

---

Konrad-Zuse-Zentrum  
für Informationstechnik Berlin

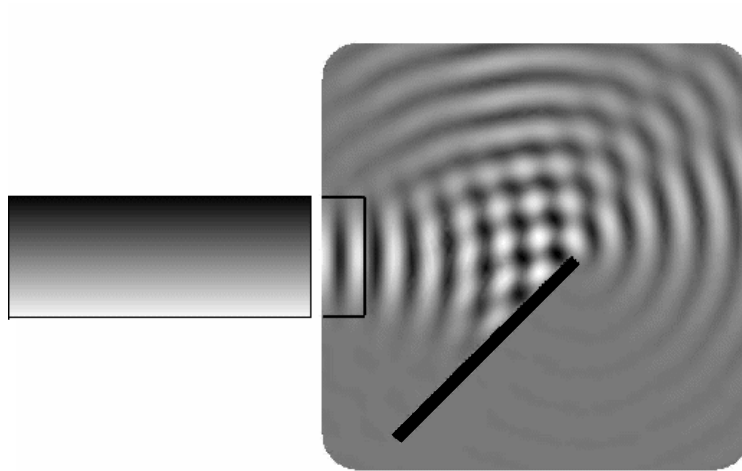
Takustraße 7  
D-14195 Berlin-Dahlem  
Germany

FRANK SCHMIDT

**A New Approach to Coupled  
Interior-Exterior Helmholtz-Type  
Problems: Theory and Algorithms**

**A New Approach to Coupled Interior-Exterior  
Helmholtz-Type Problems:  
Theory and Algorithms**

Frank Schmidt



*Key words and phrases.* Pole Condition; Helmholtz equation; Wide angle one-way equations; Schrödinger equation; Unbounded domains; Transparent boundary Conditions; Finite-element method

### **Acknowledgement**

It is a pleasure to thank all the individuals who contributed to this work. Among those I would like to thank in particular

Peter Deuffhard for making this work possible. Doubtless, without his continuous support, guidance and encouragement I would never have been in a position to reach the current state of the project. Together we overcame the first difficulties in theory and algorithms.

David Yevick for forcing the application of the theory to one-way wide angle equations.

Reinhard März for discussing the properties of periodic structures.

My colleagues Thorsten Hohage and Lin Zschiedrich for their deep interest in the subject, proofreading the manuscript, and the countless number of fruitful discussions, which intensely influenced the whole representation, especially the theoretical statements of Sections 4.2 and 4.3.

And, last but not least, my family who made it all possible in the first place.

### Notations

In physical equations like Maxwell equations, we follow the usual convention and represent two- and three-dimensional vectorial quantities like the radius vector or the magnetic field strength by lower and upper case bold roman letters, in the given case for example by  $\mathbf{x}$  and  $\mathbf{H}$ , respectively. In the algorithmic part, where the language of numerical linear algebra is used, we represent vectors by lower case roman letters, and matrices by upper case Greek or roman letters. In both cases, lower case Greek letters and lower case Italic letters represent scalars.

#### Scalars and Vectors

$x, y, z$	Cartesian coordinates
$t$	time
$\mathbf{x}$	radius vector; $\mathbf{x} \in \mathbb{R}^d$ , $d = 1, 2, 3$
$\mathbf{x}^0$	unit vector $\mathbf{x}/ \mathbf{x} $ , $\mathbf{x} \in \mathbb{R}^d$
$\mathbf{n}$	normal vector pointing to the exterior
$k$	wavenumber; $k^2$ - sometimes called potential
$\mathbf{k}$	wavevector
$\lambda$	wavelength
$\bar{s}$	complex conjugate number to $s$

#### Sets

$\mathbb{N}, \mathbb{Z}, \mathbb{R}, \mathbb{C}$	natural, real, integer, and complex numbers
$\mathbb{R}_+$	the set $\{x \in \mathbb{R} : 0 \leq x < \infty\}$
$\mathbb{N}_0$	$\mathbb{N} \cup \{0\}$
$[a, b]$	closed interval, the set $\{x \in \mathbb{R} : a \leq x \leq b\}$
$]a, b]$	half-open interval, the set $\{x \in \mathbb{R} : a < x \leq b\}$
$]a, b[$	open interval, the set $\{x \in \mathbb{R} : a < x < b\}$
$S^d$	unit sphere, $S^d = \{\mathbf{x} \in \mathbb{R}^{d+1} : \ \mathbf{x}\  = 1\}$
$\Omega$	computational domain, $\Omega \in \mathbb{R}^d$
$\partial\Omega$	boundary of $\Omega$
$\bar{\Omega}$	closure of $\Omega$
$\text{span } S$	linear hull of the set $S$
$\text{meas } S$	measure of the set $S$
$\text{ext } \Omega$	exterior of $\Omega$ , $\mathbb{R}^2 \setminus \bar{\Omega}$
$X \times Y$	product set $\{(x, y) : x \in X \text{ and } y \in Y\}$

#### Derivatives

$\partial_x, \partial_y, \partial_z$	partial derivatives, $\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z}$
$\frac{\partial^{ \alpha }}{\partial x^\alpha}$	partial derivative in multi-index notation, $\frac{\partial^{ \alpha }}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_n^{\alpha_n}}$
$\alpha$	multi-index, $\alpha = (\alpha_1, \dots, \alpha_n)$ with $\alpha_i \in \mathbb{Z}, \alpha_i \geq 0$
$ \alpha $	$\alpha_1 + \dots + \alpha_n$
$\nabla$	nabla operator $\nabla = (\partial_x, \partial_y, \partial_z)$ in Cartesian coordinates

$\text{grad } u$	gradient; in Cartesian coordinates $\nabla u = (\partial_x u, \partial_y u, \partial_z u)^T$
$\nabla_{\xi\eta} u$	gradient with respect to general $\xi\eta$ -coordinates
$\text{curl } \mathbf{u}$	curl of the vector $\mathbf{u}$ ; in Cartesian coordinates $\nabla \times \mathbf{u}$
$\text{div } \mathbf{u}$	divergence of the vector $\mathbf{u}$ ; in Cartesian coordinates $\nabla \cdot \mathbf{u}$
$\partial_n$	derivative in the direction of the exterior normal
$\Delta$	Laplacian, $\partial_x^2 + \partial_y^2 + \partial_z^2$

### Operators

$\mathcal{D}(A)$	domain of definition of the operator $A$
$\mathcal{R}(A)$	range of the operator $A$ , $\mathcal{R}(A) := \{v : v = Au \text{ for } u \in \mathcal{D}(A)\}$
$\mathcal{N}(A)$ or $\ker A$	null space (or kernel) of the operator $A$ , $\mathcal{N}(A) := \{u \in \mathcal{D}(A) : Au = 0\}$
$\sigma(A)$	spectrum of $A$
$\mathcal{N}^\perp(A)$	orthogonal complement of the null space of $A$
$M + L$	sum of the space $M$ and $L$ , $M + L = \{x + y : x \in M \text{ and } y \in L\}$
$M \oplus L$	orthogonal direct sum, $M + L$ , where $L = M^\perp$
$A \otimes B$	Kronecker tensor product of two matrices $A$ and $B$
$AB$ or $A \circ B$	product of the operators $A$ and $B$ , $(AB)(u) := A(Bu)$
$\text{Tr}$	Trace operator

### Norms and Inner Products

$ z $	absolute value of the complex number $z$
$ \cdot $	Euclidean norm
$\ \cdot\ _{L^2(\Omega)}$	$L^2$ norm on $\Omega$
$(v, u)$	inner product in $L^2(\Omega)$
$\langle v, u \rangle$	inner product in $L^2(\partial\Omega)$ or Euclidean inner product of two vectors
$\ \cdot\ _1$	$L^1$ norm
$\ \cdot\ _\infty$	$L^\infty$ norm
$\ \cdot\ _{H^s(\Omega)}$	Sobolev norm for $s \in \mathbb{R}$ ; $\ (I - \Delta)^{\frac{s}{2}} \cdot\ _{L^2(\Omega)}$

### Function Spaces

$\mathcal{S}$	space of rapidly decreasing functions
$\mathcal{S}'$	space of tempered distributions
$C^k$	space of complex, $k$ -times continuously differentiable functions
$C_0$	space of complex, differentiable functions with compact support
$C'_0$	dual space to $C_0$
$H^s$	Sobolev space with $s \in \mathbb{R}$
$L^2$	space of quadratic integrable functions

**Miscellaneous**

$\mathcal{O}(\cdot), o(\cdot)$	Landau symbols
$H_n^{(1)}, H_n^{(2)}$	Hankel functions of first and second kind and order $n$
$J_n$	Bessel function of order $n$
$\delta(\cdot)$	Dirac's delta function
$\sim$	asymptotic equivalence, $f(x) \sim g(x)$ means $f(x) = g(x)(1 + o(x))$

# Contents

Acknowledgement	3
Notations	4
Chapter 1. Introduction	9
1.1. Physical Models of Wave Phenomena	10
1.2. Examples of Coupled Interior-Exterior Helmholtz-Type Problems	12
1.3. What Are Transparent Boundary Conditions?	16
1.4. Basic Idea of the Pole Condition Approach	17
1.5. Outline of Contents	18
Chapter 2. A Survey of Different Classical Methods	20
2.1. Factorization Based on Separable Coordinates	20
2.2. Boundary Integral/Element Methods	22
2.3. Infinite Element Methods	25
2.4. Asymptotic Methods	27
2.5. Further Methods	34
Chapter 3. Characterization of Incoming and Outgoing Waves: Pole Condition	37
3.1. Helmholtz Equation in 1D with Constant Potential	37
3.2. Helmholtz Equation in 1D with Periodic Potentials	38
3.3. Radially Symmetric Helmholtz Equation in 2D	46
3.4. Generalizations	51
Chapter 4. Existence and Uniqueness Statements for Separable Problems	57
4.1. Separable Coordinates	57
4.2. Preparation of Main Theory	61
4.3. Main Theoretical Results	66
4.4. Important Consequences	80
4.4.1. Representation formula	80
4.4.2. Equivalent Formulations of the Pole Condition	81
4.4.3. Asymptotic Expansion of the Far-Field	83
4.4.4. Spectral Properties of the Dirichlet-to-Neumann Map	85
4.4.5. High-Frequency Limit of the Cut Functions	87
4.4.6. Relation to the PML Method	88
Chapter 5. Numerical Treatment of Helmholtz-Type Scattering Problems	91
5.1. Factorization Approach	91
5.1.1. Reflection at an Infinite Plane	91
5.1.2. Reflection and Diffraction by a Semi-Infinite Plane	102
5.2. Laplace domain method	105
5.2.1. Real Axis Approach	106
5.2.2. Cut Function Approach	115
5.3. Non-Separable, Discrete Problems in 2D	120
5.4. Numerical Examples	139

Chapter 6. Numerical Treatment of Schrödinger-Type Scattering Problems	144
6.1. Time-Dependent Schrödinger Equation	144
6.2. Wide Angle One-Way Equations	151
6.3. Stability	162
6.4. Numerical Experiments	164
Conclusions	170
Appendix A. Basic Properties of the Continuous and Discrete Transforms	172
A.1. Basic Properties of the Laplace Transform	172
A.2. Discrete Operational Calculus	176
Appendix. Bibliography	181
Appendix. Index	185



## CHAPTER 1

# Introduction

This work presents a new approach to the numerical solution of scattering problems on unbounded domains. The presentation is rather broad, so that it is hopefully of interest for an interdisciplinary community of mathematicians, physicists and engineers working on the simulation of wave propagation phenomena.

Our goal is to solve scattering problems modeled by the Helmholtz equation

$$\begin{aligned}
 \Delta u(\mathbf{x}) + k^2(\mathbf{x}) u(\mathbf{x}) &= f(\mathbf{x}) && \text{in } \Omega \\
 \Delta u_{\text{out}}(\mathbf{x}) + k^2(\mathbf{x}) u_{\text{out}}(\mathbf{x}) &= 0 && \text{outside } \Omega \\
 u(\mathbf{x}) &= u_{\text{src}}(\mathbf{x}) + u_{\text{out}}(\mathbf{x}) && \text{on } \partial\Omega \\
 \partial_n u(\mathbf{x}) &= \partial_n u_{\text{src}}(\mathbf{x}) + \partial_n u_{\text{out}}(\mathbf{x}) && \text{on } \partial\Omega \\
 \partial_n u_{\text{out}}(\mathbf{x}) &= Bu_{\text{out}}(\mathbf{x}) && \text{on } \partial\Omega.
 \end{aligned}
 \tag{1.0.1}$$

Here  $k(\mathbf{x})$  is the position dependent wavenumber,  $k^2$  is often also called potential, and  $f(\mathbf{x})$  is a given source term with support only in some bounded computational domain  $\Omega \subset \mathbb{R}^d$ , where  $d$  is the space dimension. Practically we consider the cases  $d = 1, 2$ ; the theoretical analysis is done for arbitrary space dimensions. We assume that  $u(\mathbf{x})$  can be decomposed into a source field  $u_{\text{src}}(\mathbf{x})$  and an outgoing part  $u_{\text{out}}(\mathbf{x})$  on the boundary of the computational domain. The source field is a field excited by sources placed in the exterior of  $\Omega$  and may contain both incoming and outgoing fields. The most prominent source field is a plane wave  $\exp(i\mathbf{k}\mathbf{x})$ , but our approach will allow for much more general source fields. A precise definition of incoming and outgoing fields will be given in Section 3.1 (1D) and Section 3.4 (higher dimensions). The above boundary operator  $B$  is called Dirichlet-to-Neumann (DtN) operator (also Calderon operator or Poincaré-Steklov operator), because it maps Dirichlet to Neumann data on the boundary. The main part of this work will focus on the definition, construction, and analysis of discrete approximations to the boundary operator. The common way to formulate the problem, for a wavenumber  $k$  constant in some exterior domain, is to pose it right from the beginning on the entire, unbounded domain, and to impose the asymptotic Sommerfeld radiation condition

$$\lim_{\rho \rightarrow \infty} \rho^{\frac{d-1}{2}} (\partial_\rho u - iku) = 0, \quad \rho = |\mathbf{x}|, \quad \text{uniformly in all directions } \frac{\mathbf{x}}{|\mathbf{x}|},
 \tag{1.0.2}$$

where  $d$  is the space dimension. For  $k = k(\mathbf{x})$  explicitly spatially dependent, however, the Sommerfeld condition does not hold and all the methods based on it must fail. Example 1.2.2 in Section 1.2 presents a practically relevant example of a scattering problem with an inhomogeneous exterior domain. Therefore our treatment is *not* based on Sommerfeld's radiation condition. We shall replace the Sommerfeld condition by a property which we call *pole condition*. This notion will be *central* in our presentation. The pole condition is applicable in a natural way to a large class of problems with position dependent wavenumbers and arbitrary space dimension. Further it offers many constructive ways to compute corresponding boundary operators  $B$ , and thus to restrict problems given on unbounded domains to such ones given on bounded domains. Moreover, the pole condition can be applied in a direct

manner to the *one-way wide angle approximations* of the Helmholtz equation which include, for example the time-dependent *Schrödinger equation*. We call the class of scattering problems consisting of the classical Helmholtz equation with possibly non-constant coefficients in the exterior domain and the one-way wide angle approximations including the time-dependent Schrödinger equation *Helmholtz-type* problems.

From the application point of view it is our goal to solve numerically scattering problems of Helmholtz-type

$$\begin{array}{l} \text{with} \\ \text{without} \end{array} \left\{ \begin{array}{l} \bullet \text{ inhomogeneities in the exterior domain,} \\ \bullet \text{ computational domains naturally defined by the problem,} \\ \\ \bullet \text{ an explicit use of Green's functions,} \\ \bullet \text{ being confined to separable coordinates,} \\ \bullet \text{ explicit knowledge of an asymptotic radiation condition.} \end{array} \right.$$

The types of the exterior inhomogeneities for the Helmholtz-type problems, which can be handled by our method, follow from the construction principle and the kind of discretization of the exterior region. Most important here is the possibility to take into account waveguide-type perturbations of homogeneous media. In Section 1.2 we will present a typical example from integrated optics.

### 1.1. Physical Models of Wave Phenomena

The Helmholtz equation on unbounded domains plays an important role in many areas of science and engineering. Most prominent here are electromagnetics, acoustics and quantum mechanics. In quantum mechanics, the Helmholtz equation usually occurs as an eigenvalue problem related to the Schrödinger equation. Despite of the fact that the numerical solution of Helmholtz-type eigenvalue problems on unbounded domains has much in common with the solution of scattering problems on unbounded domains, we will not include this in our work, because the numerical treatment of eigenproblems is a problem of its own.

