in the boundary conditions only. However, this is not true. Even if the source/sink

<sup>1</sup> **Laplacian** [laplasɪæn]: Laplace-Operator    <sup>2</sup> **stencil** [stent.səl]: Vervielfältigungsmatrix



**Figure 4.2:** Potential (left) and electric field in log-scale (right) with  $21 \times 21$  points (top) and  $101 \times 101$  points (bottom) with Dirichlet boundary conditions. The singularities occurring at the edges are marked with circles.

term  $g$  were perfectly smooth and no discontinuities of the boundary conditions occurred, the domain  $\Omega$  may consist of **reentrant corners**<sup>1</sup>, that will always provoke singularities in the solution.

To come back to our simple rectangular domain, replacing the Dirichlet boundary condition with a homogeneous Neumann boundary condition (i.e. zero flux) at the boundary section  $j=LY$  results in  $\psi[i, LY] = \psi[i, LY - 1]$  for all  $i$ . As a consequence, one row in the array is wasted, since two rows carry the same information.

```
1 equ psi[i=1..LX-1,j=LY] -> {psi[i,j] - psi[i,j-1]} / sq(DX) = 0.0;
```

source\_code/neumann\_boundary1.sg

One can get rid of this waste of memory by using a little trick: We imagine the existence of another row at  $j=LY+1$ , so that the boundary condition at  $j=LY$  is  $(\psi[i, LY + 1] - \psi[i, LY]) / DX = 0.0$ ;, which simply states that  $\psi[i, LY + 1] = \psi[i, LY]$ . Now we consider the five-point stencil

```
1 equ psi[i=1..LX-1,j=1..LY-1] -> {psi[i-1,j] + psi[i,j-1] - 4*psi[i,j] +
2 psi[i+1,j] + psi[i,j+1]} / sq(DX) = 0.0;
```

source\_code/neumann\_boundary2.sg

<sup>1</sup> **reentrant corners** [ri:en.trænt]: einspringende Ecken



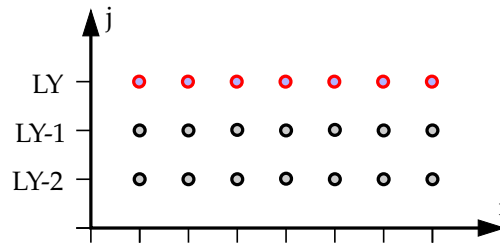


Figure 4.3: One row seems to be wasted when using Neumann boundary condition.

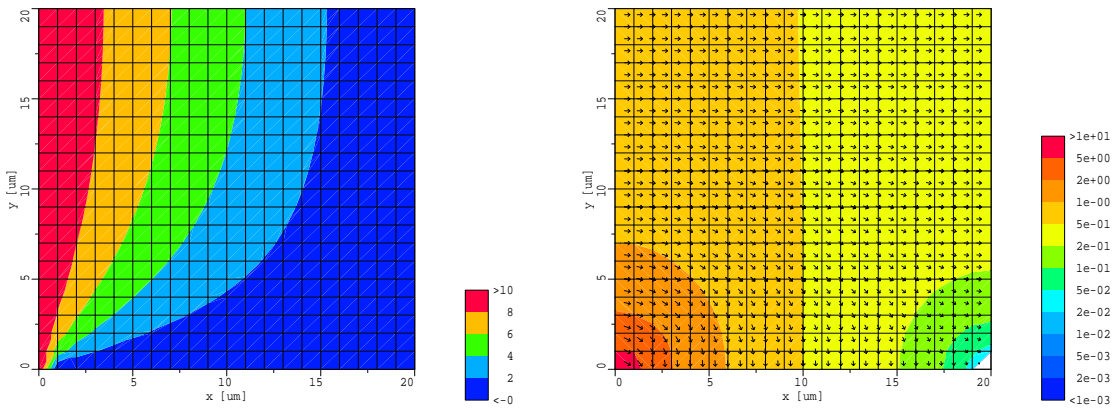


Figure 4.4: Potential (left) and electric field in log-scale (right) with  $21 \times 21$  points including Dirichlet and Neumann boundary conditions.

at  $j = LY$  and replace all occurrences of  $\psi[i, LY + 1]$  with  $\psi[i, LY]$ , as required by the Neumann boundary condition. After collection of terms, the new equation for the last row is

$$\text{equ } \psi[i=1..LX-1, j=LY] \rightarrow \{ \psi[i-1, j] - 3 * \psi[i, j] + \psi[i+1, j] + \psi[i, j-1] \} / \text{sq}(DX) = 0.0;$$

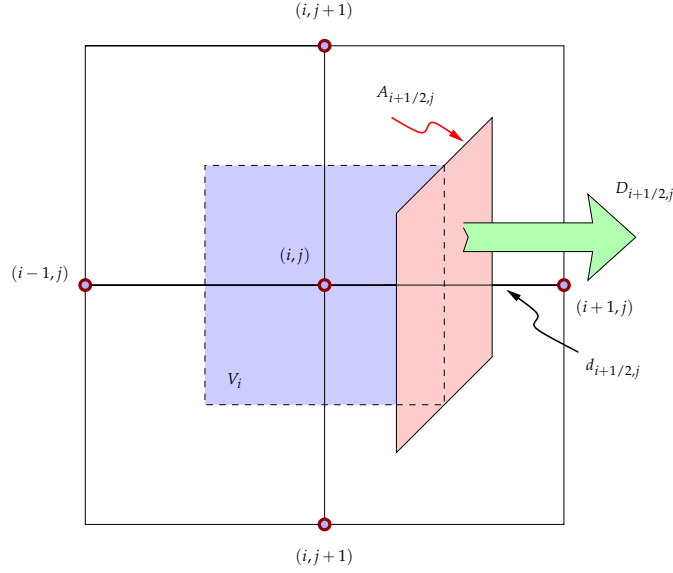
source\_code/neumann\_boundary3.sg

and has to be implemented into the `laplace.sg` SGFRAMEWORK file.

The resulting potential can be seen in Fig. 4.4. At the Neumann boundary, the solution is very smooth, the only discontinuity that shows up is due to the discontinuous Dirichlet boundary conditions. Note that the field in Fig. 4.4 resembles the field in the lower half of the problem space of Fig. 4.2. This is not a coincidence, it is a general property of homogeneous Neumann boundary conditions, which is why they are also called mirror boundary conditions (cf. method of images in electrostatics).

## 4.2 Box Integration Method

The **box integration method** (aka. *finite volume method*, suitable also for unstructured meshes) is a subset of finite difference methods, and is very important in 2D semiconductor device simulations. Box integration provides a convenient method of discretizing a large class of PDEs on both regular and irregular meshes by constructing a box around each node in such a way that



**Figure 4.5:** Approximated closed-loop path integral is made up of its four 'constant  $D$ ' edge contributions.

the boxes cover the whole simulation domain. The differential equations are then individually integrated over each of the subdomains.

#### 4.2.1 Example: The Poisson Equation

In order to demonstrate the box method, we integrate the Poisson equation  $\nabla \cdot \mathbf{D} = \rho$  over a small region (box)  $\mathcal{V}$  with volume  $V_i$  and obtain

$$\int_{\mathcal{V}} \nabla \cdot \mathbf{D} \, dx = \int_{\mathcal{V}} \rho \, dx \approx \rho_i V_i,$$

where  $\rho_i$  is the charge density evaluated at some representative point  $x_i$  inside  $\mathcal{V}$ . We will discretize the left-hand side using *Gauss' integral theorem*, which states that

$$\int_{\mathcal{V}} \nabla \cdot \mathbf{D} \, dx = \int_{\partial\mathcal{V}} \mathbf{D} \cdot \mathbf{n} \, dx,$$

where  $\mathbf{n}$  is the outward pointing local normal vector of the enclosing surface. In two dimensions,  $\partial\mathcal{V}$  is a closed-loop path.

When discretizing a closed-loop path integral, the contributions along the path are approximated and accumulated. For the example in Fig. 4.5 this means that for instance  $D_{i+1/2, j}$  is assumed constant along the edge  $(i + 1/2, j - 1/2)$  and  $(i + 1/2, j + 1/2)$ . The other three sections of the path are approximated in a similar fashion. As we make the box volume smaller, the approximation error becomes smaller. Hence, the approximation of the closed-loop path integral gives

$$\int_{\partial\mathcal{V}} \mathbf{D} \cdot \mathbf{n} \, dx \approx D_{i+1/2, j} A_{i+1/2, j} + D_{i-1/2, j} A_{i-1/2, j} + D_{i, j+1/2} A_{i, j+1/2} + D_{i, j-1/2} A_{i, j-1/2}$$

with the approximation of the fluxes

$$\begin{aligned}
 D_{(i,j) \rightarrow (l,m)} &= \mathbf{D} \Big|_{x=(i+1)/2, y=(j+m)/2} \cdot \mathbf{e}_{(i,j) \rightarrow (l,m)} \\
 &= (\varepsilon \mathbf{E}) \Big|_{x=(i+1)/2, y=(j+m)/2} \cdot \mathbf{e}_{(i,j) \rightarrow (l,m)} \\
 &= (-\varepsilon \nabla \psi) \Big|_{x=(i+1)/2, y=(j+m)/2} \cdot \mathbf{e}_{(i,j) \rightarrow (l,m)} \\
 &\approx -\frac{\varepsilon_{i,j} + \varepsilon_{l,m}}{2} \frac{\psi_{l,m} - \psi_{i,j}}{d_{(i,j) \rightarrow (l,m)}}.
 \end{aligned}$$

Assuming an equidistant mesh ( $\Delta y = \text{const}$ ,  $\Delta x = \text{const}$ ),

$$A_{i+1/2,j} = A_{i-1/2,k} = \Delta y w, A_{i,j+1/2} = A_{i,j-1/2} = \Delta x w, V_i = \Delta x \Delta y w,$$

with  $w$  as depth of the simulation domain and constant permittivity we get

$$\begin{aligned}
 \int_{\partial V} \mathbf{D} \cdot \mathbf{n} \, dx &\approx -\varepsilon \frac{\psi_{i+1,j} - \psi_{i,j}}{\Delta x} w \Delta y - \varepsilon \frac{\psi_{i-1,j} - \psi_{i,j}}{\Delta x} w \Delta y \\
 &\quad - \varepsilon \frac{\psi_{i,j+1} - \psi_{i,j}}{\Delta y} w \Delta x - \varepsilon \frac{\psi_{i,j-1} - \psi_{i,j}}{\Delta y} w \Delta x \\
 &= -\varepsilon \frac{\psi_{i+1,j} - 2\psi_{i,j} + \psi_{i-1,j}}{\Delta x} w \Delta y \\
 &\quad - \varepsilon \frac{\psi_{i,j+1} - 2\psi_{i,j} + \psi_{i,j-1}}{\Delta y} w \Delta x \\
 &= \rho_i \Delta x \Delta y w.
 \end{aligned}$$

Dividing by  $V_i$  we finally obtain

$$-\varepsilon \frac{\psi_{i+1,j} - 2\psi_{i,j} + \psi_{i-1,j}}{(\Delta x)^2} - \varepsilon \frac{\psi_{i,j+1} - 2\psi_{i,j} + \psi_{i,j-1}}{(\Delta y)^2} = \rho_i$$

which is exactly the same result as obtained by finite differences (4.2). Assuming for example an inhomogeneous permittivity  $\varepsilon = \varepsilon(y)$  one has to modify the previous equations to

$$\begin{aligned}
 \int_{\partial V} \mathbf{D} \cdot \mathbf{n} \, dx &\approx -\varepsilon_j \frac{\psi_{i+1,j} - 2\psi_{i,j} + \psi_{i-1,j}}{\Delta x} w \Delta y \\
 &\quad - \frac{\varepsilon_j + \varepsilon_{j+1}}{2} \frac{\psi_{i,j+1} - \psi_{i,j}}{\Delta y} w \Delta x \\
 &\quad - \frac{\varepsilon_j + \varepsilon_{j-1}}{2} \frac{\psi_{i,j-1} - \psi_{i,j}}{\Delta y} w \Delta x.
 \end{aligned}$$

#### 4.2.2 Example: Extraction of Capacitances

When working with Dirichlet boundary conditions for the box integration method, two implementation models are appropriate. Either a substitute equation, like  $\psi_{LY} = V_c$  is implemented or the Dirichlet boundary conditions are explicitly inserted into the .sg file in place of  $\psi_{LY}$ . While the first possibility wastes variables, the latter requires more careful index-checking.

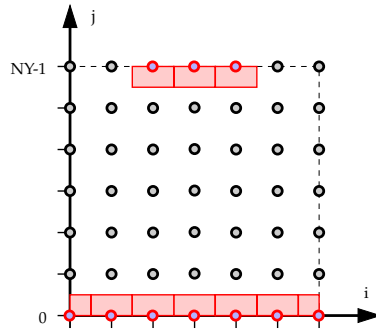
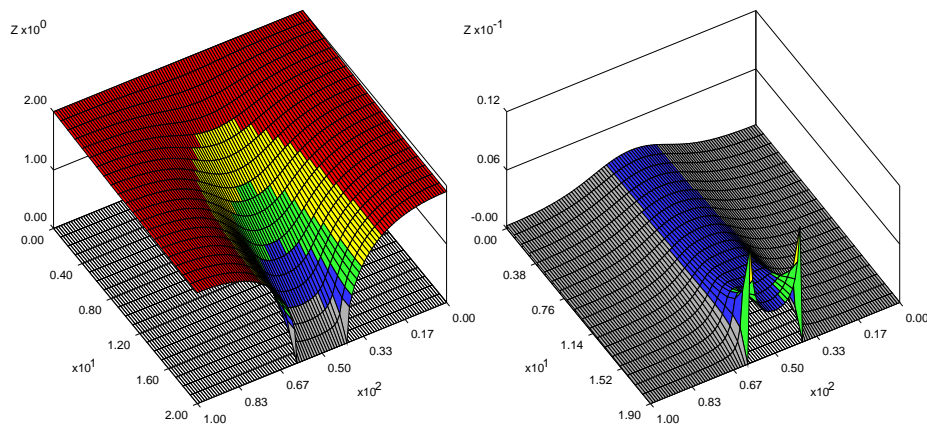


Figure 4.6: Dirichlet Boundary Conditions


 Figure 4.7: Potential  $\psi$  (left) and displacement  $D_y$  (right)

Neumann boundary conditions for zero out-flux are dealt with very simply by omitting corresponding contributions within the surface integral. For a box centered at point  $(r, s)$  with a homogeneous Neumann boundary condition along the edge at  $(r, s + 1/2)$ , we have

$$\int_{\partial V} \mathbf{D} \cdot \mathbf{n} \, dx \approx D_{r+1/2,s} A_{r+1/2,s} + D_{r-1/2,s} A_{r-1/2,s} + \underbrace{D_{r,s+1/2} A_{r,s+1/2}}_{=0} + D_{r,s-1/2} A_{r,s-1/2}.$$

In contrast to that, Neumann boundary conditions with a non-zero boundary flux  $D_c$  are added explicitly

$$\int_{\partial V} \mathbf{D} \cdot \mathbf{n} \, dx \approx D_{r+1/2,s} A_{r+1/2,s} + D_{r-1/2,s} A_{r-1/2,s} + D_c A_{r,s+1/2} + D_{r,s-1/2} A_{r,s-1/2}.$$

extraction\_capacitances.sg shows a program to simulate a capacitor with the results shown in Fig. 4.7 and Fig. 4.8.

```

1  const NX = 101, NY = 21, DX = 1.0e-9, DY = DX;
2  const Vc = 2,    A = DX;
3  const X1 = 40,  X2 = 60;
4
5  const RHO    = 1e-14;
6  const EPSr  = 3.9;           // relative permittivity region 1
7  const EPSo  = 8.854e-12;    // permittivity of vacuum (approximately)
    
```

```

8
9 var psi[NX,NY], Ex[NX-1,NY], Ey[NX,NY-1], Dx[NX-1,NY], Dy[NX,NY-1];
10 var E[NX-1,NY-1], D[NX-1,NY-1], C[2], Qs[NX,NY];
11
12 unknown psi[0..NX-1,0..NY-1];
13 known psi[X1..X2,NY-1];
14 known psi[all,0];
15
16 // Poisson Equation for the inner grid points
17 equ psi[i=1..NX-2,j=1..NY-2] ->
18   {
19     psi[i,j-1] +
20     psi[i-1,j] - 4.0*psi[i,j] + psi[i+1,j]
21     + psi[i,j+1] } / sq(DX) = RHO / ( EPSo * EPSr);
22
23 // Neumann boundary condition
24 equ psi[i=0, j=0..NY-2] -> psi[i,j] - psi[i+1,j] = 0;
25 equ psi[i=NX-1, j=0..NY-2] -> psi[i,j] - psi[i-1,j] = 0;
26 equ psi[i=0..X1-1, j=NY-1] -> psi[i,j] - psi[i,j-1] = 0;
27 equ psi[i=X2+1..NX-1, j=NY-1] -> psi[i,j] - psi[i,j-1] = 0;
28
29 begin main
30 // Contact voltages , Dirichlet boundary condition
31 assign psi[i=X1..X2, j=NY-1] = 0;
32 assign psi[i=all, j=0] = Vc;
33
34 solve;
35
36 // Find E and D from the potential psi using right sided differences
37 assign Ex[i=all, j=all] = - (psi[i+1,j] - psi[i,j])/DX ; // Ex = - dpsi/DX
38 assign Ey[i=all, j=all] = - (psi[i,j+1] - psi[i,j])/DY ; // Ey = - dpsi/dy
39 assign E [i=all, j=all] = sqrt(sq(Ex[i,j]) + sq(Ey[i,j]));
40
41 assign Dx[i=all, j=all] = EPSr * EPSo * Ex[i,j];
42 assign Dy[i=all, j=all] = EPSr * EPSo * Ey[i,j];
43 assign D [i=all, j=all] = sqrt(sq(Dx[i,j]) + sq(Dy[i,j]));
44
45 // sum up charges at Vc-electrode
46 assign Qs[i=0, j=0] = A/2 * Dy[i,j];
47 assign Qs[i=1..NX-2,j=0] = Qs[i-1,j] + A * Dy[i,j];
48 assign Qs[i=NX-1, j=0] = Qs[i-1,j] + A/2 * Dy[i,j];
49
50 // sum up charges at ground electrode (do not forget Dx-components!)
51 assign Qs[i=X1, j=NY-1] = A/2 * Dx[i-1,j] + A * Dy[i,j-1];
52 assign Qs[i=X1+1..X2-1,j=NY-1] = Qs[i-1,j] + A * Dy[i,j-1];
53 assign Qs[i=X2, j=NY-1] = Qs[i-1,j] + A/2 * Dx[i,j] + A * Dy[i,j-1];
54
55 // compute capacitances from the auxiliary charge
56 assign C[i=0] = Qs[NX-1,0] / Vc;
57 assign C[i=1] = Qs[X2,NY-1] / (- Vc);
58
59 write;
60
61 end

```

source\_code/extraction\_capacitances.sg

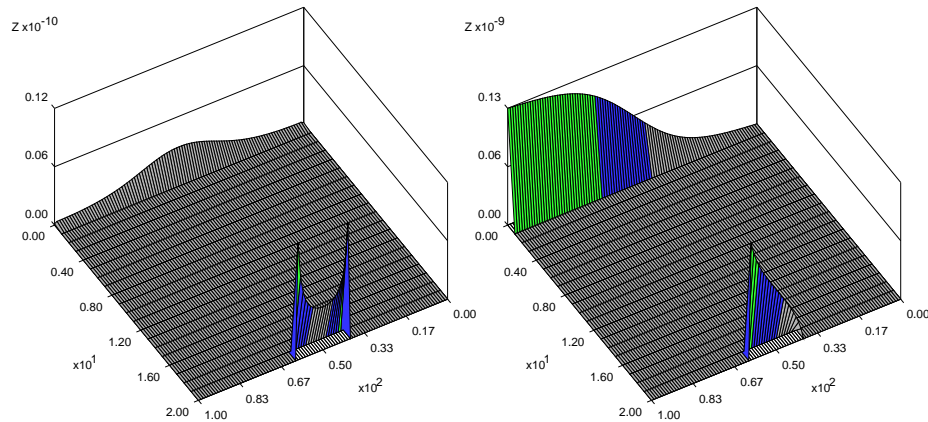


Figure 4.8: Real charge per node (left) and auxiliary charge  $Q_s$  (right)

The capacitance  $C = Q/V$  is calculated via the charge on the contact surface, which is

$$Q = \int_{V_{\text{contact}}} \rho \, dV = \int_{\partial V_{\text{contact}}} \mathbf{D} \cdot \mathbf{n} \, dx = \sum_{\text{all contact points}} D_y[i] A[i].$$

Because there is no direct support for integration of a function in SGFRAMEWORK, a workaround in the form of a sum over all contact points has to be assigned. Starting at  $i=0, j=0$  the first element is calculated (be aware that for the first and the last cell only a half of each edge belong to the contact). The charge is then added up from the first to the last cell. The implementation trick in SGFRAMEWORK can be understood when looking at Fig. 4.6 and is requires some extra attention when dealing with the last row  $j=NY-1$ . The code in `extraction_capacitances2.sg` displays the trick and Fig. 4.8 shows the result on the right hand side. Starting at  $x = 0$ , the (auxiliary-) charge is growing along the positive  $x$ -axis until its peak value at the end of the contact where it represents the approximated value of the integral. Again, the accuracy only depends on the grid spacing.

```

1  assign Qs[i=0,      j=0]      =      A/2 * Dy[i , j ];
2  assign Qs[i=1..NX-2,j=0]     = Qs[i-1,j] + A * Dy[i , j ];
3  assign Qs[i=NX-1,  j=0]     = Qs[i-1,j] + A/2 * Dy[i , j ];
4  assign C[i=0]              = Qs[NX-1,0] / Vc;
    
```

source\_code/extraction\_capacitances2.sg



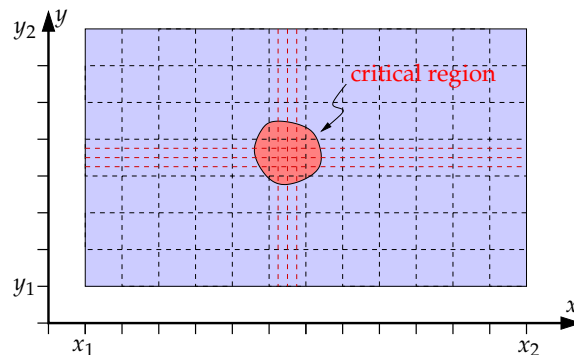
## Chapter 5

# Tessellation of Unstructured Meshes

If we need a high accuracy of the quantity of interest in a certain region of the device only, structured grids like Manhattan geometries<sup>1)</sup> or ortho-product grids require a refinement outside the simulation domain as well. For example, for a two-dimensional simulation of a MOSFET with 100 grid points in both spatial directions, we have to deal with 10 000 unknowns. Typically, we require a high resolution of the channel, which is only a very small part of the simulation domain. As can be seen in Fig. 5.1, a refinement of the channel (critical region) would also result in a refinement along strips parallel to the  $x$ - and  $y$ -axis. This way, many additional unknowns are introduced outside the channel (critical region) which are not needed.

By means of the box integration method, unstructured meshes can be used for solving partial differential equations. This allows meshes with very fine grid spacing in the area of interest (like the aforementioned channel of a MOSFET), while maintaining a larger grid spacing in regions where the quantities are known in advance to have small variations (like in the deep bulk region of a MOSFET).

This chapter deals with the generation and handling of such unstructured meshes and describes both the resulting benefits and the **subtleties**<sup>2</sup>.



**Figure 5.1:** Although only nine additional points are required in the critical region, the total number of grid points  $N_{\text{Points}}$  increases from 104 to 176.

<sup>1</sup> The upright projection of a transistor may look like the street map of Manhattan    <sup>2</sup> **subtlety** [sʌt.l.ti]: Schwierigkeit, Raffinesse



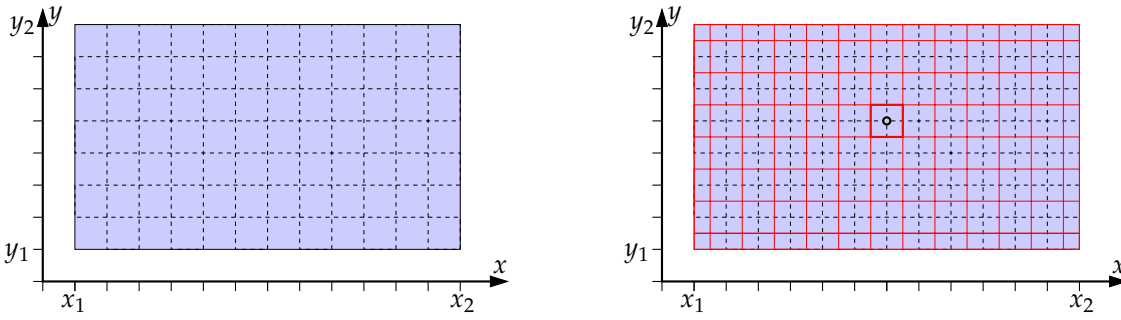


Figure 5.2: Ortho-product grid (left) and its corresponding Voronoi tessellation (right).

## 5.1 Voronoi Tessellation

When the points are located arbitrarily in space rather than on a regular grid, we need a proper generalization of the ortho-product grid, called Voronoi **tessellation**<sup>1</sup>: Consider a finite set of points  $\mathcal{D} = \{r_1, r_2, \dots, r_n\}$  in a subdomain  $\Omega$  of  $\mathbb{R}^n$ . A Voronoi region  $\Omega_i$  is the set of all points of  $\Omega$  that are closer to  $r_i$  than to any other point of  $\mathcal{D}$ :

$$\Omega_i = \{r \in \Omega \mid \|r - r_i\| < \|r - r_j\| \forall j \neq i\} \quad (5.1)$$

The resulting Voronoi regions  $\Omega_i$  form a tessellation of  $\Omega = \bigcup_i \Omega_i$ , given in Fig. 5.2.

The construction of a Voronoi tessellation works as follows: The edges of the boxes are obtained by drawing the **perpendicular bisecting**<sup>2</sup> segments of edges that connect the mesh nodes, as shown in Fig. 5.3. The bisectors are only allowed to extend until they cut another bisector.

## 5.2 Triangular Delaunay Meshes

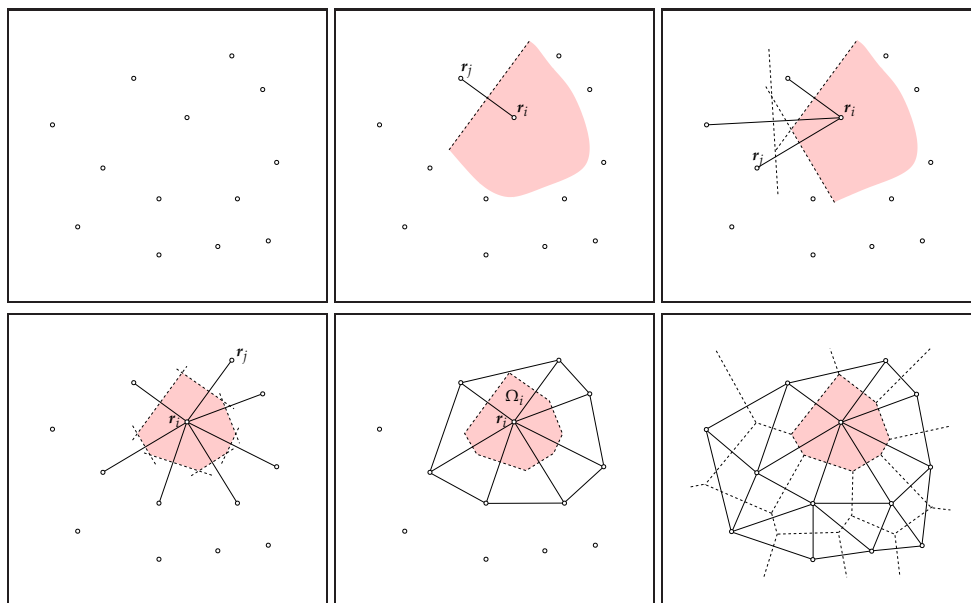
It can be shown that every Voronoi tessellation has a dual triangular mesh called **Delaunay**<sup>3</sup> mesh (can be seen in the lower right of Fig. 5.3). Since it is usually easier to construct a triangular mesh than to construct the Voronoi tessellation directly, we will use this duality to construct a Voronoi tessellation from a triangulation. Unfortunately, not every triangulation is the dual of a Voronoi tessellation, but it is possible to modify triangulations (meshes) such that they become duals of Voronoi tessellations, i.e. Delaunay meshes.

Let us have a look at the construction of a Delaunay mesh out of a point cloud. In order to obtain a Delaunay mesh, we need a criterion that allows us to determine whether a triangulation is a Delaunay mesh (Fig. 5.4):

- *Empty Sphere Criterion*: The open discs (balls) circumscribing the triangles (or tetrahedra) must not contain any other mesh point (Fig. 5.5).

It turns out that the local property of fulfilling the Empty Sphere Criterion directly extends to a global property:

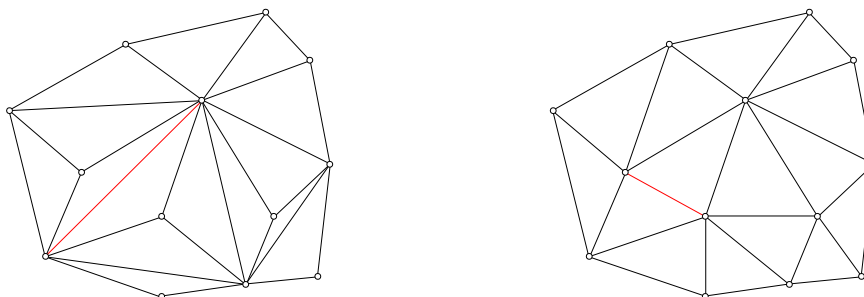
<sup>1</sup> **tessellation** [tes.ə'leɪ.ʃən]: Mosaik    <sup>2</sup> **to perpendicular bisect** [peɪ.pə'n'dɪk.jʊ.lər baɪ'sekt]: in zwei gleich große Teile teilen    <sup>3</sup> **Delaunay** [deləʊneɪ]: Delaunay



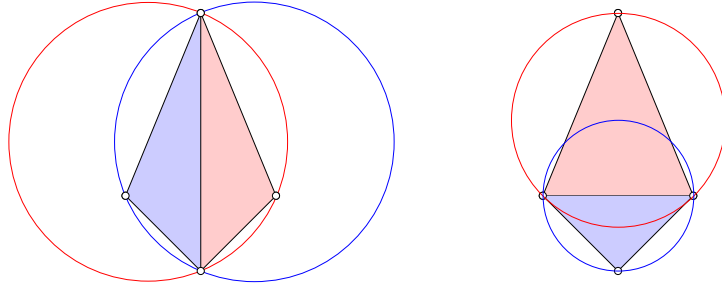
**Figure 5.3:** A set of 13 grid points and the associated Voronoi regions (bounded by the dashed lines).

**Lemma.** *If for every pair of adjacent triangles (or tetrahedrons) in a given mesh the empty sphere criterion holds, then this criterion holds globally and the mesh is a Delaunay triangulation.*

With the Empty Sphere Criterion at hand, Delaunay meshes can be constructed in many ways. One of them is the *incremental method*, where one starts with a rectangle (or any other polyhedron) made up from a few points within the simulation domain and then adds point by point (Fig. 5.6) to obtain a full mesh. Alternatively, one might take an arbitrary triangulation and swap the diagonals that violate the empty sphere criterion; this is called the *swapping method*. Moreover, as indicated in the introduction to this chapter, it is desirable that the mesh reflects the regions of interest in the device: The mesh density for example should be a function of the gradient (local change rate) of the quantities. This poses a chicken-and-egg problem: The mesh should be instantiated prior to the determination of the quantities (the simulation process it-



**Figure 5.4:** Example of two triangular grids constructed from the same point-cloud, where the right one is the dual of the Voronoi tessellation.



**Figure 5.5:** Example of two triangular meshes constructed from the same deltoid. By the swapping the diagonals in the left mesh, the mesh on the right is obtained, which is the dual of the Voronoi tessellation.

self), but the mesh generation process needs the quantities to construct a suitable mesh. This is actually the idea of adaptive schemes, where the solution is first computed on a coarse mesh, then this mesh is refined according to the behavior of the coarse solution. This procedure is then repeated on the refined mesh again, until a sufficiently good approximation to the true solution is found.

### 5.3 Skeleton Mesh

From the previous section it should have become clear that additional topological information has to be stored for unstructured grids. For example, the location of points cannot be **recovered**<sup>1</sup> by their array index anymore, as well as some kind of connectivity information must be available. In this section, the data structures used to **accomplish**<sup>2</sup> this are discussed.

Referring to the example of the Poisson equation in section 4.2.1, we write

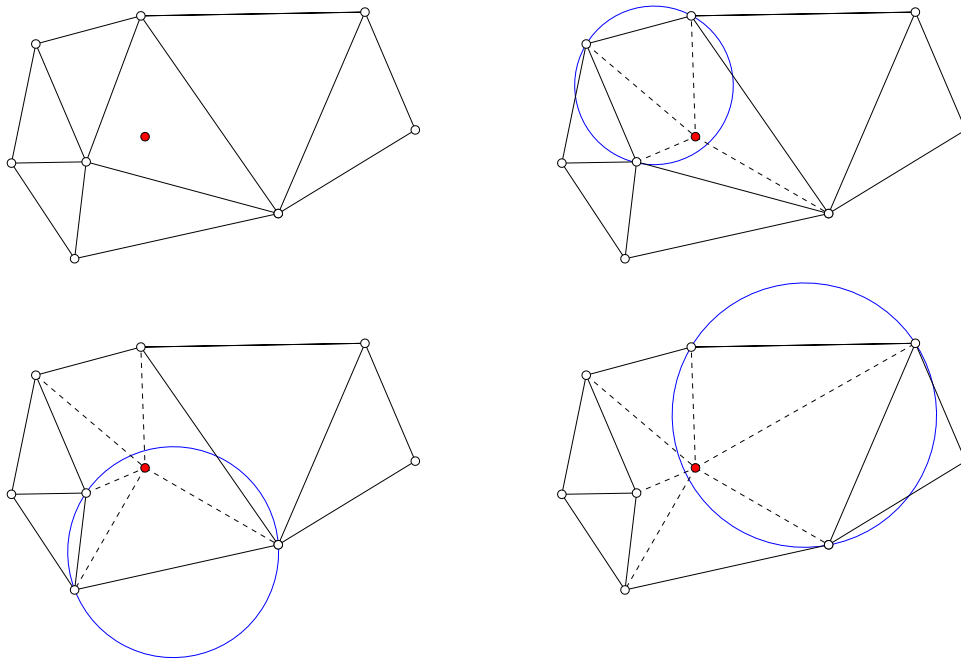
$$\sum_{j \in \mathcal{N}_i} D_{i,j} A_{i,j} = - \sum_{j \in \mathcal{N}_i} \varepsilon_{i,j} \frac{\psi_j - \psi_i}{d_{i,j}} A_{i,j} = \rho_i V_i,$$

where  $\mathcal{N}_i$  denotes the set of all nodes that are neighbors of node  $i$ . The integral is discretized by the box method using the Voronoi tessellation. In Fig. 5.8 the solid lines are the mesh edges while the dashed lines are the edges of the integration box.  $V_i$  is the volume of the box surrounding the mesh point  $i$ ,  $d_{i,j}$  is the length of the edge that connects  $i$  and its neighbor-point  $j$ , and  $A_{i,j}$  is the length of the corresponding integration area, i.e. the area of that surface element of  $\Omega_i$  that interfaces the box  $\Omega_j$  at point  $j$ .

To foster the understanding of unstructured neighborhood information, a 3 times 3 rectangular example geometry is presented in Fig. 5.8. Taking a closer look at Tab. 5.1, one can verify that the only box having a volume  $V_i = 1$  is box number 4 because of its location in the interior, while the corner boxes (0, 2, 6, 8) and the edge boxes (1, 3, 5, 7) have a volume of 1/4 and 1/2 respectively. The column names  $x_i$  and  $y_i$  represent the indices of the box.  $A_{i,j}$  and  $d_{i,j}$  have already been explained.

The information given in Tab. 5.1 is sufficient for any geometry in any dimension. It is implemented in SGFRAMEWORK using a so-called *skeleton file*. The domain specified in `sample_geometry.sk` is shown in Fig. 5.9.

<sup>1</sup> **to recover** [ri'kʌv.ə]: zurückgewinnen, wiedererlangen    <sup>2</sup> **to accomplish sth.** [ə'kʌ:m.plɪʃ]: etwas erreichen, etwas vollbringen



**Figure 5.6:** Incremental method (starting from a Delaunay mesh). First, the newly inserted point is connected with the nodes of the triangle it is located in. Then, the empty sphere criterion is checked for all adjacent triangles. If the criterion is violated for a neighbor, the interfacing edge is removed and the newly inserted point is connected with the opposite corner node of the neighbor.

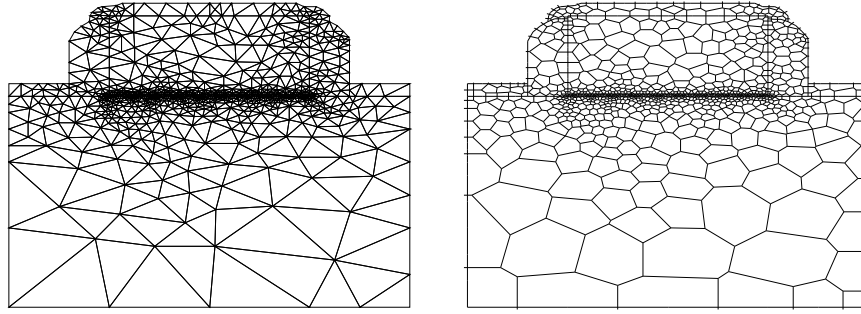
```

1 // Mesh file for simple square mesh
2 // sample_geometry.sk (will be converted to sample_geometry.msh)
3
4 const dx = 1.0, W = 20.0;
5
6 point pA = (0.0, W), pB = ( W, W);
7 point pC = ( W, 0.0), pD = (0.0, 0.0);
8
9 edge eAB = METAL1 [pA,pB] (dx, 0.0);
10 edge eBC = METAL2 [pB,pC] (dx, 0.0);
11 edge eDC = METAL2 [pD,pC] (dx, 0.0);
12 edge eAD = METAL2 [pA,pD] (dx, 0.0);
13
14 // remove keyword 'rectangles' to get a triangular mesh
15 region rABCD = AIR {eAD,eDC,eBC,eAB} rectangles;
16
17 coordinates x, y;
    
```

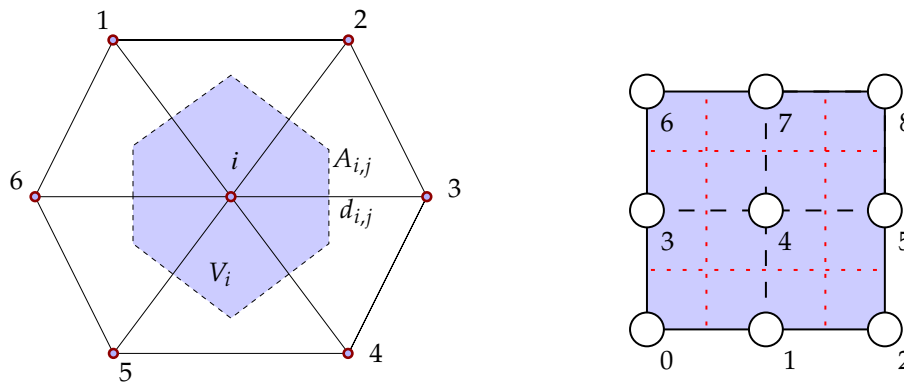
source\_code/sample\_geometry.sk

In the first section geometrical quantities such as the device width and depth are defined as constants. Next the skeleton points are defined by two coordinates enclosed in *parentheses* each. Each point is given a name that is used for further reference. Then edges are defined using the previously set points. They are described by their two endpoints enclosed in brackets. The order of the edges (running **counterclockwise**<sup>1</sup>) is important later when defining regions. If the edges shall constitute a rectangular region, opposite parallel edges must point in the same

<sup>1</sup> **counterclockwise** [kaʊn.təˈklok.waɪz]: gegen den Uhrzeigersinn



**Figure 5.7:** Unstructured mesh of a MOS transistor; triangular mesh (left) and the corresponding Voronoi regions (right).