**Maxwell's Equations.** Let us discuss first the occurrence of Helmholtz equations in electromagnetics. The electromagnetic field in charge-free and non-conducting media is governed by Maxwell's equations

$$\begin{array}{ll} \text{curl } \mathbf{H}(\mathbf{x}, t) &= \epsilon(\mathbf{x}) \partial_t \mathbf{E}(\mathbf{x}, t) & \text{div } \mu(\mathbf{x}) \mathbf{H}(\mathbf{x}, t) &= 0 \\ \text{curl } \mathbf{E}(\mathbf{x}, t) &= -\mu(\mathbf{x}) \partial_t \mathbf{H}(\mathbf{x}, t) & \text{div } \epsilon(\mathbf{x}) \mathbf{E}(\mathbf{x}, t) &= 0 \end{array}$$

where  $\mathbf{H}(\mathbf{x}, t)$  and  $\mathbf{E}(\mathbf{x}, t)$  denote the magnetic and electric field strength, respectively, in space and time. The material properties are described by the electric permittivity  $\epsilon(\mathbf{x})$  and the magnetic permeability  $\mu(\mathbf{x})$ . If we eliminate either  $\mathbf{H}$  or  $\mathbf{E}$  from one of the curl-equations, we obtain a vectorial wave equation. E.g. the elimination of  $\mathbf{E}$  from the first curl-equation results in

$$\text{curl } \frac{1}{\epsilon(\mathbf{x})} \text{curl } \mathbf{H}(\mathbf{x}, t) = -\mu(\mathbf{x}) \partial_t^2 \mathbf{H}(\mathbf{x}, t).$$

One of the oldest and most successful approximations of the vectorial wave equation is obtained, when we regard the factor  $1/\epsilon(\mathbf{x})$  as a constant with respect to the curl-operation. As long as the permittivity changes only weakly over a distance of a local wavelength, the approximation supplies an accepted physical model of optical wave propagation. It is the most essential approximation, which leads us from Maxwell's equations to geometrical optics, see [13, chapt. 3.1, pp. 111] and scalar wave optics. Following this approximation, applying the identity

$\text{curl curl } \mathbf{H} = \text{grad div } \mathbf{H} - \Delta \mathbf{H}$  and taking the divergence condition into account, we obtain the scalar wave equation

$$\Delta \mathbf{H} = \epsilon(\mathbf{x})\mu(\mathbf{x})\partial_t^2 \mathbf{H}$$

where the Laplacian is applied to each individual component of  $\mathbf{H}$  separately. Since the vectorial components are independent of each other, the equation separates into three independent scalar equations. Let us denote one of the vectorial components of  $\mathbf{H}(\mathbf{x}, t)$  by  $u(\mathbf{x}, t)$ . Further let us assume that field  $u(\mathbf{x}, t)$  is harmonic in time, i. e.  $u(\mathbf{x}, t) = \Re(\tilde{u}(\mathbf{x}) \exp(-i\omega t))$ , where  $\omega$  is the angular frequency of the wave. Set  $k^2(\mathbf{x}) := \epsilon(\mathbf{x})\mu(\mathbf{x})$ . If we introduce this complex representation into the wave equation and drop the tilde, we arrive at the scalar Helmholtz equation

$$\Delta u(\mathbf{x}) + k^2(\mathbf{x})u(\mathbf{x}) = 0.$$

Let us consider the case, where the wavenumber  $k$  is independent of the position, that is the wave propagation in homogeneous media. Then all plane waves  $u(\mathbf{x}) = \text{const} \exp(i\mathbf{k}\mathbf{x})$ , with a wavevector  $\mathbf{k}$  which satisfies the dispersion relation  $\mathbf{k} \cdot \mathbf{k} = k^2$ , are solutions of the Helmholtz equation. These plane waves correspond to time-harmonic solutions  $u(\mathbf{x}, t) = \text{const} \text{Re} \exp(i(\mathbf{k} \cdot \mathbf{x} - \omega t))$ . We consider a fixed direction  $\mathbf{x}_0$  and study the evolution of the phase  $\arg \exp(i(\mathbf{k} \cdot \mathbf{x}_0 - \omega t))$  with increasing time. Clearly, we can identify the case  $\text{Re}(\mathbf{k} \cdot \mathbf{x}_0) > 0$  with a plane wave, which propagates in the positive direction with respect to  $\mathbf{x}_0$  and the opposite case  $\text{Re}(\mathbf{k} \cdot \mathbf{x}_0) < 0$  with a plane wave, which propagates in the negative direction with respect to  $\mathbf{x}_0$ . Note that a definition of the time-dependency based on  $\exp(i\omega t)$  would cause us to define the reverse directions as positive and negative propagation direction.

**Acoustics.** We consider acoustic wave propagation in a compressible, ideal fluid. This involves the interaction of the following three physical effects [31, Vol. 1, Chapt. 47]:

- (1) The fluid moves and changes its density  $\rho(\mathbf{x}, t)$ . The corresponding physical equation is the continuity equation  $\partial_t \rho + \text{div}(\rho \mathbf{v}) = 0$ .
- (2) The change in the density results in a change of the pressure  $p(\mathbf{x}, t)$ . This is described in linear approximation by the equation  $p(\mathbf{x}, t) = p_0 + c^2(\rho(\mathbf{x}, t) - \rho_0)$ , where  $\rho_0$  and  $p_0$  characterize the constant state of the fluid with all particles at rest. The constant factor  $c$  is the velocity of sound.
- (3) Local differences in the pressure result in a change of the velocity  $\mathbf{v}(\mathbf{x}, t)$ . Since  $d(\rho \mathbf{v}) = dt \partial_t(\rho \mathbf{v}) + (d\mathbf{x} \cdot \nabla)(\rho \mathbf{v})$  it follows  $d_t(\rho \mathbf{v}) = \partial_t(\rho \mathbf{v}) + (\mathbf{v} \cdot \nabla)(\rho \mathbf{v})$ . Newton's law supplies the equation of motion  $\partial_t(\rho \mathbf{v}) + (\mathbf{v} \cdot \nabla)(\rho \mathbf{v}) = -\text{grad } p$ . We consider only problems, where the spatial variation of  $\rho \mathbf{v}$  is much smaller than its variation in time. Hence we can linearize the equation of motion to obtain Euler's equation  $\partial_t(\rho \mathbf{v}) = -\text{grad } p$ .

Euler's law, together with the continuity condition, yields  $\partial_t^2 \rho = \text{div grad } p$ . Using the linear relation between density and pressure, we obtain the wave equation for the pressure

$$\Delta p(\mathbf{x}, t) = \frac{1}{c^2} \partial_t^2 p(\mathbf{x}, t).$$

Proceeding as in Section 1.1, we restrict the equation to time-harmonic solutions and introduce the complex pressure  $\tilde{p}(\mathbf{x})$  by  $\text{Re } p(\mathbf{x}, t) = (\tilde{p}(\mathbf{x}) \exp(-i\omega t))$ . We insert this expression into the wave equation, drop the tilde and obtain the Helmholtz equation for the pressure

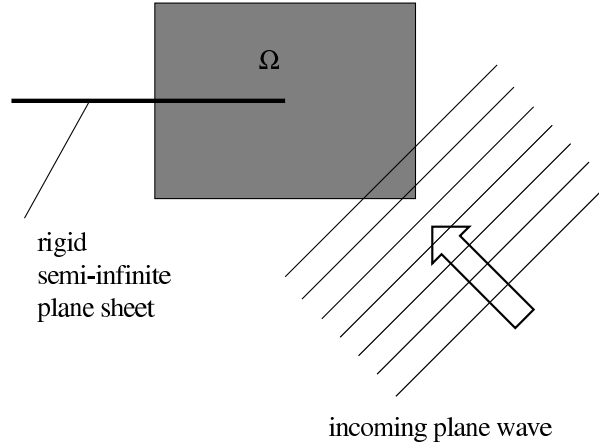


FIGURE 1.2.1. Reflection and diffraction by a semi-infinite plane sheet (Example 1.2.1)

$$\Delta p(\mathbf{x}) + k^2(\mathbf{x})p(\mathbf{x}) = 0 \quad \text{where } k^2(\mathbf{x}) = \omega^2 / c^2.$$

**Quantum Mechanics and Paraxial Beam Propagation.** The non-relativistic motion of an electron of mass  $m$  is governed by Schrödinger's equation

$$i\hbar\partial_t\psi(\mathbf{x}, t) = -\frac{\hbar^2}{2m}\Delta\psi(\mathbf{x}, t) + V(\mathbf{x}, t)\psi(\mathbf{x}, t)$$

where  $\psi(\mathbf{x}, t)$  is the quantum mechanic wave function,  $\hbar = h/(2\pi)$ , with  $h$  the Planck quantum action, and  $V(\mathbf{x}, t)$  a potential. The lowest order wide angle one-way approximation of Helmholtz's equation is of the same type. In optics it takes the form

$$-i\partial_z u(\mathbf{x}, z) = \frac{1}{2n_0k_0} (\Delta u(\mathbf{x}, z) + k_0^2 (n^2 - n_0^2) u(\mathbf{x}, z))$$

and is called Fresnel's equation or the paraxial wave equation. The  $z$ -axis is the main propagation direction of the beam, the transversal coordinates are  $\mathbf{x} = (x, y)^T$ ,  $k_0$  denotes the free-space wave number,  $n(\mathbf{x}, z)$  is the refractive index and  $n_0$  is the so-called reference refractive index. The latter constant is a model parameter and must be chosen such that the deviation with respect to the Helmholtz equation becomes as small as possible, see e. g. [84]. The function  $u(\mathbf{x}, z)$  is the slowly-varying amplitude function and is defined as  $u = \exp(-in_0k_0z)E$ , where  $E$  stands for a component of either the electric or the magnetic field.

## 1.2. Examples of Coupled Interior-Exterior Helmholtz-Type Problems

In the following we want to characterize the problem classes under consideration by typical examples from integrated optics and acoustics. Each of the examples discussed possesses at least one typical difficulty known to be a challenge to its numerical solution.

**EXAMPLE 1.2.1.** *Reflection and diffraction by a semi-infinite plane sheet.*

Fig. 1.2.1 shows schematically an arrangement, where an incoming plane wave strikes a semi-infinite plane sheet. In the acoustic case, the solution  $u(x, y)$  of the Helmholtz equation (1.0.1) is the deviation from the hydrostatic pressure, the

gradient of whose is proportional to the velocity of the particles, see Section 1.1. If the semi-infinite plane sheet is a rigid boundary, there cannot be a normal component of the velocity across it. Therefore the boundary condition is  $\partial_n u = 0$ . The plane wave itself does not satisfy this boundary condition. Therefore the solution  $u(x, y)$  must be a superposition of the incoming plane wave and a scattered field, composed of reflected and diffracted waves. A typical scattering problem of this type is: Given the incidence angle of the incoming plane wave, compute the field in the bounded computational domain  $\Omega$ . This problem belongs to the most simple ones, and we will discuss the numerical solution in Section 5.1.2. Because the whole problem can be formulated based on our one-dimensional concept, it provides an interesting and instructive application of the one-dimensional theory. Further, since an analytic solution based on the famous Wiener-Hopf technique can be derived, it allows a step-by-step comparison between our numerical method and an integral equation approach. It will turn out that the analytic solution process consists of the following main steps:

**Step 1:** First, the entire field is decomposed into the incoming field, the reflected field and the scattered field, where this decomposition is motivated by geometrical optics and needs a-priori knowledge about the solution. While we take the geometrical optics ansatz for the reflected field, the scattered field has to satisfy the Helmholtz equation and Sommerfeld's asymptotic radiation condition. Note that the reflected field *does not* satisfy the radiation condition.

**Step 2:** Apply Green's identity to the decomposed field and derive a Fredholm integral equation of the first kind with a special structure, known as Wiener-Hopf integral equation.

**Step 3:** On the other hand, the density function inside the integral is partially known. This problem can be solved by the Wiener-Hopf technique, which requires, however, an additional trick. To apply the Wiener-Hopf technique, we must extend the real wavenumber  $k$  to a new, complex wavenumber  $\tilde{k} = k + ik_2$  with a small imaginary part  $k_2$ .

In contrast to the integral equation approach, the discrete method based on transparent boundary conditions needs neither a decomposition of the field other than the decomposition into incoming and outgoing waves nor an extension of the real wavenumber to a complex number. The numerical problem here is to take into account a boundary, which extends from a finite point to infinity. This however is a problem for the numerical evaluation of the analytical integral equation too, because all integrals extend to infinity.

**EXAMPLE 1.2.2. Scattering by a finite obstacle.** This problem concerns the computation of the scattered field, excited by an incoming plane wave, in the vicinity of an obstacle with a complicated boundary. The situation, schematically displayed in Fig. 1.2.2, is similar to the preceding one, except that the boundary, which causes the scattering, can be included into a bounded computational domain  $\Omega$ . This simplification allows the application of other methods to compute the scattered field in or outside  $\Omega$ . Especially the choice of separable coordinates in the exterior of  $\Omega$ , if possible, leads to a considerable simplification of the numerical methods.

**EXAMPLE 1.2.3. Inhomogeneous exterior domain.** This example problem, see Fig. 1.2.3, presents a problem with a typical inhomogeneity in the exterior domain. The situation is: The left waveguide, which comes from infinity and ends inside the computational domain  $\Omega$ , guides a mode from infinity to the right end of the waveguide, where the field leaves the waveguide. Let us assume that the wavenumber of

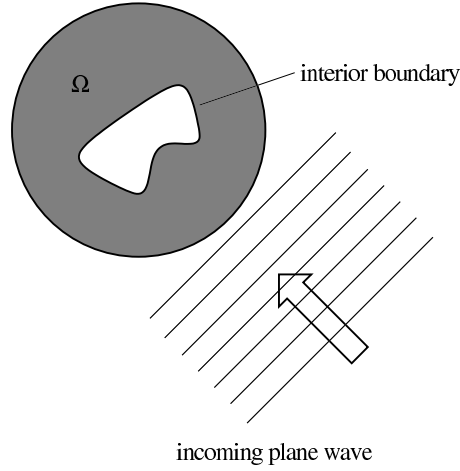


FIGURE 1.2.2. Scattering of a plane wave by a finite obstacle (Example 1.2.2)

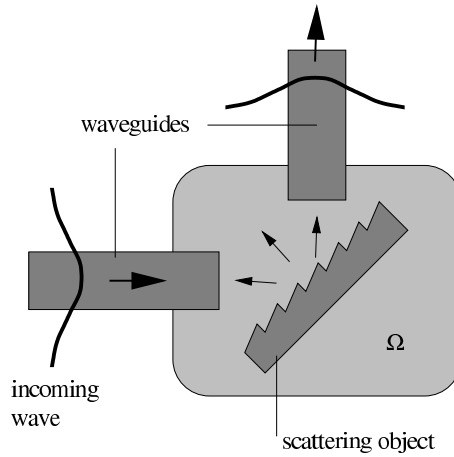


FIGURE 1.2.3. Scattering problem with an inhomogeneous exterior domain. The wavenumber corresponding to the waveguide domains (dark gray) is different from the wavenumber of the surrounding media (Example 1.2.3)

the waveguide is  $k_1$ , the wavenumber of the surrounding homogeneous medium is  $k_0$ , and that  $k_1 > k_0$ . Then the waveguide-mode travels without attenuation from infinity to the end-face of the waveguide with a phase factor  $\exp(i\beta x)$ , where  $x$  denotes the axis of the fiber, and  $k_0 < \beta < k_1$ . Hence the traveling mode *does not* satisfy Sommerfeld's radiation condition. The field radiated by the waveguide hits a scattering object. Schematically we used a wavelength selective mirror, where the field is scattered in all directions. Parts of the incoming field are scattered back into the waveguide, other parts are scattered into a second waveguide, and the remainder is scattered into the homogeneous part of the exterior domain. The part of the wave which is scattered into the second waveguide, can excite a guided mode traveling undamped towards infinity. Again, this part does not satisfy the Sommerfeld condition. We call the described inhomogeneities of the exterior domain waveguide-type inhomogeneities.

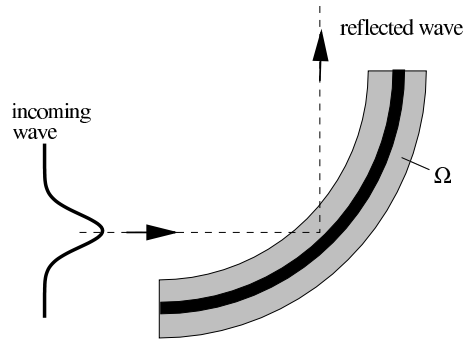


FIGURE 1.2.4. Reflection by a parabolic mirror. The computational domain is non-convex (Example 1.2.4)

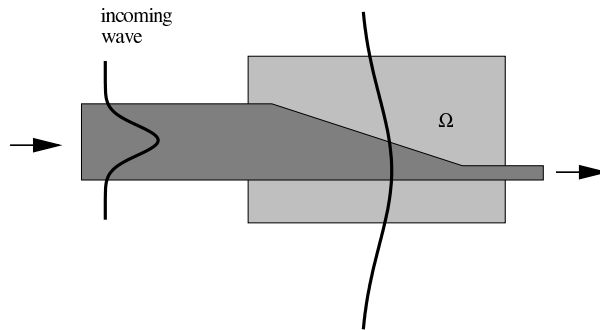


FIGURE 1.2.5. Beam propagation through a tapered waveguide (Example 1.2.5)

So far, all computational domains convex, and, regarding the physical situation, this has been the natural choice. Next, we give a problem type, where a non-convex computational domain possesses numerical advantages.

EXAMPLE 1.2.4. *Parabolic mirror.* Fig. 1.2.4 shows a typical reflection problem, where the scattering surface is of a parabolic type, embedded in a non-convex computational domain. As far as only the near-field solution of the reflection problem is needed, domains of this type might have some numerical advantages compared to domains based on the convex hull of the scattering surface, because they are much smaller than an appropriate convex domains. Thus, the numerical effort to discretize the computational domain is potentially less than in the convex case. Using one version of our Laplace domain method, the cut-function approach, we are even able to compute the far-field without additional effort.

EXAMPLE 1.2.5. *Paraxial beam propagation.* This example is concerned with the one-way wide angle approximation of the Helmholtz equation. A typical situation is shown in Fig. 1.2.5. A waveguide guides a mode from infinity to the computational domain, where the waveguide changes its shape. This change causes in general both a scattering and a deformation of the shape of the guided part of the wave. We indicated the latter schematically by a larger width of the amplitude function at the thinner end of the waveguide. Thus we have a situation similar to the previous ones, except that in this case the direction of the beam propagation is maintained which allows us to apply the one-way equation. Obviously, we could

solve this problem by the same means we will use to solve the reflection problem of Fig. 1.2.3. However, it will turn out that at the present state of the numerical methods it is much more effective to use one of the approximations of the one-way Helmholtz equation.

### 1.3. What Are Transparent Boundary Conditions?

Both our new method – the pole condition approach – as well as the classical methods for the solution of wave propagation problems on unbounded domains can be seen as methods to construct *transparent boundary conditions*. To prepare the very first presentation of our pole condition concept in the following section, we want to introduce this term here. A more detailed description of the classical methods follows in Section 2.