**Figure 5.8:** Left: The Voronoi box  $V_i$  of mesh point  $i$  is surrounded by six adjacent mesh points. For the discretization the knowledge of  $A_{i,j}$  and  $d_{i,j}$  is required. Right: in an ortho-grid geometry this information is implicitly available.

direction. Finally, the regions are defined in the last section of the mesh skeleton by a list of three or more edges. The list is enclosed in *braces*. The skeleton file exports all the user-defined constants plus three additional ones to the `.sg` file:

```

NODES      ...  number of nodes (points)
EDGES      ...  number of edges
ELEMENTS   ...  number of elements (regions; triangular and rectangular)
    
```

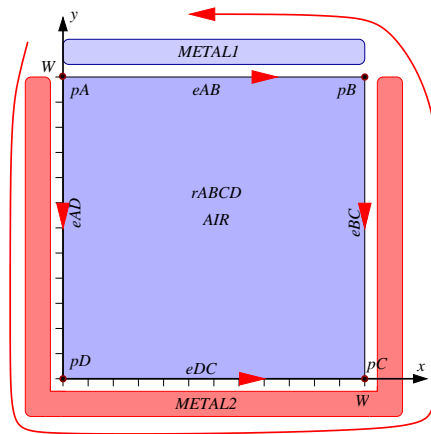
In `sample_geometry.sk` the coordinates are identified as  $(x, y)$ . The arrays  $x[j]$  and  $y[j]$  are automatically passed to `mesh_example1.sg` from the generated mesh file `sample_geometry.msh`. This is important, because otherwise there would be no way of finding the positions of the nodes. The labels representing a list of nodes are imported as well. The results are given in Fig. 5.10.

```

1 mesh "sample_geometry.msh";
2
3 var      psi[NODES]; // one-dimensional
4 unknown psi[all];
5 known   psi[METAL1], psi[METAL2];
6
7 equ psi[i=AIR] -> nsum(i,j,all,
8                      ((psi[i]-psi[node(i,j)]) / elen(i,j)) * ilen(i,j))
9                      ) = 0.0;
    
```

Boxnumber	$x_i$	$y_i$	$V_i$	Connectionnumber	$i$	$j$	$A_{i,j}$	$d_{i,j}$
0	0	0	1/4	0	0	1	1/2	1
1	1	0	1/2	1	0	3	1/2	1
2	2	0	1/4	2	1	2	1/2	1
3	0	1	1/2	3	1	4	1	1
4	1	1	1	4	2	5	1/2	1
5	2	1	1/2	5	3	4	1	1
6	0	2	1/4	6	3	6	1/2	1
7	1	2	1/2	7	4	7	1	1
8	2	2	1/4	8	4	5	1	1
				9	5	8	1/2	1

**Table 5.1:** Location and volume of Voronoi boxes (left) and connection information (right).



**Figure 5.9:** Boundary conditions have to be included into a skeleton file.

```

10
11 begin main
12   assign psi[i=METAL1] = 10.0;
13   assign psi[i=METAL2] = 0.0;
14   solve; write;
15 end
    
```

source\_code/mesh-example1.sg

## 5.4 Auxiliary Functions

Two helper functions, `nsum` and `lsum`, are described in the following as they are used in `mesh-example1.sg`. They establish a convenient means of handling neighborhood information as needed by the box integration method. The node summation function `nsum` sums over all neighbors of a given point and is declared as

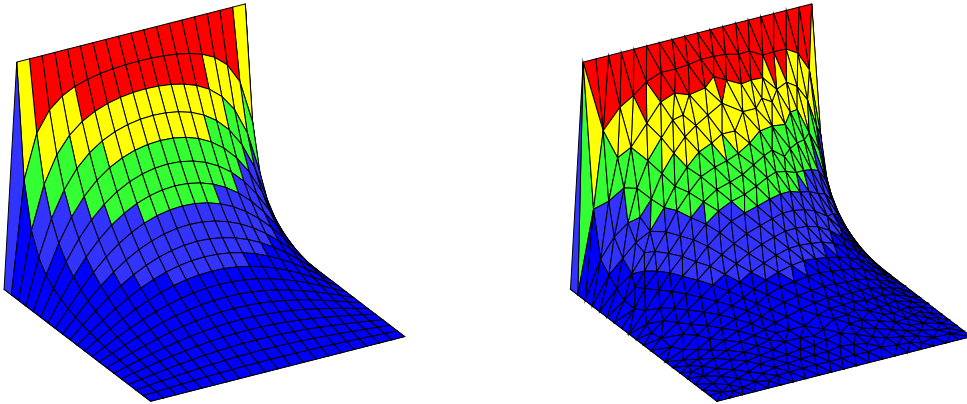


Figure 5.10: Rectangular Grid (left) and Triangular Mesh (right).

```
nsum(nodeIndex, neighborIdentifier, label, expr)
nodeIndex    ... current node
neighborIdentifier ... identifier that will run through all neighbors of nodeIndex
label        ... region label (from the skeleton file), may be all
expr         ... the expression to evaluate
```

The following `nsum_example.sg` file specifies the equation for all points inside region AIR. For each point  $i$  it sums the expression for the neighbor points  $j$ .

```
1 equ psi[i=AIR] -> nsum(i,j,all,
2                   ((psi[i]-psi[node(i,j)]) / elen(i,j)) * ilen(i,j)
3                   ) = rho[i] * area(i) / EPS;
```

source\_code/nsum\_example.sg

Here  $elen(i, j) = d_{i,j}$  is the edge length,  $ilen(i, j) = A_{i,j}$  the integration area (3D) or integration length (2D) (which is also a length in a projection, hence the name `ilen`) and  $area(i) = V_i$  the box volume. The command `node(i, j)` gives the index of  $j$ -th neighbor of  $i$ . Since we deal with a two-dimensional problem, the area the flux is passing through is made up of the perpendicular bisectors given by `ilen(i, j)`.

The label summation function `lsum` sums over all nodes with a particular label and evaluates an expression on each node. The example `lsum_example.sg` demonstrates how to sum up the expression for all points  $i$  that are labelled METAL1.

```
lsum(nodeIndex, label, expr)
nodeIndex    ... current node
label        ... group of nodes (from skeleton file)
expr         ... the expression to evaluate
```

```
1 assign Q1 = lsum(i,METAL1,
2               nsum(i,j,SIO2,
3                   D(EPSr,V[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)
4                   )
5               );
```

source\_code/lsum\_example.sg

## 5.5 Skeleton Mesh – Boundaries

In Section 5.3 the two-dimensional Laplace equation and its solution on an unstructured mesh including Dirichlet boundary conditions was considered. We have already given some hints on how to implement Neumann boundary conditions in Section 4.2.2, now we will look at it in more detail and also consider non-homogeneous Neumann boundary conditions. The implementation will be explained by the example of the discretized Poisson equation:

$$\sum_{j \in \mathcal{N}_i} D_{i,j} A_{i,j} = \rho_i V_i$$

In Fig. 5.11, a box at a domain boundary is illustrated. As discussed in section 4.2.2, Neumann boundary conditions with zero out-flux are implemented into SGFRAMEWORK by means of the implicit condition

$$\sum_{\substack{\text{all neighbors } j \\ \text{in current segment}}} D_{i,j} A_{i,j} = \rho_i V_i \quad (5.2)$$

which in SGFRAMEWORK reads as

```

1  equ psi[i=AIR] -> nsum(i,j,all,
2      ((psi[i]-psi[node(i,j)]) / elen(i,j)) * ilen(i,j)
3      ) = rho[i] * area(i) / EPS;
```

source\_code/nsum\_example.sg

Effectively, the code remains the same, all boundary handling is done within SGFRAMEWORK.

When dealing with non-zero out-flux, (5.2) has to be modified to take additional flux contributions at the boundary into account:

$$\sum_{\substack{\text{all neighbors } j \\ \text{in current segment}}} D_{i,j} A_{i,j} + D_5 \frac{l_5}{2} + D_6 \frac{l_6}{2} = \rho_i V_i.$$

```

1  nsum(i,j,all,
2      F(psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)
3      ) +
4  nsum(i,j,BND,
5      Fext * elen(i,j) / 2
6      ) = rho[i] * area(i) / EPS;
```

source\_code/nsum\_boundary\_example.sg

File `capacitor_geometry.sg` gives an example of how the structure in Fig. 5.12 is realized in SGFRAMEWORK.

Two different metal contacts METAL1 and METAL2 require further extensions to our skeleton file, as can be seen in File `capacitor_geometry.sk`. More edges and regions are needed as the METAL2 contact is area-like. The command NOFLUX has to be used for external Neumann boundaries, otherwise SGFRAMEWORK cannot triangulate the geometry.

```

1  const dx = 1.0, W1 = 8, W2 = 4, W3 = 8;
2  const W = W1 + W2 + W3, H = 20, d = 3;
```

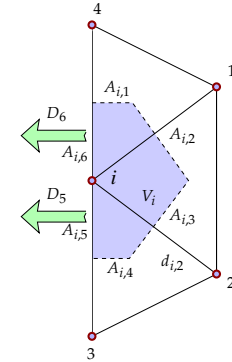


Figure 5.11: Points 4,  $i$  and 3 form a boundary.



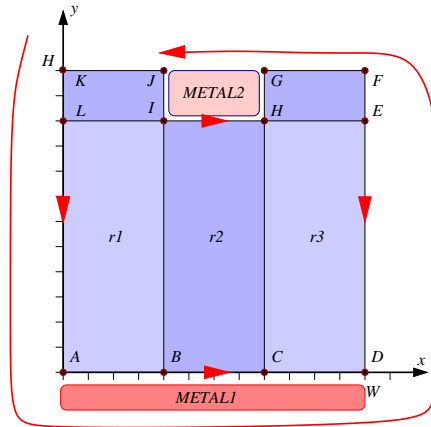
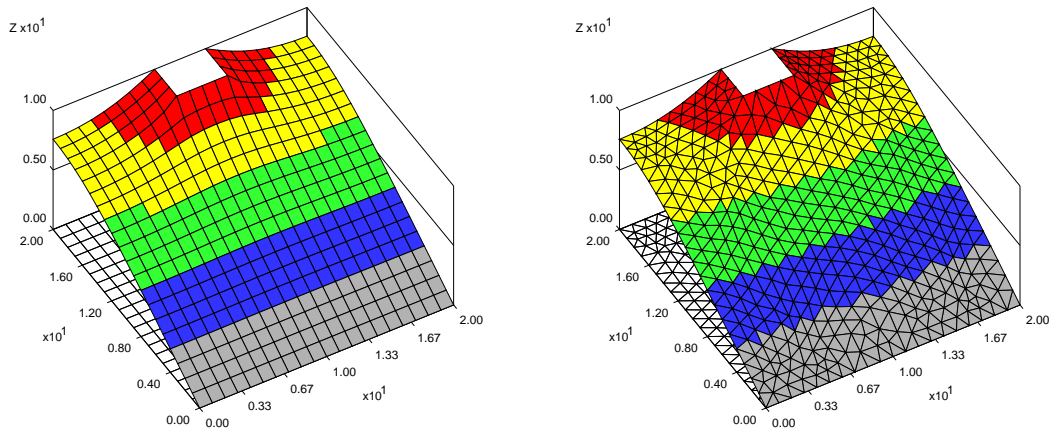


Figure 5.12: Complex geometry as METAL2 contact has three edges.

```

3
4 point pA = (0, 0);
5 point pB = (W1, 0);
6 point pC = (W1+W2, 0);
7 point pD = (W, 0);
8 point pE = (W, H-d);
9 point pF = (W, H);
10 point pG = (W1+W2, H);
11 point pH = (W1+W2, H-d);
12 point pI = (W1, H-d);
13 point pJ = (W1, H);
14 point pK = (0, H);
15 point pL = (0, H-d);
16
17 edge eAB = METAL1 [pA,pB] (dx, 0.0);
18 edge eBC = METAL1 [pB,pC] (dx, 0.0);
19 edge eCD = METAL1 [pC,pD] (dx, 0.0);
20 edge eED = NOFLUX [pE,pD] (dx, 0.0);
21 edge eFE = NOFLUX [pF,pE] (dx, 0.0);
22 edge eGF = NOFLUX [pG,pF] (dx, 0.0);
23 edge eGH = METAL2 [pG,pH] (dx, 0.0);
24 edge eIH = METAL2 [pI,pH] (dx, 0.0);
25 edge eJI = METAL2 [pJ,pI] (dx, 0.0);
26 edge eKJ = NOFLUX [pK,pJ] (dx, 0.0);
27 edge eKL = NOFLUX [pK,pL] (dx, 0.0);
28 edge eLA = NOFLUX [pL,pA] (dx, 0.0);
29
30 edge eHC = [pH,pC] (dx, 0.0);
31 edge eHE = [pH,pE] (dx, 0.0);
32 edge eIB = [pI,pB] (dx, 0.0);
33 edge eLI = [pL,pI] (dx, 0.0);
34
35 // rectangular mesh vs. triangular mesh (without 'RECTANGLES')
36 region r1 = SIO2 {eAB,eIB,eLI,eLA} RECTANGLES;
37 region r2 = SIO2 {eBC,eHC,eIH,eIB} RECTANGLES;
38 region r3 = SIO2 {eCD,eED,eHE,eHC} RECTANGLES;
39 region r4 = SIO2 {eHE,eFE,eGF,eGH} RECTANGLES;
40 region r6 = SIO2 {eLI,eJI,eKJ,eKL} RECTANGLES;
    
```

source\_code/capacitor\_geometry.sk



**Figure 5.13:** Simulation results on a rectangular grid (left) and triangular mesh (right) after a successful run of `capacitor_geometry.sg`.

```

1  mesh "capacitor_geometry.msh";
2
3  const  EPSr = 3.9;
4  const  EPSo = 8.854e-12;
5  var    psi[NODES], Q1, Q2;
6  unknown psi[all];
7  known  psi[METAL1], psi[METAL2];
8
9  func E(psi1,psi2,<h>)          return  (psi1-psi2) / h;
10 func D(<epsr>,psi1,psi2,<h>)    return  epsr * EPSo * E(psi1,psi2,h);
11
12 equ psi[i=SIO2] ->
13     nsum(i,j,all,D(EPSr,psi[i],psi[node(i,j)],elen(i,j))*ilen(i,j)) = 0.0;
14
15 begin main
16
17     assign psi[i=METAL1] = 0.0;
18     assign psi[i=METAL2] = 10.0;
19
20 solve;
21
22     assign Q1 = lsum(i,METAL1,
23         nsum(i,j,SIO2,D(EPSr,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)));
24     assign Q2 = lsum(i,METAL2,
25         nsum(i,j,SIO2,D(EPSr,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)));
26
27 write;
28
29 end
    
```

source\_code/capacitor\_geometry.sg

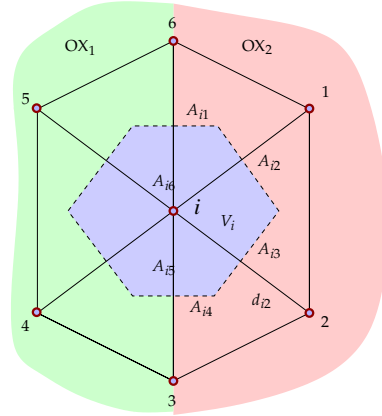


Figure 5.14: Interface conditions at the boundary between  $OX_1$  and  $OX_2$ .

## 5.6 Skeleton Mesh – Interfaces

Fig. 5.14 shows an interface condition between two different oxides. This time, we have to split the flux contributions into two terms, so that

$$\sum_{j=1,2,3,6} D_{i,j} A_{i,j}^{OX_2} + \sum_{j=3,4,5,6} D_{i,j} A_{i,j}^{OX_1} = \rho_i V_i. \quad (5.3)$$

The first term represents the contributions inside  $OX_2$  with the surrounding grid points 1, 2, 3, 6, while the second sum represents all contributions from  $OX_1$  enclosed by 3, 4, 5, 6. An implementation in SGFRAMEWORK is

```

1  equ psi[i=INT] -> nsum(i,j,OX1,F(...)*ilen(i,j)) +
2  nsum(i,j,OX2,F(...)*ilen(i,j)) = rho[i]*area(i)*EPS;
    
```

source\_code/mesh\_interfaces1.sg

Again, an example (Fig. 5.15 and 5.16) may provide a better understanding. Between two metal contacts we analyze the electrostatic field on unstructured mesh. The interface conditions have to be specified separately. The two contacts determine Dirichlet boundary conditions, while we assume homogeneous Neumann boundaries everywhere else on the boundary. The specifiers INT in the skeleton file capacitor\_example2.sk act as a label for edges that determine the interface. OX1 and OX2 as region labels are used as filters for the grid points. **Pay attention**<sup>1</sup> to specify the interface conditions before the segment conditions!

```

1  mesh "capacitor_example2.msh";
2
3  const EPSr1 = 3.9;
4  const EPSr2 = 11.8;
5  const EPSo = 8.854e-12;
6  var psi[NODES], Q1, Q2;
7  unknown psi[all];
8  known psi[METAL1], psi[METAL2];
9
10 func E(psi1,psi2,<h>) return (psi1-psi2)/h;
11 func D(<epsr>,psi1,psi2,<h>) return epsr*EPSo*E(psi1,psi2,h);
12
    
```

<sup>1</sup> to pay attention to sth. [peɪ ˈeɪ.tən.tʃən]: auf etwas achten

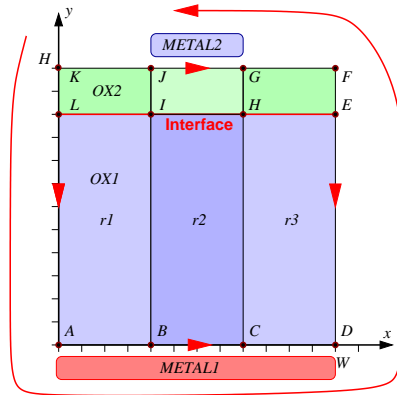


Figure 5.15: Geometry with additional interface between  $OX_1$  and  $OX_2$ .

```

13 //interface has to come first
14 equ psi[i=INT] -> nsum(i,j,OX1,
15                     D(EPSr1,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)) +
16                     nsum(i,j,OX2,
17                     D(EPSr2,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j))
18                     = 0.0;
19
20 //Laplace for OX1
21 equ psi[i=OX1] ->
22     nsum(i,j,OX1,D(EPSr1,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)) = 0.0;
23
24 //Laplace for OX2
25 equ psi[i=OX2] ->
26     nsum(i,j,OX2,D(EPSr2,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)) = 0.0;
27
28 begin main
29
30     assign psi[i=METAL1] = 0.0;
31     assign psi[i=METAL2] = 10.0;
32
33     solve;
34
35     assign Q1 = lsum(i,METAL1,
36                     nsum(i,j,OX1,D(EPSr1,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)));
37     assign Q2 = lsum(i,METAL2,
38                     nsum(i,j,OX2,D(EPSr2,psi[i],psi[node(i,j)],elen(i,j)) * ilen(i,j)));
39
40     write;
41
42 end
    
```

source\_code/capacitor\_example2.sg

```

1 const dx = 1.0, W1 = 8, W2 = 4, W3 = 8;
2 const W = W1 + W2 + W3, H = 20, d = 3;
3
4 point pA = (0, 0);
5 point pB = (W1, 0);
6 point pC = (W1+W2, 0);
7 point pD = (W, 0);
8 point pE = (W, H-d);
    
```

```

9 point pF = (W, H);
10 point pG = (W1+W2, H);
11 point pH = (W1+W2, H-d);
12 point pI = (W1, H-d);
13 point pJ = (W1, H);
14 point pK = (0, H);
15 point pL = (0, H-d);
16
17 edge eAB = METAL1 [pA,pB] (dx, 0.0);
18 edge eBC = METAL1 [pB,pC] (dx, 0.0);
19 edge eCD = METAL1 [pC,pD] (dx, 0.0);
20
21 edge eED = NOFLUX [pE,pD] (dx, 0.0);
22 edge eFE = NOFLUX [pF,pE] (dx, 0.0);
23 edge eLA = NOFLUX [pL,pA] (dx, 0.0);
24 edge eKL = NOFLUX [pK,pL] (dx, 0.0);
25
26 edge eHE = INT [pH,pE] (dx, 0.0);
27 edge eIH = INT [pI,pH] (dx, 0.0);
28 edge eLI = INT [pL,pI] (dx, 0.0);
29
30 edge eGF = NOFLUX [pG,pF] (dx, 0.0);
31 edge eJG = METAL2 [pJ,pG] (dx, 0.0);
32 edge eKJ = NOFLUX [pK,pJ] (dx, 0.0);
33
34 edge eJI = NOFLUX [pJ,pI] (dx, 0.0);
35 edge eGH = NOFLUX [pG,pH] (dx, 0.0);
36
37 edge eHC = [pH,pC] (dx, 0.0);
38 edge eIB = [pI,pB] (dx, 0.0);
39
40 //rectangluar mesh vs. triangular mesh (without 'RECTANGLES')
41 region r1 = OX1 {eAB,eIB,eLI,eLA} RECTANGLES;
42 region r2 = OX1 {eBC,eHC,eIH,eIB} RECTANGLES;
43 region r3 = OX1 {eCD,eED,eHE,eHC} RECTANGLES;
44 region r4 = OX2 {eHE,eFE,eGF,eGH} RECTANGLES;
45 region r5 = OX2 {eIH,eGH,eJG,eJI} RECTANGLES;
46 region r6 = OX2 {eLI,eJI,eKJ,eKL} RECTANGLES;
    
```

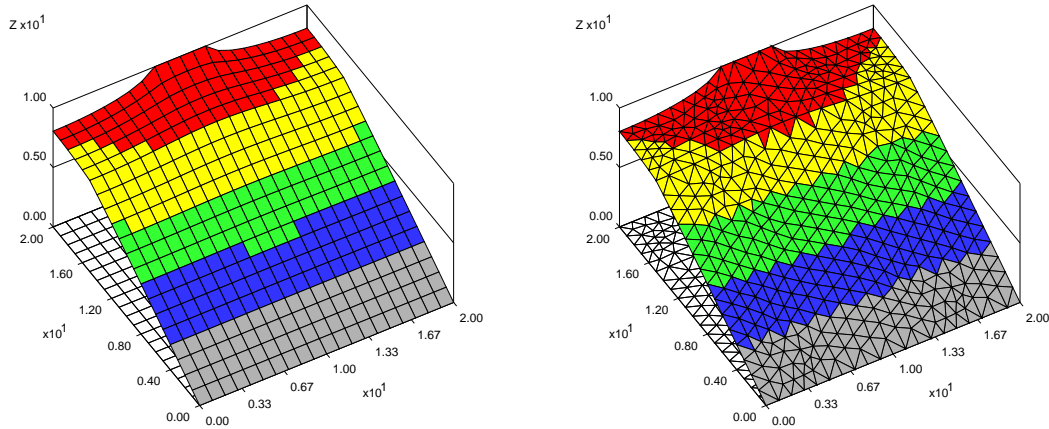
source\_code/capacitor\_example2.sk

## 5.7 Mesh Refinement

Mesh refinement is useful to increase the resolution in certain areas of the device. Sophisticated algorithms solve a given problem on a coarse mesh, apply some error indicators to this **preliminary**<sup>1</sup> solution and refine the mesh in those regions. In SGFRAMEWORK, refinement has to be done manually: The role of the error estimator is assigned to the user, who has to supply a function describing the level of refinement.

The refinement criteria can be specified in the skeleton files right after the region statements in the form of one or more refinement functions. In order to determine whether a mesh element should be refined, the mesh refinement routine will evaluate the refinement function(s) at each of the element's vertices. If the measure between two grid points  $M_{i,j}$  exceeds a reference distance  $d$ , the element will be refined, i.e. it will be split into two or more (hence smaller)

<sup>1</sup> preliminary [pri'lim.i.n<sup>o</sup>r.i]: vorläufig, vorübergehend



**Figure 5.16:** Results obtained on a rectangular grid (left) and a triangular mesh (right) for the input file `capacitor_example2.sg`.

elements. The function itself has a name, two refinement parameters and a body, using the syntax:

```

refine name(measure, refDistance) = expr
name      ... name of the criterion (arbitrary)
measure   ... refinement measure
refDistance ... reference distance
expr      ... the expression to compare with the reference distance
    
```

The refinement measure describes how the current vertex distances are weighted. It may be linear, logarithmic, or signedlog as defined by the following expressions:

$$M_{\text{lin}} = x, \quad M_{\text{log}} = \log(x), \quad M_{\text{sllog}} = \text{sign}(x) \log(1 + |x|).$$

For example, in order to determine whether a triangular element needs to be refined in a linear way,

$$\begin{aligned}
 M_{1,2} &= |M(N_1) - M(N_2)| \\
 M_{2,3} &= |M(N_2) - M(N_3)| \\
 M_{1,3} &= |M(N_1) - M(N_3)|
 \end{aligned}$$

is calculated and compared to  $d$ . If any of the  $M_{i,j}$  exceeds  $d$ , i.e. the refinement function varies strongly between the edges of the triangle, the triangle is refined. Note that refinement occurs where the refinement function has a large gradient, while the magnitude of the refinement function plays a minor role.

The last section of the refinement criterion determines the minimum and maximum number of divisions and the minimum and maximum lengths of edges. Each element in the initial grid is assigned a division level of zero. At each refinement step the division level of the current element is increased by one. The mesh refinement program loops through each element and if a refinement is necessary, a check is performed on its neighbors. If a neighboring element's division level is less than that of the considered one, the neighboring element is refined first. This is done to prevent too large variations in the dimensions of two adjacent elements. One possible refinement indicator function is shown in Fig. 5.17; it can be implemented as

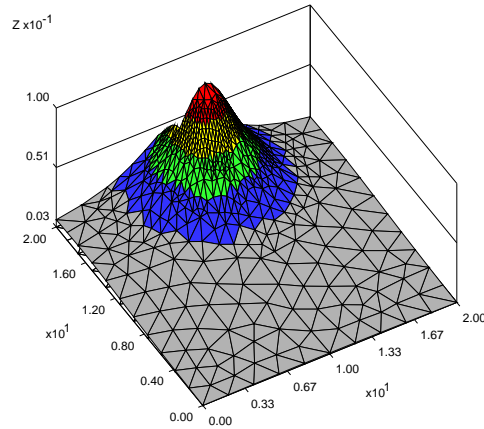


Figure 5.17: Distance Function used in `refine.sk`.

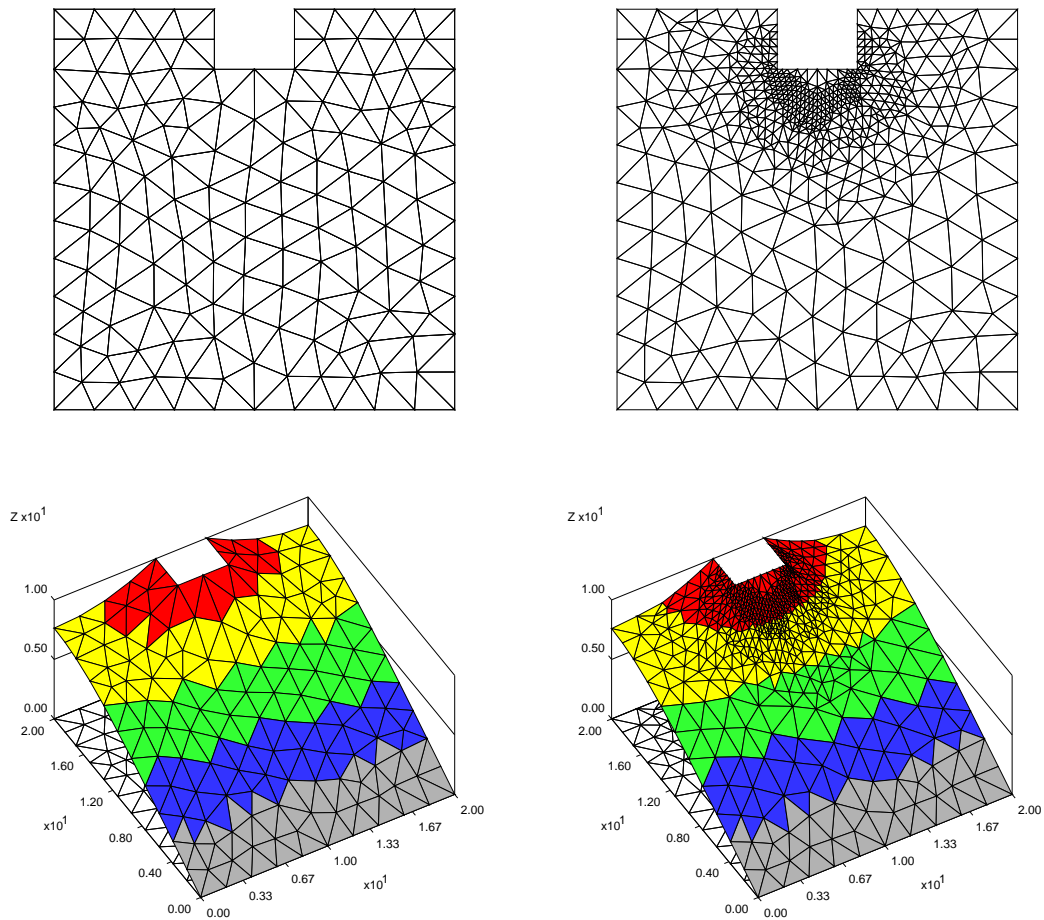
```

1  refine Dist(linear,0.01) = 1/(10 + sq(x-W/2) + sq(y-(H-d)));
2  set minimum divisions = 0;
3  set maximum divisions = 5;

```

source\_code/refinement\_example1.sk

The chosen refinement function results in a refinement at the edges of the upper metal contact, where the distance function has its maximum (Fig 5.18). The number of nodes increases from  $N = 155$  to  $N' = 541$  as the capacity  $C = 34.53$  pF changes to  $C' = 34.10$  pF, suggesting an initial error of 1.3%.



**Figure 5.18:** Mesh before (top left) and after refinement (top right). Potential without (bottom left) and with refinement (bottom right) using those meshes.





## Chapter 6

# Transport Phenomena and their Numerical Analysis

In the **preceding**<sup>1</sup> chapters we have learnt how to discretize derivatives and finally how to solve Poisson's equation. Poisson's equation describes an equilibrium distribution of some physical quantity and is numerically well-behaved. In this chapter we are going to simulate the transient behavior of a device, which means that we have to **cope**<sup>2</sup> with an additional time dependence. Our ultimate aim is the solution of the drift-diffusion model as given in (1.23) to (1.25), but before we are ready to do so, we will analyze the different contributions to each of the equations.

The equation describing conservation of a quantity  $n$  in general form is

$$\text{Increase in Time} + \text{Outflux} = \text{Production Rate} ,$$

or, mathematically

$$\frac{\partial n}{\partial t} + \nabla \cdot \Gamma = s . \quad (6.1)$$

The flux term  $\Gamma$  may consist of two distinct contributions:

- *Diffusion*: For  $\Gamma = -D\nabla n$  with (constant) diffusion coefficient  $D$ , the conservation equation becomes

$$\frac{\partial n}{\partial t} - D\nabla^2 n = s ,$$

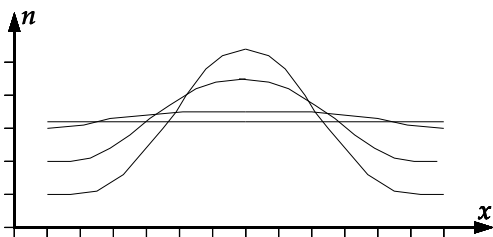
which is a parabolic partial differential equation. In pure diffusive processes, relaxation to an equilibrium determined by the boundary conditions can be observed (cf. Fig. 6.1). For example, with homogeneous Neumann boundary conditions (i.e. no out-flux),  $\int_V n(x)dV = \text{const.}$ , while this is not the case with inhomogeneous boundary conditions.

- *Convection (Drift)*:  $\Gamma = n\mu E$  with (constant) mobility  $\mu$ . The conservation equation then is a hyperbolic partial differential equation of first order,

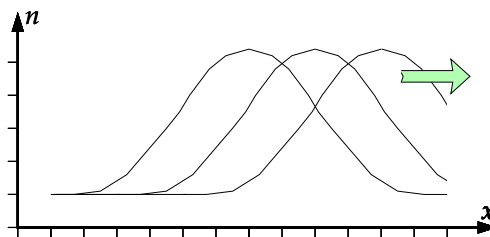
$$\frac{\partial n}{\partial t} + \mu \nabla \cdot (nE) = S .$$

---

<sup>1</sup> **preceding** [pri'si:diŋ]: vorangegangen    <sup>2</sup> **to cope** [koup]: zurechtkommen, beherrschen



**Figure 6.1:** A pure diffusive process relaxes to an equilibrium over time if there is no out-flux.



**Figure 6.2:** A pure convective process describes travelling waves.

Pure convective processes describe travelling waves, meaning that the initial shape of  $n(x)$  remains unchanged (cf. Fig. 6.2). Because of the wave-like behavior, there is a flux through boundaries.

In most systems both contributions are present, complicating the solution process. This chapter deals with discretization in the time domain and studies the numerical stability of these schemes.

## 6.1 Discretization in the Time Domain

So far we have dealt with the approximation of derivatives in spatial coordinates only. For the Poisson equation, both the five-point and the nine-point stencil lead to **satisfactory**<sup>1</sup> results without numerical stability issues. For the time discretization of parabolic partial differential equations, this is not the case anymore.

Let us consider the ordinary differential equation

$$\frac{dn(t)}{dt} = h(n(t), t)$$

for an arbitrary, smooth function  $h(\cdot, \cdot)$ . Without any further assumptions on  $h$  and knowing  $n(t_k)$  we have mainly two possibilities to obtain  $n(t_{k+1})$  after discretization of the time-derivative: We can evaluate the right hand side at  $t_k$  or at  $t_{k+1}$ .

In the first case, the scheme is called *forward Euler scheme* (aka. *explicit Euler scheme*):

$$\frac{n(t_{k+1}) - n(t_k)}{t_{k+1} - t_k} = h(n(t_k), t_k) \iff n(t_{k+1}) = n(t_k) + (t_{k+1} - t_k)h(n(t_k), t_k)$$

Since  $n(t_k)$  is known, we immediately obtain  $n(t_{k+1})$  and are able to proceed to the next time step.

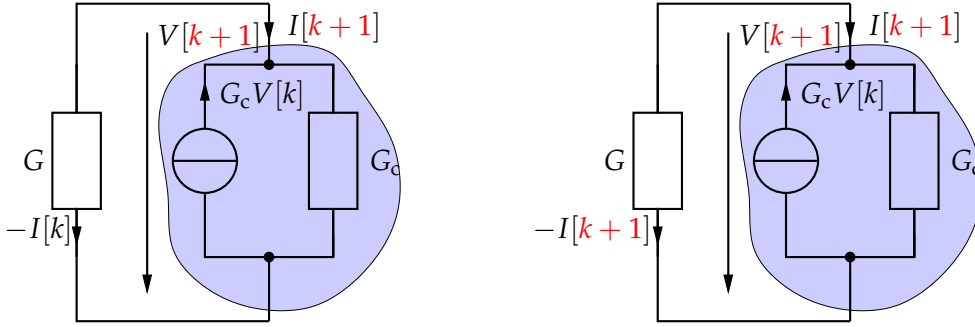
The second possibility for the evaluation of the right hand side leads to the *backward Euler scheme* (aka. *implicit Euler scheme*):

$$\frac{n(t_{k+1}) - n(t_k)}{t_{k+1} - t_k} = h(n(t_{k+1}), t_{k+1})$$

<sup>1</sup> **satisfactory** [sæt.ɪs'fækt.ɹ.i]: zufriedenstellend

The solution process for  $n(t_{k+1})$  now strongly depends on  $h$ . Especially if  $h$  is non-linear, one usually has to rely on *Newton's method* (cf. Section 3.6) and solve a (system of) non-linear equation(s) within each time step.

If we denote the step size in time with  $\Delta t := t_{k+1} - t_k$ , then both the solutions obtained by the explicit Euler scheme and the backward Euler scheme converge to the true solution for  $\Delta t \rightarrow 0$  with order  $O(\Delta t)$ .



**Figure 6.3:** Time discretization of a discharging capacitor with Forward Euler (left) and Backward Euler (right) scheme

To become **familiar**<sup>1</sup> with the Euler schemes, let us have a look at an initially charged capacitor that discharges over a resistor with conductance  $G = 1/R$ , shown in Fig. 6.3. The governing equation for the voltage over the capacitor is

$$C \frac{dV(t)}{dt} = -GV(t).$$

A discretization using the explicit Euler scheme leads to

$$C \frac{V_{k+1} - V_k}{\Delta t} = -GV_k.$$

With the auxiliary conductance  $G_C := C/\Delta t$  we get

$$V_{k+1} = V_k \frac{G_C - G}{G_C} = V_0 \left( \frac{G_C - G}{G_C} \right)^{k+1}.$$

A discretization using the implicit Euler scheme results in

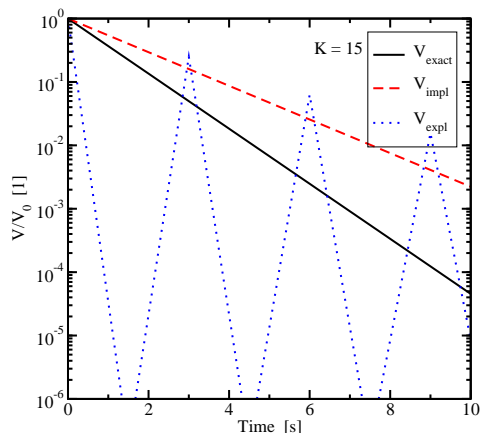
$$C \frac{V_{k+1} - V_k}{\Delta t} = -GV_{k+1} \iff V_{k+1} = V_k \frac{G_C}{G_C + G} = V_0 \left( \frac{G_C}{G_C + G} \right)^{k+1}$$

As initial condition we set  $V_0 = 1$  and compare with the analytical solution  $V(t) = \exp(-t/\tau)$ , where  $\tau = RC = \Delta t G_C / G$ . Since  $G_C / G = \tau / \Delta t$ , we can rewrite the equations obtained by the Euler schemes as

$$V_k = V_0 \left( 1 - \frac{\Delta t}{\tau} \right)^k \quad (\text{explicit Euler})$$

$$V_k = V_0 \left( 1 + \frac{\Delta t}{\tau} \right)^{-k} \quad (\text{implicit Euler})$$

<sup>1</sup> **familiar** [fə'mil.i.jə]: vertraut, geläufig



**Figure 6.4:** A discretization with the forward Euler scheme leads to numerical instabilities, whereas the backward Euler scheme remains stable.

**On closer inspection**<sup>1</sup>, one can see that the solution obtained by the implicit Euler scheme tends to zero as  $k \rightarrow \infty$  for all  $\tau, \Delta t > 0$ . However, the explicit Euler scheme leads to  $V_1 < 0$  if  $\frac{\tau}{\Delta t} > 1$ , which is illustrated in Fig. 6.4. This behavior is not just specific for this example only, as the next section will show.

Note that for fixed  $t = k\Delta t$  the explicit Euler scheme gives us

$$\begin{aligned} \lim_{\Delta t \rightarrow 0} V(t) &= V(0) \lim_{k \rightarrow \infty} \left(1 - \frac{t}{\tau k}\right)^k, \\ &= V(0) e^{-\frac{t}{\tau}} \\ &= V(t), \end{aligned}$$

while for the implicit Euler scheme we find

$$\begin{aligned} \lim_{\Delta t \rightarrow 0} V(t) &= V(0) \lim_{k \rightarrow \infty} \left(1 + \frac{t}{\tau k}\right)^{-k} && \text{(implicit Euler)} \\ &= V(0) 1/e^{\frac{t}{\tau}} \\ &= V(t). \end{aligned}$$

An implementation in SGFRAMEWORK reads

```

1  const C = 0.1,    R = 10.0,    G = 1/R;
2  const K = 1,     dt = 0.1*K,   tmax = 10;
3  const GC = C/dt, N = 100/K,   V0 = 1;
4
5  var t[N+1], Vimpl[N+1], Vexpl[N+1], Vexact[N+1];
6
7  begin main
8    assign t[i=all]      = i*dt;
9    assign Vimpl[i=0]    = V0;
10   assign Vimpl[i=1..N] = Vimpl[i-1] * GC/(GC+G);
11
12   assign Vexpl[i=0]    = V0;
13   assign Vexpl[i=1..N] = Vexpl[i-1] * (GC-G)/GC;
    
```

<sup>1</sup> **on closer inspection** [ɔn kloʊsər m'spek.ʃən]: bei näherer Betrachtung

```

14
15 assign Vexact[i=0..N] = V0 * exp(-t[i]/(R*C));
16 write;
17 end

```

source\_code/RC\_circuit.sg

## 6.2 Stability of Discretization Schemes

The Nyquist-Shannon sampling theorem states that if a signal is band-limited to  $W$ , a sampling frequency of  $2W$  is sufficient for exact reconstruction. A related question arises for numerical discretization schemes: Will the chosen discretization yield meaningful results? Which schemes are “better” than others? Which is the maximum allowed distance between grid points?

We will investigate stability using a technique introduced by John von Neumann and start with the continuous *Fourier transform*,

$$N(\chi, t) = \mathcal{F}\{n(x, t)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-j\chi x} n(x, t) dx,$$

$$n(x, t) = \mathcal{F}^{-1}\{N(\chi, t)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{j\chi x} N(\chi, t) d\chi$$

and the Fourier transform for discrete functions,

$$N(\xi, k) = \mathcal{F}\{n_{i,k}\} = \frac{1}{\sqrt{2\pi}} \sum_{i=-\infty}^{\infty} e^{-j\Delta x \xi i} n_{i,k} \Delta x,$$

$$n_{i,k} = \mathcal{F}^{-1}\{N(\xi, k)\} = \frac{1}{\sqrt{2\pi}} \int_{-\pi/\Delta x}^{\pi/\Delta x} e^{j\Delta x \xi i} N(\xi, k) d\xi,$$

where  $\xi \in [-\pi/\Delta x, \pi/\Delta x]$ . For displacement and differentiation in space in the continuous case,

$$\mathcal{F}\{n(x + x_0, t)\} = e^{j\chi x_0} N(\chi, t), \quad \mathcal{F}\left\{\frac{\partial n(x, t)}{\partial x}\right\} = j\chi N(\chi, t)$$

holds, whereas for the discrete version we have

$$\mathcal{F}\{n(i + i_0, k)\} = e^{j\Delta x \xi i_0} N(\xi, t). \tag{6.2}$$

Let us first consider a Forward Euler discretization of the diffusion equation,

$$n_{i,k+1} = n_{i,k} + \eta D[n_{i+1,k} - 2n_{i,k} + n_{i-1,k}], \tag{6.3}$$

where  $\eta = \frac{\Delta t}{(\Delta x)^2}$ . We use as initial condition  $n_{i,0} = N_0 \delta_{i-i_0}$  (think about why the choice  $n_{i,0} = N_0$  does not yield a too meaningful example!). At  $k = 1$  we get

$$\begin{aligned} n_{i_0-1,1} &= 0 + \eta D(N_0 - 0 + 0) = \eta D N_0 \\ n_{i_0,1} &= N_0 + \eta D(0 - 2N_0 + 0) = N_0(1 - 2\eta D) \\ n_{i_0+1,1} &= 0 + \eta D(0 - 0 + N_0) = \eta D N_0 \end{aligned} \tag{6.4}$$

In particular, if  $1 - 2\eta D < 0 \Leftrightarrow 2\eta D > 1 \Leftrightarrow \Delta t > (\Delta x)^2/2D$  we get  $n_{i_0,1} < 0$ , which is meaningless from a physical point of view.

A more stringent investigation of the stability properties of (6.3) is as follows: A transformation of (6.3) into the frequency domain and application of the displacement properties yields

$$\begin{aligned} N(\xi, k+1) &= N(\xi, k) + \eta D N(\xi, k) \left( e^{j\Delta x \xi} - 2 + e^{-j\Delta x \xi} \right) \\ &= N(\xi, k) (1 + 2\eta D (\cos(\Delta x \xi) - 1)). \end{aligned}$$

If we now define the *amplification factor*  $g(\Delta x \xi) := 1 + 2\eta D (\cos(\Delta x \xi) - 1)$  we immediately see that

$$N(\xi, k+1) = N(\xi, k) g(\Delta x \xi) = N(\xi, 0) (g(\Delta x \xi))^{k+1},$$

so if we require that the solution stays bounded as  $k \rightarrow \infty$ , we require

$$|g(\Delta x \xi)| \leq 1.$$

For our example, we have  $g(\Delta x \xi) = 1 + 2\eta D (\cos(\Delta x \xi) - 1) = 1 - 4\eta D \sin^2\left(\frac{\Delta x \xi}{2}\right)$ , so in order to fulfill  $-1 \leq g(\Delta x \xi) \leq 1$ ,

$$0 \leq 4\eta D \sin^2\left(\frac{\Delta x \xi}{2}\right) \leq 2.$$

The lower inequality is trivially fulfilled, but the upper requires

$$2\eta D \leq 1 \quad \overset{\eta = \frac{\Delta t}{(\Delta x)^2}}{\Leftrightarrow} \quad \Delta t \leq \frac{(\Delta x)^2}{2D}. \quad (6.5)$$

Therefore, discretization in time and space are tightly connected for the Forward Euler scheme. Note that we have already found such a condition in our simple example (6.4).

Let us have a look at the Backward Euler scheme as well. We will use the general scheme of replacing  $n_{i,k}$  by  $(g(\Delta x \xi))^k e^{-ji\Delta x \xi}$  in

$$n_{i,k} - n_{i,k-1} = \eta D (n_{i+1,k} - 2n_{i,k} + n_{i-1,k})$$

and obtain

$$g(\Delta x \xi) - 1 = \eta D g(\Delta x \xi) \left( e^{-j\Delta x \xi} - 2 + e^{j\Delta x \xi} \right) = 2\eta D g(\Delta x \xi) (\cos(\Delta x \xi) - 1).$$

With this, we find  $g(\Delta x \xi)$  to be

$$g(\Delta x \xi) = \frac{1}{1 + 4\eta D \sin^2\left(\frac{\Delta x \xi}{2}\right)},$$

so  $-1 \leq g(\Delta x \xi) \leq 1$  is always fulfilled. Therefore, the Backward Euler scheme is said to be *unconditionally stable*: The numerical results will be stable (which does not necessarily mean that they are good), no matter which discretization is chosen.

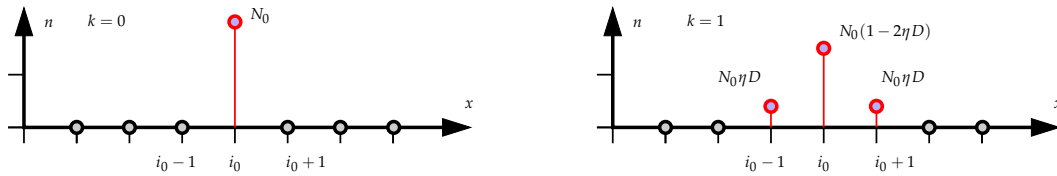


Figure 6.5: A Dirac impulse as initial condition confirms the results of the von Neumann analysis.

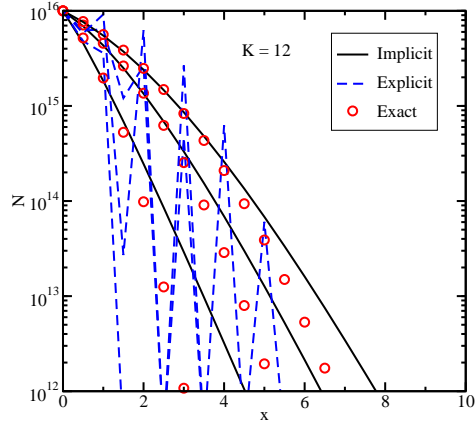


Figure 6.6: An explicit Euler scheme for the time discretization of the diffusion equation leads to an unstable numerical solution for large (time-)step sizes, whereas the implicit Euler scheme is unconditionally stable.

### 6.3 Diffusive Problems

For  $\Gamma = -D\nabla n$  and  $s = 0$  in (6.1), the *diffusion equation*

$$\frac{\partial n}{\partial t} = D\nabla^2 n,$$

$$n(x, t = 0) = n_0(x)$$

without source term is obtained. It is also known as *heat equation*, since the spreading of heat is also described by this equation.

Unlike the wave equation, the diffusion equation has very strong smoothing properties. This means that for  $t > 0$  the solution  $n(x, t)$  is smooth (infinitely differentiable), even if this is not true for  $n_0$ .

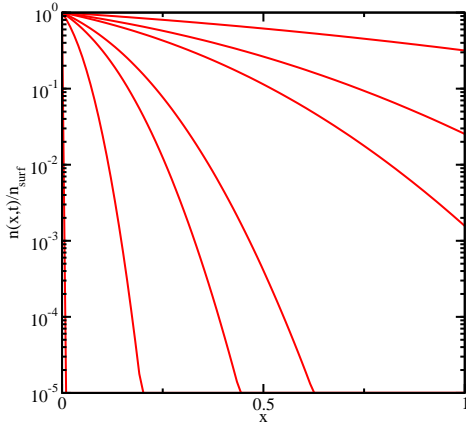
We are interested in the description of dopant diffusion, which will serve as an explanatory framework in this section. For this we assume the interface between the dopant source ( $x < 0$ ) and the semiconductor ( $x > 0$ ) at  $x = 0$ .

The diffusion coefficient  $D$  in a semiconductor is given by an *Arrhenius law*

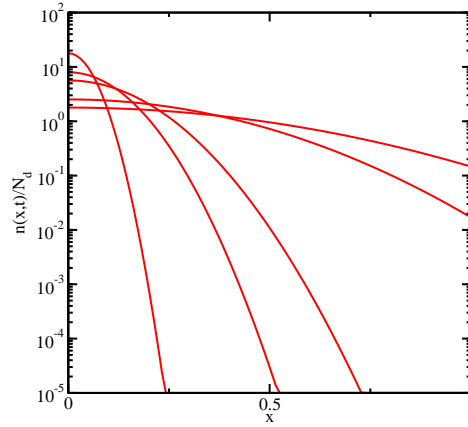
$$D = D_0 \exp\left(-\frac{E_{\text{act}}}{k_B T_L}\right).$$

The parameters for phosphorus are  $D_0 = 3.85 \text{ cm}^2/\text{s}$  and  $E_{\text{act}} = 3.66 \text{ eV}$ , which results in





**Figure 6.7:** Dopant diffusion from an inexhaustible source (i.e. Dirichlet boundary condition) at  $x = 0$ .



**Figure 6.8:** Dopant drive-in diffusion with a constant number of dopants (i.e. Neumann boundary condition) at  $x = 0$ .

$D = D_0 e^{-142}$  at  $T_L = 300$  K. The diffusion front (which can be defined via the variance of  $u$ ) at time  $t$  is located at depth  $\sqrt{Dt}$ . For  $t > 0$ , there are two types of boundary conditions possible:

- *Constant surface concentration:* The source of dopants is **inexhaustible**<sup>1</sup>, thus  $n(0, t) = n_{\text{surf}}$ . Since dopants cannot spread over the whole device instantly, we have for every finite  $t$

$$\lim_{x \rightarrow \infty} n(x, t) = 0.$$

In one spatial dimension, the analytical solution for the dopant diffusion problem is

$$n(x, t) = n_{\text{surf}} \operatorname{erfc}\left(\frac{x}{\sqrt{Dt}}\right).$$

As  $t \rightarrow \infty$ , the equilibrium (which is formally obtained from the solution of the equation that remains if all time derivatives are set to zero) is

$$\lim_{t \rightarrow \infty} n(x, t) = n_{\text{surf}}.$$

- *Constant dose (drive-in diffusion):* Only a fixed amount of dopants is available, so

$$N_d(t) = \int_0^\infty n(x, t) dx = \text{const.}$$

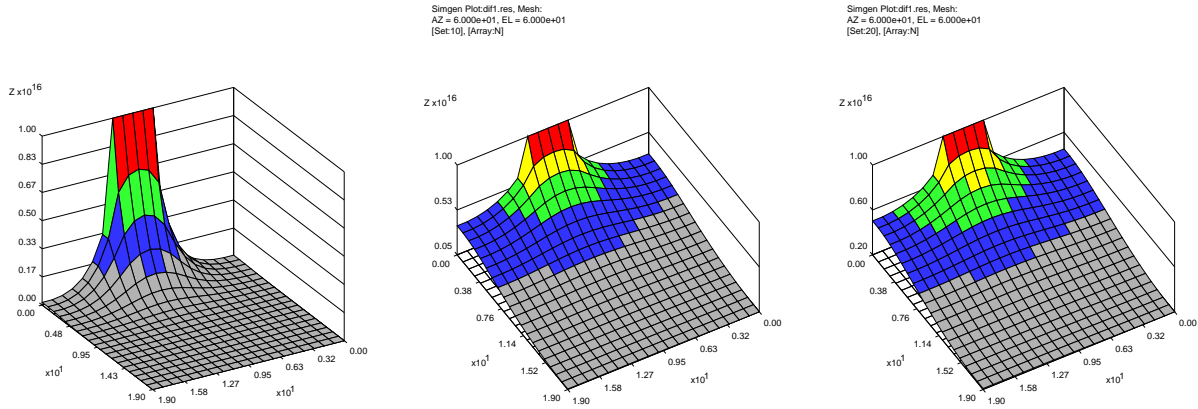
As before,

$$\lim_{x \rightarrow \infty} n(x, t) = 0$$

for every  $t < \infty$ . An analytical solution in one spatial dimension is

$$n(x, t) = \frac{N_d}{\sqrt{\pi Dt}} \exp\left(-\frac{x^2}{4Dt}\right),$$

<sup>1</sup> **inexhaustible** [m.ig'zɔ:.stɪ.b]: unerschöpflich



**Figure 6.9:** Results for dopant diffusion with constant surface concentration for three equidistant time-steps.

so the equilibrium solution for a device with infinite length is

$$\lim_{t \rightarrow \infty} n(x, t) = 0.$$

However, real devices have a finite length  $l$ , so that

$$\lim_{t \rightarrow \infty} n(x, t) = \frac{N_d}{l}.$$

Let us consider an example for each boundary condition. In the case of constant surface concentration at  $\Gamma_{\text{surf}}$ , the initial condition in our simulation domain  $\Omega$  is

$$n(x, t = 0) = \begin{cases} n_{\text{surf}}, & x \in \Gamma_{\text{surf}} \\ 0, & x \in \Omega \setminus \Gamma_{\text{surf}} \end{cases}. \quad (6.6)$$

The boundary condition is  $n(x, t) = n_{\text{surf}}$  for  $x \in \Gamma_{\text{surf}}$ , and there is no out-flux elsewhere, hence  $\Gamma \cdot \nabla n = 0$ , where  $\Gamma$  is the outer normal vector at each point on the boundary arc  $\partial\Omega \setminus \Gamma_{\text{surf}}$ .

An implementation in SGFRAMEWORK reads

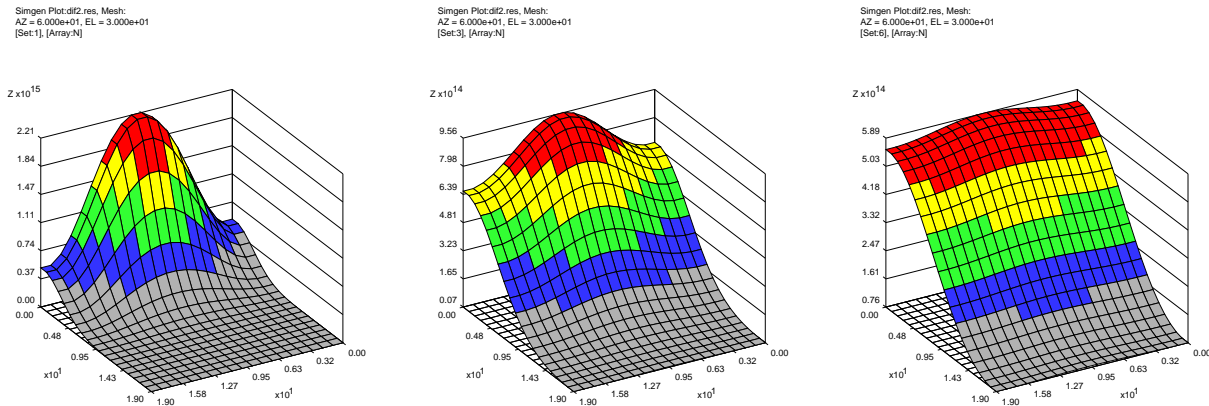
```

1 // constants
2 const NX    = 20;           // number of mesh points, x direction
3 const NY    = 20;           // number of mesh points, y direction
4 const L     = 10.0e-4;      // distance scaling (cm)
5 const dx    = L/NX;
6 const dy    = L/NY;
7 const T     = 300.0;        // operating temperature
8 const D     = 1.0e-2;        // diffusion constant
9 const dt    = 0.1*dx*dx/(2*D); // time step — obeys von Neumann condition
10 const Xmin  = 8;
11 const Xmax  = 12;
12
13 // declare variables and specify the unknowns
14 var itime;
15 var iwrite;
16 var x[NX];
17 var y[NY];
18 var N[NX,NY], Nold[NX,NY];
19 unknown N[all, all];
20 known N[Xmin..Xmax, 0];
    
```

```

21
22 // Diffusion equation
23 equ N[i=1..NX-2,j=1..NY-2] ->
24     D*(N[i+1,j ]-2*N[i,j ]+N[i-1,j ])/(dx*dx) +
25     D*(N[i ,j+1]-2*N[i,j ]+N[i ,j-1])/(dy*dy) -
26     (N[i,j]-Nold[i,j])/dt = 0.0;
27
28 //boundary at i=NX-1
29 equ N[i=NX-1,j=1..NY-2] -> N[i,j] -N[i-1,j] = 0.0;
30
31 // boundary at i=0
32 equ N[i=0,j=1..NY-2] -> N[i,j] - N[i+1,j] = 0.0;
33
34 // boundary at j=NY-1
35 equ N[i=1..NX-2,j=NY-1] -> N[i,j] - N[i,j-1] = 0.0;
36
37 // first boundary at j=0
38 equ N[i=1..Xmin-1,j=0] -> N[i,j] - N[i,j+1] = 0.0;
39
40 // second boundary at j=0
41 equ N[i=Xmax+1..NX-2,j=0] -> N[i,j] - N[i,j+1] = 0.0;
42
43 // corner boundary conditions:
44 equ N[i=0,j=0] -> N[i,j] - N[i+1,j+1] = 0.0;
45 equ N[i=0,j=NY-1] -> N[i,j] - N[i+1,j-1] = 0.0;
46 equ N[i=NX-1,j=0] -> N[i,j] - N[i-1,j+1] = 0.0;
47 equ N[i=NX-1,j=NY-1] -> N[i,j] - N[i-1,j-1] = 0.0;
48
49 // set the numerical algorithm parameters
50 set NEWTON DAMPING = 3;
51 set NEWTON ACCURACY = 1.0e+4;
52 set NEWTON ITERATIONS = 100;
53 set LINSOL ALGORITHM = GAUSSELIM;
54 set LINSOL FILL = INFINITY;
55
56 begin InitVars
57     assign x[i=all] = i*dx;
58     assign y[i=all] = i*dy;
59     assign N[i=5..15,j=0] = 1.0e16;
60     assign Nold[i=all,j=all] = N[i,j];
61 end
62
63 begin main
64     assign itime = 0;
65     call InitVars;
66     while (itime < 2000) begin
67         assign Nold[i=all,j=all] = N[i,j];
68         solve;
69         assign itime = itime + 1;
70         assign iwrite = iwrite + 1;
71         if (iwrite == 100) begin
72             write;
73             assign iwrite = 0;
74         end
75     end
76 end
    
```

source\_code/diffusion\_example1.sg



**Figure 6.10:** Results for drive-in dopant diffusion for three equidistant time-steps. The initial local peak concentration spreads over the whole domain and ultimately results in a constant (and smaller compared to the initial peak) concentration over the whole device.

Results are shown in Fig. 6.9. It can be seen that initially there is a steep gradient of the dopant concentration, because they only penetrate the surface. As time advances, the dopants spread over the whole simulation domain until finally after an infinite amount of time a constant dopant level  $n_{\text{surf}}$  in the whole material is reached.

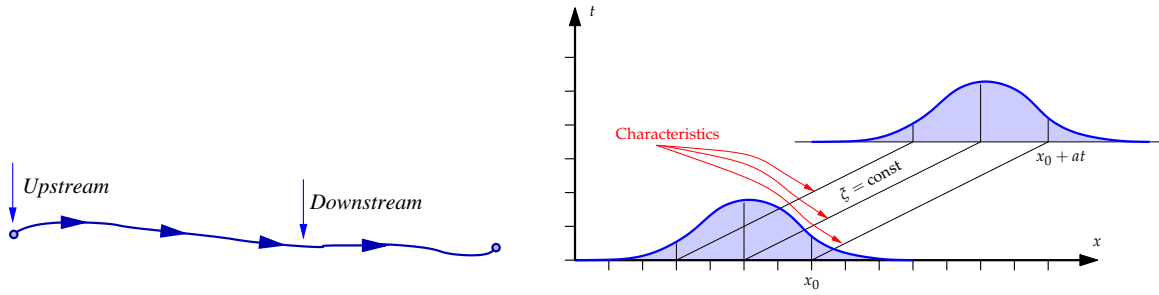
For the second example, we choose the same initial conditions (6.6), but the boundary conditions are  $\Gamma \cdot \nabla n = 0$ , where  $\Gamma$  is the outer normal vector at each point on the whole boundary  $\Gamma$ . This means that the number of dopant is fixed in the whole domain, because no particles enter or leave the simulation domain. The changes in the code for the previous example are

```

1 // changes for fixed number of dopants
2 unknown N[all,all];
3 // all variables are known now
4 // known N[Xmin..Xmax,0];
5
6 // contact side
7 // equ N[i=1..Xmin-1, j=0] -> N[i,j] - N[i,j+1] = 0.0;
8 // equ N[i=Xmax+1..NX-2,j=0] -> N[i,j] - N[i,j+1] = 0.0;
9 // no outflux anywhere
10 equ N[i=1..NX-2,j=0] -> N[i,j] - N[i,j+1] = 0.0;
11
12 begin InitVars
13   assign N[i=5..15,j=1] = 1.0e16; // initial condition
14 end
    
```

source\_code/diffusion\_example2.sg

The numerical result are depicted in Fig. 6.10. Compared to the first example, the dopant concentration smoothes out rather quickly, because there is no more dopant source. Ultimately, the dopant concentration becomes constant in the whole domain, however, this time the equilibrium concentration is not  $n_{\text{surf}}$ , but determined by the number of particles initially available in the simulation domain.



**Figure 6.11:** Boundary conditions have to be specified “upstream”, so that they are initial conditions along a characteristic.

## 6.4 Convective Problems

In the introductory part of this chapter we have already discussed the convective case with  $\Gamma = n\mu E$ . Since the resulting partial differential equation is of first order, we have to specify one boundary condition.

The most prominent type of hyperbolic equations are wave equations, the simplest one being

$$\begin{aligned}\frac{\partial n}{\partial t} + a \frac{\partial n}{\partial x} &= 0, \\ n(x, 0) &= n_0(x).\end{aligned}$$

The solution of this equation is  $n_0(x - at)$ , which is a wave with initial shape  $n_0(x)$  traveling with velocity  $a$ . Those lines in the  $(x, t)$ -plane with constant  $x - at$  are called *characteristics* or *characteristic lines* and are illustrated in Fig. 6.11. The solution along a characteristic is determined by the initial conditions given at a single point on it. If, however, two points on a characteristic are specified by initial conditions, a solution cannot exist anymore, unless the two solutions induced by the two points happen to **coincide**<sup>1</sup>. This can be seen from Fig. 6.11: If we specify an initial condition at  $x_0$  at  $t = 0$ , the value at all points  $x_0 + at$  is determined. Thus, it is meaningless to specify at, say,  $t = 2$  another condition at  $x_0 + 2a$ .

Let us have a closer look at the simple example

$$\frac{\partial n}{\partial t} + a \frac{\partial n}{\partial x} = 0, \quad 0 \leq x \leq 1, t \geq 0$$

If  $a > 0$ , the wave travels from left to right, therefore we specify the boundary condition on the left boundary:  $n(x, 0) = n_0(x)$ ,  $n(0, t) = g(t)$ . The analytical solution for  $a > 0$  is then given as

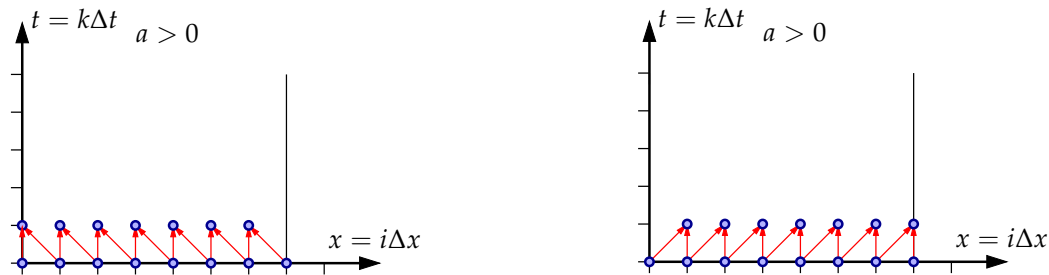
$$n(x, t) = \begin{cases} n_0(x - at), & x - at > 0 \\ g(t - x/a), & x - at < 0 \end{cases}$$

and vice versa for  $a < 0$ .

A numerical discretization forward in time with right-sided differences in space results in

$$\frac{n_{i,k+1} - n_{i,k}}{\Delta t} + a \frac{n_{i+1,k} - n_{i,k}}{\Delta x} = 0,$$

<sup>1</sup> to coincide [kou.m'said]: übereinstimmen



**Figure 6.12:** For a wave travelling from the left to the right, right-sided differences cannot reflect the wave-nature of the solution, while left-sided differences can. The opposite is true for a wave travelling from the right to the left.

leading to

$$n_{i,k+1} = n_{i,k}(1 + a\lambda) - a\lambda n_{i+1,k}, \quad \lambda = \frac{\Delta t}{\Delta x}.$$

We see that the discretization is unable to fully reflect the wave type of the solution, because at point  $i$  no information from the left neighbor  $i - 1$  is used.

Using left-sided differences instead, we obtain

$$n_{i,k+1} = n_{i,k}(1 - a\lambda) + a\lambda n_{i-1,k}, \quad \lambda = \frac{\Delta t}{\Delta x}. \quad (6.7)$$

Now data from the left neighbor is used leading to good results as the next example will show. Nevertheless, this does not mean that right-sided differences are worse than left-sided differences in general! Replacing  $a$  with  $-a$ , right-sided differences perform acceptably, while left-sided differences fail (Fig. 6.12).

For demonstration purposes we will solve

$$\frac{\partial n}{\partial t} + a \frac{\partial n}{\partial x} = 0 \quad (6.8)$$

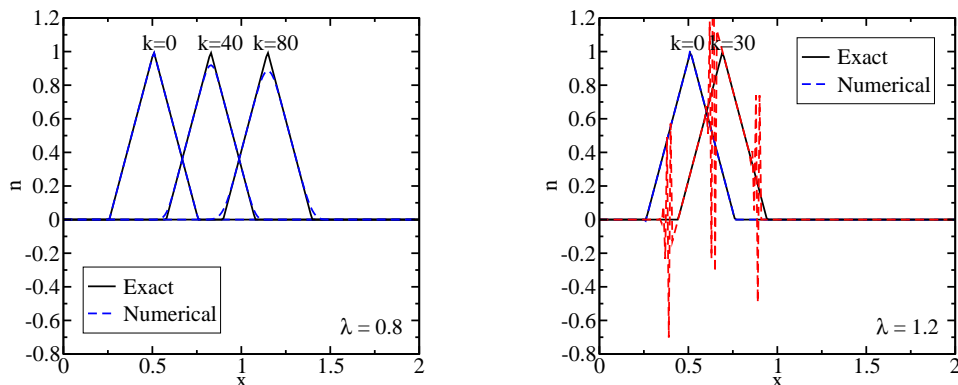
numerically with initial condition

$$n(x,0) = \begin{cases} 1 - |x - at - 1/2|, & |x| < 1 \\ 0, & \text{otherwise} \end{cases}$$

The results can be seen in Fig. 6.13 and illustrate the stability of left-sided differences and the instability of right-sided differences for the given equation.

```

1  const dt = 0.008, dx = 0.01, a = 1;
2  equ V[i=1..Nx] -> (V[i] - Vold[i]) / dt
3      + a * (Vold[i] - Vold[i-1]) / dx = 0;
4  begin calcA
5      assign A[i=0..Nx] = max(0,
6          1 - abs(i*dx - a*time*dt - Nx/4*dx)*k);
7  end
8  begin main
9      ...
10     call calcA;
    
```



**Figure 6.13:** An appropriate difference scheme has to be chosen for hyperbolic equations (here: left-sided for a wave travelling from the left to the right), otherwise instabilities occur.

```

11  assign V[i=all] = A[i];
12  while (...) begin
13      assign Vold[i=all] = V[i];
14      solve;
15      call calcA;
16      assign Vex[i=all] = A[i];
17  end
18  end
    
```

source\_code/hyperbolic\_unstable.sg

We will use the von Neumann analysis to investigate the reason for instability a little bit further. Replacing  $n_{i,k} = (g(\theta))^k e^{-j i \theta}$ ,  $\theta = \Delta x \xi$  in (6.7) yields

$$g(\theta) = 1 + a\lambda(e^{-j\theta} - 1)$$

and with  $e^{-j\theta} = \cos \theta - j \sin \theta$  the **modulus**<sup>1</sup> is

$$|g(\theta)|^2 = (1 + a\lambda(\cos(\theta) - 1))^2 + (a\lambda)^2 \sin^2(\theta) = 1 - 4a\lambda(1 - a\lambda) \sin^2(\theta/2).$$

In order to satisfy  $|g(\theta)| \leq 1$ , the *Courant-Friedrichs-Lewy (CFL) Condition*

$$|a\lambda| < 1 \tag{6.9}$$

is required for stability. This is quite intuitive: The discretized information “travels” at speed  $1/\lambda$  so information at point  $i$  accesses information at point  $i - 1/\lambda$ . Thus, for a proper resolution of the information propagation, the grid size  $a$  should be smaller than  $1/\lambda$  to resolve the flow of information.

Apart from stability issues with “wrong-sided” differences, numerical dispersion occurs even for “correct-sided” differences (cf. left illustration in Fig. 6.13). While the exact solution is a travelling hat, the numerical solution is smoothed out, which can be explained by realizing that different frequencies travel (numerically) at different speeds. This becomes – again – clearer

<sup>1</sup> **modulus** [mɔd.ju:lɔs]: Absolutbetrag

with the von Neumann analysis. Setting  $n_{i,k} = (g(\theta))^k e^{-jm\theta}$  and using (6.2) and  $N(\xi, k+1) = g(\Delta x \xi) N(\xi, k)$  in (6.7) gives

$$g(\Delta x \xi) = (1 - a\lambda) + a\lambda e^{-j\Delta x \xi}.$$

The exact result  $n(x, t) = n_0(x - at)$  has the Fourier transform

$$\begin{aligned} N(\chi, t) = \mathcal{F}\{u(x, t)\} &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-j\chi x} u(x, t) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-j\chi x} u_0(x - at) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-j\chi x'} e^{-j\chi at} u_0(x') dx' \\ &= e^{-j\chi at} N_0(\chi). \end{aligned}$$

Consequently, at a given time  $t$  the Fourier transform at time  $t + \Delta t$  can be computed via

$$N(\chi, t + \Delta t) = e^{-j\chi a \Delta t} N(\chi, t). \quad (6.10)$$

Since the discretized solution is an approximation to the real solution, the amplification factor  $g(\Delta x \xi)$  is an approximation of  $e^{-j\chi a \Delta t}$ . However, the difference scheme yields

$$\begin{aligned} g(\Delta x \xi) &= (1 - a\lambda) + a\lambda e^{-j\Delta x \xi} \\ &= |g(\Delta x \xi)| e^{-j\alpha \Delta t} \end{aligned}$$

for some  $\alpha$  that we are not going to compute explicitly. Now,  $|g(\Delta x \xi)|$  is the damping factor of a wave with frequency  $\xi$ , while  $\alpha(\Delta x \xi)$  is proportional to the phase speed. Comparing with (6.10), the term  $a - \alpha(\Delta x \xi)$  is responsible for dispersion. Fig. 6.14 summarizes the observed effects.

## 6.5 Diffusive and Convective Problems

As mentioned in the introductory part of this chapter, both diffusive and convective contributions are often present. Since a diffusive contribution leads to parabolic PDEs and convective contribution to hyperbolic PDEs, reliable numerical results are harder to obtain in the combined case. This section emphasizes the additional problems and aims at finding a robust discretization scheme for the box integration method.

In the *drift-diffusion model* (1.23) to (1.25), we have both diffusive and convective contributions:

$$J_n = qn\mu_n E + qD_n \nabla n, \quad (6.11)$$

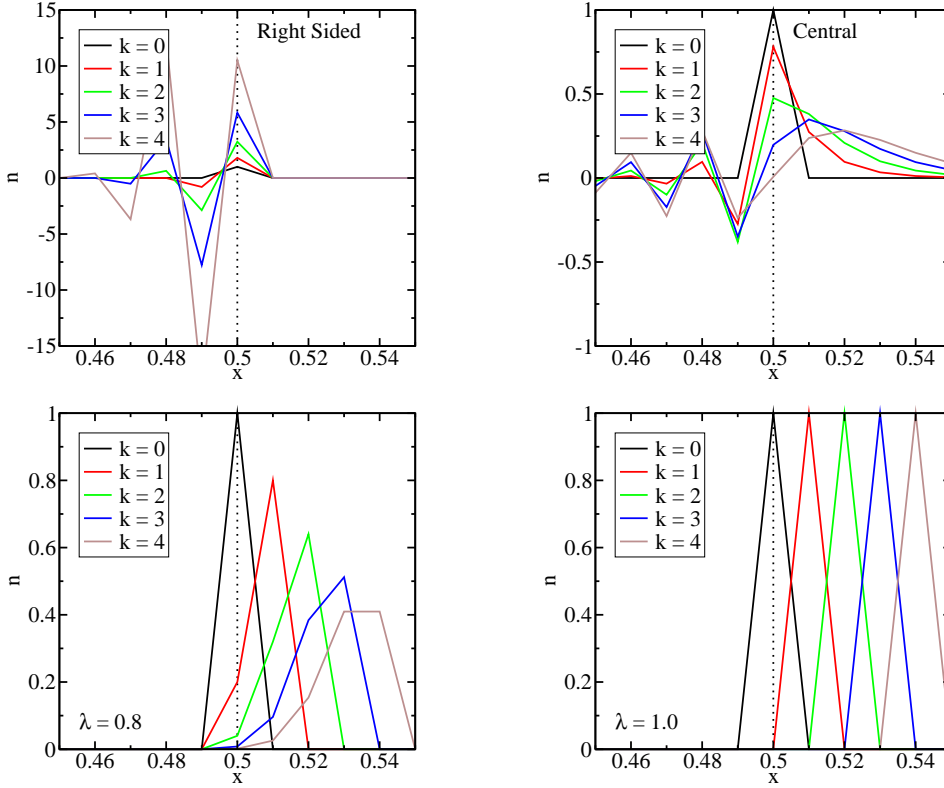
$$J_p = qp\mu_p E - qD_p \nabla p. \quad (6.12)$$

These expressions can now be plugged into the continuity equations

$$\nabla \cdot J_n - q \frac{\partial n}{\partial t} = +qR, \quad (6.13)$$

$$\nabla \cdot J_p + q \frac{\partial p}{\partial t} = -qR. \quad (6.14)$$





**Figure 6.14:** For a wave traveling to the right, neither right-sided differences (upper left) nor central differences (upper right) yield satisfactory results. Even left-sided differences may (lower left) or may not (lower right) show dispersion.

Substituting (6.11) into (6.13) leads to

$$\nabla \cdot (q n \mu_n \mathbf{E} + q D_n \nabla n) - q \frac{\partial n}{\partial t} = +q R.$$

Expanding terms, neglecting recombination, cancelling  $q$  and introducing the *thermal voltage*  $V_T = k_B T / q = D_n / \mu_n$  finally results in

$$\frac{\partial n}{\partial t} - \mu_n \mathbf{E} \cdot \nabla n = \mu_n V_T \Delta n. \quad (6.15)$$

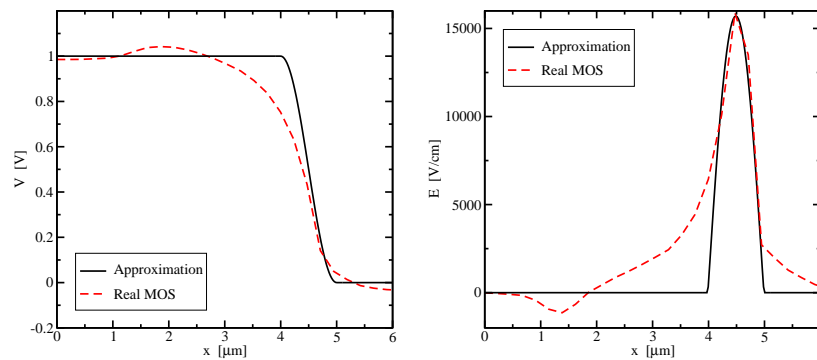
This is a second order parabolic partial differential equation. Depending on  $\mathbf{E}$ , the drift term can dominate in certain device regions so that one can expect a (more or less) dominant propagation in one specific direction with velocity  $v = -\mu \mathbf{E}$ . Diffusion takes place due to the term on the right hand side with coefficient  $D_n = \mu_n V_T$ .