Many methods to solve the discussed wave propagation and scattering problems rely on a decomposition of the unbounded domain, on which the problem is posed, into a bounded – possibly large – interior domain  $\Omega_{\text{int}}$ , which contains the scattering objects, and an exterior domain. Now the following question which quote from Grote and Keller [42] arises: “Does there exist a boundary condition such that the solution of the problem in  $\Omega_{\text{int}}$ , with this boundary condition, coincides exactly with the restriction to  $\Omega_{\text{int}}$  of the solution in the unbounded domain?” Such boundary conditions are called exact or *transparent* boundary conditions, sometimes also radiating or absorbing boundary conditions. There exist a large number of applications, where one is interested in the solution of the scattering problem only in the vicinity of the scatterer, i. e. in the solution on the interior domain. In these cases, the transparent boundary conditions plus a consistent and stable discretization of the interior problem provide the basis of a numerical solution. It is the goal of this work to supply a general approach to construct such transparent boundary conditions.

The above discussion suggests to decompose the problem into two sub-problems: (1) The derivation of the transparent boundary conditions by some method, and, independently of it, (2) the solution of the interior problem by another method. In fact, this is an approach followed successfully by many researchers, cf. e.g. the work of Keller, Givoli, and Grote [64, 38] for the Laplace and the Helmholtz equation, [44, 45] for the wave and time-dependent Maxwell equations. Nevertheless, this strategy has the disadvantage that an analytic solution representation in the exterior domain *must be known in advance*. The same holds true for the already mentioned class of boundary integral methods. Our approach is quite different. Within our methods the transparent boundary conditions and the solution of the interior problem are computed *simultaneously*.

There is another important aspect. Usually the transparent boundary conditions are derived from an exact analytic representation formula. Then they are discretized and applied to the discretized interior problem. The discretization violates the transparency of the boundary which becomes transparent only in the limiting case of a vanishing discretization step width. The pole condition approach, however, which is first given for the continuous case too, has a natural counterpart in the discrete case. Naturally enough, we require for a given discrete scheme that the pole condition has to be satisfied, thus realizing the transparency even for the discrete case. Therefore we call our boundary conditions sometimes also *discrete transparent boundary conditions*.

The next section supplies a first impression of our concept.



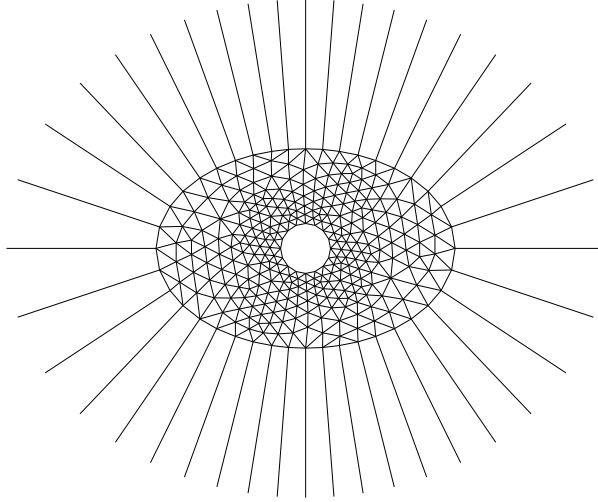


FIGURE 1.4.1. Inner discretization by triangulation and outer discretization by rays used by the methods developed in this work

#### 1.4. Basic Idea of the Pole Condition Approach

We refer to Fig. 1.4.1. The interior domain, which is bounded by an ellipse, is discretized by a triangulation. Any other discretization technique might be applied as well. However, the numerical implementation of the algorithms, which we shall develop, is always based on finite elements and triangulations. This results both in a high flexibility with respect to the discretization of scatterers with complicated shapes and a very natural and simple coupling to the exterior domain. The exterior domain itself is discretized by linear rays with initial points on the artificial boundary. The exterior discretization is chosen such that each node on the artificial boundary is connected with infinity such that the rays do not intersect each other. The transparent boundary condition is realized automatically on the initial points of the rays, but not independently of the interior solution. We compute simultaneously

- (1) A discrete interior solution on the triangulation.
- (2) A discrete solution on the rays.

However, the solution on the rays is not obtained in the given spatial form, instead its Laplace transform is computed. From all possible exterior solutions we accept only those as *outgoing* solutions, whose Laplace transform has no singularity in the lower half of the complex plane. We call this selection rule *pole condition*. The restriction to outgoing solutions supplies the desired transparent boundary conditions.

We can carry out the concept without any analytic pre-treatment. The Laplace transform is computed numerically by means of an Volterra integral equation. This is done very effectively by different types of collocation methods, among them explicit and implicit Runge-Kutta methods and a classic spline approximation method. Note, that no *inverse* Laplace transform is needed. The restriction to linear rays lies not in the nature of the method, however, it supplies the easiest treatment of the exterior discretization.

Our main tool will be the Laplace transform of the fields outside of  $\Omega$ . In this uniform framework the following three principles are of fundamental importance:

- (1) The singularities of the Laplace transform of a complex function defined on a semi-infinite interval determine its far-field. In short: The near-field in the Laplace domain determines the far-field in the physical domain.
- (2) The near field of this function is determined by its Laplace transform for large values of the Laplace variable. In short: The near-field in the physical domain is determined by the far-field in the Laplace domain.
- (3) Problems in higher space dimensions can be composed from mutually interacting one-dimensional problems.

### 1.5. Outline of Contents

The structure of the individual chapters and their mutual relation is the following: In Chapter 2 we give a survey of important classical methods used to derive transparent boundary conditions.

In Chapter 3 we introduce the pole condition. The point of departure is a one-dimensional theory, whose basic point is a constructive definition of incoming and outgoing waves. We develop a tool, which enables us to decompose complex functions given on a semi-infinite interval into incoming and outgoing parts, and we give precise definitions of the notions “incoming” and “outgoing”. We apply the one-dimensional theory to several illustrating and non-trivial examples, including e. g. periodic potentials. Next we generalize the concept to higher space dimensions. It turns out that this can be done completely in terms of the one-dimensional decomposition of waves. Finally we discuss the extension of the one-dimensional concept to the time-dependent Schrödinger equation.

In Chapter 4 we develop the general existence and uniqueness theory for a model problem: the radially symmetric Helmholtz equation. The main results are:

- The pole condition is equivalent to Sommerfeld’s asymptotic radiation condition, at least in the considered finite-dimensional setting.
- The pole condition generalizes Sommerfeld’s radiation condition for cases with potentials depending on the radial distance.
- A new representation formula for the exterior solutions is obtained.
- New asymptotic series representations for the exterior solutions are derived.

In Chapter 5 we derive both the principal structure as well as the essential details of our algorithms to solve scattering problems based on the Helmholtz equation. The emphasis here lies on the development of a new class of algorithms, which we call *Laplace domain methods* since they approximate the essential quantities in the spectral (Laplace) domain. To the best of our knowledge, these Laplace domain methods do not have a predecessor in the numerical literature. They follow naturally from our theory. The algorithms can be divided into two sub-classes: The real axis integral approach and the cut-function approach. We study in detail the numerical properties of these algorithms based on our basic model equation, the Bessel equation. We demonstrate that accuracy in the order of  $10^{-10}$  can be obtained even close to singular cases with a very moderate numerical effort, cf. Fig.’s 5.2.3 and 5.2.7. These studies are performed based on the 1D Bessel model, since corresponding studies in higher dimensions would cause an additional preparation effort. Along with the most general algorithmic realization we discuss briefly the asymptotic approximation of the derived general formulation. We point out its direct correlation to known analytic approximations.

In Chapter 6 we derive the algorithms for the time-dependent Schrödinger equation and the wide angle one-way equations. Since our focus here is on the time-evolution, we consider only problems with one transversal (spatial) direction. First we study the Schrödinger equation and derive the basic technique which allows

us to extend our method to time-discrete problems. In Section 6.2 we generalize this concept to the wide angle one-way approximations of the Helmholtz equation. Finally, we give some typical applications of Schrödinger-type scattering problems along with a study of their numerical properties.

## A Survey of Different Classical Methods

A number of quite different methods are available to solve scattering problems of the described types. It seems to be impossible to give a single satisfactory classification scheme which sets all the methods in relation to each other, with all their pros and cons. However, most important here is to discuss the relation between the discrete method developed in this paper and the most successful classical methods. Here use the term “classical” even for relatively new methods such as the the infinite-element method. Fig. 2.0.1 attempts to show the relation between the different methods postulating that a continuous form of a radiation condition provides the starting point for all methods. According to this classification the common strategy of all methods is to construct a representation scheme for the exterior solutions such that all solutions both satisfy the continuous Helmholtz equation and the continuous radiation condition. Here the radiation condition is of fundamental importance. It is constructed such that it allows only for “outgoing” or “incoming” solutions.

### 2.1. Factorization Based on Separable Coordinates

For simplicity, we consider only the solution of the Helmholtz equation exterior to a circle with radius  $a$ . To get the point, we leave out some details needed to justify the operations. Every function  $u \in L^2(\partial\Omega)$ ,  $\partial\Omega = \{x \in \mathbb{R}^2 : |x| = a\}$  can be expressed via its Fourier series expansion

$$u(a, \phi) = \sum_{\nu=-\infty}^{\infty} a_{\nu} e^{i\nu\phi}, \quad a_{\nu} = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-i\nu\phi} u(a, \phi) d\phi.$$

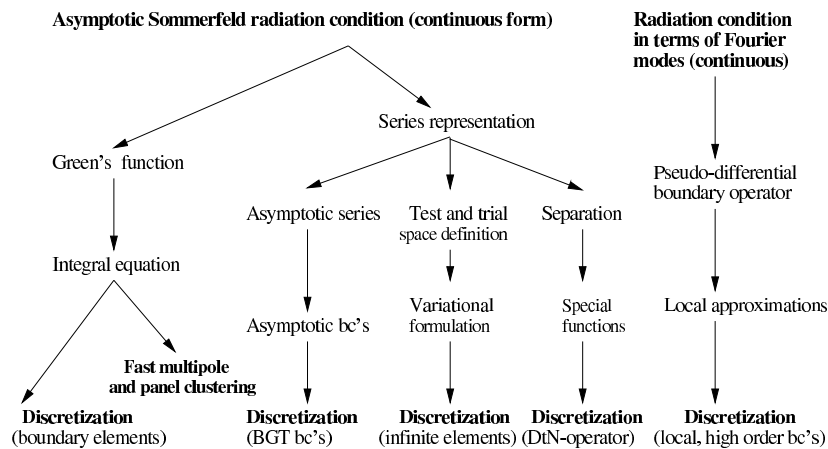


FIGURE 2.0.1. Classical approaches to derive transparent boundary conditions

A separation ansatz in cylindric coordinates  $u(r, \phi) = \sum_{\nu} a_{\nu} e^{i\nu\phi} u_{\nu}(r)$  shows that  $u_{\nu}(r)$  must obey Bessel's differential equation

$$\partial_r^2 u_{\nu} + \frac{1}{r} \partial_r u_{\nu} + \left( k^2 - \frac{\nu^2}{r^2} \right) u_{\nu} = 0, \quad r > a.$$

The fundamental solution and its derivative are

$$\begin{aligned} u_{\nu}(r) &= c_1 H_{\nu}^{(1)}(kr) + c_2 H_{\nu}^{(2)}(kr) \\ u'_{\nu}(r) &= c_1 H_{\nu}^{(1)'}(kr) + c_2 H_{\nu}^{(2)'}(kr) \end{aligned}$$

where  $H_{\nu}^{(1)}, H_{\nu}^{(2)}$  are Hankel's functions of the first and the second kind, respectively, and the prime denotes the derivative with respect to the argument. From an asymptotic study of Hankel's functions we know that only the first kind obeys the Sommerfeld condition. In order to drop the Hankel function of the second kind, we must establish an additional condition, namely,

$$u'_{\nu}(r) = k \frac{H_{\nu}^{(1)'}(kr)}{H_{\nu}^{(1)}(kr)} u_{\nu}(r),$$

or, equivalently,

$$u'_{\nu}(r) = \left( \frac{\nu}{r} - k \frac{H_{\nu+1}^{(1)}(kr)}{H_{\nu}^{(1)}(kr)} \right) u_{\nu}(r).$$

Taking the angular-dependent part into account, we obtain the DtN-map

$$\partial_n u(r, \phi)|_{r=a} = \sum_{\nu=-\infty}^{\infty} \left( \frac{\nu}{r} - k \frac{H_{\nu+1}^{(1)}(kr)}{H_{\nu}^{(1)}(kr)} \right)_{r=a} \frac{e^{i\nu\phi}}{2\pi} \langle e^{i\nu\phi}, u(a, \phi) \rangle_{L^2(\partial\Omega)}.$$

The main ingredients, with respect to this classical approach, have been the expression of the fundamental solution in terms of a superposition of Hankel's functions and the decision, which of them satisfy Sommerfeld's radiation condition. For example, if we would start with another pair of fundamental solutions, say with Bessel and Weber functions, none of them would satisfy the Sommerfeld condition and we must try to find a superposition of both, namely the Hankel functions, which obey the radiation condition.

This technique has been introduced by Givoli and Keller for the Laplace equation, the Helmholtz equation and the elastic wave equation in a series of papers, cf. e. g. [38, 37, 35, 64], and has been extended later by Givoli, Grote and Keller to time-dependent problems (e. g. wave equation, Maxwell's equations), see the corresponding articles [36, 43, 42, 44, 45].

Some of the features of the method are:

- The DtN-number  $u'_{\nu}(a)/u_{\nu}(a)$  corresponding to a given mode number  $\nu$  and a distance  $a$  is derived in straightforward way.
- In standard cases like the one above the implementation does not cause a problem, since the necessary special functions are available from standard libraries.
- Efficient tools to couple the non-local Fourier modes to local interior discretizations are available.

The obvious drawbacks are:

- The problem must have a separable structure. In general, this is not the case for variable coefficient problems.
- The generalization to other than cylindric and elliptic coordinate systems and their 3D generalizations is difficult.

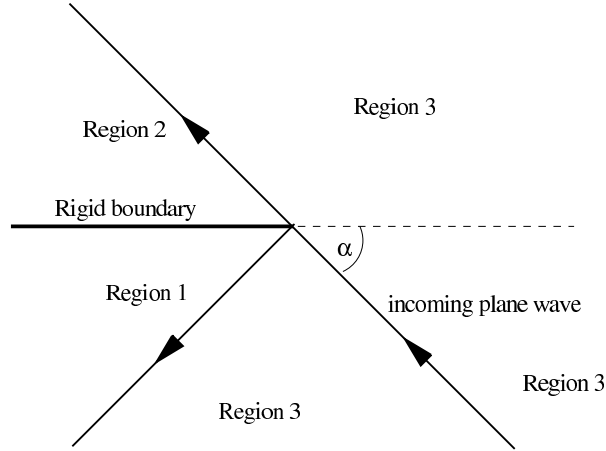


FIGURE 2.2.1. Diffraction problem

## 2.2. Boundary Integral/Element Methods

We describe the concept of boundary integral methods (boundary element methods) by means of a particular example: the above mentioned classical diffraction problem of a plane wave by a rigid semi-infinite plane sheet, Example 1.2.1, p. 12. A representation of boundary integral methods for the Helmholtz equation for problems with bounded scatterers may be found in [3], for the coupling between finite elements and boundary elements see, e.g. [2]. Our presentation follows the ones given in [24, Sec. 8.7] and [23, Sec. 5.6, Chap. 8]. An alternative operator theoretical approach of the same subject can be found in [73].

We decompose the whole procedure into the three main steps discussed briefly in connection with Example 1.2.1:

- (1) Decomposition of the full field into three parts: The incoming plane wave  $u_{\text{in}} = \exp(i\mathbf{k}\mathbf{x})$ , the reflected wave  $u_{\text{ref}}$  and the diffracted part  $u_{\text{diff}}$ .
- (2) Solution representation by a properly constructed Green's function.
- (3) Solution of the integral equation.

Without going too much into the details, we describe each of the steps.

**Field decomposition.** The situation is illustrated in Fig. 2.2.1. The incoming plane wave strikes the semi-infinite plane. The incoming field already satisfies the Helmholtz equation. It excites a reflected field, which is assumed to live mainly in Region 1, and a diffracted field, which lives in all the regions 1, 2, and 3. The entire field, the incoming plus the fields excited by the rigid boundary (the scattered field), must satisfy the Helmholtz equation. The problem to formulate a Green's function representation of the scattered field now is that it obviously does not meet Sommerfeld's radiation condition (1.0.2), since the reflected field satisfies it only in one but not in all directions. In order to circumvent this difficulty, that is to extract the part of the field which in fact satisfies the Sommerfeld condition, one considers the solutions predicted by geometrical optics first, see Fig. 2.2.1:

$$\begin{aligned} \text{Region 1} & : u = u_{\text{in}} + u_{\text{ref}} \\ \text{Region 2} & : u = 0 \\ \text{Region 3} & : u = u_{\text{in}}. \end{aligned}$$

From geometrical optics we know the explicit formula for the reflected field. If the incoming plane wave is  $u_{\text{in}} = \exp(ik_x x + ik_y y)$  then  $u_{\text{ref}} = \exp(ik_x x - ik_y y)$  in Region 1. We use these fields as initial guess of the solution of the Helmholtz

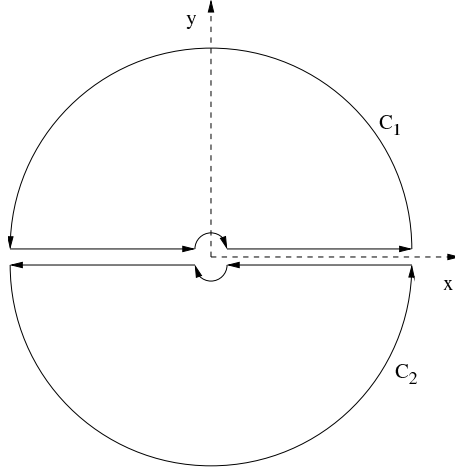


FIGURE 2.2.2. Semi-circles used for Green's identity

equation, which is corrected by the diffracted part  $u_{\text{diff}}$  according to

$$\begin{aligned} \text{Region 1} & : u = u_{\text{in}} + u_{\text{ref}} + u_{\text{diff}} \\ \text{Region 2} & : u = u_{\text{diff}} \\ \text{Region 3} & : u = u_{\text{in}} + u_{\text{diff}}. \end{aligned}$$

Now, we require that  $u_{\text{diff}}$  satisfies Sommerfeld's radiation condition. Additionally we require  $\partial_y u = 0$  on the rigid boundary, since the normal derivative of the acoustic velocity potential supplies the normal velocity of the particles, which must be zero on a sound-hard boundary. If we compare this decomposition step with the corresponding discrete solution based on the pole condition, which we will present in 5.1.2, pp. 102, we find that no comparable procedure is necessary. The only point, where in our approach some a-priori knowledge about the solution behavior plays a role, is in the definition of the fields along artificial vertical boundaries far away from the origin. Here, we will use again informations from geometrical optics.