For a box integration discretization, the integral form of (6.13) is needed. Therefore we integrate the equation over a box  $\mathcal{V}$  and obtain

$$\int_{\mathcal{V}} \nabla \cdot \mathbf{J}_n \, dV - q \int_{\mathcal{V}} \frac{\partial n}{\partial t} \, dV = q \int_{\mathcal{V}} R \, dV.$$

Using Gauss' integral theorem on the first integral, we find

$$\int_{\partial \mathcal{V}} \mathbf{J}_n \cdot d\mathbf{A} - q \int_{\mathcal{V}} \frac{\partial n}{\partial t} \, dV = q \int_{\mathcal{V}} R \, dV,$$



**Figure 6.15:** Potential (left) and electric field (right) of an approximate potential barrier (solid lines) designed to mimic the potential distribution inside a MOS transistor.

which can be approximated as

$$\sum_{j \in \mathcal{N}_i} J_{i,j} A_{i,j} - q \int_{\mathcal{V}} \frac{\partial n}{\partial t} dV = q \int_{\mathcal{V}} R dV.$$

A naive discretization for  $J_{i,j}$  recalling  $J_n = qn\mu_n E + qD_n \nabla n = q\mu_n (-n \nabla V + V_T \nabla n)$  would be

$$J_{i,j} = q\mu_n \left( -\frac{n_i + n_j}{2} \frac{V_j - V_i}{d_{i,j}} + V_T \frac{n_j - n_i}{d_{i,j}} \right). \quad (6.16)$$

Recall that in the box discretization the fluxes are required in the middle of a connection  $i, j$  (see Fig. 5.8), while the quantities  $E$  and  $\nabla n$  can be expressed by central differences valid at this mid point. The carrier concentration is only available on the mesh points  $i$  and  $j$ . Naively, it has been assumed that

$$n_{i,j}^{\text{mid}} \approx \frac{n_i + n_j}{2}.$$

The discretization (6.16) is used in the following example: A potential barrier is given (Fig. 6.15) and the resulting electron concentration has to be computed.

```

1  func window(min,max,x) return step(x-min)*step(max-x);
2
3
4  assign x[i=all] = -l2 + i*DX;
5  assign V[i=0..NX] = window(-l2, 0, x[i]) * Vc +
6  window( 0, l1, x[i]) * Vc * sq(cos(x[i]/l1*PI/2));
    
```

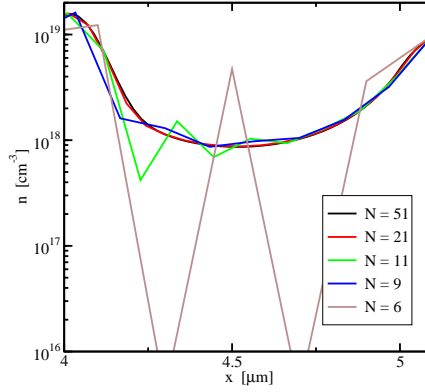
source\_code/instability1.sg

The current continuity equation is implemented using the box integration method as

```

1  func Jsimple(ni,nj,Vi,Vj)
2  return q0 * mu * Vt / DX * (- (ni+nj)/2 * (Vj-Vi)/Vt + (nj-ni));
3
4
5  equ n[i=1..NX-1] -> (Jsimple(n[i],n[i-1],V[i],V[i-1]) +
6  Jsimple(n[i],n[i+1],V[i],V[i+1]))*A = 0;
    
```

source\_code/instability2.sg



**Figure 6.16:** Large grid spacings show instabilities of the numerical solution.

The simulation results shown in Fig. 6.16 are poor: There is a strong dependence on the grid spacing. The scheme does not work for a large grid spacing, making simulations in two or three dimensions **unfeasible**<sup>1</sup>.

We have to find out why the discretization (6.16) performs so badly. Let us therefore consider a one-dimensional example with constant electric field. For a conductor of width  $w$  and height  $h$  we have  $V_i = \Delta x wh$ ,  $A = wh$ . Neglecting recombination, we have to solve

$$\sum_{j \in \mathcal{N}_i} J_{i,j} A = q \frac{\partial n_i}{\partial t} V_i. \quad (6.17)$$

At node  $i$ , the discretization reads

$$(J_{i,i-1} + J_{i,i+1}) A = q \frac{\partial n}{\partial t} V_i \Leftrightarrow (J_{i,i-1} + J_{i,i+1}) = q \frac{\partial n}{\partial t} \Delta x. \quad (6.18)$$

With the naive discretization (6.16) and a constant electric field we obtain

$$J_{i,i\pm 1} = q \mu_n \left( \frac{n_i + n_{i-1}}{2} (\pm E) + V_T \frac{n_{i\pm 1} - n_i}{\Delta x} \right),$$

so all together

$$J_{i,i-1} + J_{i,i+1} = q \mu_n \left( \frac{n_{i+1} - n_{i-1}}{2} E + V_T \frac{n_{i+1} - 2n_i + n_{i-1}}{\Delta x} \right).$$

Inserting into (6.18) with the time-derivative discretized with a forward Euler scheme and dividing by  $\mu_n \Delta x$ , (6.17) becomes

$$\frac{n_{i+1,k} - n_{i-1,k}}{2\Delta x} E + V_T \frac{n_{i+1,k} - 2n_{i,k} + n_{i-1,k}}{(\Delta x)^2} = \frac{n_{i,k+1} - n_{i,k}}{\mu_n \Delta t}. \quad (6.19)$$

Let us further analyze the discretization: With  $a := -\mu_n E$ ,  $b := \mu_n V_T$ ,  $\eta := \Delta t / (\Delta x)^2$  and  $\alpha := a \Delta x / (2b)$ , equation (6.19) can be rearranged to

$$n_{i,k+1} = b\eta(1 + \alpha)n_{i-1} + (1 - 2b\eta)n_i + b\eta(1 - \alpha)n_{i+1}.$$

<sup>1</sup> **unfeasible** [ʌn'fi:zɪ.b]: undurchführbar

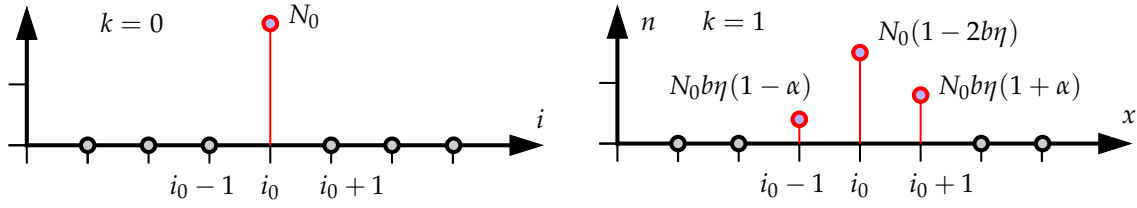


Figure 6.17: Motivation for the derivation of the Reynolds number.

We will assume  $\alpha > 0$  and test with a Dirac pulse  $n_{i,0} = N_0\delta_{i-i_0}$ . The values for the first time step are

$$\begin{aligned} n_{i_0-1,1} &= 0 & + 0 & + N_0b\eta(1 - \alpha), \\ n_{i_0,1} &= 0 & + N_0(1 - 2b\eta) & + 0, \\ n_{i_0+1,1} &= N_0b\eta(1 + \alpha) & + 0 & + 0. \end{aligned}$$

To obtain positive electron concentrations,  $\alpha = a\Delta x/2b < 1$  must hold, or equivalently,  $\Delta x < 2b/a = -2V_T/E$ , with  $E < 0$  because of  $\alpha > 0$ . We have just derived the *Reynolds-* or *Peclet-number*  $|2V_T/E|$  of the cell. Keep in mind that it is *not* a stability criterion in the sense of (6.9), since there is *no time dependence!* Nevertheless, when  $\Delta x$  is larger than the Peclet-number, oscillations in the numerical solution will occur.

Please note that – just as the original PDE – the resulting discretization scheme is not symmetric: There is a net flow of electrons into one direction, driven by the electric field. Thus, for a fixed spatial point  $x$ , the concentration in upstream direction is more important than the concentration in downstream direction. Or to be (much more) pictorial: If you are sitting in a river and do not want to be hit by drifting wood, you have to look carefully for such dangerous wood in the opposite direction of the water flow.

Another problem of the naive discretization (6.16) is that the two contributions from diffusion current and drift current might be of similar magnitude and cancel each other. The main reason for the instability inherent to (6.16) is that, since the carrier concentration depends exponentially on the potential, a discretization of the electron concentration at the midpoint of an edge (i.e. the term  $(n_i + n_j)/2$ ) is definitely a bad one. Historically, this instability prevented efficient application of numerical methods in the analysis of semiconductor devices until the pioneering work of Scharfetter and Gummel [?].

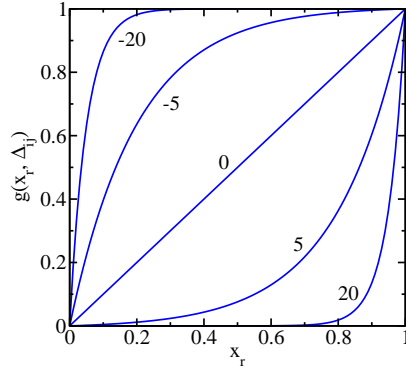
### 6.5.1 Scharfetter-Gummel Discretization

To overcome the issues mentioned above, an alternative discretization was suggested by *Scharfetter* and *Gummel*, which will be summarized in the following. Consider the projection of the current on the edge  $e_{i,j}$ :

$$J_{n,i,j} = e_{i,j} \cdot J_n.$$

With the local coordinate  $x_r = (x - x_i)/d_{i,j}$  along edge  $e_{i,j}$  the current relation (6.11) becomes

$$\frac{J_{n,i,j}}{q\mu_n} = nE_{i,j} + V_T \frac{dn}{dx_r} d_{i,j}, \quad (6.20)$$



**Figure 6.18:** The growth function  $g(x_r, \Delta_{i,j})$  plotted for different values of  $\Delta_{i,j}$ .

where the additional term  $d_{i,j}$  shows up because of the change of variables. We will now assume the three quantities  $J_{n,i,j}$ ,  $E_{i,j} = -\frac{dV}{d\xi}$  and  $\mu$  to be constant along the edge  $e_{i,j}$ , leading to an ordinary differential equation for the unknown function  $n(x_r)$ . The boundary conditions are  $n(0) = n_i$  and  $n(1) = n_j$ . The attentive reader may object that the ODE is of first order only, still we have specified two boundary conditions. This is not **contradictory**<sup>1</sup>, because there is an additional degree of freedom: We have chosen  $J_{n,i,j}$  to be *constant*, but we can still choose its *value*!

The solution of (6.20) is

$$n(x, V) = (1 - \tilde{g}(x, V))n_i + \tilde{g}(x, V)n_j, \quad x_i \leq x \leq x_j, \quad (6.21)$$

with the growth function

$$\tilde{g}(x, V) = g(x_r, \Delta_{i,j}) = \frac{1 - \exp(\Delta_{i,j}x_r)}{1 - \exp(\Delta_{i,j})}, \quad x_r = \frac{x - x_i}{d_{i,j}}, \quad \Delta_{i,j} = \frac{V_j - V_i}{V_T}.$$

For the discretization, we need the values of  $n$  and  $dn/dx_r$  at the midpoint:

$$\frac{J_{n,i,j}}{q\mu_n} = n|_{\text{midpoint}}E_{i,j} + V_T \frac{dn}{dx_r} \Big|_{\text{midpoint}}. \quad (6.22)$$

They can be obtained from the solution of the ODE with  $x_r = 1/2$ :

$$n|_{\text{midpoint}} = \frac{n_i}{1 + \exp(-\Delta_{i,j}/2)} + \frac{n_j}{1 + \exp(\Delta_{i,j}/2)}, \quad (6.23)$$

$$\frac{dn}{dx_r} \Big|_{\text{midpoint}} = \frac{\Delta_{i,j}/2}{\sinh(\Delta_{i,j}/2)} \frac{n_j - n_i}{2}. \quad (6.24)$$

After some formal rearrangement of (6.23) and (6.24), the *Scharfetter-Gummel discretization* for the current is finally obtained as

$$J_{n,i,j} = \frac{q\mu_n V_T}{d_{i,j}} (n_j \mathcal{B}(\Delta_{i,j}) - n_i \mathcal{B}(-\Delta_{i,j})), \quad (6.25)$$

$$J_{p,i,j} = -\frac{q\mu_p V_T}{d_{i,j}} (p_j \mathcal{B}(-\Delta_{i,j}) - p_i \mathcal{B}(\Delta_{i,j})), \quad (6.26)$$

<sup>1</sup> **contradictory** [kɔn.trəˈdɪk.tɔr.i]: widersprüchlich

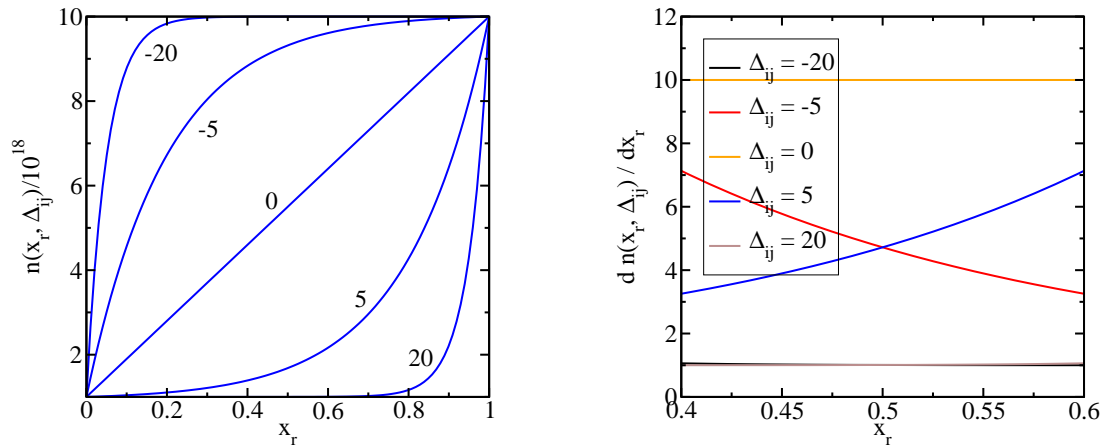


Figure 6.19: The concentration  $n(x_r, \Delta_{ij})$  given by (6.23) and its derivative  $dn/dx_r$  given by (6.24).

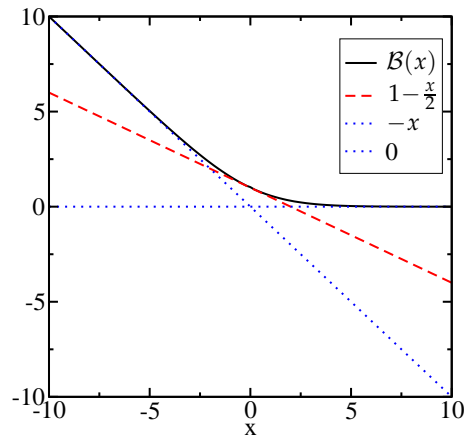


Figure 6.20: The Bernoulli function  $\mathcal{B}(x)$ .

with

$$\Delta_{i,j} = \frac{V_j - V_i}{V_T}, \quad \mathcal{B}(x) = \frac{x}{e^x - 1}.$$

The Bernoulli function  $\mathcal{B}(x)$  is depicted in Fig. 6.20. Its asymptotic behavior is

$$\mathcal{B}(x) \approx \begin{cases} -x, & x \ll 0, \\ 1 - \frac{x}{2}, & x \approx 0, \\ 0, & x \gg 0. \end{cases}$$

How does the Scharfetter-Gummel discretization behave for the large electric fields occurring in the previous example? Let us assume  $E \gg 0$ , then (6.25) used in (6.17) yields

$$J_{i,i-1} + J_{i,i+1} = \frac{q\mu_n V_T}{\Delta x} (n_{i-1} \mathcal{B}(\Delta_{i,i-1}) - n_i \mathcal{B}(-\Delta_{i,i-1}) + n_{i+1} \mathcal{B}(\Delta_{i,i+1}) - n_i \mathcal{B}(-\Delta_{i,i+1})).$$

Now,

$$\begin{aligned}\Delta_{i,i-1} &= \frac{V_{i-1} - V_i}{V_T} = -E \frac{x_{i-1} - x_i}{V_T} = +E \frac{\Delta x}{V_T} \gg 0, \\ \Delta_{i,i+1} &= \frac{V_{i+1} - V_i}{V_T} = -E \frac{x_{i+1} - x_i}{V_T} = -E \frac{\Delta x}{V_T} \ll 0,\end{aligned}$$

so we get

$$J_{i,i-1} + J_{i,i+1} = \frac{q\mu_n V_T}{\Delta x} (-n_i(\Delta_{i,i-1}) + n_{i+1}(-\Delta_{i,i+1})).$$

This finally gives for  $E \gg 0$

$$\begin{aligned}J_{i,i-1} + J_{i,i+1} &= \frac{q\mu_n V_T}{\Delta x} (-n_i \Delta_{i,i-1} - n_{i+1} \Delta_{i,i+1}) \\ &= \frac{q\mu_n V_T}{\Delta x} \left( -n_i E \frac{\Delta x}{V_T} + n_{i+1} E \frac{\Delta x}{V_T} \right) \\ &= q\mu_n E (n_{i+1} - n_i) \\ &\stackrel{!}{=} q \frac{\partial n_i}{\partial t} \Delta x.\end{aligned}$$

With this,

$$\frac{\partial n_i}{\partial t} = \mu_n E \frac{n_{i+1} - n_i}{\Delta x},$$

which is just the right-sided difference with the diffusive term neglected. Analogously, for  $E \ll 0$ , one finds the left-sided difference

$$\frac{\partial n_i}{\partial t} = \mu_n E \frac{n_i - n_{i-1}}{\Delta x}.$$

Comparing these results with (6.19), the diffusion term is gone and the central difference for the drift term is replaced by a right-sided difference and a left-sided difference, respectively.

However, these modifications of the simple discretization are not surprising: Since  $|E| \gg 0$ , the drift term dominates. Recalling the phenomena observed in the previous section about hyperbolic problems, Scharfetter-Gummel adapts itself to the dominant direction of convection and uses right-sided differences for electrons moving from the right to the left and vice versa.

Finally, let us consider the case of small electric fields: Rearranging the simple discretization (6.16) gives

$$\begin{aligned}J_{i,j} &= q\mu_n \left( -\frac{n_i + n_j}{2} \frac{V_j - V_i}{\Delta x} + V_T \frac{n_j - n_i}{\Delta x} \right) \\ &= \frac{q\mu_n V_T}{\Delta x} \left( -\frac{n_i + n_j}{2} \frac{V_j - V_i}{V_T} + n_j - n_i \right) \\ &= \frac{q\mu_n V_T}{\Delta x} \left( -\frac{n_i + n_j}{2} \Delta_{i,j} + n_j - n_i \right) \\ &= \frac{q\mu_n V_T}{\Delta x} \left( n_j \left( 1 - \frac{\Delta_{i,j}}{2} \right) - n_i \left( 1 + \frac{\Delta_{i,j}}{2} \right) \right)\end{aligned}$$

Comparison with (6.25) shows that the simple discretization can be obtained from the Scharfetter-Gummel discretization by linearizing the Bernoulli function around  $E = 0$ . This implies that the naive discretization is valid for **negligible**<sup>1</sup> fields only. Only under these circumstances can the carrier concentration at the midpoint be expressed by their arithmetic average, see also Fig. 6.19.

<sup>1</sup> **negligible** [neg.l.dʒə.b]: vernachlässigbar

# Chapter 7

## Parameter Modeling

The basic semiconductor equations discussed in Chapter 1 contain several physical parameters (in particular: the mobilities  $\mu_n$  and  $\mu_p$  and the recombination term  $R$ ) which have to be accurately modeled for the purpose of reliable device simulation. For accurate results, these parameters cannot be assumed to be constant because they do not reflect the underlying physics to a sufficient degree of accuracy. Unfortunately this leads to an increased computational effort, so a good compromise between (physical) model accuracy and computational effort has to be found.

In many cases, sufficiently accurate and simple analytic models derived from theoretical considerations are not available and parameters such as the mobility have to be extracted from measurement results. One can then fit an analytic curve through the available data points and use this empirical model for the simulation.

### 7.1 Carrier Mobilities

Carrier mobilities in semiconductors are influenced by a variety of physical mechanisms. We will focus on the following dominant effects:

- Scattering at lattice atoms or defects:  $\mu^L$
- Scattering at charges or neutral impurities:  $\mu^I$
- Surface **roughness**<sup>1</sup> scattering:  $\mu^S$
- Increased scattering due to heating:  $\mu^F$

In order to obtain a **tractable**<sup>2</sup> model, scattering effects are frequently assumed to be independent. A common assumption is that the effective mobility is given by

$$\mu^{\text{LISF}} = \mu^{\text{LISF}}(\mu^{\text{LIS}}(\mu^{\text{LI}}(\mu^{\text{L}}))) ,$$

which means that one first computes a mobility due to scattering at lattice atoms  $\mu^L$  from an initially constant mobility  $\mu^0$ . Then, a modified mobility from  $\mu^L$  by taking a possible modification due to scattering at charges or neutral impurities into account. After that, surface roughness

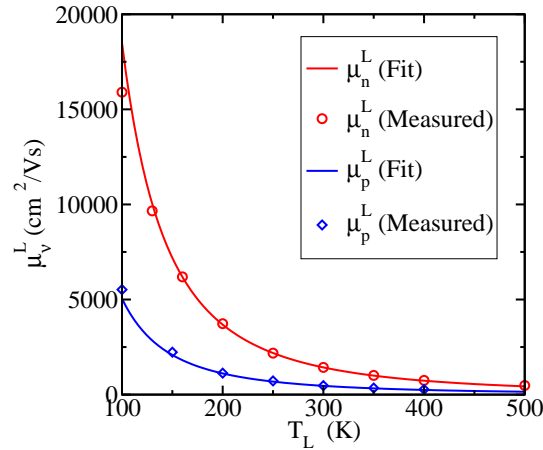
---

<sup>1</sup> **roughness** [rʌf.nəs]: Rauheit, Unebenheit    <sup>2</sup> **tractable** [træk.tə.bɪl]: handhabbar, lenkbar



Material	$\mu_n^0$ (cm <sup>2</sup> /Vs)	$\alpha_n$ (1)	$\mu_p^0$ (cm <sup>2</sup> /Vs)	$\alpha_p$ (1)
Si	1430	2.33	460	2.18
Ge	3800	1.66	1800	2.33
GaAs	8500	1	400	2.1

**Table 7.1:** Lattice scattering parameters for several semiconductors. Note that electron mobilities are higher than hole mobilities, therefore an *n*-MOS typically has a higher transconductance than a *p*-MOS (see Chapter 8).



**Figure 7.1:** Comparison of measured and fitted mobilities due to lattice scattering in a homogeneous semiconductor.

scattering is considered to find the mobility  $\mu^{\text{LIS}}$  including scattering at lattice atoms or charges and surface scattering into account. Finally, increased scattering due to heating yields a mobility  $\mu^{\text{LISF}}$  that includes the physical effects mentioned above.

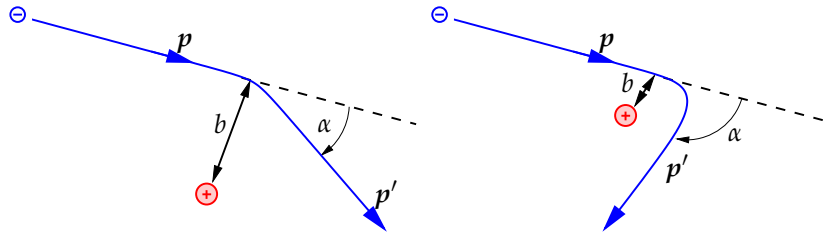
### 7.1.1 Lattice Scattering

At non-zero temperature the lattice atoms oscillate about their equilibrium sites. Even in pure, perfectly crystalline semiconductors, carriers are scattered because of their interaction with the vibrating lattice. This results in the lattice mobility  $\mu^L$  which is a function of the lattice temperature. Theoretical models describing lattice scattering are complicated and results often unsatisfactory. However, for the purpose of simulation an empirical power law is conventionally assumed:

$$\mu_v^L = \mu_v^0 \cdot \left( \frac{T}{300 \text{ K}} \right)^{-\alpha_v}, \quad v = n, p$$

This power law shows good agreement with experimental data (Fig. 7.1) in the temperature range between 200 K and 500 K (cf. Tab. 7.1), the typical operation temperature of semiconductor devices.

When comparing experimental data given in the literature, one finds that such data can show considerable scatter. For example, the parameters for the electron mobility are often in the (pretty



**Figure 7.2:** A carrier trajectory is influenced by ionized impurities. The potential (attraction or repulsion) “felt” by electrons is inversely proportional to the distance to the impurity, so trajectories far away from the impurity are less influenced (left) compared to those closer to the impurity (right).

large intervals)  $1240 \text{ cm}^2/\text{Vs} < \mu_n^0 < 1600 \text{ cm}^2/\text{Vs}$  and  $2.2 < \alpha_n < 2.6$  for Si. One reason is that each author uses their own measurement range, so that the available data puts some implicit weight on the measured interval. Furthermore, measurements and device fabrications are (at least to some extent) stochastic processes, resulting in many different sets of measurements.

### 7.1.2 Ionized Impurity Scattering

In semiconductor devices mobility reduction due to scattering by charged impurities (Fig. 7.2) is a dominant effect. The impact of lattice and impurity scattering has to be combined in some way to obtain an effective mobility. From a theoretical point of view, ionized impurity scattering is very difficult to model, because there is a strong dependence on doping concentration and mobile charges  $n$  and  $p$  that **screen**<sup>1</sup> the impurities.

An empirical model for the combined lattice and ionized impurity mobility was introduced by Caughey and Thomas [?]. To fit experimental data they used the following (empirical) expression:

$$\mu_v^{LI} = \mu_v^{\min} + \frac{\mu_v^L - \mu_v^{\min}}{1 + \left(\frac{N_I}{N_v^{\text{ref}}}\right)^{\alpha_v}}, \quad v = n, p, \quad (7.1)$$

where

$$N_I = \sum_i |Z_i| N_i$$

is a sum over all charged impurities,  $Z_i$  is the charge state of the impurity. For (singly ionized) impurities (such as **boron**<sup>2</sup>, **arsenic**<sup>3</sup> and **phosphorus**<sup>4</sup> in Si),  $|Z_i| = 1$ . The above expression has three free parameters that need to be calibrated to experimental data. Typical values for silicon at room temperature are  $\mu_n^{\min} = 80 \text{ cm}^2/\text{Vs}$ ,  $N_n^{\text{ref}} = 1.12 \times 10^{17} \text{ cm}^{-3}$ ,  $\alpha_n = 0.72$  for electrons and  $\mu_p^{\min} = 45 \text{ cm}^2/\text{Vs}$ ,  $N_p^{\text{ref}} = 2.23 \times 10^{17} \text{ cm}^{-3}$ ,  $\alpha_p = 0.72$  for holes.

### 7.1.3 Surface/Interface Scattering

The perfect periodicity of a semiconductor crystal is due to finite spatial dimensions broken by crystal surfaces. On interfaces between different materials, different lattice constants lead

<sup>1</sup> to screen [skri:n]: abschirmen, schützen, filtern    <sup>2</sup> boron ['bɔ:ɹən]: Bor    <sup>3</sup> arsenic ['ɑ:seɪn.ɪk]: Arsen    <sup>4</sup> phosphorus [fɔs.fɹ.əs]: Phosphor

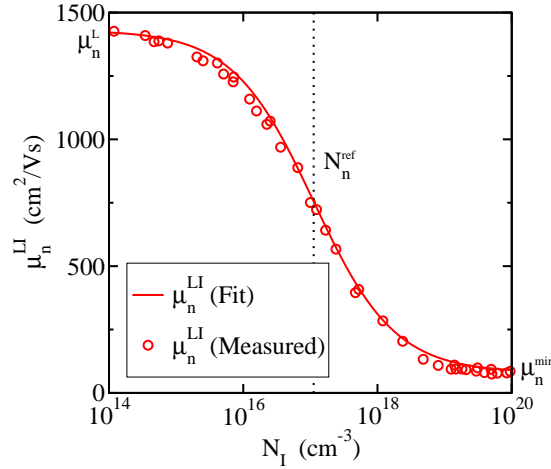


Figure 7.3: Comparison of measured and fitted mobilities  $\mu_n^{LI}$ .

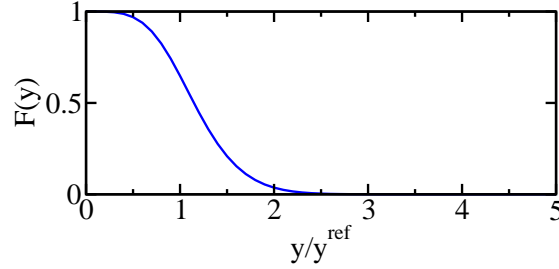


Figure 7.4: Depth dependence function  $F(y)$  for modeling surface scattering in (7.2).

to unavoidable imperfections. This kind of imperfection is a severe concern if current flows primarily along the interface, as it is the case in a MOSFET. Typically, the mobility along a surface is considerably smaller than in the center of a crystal. There is no sharp transition from the high mobility in the crystal to the low mobility on the surface, instead, a smooth transition occurs.

An empirical expression for the description of such a depth dependence of the mobility suggested by Selberherr [?] is

$$\mu_v^{LIS} = \frac{\mu_v^{\text{ref}} + (\mu_v^{LI} - \mu_v^{\text{ref}})(1 - F(y))}{1 + F(y) \left( \frac{S_v}{S_v^{\text{ref}}} \right)^{\gamma_v}}, \quad v = n, p. \quad (7.2)$$

The depth dependence  $F(y)$  is given by

$$F(y) = \frac{2 \exp\left(- (y/y^{\text{ref}})^2\right)}{1 + \exp\left(-2 (y/y^{\text{ref}})^2\right)},$$

where the depth parameter  $y^{\text{ref}}$  is typically 2...10 nm. The *pressing forces*  $S_n$  and  $S_p$  are equal to the magnitude of the normal field strength at the interface ( $E_{\perp} = \mathbf{E} \cdot \mathbf{n}$ ,  $\mathbf{n}$  is the surface normal vector), if the carriers are attracted by it, otherwise they are zero. The remaining parameters are fit-parameters.

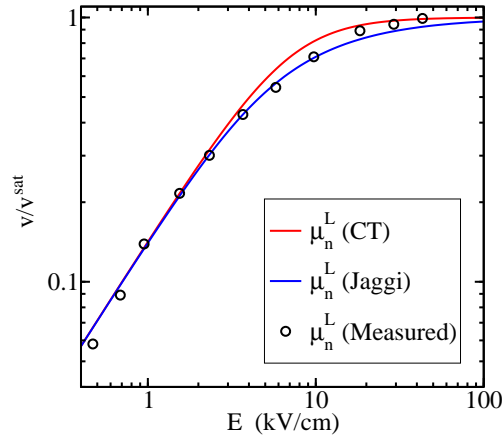


Figure 7.5: Comparison of carrier heating models. (CT = Caughey-Thomas)

### 7.1.4 Carrier Heating

The carrier energy consists of two contributions: one due to random thermal motion, which corresponds to the temperature, and a kinetic contribution given by the kinetic energy  $mv^2/2$ . The average energy  $w$  per particle is thus given by

$$w = \frac{3k_B T_n}{2} + \frac{mv^2}{2}, \quad (7.3)$$

where  $T_n$  is the carrier temperature. Application of a field first increases the kinetic energy, while scattering eventually transfers the kinetic contribution to thermal energy by increasing the carrier temperature  $T_n$ .

At low electric fields the mobility is constant with respect to the field and the velocity-field relationship is consequently linear. The thermal velocity of electrons and holes is large compared to the movement due to the externally applied field and the lattice temperature matches the carrier temperature. At high fields, however, scattering events occur more frequently, converting a considerable amount of kinetic energy to thermal energy. As a result, the velocity-field relationship becomes non-linear (cf. Fig. 7.5). This effect is modeled in the framework of the drift-diffusion model by a field-dependent mobility.

Usually empirical mobility expressions are chosen where parameters are determined by fitting experimental data. One widely used expression is due to Caughey and Thomas [?]:

$$\mu_v^{\text{LISF}}(E) = \frac{\mu_v^{\text{LIS}}}{\left(1 + \left(\frac{\mu_v^{\text{LIS}} E}{v_v^{\text{sat}}}\right)^{\beta_v}\right)^{1/\beta_v}}, \quad v = n, p \quad (7.4)$$

and another one is due to Jaggi [?, ?]:

$$\mu_v^{\text{LISF}}(E) = \frac{2\mu_v^{\text{LIS}}}{1 + \left(1 + \left(\frac{2\mu_v^{\text{LIS}} E}{v_v^{\text{sat}}}\right)^{\beta_v}\right)^{1/\beta_v}}, \quad v = n, p. \quad (7.5)$$

Both expressions contain as parameters the low-field mobility  $\mu_v^{\text{LIS}}$  and the saturation velocity  $v_v^{\text{sat}}$ . These parameters determine the low-field and high-field limits of the  $v(E)$  curve:

$$\lim_{E \rightarrow 0} \mu_v^{\text{LISF}}(E) = \mu_v^{\text{LIS}}, \quad \lim_{E \rightarrow \infty} v_v(E) = \lim_{E \rightarrow \infty} \mu_v^{\text{LISF}}(E) \times E = v_v^{\text{sat}}.$$

The following parameters can be used for silicon at room temperature:  $v_n^{\text{sat}} = 1 \times 10^7$  cm/s,  $\beta_n = 2$ ,  $v_p^{\text{sat}} = 8 \times 10^6$  cm/s,  $\beta_p = 1$ . Note that for high electric fields, both (7.4) and (7.5) asymptotically reach  $\mu_v^{\text{LISF}}(E) \sim 1/E$ , as expected.

## 7.2 Carrier Generation and Recombination

Generation-recombination phenomena determine many essential effects such as leakage current and device breakdown. For a semiconductor in thermal equilibrium there is a dynamic balance between generation and recombination processes, which leads to an equilibrium concentration  $n_0$  for electrons and  $p_0$  for holes:

$$n_0 = N_c \exp\left(\frac{E_F - E_c}{k_B T}\right) = n_i \exp\left(\frac{E_F - E_i}{k_B T}\right), \quad (7.6)$$

$$p_0 = N_v \exp\left(\frac{E_v - E_F}{k_B T}\right) = n_i \exp\left(\frac{E_i - E_F}{k_B T}\right). \quad (7.7)$$

For the product we get

$$n_0 p_0 = N_c N_v \exp\left(-\frac{E_c - E_v}{k_B T}\right) = N_c N_v \exp\left(-\frac{E_g}{k_B T}\right) = n_i^2, \quad (7.8)$$

with the intrinsic concentration

$$n_i = \sqrt{N_c N_v} \exp\left(-\frac{E_g}{2k_B T}\right). \quad (7.9)$$

Note that these equations are valid for non-degenerate semiconductors only, because we assumed Boltzmann statistics in (7.6) and (7.7).

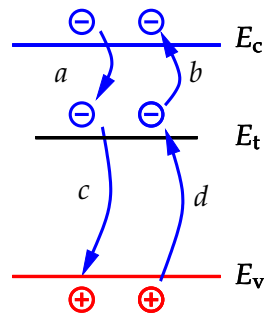
In the case of a deviation of the electron and hole concentrations from their equilibrium values the balance of generation and recombination rates is disturbed. In regions with **excess**<sup>1</sup> carriers ( $np > n_i^2$ ), recombination will **prevail**<sup>2</sup>, whereas in the opposite case ( $np < n_i^2$ ) generation will dominate.

Generation-recombination can be induced by various physical mechanisms, such as absorption or emission of a photon, absorption or emission of a phonon (i.e. the quantum of lattice vibrations), three-particle transitions, and transitions assisted by recombination centers. The relative importance of these mechanisms depends on material properties and operating conditions.

Transition from the valence band to the conduction band and vice versa requires energy. If an electron is **elevated**<sup>3</sup> from the valence band to the conduction band (which corresponds to a hole going from the conduction band to the valence band), its energy increase approximately equals the band gap energy  $E_g$ . This amount of energy can be supplied by several sources:

- *Phonons*: Due to thermal lattice vibration, which is quantified by means of phonons, electrons can gain energy.

<sup>1</sup> **excess** [ek'ses]: überschüssig    <sup>2</sup> **to prevail** [pri'veil]: überwiegen, vorherrschen    <sup>3</sup> **to elevate** [el.i.veit]: emporheben, erhöhen



**Figure 7.6:** The complex process of phonon transition can be split up into four partial processes.

- *Photons:* When exposed to light, photons hit the semiconductor, each carrying an energy  $\hbar\omega$ . If this energy is larger or equal to the band gap, photons can be absorbed by electrons to pass the band gap.
- *Collisions:* The random motion of electrons inevitably leads to collisions among electrons. An electron from the conduction band with high energy (i.e. much higher than the conduction edge) can transfer enough energy to its collision partner in the valence band so that both electrons can finally be found in the conduction band.

The first two possibilities can also act as energy sinks: The first case usually leads to the generation of heat, while the second occurs in direct bandgap materials only and is exploited in lasers and LEDs. In the following we will look at each of these energy supplies separately.

### 7.2.1 Phonon Assisted Recombination and Generation

In indirect band gap semiconductors, such as silicon and germanium, it was found experimentally that generation-recombination phenomena occur primarily via trap centers. This mechanism is commonly termed Shockley-Read-Hall generation-recombination after the authors who established the theory [?, ?]. Indirect generation-recombination is a non-radiative process. In detail, four partial processes are involved (Fig. 7.6):

- Electron capture: An electron from the conduction band is trapped by an unoccupied defect which becomes occupied.
- Electron emission: An electron from an occupied trap moves to the conduction band. The trap becomes unoccupied.
- Hole capture: An electron from an occupied trap moves to the valence band and neutralizes a hole. The trap becomes unoccupied.
- Hole emission: An electron from the valence band is trapped by a defect, thus leaving a hole in the valence band and an occupied trap.

With  $k_a, k_b, k_c, k_d$  the corresponding reaction rates are given by

$$\begin{aligned} v_a &= k_a n N_t^0 && \text{(electron capture) ,} \\ v_b &= k_b N_t^- && \text{(electron emission) ,} \\ v_c &= k_c p N_t^- && \text{(hole capture) ,} \\ v_d &= k_d N_t^0 && \text{(hole emission) ,} \end{aligned}$$

where  $N_t^0$  is the concentration of neutral traps and  $N_t^-$  is the concentration of occupied traps. The total trap concentration  $N_t$  is given by  $N_t = N_t^0 + N_t^-$ . The fraction of occupied traps is given by  $f_t = N_t^- / N_t$ ,  $1 - f_t = N_t^0 / N_t$ . The rate equation for  $v_a$  for electron capture states that the transmission rate is proportional to the number of carriers in the conduction band  $n$  and the number of neutral (free) traps  $N_t^0$ . For electron emission, the reaction rate  $v_b$  is proportional to the number of electrons  $N_t^-$  in the traps only, because most states in the conduction band are free (i.e. the distribution function  $f$  is close to zero, hence  $1 - f$  is close to 1). A similar reasoning applies to holes.

In thermal equilibrium the principle of detailed balance holds, which implies in this particular case

$$v_a^{\text{eq}} = v_b^{\text{eq}} , \quad v_c^{\text{eq}} = v_d^{\text{eq}} .$$

Thus, the four rate constants are not independent:

$$k_b = k_a n_0 \underbrace{\frac{1 - f_{t,0}}{f_{t,0}}}_{=:n_1} , \quad (7.10)$$

$$k_d = k_c p_0 \underbrace{\frac{f_{t,0}}{1 - f_{t,0}}}_{=:p_1} , \quad (7.11)$$

with the auxiliary concentrations  $n_1$  and  $p_1$ . Here,  $f_{t,0}$  denotes the fraction of occupied traps in thermal equilibrium. The meaning of the concentrations  $n_1$  and  $p_1$  will be discussed below.

With these definitions the net recombination rates become

$$\begin{aligned} R_n^{\text{SRH}} &= v_a - v_b = k_a N_t (n(1 - f_t) - n_1 f_t) \\ R_p^{\text{SRH}} &= v_c - v_d = k_c N_t (p f_t - p_1 (1 - f_t)) . \end{aligned}$$

In the general dynamic case,  $R_n^{\text{SRH}}$  is not necessarily equal to  $R_p^{\text{SRH}}$ . This is modeled by adding a conservation equation to the basic semiconductor equations:

$$\frac{\partial N_t^-}{\partial t} = R_n^{\text{SRH}} - R_p^{\text{SRH}} \quad (7.12)$$

that has to be solved in the whole domain. Under steady-state conditions, the time derivative vanishes and the net recombination rates of electrons and holes are equal. This allows the trap occupancy function to be calculated explicitly:

$$f_t = \frac{k_a n + k_c p_1}{k_a (n + n_1) + k_c (p + p_1)} .$$

The net recombination rate after Shockley, Read [?] and Hall [?] thus evaluates to

$$\begin{aligned} R^{\text{SRH}} &= N_t k_a k_c \frac{np - n_1 p_1}{k_a(n + n_1) + k_c(p + p_1)} \\ &= N_t k_a k_c \frac{np - n_i^2}{k_a(n + n_1) + k_c(p + p_1)}. \end{aligned}$$

Introducing the carrier lifetimes  $\tau_n^{-1} = k_a N_t$  and  $\tau_p^{-1} = k_c N_t$  one finally obtains

$$R^{\text{SRH}} = \frac{np - n_i^2}{\tau_p(n + n_1) + \tau_n(p + p_1)}. \quad (7.13)$$

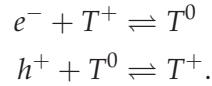
In case  $np > n_i^2$  we have  $R > 0$ , which means that recombination takes place until  $np = n_i^2$ . On the contrary, if  $np < n_i^2$ , generation is dominant, such that carrier concentrations increase until again  $np = n_i^2$ . The carrier lifetimes  $\tau_n$  and  $\tau_p$  control the transient response of the material in non-equilibrium situations. The shorter the carrier lifetimes are, the larger is  $|R|$ , thus the material approaches its equilibrium more rapidly.

Traps are defects with an energy level  $E_t$  and a concentration  $N_t$ . The characteristic parameters describing the interaction of carriers and trap centers are the capture cross sections  $\sigma_n$  and  $\sigma_p$  for electrons and holes, respectively and the rate constants and the carrier lifetimes are conventionally expressed as

$$k_a = \sigma_n v_{\text{th}}^n, \quad \tau_b^{-1} = \sigma_n v_{\text{th}}^n N_t, \quad k_c = \sigma_p v_{\text{th}}^p, \quad \tau_d^{-1} = \sigma_p v_{\text{th}}^p N_t,$$

with the thermal velocities  $v_{\text{th}}^n$  and  $v_{\text{th}}^p$  given by (2.14) for electrons and holes respectively.

So far, the model was derived for *acceptor-like* traps. These can exist in a neutral or a negatively charged state ( $N_t^0, N_t^-$ ). *Donor-like* traps have a neutral and a positively charged state ( $N_t^0, N_t^+$ ). For the latter type of traps the reaction rates are of the form



Following an analogous derivation as the one leading to (7.13) it can be shown that the net recombination rate for donor-like traps is the same as for acceptor-like traps.

We still have to interpret the auxiliary concentrations  $n_1$  and  $p_1$ . In equilibrium, the trap occupancy function  $f_{t,0}$  is determined by the Fermi-Dirac distribution function:

$$f_{t,0} = \left( 1 + \exp\left(\frac{E_t - E_F}{k_B T}\right) \right)^{-1}.$$

It has to be emphasized that one must not use the Boltzmann distribution at this point, because the traps are located in the band gap, where the Boltzmann distribution is not valid anymore. We find

$$\frac{1 - f_{t,0}}{f_{t,0}} = \exp\left(\frac{E_t - E_F}{k_B T}\right).$$



With this we find that for the equilibrium carrier concentrations (cf. (7.6) and (7.7)), the auxiliary concentrations  $n_1$  and  $p_1$  introduced in (7.10) and (7.11) can be written as

$$n_1 = n_0 \frac{1 - f_{t,0}}{f_{t,0}} = n_i \exp\left(\frac{E_t - E_i}{k_B T}\right), \quad (7.14)$$

$$p_1 = p_0 \frac{f_{t,0}}{1 - f_{t,0}} = n_i \exp\left(\frac{E_i - E_t}{k_B T}\right). \quad (7.15)$$

Therefore,  $n_1$  and  $p_1$  are those carrier concentrations that correspond to a Fermi level equal to the trap level ( $E_F = E_t$ ).

A trap acts most effectively as recombination center if its energy level  $E_t$  lies in the middle of the band gap ( $E_t \approx E_i$ ), because in this case both  $n_1$  and  $p_1$  in (7.14) and (7.15) take moderate values and the recombination rate in (7.13) is large. These traps are termed *deep traps*. If the dependence of the helper concentrations  $n_1$  and  $p_1$  on the energy difference in (7.14) and (7.15) were linear, the trap location would be irrelevant, but since there is an exponential dependence, the convex function  $n_1 + p_1$  has a unique minimum at  $E_t = E_i$ .

Dopants are impurities with energy levels close to the band edges and are termed *shallow traps*. For shallow traps  $|E_t - E_i| \gg k_B T$ . Thus either  $n_1$  or  $p_1$  takes a large value because of the exponential dependence, which results in a negligible net recombination rate according to (7.13). Therefore, dopants are very ineffective recombination centers.

Let us consider the special case of the space charge region of a *pn*-diode under reverse bias. We can make the assumptions  $n \ll n_i$  and  $p \ll n_i$  as well as  $\tau_n = \tau_p = \tau$ . The *Shockley-Read-Hall recombination* is then computed as

$$\begin{aligned} R^{\text{SRH}} &= \frac{np - n_i^2}{\tau_p(n + n_1) + \tau_n(p + p_1)} \\ &\approx -\frac{n_i^2}{\tau_p n_1 + \tau_n p_1} \\ &= -\frac{n_i}{\tau} \frac{1}{\exp\left(\frac{E_t - E_i}{k_B T}\right) + \exp\left(\frac{E_i - E_t}{k_B T}\right)} \\ &= -\frac{n_i}{\tau} \frac{1}{2 \cosh\left(\frac{E_t - E_i}{k_B T}\right)}. \end{aligned}$$

Since  $\cosh(x) \geq 1$  for all  $x$ , the maximum recombination rate is

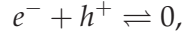
$$R^{\text{SRH}} \Big|_{\text{max}} = -\frac{n_i}{2\tau}$$

for  $E_t = E_i$ . In general, the carrier lifetimes differ ( $\tau_n \neq \tau_p$ ) and this energy is slightly shifted but still close to  $E_i$  under practical circumstances.

## 7.2.2 Photon Transition

The mechanism of direct generation/recombination can also be accompanied by photon emission or absorption. Due to the small photon momentum direct band-to-band transitions are only important in direct gap semiconductors such as GaAs. In silicon and germanium the direct generation-recombination mechanism is insignificant.

According to the reaction



two partial processes are involved:

- a) Electron-hole recombination: An electron moves from the conduction band to the valence band where it neutralizes a hole. A photon with energy of approximately the band gap energy is emitted (radiative recombination).
- b) Electron-hole pair generation: An electron from the valence band absorbs a photon, whose energy is larger than the band gap energy, and moves to the conduction band. A hole is left behind in the valence band (optical generation).

With the rate constants  $k_a^{\text{opt}}$  and  $k_b^{\text{opt}}$  the rate equations become (note that in case **a**) both an electron and a hole are required)

$$\begin{aligned} v_a &= k_a^{\text{opt}}(T)np, \\ v_b &= k_b^{\text{opt}}(T). \end{aligned}$$

These two rates must be equal in thermal equilibrium:

$$v_{a,0} = v_{b,0} \Rightarrow k_a^{\text{opt}} n_i^2 = k_b^{\text{opt}}.$$

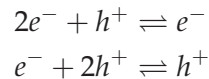
Therefore, the net recombination rate evaluates to

$$R^{\text{opt}} = v_{a,0} - v_{b,0} = k_a^{\text{opt}}(np - n_i^2).$$

Note, again, the term  $np - n_i^2$ , which expresses the tendency to finally reach an equilibrium, where  $np = n_i^2$  holds.

### 7.2.3 Auger Generation-Recombination

In this process three particles are involved, but only two move from one band to another. The third one, which provides or receives the excess energy, moves to another energy level within the same band, where it ultimately loses (in case of recombination) its energy to thermal vibrations. We consider the direct band-to-band Auger process which is sometimes also referred to as *phonon-assisted* Auger process. Four partial processes are involved:



- a) Electron capture: An electron from the conduction band moves to the valence band. The excess energy is transmitted to another conduction electron. In the valence band the electron neutralizes a hole.
- b) Electron emission: A valence electron moves to the conduction band by consuming the energy of a high energetic electron in the conduction band. A hole is left behind in the valence band.

- c) Hole capture: An electron from the conduction band moves to the valence band. The excess energy is transmitted to another hole. In the valence band the electron neutralizes a hole.
- d) Hole emission: A valence electron moves to the conduction band by consuming the energy of a high energetic hole in the valence band. A hole is left behind in the valence band.

With the rate constants  $c_n$ ,  $e_n$ ,  $c_p$  and  $e_p$  the reaction rates read

$$v_a = c_n n^2 p \quad (\text{electron capture}) , \quad (7.16)$$

$$v_b = e_n n \quad (\text{electron emission}) , \quad (7.17)$$

$$v_c = c_p n p^2 \quad (\text{hole capture}) , \quad (7.18)$$

$$v_d = e_p p \quad (\text{hole emission}) . \quad (7.19)$$

The term  $n^2 p$  in (7.16) is due to the need for two electrons and one hole from the conduction and valence band respectively. Even though there are two electrons involved in the electron emission process in (7.17), only one electron from the conduction band (recall that  $n$  represents the number of electrons in the conduction band only) participates. A similar reasoning applies to (7.18) and (7.19).

In thermal equilibrium the principle of detailed balance holds, which implies

$$v_{a,0} = v_{b,0} \Rightarrow c_n n_i^2 = e_n ,$$

$$v_{c,0} = v_{d,0} \Rightarrow c_p n_i^2 = e_p .$$

The remaining rate constants  $c_n$  and  $c_p$  are referred to as *Auger coefficients*. The net recombination rate for the Auger process becomes

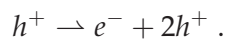
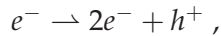
$$R^{\text{AU}} = v_a - v_b + v_c - v_d = (c_n n + c_p p)(np - n_i^2) .$$

Note again the term  $np - n_i^2$ , which expresses the tendency to finally reach an equilibrium. The Auger coefficients vary slightly with temperature between 77 K and 400 K. Typical values for silicon at room temperature are:  $c_n = 2.9 \times 10^{-31} \text{ cm}^6/\text{s}$  and  $c_p = 9.9 \times 10^{-32} \text{ cm}^6/\text{s}$ .

## 7.2.4 Impact Ionization

Impact ionization is a pure generation process in which a high energetic carrier generates an electron-hole pair. From a microscopic point of view there is no difference between impact ionization and Auger generation. The only difference lies in the energy sources of the two processes: the Auger generation rate was evaluated by making use of the principle of detailed balance, which holds in equilibrium. Impact ionization, however, is a typical non-equilibrium process requiring large electric fields.

Two partial processes have to be considered:



- a) Electron emission: A valence electron moves to the conduction band by consuming the energy of a high energetic electron in the conduction band. A hole is left behind in the valence band.

- b) Hole emission: A valence electron moves to the conduction band by consuming the energy of a high energetic hole in the valence band. A hole is left behind in the valence band.

Obviously, these two partial processes are identical to the Auger partial processes b) and d). However, for impact ionization the reaction rates are modeled differently in the framework of the drift-diffusion model:

$$v_a = \alpha_n \frac{|J_n|}{q},$$

$$v_b = \alpha_p \frac{|J_p|}{q}.$$

Here,  $\alpha_n$  and  $\alpha_p$  are the ionization coefficients for electrons and holes, respectively. They are defined as the reciprocals of the average distances traveled by the carriers between consecutive ionization events. For instance, over a distance  $1/\alpha_n$  an electron generates on average one electron-hole pair. The total generation rate can thus be written as

$$G^{\text{II}} = v_a + v_b = \alpha_n \frac{|J_n|}{q} + \alpha_p \frac{|J_p|}{q}. \quad (7.20)$$

Generation due to impact ionization is proportional to the current densities, while Auger generation is proportional to the carrier concentrations (cf. (7.17) and (7.19)). This means that Auger generation may take place in regions with a high concentration of mobile carriers with negligible current flow, whereas impact ionization requires non-negligible current flow as **prerequisite**<sup>1</sup>.

Both theoretical and experimental investigations indicate an exponential dependence of the ionization coefficients on the electric field:

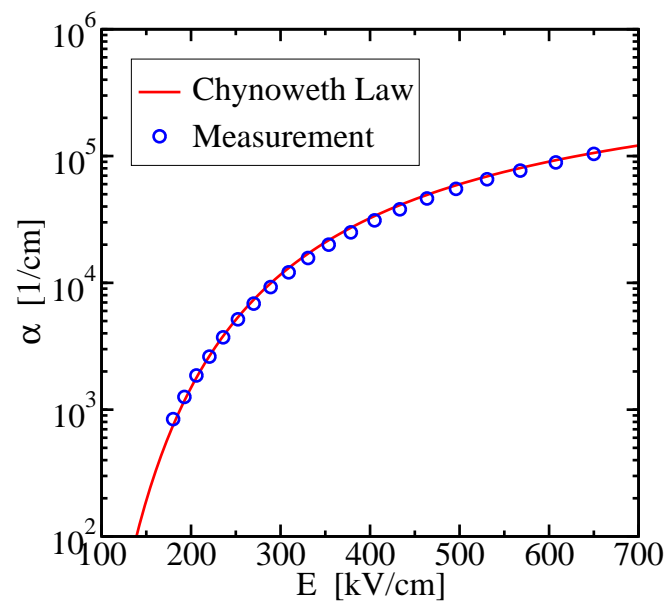
$$\alpha_n = A_n \exp\left(- (B_n/E)^{\beta_n}\right), \quad \alpha_p = A_p \exp\left(- (B_p/E)^{\beta_p}\right).$$

$E = \mathbf{E} \cdot \mathbf{J}/|\mathbf{J}|$  is the field component in direction of current flow. By interpreting a lot of experimental data, Chynoweth found the exponents  $\beta_n$  and  $\beta_p$  to be unity. Theoretical investigations by Shockley predicted the same exponents. A different treatment by Wolff predicts the exponents to be two. Fig. 7.2.4 shows the good agreement of theoretical and measurement results.

Practically,  $\beta_n$  and  $\beta_p$  are chosen between one and two in order to obtain good agreement with experimental data. For silicon at room temperature typical parameters are  $A_n = 7.03 \times 10^5 \text{ cm}^{-1}$ ,  $B_n = 1.231 \times 10^6 \text{ V/cm}$ ,  $\beta_n = 1$ ,  $A_p = 6.71 \times 10^5 \text{ cm}^{-1}$ ,  $B_p = 1.693 \times 10^6 \text{ V/cm}$  and  $\beta_p = 1$ .

---

<sup>1</sup> **prerequisite** [pri:'rek.wi.zit]: Voraussetzung, Bedingung



**Figure 7.7:** Verification of lattice scattering parameters for a homogeneous (bulk) sample.

## Chapter 8

# Devices in Detail

While the 1950s and 1960s were dominated by bipolar transistor technology, the 1970s saw Metal-Oxide-Semiconductor (MOS) technology begin to overtake bipolar technology in terms of functional complexity and level of integration. Geometrical scaling in MOS technology made it possible to increase chip complexity and to challenge bipolar circuits in the area of high-speed applications. In the 1980s, complementary MOS (CMOS) technology began to replace bipolar technologies such as transistor-transistor-logic (TTL) in many system applications [?, ?].

This chapter starts with the analytical description of a diode, after that the bipolar transistor is discussed. The main motivation is to highlight the approximations necessary to obtain closed form results. These approximations are not required in a numerical scheme and the results will be compared. The second half of this chapter deals with the MOS capacitor, leading to the MOS field effect transistor (MOSFET). Finally, scaling issues in CMOS technology will be discussed.

### 8.1 Analytical Diode Model

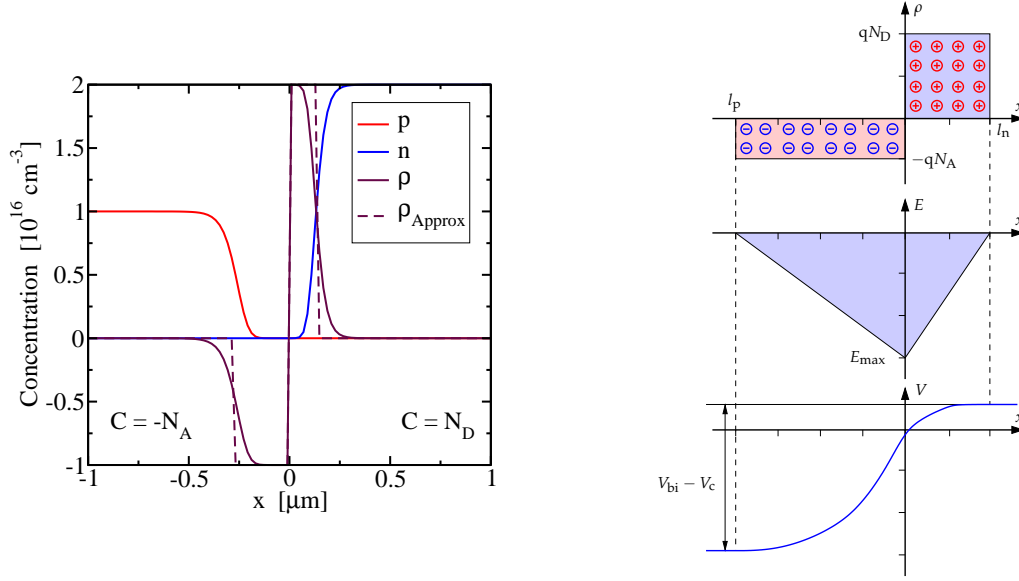
In 1944, Shockley theoretically investigated the junction of two semiconductor segments with different doping. When a  $p$ - and an  $n$ -doped region are joined, a *space-charge* will build up, resulting from diffusion of carriers from one material into the other caused by the difference of Fermi levels, leaving behind their ionized atoms they originate from. This results in an electric field that opposes the diffusion and drives the carriers back to their origin, until an equilibrium is reached (Fig. 8.1). In the space charge region, carriers of opposite charge recombine, so that the region is **depleted**<sup>1</sup> of carries, therefore this region is also called *depletion region*.

Inside the space charge region, we assume  $\rho = qC$ , while  $\rho = 0$  in the outside. The width of transition between these two regions is assumed to be negligible for analytical purposes. This *depletion approximation* is justified if the width of the space charge region is much larger than the Debye length  $L_D = \sqrt{\epsilon V_T / qN}$ . With this we can immediately solve Poisson's equation, because it is now decoupled from the carrier concentrations: Outside the depletion region,  $\rho = 0$  holds and thus  $E$  is constant. Inside the space-charge region,  $\rho = qC$  holds, so

$$\frac{dE}{dx} = \frac{q}{\epsilon}(p - n + C) \approx \frac{qC}{\epsilon}. \quad (8.1)$$

---

<sup>1</sup> to deplete [di'pli:t]: verringern, aufbrauchen, verarmen



**Figure 8.1:** Concentrations of electrons and holes around the  $pn$ -junction of a diode (left). The depletion approximation enables an analytical calculation of the electric field and the potential.

This gives

$$E(x) = \begin{cases} -\frac{q}{\epsilon} N_A (x + l_p), & -l_p \leq x \leq 0, \\ \frac{q}{\epsilon} (N_D x - N_A l_p), & 0 \leq x \leq l_n. \end{cases}$$

The maximum electric field is found at the interface:

$$E_{\text{max}} = E(0) = -\frac{q}{\epsilon} N_A l_p,$$

so the width of the negatively charged space-charge region  $l_p$  is

$$l_p = -\frac{\epsilon E_{\text{max}}}{q N_A}.$$

Since  $E = -dV/dx$ , the potential is found as (left contact at potential  $V_c$ )

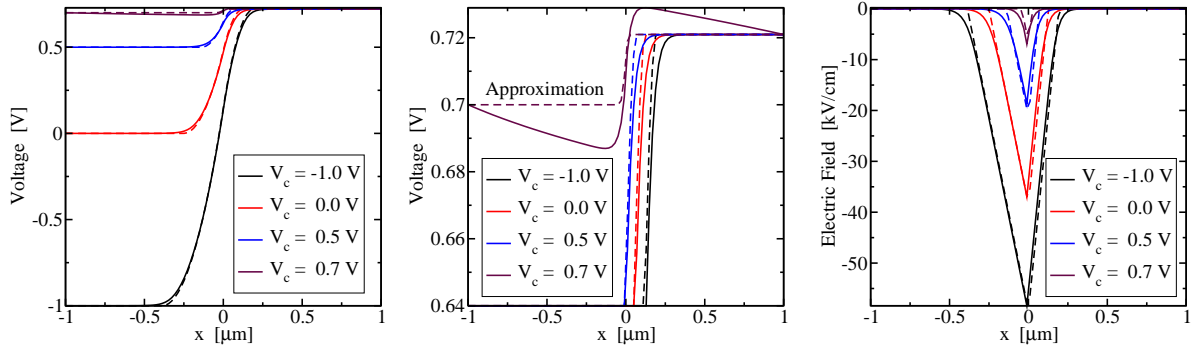
$$V(x) = \begin{cases} V_c & x \leq -l_p, \\ V_c + \frac{q}{2\epsilon} N_A (x - l_p)^2, & -l_p \leq x \leq 0, \\ V_c + \frac{q}{2\epsilon} (N_A l_p^2 - N_D x^2 + N_A x l_p), & 0 \leq x \leq l_n. \end{cases}$$

The total charge must vanish due to charge neutrality, therefore

$$N_A l_p = N_D l_n \quad \Rightarrow \quad l_n = -\frac{\epsilon E_{\text{max}}}{q N_D}.$$

At  $x = l_n$ , we find the potential

$$\begin{aligned} V(l_n) &= V_c + \frac{q}{2\epsilon} (N_A l_p^2 - N_D l_n^2 + N_A l_n l_p) \\ &= V_c + E_{\text{max}}^2 \frac{\epsilon}{2q} \left( \frac{1}{N_A} + \frac{1}{N_D} \right) \\ &= \psi_{bi}. \end{aligned}$$



**Figure 8.2:** Comparison of the electrostatic potential and the electric field obtained by numerical simulation with the (analytical) depletion approximation. Results are poor for  $V_c \gtrsim \psi_{bi}$  and there is a considerable voltage drop outside the space charge region.

At this point, the built-in potential from Section 2.3 turns up again: If two semiconductors with different Fermi level are joined, we must not forget about the potential difference at the contacts. We find

$$\psi_{bi} - V_c = E_{\max}^2 \frac{\epsilon}{2q} \left( \frac{1}{N_A} + \frac{1}{N_D} \right).$$

Conversely, the space charge region is defined via

$$E_{\max} = -\sqrt{\frac{2q}{\epsilon} \frac{\psi_{bi} - V_c}{\frac{1}{N_A} + \frac{1}{N_D}}},$$

which yields

$$l = \sqrt{\frac{2\epsilon}{q} (\psi_{bi} - V_c) \left( \frac{1}{N_A} + \frac{1}{N_D} \right)}, \quad l_p = l N_D / (N_A + N_D), \quad l_n = l N_A / (N_A + N_D). \quad (8.2)$$

The expression in the square root has to be non-negative, so we can immediately see that this model fails for the case  $V_c > \psi_{bi}$ ! A more physical explanation is that as  $V_c$  approaches  $\psi_{bi}$ , the opposing electric field becomes smaller and smaller. Therefore, diffusion acting on the carriers is only partially compensated by the force resulting from the junction potential variation. Thus, if  $V_c \approx \psi_{bi}$ , there is no clearly defined space-charge region anymore, so that our model assumptions cannot be justified anymore.

For  $V_c \ll \psi_{bi}$ , the simulation results for the potential and the electric field (cf. Fig. 8.3) are in good agreement with our analytical results. Let us consider the equilibrium case first: In the  $p$  region, outside the space charge region, there holds

$$p_{p0} = N_A, \quad n_{p0} = n_i^2 / N_A,$$

while in the  $n$  region we have

$$n_{n0} = N_D, \quad p_{n0} = n_i^2 / N_D.$$

Away from thermal equilibrium we have according to Chapter 2

$$n = n_i \exp\left(\frac{E_{Fn} - E_i}{k_B T_L}\right), \quad p = n_i \exp\left(\frac{E_i - E_{Fp}}{k_B T_L}\right),$$



so we get

$$np = n_i^2 \exp\left(\frac{E_{Fn} - E_{Fp}}{k_B T_L}\right) = n_i^2 \exp\left(\frac{V_c}{V_T}\right). \quad (8.3)$$

Shockley assumed in his derivation (1949) a so called *low-level injection*, which means that the excess minority carrier concentrations injected in a quasi-neutral region are low compared to the majority carrier concentration. With this assumption, the majority carrier concentrations remain at their equilibrium values:

$$n_n = n_{n0}, \quad p_p = p_{p0}. \quad (8.4)$$

From (8.3) we deduce

$$n_p = n_{p0} \exp\left(\frac{V_c}{V_T}\right), \quad p_n = p_{n0} \exp\left(\frac{V_c}{V_T}\right), \quad (8.5)$$

for the distribution of minority carriers concentrations at the edge of the space charge region.

Since the electric field outside the space-charge region is neglected, the resulting current in these neutral regions is purely diffusive. The majority carrier concentrations are due to (8.4) constant in the quasi-neutral regions, but (8.5) tells us that this is not true for the minority concentrations. Consequently, the minority carrier concentrations determine the (purely diffusive) current.

Thus, the minority carrier current inside the  $p$ -region is

$$J_n \approx \mu_n k_B T_L \nabla n. \quad (8.6)$$

We will further neglect recombination both inside and outside the space-charge region, thus  $J_n \approx \text{const.}$  With  $n(-l_p) = n_p$  as first boundary condition we can write

$$n(x) = n_p + \frac{J_n}{\mu_n k_B T_L} (x + l_p)$$

Let  $L$  denote the width of the  $p$ -region and with  $n(-L) = n_{p0}$  and  $L \gg l_p$  we obtain for the current

$$J_n = \mu_n k_B T_L \frac{n_p - n_{p0}}{L - l_p} \approx \mu_n k_B T_L \frac{n_p - n_{p0}}{L} = \frac{\mu_n k_B T_L n_i^2}{N_A L} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right).$$

An analogous procedure for holes yields

$$J_p = \mu_p k_B T_L \frac{p_n - p_{n0}}{L - l_n} \approx \mu_p k_B T_L \frac{p_n - p_{n0}}{L} = \frac{\mu_p k_B T_L n_i^2}{N_D L} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right).$$

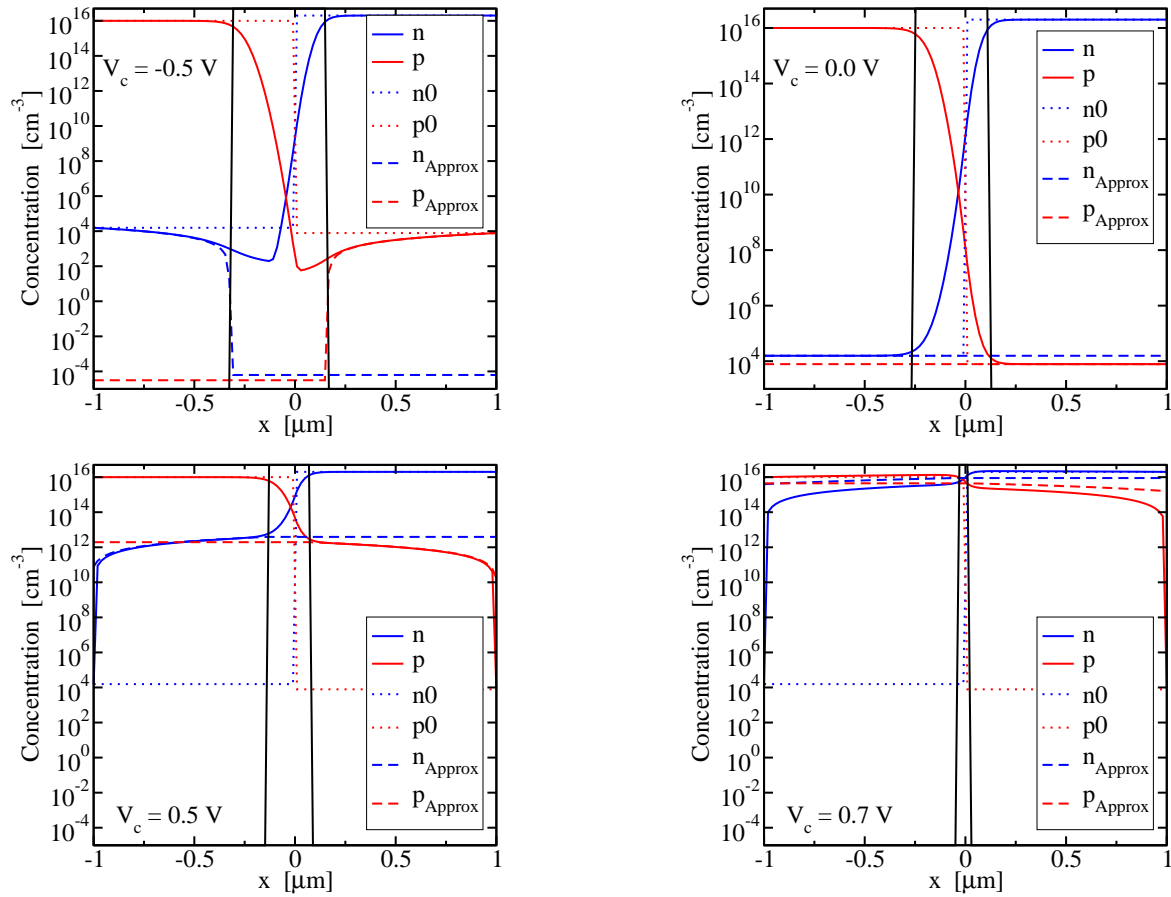
The total diode current is the sum of the electron and hole currents (which do not interact because we have neglected  $R$ ):

$$J = J_n + J_p = \frac{\mu_n k_B T_L n_i^2}{N_A L} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right) + \frac{\mu_p k_B T_L n_i^2}{N_D L} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right) \quad (8.7)$$

$$= \underbrace{\frac{k_B T_L n_i^2}{L} \left( \frac{\mu_n}{N_A} + \frac{\mu_p}{N_D} \right)}_{=: J_s} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right), \quad (8.8)$$

which is the characteristic current relation for the diode.

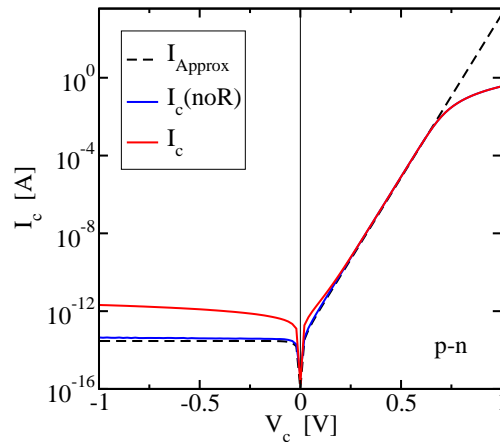
Let us recapitulate the assumptions leading to the analytic diode model:



**Figure 8.3:** Comparison of the compact model with numerical results. At  $V_c \neq 0.0$  V we observe a slope in the minority carrier concentrations, which determines the current-voltage relation. At  $V_c \approx 0.7$  V, the space charge region becomes very thin and the compact model does not show good agreement with the simulation results anymore.

- We have assumed *abrupt doping profiles*. In reality, such profiles do not exist. This is not a problem as long as the doping transition length is small compared to the space charge region. For small diodes with high (asymmetric) doping concentrations, space charge regions become very thin (cf. Eq. 8.2) and realistic shapes of doping profiles have to be considered.
- The *depletion approximation* assumes an abrupt space charge region. In reality, there is a smooth (but still steep) transition between the quasi neutral and the space charge regions in the order of the Debye length, which has to be considered for highly asymmetric doping profiles.
- *No recombination* within the space charge region was assumed. Real diodes are, however, non-ideal and the effects of generation and recombination have to be taken into account. Doing so, the exponential in (8.7) has to be extended by an *ideality factor*  $1 \leq n \leq 2$ , so that the current relation becomes

$$J = J_s \left( \exp \left( \frac{V_c}{nV_T} \right) - 1 \right).$$



**Figure 8.4:** A real diode differs from the compact model especially through high-level injection and increased reverse-bias current due to generation and recombination.

- *Short diode:* On the one hand, recombination outside the space charge region can be neglected (diode is relatively short), but the diode is still long enough to neglect the space charge width in the gradient at (8.6).
- *Low-level injection:* We assumed that majority carriers were not influenced and that a pure minority diffusion current dominated.

A comparison of the compact model with numerical results in Fig. 8.3 shows good agreement of the minority concentrations outside the space charge region. However, as the applied voltage approaches  $\psi_{bi}$ , the approximation becomes poorer.

Finally, some further effects in real diodes shall be mentioned: First, *high-level injection* decreases the current for higher bias voltages. Second, the voltage drop in the quasi neutral regions decreases the current at high bias voltages (*high bias effect*). The third effect we mention here is a higher current under reverse-bias compared to the predictions of our model. This can be traced back to generation in the space charge region.

Apart from the physical effects mentioned, a real diode is at least two-dimensional, resulting in additional corner effects (larger electric fields!) that cannot be tackled by a simple one-dimensional approximation. Furthermore, doping profiles are more complex and the mobility is not constant either. Fig. 8.4 compares numerically simulated diode characteristics with the results from the compact model.

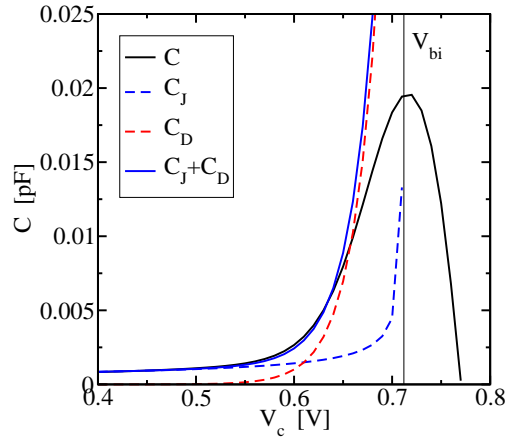
## 8.2 Small Signal Analysis of a Diode

In (8.7) the DC-characteristic of a diode is given. The transient description, however, requires additional considerations. In particular, charge storage effects have to be taken into account, otherwise the device would be infinitely fast, which contradicts the **inertia**<sup>1</sup> of charges.

There are two dominant forms of charge-storage:

- Charge is stored in the depletion region due to the dopants, resulting in a *junction capacitance* (or *depletion capacitance*). The junction capacitance is caused by charge dipoles

<sup>1</sup> **inertia** [ɪˈnɜːrjə]: Trägheit



**Figure 8.5:** Comparison of compact model with numerical simulation with small-signal mode of *Minimos – NT*. There is a good agreement for  $V_c$  not too close to  $\psi_{bi}$ .

formed by ionized dopants in the depletion region so the depletion region width and thus the charge changes with the applied voltage. Let the contact potential change by a small amount  $dV_c$ . Then the depletion region charge per unit area changes by  $dQ'_J$ . The junction capacitance per unit area is therefore defined as

$$C'_J = \frac{dQ'_J}{dV_c}.$$

Since  $Q'_J = qN_D l_n = qN_A l_p$  we get with (8.2)

$$C'_J(V_c) = qN_D \frac{dl_n}{dV_c} = \sqrt{\frac{q\epsilon}{2(\psi_{bi} - V_c)(1/N_A + 1/N_D)}} = \frac{C'_J(0)}{\sqrt{1 - V_c/\psi_{bi}}}$$

Using again (8.2), we get

$$C'_J(V_c) = \frac{\epsilon}{l},$$

which is just the same as for a parallel-plate capacitor!

- Charge is stored due to injected minority-carriers in the space-charge regions, leading to the *diffusion capacitance* (or *storage capacity*)  $C'_D$ .

The hole charge stored in the  $n$ -region is

$$Q'_p = q \int_{l_n}^L (p_n(x) - p_{n0}) dx = \frac{qL}{2} p_{n0} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right) \quad \text{for } L \gg l_n.$$

Similarly, the electron charge stored in the  $p$ -region is

$$Q'_n = q \int_{-l_p}^{-L} (n_p(x) - n_{p0}) dx = \frac{qL}{2} n_{p0} \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right) \quad \text{for } L \gg l_p.$$

In total we have

$$C'_D = \frac{dQ'_D}{dV_c} = \frac{dQ'_p + dQ'_n}{dV_c} = \frac{qLn_i^2}{2V_T} \left( \frac{1}{N_A} + \frac{1}{N_D} \right) \left( \exp\left(\frac{V_c}{V_T}\right) - 1 \right).$$

Due to the exponential term, the contribution is negligible under reverse bias. However, as  $V_c$  approaches  $\psi_{bi}$ , the capacitance is overestimated (cf. Fig. 8.5).

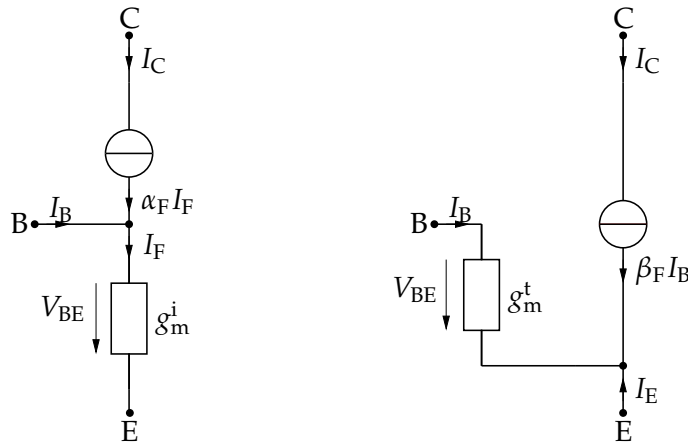


Figure 8.6: Frequently used equivalent circuit diagrams for a BJT.

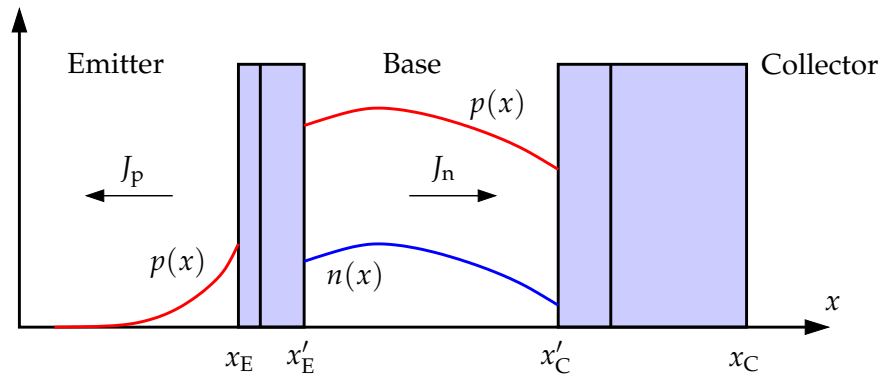


Figure 8.7: Electron and hole densities in a  $npn$ -transistor.

### 8.3 Analytical Bipolar Junction Transistor Model

As for the diode, we will derive a compact model for *bipolar junction transistors* (BJTs). In the literature, two equivalent circuit diagrams for such a model are often used, shown in Fig. 8.6. They are indeed equivalent, which is not obvious at first sight, but shall be shown now.