**Integral representation of the solution.** We construct a suitable Green's function from the reflection principle, based on Green's function for the free-space problem:

$$\begin{aligned} G(\xi, \eta, x, y) &= \frac{i}{4} H_0^{(1)} \left( k \sqrt{(x - \xi)^2 + (y - \eta)^2} \right) \\ &\quad + \frac{i}{4} H_0^{(1)} \left( k \sqrt{(x - \xi)^2 + (y + \eta)^2} \right). \end{aligned}$$

This Green's function satisfies the Helmholtz equation except at  $(x, y) = (\xi, \pm\eta)$ , due to the Hankel functions Sommerfeld's radiation condition and  $\partial_y G(\xi, \eta, x, 0) = 0$ , which is the imposed boundary condition for  $x \leq 0$ . Next, we apply Green's identity to the whole field in the upper half of the plane  $u^{(1)}$  using the contour  $C_1$ , see Fig. 2.2.2. This yields

$$\begin{aligned} u^{(1)}(\xi, \eta) &= \oint_{c_1} ds \left( G \frac{du^{(1)}}{dn} - u^{(1)} \frac{dG}{dn} \right) \\ &= - \int_{-\infty}^{\infty} dx G \frac{du^{(1)}}{dy} \Big|_{y=+0} \end{aligned}$$

The latter equation follows because the integrals on the small and the large semi-circles tend to zero if their radii go to zero and infinity, respectively, see [23, Sec. 5.6, pp. 264-265]. Next, we derive a second integral equation, applying Green's identity to the whole field in the lower half of the plane minus the incoming plane wave and the reflected wave. We obtain

$$u^{(2)}(\xi, \eta) - u_{\text{in}}(\xi, \eta) - u_{\text{refl}}(\xi, \eta) = \int_{-\infty}^{\infty} dx G \frac{du^{(2)}}{dy} \Big|_{y=-0}.$$

We let  $\eta$  approach  $\pm 0$ , take  $u^{(1)}(\xi, 0) = u^{(2)}(\xi, 0)$  for  $\xi > 0$  into account, and take the difference between the two integral equations into account to obtain

$$f(\xi) = i \int_{-\infty}^{\infty} dx H_0^{(1)}(k|x - \xi|) \partial_y u(x, 0),$$

where

$$f(\xi) = \begin{cases} -2e^{ik_x \xi}, & \xi > 0 \\ u(\xi, -0) - u(\xi, +0) - 2e^{ik_x \xi}, & \xi < 0. \end{cases}$$

Integral equations of this type are called Wiener-Hopf integral equations. Their difficulty lies in the fact that the function  $f$  is not completely known (namely for  $\xi < 0$ ). However, the function  $\partial_y u(\xi, 0)$  under the integral sign is already known in parts (namely for  $\xi > 0$ ).

If we compare this step with our alternative approach from Section 5.1.2, we find that we obtain an algebraic system of exactly the same structure, a linear system of the type  $Ax = b$ , where the right-hand side  $b$  is not completely known, whereas  $x$  is already partially known. However, we obtain this system without applying the techniques described above.

**Solution of the integral equation.** There are several ways to solve the Wiener-Hopf equation. The first one is to reformulate this integral equation by analytical means such that a classical Fredholm integral equation of the first kind results. The technique to do this is known as Wiener-Hopf technique. It is based on a Fourier transform of the convolution integral and certain factorizations of the Fourier-transformed kernel and the known part of the right-hand side of the integral equation. Applying analytic continuation techniques and Liouville's theorem supplies the desired representation. However, finally we are left with an integral equation, which must be solved, preferably by numerical methods. Therefore we can, as a second way to solve the problem, directly apply numerical discretization techniques to the Wiener-Hopf equation.

To formalize this approach, we restrict the infinite interval, on which the integral equation is defined to the finite interval  $[-a, a]$  and set the approximate boundary conditions known from the geometrical optics model to  $u^{(1)}(-a) = 0$ ,  $u^{(1)}(a) = u^{(2)}(a) = u_{\text{in}}(a)$ ,  $u^{(2)}(-a) = u_{\text{in}}(-a) + u_{\text{refl}}(-a)$ . Setting  $v(x) := \partial_y u(x, 0)$  we have to solve the linear Fredholm integral equation

$$\int_{-a}^a dx K(x, \xi) v(x) = f(\xi),$$

where the kernel  $K(x, \xi) = H_0^1(k|x - \xi|)$  is in  $L^2[-a, a]^2$  due to the logarithmic singularity  $H_0^1(k|x - \xi|) \rightarrow \text{const} \ln|x - \xi|$  as  $(x - \xi) \rightarrow 0$ . The integral operator is given by the formula

$$(\mathcal{K}u)(\xi) := \int_{-a}^a dx K(x, \xi) u(x)$$



for all  $u \in L^2[-a, a]$  and all  $\xi \in [-a, a]$ . The Nyström, or quadrature, method replaces this integral operator by a finite-dimensional approximation

$$(\mathcal{K}_h u)(\xi) := \sum_{j=1}^n w_j K(x_j, \xi) u(x_j)$$

for all  $u \in L^2[-a, a]$  and all  $\xi \in [-a, a]$  and given points  $x_1, \dots, x_n \in [-a, a]$ . Ideally, the real weights  $w_1, \dots, w_n$  are chosen such that integration along the desired solution is approximated very well.

We require that the approximation be exact at all the points  $x_j$  (collocation method), which supplies the algebraic system

$$(\mathcal{K}_h u)(\xi_j) = f(\xi_j), \quad j = 1, \dots, n.$$

If we again compare this last step with with the formulation of our discrete method, we find a perfect coincidence with respect to the structure of the algebraic system. However, the derivation is quite different. Instead of using a discretization of an integral equation, we will use a direct discretization of the Helmholtz operator to derive the discrete system.

REMARK. In general, integral equations of the first kind are ill-posed problems. To solve such problems numerically, regularization techniques has to be applied. Since this has no effect on our elaboration in the sequel, we will not discuss this point further. The special Wiener-Hopf structure is not characteristic for the class of integral equation methods. It is a property which results from our diffraction problem. On the other hand, the main steps of the above procedure are characteristic: (1) The construction of a representation formula, (2) The restriction of the representation formula to the boundary of an obstacle or to an artificial boundary, (3) Numerical solution of an integral equation.

### 2.3. Infinite Element Methods

The infinite element method belongs to the most attractive methods for the solution of exterior Helmholtz scattering problems.

To keep the setting simple, we consider only the scattering exterior to the unit sphere, i. e., we set  $\Omega = \{\mathbf{x} \in \mathbb{R}^3 : |\mathbf{x}| < a\}$  and assume given Neumann data on the boundary:

$$(2.3.1) \quad \begin{aligned} \Delta u + k^2 u &= 0 && \text{in ext } \Omega \\ \partial_n u &= g && \text{on } \partial\Omega \\ \partial_r u - iku &= o\left(\frac{1}{r}\right) && (r = |\mathbf{x}| \rightarrow \infty) \text{ uniformly in all directions.} \end{aligned}$$

In our representation here we follow the infinite element method proposed by Astley et. al., following the presentation given by Ihlenburg in [61].

First, the scattering problem is rewritten in variational form: Find  $u \in H_{\mathbb{W}}^1(\text{ext } \Omega)$  such that

$$b(v, u) = \langle v, g \rangle_{\partial\Omega} \quad \text{for all } v \in H_{\mathbb{W}^*}^1(\text{ext } \Omega)$$

Here, the sesquilinear forms  $b(\cdot, \cdot)$  and  $\langle \cdot, \cdot \rangle_{\partial\Omega}$  are defined by

$$\begin{aligned} b(v, u) &:= \int_{\text{ext } \Omega} (\overline{\nabla v} \cdot \nabla u - k^2 \overline{v} u) \, dx \\ \langle v, g \rangle_{\partial\Omega} &:= \int_{\partial\Omega} \overline{v} g \, dx \end{aligned}$$

and

$$\begin{aligned} H_w^1(\text{ext } \Omega) &= \left\{ u : \int_{\text{ext } \Omega} \frac{1}{r^2} (|\nabla u|^2 + |u|^2) dx + \int_{\text{ext } \Omega} |\partial_r u - iku|^2 dx < 0 \right\} \\ H_{w^*}^1(\text{ext } \Omega) &= \left\{ u : \int_{\text{ext } \Omega} r^2 (|\nabla u|^2 + |u|^2) dx < 0 \right\} \end{aligned}$$

are the trial and the test space, respectively. The subscript  $w$  indicates that the test and trial space defined this way are *weighted* Sobolev spaces. Based on the theoretical framework of Leis [68] this variational form was proposed by Demkovicz and Gerdes [22]. The existence and uniqueness of a solution of the variational problem is known. To construct a convergent discrete method, sequences of finite dimensional subspaces  $V_w^1 \subset V_w^2 \subset \dots \subset V_w^N \subset \dots \subset H_w^1(\text{ext } \Omega)$  and  $V_{w^*}^1 \subset V_{w^*}^2 \subset \dots \subset V_{w^*}^N \subset \dots \subset H_{w^*}^1(\text{ext } \Omega)$  must be constructed. The trial space must be such that the approximation property  $\inf_{v^N \in V_w^N} \|u - v^N\|_{H_w^1} \rightarrow 0$  as  $N \rightarrow \infty$  for each  $u \in H_w^1(\text{ext } \Omega)$  holds true. The test space, together with the trial space and the variational formulation must guaranty that the discrete solution process is stable, i. e., the condition number of the discrete problem must be bounded. The practical construction of the trial and test spaces is based on the Wilcox expansion theorem [90], which also supplies the basis for the Bayliss-Gunzburger-Turkel boundary conditions. Since we will come back several times to the Wilcox theorem and its 2D counterpart, the Karp theorem, we state them here explicitly.

**THEOREM 2.3.1.** *(The Wilcox Expansion Theorem [90].) Let  $u \in C^2$  be a solution of (2.3.1) for the region exterior to a sphere  $|\mathbf{x}| = a$ , and let  $(r, \theta, \phi)$  be the spherical coordinates for  $\mathbf{x}$ . Then*

$$u(r, \theta, \phi) = \frac{e^{ikr}}{kr} \sum_{j=0}^{\infty} \frac{F_j(\theta, \phi)}{(kr)^j}$$

where the series converges for  $r > a$  and converges absolutely and uniformly in  $r, \theta$  and  $\phi$  in any region  $|\mathbf{x}| \geq a + \epsilon > a$ . The series may be differentiated term by term with respect to  $r, \theta$  and  $\phi$  any number of times and the resulting series all converge absolutely and uniformly.

The corresponding theorem for two dimensions has been proved by Karp in 1961:

**THEOREM 2.3.2.** *(The Karp Expansion Theorem [63].) Let  $u \in C^2$  be a solution of (2.3.1) for the region exterior to a circle  $|\mathbf{x}| = a$ , and let  $(r, \theta)$  be the cylindrical coordinates for  $\mathbf{x}$ . Then there exists a convergent expansion, valid for  $r > a$ , in the form*

$$u = H_0^{(1)}(kr) \sum_{j=0}^{\infty} \frac{F_j(\theta)}{(kr)^j} + H_1^{(1)}(kr) \sum_{j=0}^{\infty} \frac{G_j(\theta)}{(kr)^j}$$

where the series converges uniformly and absolutely for  $|\mathbf{x}| \geq a + \epsilon > a$ , and may be differentiated term-wise with respect to  $r$  as often as desired.

Thus, the key-point in infinite element methods is to have an expansion theorem valid on the unbounded domain, together with a variational formulation which supplies stable discrete problems. In contrast, our Laplace domain methods derived from the pole condition are not based on series representations like the ones above (even not theoretically). Instead, the field representation on the unbounded domain is computed simultaneously with the interior solution.

## 2.4. Asymptotic Methods

In practice asymptotic approximations are used very often. The two main categories of asymptotic methods are the one where the distance between the obstacle and the artificial boundary tends to infinity and the one where the wavenumber tends to infinity. Especially the latter class allows to include problems with variable coefficients. We emphasize this here, since this is a key feature of the pole condition approach, developed in the main part of this work, too. The other key feature of the asymptotic methods is that the derived asymptotic boundary conditions are local conditions.

**Far-Field Approximation (Bayliss-Gunzburger-Turkel).** In [10] Bayliss and Turkel and in [9] Bayliss, Gunzburger and Turkel developed a concept of asymptotic boundary conditions (BGT boundary conditions) characterized by the central property that these boundary conditions can be associated with local boundary conditions on the artificial boundary. Especially when used with direct sparse solvers, the BGT-type boundary conditions have the advantage as opposed to boundary conditions based on standard integral formulations because that they do not destroy the sparsity (at least in their lowest order form) of the FEM-discretization of the interior problem.

The starting point of the derivation of the BGT boundary conditions is, exactly as in the case of infinite elements, the existence of a convergent Wilcox expansion in 3D [90], cf. Theorem 2.3.1 and the Karp expansion in 2D [63], cf. Theorem 2.3.2. We consider the derivation of the BGT boundary conditions for the 2D case. Since the series expansion from Theorem 2.3.2 is difficult to work with, it is replaced by a series of the simpler 3D type of Theorem 2.3.1 using the well-known asymptotic expansion of the Hankel functions (cf. Corollary 4.4.5 in Section 4.2, p. 84)

$$(2.4.1) \quad u(r, \phi) \sim \sqrt{\frac{2}{\pi kr}} e^{i(kr - \pi/2)} \sum_{j=0}^{\infty} \frac{f_j(\phi)}{r^j}.$$

Note, that this manipulation results in a series which is no longer convergent. Now the central idea of the BGT boundary conditions is the construction of a family of boundary operators  $B_m$ ,  $m = 1, 2, \dots$ , such that the action of the operator  $B_m(r, \partial_r)$  applied to the series (2.4.1) results in

$$(2.4.2) \quad B_m(r, \partial_r) u(r, \phi) \sim \sqrt{\frac{2}{\pi kr}} e^{i(kr - \pi/2)} \sum_{j=2m}^{\infty} \frac{f_j(\phi)}{r^j}.$$

Hence the operator is designed such that it annihilates the leading  $2m$  terms (the terms with  $j = 0, \dots, 2m - 1$ ) of the series expansion. It is not immediately clear that the operator  $B_m$  defines an asymptotic boundary condition. To discuss this point we consider a parameterization

$$\begin{aligned} \gamma : [0, 1] &\rightarrow \mathbb{R}_+ \times [-\pi, \pi] \\ s &\mapsto (r, \phi) \in \partial\Omega \end{aligned}$$

of the artificial boundary  $\partial\Omega$  in polar coordinates. We set  $B_m(r, \partial_r) u(r, \phi)|_{r=r(s)} = 0$  for each parameter  $s \in [0, 1]$  which supplies a relation between the function  $u(r(s), \phi(s))$  at each point on the boundary (the Dirichlet data) and radial derivatives at this point. Thus, asymptotic boundary conditions of the form  $B_m u = 0$  yield local boundary conditions containing, in general, higher order derivatives in radial direction.

To construct the family of operators  $B_m$  we use the identities

$$(2.4.3) \quad \left( \partial_r - ik + \frac{j + \frac{1}{2}}{r} \right) \frac{e^{ikr}}{r^{j+1/2}} = 0$$

$$(2.4.4) \quad \left( \partial_r - ik + \frac{j + \frac{1}{2}}{r} \right) \frac{e^{ikr}}{r^{j+1+1/2}} = -\frac{e^{ikr}}{r^{j+2+1/2}}.$$

We set

$$B_1 = \partial_r - ik + \frac{1}{2r}$$

which annihilates the zeroth and the first order term of the asymptotic expansion (2.4.1), i.e., it holds  $B_1 u = O(1/r^{2+1/2})$  for  $r \rightarrow \infty$ . Based on (2.4.3) and (2.4.4) we find that the application of  $(\partial_r - ik + (j + 1/2)/r)$  annihilates both the  $j$ th and the  $j + 1$ th terms in the series (2.4.1) yielding the recurrence

$$(2.4.5) \quad B_j = \left( \partial_r - ik + \frac{2j - \frac{3}{2}}{r} \right) B_{j-1}, \quad j \geq 1.$$

with an asymptotic behavior

$$B_m u = \mathcal{O}\left(\frac{1}{r^{2m+1/2}}\right), \quad m \geq 1 \quad (r \rightarrow \infty).$$

We add the results for three dimensions. The same procedure applied to the Atkinson-Wilcox expansion supplies the recurrence

$$B_j = \left( \partial_r - ik + \frac{2j - 1}{r} \right) B_{j-1}, \quad j \geq 1$$

with  $B_0 := 1$  and an asymptotic behavior

$$B_m u = \mathcal{O}\left(\frac{1}{r^{2m+1}}\right), \quad m \geq 1 \quad (r \rightarrow \infty).$$

**Operator Factorization.** In two important papers, Engquist and Majda [28, 29] constructed a family of transparent boundary conditions for the wave equation. An error analysis of this scheme specialized to the constant coefficient scalar wave equation, carried out by Halpern and Rauch [54], shows that the error consists of two terms. One is proportional to the largest reflection coefficient for the artificial boundary condition, and can be made arbitrary small in case of a planar boundary. The second term is inverse proportional to the average frequency present in the solution. Roughly speaking, the Engquist-Majda construction becomes the better the smoother the artificial boundary and the higher the frequency of the wave. Since their method is very general and can take into account both variable coefficients and general curved boundaries, we repeat the main steps of the construction here. Instead of summarizing the full theory in the following, we will introduce only the very basic quantities needed to motivate and develop the desired technique of operator factorization. The applicability and limitations of the method are studied using the Helmholtz equation with constant coefficients and a fixed angular frequency in the exterior of a circle as a model problem. We introduce the method here for two reasons:

- (1) The method is applicable to very general situations. It allows arbitrary curved, smooth boundaries and variable coefficients in the exterior domain.
- (2) The method is directly applicable to the semi-discrete formulation of the scattering problem to be derived in Section 5.3. Hence it might be considered as one of the several possible approximate realizations of our exact solution of the basic semi-discrete formulation.