We start with the electron current relations in the base of a  $npn$ -transistor:

$$J_n = q\mu_n nE + qD_n \nabla n, \quad (8.9)$$

$$\nabla \cdot J_n = qR. \quad (8.10)$$

The Shockley-Read-Hall recombination rate (7.13) in the base with  $p \gg n$  and  $p \approx N_A$  becomes

$$R^{\text{SRH}} \approx \frac{nN_A - n_i^2}{\tau_n N_A} = \frac{n - n_{p0}}{\tau_n} =: \frac{\Delta n}{\tau_n},$$

where  $n_{p0} = n_i^2/N_A$  is the (constant) equilibrium concentrations of electrons in the  $p$ -type base.

Assuming that the diffusion current dominates in (8.9), we obtain

$$qD_n \nabla^2 n = q \frac{\Delta n}{\tau_n}.$$

Cancelling  $q$  and using the identity  $\nabla^2 n = \nabla^2(n - n_{p0}) = \nabla^2(\Delta n)$  (because we assumed  $n_{p0}$  to be constant) leads to

$$D_n \nabla^2(\Delta n) = \frac{\Delta n}{\tau_n}.$$

In one dimension, the general solution of this differential equation is

$$\Delta n = C_1 \exp\left(\frac{x}{L_B}\right) + C_2 \exp\left(-\frac{x}{L_B}\right),$$

where  $L_B = \sqrt{D_n \tau_n}$  is the diffusion length. From the base-emitter-diode we find as boundary condition

$$n(0) = n_{p0} \exp\left(\frac{V_{BE}}{V_T}\right), \quad \text{hence} \quad \Delta n(0) = n_{p0} \left( \exp\left(\frac{V_{BE}}{V_T}\right) - 1 \right) \quad (8.11)$$

and from the base-collector-diode

$$n(W) = n_{p0} \exp(V_{BC}/V_T), \quad \text{hence} \quad \Delta n(W) = n_{p0} (\exp(V_{BC}/V_T) - 1). \quad (8.12)$$

The recombination in the base is small, therefore  $L_B$  is large,  $L_B \gg W$ , where  $W$  is the base width. For this reason, we can linearize  $\exp(x) \approx 1 + x$  around zero and can now write

$$\Delta n \approx \Delta n(0) \left(1 - \frac{x}{W}\right) + \Delta n(W) \frac{x}{W}.$$

The current density at the beginning of the base is

$$J_n(0) = qD_n \frac{d\Delta n}{dx} = \frac{qD_n}{W} (\Delta n(W) - \Delta n(0)) \quad (8.13)$$

and with (8.11) and (8.12) this becomes

$$I_n = I_{n0} ((\exp(V_{BE}/V_T) - 1) - (\exp(V_{BC}/V_T) - 1)) = I_{n1} - I_{n2},$$

where  $I_{n0} = qD_n A n_i^2 / (N_A W)$ . Note that the minority carriers in the emitter and collector are not affected:

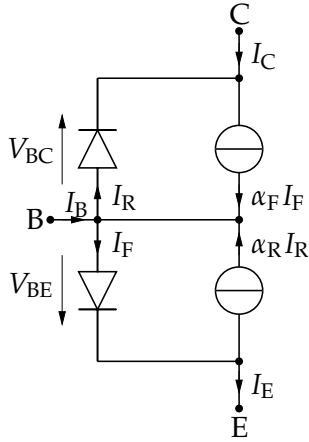
$$I_{Ep} = \frac{qD_p A}{W} \frac{n_i^2}{N_{DE}} (\exp(V_{BE}/V_T) - 1), \quad I_{Cp} = \frac{qD_p A}{W} \frac{n_i^2}{N_{DC}} (\exp(V_{BC}/V_T) - 1).$$

For later use we define:

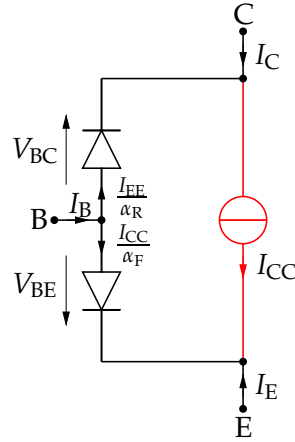
$$\frac{I_{n2}}{I_{Ep}} = \frac{D_n L_p N_{DE}}{D_p W N_A} =: g_F, \quad \frac{I_{n1}}{I_{Cp}} = \frac{D_n L_p N_{DC}}{D_p W N_A} =: g_R.$$

The total emitter current is then

$$I_E = -I_n + I_{Ep} = -I_{n1} + I_{n2} + I_{Ep} = -I_{Cp} g_R + \underbrace{I_{Ep}(1 + g_F)}_{I_F};$$



**Figure 8.8:** Injection version of the Ebers-Moll model.



**Figure 8.9:** Transport version of the Ebers-Moll model.

Similarly, the collector current can be written as

$$I_C = -I_n - I_{Cp} = -I_{n1} + I_{n2} - I_{Cp} = -\underbrace{I_{Cp}(1 + g_R)}_{I_R} + I_{Ep}g_F.$$

This gives for the controlled currents

$$g_F I_{Ep} = \frac{1 + g_F}{g_F} I_F =: \alpha_F I_F, \quad g_R I_{Cp} = \frac{1 + g_R}{g_R} I_R =: \alpha_R I_R$$

and so we arrive at

$$I_E = -\alpha_R I_R + I_F, \quad I_C = -I_R + \alpha_F I_F.$$

This is the *injection version* of the *static Ebers-Moll model*, where

$$I_F = I_{ES} \left( \exp \left( \frac{V_{BE}}{V_T} \right) - 1 \right), \quad I_R = I_{CS} \left( \exp \left( \frac{V_{BC}}{V_T} \right) - 1 \right). \quad (8.14)$$

If we define

$$I_{CC} = \alpha_F I_F = I_S \left( \exp \left( \frac{V_{BE}}{V_T} \right) - 1 \right), \quad (8.15)$$

$$I_{EE} = \alpha_R I_R = I_S \left( \exp \left( \frac{V_{BC}}{V_T} \right) - 1 \right), \quad (8.16)$$

with  $I_S = \alpha_F I_{ES} = \alpha_R I_{CS}$ , we arrive at the *transport version* of the static Ebers-Moll model with

$$I_{CT} = I_{CC} - I_{EE}, \quad \beta_F = \frac{1 - \alpha_F}{\alpha_F}, \quad \beta_R = \frac{1 - \alpha_R}{\alpha_R}.$$

Let us come back to the injection version of the model. In the active region there holds  $I_R \approx 0$ , so linearization of (8.14) around the operating point  $(V_{BE}^0, I_F^0)$  yields

$$\begin{aligned} I_F(V_{BE}) &= I_F(V_{BE}^0) + \underbrace{\frac{\partial I_F}{\partial V_{BE}} \Big|_{V_{BE}^0}}_{=:g_m^i} (V_{BE} - V_{BE}^0) \\ &= I_F^0 + g_m^i (V_{BE} - V_{BE}^0), \end{aligned}$$

so in small signal analysis (i.e.  $V_{BE} = V_{BE}^0 + v_{BE}$ ,  $I_F = I_F^0 + i_F$ ) this results in

$$i_F = g_m^i v_{BE}, \quad (8.17)$$

which results in the equivalent circuit shown on the left of Fig. 8.6.

We can do similar simplifications for the transport version of the model: In the active region,  $I_{EE} \approx 0$  and with (8.15) we can write the current at the basis as

$$I_{CC} = I_{BS} \left( \exp\left(\frac{V_{BE}}{V_T}\right) - 1 \right).$$

Linearization around the operating point yields

$$\begin{aligned} I_B(V_{BE}) &= I_B(V_{BE}^0) + \underbrace{\frac{\partial I_B}{\partial V_{BE}} \Big|_{V_{BE}^0}}_{=:g_m^t} (V_{BE} - V_{BE}^0) \\ &= I_B^0 + g_m^t (V_{BE} - V_{BE}^0), \end{aligned}$$

so in small signal analysis we find in analogy to (8.17) the relation

$$i_B = g_m^t v_{BE},$$

thus the equivalent circuit shown on the right of Fig. 8.6 is justified.

Comparing both models, the conductivities  $g_m^i = \partial I_F / \partial V_{BE}$  and  $g_m^t = \partial I_B / \partial V_{BE}$  are not the same. However, they are related via

$$g_m^t = (1 - \alpha_F) g_m^i, \quad \alpha_F \leq 1,$$

as can easily be seen from Fig. 8.6 after expressing  $I_B$  as a function of  $I_F$ .

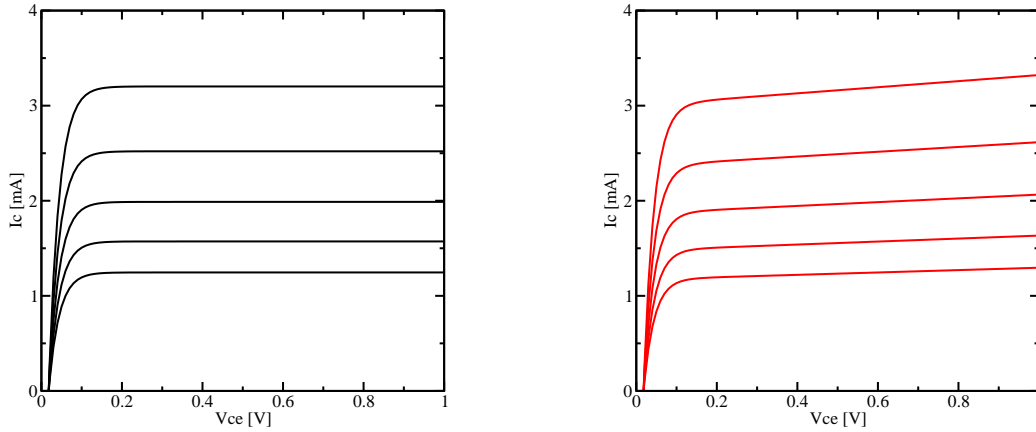
How well do the derived compact models reflect a real transistor? The most obvious difference is that the *transfer characteristics*  $I_C(V_{BE})$  of our compact model show an output conductivity of zero. The reason for this is that we assumed a constant width of the base in (8.13). However, this is not true: The neutral base width is modulated by the applied collector-base reverse bias  $V_{CB}$ , which gives rise to the so-called *Early effect*: The off-state current  $I_S$  as well as the amplification factor  $\beta_F$  depend on  $V_{CB}$ , in linear approximation

$$\begin{aligned} I_S(V_{CB}) &= I_S(0) \left( 1 - \frac{V_{CB}}{V_A} \right), \\ \beta_F(V_{CB}) &= \beta_F(0) \left( 1 - \frac{V_{CB}}{V_A} \right), \end{aligned}$$

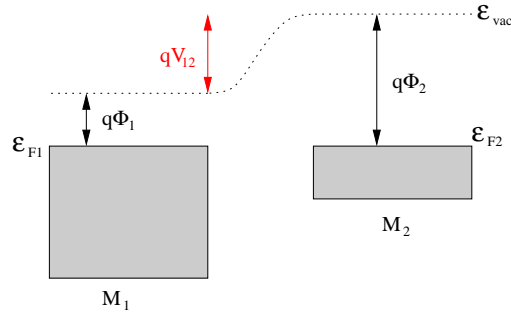
where  $V_A$  is the *Early voltage*. The output conductance is then

$$g_c = \frac{\partial I_C}{\partial V_{CE}} \Big|_{V_{BE}=\text{const.}} \approx \frac{I_C(0)}{V_A}.$$





**Figure 8.10:** Transfer characteristic of our compact model (left) and after taking the Early effect into account (right).



**Figure 8.11:** Work functions are used to describe a potential drop between two materials.

## 8.4 The Metal-Oxide-Semiconductor Capacitor

The work function of a material is the potential difference from the vacuum level to the Fermi level. Thus, for two materials

$$q\Phi_1 = E_{\text{vac}} - E_{F1}, \quad q\Phi_2 = E_{\text{vac}} - E_{F2}.$$

The potential drop between material  $M_1$  and material  $M_2$  is the *contact potential*

$$\psi_{12} = \Phi_2 - \Phi_1.$$

With the concept of work functions in hand, we proceed to the simplest layout of a *MOS capacitor* (Fig. 8.12). Going from the gate to the substrate, the potentials are

$$\psi_{M1} + \psi_{1S} = (\Phi_1 - \Phi_M) + (\Phi_S - \Phi_1) = \Phi_S - \Phi_M = \Phi_{MS}. \quad (8.18)$$

Metal	Al	Pt	W	Mg	Ag	Au	Cu	Cr
$q\Phi_M$ (eV)	4.28	5.65	4.63	3.66	4.30	4.80	4.25	4.50

**Table 8.1:** Work function values for several metals.

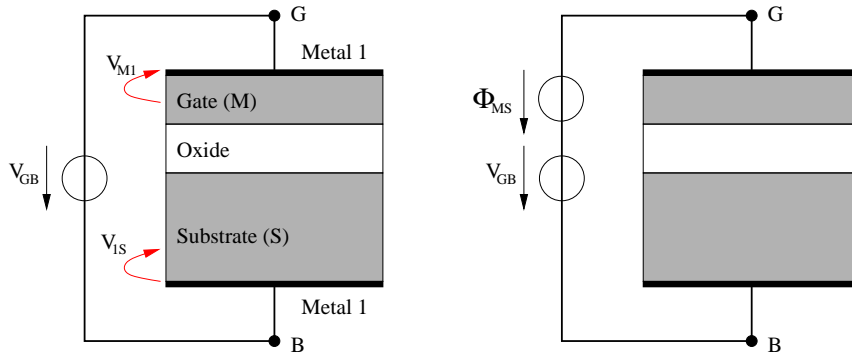


Figure 8.12: Schematics of a MOS capacitor.

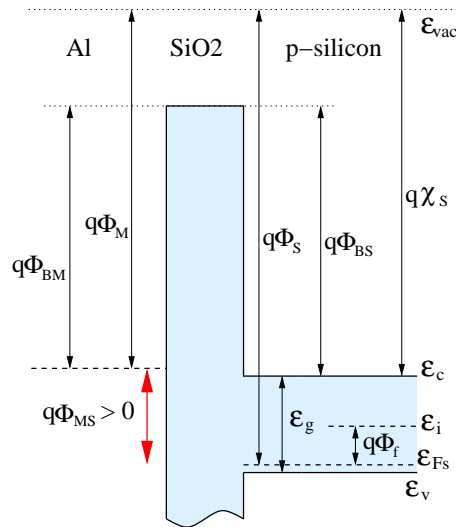


Figure 8.13: Band diagram of an ideal MOS capacitor.

Taking the work function difference  $\Phi_{MS}$  into account allows us to assume ideal contacts in the sense that no additional built-in potentials occur. Conversely, if the voltage applied from the outside equals the work function difference, carrier concentrations are at their equilibrium values, so there is no voltage drop. From Kirchhoff's voltage law,  $\Phi_{MS} + V_{GB} = 0$ , so the *flat band voltage* is  $\psi_{FB} = -\Phi_{MS}$ .

Let us consider a more realistic setting, where the gate is made of aluminum, the oxide is  $\text{SiO}_2$  and the substrate is  $p$ -doped silicon (Fig. 8.13). For the flat band voltage there holds  $V_{FB} = -\Phi_{MS} > 0$ . Similar to the work function in metals, we introduce the *electron affinity*  $\chi_S$  of a semiconductor: It is the energy required to detach an electron from a singly charged negative ion. For our purposes, it is proportional to the difference between the vacuum potential and the conduction band edge:  $q\chi_S = E_{vac} - E_c$ .

The *Fermi potential* is given as  $q\Phi_F = E_i - E_{Fs}$ , thus we can write the work function as

$$q\Phi_S = E_{vac} - E_{Fs} = q\chi_S + E_g - (E_i - E_v) + q\Phi_F. \quad (8.19)$$

Semiconductor	Si	Ge	GaAs	GaP	GaSb	InAs	InP	InSb
$\chi_s$ [V]	4.05	4.00	4.07	3.80	4.06	4.90	4.38	4.59

Table 8.2: Electron affinity for several semiconductors.

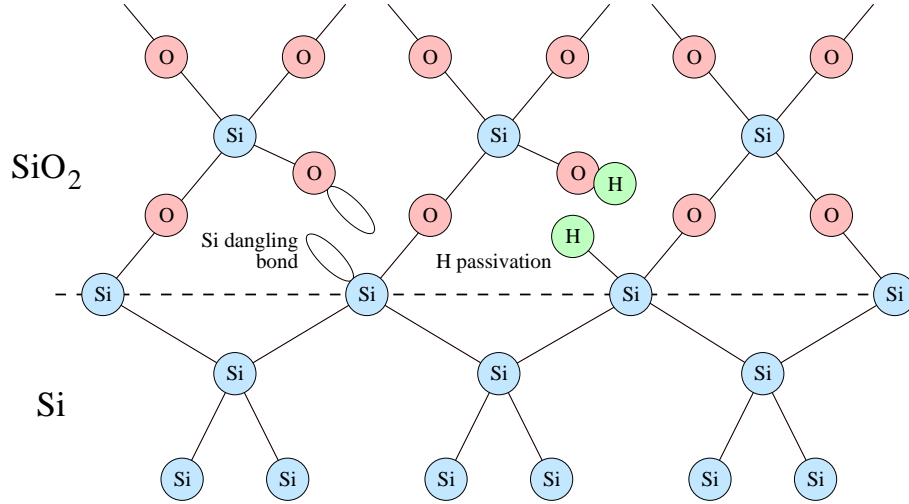


Figure 8.15: Dangling bonds on the interface can be passivated with hydrogen.

In Fig. 8.13 all quantities are illustrated as well as the barrier heights  $q\Phi_{BM}$  and  $q\Phi_{BS}$ . The Fermi potentials depend on the doping level:

$$p\text{-type} : \Phi_F = +\frac{k_B T}{q} \ln\left(\frac{N_A}{n_i}\right) > 0,$$

$$n\text{-type} : \Phi_F = -\frac{k_B T}{q} \ln\left(\frac{N_D}{n_i}\right) < 0.$$

Since the Fermi potential is the only contribution that depends on the doping, a doping-independent part

$$q\Phi'_{MS} = q\Phi_M - q\chi_s - E_g + E_i - E_v.$$

Typical values are  $\Phi'_{MS} = -0.6$  V for an Al-gate and a silicon substrate and  $\Phi'_{MS} = 0.3$  V for an Au-gate and a silicon substrate. For a silicon gate, one finds  $\Phi_{MS} = \Phi_F^G - \Phi_F^B$ . In the case of a  $n^+$ -gate,  $\Phi_F^G = -0.56$  V.

In a non-ideal MOS capacitor, positively charged *dangling bonds* at the Si-SiO<sub>2</sub>-interface occur (cf. Fig. 8.15). Most of them are passivated by hydrogen. Apart from charges caused by dangling bonds, there may be interface trapped charges, oxide trapped charges or mobile ionic charges such as sodium, but they are usually negligible compared to  $10^9 \dots 10^{10} \text{ m}^{-2}$  caused by fixed oxide charges (that normally carry a positive charge) and

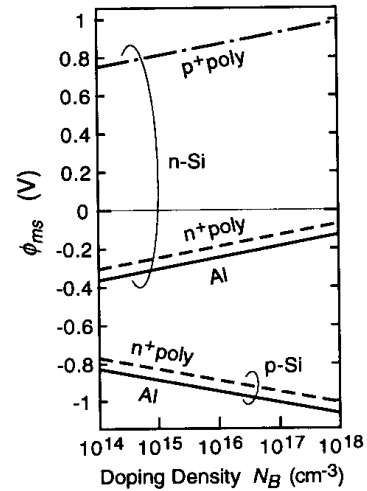


Figure 8.14: Work functions for several materials and doping densities.

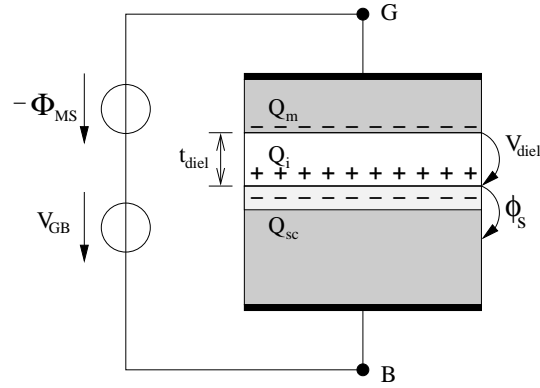


Figure 8.16: Charges in a MOS capacitor.

dangling bonds (positive charges in a pMOS and negative charges in an nMOS). These charges cause a shift in the flat band voltage. With the notations used in Fig. 8.16, the balance of charges requires

$$Q_m + Q_i + Q_{sc} = 0. \quad (8.20)$$

For the dielectric displacement,  $\epsilon_{diel} E_{diel} = Q_m$  holds. Using Kirchhoff's voltage law again, one finds

$$\psi_{GB} - \Phi_S - \psi_{diel} + \Phi_{MS} = 0. \quad (8.21)$$

The flat band condition requires  $Q_{sc} = 0$  and  $\Phi_S = 0$ . Putting all equations together, we find

$$\psi_{diel} = E_{diel} t_{diel} = \frac{Q_m}{\epsilon_m} t_{diel} = -\frac{Q_i}{C_{diel}}. \quad (8.22)$$

So, the new flat band voltage is

$$\psi_{FB} = \frac{Q_i}{C_{diel}} - \Phi_{MS}. \quad (8.23)$$

Next we are going to investigate the case  $V_{GB} \neq V_{FB}$ . In this case, a space-charge region forms in the semiconductor near the interface. The total potential drop caused by the space-charge region is called the *surface potential*  $\phi_S$  and induces a shift of the band edges:

$$\begin{aligned} E_c(x) &= E_{c,0} - q\phi(x), \\ E_v(x) &= E_{v,0} - q\phi(x), \end{aligned}$$

where the potential  $\phi|_{bulk} = 0$  and the surface potential is  $\phi_S = \phi(0)$ .

In Tab. 8.3 and Fig. 8.17 the band edge energies and electrostatic potential for a nMOS capacitor for several bias voltages are shown. Very important with respect to the MOSFET is the case of inversion: In equilibrium, the carrier concentrations are given by (cf. Chapter 2)

$$n = n_i \exp\left(\frac{E_{Fn} - E_i}{k_B T_L}\right), \quad p = n_i \exp\left(\frac{E_i - E_{Fp}}{k_B T_L}\right). \quad (8.24)$$

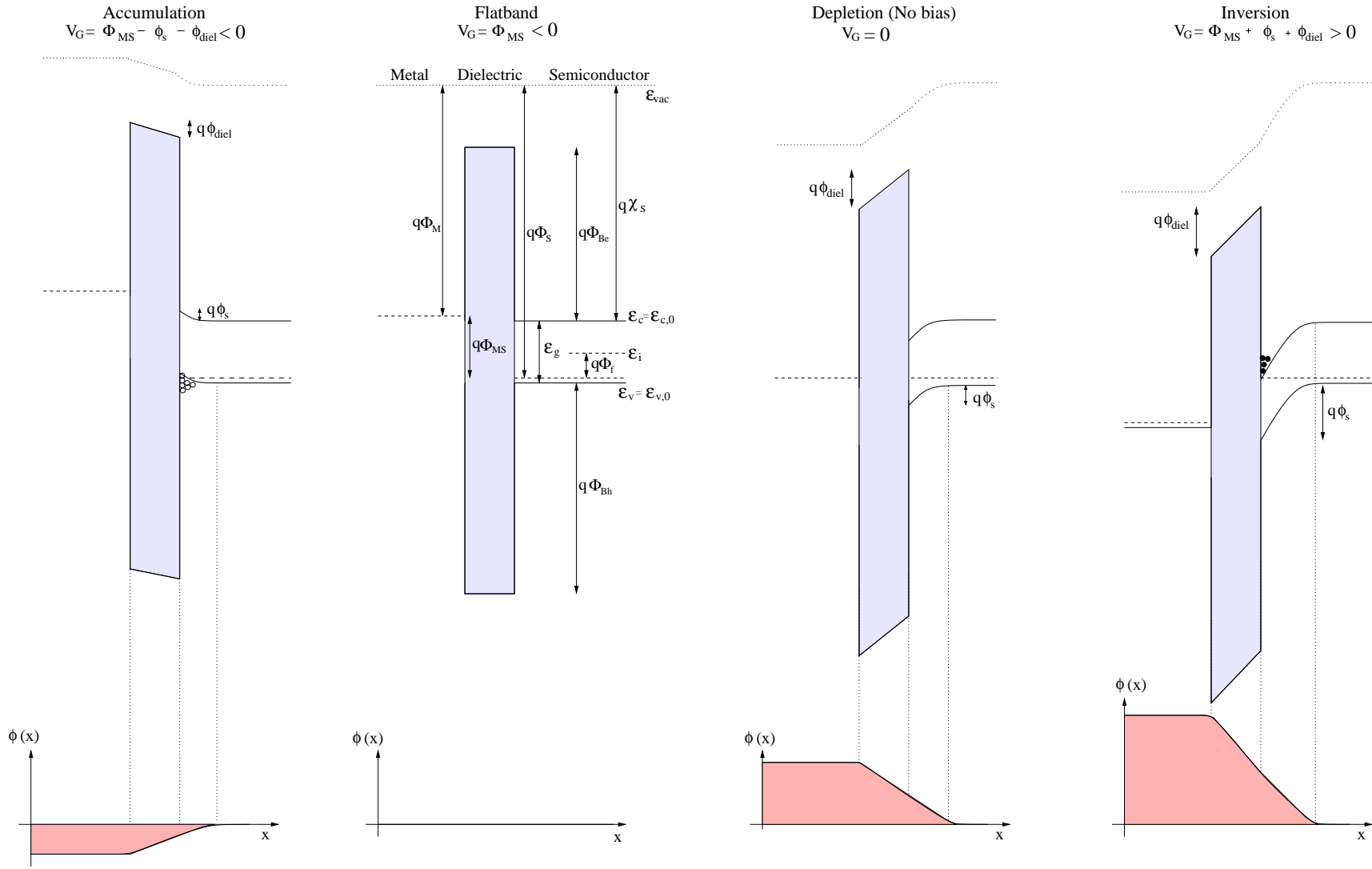
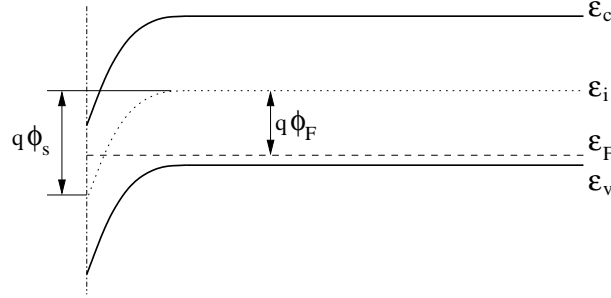


Figure 8.17: Band energies and electrostatic potential for an nMOS in different operation conditions.

	Bias	Surface potential	Charges in substrate
Accumulation	$V_{GB} < V_{FB}$	$\phi_S < 0$	$Q_{sc} > 0$
Flat band	$V_{GB} = V_{FB}$	$\phi_S = 0$	$Q_{sc} = 0$
Depletion	$V_{GB} > V_{FB}$	$\phi_S > 0$	$Q_{sc} < 0$
Inversion	$V_{GB} \gg V_{FB}$	$\phi_S > 0$	$Q_{sc} < 0$

**Table 8.3:** Operation ranges of a nMOS (i.e. p-type substrate,  $V_{FB} < 0$ ).



**Figure 8.18:** Inversion occurs when the conduction band-edge falls below the Fermi level.

The onset of strong inversion is (somewhat arbitrarily) defined as the point at which  $n_{\text{surf}} = p_{\text{bulk}}$ , where

$$n_{\text{surf}} = n_i \exp\left(\frac{E_{Fn} - E_{i,\text{surf}}}{k_B T_L}\right), \quad E_{i,\text{surf}} = E_i - q\phi_S. \quad (8.25)$$

According to Fig. 8.18, the conduction band-edge falls below the Fermi-level as soon as  $\phi_S = 2\Phi_F$ .

Still, charge neutrality has to hold in inversion as well. According to Fig. 8.19, charges in the space charge region consist of a depletion layer charge  $Q_d$  and an inversion (channel) charge  $Q_{\text{ch}}$ , thus

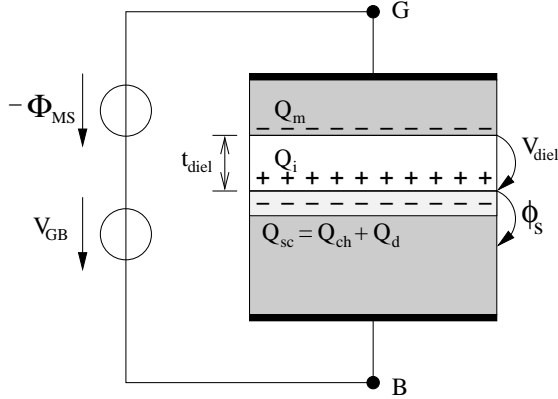
$$Q_m + Q_i + Q_{\text{ch}} + Q_d = 0. \quad (8.26)$$

Kirchhoff's voltage law requires  $V_{GB} = \Phi_{MS} + V_{\text{diel}} + \phi_S$ . As before, the dielectric displacement is  $\epsilon_{\text{diel}} E_{\text{diel}} = Q_m$ , which together with the conditions  $\phi_S = 2\Phi_{MS}$  and  $Q_{\text{ch}} \approx 0$  for strong inversion yields the *threshold voltage*

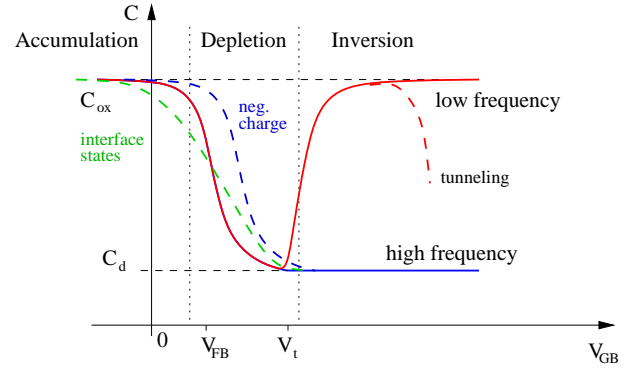
$$V_t = \Phi_{MS} - \frac{Q_i}{C_{\text{diel}}} + 2\Phi_F - \frac{Q_d}{C_{\text{diel}}}. \quad (8.27)$$

For  $n$ -channel devices (i.e.  $p$ -bulk), normally  $V_t \geq 0$ , while for  $p$ -channel devices (i.e.  $n$ -bulk),  $V_t < 0$ . One can express  $Q_d$  in (8.27) by

$$Q_d = \begin{cases} -qN_A d, & p\text{-type,} \\ qN_D d, & n\text{-type.} \end{cases}$$



**Figure 8.19:** Charges in a MOS channel in case of inversion.



**Figure 8.20:**  $C(V)$ -curve of a MOS capacitor.

With the *depletion layer width*

$$d = \begin{cases} \sqrt{\frac{2\epsilon_S\phi_S}{qN_A}}, & p\text{-type,} \\ \sqrt{\frac{2\epsilon_S|\phi_S|}{qN_D}}, & n\text{-type,} \end{cases} \quad (8.28)$$

one finally obtains

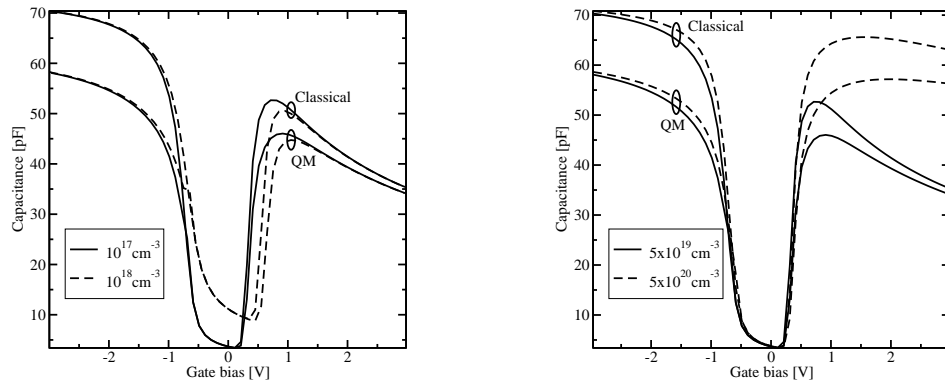
$$V_t = \begin{cases} \Phi_{MS} - \frac{Q_i}{C_{diel}} + 2\Phi_F - \frac{\sqrt{4\epsilon_S q N_A \Phi_F}}{C_{diel}}, & p\text{-MOS } (\Phi_F > 0), \\ \Phi_{MS} - \frac{Q_i}{C_{diel}} + 2\Phi_F - \frac{\sqrt{4\epsilon_S q N_D |\Phi_F|}}{C_{diel}}, & n\text{-MOS } (\Phi_F < 0). \end{cases} \quad (8.29)$$

The main electrical measurement to determine the insulator quality in a MOSFET is the *capacity-voltage curve* ( $C(V)$ -curve) of the MOS-capacitor between gate and bulk. It enables the determination of insulator thickness (magnitude of capacity in accumulation), the flat band voltage as well as threshold voltage, the bulk doping (from the inversion capacitance) and the interface trap density (shape of  $C(V)$ -curve), as can be seen in Fig. 8.20. It is important to note that there are two distinct behaviors at inversion, depending on the frequency of the applied signal. For high frequencies, the inversion channel cannot be built up properly, thus only  $C_d$  contributes. In a real MOSFET, the high frequency branch is not observed as carriers are supplied by the source and drain regions.

Despite all the details we put into our model, it does not include any quantum-mechanical effects, which leads to significant deviations as devices are scaled down to several nanometers. In Fig. 8.21 a comparison of classical numerical simulation with quantum-mechanical simulation for an nMOS with  $t_{diel} = 15\text{\AA}$  is shown.

## 8.5 The Metal-Oxide-Semiconductor Field-Effect-Transistor (MOS-FET)

Already back in 1926, Julius Edgar Lilienfeld proposed a *Method and Apparatus for Controlling Electric Currents* (Fig. 8.22). However, fabrication was not possible due to material-related problems. In 1960, Kahng and Attala realized the first field-effect transistor in MOS technology. We will summarize the analytical compact model given by Sah in 1964 [?].



**Figure 8.21:** Influence on substrate doping (left) and polydepletion (right) using a classical and quantum-mechanical simulation of a nMOS with  $t_{\text{diel}} = 15 \text{ \AA}$ .

In the *linear region*, where  $V_{\text{GS}} > V_t$  and  $V_{\text{GD}} > V_t$ , the drain current is given as

$$I_D = \frac{W}{L} \mu_n C_{\text{diel}} \left( (V_{\text{GS}} - V_t) V_{\text{DS}} - \frac{1}{2} V_{\text{DS}}^2 \right),$$

with the *threshold voltage*

$$V_t = V_{\text{FB}} + 2\Phi_F + \gamma \sqrt{2\Phi_F - V_{\text{BS}}} \quad (8.30)$$

including the *body factor*

$$\gamma = \frac{1}{C_{\text{diel}}} \sqrt{2\epsilon_{\text{Si}} q N_A}.$$

This is the level 1 MOSFET model of the circuit simulator SPICE.

In the *saturation region* (also referred to as *pinch-off region*),  $V_{\text{GS}} > V_t$ , but  $V_{\text{GD}} < V_t$ . The current remains at constant level as soon as the saturation voltage is reached:

$$\frac{dI_D}{dV_{\text{DS}}} = 0 \implies V_{\text{DS}}^{\text{sat}} = V_{\text{GS}} - V_t.$$

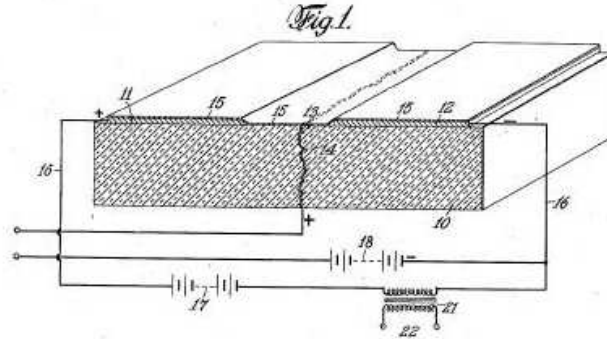
The idealized MOSFET then behaves like an ideal current source (i.e. no dependence on  $V_{\text{DS}}$ ):

$$I_D = \frac{W}{L} \mu_n C_{\text{diel}} \frac{1}{2} (V_{\text{GS}} - V_t)^2,$$

For example, with technological parameters  $N_A = 5 \times 10^{17} \text{ cm}^{-3}$ ,  $T = 300 \text{ K}$ ,  $Q_i = 10^{10} \text{ q/cm}^2$ ,  $t_{\text{ox}} = 10 \text{ nm}$ ,  $\mu_n = 300 \text{ cm}^2/\text{V}$ ,  $W/L = 1$  and  $\Phi_F^G = -0.56 \text{ V}$  the calculated parameters that determine the output characteristics are  $\phi_S = 2\Phi_F = 0.896 \text{ V}$ ,  $V_{\text{FB}} = -1.01 \text{ V}$ ,  $V_t = 1.00 \text{ V}$  and  $\gamma = 1.18 \text{ V}^{1/2}$ . In Fig. 8.23 a plot of the resulting output characteristics is shown.



Jan. 28, 1930. J. E. LILIENTELD 1,745,175  
 METHOD AND APPARATUS FOR CONTROLLING ELECTRIC CURRENTS  
 Filed Oct. 8, 1926



**Figure 8.22:** Schematic of a *Method and Apparatus for Controlling Electric Currents*, proposed by Lilienfeld in 1926.

Even though the output characteristics reflect the physical device to a certain extent, the model still yields infinite output conductance. In real MOSFETs, this is not the case, because the channel is pinched off at  $x = L - l_d$ , where  $l_d$  depends on  $V_{DS}$  (*channel length modulation*). The drain saturation current now becomes

$$I_D^{\text{sat}} = \frac{W}{L - l_d} \mu_n C_{\text{diel}} \frac{1}{2} (V_{GS} - V_t)^2.$$

Therefore, with the approximation

$$I_D = I_D^{\text{sat}} \frac{L}{L - l_d},$$

a linearization (recall  $1/(1 - x) \approx 1 + x$  for small  $x$ ) yields

$$I_D \approx \left(1 + \frac{l_d}{L}\right) I_D^{\text{sat}} \approx (1 + \lambda V_{DS}) I_D^{\text{sat}},$$

where the *channel length parameter*  $\lambda$  [ $V^{-1}$ ] was introduced. This approximation now results in a finite output conductance.

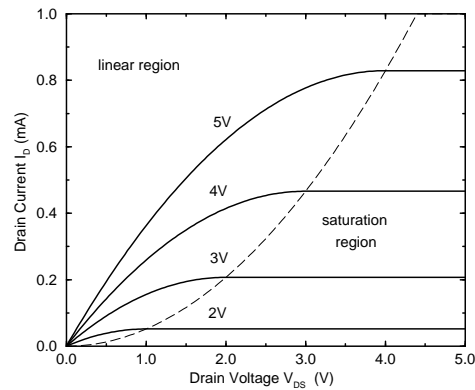
In Fig. 8.24, simulation results with MINIMOS-NT for a MOSFETs are shown in two different scales. In the subthreshold region, where  $V_{GS} < V_t$ , the *subthreshold slope* is

$$S = \frac{dV_{GS}}{d(\log I_D)} \leq \ln(10) \frac{k_B T}{q}.$$

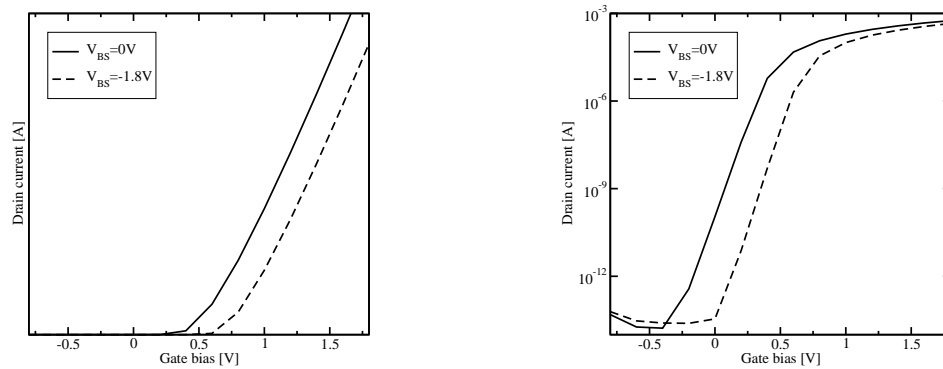
This raises the question of a proper definition of the *threshold current*, so that values for the threshold voltage are comparable. A practical definition (or convention) is

$$I_D^{\text{th}} = 0.1 \mu\text{A} \frac{W}{L}.$$

Finally, we must not skip the *body effect*: The threshold voltage  $V_t$  depends on the substrate bias  $V_{SB}$ , but we are not going into further details.