Since we are interested in the solution of the Helmholtz equation, we specialize the general consideration of [28, 29] to the time-harmonic case. Further we use the Helmholtz operator  $L$  in cylindric coordinates in 2D by

$$L := \partial_r^2 + \frac{1}{r} \partial_r + \frac{1}{r^2} \partial_\phi^2 + n^2(\phi) k^2.$$

In the special variable coefficient case which we want to consider here, we allow a variation of the wavenumber with respect to the angular coordinate  $\phi$ , denoted by the scalar, real refractive index function  $n(\phi)$ . We want to compute a family of boundary conditions at  $r = r_0$ . The key tool to derive the boundary conditions is a factorization of the differential operator  $L$  by methods arising on the theory of pseudodifferential operators. Before we study the method, we recall some basic facts about pseudodifferential operators which are used by Engquist and Majda, see the book of Taylor [87]. Our presentation here follows the one given in [29] for the variable coefficient case. We need the definition of a symbol class, the construction of the associated pseudodifferential operators and the asymptotic expansion of the symbol of composite pseudodifferential operators.

Let  $\Omega$  be an open subset of  $\mathbb{R}^n$ . The complex, infinitely differentiable function  $b(\mathbf{x}, \boldsymbol{\xi}) \in C^\infty(\Omega \times \mathbb{R}^n)$  is an element of the symbol class  $S^m(\Omega)$  provided there exists constants  $C_{\alpha, \beta}$  such that

$$(2.4.6) \quad \left| \frac{\partial^{|\alpha|}}{\partial \mathbf{x}^\alpha} \frac{\partial^{|\beta|}}{\partial \boldsymbol{\xi}^\beta} b(\mathbf{x}, \boldsymbol{\xi}) \right| \leq C_{\alpha, \beta} (1 + |\boldsymbol{\xi}|)^{m - |\beta|}$$

for any multi-indices  $\alpha, \beta$ , ([87, Chapt. II, Def. 1.1]). Associated with any symbol  $b(\mathbf{x}, \boldsymbol{\xi})$  is a pseudodifferential operator  $b(\mathbf{x}, (1/i) \partial/\partial \mathbf{x}) \in OP(S^m)$  defined via the Fourier transform by

$$(2.4.7) \quad b \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) v(\mathbf{x}) = \int \exp(i\mathbf{x} \cdot \boldsymbol{\xi}) b(\mathbf{x}, \boldsymbol{\xi}) \widehat{v}(\boldsymbol{\xi}) d\boldsymbol{\xi}.$$

Conversely, each pseudodifferential operator  $b(\mathbf{x}, (1/i) \partial/\partial \mathbf{x})$  has a symbol  $\sigma(b(\mathbf{x}, (1/i) \partial/\partial \mathbf{x})) = b(\mathbf{x}, \boldsymbol{\xi})$ . Note that the symbol of a pseudodifferential operator does not contain any differential operator. Given two pseudodifferential operators  $a(\mathbf{x}, (1/i) \partial/\partial \mathbf{x}) \in OP(S^{m_1})$  and  $b(\mathbf{x}, (1/i) \partial/\partial \mathbf{x}) \in OP(S^{m_2})$  it holds  $a \circ b \in OP(S^{m_1+m_2})$  and the symbol of the composite pseudodifferential operator  $a \circ b$  possesses the asymptotic expansion ([87, Chapt. II, Theorem 4.4])

$$(2.4.8) \quad \sigma(a \circ b) \sim \sum_{|\alpha| \geq 0} \frac{(-i)^{|\alpha|}}{\alpha!} \frac{\partial^{|\alpha|}}{\partial \boldsymbol{\xi}^\alpha} a(\mathbf{x}, \boldsymbol{\xi}) \frac{\partial^{|\alpha|}}{\partial \mathbf{x}^\alpha} b(\mathbf{x}, \boldsymbol{\xi}) \quad (|\boldsymbol{\xi}| \rightarrow \infty).$$

In our special case we define  $\mathbf{x}$  by 2D cylindric coordinates  $\mathbf{x} = (r, \phi)$  and consider the exterior of a circle with radius  $r_0$ , that is  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : r > r_0, -\pi < \phi \leq \pi\}$ . Let us denote the Fourier coordinates dual to the original coordinates by  $\boldsymbol{\xi} = (\rho, \nu)$ .

The key of the Engquist-Majda approach is the construction of pseudodifferential operators  $\lambda_+(\mathbf{x}, (1/i) \partial/\partial \mathbf{x})$  and  $\lambda_-(\mathbf{x}, (1/i) \partial/\partial \mathbf{x})$  such that the factorization

$$(2.4.9) \quad L = (\partial_r - \lambda_+) (\partial_r - \lambda_-) + (\text{smooth error})$$

holds true. This factorization requires the computation of the symbol of the composite pseudodifferential operators  $\sigma(\lambda_+(\mathbf{x}, (1/i) \partial/\partial \mathbf{x}) \partial_r)$  and  $\sigma(\partial_r \lambda_-(\mathbf{x}, (1/i) \partial/\partial \mathbf{x}))$ ,

where  $\lambda$  is either  $\lambda_+$  or  $\lambda_-$ . It follows directly from (2.4.7) that

$$\begin{aligned} \left[ \lambda \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \partial_r \right] u(\mathbf{x}) &= \int \exp(i\mathbf{x} \cdot \boldsymbol{\xi}) \lambda(\mathbf{x}, \boldsymbol{\xi}) \widehat{\partial_r u} d\boldsymbol{\xi} \\ &= \int \exp(i\mathbf{x} \cdot \boldsymbol{\xi}) [\lambda(\mathbf{x}, \boldsymbol{\xi}) i\xi] \widehat{u}(\boldsymbol{\xi}) d\boldsymbol{\xi}. \end{aligned}$$

Similarly, we have

$$\begin{aligned} \partial_r \left[ \lambda \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) u(\mathbf{x}) \right] &= \left[ \partial_r \lambda \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \right] u(\mathbf{x}) + \left[ \lambda \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \partial_r \right] u(\mathbf{x}) \\ &= \int \exp(i\mathbf{x} \cdot \boldsymbol{\xi}) [\partial_r \lambda(\mathbf{x}, \boldsymbol{\xi}) + \lambda(\mathbf{x}, \boldsymbol{\xi}) i\xi] \widehat{u}(\boldsymbol{\xi}) d\boldsymbol{\xi}. \end{aligned}$$

This shows the identities

$$(2.4.10) \quad \sigma \left( \lambda \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \partial_r \right) = \lambda(\mathbf{x}, \boldsymbol{\xi}) i\xi$$

$$(2.4.11) \quad \sigma \left( \partial_r \lambda \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \right) = \lambda(\mathbf{x}, \boldsymbol{\xi}) i\xi + \partial_r \lambda(\mathbf{x}, \boldsymbol{\xi}).$$

Hence we have recovered a very special form of the asymptotic expansion (2.4.8), where we can replace the asymptotic equality by the equality itself.

We shall show experimentally that we can construct sequences of pseudodifferential operators  $\lambda_+^1, \lambda_+^0, \dots, \lambda_+^{-j}$  and  $\lambda_-^1, \lambda_-^0, \dots, \lambda_-^{-j}$ ,  $1 \geq j \geq -1$  such that the following crucial asymptotic properties hold true, with  $u(r) = H_\nu^{(1)}(knr)$  the exact solution of Bessel's equation:

i) *Far field asymptotics.*

$$\frac{\left( \partial_r - \sum_{k=-1}^j \lambda_-^{-k} \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \right) u(r)}{u(r)} = \mathcal{O} \left( \frac{1}{r^{j+3}} \right) \quad \text{for } r \rightarrow \infty \text{ and fixed } \nu.$$

ii) *High frequency asymptotics in  $\nu$ .*

$$\frac{\left( \partial_r - \sum_{k=-1}^j \lambda_-^{-k} \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \right) u(r)}{u(r)} = \mathcal{O} \left( \frac{1}{|\nu|^{2j+1}} \right) \quad \text{for } |\nu| \rightarrow \infty \text{ and fixed } r,$$

and  $\text{Im} \left\{ \sigma \left( \lambda_-^j \right) \right\} > 0$  as  $r \rightarrow \infty$  and  $\nu$  fixed. Once a function  $u(\mathbf{x})$  in  $\Omega$  is known which satisfies  $(\partial_r - \lambda_-)u = 0$  it follows  $Lu \cong 0$ , where  $\cong$  means equal up to a smooth error. Viewing  $\lambda_-$  as a constant with positive imaginary part, the equation  $(\partial_r - \lambda_-)u = 0$ , derived from (2.4.9), possesses a solution  $\exp(\lambda_- r)$  showing that  $u(\mathbf{x})$  is an outgoing solution with respect to  $r$ . Hence, if a function  $u(\mathbf{x})$  satisfies  $(\partial_r - \lambda_-)u = 0$ , it satisfies  $Lu \cong 0$  and consists only of outgoing functions. Consequently, the restriction of  $\mathbf{x}$  to the boundary supplies the desired boundary condition (up to the factorization error)

$$\left( \partial_r - \lambda_- \left( \mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}} \right) \right) u(\mathbf{x}) = 0 \quad \mathbf{x} \in \partial\Omega.$$

To obtain a family of factorizations of ascending degree of accuracy, we construct sequences of pseudodifferential operators  $\lambda_+^1, \lambda_+^0, \dots, \lambda_+^{-j}$  and  $\lambda_-^1, \lambda_-^0, \dots, \lambda_-^{-j}$ ,

$j \geq -1$ , such that  $\sigma(\lambda_{\pm}^{-j}) \in S^{-j}$  and the symbol of the residual satisfies

$$(2.4.12) \quad \sigma\left(L - \left(\partial_r - \sum_{k=-1}^j \lambda_+^{-k}\right) \left(\partial_r - \sum_{k=-1}^j \lambda_-^{-k}\right)\right) \in S^{-j}.$$

The idea behind this is to approximate  $\lambda_-$  by  $\sum_{k=-1}^j \lambda_-^{-k}$  such that the terms  $\geq j$  are all smoother than the initial terms since  $\sigma(\lambda_-^{-j-k}) = \mathcal{O}(|\xi|^{-j-k})$  for  $k \geq 1$ , by definition of the symbol class (2.4.6). Consequently we construct a method adapted to situations where high frequencies play the dominant role. First we compute the symbol of the Helmholtz operator

$$(2.4.13) \quad \sigma(L) = -\rho^2 + \frac{1}{r}i\rho - \frac{1}{r^2}\nu^2 + n^2k^2.$$

Its principal part consists of the terms quadratic in  $\rho, \nu, k$  (since we view it as time-harmonic wave equation) and possesses a factorization

$$-\rho^2 - [(1/r^2)\nu^2 - n^2k^2] = \left(i\rho + \sqrt{(1/r^2)\nu^2 - n^2k^2}\right) \left(i\rho - \sqrt{(1/r^2)\nu^2 - n^2k^2}\right),$$

where we take the branch of the square root such that its real part is positive and  $\sqrt{-1} = +i$ . Motivated by this factorization we choose

$$\lambda_{\pm}^1\left(\mathbf{x}, \frac{1}{i} \frac{\partial}{\partial \mathbf{x}}\right) = \mp \sqrt{-\frac{1}{r^2} \partial_{\phi}^2 - n^2k^2}.$$

Its symbol is

$$(2.4.14) \quad \lambda_{\pm}^1(\mathbf{x}, \boldsymbol{\xi}) = \mp \sqrt{\frac{1}{r^2}\nu^2 - n^2k^2}$$

and is bounded by

$$|\lambda_{\pm}^1(\mathbf{x}, \boldsymbol{\xi})| \leq C(1 + |\boldsymbol{\xi}|), \quad C = \max\left(nk, \frac{1}{r}\right).$$

However, already the first derivative  $\partial_{\nu} \lambda_{\pm}^1(\mathbf{x}, \boldsymbol{\xi})$  generates a factor  $1/\sqrt{\nu^2/r^2 - n^2k^2}$  which prevents to compute a constant  $C_{\alpha,1}$  as required by (2.4.6) in the definition of the symbol class. Hence the concept to realize smoother and smoother approximations can be realized only in the cones

$$(2.4.15) \quad \frac{|\nu^2|}{r^2} < n^2k^2$$

$$(2.4.16) \quad \frac{|\nu^2|}{r^2} > n^2k^2.$$

We refer to the first one as *far-field regime* and to the latter one as *high-frequency regime* in  $\nu$ . Note that the high-frequency regime in  $k$  is identical with the far-field regime. We show that the choice (2.4.14) satisfies condition (2.4.12) with  $j = -1$ . Based on (2.4.10) and (2.4.11) we compute first the symbol of the composite operators

$$\begin{aligned} \sigma(\partial_r \lambda_{\pm}^1) &= i\rho \lambda_{\pm}^1(\mathbf{x}, \boldsymbol{\xi}) + \frac{\partial}{\partial r} \lambda_{\pm}^1(\mathbf{x}, \boldsymbol{\xi}) \\ \sigma(\lambda_{\pm}^1 \partial_r) &= i\rho \lambda_{\pm}^1(\mathbf{x}, \boldsymbol{\xi}). \end{aligned}$$

This shows that the highest order terms of the residual

$$\begin{aligned} \sigma(L - (\partial_r - \lambda_+^1)(\partial_r - \lambda_-^1)) &= \sigma(L - \partial_r^2 + \partial_r \lambda_-^1 + \lambda_+^1 \partial_r - \lambda_+^1 \lambda_-^1) \\ &= \frac{1}{r}i\rho + \frac{\partial}{\partial r} \lambda_-^1(\mathbf{x}, \boldsymbol{\xi}) \end{aligned}$$

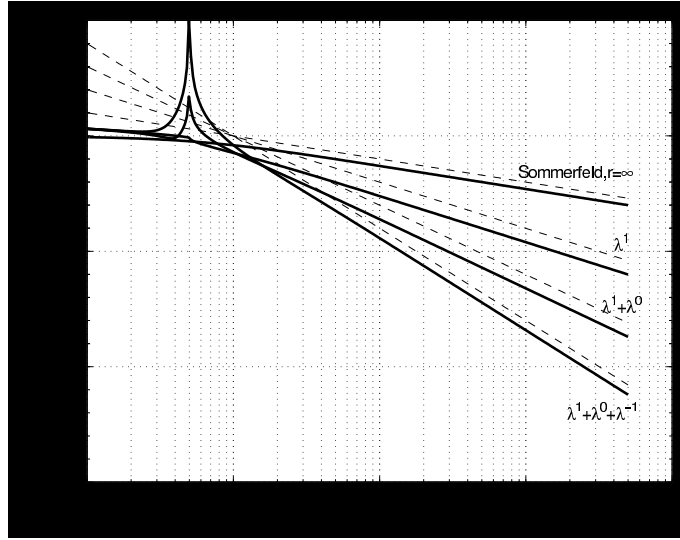


FIGURE 2.4.1. Error of the asymptotic boundary conditions ( $\nu = 0$ ). The dashed lines represent the functions  $r^{-1}$ ,  $r^{-2}$ ,  $r^{-3}$ ,  $r^{-4}$ .

are in fact first order terms, hence the residual is in  $S^1$ , if we modify the definition of  $S^1$ , see (2.4.6), such that it contains elements defined on the cones (2.4.15) or (2.4.16). Engquist and Majda extend the procedure up to pseudodifferential operators  $\lambda_{\pm}^{-1}$ . That is, we approximate  $\lambda_{\pm} \simeq \lambda_{\pm}^1 + \lambda_{\pm}^0 + \lambda_{\pm}^{-1}$  and set the appropriate terms in the residual equation to zero. This yields

$$\begin{aligned} 0 &= \frac{1}{r}i\rho + i\rho(\lambda_+^0 + \lambda_-^0) + \partial_r \lambda_{\pm}^1 - (\lambda_+^1 \lambda_-^0 + \lambda_+^0 \lambda_-^1) && \text{for } \mathcal{O}(|\xi|) \\ 0 &= i\rho(\lambda_+^{-1} + \lambda_-^{-1}) + \partial_r \lambda_{\pm}^0 - (\lambda_+^1 \lambda_-^{-1} + \lambda_+^0 \lambda_-^0 + \lambda_+^{-1} \lambda_-^1) && \text{for } \mathcal{O}(1). \end{aligned}$$

We have the freedom to prescribe a relation between  $\lambda_+^j$  and  $\lambda_-^j$  as long as the order condition is not violated, and following [29] we use  $\lambda_+^{-j} = -\lambda_-^{-j}$  for  $j \geq 0$  to obtain two equations for the symbols  $\lambda_{\pm}^0(\mathbf{x}, \xi)$  and  $\lambda_{\pm}^{-1}(\mathbf{x}, \xi)$

$$\begin{aligned} 0 &= \frac{1}{r}i\rho + \partial_r \lambda_{\pm}^1 - \lambda_{\pm}^0 (\lambda_+^1 - \lambda_-^1) \\ 0 &= \quad \quad \quad + \partial_r \lambda_{\pm}^0 - \lambda_{\pm}^{-1} (\lambda_+^1 - \lambda_-^1) + \lambda_{\pm}^0 \lambda_{\pm}^0. \end{aligned}$$

The first of these equations supplies  $\lambda_{\pm}^0$  using  $\lambda_{\pm}^1$  and the second supplies  $\lambda_{\pm}^{-1}$  using  $\lambda_{\pm}^1$  and  $\lambda_{\pm}^0$ .

**Example.** We study Bessel's equations with  $nk = 1$ . First we investigate the far-field asymptotic for the two fixed frequencies  $\nu = 0$  and  $\nu = 20$ . Fig. 2.4.1 shows the relative error of the normal derivative as function of the radius  $r$ . Clearly, the result becomes better if the order of the approximation is increased. Additionally we observe that this behavior is *not* uniform. Close to the resonance case  $nk = 1 = |\nu/r|$ , which separates the cones (2.4.15) and (2.4.16), the approximation fails.

The same observation can be made in Fig. 2.4.2 for  $\nu = 20$ . Again we find that the approximation fails close to  $nk = 1 = |\nu/r|$ . To justify this, Fig. 2.4.3 displays a zoom into Fig. 2.4.2 around the position  $\nu/r = 20$ .

Next, we investigate the the high-frequency asymptotic. Fig. 2.4.4 shows the error for a fixed radius  $r = 0.1$  as function of the frequency  $\nu$ . Again, we find that



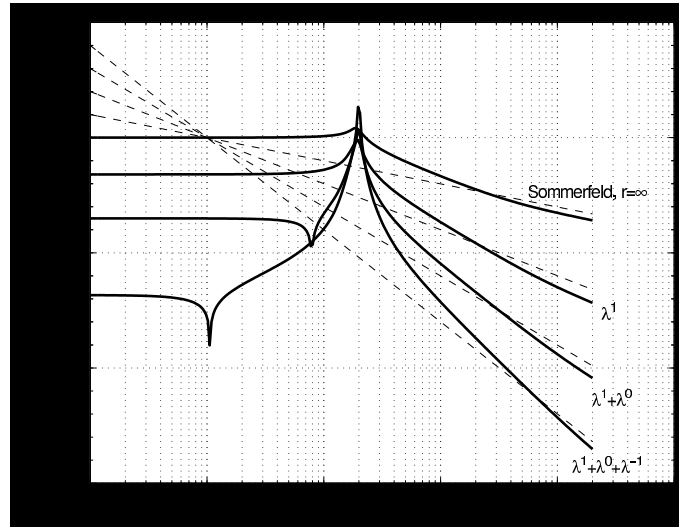


FIGURE 2.4.2. Error of the asymptotic boundary conditions ( $\nu = 20$ ). The dashed lines represent the functions  $r^{-1}$ ,  $r^{-2}$ ,  $r^{-3}$ ,  $r^{-4}$ .

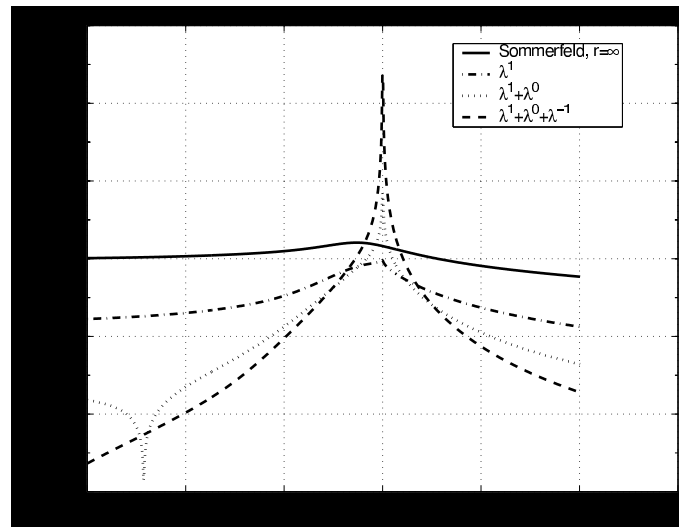


FIGURE 2.4.3. Zoom of the error plot Fig. 2.4.2 into the interval  $5 \leq r \leq 30$

the approximation is the better the higher the order of the approximation for  $\nu$  large. As in the far-field case above, the approximation fails for  $nk = 1 = |\nu/r|$ , as it becomes apparent from Fig. 2.4.5, which shows the error for  $r = 10$ , fixed. In conclusion, the experiments show that the factorization technique works well both in the far-field and in the high-frequency regime. In Section 5.3, p. 138, we will apply this technique to derive a high frequency approximation for the boundary condition in a general case.

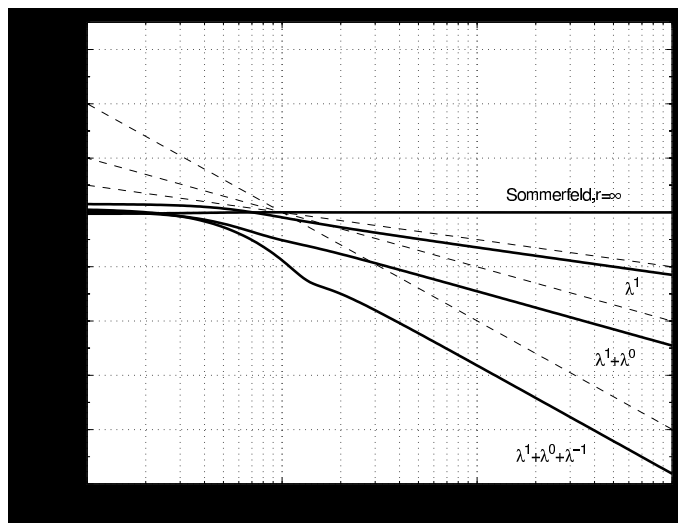


FIGURE 2.4.4. Error of the asymptotic boundary condition ( $r = 0.1$ ). The dashed lines indicate the functions  $\nu^0, \nu^{-1}, \nu^{-2}, \nu^{-4}$ .

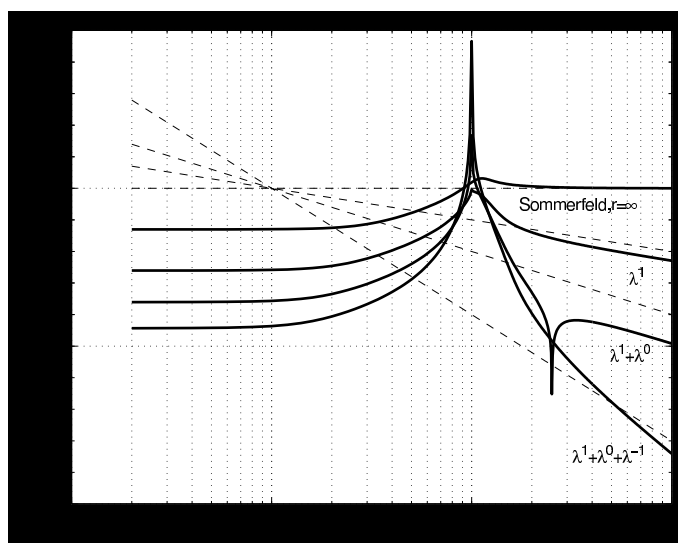


FIGURE 2.4.5. Error of the asymptotic boundary condition ( $r = 10$ ). The dashed lines indicate the functions  $\nu^0, \nu^{-1}, \nu^{-2}, \nu^{-4}$ .

## 2.5. Further Methods

There are many other important methods tailored to solve scattering problems in exterior domains. Starting from an integral representation based on explicitly known Green's functions, both *multipole methods* and *panel clustering* seek to minimize the numerical effort to obtain a solution with prescribed accuracy.

**Multipole Method.** The multipole methods have been developed to have fast summation algorithms for situations, where the potential resulting from many charges has to be computed simultaneously at a number of positions. The first

multipole method has been proposed by Rokhlin [76], where he applied the idea to discretized integral equations with a non-oscillatory kernel. Since the contributions of Greengard [41], [40] and [15] the method has been established as standard tool. In 1993 Rokhlin extended the concept to the 3D Helmholtz equation [77].

**Panel Clustering.** The panel clustering technique goes back to Hackbusch and Nowak [47] and has been improved by Hackbusch, Sauter and Lage in [46, 79]. As opposed to the multipole method, the panel clustering aims to make the boundary integral method efficient. The thesis of Giebermann [34] gives both an overview and many details of these fast summation methods.

**Fictitious domain method.** A further method for  $k \geq 1$  is the so-called fictitious domain method, cf. Ernst [30] and Heikkola et. al. [56]. The idea is to surround the computational domain by sufficient large artificial domain, for which a good asymptotic approximation of the exact transparent boundary condition is known. The technique becomes fast, if the artificial layer can be efficiently discretized by tensorial meshes.

The difference of all these methods to our approach is that they are either based on an explicitly known Green's function or radiation condition, whereas our approaches do not need such a-priori knowledge. The last class of methods which we want to mention here works without using Green's functions.

**Absorbing Layers.** The currently perhaps most widely used and most successful approach is the perfectly matched layer technique (abbreviated to PML afterwards) proposed by Bérenger [11] in 1994 to solve the time-dependent scattering problem of Maxwell's equations in 2D, and extended in 1996 by him to 3D problems [12]. After Bérenger's initial proposal a large number of modifications and extensions to other problem classes than time-dependent Maxwell's equations has been made. Among them are the generalization to time-harmonic and curvilinear coordinates (instead of the usual Cartesian grids) by Collino and Monk [17], the application to paraxial equations like the Fresnel equation by Collino [18], the analysis of different possible auxiliary forms (so-called split-field, biaxial and uniaxial forms) [91] by Yang and Petropoulos and Petropoulos in [75], and the convergence proof of Lassas and Somersalo [67], to mention only a very few of the many interesting articles.

Before Bérenger's publication appeared, many other forms of absorbing layers had been developed, mainly motivated by the modeling of a physical absorption. However, caused by the success of the PML-technique, most of them are no longer in use.

The main idea of absorbing layers is to surround the artificial boundary by an artificial layer of finite thickness. The goal in the design of these layers is to avoid spurious reflections from the artificial boundary and the interior of the absorbing layer to the greatest possible extend. Bérenger showed how to modify Maxwell's equations to provide a perfectly matched absorber. That is, the waves are absorbed independently from their frequency and incidence direction. Moreover, the waves decay exponentially with the distance into the layer. Thus the PML layer can be truncated after a relatively short distance. The PML layer must be discretized like the interior problem and the discretized interior problem and the modified problem on the PML domain have to be solved simultaneously.

Surprisingly, it turned out that the PML method has much in common with the Laplace domain methods which forms the main topic of this work. In fact, from a certain point of view we can see the PML technique as one (among other possible)

realization of the pole condition, compare Section 4.4.6 for a short discussion of this issue.

## Characterization of Incoming and Outgoing Waves: Pole Condition

This chapter introduces our central condition, the pole condition. We start with the analysis of non-trivial examples in 1D and show that their far-field is characterized by the singularities of the Laplace transformed solutions. We discuss the physical meaning of the solution representation in terms of the Laplace transform, which shows that the Laplace transform is a convenient tool to represent functions on semi-infinite intervals. The essence of the examples leads to the definition of the pole condition. We extend the pole condition to both higher space dimensions and time-dependent problems. The crucial definitions are the following:

- (1) Pole condition in 1D for continuous functions, Definition 3.1.1. This refers to potentials which become constant at infinity.
- (2) Pole condition in arbitrary space dimensions for continuous functions, Definition 3.4.5. This is a generalization of the first point.
- (3) Pole condition for time-dependent problems. A direct generalization to Schrödinger's equation.

In this section we present and analyze three characteristic examples, which can be treated almost analytically. These are

- The Helmholtz equation in 1D with constant coefficients.
- The Helmholtz equation in 1D with periodic coefficients.
- The radially symmetric Helmholtz equation in 2D with constant coefficients.

The Helmholtz equation in 1D is the most simple equation and therefore well suited to introduce the concept. The Bessel equation is a 1D equation, too, nevertheless it characterizes the situation in higher space dimension and is therefore of fundamental interest. We present an analytic solution, up to a convolution integral, of Bessel's equation in Laplace domain. This will motivate our treatment of the generalized Bessel (Hankel) functions, which provides the basis of the general theory for problems with variable coefficients. Finally, we consider the case of periodic coefficients. This is of interest for a number of practical problems and shows, how theory and algorithms can be extended to cover this problem class.

### 3.1. Helmholtz Equation in 1D with Constant Potential

Let us consider the Helmholtz equation in the semi-infinite interval  $[a, \infty)$  and let us introduce the shifted coordinate-system with the independent variable  $x$  via  $r = a + x, x \geq 0$ ,

$$\frac{d^2}{dx^2}u(x) + u(x) = 0, \quad x > 0.$$

Here we set w. r. o. g.  $k^2 = 1$ . Laplace transform

$$(3.1.1) \quad \hat{u}(s) = (Lu(x))(s) = \int_0^\infty dx e^{-sx}u(x), \quad \operatorname{Re} s > 0$$

yields

$$\hat{u}(s) = \frac{su(0) + \partial_x u|_{x=0}}{s^2 + 1}.$$

In the following we abbreviate always:  $u'(0) := \partial_x u|_{x=0}$ . The complex valued function  $\hat{u}(s)$  is holomorphic in  $\mathbb{C} \setminus \{\pm i\}$  and possesses two isolated singularities, namely the points  $\{\pm i\}$ . Hence  $\hat{u}(s)$ , can be extended from  $\text{Re } s > 0$ , where it has been defined to the whole complex plane except the singular points. A partial fraction decomposition supplies

$$\hat{u}(s) = \frac{1}{2} \frac{u(0) + iu'(0)}{s+i} + \frac{1}{2} \frac{u(0) - iu'(0)}{s-i}.$$

Obviously, the pole at  $s = -i$  vanishes, if the 1D Sommerfeld condition holds:  $u(0) + iu'(0) = 0$ . We will call this condition pole condition. Thus, the Laplace-technique supplies two essential informations:

- (1) It gives a natural splitting of a given function into asymptotically incoming and outgoing waves, because we can identify a complex function  $\hat{u}_-(s) := \frac{1}{s+i}$  with an incoming wave  $u_-(x) = \exp(-ix)$  and a function  $\hat{u}_+(s) := \frac{1}{s-i}$  with an outgoing wave  $u_+(x) = \exp(ix)$ .
- (2) It supplies automatically the relation between Neumann data  $u'(0)$  and the Dirichlet data  $u(0)$  needed to drop one of them.

**DEFINITION 3.1.1.** *Pole condition in 1D.* A function  $u : \mathbb{R}_+ \rightarrow \mathbb{C}$  satisfies the pole condition if the function  $\hat{u}$  defined by the Laplace transform (3.1.1) has a holomorphic extension to the lower half of the complex plane  $\{s \in \mathbb{C} : \text{Im } s < 0\}$ .

**REMARK 3.1.2.** The pole condition is a condition concerning the behavior of  $u$  at infinity, i.e. it is independent of the choice of the boundary of the computational domain. If we replace the boundary at  $x = a$  by a boundary at  $a + \delta$ , this property follows from

$$\begin{aligned} \int_0^\infty dx e^{-sx} u(x) &= \int_0^\delta dx e^{-sx} u(x) + \int_{-\delta}^\infty dx e^{-sx} u(x) \\ &= \int_0^\delta dx e^{-sx} u(x) + e^{-s\delta} \int_0^\infty dx e^{-sx} u(x + \delta) \end{aligned}$$

and the fact that both  $s \mapsto \int_0^\delta dx e^{-sx} u(x)$  and  $s \mapsto e^{-s\delta}$  are entire functions.

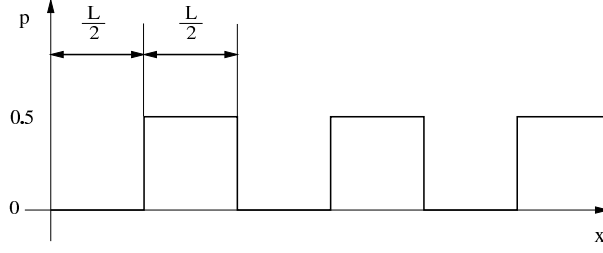
### 3.2. Helmholtz Equation in 1D with Periodic Potentials

We consider real, positive and piecewise continuous potentials with period  $L > 0$ , that is  $p(x + L) = p(x)$ ,  $x \geq 0$  along with the 1D Helmholtz equation

$$(3.2.1) \quad u''(x) + (k^2 + p(x)) u(x) = 0, \quad x > 0$$

with initial data  $u(0)$  and  $u'(0)$ . If  $u(x)$  is a solution to (3.2.1), it satisfies the *conservation property*  $\text{Im}(\bar{u}(x)u'(x)) = \text{const}$  for all  $x \geq 0$ . The quantity  $P := \text{Im}(\bar{u}(x)u'(x))$  is called power flux. The conservation property is derived by a multiplication of (3.2.1) by  $\bar{u}$ , followed by an integration by parts and a computation of the imaginary part of the resulting expression.