**Figure 8.23:** Characteristics of a MOSFET in the linear and in the saturation region.



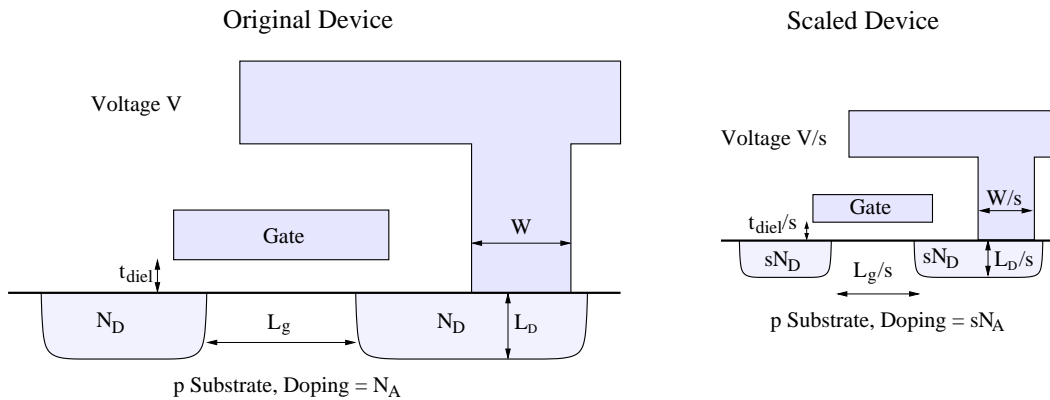
**Figure 8.24:** Transfer characteristics of a MOSFET simulated with MINIMOS-NT and plotted in linear scale (left) and logarithmic scale (right) for the drain current.

## 8.6 CMOS Design Issues

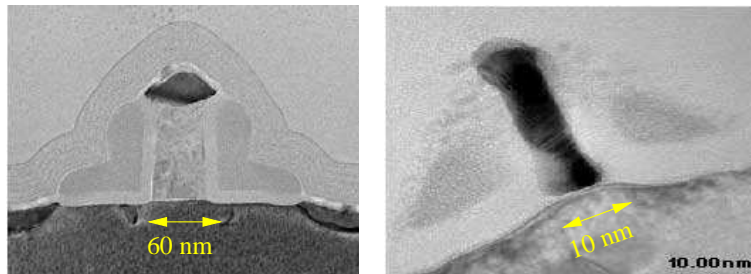
Current MOSFETs are scaled down to gate lengths of several nanometers only. In 1974, Denard analyzed the effects of scaling on physical and electrical quantities in devices. We will have a look at constant-field scaling, where design parameters are chosen such that the electrical field within the device remains the same.

In Fig. 8.25, the effects of (*constant-field*) scaling of a device by a factor  $s > 1$  are demonstrated. Some problems and deviations from the realizations in practice are as follows:

- Within one product family, supply voltages remain constant for reasons of compatibility with other peripheral devices. Thus, the electric field inside a device actually increases even for a constant-field scaling.
- In modern MOSFETs, the gate oxide thickness approaches one nanometer, which corresponds to a few atomic layers only. The atomic structure of matter comes into play: It is not possible to decrease the oxide thickness by a constant factor anymore. Instead,



**Figure 8.25:** Effect of constant-field scaling on device dimensions.



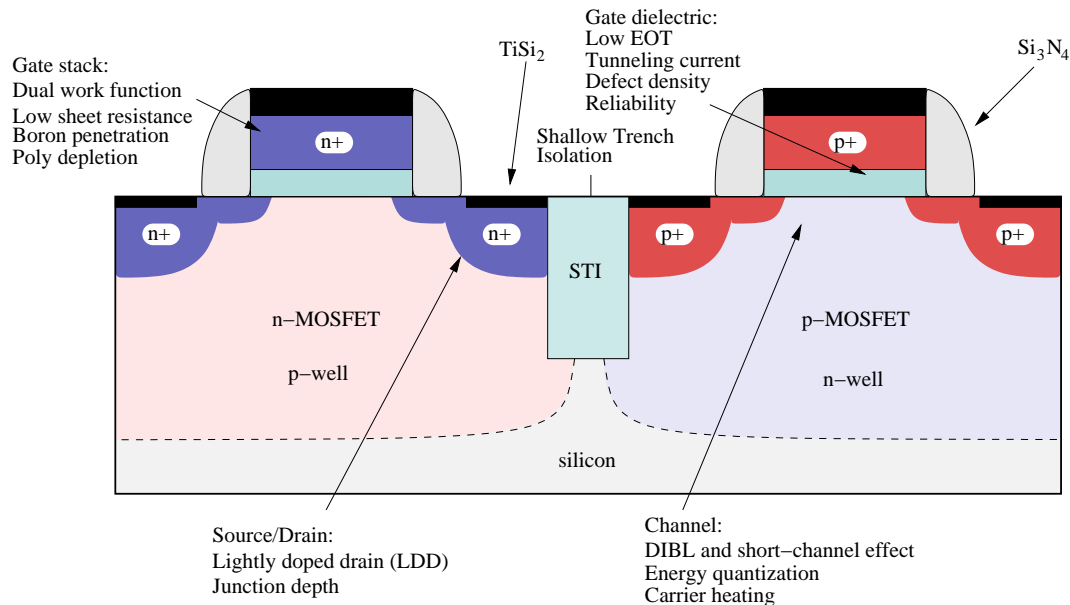
**Figure 8.26:** Two MOSFETs (Intel) with 60nm (left) and 10nm (right) gate length.

there are only discrete multiples of the width of one atomic layer possible. Additionally, quantum mechanical effects have to be taken into account.

- A decrease of the diameter of metallic interconnects leads to an increase of the current densities in the interconnect. High currents may destroy the metallic lattice of the interconnect over time, leading to a breakdown of the circuit. This effect is known as *electromigration*.
- Doping concentrations face natural limits: The impurities have to be incorporated into the semiconductor lattice. If their concentration is too high, this is not possible anymore. For instance, it is difficult in practice to obtain electrically active concentrations larger than  $e \times 10^{20} \text{ cm}^{-3}$ .

In Fig. 8.26, pictures of real MOSFETs from Intel are shown. In contrast to clearly-defined rectangular domains in a two-dimensional cross-section of a MOSFET, the real device does not have such clearly defined sharp transitions between different materials. One has to keep in mind that such small structures are currently fabricated using ultra-violet (UV) light with a wavelength of 193nm (extreme UV light with smaller wavelengths is in preparation, but not technically mature yet), thus posing highest quality requirements (and costs) on the fabrication equipment.

For a standard CMOS inverter, Fig. 8.27 shows some of the requirements and problems that

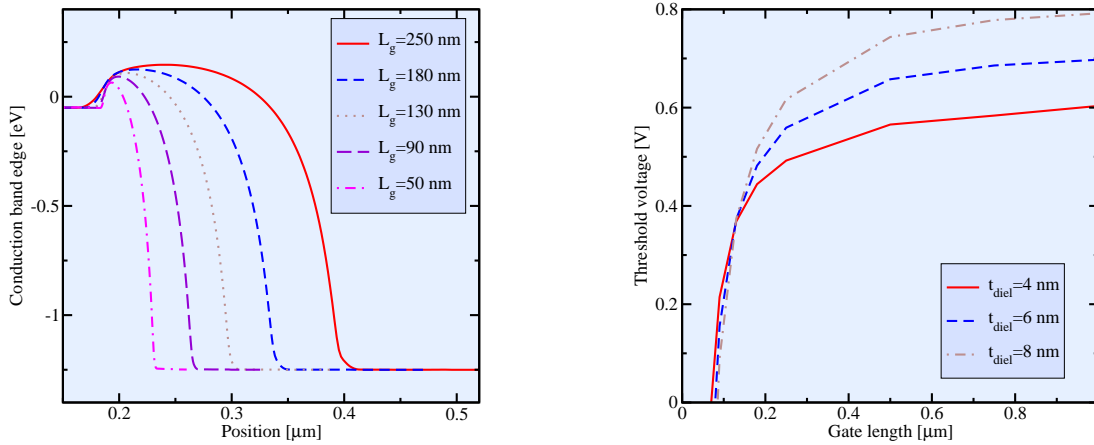


**Figure 8.27:** Requirements and problems with CMOS scaling.

arise from scaling. We will list some of the electrical and fabrication issues in the following:

- *Shallow trench isolation:* As devices are packed closer to each other, the surface leakage current between adjacent wells increases dramatically. This effect can be suppressed (though not completely eliminated) by fabricating shallow trenches filled with a dielectric between wells.
- *Threshold voltage roll-off:* The threshold voltage of a MOSFET is (besides other effects) determined by the depletion charge  $Q_d$  as given in (8.27). The zone of the gate-controlled depletion charge is determined by the gate length and has a trapezoidal shape due to the **encroachment**<sup>1</sup> of the depletion regions from the source and drain reversed-bias junctions into the depletion zone created by the gate electrode. If the channel is long, this encroachment can be neglected. However, for short channels, the gate charge is considerably reduced and the threshold voltage changes (cf. Fig. 8.28). The problem associated with the short-channel effect is not that devices with different channel lengths have different threshold voltages, rather, the problem is that in short devices small statistical variations in the gate length give rise to larger statistical variations of the threshold voltage, which poses a reproducibility problem.
- *Drain-induced barrier lowering (DIBL):* When the drain voltage of a turned-on short-channel MOSFET increases from the linear region toward the saturation region, its threshold voltage roll-off becomes larger (cf. Fig. 8.28, where the conduction band edge is shown on the left). In a turned-on MOSFET, carriers face an energy barrier between source and drain that is controlled by  $V_{GS}$ . In short-channel MOSFETs, the barrier height is additionally influenced by  $V_{DS}$ : In the saturation region, the depletion-layer width of the reverse-biased  $pn$ -junctions increases and reduces the effective potential barrier width. In long-channel devices, this reduction of the barrier width is not significant, but at short channels, the

<sup>1</sup> **encroachment** [ɪn'krəʊtʃ.mənt]: Beeinträchtigung



**Figure 8.28:** Drain induced barrier lowering (left) and roll-off-curve resulting from short-channel  $n$ -MOSFETS (right).

maximum barrier height is also reduced and leads to a substantial increase in electron injection from the source to the drain. As a result, the subthreshold current increases. In other words: An increase of the drain voltage leads to a decrease of the threshold-voltage. Therefore, the threshold voltage increases rapidly for small gate lengths, so that the MOSFETs cannot be turned off properly anymore.

- *Bulk punch-through:* DIBL causes the formation of a leakage path at the SiO<sub>2</sub>/Si interface. If the drain voltage is large enough, significant leakage current may also flow from drain to source via the bulk of the substrate in a short-channel MOSFET due to the increasing depletion-layer width at the drain junction. Consequently, the gate can no longer turn the device completely off and loses control of the drain current at high drain voltages. Even worse, when two adjacent depletion regions touch, *punchthrough* occurs, which will destroy the device. To overcome this problem, anti-punchthrough implants (*retrograde wall, pocket, halo*) exist.

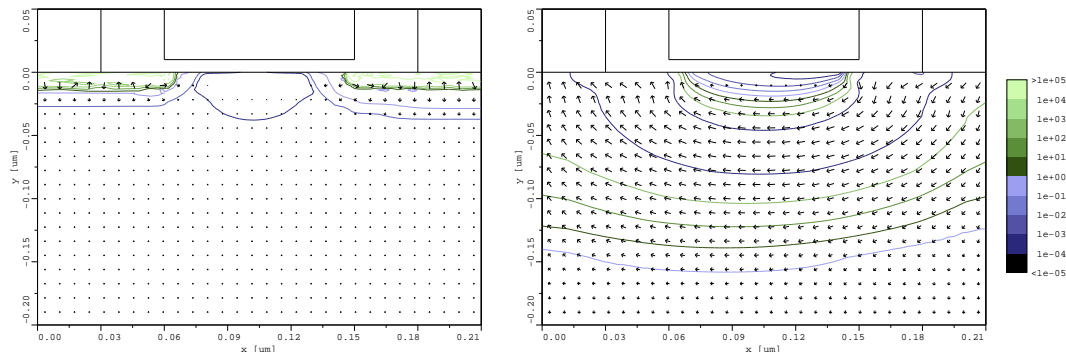
In CMOS inverters with small size, the problem of *polydepletion* comes into play: Polysilicon gates allow the adjustment of work functions and hence threshold voltages by doping. However, a depletion layer at the interface to the dielectric can show up at low electron concentrations (Fig. 8.30), resulting in a voltage drop

$$V_{\text{poly}} \approx \frac{\epsilon_{\text{diel}}^2 E_{\text{diel}}^2}{2q\epsilon_{\text{Si}} N_{\text{poly}}},$$

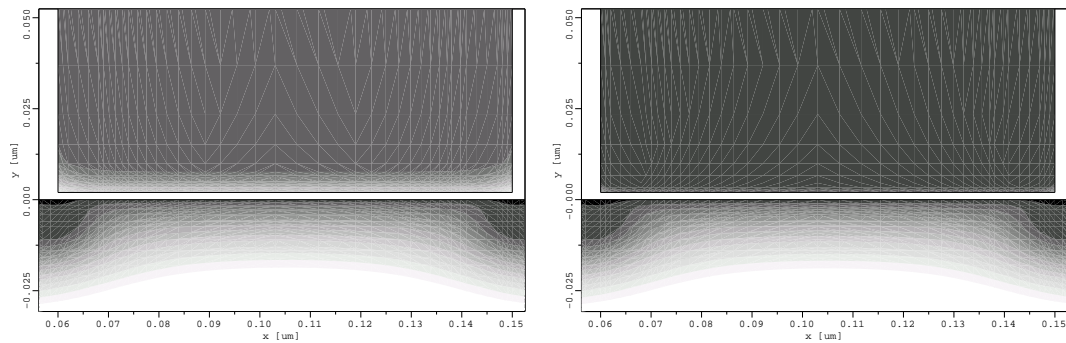
This voltage drop leads to an effective increase of the dielectric thickness, which is not at all desired. However, the problem can be avoided by the use of metal gates, which, unfortunately, introduce a **plethora**<sup>1</sup> of additional problems into the process.

Another quantum-mechanical issue is *carrier quantization*. The distribution of carriers in the channel taking quantum mechanics into account differs compared to classical models. In a classical view, the peak concentration of carriers is directly at the interface to the gate-dielectric,

<sup>1</sup> **plethora** [ˈplɛθ.ə.ə]: Fülle, Vielzahl



**Figure 8.29:** Current density in a turned-off MOSFET with (left) and without (right) retrograde wall. In the latter case the gate cannot control the current flow anymore, thus the device cannot be switched off properly.



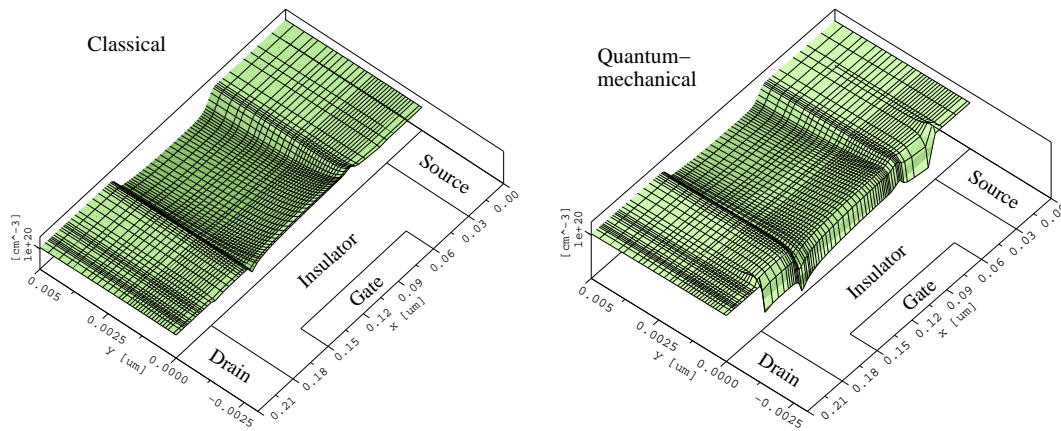
**Figure 8.30:** Electron concentrations for  $N_{\text{poly}} = 10^{19} \text{ cm}^{-3}$  (left) and  $N_{\text{poly}} = 10^{20} \text{ cm}^{-3}$ . In the former case, polydepletion in the gate occurs.

but quantum-mechanics **prohibits**<sup>1</sup> a peak at the interface, so that the concentration maximum is in fact inside the channel (Fig. 8.31). This modifies the output-characteristics (cf. Fig. 8.32)

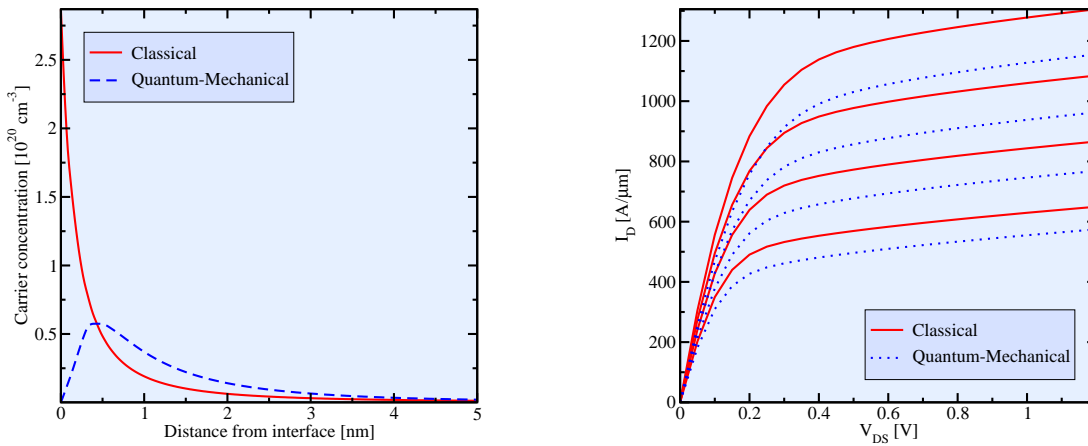
Single- and double-gate SOI (silicon on insulator) MOSFETs offer a superior control of the inversion charge. In single-gate SOI MOSFETs, the undoped channels offer high mobility (recall that impurity scattering at dopants decreases the mobility!), whereas the insulator at the bottom prevents punchthrough. Double-gate SOI result in a symmetric concentration profile of the inversion charge thanks to their symmetric layout with a gate electrode on both ends of the channel. However, SOI devices behave worse when it comes to heat conductance, because the insulator typically has a much lower heat conductivity than the bulk silicon.

The scaling issues discussed in this section are only a selection of the challenges. For more details, the reader is referred to the literature [?, ?].

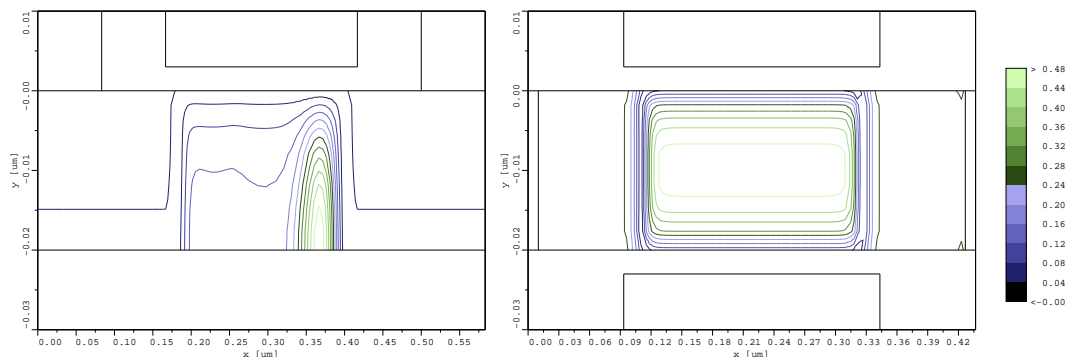
<sup>1</sup> to prohibit [prəˈhɪb.ɪt]: verhindern, unterbinden



**Figure 8.31:** Carrier quantization is a quantum-mechanical effect and results in a peak concentration in the interior of the channel.



**Figure 8.32:** Carrier quantization and its effect on output-characteristics.



**Figure 8.33:** Inversion charge in a single-gate (left) and a double-gate (right) SOI MOSFET, illustrating the superior control of the inversion charge compared to traditional bulk-MOSFETs.

# Appendix A

## Partial Differential Equations

As we have seen in the introductory chapter, semiconductors can be described by a nonlinear set of coupled partial differential equations (PDEs), the drift-diffusion model given in (1.23), (1.24) and (1.25). In contrast to systems of *linear* equations, a more sophisticated solution procedure is required here, which must be carefully chosen depending on the type of the underlying PDEs. Governed by the structure of these equations, **archetypical**<sup>1</sup> behavior of the field quantities can be observed.

### A.1 Boundary and Initial Conditions

In this lecture we will be dealing mainly with *partial differential equations* of the form

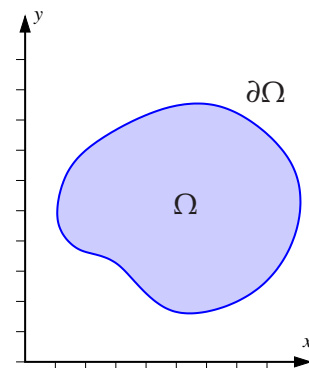
$$F(u, \partial u / \partial x_1, \partial u / \partial x_2, \partial u / \partial x_3, \dots) = G(x) .$$

It is a difficult topic to prove the existence of a solution of equations of this type. Similarly to ordinary differential equations, there are typically degrees of freedom in the solution of partial differential equations (provided it exists). To narrow the solution down, boundary and initial conditions must be provided. For the purpose of device simulation, usually a finite problem domain is assumed. Denoting the problem domain with  $\Omega$  and its boundary with  $\partial\Omega$ , either the values of  $u$  at  $\partial\Omega$  or the values of the normal derivative of  $u$  at  $\partial\Omega$  may be given for a second order partial differential equation.

In the first case, the boundary condition is termed *Dirichlet condition*, and the whole problem is formulated as

“For a given problem space  $\Omega$  and a given function  $\phi$ , find a suitable function  $u$  that fulfills  $F(u, \partial u / \partial x_1, \partial u / \partial x_2, \partial u / \partial x_3, \dots) = G(x)$  in the interior of  $\Omega$  with  $u(x) = \phi(x)$  on the boundary  $\partial\Omega$ .”

Dirichlet conditions are quite intuitive when considering for example Poisson’s equation  $\nabla^2\psi = \rho/\epsilon$ : The function  $u$  is then identified with the electrostatic potential  $\psi$ , the boundary  $\partial\Omega$  is formed by electrodes, and the function  $\phi$  is the voltage at these electrodes.



**Figure A.1:** A region  $\Omega$  and its boundary  $\partial\Omega$ .

<sup>1</sup> **archetypical** [ɑ:krɪ'taɪ.p.ə]: urbildlich, musterhaft, typisch



The second important type are *Neumann conditions*, where the change of the field in the direction perpendicular to the boundary is given,

$$\frac{\partial u(\mathbf{x})}{\partial \mathbf{n}} := \mathbf{n} \cdot \nabla u(\mathbf{x}) = \phi(\mathbf{x}) \quad \text{for } \mathbf{x} \in \partial\Omega,$$

with the unit vector  $\mathbf{n}$  locally perpendicular to  $\partial\Omega$ . Continuing the example from above, where  $u$  is the electrostatic potential  $\psi$ , prescribing the normal projection of  $\nabla u$  effectively prescribes the electric *flux* coming out of  $\Omega$  (remember,  $\nabla u$  reflects the electric field, which is proportional to the flux). In most cases,  $\partial u / \partial \mathbf{n} = 0$ , i.e. *no flux* out of  $\Omega$  at all, is prescribed. This is, however, in many cases true for infinitely large simulation domains only, while usually only finite simulation domains are realized numerically. As a **remedy**<sup>1</sup>, a sufficiently large simulation domain is chosen with (at least to some extent) artificial homogeneous Neumann boundary conditions.

Be aware that it is not sufficient for Laplace's (and Poisson's) equation to prescribe Neumann boundary conditions only: For a solution  $u_0(\mathbf{x})$  of the system

$$\begin{aligned} \nabla^2 u &= 0, & \text{in } \Omega, \\ \frac{\partial u(\mathbf{x})}{\partial \mathbf{n}} &= \phi(\mathbf{x}), & \text{on } \partial\Omega, \end{aligned}$$

every function  $u(\mathbf{x}) = u_0(\mathbf{x}) + C$  with constant  $C$  is also a solution. Therefore, in order to assure uniqueness of a solution of Laplace's (or Poisson's) equation, Dirichlet boundary conditions on at least part of the boundary  $\partial\Omega$  must be given.

As indicated, it is also possible to have the boundary  $\partial\Omega$  decomposed into several (**disjoint**)<sup>2</sup> parts  $\Gamma_1, \Gamma_2, \dots, \Gamma_r$  and impose boundary conditions of different types on each part:

$$\begin{aligned} u(\mathbf{x}) &= g_1(\mathbf{x}) & \text{on } \Gamma_1, \\ \frac{\partial u(\mathbf{x})}{\partial \mathbf{n}} &= g_2(\mathbf{x}) & \text{on } \Gamma_2, \\ & \vdots \\ u(\mathbf{x}) &= g_r(\mathbf{x}) & \text{on } \Gamma_r. \end{aligned}$$

Actually, this is the regular case; think of a simple parallel-plate capacitor: The electrodes, of course, form a Dirichlet boundary, while at the remaining boundaries the potential can not possibly be fixed. Instead, it is assumed that no electric flux leaves the problem space at these boundaries, an approximation which has less impact the farther the boundaries are away from the electrodes. Such a combination of boundary conditions is called *mixed boundary conditions*.

There is a third type of boundary conditions possible, called *Robin conditions*. They can be considered a linear combination of Dirichlet and Neumann boundary conditions:

$$\alpha(\mathbf{x})u(\mathbf{x}) + \beta(\mathbf{x})\frac{\partial u(\mathbf{x})}{\partial \mathbf{n}} = \gamma(\mathbf{x}) \quad \text{on } \partial\Omega.$$

Depending on the problem formulation and the chosen solution method, it can be quite difficult to incorporate this type of boundary conditions into numerical simulations.

If the region  $\Omega$  is unbounded (e.g.  $\Omega = \mathbb{R}^3$ ), other conditions must be given, e.g. the value of  $u(\mathbf{x})$  for  $|\mathbf{x}| \rightarrow \infty$ . But these unbounded problem spaces pose another problem when tackled by numeric simulation: How should an infinite region be represented in a finite computer

<sup>1</sup> **remedy** [rem.ə.di]: Abhilfe, Behelf    <sup>2</sup> **disjoint** [dis'dʒɔɪn.t]: disjunkt, einander nicht überschneidend

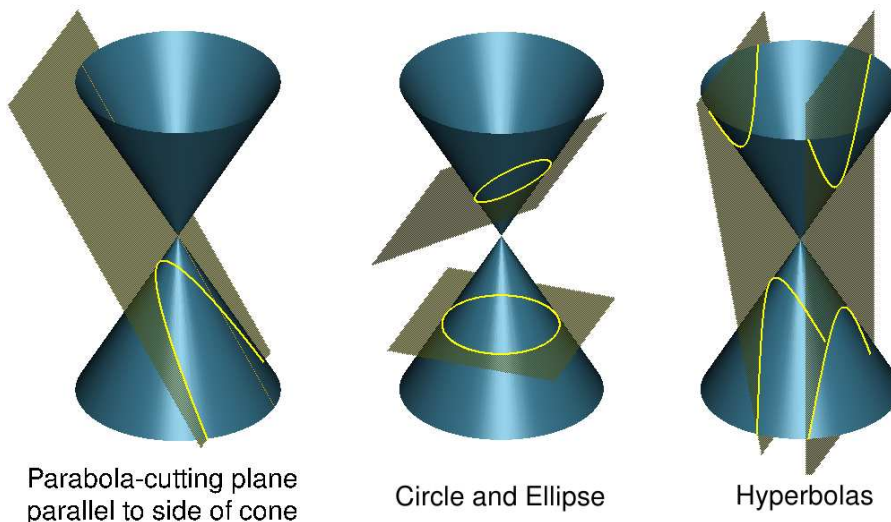


Figure A.2: Conic sections [?].

memory? Therefore, this case will not be discussed in this lecture (although, of course, there are methods around to solve even this kind of problem).

Note that if the stated problem is non-stationary (time-dependent), the time  $t$  is one of the independent variables  $x_i$ . In this case, the boundary condition is also termed *initial condition*. But from a mathematical point of view it is in no way special—time is a variable just as the spatial coordinates are.

## A.2 Classification

Every linear, second order partial differential equation with constant coefficients belongs to one of three groups, depending on the coefficients in the equation. The members of these groups share a lot of common properties and have the same **peculiarities**<sup>1</sup> when it comes to discretization and numerical solution.

Specializing the general equation scheme from the beginning of this chapter to a linear, second order type in two variables  $x$  and  $y$  we have

$$a \frac{\partial^2 u}{\partial x^2} + 2b \frac{\partial^2 u}{\partial x \partial y} + c \frac{\partial^2 u}{\partial y^2} + d \frac{\partial u}{\partial x} + e \frac{\partial u}{\partial y} + fu = G, \quad (\text{A.1})$$

where  $a$  through  $f$  are constants and  $G$  is a function in  $x$  and  $y$ . Now, every differentiation with respect to  $x$  is formally replaced by a multiplication with  $\alpha$ , every  $\partial/\partial y$  by  $\beta$ , yielding the **bivariate**<sup>2</sup> polynomial

$$P(\alpha, \beta) = a\alpha^2 + 2b\alpha\beta + c\beta^2 + d\alpha + e\beta + f. \quad (\text{A.2})$$

The equation  $P(\alpha, \beta) = 0$  describes *conic sections*. Which conic section it describes is determined by the **discriminant**<sup>3</sup>  $ac - b^2$ : For  $ac - b^2 > 0$  the equation describes an **ellipse**<sup>4</sup>, for  $ac - b^2 = 0$

<sup>1</sup> **peculiarity** [pɪˌkjuː.liˈer.ə.ti]: Eigenheit, Ausprägung    <sup>2</sup> **bivariate** [baɪvəriːt]: von zwei Variablen abhängig  
<sup>3</sup> **discriminant** [dɪskrɪmɪnənt]: Diskriminante    <sup>4</sup> **ellipse** [ɪˈlɪps]: Ellipse

a **parabola**<sup>1</sup> and for  $ac - b^2 < 0$  a **hyperbola**<sup>2</sup>. For partial differential equations in three or more independent variables, the above scheme can be generalized analogously; the derivatives with respect to the different spatial coordinates represent the variables  $\alpha$ ,  $\beta$  and so on in our polynomial.

Moreover, if the coefficients  $a$  through  $f$  are allowed to be functions of  $x$  and  $y$  instead of constants, the type of the partial differential equation may change over the domain  $\Omega$ , leading to additional numerical subtleties.

- *Elliptic Partial Differential Equations:* The elliptic type is the ‘nicest’ type in terms of sensitivity to discretization errors and numerical stability. A typical example is the Poisson equation  $\nabla^2 u = \phi$ . In two dimensions it has the form  $\partial^2 u / \partial x^2 + \partial^2 u / \partial y^2 = \phi$ , giving the characteristic equation’s coefficients  $a = c = 1$  and  $b = d = e = f = 0$ , which results in a discriminant greater than zero. Elliptic differential equations describe stationary processes and hence do not have time derivatives.
- *Parabolic Partial Differential Equations:* The diffusion equation  $\partial u / \partial t - \partial^2 u / \partial x^2 = \phi$  is a typical example for a parabolic differential equation. It contains a first order derivative with respect to time and a second order derivative with respect to the spatial coordinate, giving  $a = b = e = f = 0$ ,  $c = -1$ , and  $d = 1$ , which results in a discriminant equal to zero. The combined drift-diffusion and continuity equations (1.24) and (1.25) on page 5 as well as the heat-flow equation (1.26) are examples for parabolic differential equation.
- *Hyperbolic Partial Differential Equations:* Describing **propagating**<sup>3</sup> ‘distortions’ like waves, the hyperbolic type is the most cumbersome for numerical simulation. First of all, in contrast to the parabolic type, where discontinuous distortions are smoothed out over time, discontinuities (‘shocks’) propagate *as they are* through a hyperbolic system. Secondly, the fact that propagation takes place lets us observe a preferred direction of propagation in the system, making the discretization more complicated. (Imagine a plane wave propagating in a medium. It obviously does not make sense to discretize the space along the direction of propagation in the same way as perpendicular to the direction of propagation.) Unfortunately, the hyperbolic differential equations associated with semiconductors may change their direction of propagation, requiring an adaptive discretization scheme.

The homogeneous wave equation  $\partial^2 u / \partial t^2 - \partial^2 u / \partial x^2 = 0$ , with  $a = 1$ ,  $b = d = e = f = 0$ , and  $c = -1$  results in a discriminant of  $-1$ , and is therefore classified as a hyperbolic differential equation.

<sup>1</sup> **parabola** [pəˈræb.ə.lə]: Parabel    <sup>2</sup> **hyperbola** [haɪˈpɛɪ.bə.lə]: Hyperbel    <sup>3</sup> **to propagate** [prɒp.ə.ɡeɪt]: ausbreiten, fortpflanzen

## Appendix B

# Vector Analysis and Its Implementation in SGFramework

One quantity of interest in the simulation of a semiconductor device is the electromagnetic field inside the device. This field is **governed**<sup>1</sup> by Maxwell's equations, which are coupled with the semiconductor equations. Since the mathematical representation of scalar and vectorial fields, i.e. mappings from  $\mathbb{R}^3 \mapsto \mathbb{R}$  and  $\mathbb{R}^3 \mapsto \mathbb{R}^3$ , respectively, a few things you actually should remember from your mathematics- or electrodynamics lectures are summarized in this section.

### B.1 Divergence

The *divergence* of a vector field  $\boldsymbol{v}(\boldsymbol{r})$  is a scalar field, which, if  $\boldsymbol{v}$  is interpreted as a physical flow, is a measure for the source density of this flow. If

$$\boldsymbol{v}(\boldsymbol{r}) = \begin{pmatrix} v_x(\boldsymbol{r}) \\ v_y(\boldsymbol{r}) \\ v_z(\boldsymbol{r}) \end{pmatrix}, \quad \boldsymbol{r} = \begin{pmatrix} x \\ y \\ z \end{pmatrix}, \quad (\text{B.1})$$

then in Cartesian coordinates in three dimensions,

$$\text{div } \boldsymbol{v}(\boldsymbol{r}) = \frac{\partial v_x(\boldsymbol{r})}{\partial x} + \frac{\partial v_y(\boldsymbol{r})}{\partial y} + \frac{\partial v_z(\boldsymbol{r})}{\partial z}. \quad (\text{B.2})$$

An example is  $\text{div } \boldsymbol{D} = \rho$ , where  $\boldsymbol{D}$  is the electric flux density (electric displacement field) and  $\rho$  is the charge density.

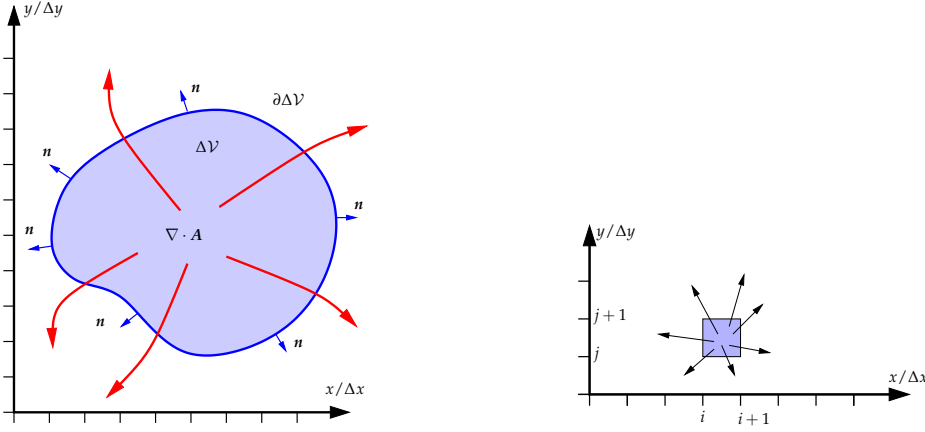
The physical meaning of the divergence can be understood as follows: Imagine a spatial point  $\boldsymbol{r}$  and a small volume  $\Delta\mathcal{V}$  (with surface  $\partial\Delta\mathcal{V}$ ) around  $\boldsymbol{r}$ . Moreover, imagine a vector field  $\boldsymbol{A}$ , defined at least in the point  $\boldsymbol{r}$  and its vicinity. Then, the *flux*  $F$  of  $\boldsymbol{A}$  leaving  $\partial\Delta\mathcal{V}$  is

$$F = \int_{\partial\Delta\mathcal{V}} \boldsymbol{A} \cdot d\boldsymbol{S} = \int_{\partial\Delta\mathcal{V}} \boldsymbol{A} \cdot \boldsymbol{n} \, dS. \quad (\text{B.3})$$

Many semiconductor devices can be characterized by a two-dimensional layout, for example a MOSFET. This allows a two-dimensional simulation such that the quantities of interest in the real device are found by scaling along the third axis, because the variation along the third axis

---

<sup>1</sup> to be governed [gAV.°nd]: bestimmt sein



**Figure B.1:** Flux of vector field  $\mathbf{A}$  through the surface  $\partial\Delta\mathcal{V}$  (left) and its box discretization (right).

is small. In mathematical terms, we will compute an averaged flux  $\bar{F}$  in two dimensions, which is linked to the true flux  $F$  by a multiplication with the width  $w$  along the third dimension:

$$\bar{F} := \frac{F}{w} = \int_{\partial\Delta\mathcal{V}_{2d}} \mathbf{A} \cdot \mathbf{n} \, ds . \quad (\text{B.4})$$

In order to discretize the divergence operator in two dimensions, we assume a box between  $(i\Delta x, j\Delta y)$  and  $((i+1)\Delta x, (j+1)\Delta y)$  (cf. Fig. B.1). With  $\mathbf{A} = (A_x, A_y)^T$ , the fluxes leaving the top and the bottom of the box are

$$\bar{F}^{\text{top}} = \int_{i\Delta x}^{(i+1)\Delta x} A_y(x, y = (j+1)\Delta y) \, dx , \quad \bar{F}^{\text{bottom}} = - \int_{i\Delta x}^{(i+1)\Delta x} A_y(x, y = j\Delta y) \, dx . \quad (\text{B.5})$$

The net flux leaving in  $y$ -direction thus is

$$\bar{F}^y = \bar{F}^{\text{top}} + \bar{F}^{\text{bottom}} = \int_{i\Delta x}^{(i+1)\Delta x} (A_y(x, y = (j+1)\Delta y) - A_y(x, y = j\Delta y)) \, dx . \quad (\text{B.6})$$

For sufficiently small  $\Delta y$  we can approximate

$$A_y(x, y = (j+1)\Delta y) \approx A_y(x, y = j\Delta y) + \frac{\partial A_y}{\partial y} \Delta y . \quad (\text{B.7})$$

Inserting (B.7) into (B.6) yields

$$\begin{aligned} \bar{F}^y &= \int_{i\Delta x}^{(i+1)\Delta x} (A_y(x, y = (j+1)\Delta y) - A_y(x, y = j\Delta y)) \, dx \\ &\approx \int_{i\Delta x}^{(i+1)\Delta x} \frac{\partial A_y}{\partial y} \Delta y \, dx \\ &\approx \frac{\partial A_y}{\partial y} \Delta x \Delta y . \end{aligned} \quad (\text{B.8})$$

Approaching the  $x$ -direction analogously,

$$\bar{F}^x \approx \frac{\partial A_x}{\partial x} \Delta x \Delta y, \quad (\text{B.9})$$

and using  $\Delta V = w \Delta x \Delta y$  we finally get

$$\frac{F^x + F^y}{\Delta V} = \frac{\bar{F}^x + \bar{F}^y}{\Delta x \Delta y} = \frac{\partial A_x}{\partial x} + \frac{\partial A_y}{\partial y} = \left( \begin{array}{c} \frac{\partial}{\partial x} \\ \frac{\partial}{\partial y} \end{array} \right) \cdot \left( \begin{array}{c} A_x \\ A_y \end{array} \right) = \nabla \cdot \mathbf{A}. \quad (\text{B.10})$$

We have just confirmed Gauss' theorem, which states the boundary integral can be translated into an integral over the domain:

$$F = \int_{\partial \Delta V} \mathbf{A} \cdot \mathbf{n} \, dS = \int_{\Delta V} \text{div } \mathbf{A} \, dV. \quad (\text{B.11})$$

Thus,  $\text{div } \mathbf{A}$  can be interpreted as a sink/source density. The above equation can be recast in a more physical setting as follows: The net number of particles leaving the volume  $\Delta V$  has to be equal to the number of particles generated within  $\Delta V$  to obtain an equilibrium situation.