It is convenient to rewrite the Helmholtz equation (3.2.1) into a first order system

FIGURE 3.2.1. Example of a periodic potential  $p(x)$ 

$$(3.2.2) \quad \begin{pmatrix} u \\ v \end{pmatrix}' = \begin{pmatrix} 0 & 1 \\ -k^2 - p & 0 \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix}, \quad x > 0,$$

with initial data  $u(0)$  and  $v(0) = u'(0)$ . The points  $\{x_j = (j-1)L : j = 1, 2, \dots\}$  define a partition of  $\mathbb{R}_+$  into intervals  $[x_j, x_{j+1}]$ . We consider a solution of (3.2.2) on the interval  $[x_j, x_{j+1}]$  with respect to the initial data  $u_j := u(x_j)$  and  $v_j := v(x_j)$ . Let us denote the map from the initial data at  $x_j$  to the data at  $x_{j+1}$  as

$$(3.2.3) \quad \begin{pmatrix} u \\ v \end{pmatrix}_{j+1} = W \begin{pmatrix} u \\ v \end{pmatrix}_j, \quad W \in \mathbb{R}^{2 \times 2}, j = 1, 2, \dots$$

$$\begin{pmatrix} u \\ v \end{pmatrix}_1 = \begin{pmatrix} u_1 \\ v_1 \end{pmatrix} = \begin{pmatrix} u(0) \\ v(0) \end{pmatrix}.$$

Here, we understood that  $(u, v)^T$  has a continuous continuation to the boundary of the interval. The matrix  $W$  is regular since a backward integration exists on the same interval. Further,  $W$  does not depend on  $j$ , due to the periodicity of  $p$ . The repeated application of (3.2.3) yields the sequences

$$U = \{u_j : j = 1, 2, \dots\} \quad \text{and} \quad V = \{v_j : j = 1, 2, \dots\}.$$

Let us introduce the z-transforms of  $U$  and  $V$  by

$$\widehat{U}(z) = \sum_{j=1}^{\infty} \frac{u_j}{z^j}, \quad \text{and} \quad \widehat{V}(z) = \sum_{j=1}^{\infty} \frac{v_j}{z^j} \quad |z| \geq a > 0$$

with a constant  $a$  such that the sums converge. Furthermore, we define a two element vector of sequences  $(U, V)^T$  and a corresponding matrix-vector multiplication by

$$W \begin{pmatrix} U \\ V \end{pmatrix} = \left\{ W \begin{pmatrix} u \\ v \end{pmatrix}_j : j = 1, 2, \dots \right\}.$$

The product  $W \cdot (\widehat{U}, \widehat{V})^T$  is the conventional matrix-vector product. Obviously it holds

$$z\widehat{U}(z) = u(0) + \sum_{j=1}^{\infty} \frac{u_{j+1}}{z^j}$$

$$z\widehat{V}(z) = v(0) + \sum_{j=1}^{\infty} \frac{v_{j+1}}{z^j},$$

which allows us to write the z-transform of the left-hand side of (3.2.3) as

$$\begin{pmatrix} z\widehat{U}(z) - u(0) \\ z\widehat{V}(z) - v(0) \end{pmatrix}.$$

Consequently, we obtain the  $z$ -transform of the complete equation

$$(3.2.4) \quad (z - W) \begin{pmatrix} \widehat{U}(z) \\ \widehat{V}(z) \end{pmatrix} = \begin{pmatrix} u(0) \\ v(0) \end{pmatrix},$$

where here and in the following  $z - W$  denotes  $zI - W$ . Apparently, this is the direct analogue to the Laplace transform of the Helmholtz equation with constant coefficients. Obviously, the solution  $(\widehat{U}(z), \widehat{V}(z))^T$  has in general poles if  $z$  approaches the eigenvalues of  $W$ . We assign to each pole a special solution of the homogeneous equation, where one corresponds to an incoming and the other to an outgoing function. Our goal will be to find the conditions with respect to  $(u(0), v(0))^T$  such that  $(\widehat{U}(z), \widehat{V}(z))^T$  has only one pole and the corresponding functions are outgoing functions. By means of the  $z$ -transform we will show that the sequence  $(u, v)_{j=1, \dots, \infty}^T$  has a representation

$$(3.2.5) \quad \begin{pmatrix} u \\ v \end{pmatrix}_j = C_1 \mathbf{q}^{(1)} t^j + C_2 \mathbf{q}^{(2)} \bar{t}^j \quad \text{for all } j \geq 1$$

with  $C_1, C_2, t \in \mathbb{C}$  and  $\mathbf{q}^{(1)}, \mathbf{q}^{(2)} \in \mathbb{C}^2$ . In case of a vanishing periodic potential, that is  $p = 0$ , we have a representation:

$$(3.2.6) \quad \begin{pmatrix} u \\ v \end{pmatrix}_j = C_1 \begin{pmatrix} 1 \\ ik \end{pmatrix} e^{ikx_j} + C_2 \begin{pmatrix} 1 \\ -ik \end{pmatrix} e^{-ikx_j} \quad \text{for all } j \geq 1.$$

Obviously  $t^j$  plays in the periodic case a role similar to  $\exp(ikx_j)$  in the standard case. In fact, Lemma 3.2.3 at the end of this section, p. 45, shows that  $|t| = 1$ . Moreover, we will show that one term of (3.2.5) has positive power flux whereas the other must have a negative power flux. This corresponds exactly to the fact that the power flux  $P_+$  which we can assign to the first term of (3.2.6) is  $P_+ = k |C_1|^2 > 0$ , whereas the power flux of the second term is  $P_- = -k |C_2|^2 < 0$ . The representation (3.2.5) can be seen as a variant of Floquet's theorem. Additionally to the existence of at least one quasi-periodic solution  $C \mathbf{q} t^j$  stated by Floquet's theorem, cf. e.g. [16], Chapters 31 and 32, pp. 121, we obtained a complete representation formula due to exploitation of the special structure of the system (39).

The first idea to apply the pole condition technique to sequences is to require the holomorphic extension of  $\widehat{U}(z)$  to the lower half of the complex plane. Of course, this is the right idea in the limiting case  $p$  small. For  $L$  sufficiently large, however,  $u(x)$  may have a large phase increment over one period  $L$ . Therefore the simple correspondence: negative phase increment of  $u(x) \longleftrightarrow$  singularity of  $\widehat{U}(z)$  in the lower half of the complex plane does not hold true in general. As an example let us consider

$$u_j = e^{ij(2\pi-\alpha)} + e^{-ij(2\pi-\alpha)} = e^{-ij\alpha} + e^{ij\alpha}, \quad j \geq 1$$

with an angle  $\alpha$ ,  $0 < \alpha < \pi$ . The sequence composed from  $u_j$  has a  $z$ -transform

$$\frac{e^{i(2\pi-\alpha)}}{z - e^{i(2\pi-\alpha)}} + \frac{e^{-i(2\pi-\alpha)}}{z - e^{-i(2\pi-\alpha)}}.$$



Hence, the first term which corresponds to a *positive* phase increment has its pole in the *lower* half plane, whereas the second term which corresponds to the *negative* phase increment has its pole in the *upper* half of the complex plane. Thus the situation known from the continuous 1D Helmholtz equation is reversed. Nevertheless, it turned out that there is another tool which allows us to decide which of the two terms of (3.2.5) corresponds to an *outgoing* wave. The important fact which is proved in Lemma 3.2.2 at the end of this section, p. 44, is that

$$(3.2.7) \quad \operatorname{Im} \left( \frac{\mathbf{q}^{(1)}(2)}{\mathbf{q}^{(1)}(1)} \right) \cdot \operatorname{Im} \left( \frac{\mathbf{q}^{(2)}(2)}{\mathbf{q}^{(2)}(1)} \right) < 0$$

holds true under general conditions. Consider, e. g. the case  $\operatorname{Im} (\mathbf{q}^{(1)}(2)/\mathbf{q}^{(1)}(1)) > 0$ , consequently  $\operatorname{Im} (\mathbf{q}^{(2)}(2)/\mathbf{q}^{(2)}(1)) < 0$ . If we set  $C_2 = 0$  in (3.2.5), we have the simple series

$$(3.2.8) \quad \begin{pmatrix} u \\ v \end{pmatrix}_j = C_1 t^j \mathbf{q}^{(1)} \quad \text{for all } j \geq 1.$$

It follows

$$\operatorname{Im} \left( \frac{v_j}{u_j} \right) = \operatorname{Im} \left( \frac{u'_j}{u_j} \right) = \operatorname{Im} \left( \frac{\mathbf{q}^{(1)}(2)}{\mathbf{q}^{(1)}(1)} \right) > 0 \quad \text{for all } j \geq 1.$$

Hence the *local* phase increment  $\operatorname{Im} (u'_j/u_j)$  is positive at *all* sample points  $x_j$  and the sequence  $(u, v)_{j=1, \dots, \infty}^T$  is the sequence of sampled points of an *outgoing* function. Reversely, we characterize the sequence

$$\begin{pmatrix} u \\ v \end{pmatrix}_j = C_2 \bar{t}^j \mathbf{q}^{(2)} \quad \text{for all } j \geq 1$$

with  $\operatorname{Im} (\mathbf{q}^{(2)}(2)/\mathbf{q}^{(2)}(1)) < 0$  as the sequence of sampled points of an *incoming* function. Note that together with a positive phase increment the power flux becomes positive. This follows trivially from

$$\operatorname{Im} \left( \frac{u'_j \bar{u}_j}{u_j \bar{u}_j} \right) = \frac{P}{|u_j|^2} > 0.$$

Assuming the property (3.2.7), we can give explicitly the Dirichlet-to-Neumann number  $u'(0)/u(0)$  which guarantees that all local phase increments at the sample points are positive:

$$(3.2.9) \quad u'(0)/u(0) = \begin{cases} \frac{\mathbf{q}^{(1)}(2)}{\mathbf{q}^{(1)}(1)} & \text{if } \operatorname{Im} \left( \frac{\mathbf{q}^{(1)}(2)}{\mathbf{q}^{(1)}(1)} \right) > 0 \\ \frac{\mathbf{q}^{(2)}(2)}{\mathbf{q}^{(2)}(1)} & \text{if } \operatorname{Im} \left( \frac{\mathbf{q}^{(1)}(2)}{\mathbf{q}^{(1)}(1)} \right) < 0. \end{cases}$$

This is the main result with respect to our 1D example with periodic coefficients.

Having outlined the characterization of sequences, we want to show how to compute the quantities  $\mathbf{q}^{(1)}$ ,  $\mathbf{q}^{(2)}$  and  $t$  from our basic equation (3.2.4). The matrix  $W$  has a Schur decomposition with a unitary matrix  $Q^{(1)} \in \mathbb{C}^{2 \times 2}$  such that

$$Q^{(1)H} W Q^{(1)} = T^{(1)} = D^{(1)} + N^{(1)}$$

where  $D^{(1)} = \operatorname{diag}(\lambda_1, \lambda_2)$  and  $N^{(1)} \in \mathbb{C}^{2 \times 2}$  is strictly upper triangular, cf. e. g. [39, Chap. 7, Theorem 7.1.3, p. 313]. To keep the representation simple, we consider only the case that the matrix  $W$  has complex eigenvalues with non-vanishing imaginary part. Since  $W$  is real, these eigenvalues are complex conjugates. Further,

because  $W$  is a  $2 \times 2$  matrix, there exists a second unitary matrix  $Q^{(2)} \in \mathbb{C}^{2 \times 2}$  such that

$$Q^{(2)H} W Q^{(2)} = T^{(2)}$$

and

$$T^{(1)} = \begin{pmatrix} t_{11} & t_{12} \\ & t_{22} \end{pmatrix}, \quad T^{(2)} = \begin{pmatrix} t_{22} & \tilde{t}_{12} \\ & t_{11} \end{pmatrix},$$

that is, the main diagonal elements of  $T^{(1)}$  and  $T^{(2)}$  are reordered. This reordering can be obtained by the application of a complex Givens matrix  $U$ , cf. e.g. [78], which realizes the transforms

$$(3.2.10) \quad \begin{aligned} Q^{(2)} &= Q^{(1)} U \\ T^{(2)} &= U^H T^{(1)} U. \end{aligned}$$

Let us denote in the following the two column vectors of  $Q^{(1,2)}$  by  $\mathbf{q}_1^{(1,2)}$  and  $\mathbf{q}_2^{(1,2)}$ . A unitary transform of (3.2.4) yields

$$(3.2.11) \quad (zI - T) Q^{(1)H} \begin{pmatrix} \widehat{U}(z) \\ \widehat{V}(z) \end{pmatrix} = Q^{(1)H} \begin{pmatrix} u(0) \\ v(0) \end{pmatrix}.$$

This equation has two singularities, namely for  $z \in \{t_{11}, t_{22}\}$ . Obviously, the solution of (3.2.11) has no pole in  $z = t_{22}$  if

$$\begin{pmatrix} u(0) \\ v(0) \end{pmatrix} \in \text{span} \{ \mathbf{q}_1^{(1)} \}.$$

In this case, we obtain a first special solution to (3.2.11):

$$\begin{aligned} Q^{(1)H} \begin{pmatrix} \widehat{U}(z) \\ \widehat{V}(z) \end{pmatrix}^{(1)} &= (zI - T^{(1)})^{-1} \begin{pmatrix} C_1 \\ 0 \end{pmatrix} \\ &= \frac{1}{z} \left( I + \frac{1}{z} T^{(1)} + \frac{1}{z^2} (T^{(1)})^2 + \dots \right) \begin{pmatrix} C_1 \\ 0 \end{pmatrix} \\ &= \frac{1}{z} \left( 1 + \frac{1}{z} t_{11} + \frac{1}{z^2} t_{11}^2 + \dots \right) \begin{pmatrix} C_1 \\ 0 \end{pmatrix} \end{aligned}$$

which yields

$$\begin{pmatrix} \widehat{U}(z) \\ \widehat{V}(z) \end{pmatrix}^{(1)} = \frac{C_1}{z} \left( 1 + \frac{1}{z} t_{11} + \frac{1}{z^2} t_{11}^2 + \dots \right).$$

Inverse  $z$ -transform supplies in fact the first term of (3.2.5). We may repeat the whole procedure replacing  $Q^{(1)}$  by  $Q^{(2)}$  and  $T^{(1)}$  by  $T^{(2)}$  to obtain

$$\begin{pmatrix} \widehat{U}(z) \\ \widehat{V}(z) \end{pmatrix}^{(2)} = \frac{C_2}{z} \left( 1 + \frac{1}{z} t_{22} + \frac{1}{z^2} t_{22}^2 + \dots \right).$$

Thus we have proved the representation formula (3.2.5).

**REMARK 3.2.1.** Following Ruhe [78] we can compute the unitary Givens matrix  $U$  using the ansatz

$$U = \begin{pmatrix} \alpha & \bar{\beta} \\ -\beta & \bar{\alpha} \end{pmatrix}.$$

The coefficients are obtained as follows:

$$\begin{aligned}
 (i) \quad |t_{11} - t_{22}| \leq |t_{12}| : \quad & \alpha = \frac{1}{\sqrt{1 + |\gamma|^2}} \\
 & \gamma = \frac{t_{11} - t_{22}}{t_{12}} \\
 & \beta = \gamma\alpha \\
 (ii) \quad |t_{11} - t_{22}| > |t_{12}| : \quad & \beta = \frac{1}{\sqrt{1 + |\gamma|^2}} \\
 & \gamma = \frac{t_{12}}{t_{11} - t_{22}} \\
 & \alpha = \gamma\beta.
 \end{aligned}$$

To illustrate the properties arising from periodic media, we want to study, in brief, the reflection and transmission of the periodic potential given in Fig. 3.2.1 with  $L = \pi/2$  and

$$\text{First example, low contrast :} \quad p = \begin{cases} 0 & \text{for } 0 \leq x < \frac{L}{2} \\ 0.05 & \text{for } \frac{L}{2} \leq x < L \end{cases}, \quad p(x+L) = x$$

as well as

$$\text{Second example, high contrast :} \quad p = \begin{cases} 0 & \text{for } 0 \leq x < \frac{L}{2} \\ 0.5 & \text{for } \frac{L}{2} \leq x < L \end{cases}, \quad p(x+L) = x$$

Corresponding to the layers, we set  $k_1 = 1$ ,  $k_2 = 1.05$  in the first and  $k_1 = 1$ ,  $k_2 = 1.5$  in the second example. To obtain the transmission matrix  $W$ , compare (3.2.3), we write the general solution of the Helmholtz equation in each of the two layers as

$$(3.2.12) \quad u_{1,2}(x) = c_1 e^{ik_{1,2}x} + c_2 e^{-ik_{1,2}x}.$$

Accordingly, we obtain the matrices

$$W_{1,2} = \begin{pmatrix} \cosh\left(ik_{1,2}\frac{L}{2}\right) & \frac{1}{ik_{1,2}} \sinh\left(ik_{1,2}\frac{L}{2}\right) \\ ik_{1,2} \sinh\left(ik_{1,2}\frac{L}{2}\right) & \cosh\left(ik_{1,2}\frac{L}{2}\right) \end{pmatrix}$$

supplying the composite matrix  $W = W_1 W_2$ . By (3.2.9) we compute the DtN-number  $\text{DtN} = u'(0)/u(0)$ . Further, we define an amplitude reflection factor  $r$  with respect to the *first* layer

$$r := \frac{c_2}{c_1} = \frac{ik_1 - \text{DtN}}{ik_1 + \text{DtN}}$$

where  $c_1$  and  $c_2$  are the constants from general representation (3.2.12). The reflection factor is the physical quantity of interest. Fig. 3.2.2 shows the computed reflection coefficients if the wavelength  $\lambda$  varies from  $\pi/8$  to  $\pi/8 + 2\pi$ . The figures show the typical resonance effects of periodic layered media. The intervals where the reflection coefficient takes its maximum value, i. e.  $r = 1$ , have finite widths. In physics these intervals are called stopbands and are of particular interest for many applications. Among others, it can be seen from Fig. 3.2.2 that the larger the index contrast between the two layers, the larger the width of the stopband. These properties of 1D periodic layered media are well documented in the literature, cf. the book of März [71], Chapters 2.2.4 and 7.2. The merit of our approach via the pole condition lies in the fact that no limiting process as e. g. number of layers  $\rightarrow \infty$

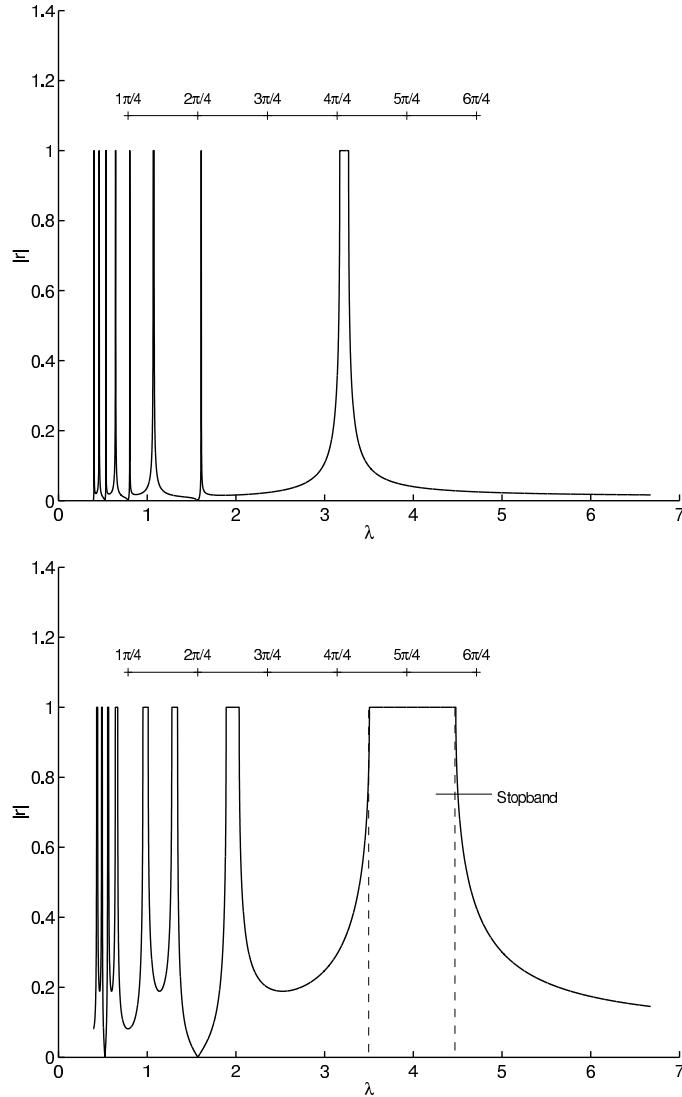


FIGURE 3.2.2. semi-infinite Bragg grating. Amplitude reflection coefficients vs. the wavelength for two different index contrasts. The upper figure refers to a layer stack with low index contrast ( $k_1 = 2\pi/\lambda$ ,  $k_2 = 1.05 \cdot 2\pi/\lambda$ ), the lower figure to a stack with high index contrast ( $k_1 = 2\pi/\lambda$ ,  $k_2 = 1.5 \cdot 2\pi/\lambda$ )

has to be considered. Moreover, there are several ways to extend this to higher dimensions.