As an example, for  $\mathbf{A} = (3x - 0.2y^2)\mathbf{e}_x + (y + 4y^2)\mathbf{e}_y$  we get  $\nabla \cdot \mathbf{A} = 4 + 8y$ , displayed in the following program and visualized in Figure B.2.

```

1  const DIM = 30;           // number of mesh points in x and y
2  const DX = 0.2;         // mesh spacing in x
3  const DY = 0.1;         // mesh spacing in y
4  var x[DIM], y[DIM], Ax[DIM,DIM], Ay[DIM,DIM],
5      divA[DIM-1,DIM-1];
6
7  begin main
8
9      // initialize the x and y components of vector A
10     assign x [i=all]      = i*DX;
11     assign y [j=all]      = j*DY;
12     assign Ax[i=all, j=all] = 3.0*x[i]-0.2*sq(y[j]);
13     assign Ay[i=all, j=all] = y[j]+4.0*sq(y[j]);
14
15     // compute the divergence and write the results
16     assign divA[i=all, j=all] = (Ax[i+1,j]-Ax[i,j])/DX +
17                                   (Ay[i,j+1]-Ay[i,j])/DY;
18
19     write;
20 end
    
```

source\_code/divergence\_example.sg

## B.2 Curl (Rotation)

Just like the divergence, the curl operator (aka. rot operator) is applied to vector fields; but it results in a vector field whose magnitude is a measure for the rate of rotation, and whose direction is perpendicular to the rotation plane of the operand.

Unlike the divergence, the curl operator is only meaningful for three-dimensional spaces,<sup>1</sup>

<sup>1</sup> For spaces of dimension  $n > 3$ , the framework of *alternating differential forms*, of which the vector calculus presented here is a special case, has to be employed.

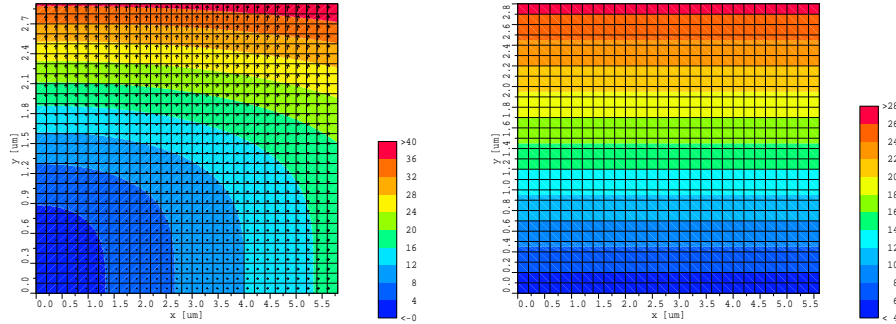


Figure B.2: The vector field  $A$  (left) and its divergence  $\nabla \cdot A$  (right).

where it is defined as

$$\text{curl } v(\mathbf{r}) = \begin{pmatrix} \frac{\partial v_z(\mathbf{r})}{\partial y} - \frac{\partial v_y(\mathbf{r})}{\partial z} \\ \frac{\partial v_x(\mathbf{r})}{\partial z} - \frac{\partial v_z(\mathbf{r})}{\partial x} \\ \frac{\partial v_y(\mathbf{r})}{\partial x} - \frac{\partial v_x(\mathbf{r})}{\partial y} \end{pmatrix}. \quad (\text{B.12})$$

### B.3 Gradient

While the curl and divergence operators act on vector fields, the *gradient* has to be used on scalar fields, returning a vector field that points in the direction of the **steepest**<sup>1</sup> **ascent**<sup>2</sup> of the operand. The definition reads

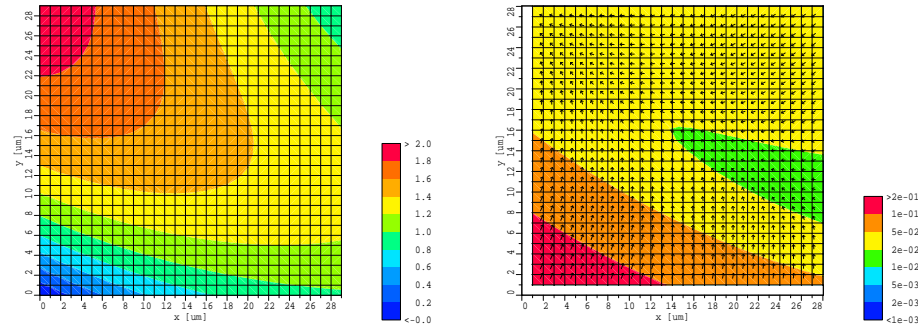
$$\text{grad } f(\mathbf{r}) = \begin{pmatrix} \frac{\partial f(\mathbf{r})}{\partial x} \\ \frac{\partial f(\mathbf{r})}{\partial y} \\ \frac{\partial f(\mathbf{r})}{\partial z} \end{pmatrix}. \quad (\text{B.13})$$

Program `exgrad.sg` gives a short example of how (B.13) is implemented in SGFRAMEWORK, while Fig. B.3 displays the result.

```

1  const DIM = 30;           // number of mesh points in x and y
2  const DX = 1.0;          // mesh spacing in x
3  const DY = 1.0;          // mesh spacing in y
4
5  var x[DIM], y[DIM], phi[DIM,DIM];
6  var gphix[DIM,DIM], gphiy[DIM,DIM], magphi[DIM,DIM];
7
8  const A = 1.0, B = 2.0, X0 = 50.0, Y0 = 50.0;
    
```

<sup>1</sup> **steep** [sti:p]: steil    <sup>2</sup> **ascent** [ə'sent]: Anstieg



**Figure B.3:** Solution of `exgrad.sg`:  $\Phi$  (left),  $\nabla\Phi$  (right). The arrows in the right figure are perpendicular to isolines (i.e. the interfaces between different colors) on the left, showing the direction of steepest ascent.

```

9  begin main
10  assign x[i=all] = i*DX;
11  assign y[j=all] = j*DY;
12  assign phi[i=all, j=all]=
13    {A*[1.0 - sq(x[i]/X0-1.0)]+B*[1.0 - sq(y[j]/Y0-1.0)]} *
14    {sq(x[i]/X0-1.0)+sq(y[j]/Y0-1.0)};
15
16  // determine the x and y components of phi and its magnitude
17  assign gphix[i=1..DIM-2, j=1..DIM-2] = (phi[i+1, j]-phi[i-1, j]) / (2.0*DX);
18  assign gphiy[i=1..DIM-2, j=1..DIM-2] = (phi[i, j+1]-phi[i, j-1]) / (2.0*DY);
19  assign magphi[i=all, j=all] = sqrt(sq(gphix[i, j])+sq(gphiy[i, j]));
20
21  // write the results
22  write;
23  end
    
```

source\_code/grad\_example.sg

## B.4 Nabla

As an aid when manipulating expressions using the differential operators above, the linear differential operator *Nabla* is defined as

$$\nabla = \begin{pmatrix} \frac{\partial}{\partial x} \\ \frac{\partial}{\partial y} \\ \frac{\partial}{\partial z} \end{pmatrix}. \quad (\text{B.14})$$



Using the Nabla operator and the two basic product operations for vectors, the previously discussed operators can simply be written as

$$\text{grad } f(\mathbf{r}) = \nabla f(\mathbf{r}) \ , \quad \text{div } \mathbf{v}(\mathbf{r}) = \nabla \cdot \mathbf{v}(\mathbf{r}) \ , \quad \text{and} \quad \text{curl } \mathbf{v}(\mathbf{r}) = \nabla \times \mathbf{v}(\mathbf{r}) \ . \quad (\text{B.15})$$

The fact that the curl is restricted to three-dimensional spaces can be **attributed**<sup>1</sup> to the fact that the cross product of two vectors is only defined for these spaces.

## B.5 Manipulating Expressions

When dealing with more complicated expressions involving vector calculus operators, the Nabla-concept **comes in handy**<sup>2</sup>. But one has to be careful and keep in mind that Nabla is *both a vector and a differential operator*. From the former it follows that all identities known for vectors can be applied to expressions using Nabla, but the latter tells us that if such an identity involves rearranging the variables following Nabla, the product-rule of differentiation must be obeyed, since Nabla's differentiating nature acts on everything following it!

For example, if we were to simplify  $\text{div}(\mathbf{u} \times \mathbf{v})$  (the dependence on  $\mathbf{r}$  will be suppressed from now on **for the sake of**<sup>3</sup> a **concise**<sup>4</sup> notation), we translate

$$\text{div}(\mathbf{u} \times \mathbf{v}) = \nabla \cdot (\mathbf{u} \times \mathbf{v}) \quad (\text{B.16})$$

and use the fact that the triple product  $\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c})$  allows **cyclic**<sup>5</sup> permutation:

$$\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = \mathbf{c} \cdot (\mathbf{a} \times \mathbf{b}) = \mathbf{b} \cdot (\mathbf{c} \times \mathbf{a}) = -\mathbf{b} \cdot (\mathbf{a} \times \mathbf{c}) = \dots \ . \quad (\text{B.17})$$

Before applying (B.17) to (B.16), however, one has to take care of Nabla's differentiating property, which says that a variable can only be moved out of Nabla's scope, i.e. pulled in front of Nabla, if it is constant (cf. the rule of ordinary differentiation, where  $d(cf) = cd f$  if and only if  $c$  is constant). By writing  $\nabla_u$  and  $\nabla_v$  for Nablas that **act only upon**<sup>6</sup>  $\mathbf{u}$  and  $\mathbf{v}$ , respectively, and therefore treat all other variables as constants, the product rule of differentiation yields

$$\nabla \cdot (\mathbf{u} \times \mathbf{v}) = \nabla_u \cdot (\mathbf{u} \times \mathbf{v}) + \nabla_v \cdot (\mathbf{u} \times \mathbf{v}) \ . \quad (\text{B.18})$$

Now the vector identities from above can be used to rearrange the individual terms,

$$\nabla_u \cdot (\mathbf{u} \times \mathbf{v}) = \mathbf{v} \cdot (\nabla_u \times \mathbf{u}) \ , \quad \nabla_v \cdot (\mathbf{u} \times \mathbf{v}) = -\mathbf{u} \cdot (\nabla_v \times \mathbf{v}) \ , \quad (\text{B.19})$$

leading to the result

$$\text{div}(\mathbf{u} \times \mathbf{v}) = \mathbf{v} \cdot \text{curl } \mathbf{u} - \mathbf{u} \cdot \text{curl } \mathbf{v} \ . \quad (\text{B.20})$$

**Caveat**<sup>7</sup>: This identity is part of the derivation of the lemma of Poynting. You may enjoy guessing what  $\mathbf{u}$  and  $\mathbf{v}$  stand for in this case, and how the expression is manipulated further. Of course, you may also look it up in the appropriate lecture notes or books, or even don't bother thinking about it at all—but that way, you'll miss a lot of fun!

---

<sup>1</sup> **to attribute** [ˈæt.rɪ.bju:t]: zuschreiben, zurückführen auf    <sup>2</sup> **to come in handy** [kʌm ɪn hæŋ.dɪ]: sich gut treffen, gelegen kommen    <sup>3</sup> **for the sake of** [fɔːr ðe seɪk əv]: um ... Willen    <sup>4</sup> **concise** [kən'saɪs]: übersichtlich    <sup>5</sup> **cyclic** [ˈsaɪ.klɪ.k]: zyklisch, periodisch    <sup>6</sup> **act upon** [ækt ə'pɒn]: einwirken    <sup>7</sup> **caveat** [ˈkæv.i.æt]: Hinweis, Warnung

## B.6 Identities

The most important identity in vector calculus involves the curl of a gradient field,

$$\text{curl grad } f = \nabla \times (\nabla f) = \begin{pmatrix} \frac{\partial}{\partial y} \frac{\partial f}{\partial z} - \frac{\partial}{\partial z} \frac{\partial f}{\partial y} \\ \frac{\partial}{\partial z} \frac{\partial f}{\partial x} - \frac{\partial}{\partial x} \frac{\partial f}{\partial z} \\ \frac{\partial}{\partial x} \frac{\partial f}{\partial y} - \frac{\partial}{\partial y} \frac{\partial f}{\partial x} \end{pmatrix}. \quad (\text{B.21})$$

If  $f$  is sufficiently smooth (which will be implicitly assumed, otherwise the whole analysis would be pointless), the lemma of Schwartz states that the order of differentiation in the mixed second partial derivatives is arbitrary. That means that all terms cancel. A gradient field therefore never has a curl component,

$$\text{curl grad } f = \nabla \times (\nabla f) \equiv 0. \quad (\text{B.22})$$

Even more important is the reverse statement: If a vector field is rotation-free, it can always be represented as the gradient of a suitable scalar field, called a *scalar potential*. The use of a scalar potential not only makes life a bit easier (remember, a scalar field is *one* function dependent on three coordinates, while a vector field is *three* functions dependent on three coordinates), but also ensures that the claim of being rotation-free is automatically fulfilled. A scalar potential is never unique, since the gradient of a constant field  $c$  vanishes,  $\text{grad}(f + c) = \text{grad } f + \text{grad } c = \text{grad } f$ , and therefore also  $f + c$  is a valid potential.

A similar observation (with basically the same proof) can be made with the divergence of a curl field,

$$\text{div curl } \mathbf{v} = \nabla \cdot (\nabla \times \mathbf{v}) \equiv 0. \quad (\text{B.23})$$

This identity enables the representation of a field that is free of sources by the curl of a vector field, called its *vector potential*. In terms of complexity, nothing is gained, but again the requirement of being source-free is automatically fulfilled. Note that **analogous**<sup>1</sup> to the scalar potential, a vector potential can always be ‘shifted’: Following (B.22), the identity  $\text{div curl}(\mathbf{v} + \text{grad } f) = \text{div}(\text{curl } \mathbf{v} + \text{curl grad } f) = \text{div curl } \mathbf{v}$  shows that it is possible to add an arbitrary gradient field to a vector potential.

The last identity makes use of the rules stated in the last section. By noting

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = \mathbf{b}(\mathbf{a} \cdot \mathbf{c}) - \mathbf{c}(\mathbf{a} \cdot \mathbf{b}) = \mathbf{b}(\mathbf{a} \cdot \mathbf{c}) - (\mathbf{a} \cdot \mathbf{b})\mathbf{c}, \quad (\text{B.24})$$

the ‘double curl’  $\text{curl curl } \mathbf{v}$  is transformed into

$$\text{curl curl } \mathbf{v} = \nabla \times (\nabla \times \mathbf{v}) = \nabla(\nabla \cdot \mathbf{v}) - (\nabla \cdot \nabla)\mathbf{v} = \nabla(\nabla \cdot \mathbf{v}) - \nabla^2 \mathbf{v} = \text{grad div } \mathbf{v} - \text{div grad } \mathbf{v}. \quad (\text{B.25})$$

This identity can be used to derive the electromagnetic wave equations from Maxwell’s equations. We won’t make use of it, though.

## B.7 Integral Theorems of Stokes and Gauss

The theorems of Stokes and Gauss are used to transform integrals of the divergence or curl of vector fields into integrals of the fields themselves. A particularly interesting point in these transformations is how the integration domain is transformed.

<sup>1</sup> **analogous** [əˈnæl.ə.gəs]: analog, entsprechend, sinngemäß

Consider a surface  $\mathcal{A}$ <sup>1</sup> and a vector field  $v$  defined at least on every point that is a member of  $\mathcal{A}$ . The surface  $\mathcal{A}$  is confined by its boundary curve  $\partial\mathcal{A}$ ; the field  $v$  needs to be defined on the boundary, too. When integrating over surfaces, the differential in the integrand is vector-valued; its magnitude is the area of the infinitesimal area element, and its direction is perpendicular to the plane spanned by the area element. Integrals along a curve also involve vector-valued differentials, but in this case their direction is the tangential to the curve. Denoting the differential area element by  $d\mathbf{A}$  and the differential curve element by  $d\mathbf{s}$ , the theorem of Stokes states that

$$\int_{\mathcal{A}} \text{curl } v \cdot d\mathbf{A} = \int_{\mathcal{A}} (\nabla \times v) \cdot d\mathbf{A} = \int_{\partial\mathcal{A}} v \cdot d\mathbf{s} \quad . \quad (\text{B.26})$$

Similarly, the integral theorem of Gauss states that the volume integral of a divergence is equal to the integral of the original field over the surface of the volume:

$$\int_{\mathcal{V}} \text{div } v \, dV = \int_{\mathcal{V}} \nabla \cdot v \, dV = \int_{\partial\mathcal{V}} v \cdot d\mathbf{A} \quad (\text{B.27})$$

---

<sup>1</sup> Be careful not to mix up the *surface*  $\mathcal{A}$ , which is an abstract geometrical object, and its *surface area*  $A(\mathcal{A})$ , which is a physical quantity, usually measured in square meters! The same applies to *volumes*.

## Appendix C

# Basics of Electromagnetism – Maxwell's Equations

Maxwell's equations describe the *structure* of the electromagnetic field. In differential form they are formulated as

$$\operatorname{curl} \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t} \quad \text{Faraday's Law of Induction} \quad (\text{C.1})$$

$$\operatorname{curl} \mathbf{H} = \mathbf{J} + \frac{\partial \mathbf{D}}{\partial t} \quad \text{Ampère's Circuital Law with Maxwell's Extension} \quad (\text{C.2})$$

$$\operatorname{div} \mathbf{D} = \rho \quad \text{Gauss' Law} \quad (\text{C.3})$$

$$\operatorname{div} \mathbf{B} = 0 \quad \text{Gauss' Law for Magnetism} \quad (\text{C.4})$$

There is another form, the integral form, which some find easier to interpret, because in this form the equations are directly related to physically observable global quantities like charge, voltage and current. By using Stoke's integral theorem, Ampère's law is integrated over a surface  $\mathcal{A}$  fixed in space to yield the magnetic voltage  $\psi_m$

$$\psi_m(\partial\mathcal{A}) = \int_{\partial\mathcal{A}} \mathbf{H} \cdot d\mathbf{s} = \int_{\mathcal{A}} \operatorname{curl} \mathbf{H} \cdot d\mathbf{A} = \int_{\mathcal{A}} \left( \mathbf{J} + \frac{\partial \mathbf{D}}{\partial t} \right) \cdot d\mathbf{A} = I(\mathcal{A}) + \frac{\partial \Psi(\mathcal{A})}{\partial t} . \quad (\text{C.5})$$

The law of induction is transformed into

$$\psi(\partial\mathcal{A}) = \int_{\partial\mathcal{A}} \mathbf{E} \cdot d\mathbf{s} = \int_{\mathcal{A}} \operatorname{curl} \mathbf{E} \cdot d\mathbf{A} = - \int_{\mathcal{A}} \frac{\partial \mathbf{B}}{\partial t} \cdot d\mathbf{A} = \frac{\partial \Phi(\mathcal{A})}{\partial t} . \quad (\text{C.6})$$

The two divergence equations thus read

$$\Psi(\partial\mathcal{V}) = \int_{\partial\mathcal{V}} \mathbf{D} \cdot d\mathbf{A} = \int_{\mathcal{V}} \operatorname{div} \mathbf{D} \, dV = \int_{\mathcal{V}} \rho \, dV = Q(\mathcal{V}) \quad (\text{C.7})$$

and

$$\Phi(\partial\mathcal{V}) = \int_{\partial\mathcal{V}} \mathbf{B} \cdot d\mathbf{A} = \int_{\mathcal{V}} \operatorname{div} \mathbf{B} \, dV = 0 . \quad (\text{C.8})$$

Maxwell's equations introduce six quantities: A flux density and field strength for the electric and magnetic part of the field, each, an electrical current density, and an electrical charge density. Therefore, to **unambiguously**<sup>1</sup> solve a problem a total of six equations is required. But

<sup>1</sup> **unambiguously** [ʌn.æm'big.ju.ə.sli]: eindeutig

Maxwell himself only provided four equations, and they are only ‘half-equations’ in the sense that they only assert properties of either the divergence- or curl part of the respective field, so we are in need for another four relations. Moreover, we have not considered the presence of matter so far, which is composed of a large number of charged particles that may either be bound to the matter’s crystal structure or may possibly move around freely. All charge carriers, bound or free, are subject to the Lorentz force

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}) , \quad (\text{C.9})$$

where  $\mathbf{v}$  is the velocity of the individual particle (referenced to the same coordinate system where  $\mathbf{E}$  and  $\mathbf{B}$  are referenced to). Depending on their bond state, they contribute to different quantities of the electromagnetic field.

The bound charge particles that **constitute**<sup>1</sup> matter (the protons and the electrons in the inner hulls) may be ‘electromagnetically active’ by contributing to the electrical (*polarization*) or magnetical (*magnetization*) flux. In simple matter (homogeneous, linear, and isotropic), the flux densities of the two field components are proportional to the respective field strengths,

$$\mathbf{D} = \epsilon\mathbf{E} \quad \text{and} \quad \mathbf{B} = \mu\mathbf{H} . \quad (\text{C.10})$$

Bound charge carriers inserted artificially cause a **net**<sup>2</sup> charge of the body, which is represented by a non-vanishing charge density  $\rho$ .

The free charges are the ones responsible for conduction. The Lorentz force fully acts on them, but frequent collisions with other particles ensure that their velocity does not exceed a certain material-dependent limit. In simple matter, this behavior is expressed as a linear relation between the electric field strength and the electric current density, which is known as *Ohm’s law*:

$$\mathbf{J} = \sigma\mathbf{E} . \quad (\text{C.11})$$

Note that depending on the properties of the matter in question, other transport mechanisms may need to be added to the model, effectively *replacing* Ohm’s law. The semiconductors the whole lecture is about are an example for this.

## C.1 Interface Conditions

When considering the interface of two adjacent bodies with different material parameters, it is obvious that not all field quantities can be continuous; e.g. consider two bodies composed of simple matter in the sense of the previous section having different permittivities. Since in both bodies the respective linear relations between  $\mathbf{D}$  and  $\mathbf{E}$  are valid, both fields can not be continuous simultaneously across the interface.

The remedy is to allow step-like discontinuities in the fields at the interface. Clearly, at these discontinuities the fields are not differentiable any more, but the integral formulation of Maxwell’s equations provides clues on how the discontinuities behave. When integrating over a surface or volume that includes an interface, the integral is divided in two parts, one ‘before’ and one ‘after’ the interface, and the integral theorems of Stokes and Gauss are extended by expressions that handle the discontinuities at the interface. For the electric field, the relations are

$$\mathbf{D}_2 \cdot \mathbf{n} - \mathbf{D}_1 \cdot \mathbf{n} = \rho_s \quad \text{and} \quad \mathbf{E}_2 \times \mathbf{n} - \mathbf{E}_1 \times \mathbf{n} = 0 , \quad (\text{C.12})$$

---

<sup>1</sup> **to constitute** [ˈkɑːn.stɪ.tuːt]: bilden, ausmachen    <sup>2</sup> **net** [net]: netto (auch: Netz)

where  $\mathbf{n}$  is the unit vector locally perpendicular to the interface when **traversing**<sup>1</sup> from body ‘1’ (with  $\mathbf{D}_1$  and  $\mathbf{E}_1$ ) to body ‘2’ (with  $\mathbf{D}_2$  and  $\mathbf{E}_2$ ). The equations state that the electric flux component *perpendicular to the interface* has a discontinuity through the sheet charge density  $\rho_s$  at the interface,<sup>2</sup> and the electric field strength’s component *tangential to the interface* is continuous. The respective other component (the tangential flux and the perpendicular field strength) can be calculated using the respective permittivity relation in the body in question,

$$\mathbf{D}_{1,2} \times \mathbf{n} = \varepsilon_{1,2} \mathbf{E}_{1,2} \times \mathbf{n} \quad \text{and} \quad \mathbf{E}_{1,2} \cdot \mathbf{n} = \frac{1}{\varepsilon_{1,2}} \mathbf{D}_{1,2} \cdot \mathbf{n} . \quad (\text{C.13})$$

## C.2 Continuity Equation

One of the most fundamental axioms in physics is the assumption that electric charges can not be generated or destroyed. Being a prerequisite to Maxwell’s equations, it must be possible to deduce it from them. In fact, taking the divergence of Ampère’s Law,

$$\text{div rot } \mathbf{H} = \text{div } \mathbf{J} + \text{div} \left( \frac{\partial \mathbf{D}}{\partial t} \right) = \text{div } \mathbf{J} + \frac{\partial \text{div } \mathbf{D}}{\partial t} = \text{div } \mathbf{J} + \frac{\partial \rho}{\partial t} . \quad (\text{C.14})$$

But at the other hand,  $\text{div rot} \equiv 0$ , and therefore we get the so called *charge continuity equation*

$$\text{div } \mathbf{J} = -\frac{\partial \rho}{\partial t} . \quad (\text{C.15})$$

---

<sup>1</sup> **to traverse** [trə'vɜ:z]: überschreiten, durchlaufen    <sup>2</sup> Since at interfaces the crystal structure of matter is perturbed, a net charge is assumed to be present there in general.



## Appendix D

### Vocabulary

abbreviation	[ə.bri:vi'eɪ.ʃən]	Abkürzung, Kurzwort
accordance	[ə'kɔ:r.dənts]	Übereinstimmung
accurate	[æk.jʊ.rət]	genau
act upon	[ækt ə'pɒn]	einwirken
adjacent	[ə'dʒeɪ.sənt]	benachbart, angrenzend
aka, also known as	[ɔ:lsoʊ noʊŋ æz]	auch bekannt unter, so genannt
analogous	[ə'næl.ə.gəs]	analog, entsprechend, sinngemäß
arbitrary	['ɑ:rbətəri]	willkürlich, beliebig
archetypical	[ɑ:kɪ'taɪ.p.əl]	urbildlich, musterhaft, typisch
arsenic	['ɑ:sən.ɪk]	Arsen
ascent	[ə'sent]	Anstieg
assumption	[ə'sʌmp.ʃən]	Annahme
big picture	[bɪg pɪk.tʃə]	Das große Ganze, ein erster Einstieg
bivariate	[baɪvəri:ət]	von zwei Variablen abhängig
boron	['bɔ:ʀən]	Bor
brevity	['brev.ə.ti]	Kürze
calculation	[kæl.kjʊ'leɪ.ʃən]	Berechnung
Cartesian	[kɑ:'ti:zi.ən]	kartesisch
caveat	['kæv.i.æt]	Hinweis, Warnung
cf, to confer	[kən'fɜ:r]	vergleichen, konsultieren
concise	[kən'saɪs]	übersichtlich
contradictory	[kɒn.trə'dɪk.tɔ:ri]	widersprüchlich
counterclockwise	[kaʊn.tə'klok.waɪz]	gegen den Uhrzeigersinn
crude	[kru:d]	grob, ungehobelt
cyclic	['saɪ.kli.k]	zyklisch, periodisch
Delaunay	[deləʊneɪ]	Delaunay
deliberately	[dɪ'lɪb.ə'r.ət.li]	absichtlich
denominator	[dɪ'na:mə.neɪ.tɔ:]	Nenner
derivative	[dɪ'rɪv.ə.tɪv]	Ableitung
deviation	[di:vi'eɪ.ʃən]	Abweichung



dipole	['dɑ:pəʊl]	Dipol
discretization	[dɪ'skri:tɪ'seɪ.ʃən]	Diskretisierung
discriminant	[dɪskrɪmɪnənt]	Diskriminante
disjoint	[dɪs'dʒɔɪn.t]	disjunkt, einander nicht überschneidend
distinction	[dɪ'stɪŋk.ʃən]	Unterscheidung
distribution	[dɪ'strɪb.ju:ʃən]	Verteilung
dumb	[dʌm]	einfältig, primitiv
electron	[ɪ'lektɹɑ:n], NOT ['ələktɹɑ:n]	Elektron
ellipse	[ɪ'lɪps]	Ellipse
encroachment	[ɪn'krəʊtʃ.mənt]	Beeinträchtigung
evenly	[i:vən.li]	gleichmäßig
excess	[ek'ses]	überschüssig
expansion	[ɪk'spæn.tʃən]	hier: Entwicklung
familiar	[fə'mɪl.i.jə]	vertraut, geläufig
fictitious	[fɪk'tɪʃ.əs]	fiktiv
for the sake of sth.	[fɔ:ʀ ðe seɪk əv]	um ... Willen
handy	[hændi]	praktisch, geschickt
hyperbola	[haɪ'pɛɪ.bəl.ə]	Hyperbel
impurity	[ɪm'pjʊərɪti]	Störstelle, Störatom
inertia	[ɪ'nɜ:ʃə]	Trägheit
inexhaustible	[ɪn.ɪg'zɔ:stɪ.bəl]	unerschöpflich
intersection	[ɪn.tə'sek.ʃən]	Schnittmenge, Schnittpunkt
ionic conductor	[aɪ'ɔn.ɪk kən'dʌk.tə]	Ionenleiter
junction	['dʒʌŋkʃən]	die Sperrschicht
Laplacian	[ləpləsiən]	Laplace-Operator
lattice	[læt.ɪs]	Kristallgitter
modulus	[mɒd.ju:ləs]	Absolutbetrag
negligible	['neg.lɪ.dʒə.bəl]	vernachlässigbar
net	[net]	netto (auch: Netz)
no matter	[nəʊ 'mæt.ə]	ganz egal
obsolete	[əb.sə'li:t]	hinfällig
on closer inspection	[ɔn kləʊsər ɪn'spek.ʃən]	bei näherer Betrachtung
overshoot	[əʊ.və'ʃu:t]	die Überhöhung
painstakingly	[peɪnz'teɪ.kɪŋ.li]	sorgfältig
parabola	[pə'ræb.əl.ə]	Parabel
peculiarity	[pɪ.kju:li'eri.ə.ti]	Eigenheit, Ausprägung
phosphorus	[fɒs.fə.əs]	Phosphor
pitfall	[pɪt.fa:l]	Fallgrube, Fallstrick, Fehler

VOCABULARY

plethora	[ˈpleθ.ə.ə]	Fülle, Vielzahl
preceding	[priːsiː.dɪŋ]	vorangegangen
preface	[ˈprefɪs], NOT [ˈpriːfeɪs]	Vorwort
preliminary	[priːlɪm.ɪ.nɪ.ə.ri]	vorläufig, vorübergehend
prerequisite	[priːˈrek.wɪ.zɪt]	Voraussetzung, Bedingung
principal minors	[prɪnt.sɪ.pəl maɪ.nəʳ]	Hauptminoren
pronunciation	[prəˌnʌntsiˈeɪʃən]	Ausprache
quantity	[ˈkwɑːn.tə.ti]	Größe, hier speziell: Matrixelemente
reentrant corners	[riː.en.trənt]	einspringende Ecken
remedy	[rem.ə.di]	Abhilfe, Behelf
root	[ruːt]	Wurzel, Ursprung, hier: Nullstelle
roughness	[rʌf.nəs]	Rauheit, Unebenheit
rudimentary	[ruːdiˈmen.tɪ.ə.ri]	elementar
satisfactory	[sæt.ɪsˈfækt.ɪ.ə.ri]	zufriedenstellend
solely	[səʊl.li]	lediglich
sound	[saʊnd]	auch: vernünftig
sophistication	[səˌfɪs.tɪˈkeɪ.ʃən]	Raffinesse
steep	[stiːp]	steil
stencil	[stent.səl]	Vervielfältigungsmatrix
straightforward	[streɪtˈfɔː.wəd]	unkompliziert
subtlety	[sʌt.l.ti]	Schwierigkeit, Raffinesse
sufficient	[səf.ɪ.ʃ.ənt]	ausreichend
tessellation	[tes.əl.ɪ.ʃən]	Mosaik
throughout	[θruːˈaʊt]	durchweg, hindurch
to accomplish	[əˈkɑːm.plɪʃ]	etwas erreichen, etwas vollbringen
to attribute	[ˈæt.rɪ.bjuːt]	zuschreiben, zurückführen auf
to be governed	[gʌv.ənd]	bestimmt sein
to be subjected to sth.	[sʌb.dʒekt]	etwas ausgesetzt werden
to cope	[kəʊp]	zurechtkommen, beherrschen
to coincide	[kəʊ.ɪnˈsaɪd]	übereinstimmen
to come in handy	[kʌm ɪn hæ.n.dɪ]	sich gut treffen, gelegen kommen
to constitute	[ˈkɑːn.stɪ.tuːt]	bilden, ausmachen
to decompose	[diː.kəmˈpəʊz]	aufteilen, spalten
to decouple	[dɪkʌp.l]	entkoppeln
to deplete	[dɪˈpliːt]	verringern, aufbrauchen, verarmen
to determine	[dɪˈteɪ.mɪn], NOT [determɪn]	bestimmen, festlegen
to devote	[dɪˈvəʊ.tɪd]	widmen
to diminish	[dɪˈmɪn.ɪʃ]	abnehmen, abklingen
to disturb	[dɪˈstɜːb]	stören, durcheinanderbringen
to elevate	[el.ɪ.veɪt]	emporheben, erhöhen
to employ	[ɪmˈplɔɪ]	einführen, einsetzen
to emphasize	[emp.fəˈsaɪz]	betonen, hervorheben
to encourage	[ɪmˈkʌr.ɪdʒ]	animieren, ermuntern
to evaluate	[ɪˈvæl.ju.ert]	etwas auswerten

to justify	[dʒʌs.tɪ.fɑɪ]	rechtfertigen
to neglect	[nɪ'glekt]	vernachlässigen
to omit	[oʊ'mɪt]	auslassen
to pay attention to sth.	[peɪ ɛ'ten.tʃən]	auf etwas achten
to perpendicular bisect	[pɛr.pə'n'dɪk.jʊ.lər bar'sekt]	in zwei gleich große Teile teilen
to prevail	[prɪ'veɪl]	überwiegen, vorherrschen
to prohibit	[prə'hɪb.ɪt]	verhindern, unterbinden
to propagate	[prɒp.ə.gert]	ausbreiten, fortpflanzen
to provoke	[prə'vəʊk]	auslösen, bewirken
to readdress	[,rɪ:ə'dres]	sich noch einmal zuwenden
to reason	[ri:zən]	begründen, überlegen
to recover	[rɪ'kʌv.ə]	zurückgewinnen, wiedererlangen
to resort to sth.	[rɪ'zɔ:rt]	auf etw. zurückgreifen
to scatter	[skæɪ.ə]	streuen, zerstreuen
to screen	[skri:n]	abschirmen, schützen, filtern
to seek sth.	[si:k]	etwas suchen
to surmount	[sə'maʊnt]	bewältigen, überwinden
to traverse	[trə'vɜ:s]	überschreiten, durchlaufen
to vanish	[væn.ɪʃ]	verschwinden
to yield sth.	[ji:ld]	etwas ergeben, etwas hervorbringen
tractable	[træk.tə.bəl]	handhabbar, lenkbar
truncation error	[trʌŋ'keɪ.ʃən 'er.ə]	Abschneidefehler
unambiguously	[ʌn.æm'big.ju.ə.sli]	eindeutig
unfeasible	[ʌn'fi:zɪ.bəl]	undurchführbar
vacancy	[veɪ.kənt.si]	freie Stelle, insbes. Gitterfreistelle
vicinity	[və'sɪnəti]	Umgebung
worth noticing	[wɜ:θ nɒv.tɪsɪŋ]	erwähnenswert

# Bibliography

- [1] D.M. Caughey and R.E. Thomas. Carrier mobilities in silicon empirically related to doping and field. *Proceedings of the IEEE*, 55(12):2192–2193, Dec. 1967.
- [2] J.-P. Colinge and C. A. Colinge. *Physics of Semiconductor Devices*. Springer Science+Business Media, Inc., New York, 2002.
- [3] Conic Sections — Wikipedia. [http://en.wikipedia.org/wiki/Conic\\_sections](http://en.wikipedia.org/wiki/Conic_sections).
- [4] R. N. Hall. Electron-Hole Recombination in Germanium. *Phys. Rev.*, 87(2):387, Jul 1952.
- [5] R. Jaggi. High-Field Drift Velocities in Silicon and Germanium. *Helvetica Physica Acta*, 42:941–943, 1969.
- [6] R. Jaggi and H. Weibel. High-Field Electron Drift Velocities and Current Densities in Silicon. *Helvetica Physica Acta*, 42:631–632, 1969.
- [7] G. Massobrio and P. Antognetti. *Semiconductor Device Modeling*. McGraw-Hill Professional, 1998.
- [8] E. H. Nicollian and J. R. Brews. *MOS (Metal Oxide Semiconductor) Physics and Technology*. Wiley-Interscience, 1982.
- [9] C. T. Sah. Characteristics of the metal-Oxide-semiconductor transistors. *IEEE Transactions on Electron Devices*, 11:324–345, 1964.
- [10] D. L. Scharfetter and H. K. Gummel. Large-Signal Analysis of a Silicon Read Diode Oscillator. *IEEE Transactions on Electron Devices*, 16:64–77, 1969.
- [11] S. Selberherr. *Analysis and Simulation of Semiconductor Devices*. Springer-Verlag, New York, 1984.
- [12] W. Shockley and W. T. Read. Statistics of the Recombinations of Holes and Electrons. *Phys. Rev.*, 87(5):835–842, Sep 1952.
- [13] John C. Strikwerda. *Finite Difference Schemes and Partial Differential Equations*. Mathematics Series. Wadsworth & Brooks/Cole, Pacific Grove, California, 1989.
- [14] Simon M. Sze. *Physics of Semiconductor Devices*. Wiley-Interscience, November 1981.
- [15] W. V. Van Roosbroeck. Theory of Flow of Electrons and Holes in Germanium and other Semiconductors. *Bell Syst. Techn. J.*, 29:560–607, 1950.