LEMMA 3.2.2. *Let the matrix  $W \in \mathbb{R}^2$  defined in (3.2.3) have a Schur transform  $Q^{(1)H}WQ^{(1)} = T^{(1)}$  with a unitary matrix  $Q^{(1)} := [\mathbf{q}_1^{(1)}, \mathbf{q}_2^{(1)}]$  and an upper triangular matrix  $T^{(1)}$  such that  $\text{Im}(t_{11}) \neq 0$ . Then there exists a second Schur transform  $Q^{(2)H}WQ^{(2)} = T^{(2)}$  with  $Q^{(1)} \neq Q^{(2)} := [\mathbf{q}_1^{(2)}, \mathbf{q}_2^{(2)}]$  and the following holds true:*

$$\text{Im} \left( \frac{\mathbf{q}_1^{(1)}(2)}{\mathbf{q}_1^{(1)}(1)} \right) \cdot \text{Im} \left( \frac{\mathbf{q}_1^{(2)}(2)}{\mathbf{q}_1^{(2)}(1)} \right) \leq 0.$$

In particular, if  $w_{21} \neq 0$ , then the expression becomes strictly negative.

PROOF. Since  $\text{Im}(t_{11}) \neq 0$  and  $W$  is real we have  $t_{11} = \bar{t}_{22}$ , hence  $t_{11} \neq t_{22}$  and  $t_{22} \neq w_{22}$ . Thus the unitary matrix  $U$  defined in Remark 3.2.1 is different from the identity matrix. Hence it follows from the transform (3.2.10) that  $Q^{(1)} \neq Q^{(2)}$ , in particular  $\mathbf{q}_1^{(1)} \neq \mathbf{q}_1^{(2)}$ .

Taking the first column of each side of  $Q^{(1)H}WQ^{(1)} = T^{(1)}$  we obtain

$$Wq_1^{(1)} = t_{11}q_1^{(1)}.$$

An evaluation in component-wise form yields

$$\frac{q_1^{(1)}(2)}{q_1^{(1)}(1)} = \frac{w_{21}}{t_{11} - w_{22}}.$$

Consequently we have

$$(3.2.13) \quad \text{Im} \left( \frac{q_1^{(1)}(2)}{q_1^{(1)}(1)} \right) = \text{Im} \left( \frac{w_{21}}{t_{11} - w_{22}} \right)$$

which is the first factor of the expression under investigation. We repeat the procedure with the second form of the Schur transform and obtain

$$(3.2.14) \quad \begin{aligned} \text{Im} \left( \frac{q_1^{(2)}(2)}{q_1^{(2)}(1)} \right) &= \text{Im} \left( \frac{w_{21}}{t_{22} - w_{22}} \right) \\ &= \text{Im} \left( \frac{w_{21}}{\bar{t}_{11} - w_{22}} \right) \\ &= \text{Im} \left( \frac{w_{21}}{t_{11} - w_{22}} \right) \\ &= -\text{Im} \left( \frac{w_{21}}{t_{22} - w_{22}} \right). \end{aligned}$$

Computing the product of (3.2.13) and (3.2.14) we obtain finally the desired result

$$\text{Im} \left( \frac{q_1^{(1)}(2)}{q_1^{(1)}(1)} \right) \cdot \text{Im} \left( \frac{q_1^{(2)}(2)}{q_1^{(2)}(1)} \right) = -w_{21}^2 \left( \text{Im} \left( \frac{1}{t_{22} - w_{22}} \right) \right)^2 \leq 0.$$

□

LEMMA 3.2.3. *Let the matrix  $W$  satisfy the same assumptions as in Lemma 3.2.2. Then the complex number  $t$  from the representation formula (3.2.5) satisfies  $|t| = 1$ .*

PROOF. The representation formula holds true for any initial data  $(u(0), v(0))$ . Let the initial data be chosen such that  $C_2 = 0$  which means

$$\begin{pmatrix} u \\ v \end{pmatrix}_j = \begin{pmatrix} u \\ u' \end{pmatrix}_j = C_1 t^j \mathbf{q}^{(1)} \quad j \geq 1.$$

From this we compute the power flux

$$\text{Im}(\bar{u}u')_j = |C_1|^2 |t|^{2j} \overline{\mathbf{q}_1^{(1)}(1)} q_1^{(1)}(2).$$

It follows from the conservation property  $\text{Im}(\bar{u}u')_j = \text{const}$  for all  $j \geq 1$  that

$$\frac{\text{Im}(\bar{u}u')_{j+1}}{\text{Im}(\bar{u}u')_j} = |t|^2 = 1.$$

□

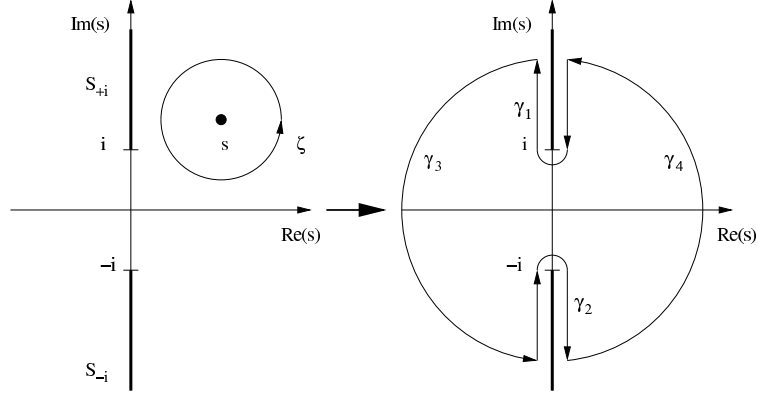


FIGURE 3.3.1. Definition of the cuts  $S_{\pm i}$  and the contours  $\gamma_1, \dots, \gamma_4$ . The contours around the cuts can be closed at infinity.

### 3.3. Radially Symmetric Helmholtz Equation in 2D

The following example is essentially for the understanding of the full 2D problem. Let us consider the complex function  $\hat{u}(s)$

$$(3.3.1) \quad \hat{u}(s) := \frac{g(s)}{\sqrt{s^2 + 1}}.$$

Here, the square-root is understood as the square root on the principal plane, i. e. , for  $a \in \mathbb{C}$  we define  $\sqrt{a} = |a|^{1/2}e^{i\phi/2}$ ,  $\phi = \arg(a)$ ,  $\phi \in (-\pi, \pi]$ . With respect to the definition of  $\hat{u}(s)$  this means, that we define the two branch cuts  $S_{\pm i}$

$$(3.3.2) \quad S_{\pm i} = \{\pm i(1 + \tau) : \tau \in \mathbb{R}_+\}.$$

In the more general cases, which we shall consider later,  $\hat{u}$  is not given explicitly, but by differential or in the most general cases by integral equations. Then it will turn out to be more convenient to define the cuts parallel to the real axis, in order to construct the analytic continuation. For the present case, however, we will use the above given standard cuts.

The function  $g(s)$  in the numerator of (3.3.1) is defined as follows:  $g : \mathbb{C} \setminus (S_{+i} \cup S_{-i}) \rightarrow \mathbb{C}$ ,  $g(s)$  is holomorphic in  $\mathbb{C} \setminus (S_{+i} \cup S_{-i})$  and  $\lim_{|s| \rightarrow \infty} |g(s)| = \text{const.}$  The two singularities  $\pm i$  and the infinite point are branch points of  $\hat{u}(s)$  of order 1. The infinite point is a zero of  $\hat{u}$ , and the derivative  $\hat{u}'(s)$  vanishes for  $|s| \rightarrow \infty$ , too. Hence, the infinite point is a regular point. In the following, we want to analyze the function  $\hat{u}(s)$  in the same way as we analyzed the function  $1/(s^2 + 1)$  occurring in the 1D example. The following lemma supplies the key for this investigation.

LEMMA 3.3.1. *The complex function  $\hat{u}$  defined by (3.3.1) possesses a decomposition*

$$\hat{u}(s) = \frac{1}{2\pi i} \oint_{\gamma_1} \frac{g(\zeta)}{\sqrt{\zeta^2 + 1}(s - \zeta)} d\zeta + \frac{1}{2\pi i} \oint_{\gamma_2} \frac{g(\zeta)}{\sqrt{\zeta^2 + 1}(s - \zeta)} d\zeta,$$

where the paths  $\gamma_1$  and  $\gamma_2$  are contours in  $\mathbb{C} \setminus (S_{+i} \cup S_{-i})$  enclosing the cuts  $S_{\pm i}$  (see Fig. 3.3.1).

PROOF. For a small disc with radius  $\epsilon$  we have, according to Cauchy's integral formula,

$$\hat{u}(s) = \frac{1}{2\pi i} \oint_{|\zeta-s|=\epsilon} \frac{\hat{u}(\zeta)}{\zeta-s} d\zeta.$$

Without to change the result we deform the path, as it is shown in Fig. 3.3.1, and obtain

$$\hat{u}(s) = \frac{1}{2\pi i} \oint_{\gamma} \frac{\hat{u}(\zeta)}{\zeta-s} d\zeta.$$

Here,  $\gamma = \gamma_1 \cup \gamma_2 \cup \gamma_3 \cup \gamma_4$ . Choosing for  $\gamma_3, \gamma_4$  half-circles with radius  $r$ , we have only to show that the integrals along these paths vanish for  $r \rightarrow \infty$ . In polar coordinates  $r, \phi, \frac{\pi}{2} \leq \phi \leq \frac{3\pi}{2}$ , we obtain on  $\gamma_3$

$$\begin{aligned} \left| \int_{\gamma_3} \frac{\hat{u}(\zeta)}{\zeta-s} d\zeta \right| &\leq \int_{\frac{\pi}{2}}^{\frac{3\pi}{2}} \left| \frac{|g(\zeta)|}{\sqrt{r^2 e^{2i\phi} + 1} (re^{i\phi} - s)} ire^{i\phi} \right| d\phi \\ &\leq \int_{\frac{\pi}{2}}^{\frac{3\pi}{2}} \frac{r|g|}{\sqrt{|r^2-1|} |r-|s||} d\phi. \end{aligned}$$

Hence, the integral tends to zero for  $r \rightarrow \infty$ . The same applies to the integral on  $\gamma_4$ .  $\square$

We discuss two examples. First, we consider the more simpler treat case  $g = 1$  on the entire complex plane. It will turn out that this corresponds to the Bessel function  $J_0$ . The second example deals with the Laplace transform of Bessel's equation, for which we will present a nearly analytic solution. We will show that a decomposition of the solution of the Laplace transformed Bessel equations according to Lemma 3.3.1 can be obtained.

We want to analyze the first term of the partial fraction decomposition of Lemma 3.3.1

$$\hat{u}_+(s) := \frac{1}{2\pi i} \oint_{\gamma_1} \frac{d\zeta}{\sqrt{\zeta^2+1}(s-\zeta)}.$$

We split the contour  $\gamma_1$  into two parts  $\gamma_1 = \gamma'_1 + \gamma''_1$ , where  $\gamma'_1$  lies to the left and  $\gamma''_1$  to the right of the imaginary axis and  $\zeta_1 \in \gamma'_1$  and  $\zeta_2 \in \gamma''_1$

$$\hat{u}_+(s) = \frac{1}{2\pi i} \int_i^{i\infty} \frac{d\zeta_1}{\sqrt{\zeta_1^2+1}(s-\zeta_1)} + \frac{1}{2\pi i} \int_{i\infty}^i \frac{d\zeta_2}{\sqrt{\zeta_2^2+1}(s-\zeta_2)}.$$

Next, we let both contours get arbitrary close to the imaginary axis. Taking the definition of the square root into account and performing the limiting process  $\zeta_{1,2} \rightarrow \zeta$ , such that  $\text{Re}(\zeta_{1,2}) \rightarrow 0$ , with

$$\begin{aligned} \zeta_1 &\rightarrow \zeta & \text{and} & \quad \sqrt{\zeta_1^2+1} \rightarrow -\sqrt{\zeta^2+1} \\ \zeta_2 &\rightarrow \zeta & \text{and} & \quad \sqrt{\zeta_2^2+1} \rightarrow \sqrt{\zeta^2+1} \end{aligned}$$

we obtain

$$(3.3.3) \quad \hat{u}_+(s) = -\frac{1}{\pi i} \int_i^{i\infty} \frac{d\zeta}{\sqrt{\zeta^2+1}(s-\zeta)}.$$

In order to find the standard form of the integral along the real axis, we compute

$$(3.3.4) \quad \hat{u}_+(s) = \frac{1}{\pi i} \int_{i-\infty}^{i-0} \frac{d\zeta}{\sqrt{\zeta^2+1}(s-\zeta)}$$

$$(3.3.5) \quad = \frac{1}{\pi i} \int_0^\infty \frac{d\tau}{\sqrt{\tau^2-2i\tau}(s-(i-\tau))}.$$

In (3.3.4) we applied Cauchy's residue theorem in the left upper part of the complex plane and in (3.3.5) we introduced  $\zeta \mapsto \tau = \zeta - i$ ,  $\tau \in \mathbb{R}_+$  and reversed the direction of integration. Both equations (??) and (3.3.5) possess an interesting interpretation. Using  $\zeta = i\tau$ , applying  $L^{-1}\left(\frac{1}{s+\tau}\right) = \exp(-\tau x)$  to (??), and reversing the order of integration we find

$$(3.3.6) \quad u_+(x) = \frac{1}{\pi} \int_1^\infty \frac{e^{i\tau x} d\tau}{\sqrt{\tau^2 - 1}},$$

i. e. , we can represent  $u(x)$  as a superposition of *undamped* Fourier modes with positive frequencies (i. e. wavenumbers)  $\tau \geq 1$ . On the other hand, (3.3.5) allows a further representation of the same function  $u(x)$ , namely,

$$(3.3.7) \quad \begin{aligned} u_+(x) &= \frac{1}{i\pi} \int_0^\infty \frac{e^{(i-\tau)x} d\tau}{\sqrt{\tau^2 - 2i\tau}}, \\ &= \frac{e^{i(x-\frac{\pi}{4})}}{\pi\sqrt{2x}} \int_0^\infty \frac{e^{-\tau} d\tau}{\sqrt{\tau} \sqrt{1 - \frac{\tau}{2ix}}}. \end{aligned}$$

That is, we have equally valid  $u(x)$  as a superposition of *damped* Fourier modes of a single frequency, due to the factor  $\exp(ix)$ , and with positive damping parts  $\tau \geq 0$ . Observe, that for  $x \rightarrow \infty$  the integral converges to  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ , which in turn supplies the well-known asymptotic formula for the Hankel function

$$u_+(x) = H_0^{(1)}(x) = \frac{e^{i(x-\pi/4)}}{\sqrt{2\pi x}} \left(1 + O\left(\frac{1}{x}\right)\right).$$

Differentiating (3.3.7) with respect to  $x$  and expanding the square-root term into a Taylor series supplies the asymptotic behavior

$$(3.3.8) \quad \frac{\partial_r u_+(r)}{u(r)} = i + O\left(\frac{1}{r}\right) \quad \text{for } r \rightarrow \infty.$$

Obviously, (3.3.8) tells us that  $u_+$  behaves asymptotically and locally like an outgoing plane wave.

REMARK. Each of the integral representations (3.3.6) and (3.3.7) coincides, up to normalization constants, to one of the two types of the classical integral representation formulas of Hankel's function. Referring to the book of Taylor, Eq. (3.3.6) corresponds to Eq. (6.32) in [88, Chap. 3.6, p. 230] and (3.3.7) corresponds to Eq. (6.33) [88, Chap. 3.6, p. 230].

Despite the simplicity of the above example, it contains all characteristic properties of the general Helmholtz-type scattering problem with variable coefficients, which we will analyze in Chap. 4. These basic properties are:

- (1) The Laplace transformed solution contains two singularities. The singularities determine the far-field behavior of the function in the spatial domain. The singularity with positive imaginary part corresponds to an *outgoing* field, and the singularity with negative imaginary part corresponds to an *incoming* field. The "far-field" property of the Laplace transformed function, that is the behavior of  $\hat{u}$  for large  $|s|$  determines the boundary condition at the artificial boundary.
- (2) Different cuts connecting the singularities to infinity correspond to different forms of solution representations.

**Bessel's equation.** The radially symmetric Helmholtz equation in 2D separates into Bessel's equation for the distance variable  $r$  and an eigenvalue problem for the



angular variable, cf. Section 4.1, p. 59. The starting point of the classical derivation of transparent boundary conditions is always to find proper solutions of the Bessel equation

$$u''(r) + \frac{1}{r}u'(r) + \left(1 - \frac{\nu^2}{r^2}\right)u = 0, \quad r > a,$$

with initial data  $u(a)$  and  $u'(a)$  exploiting the known properties of the Hankel functions at infinity. For the alternative derivation, we start again from the Bessel equation. We introduce a new coordinate system with the independent variable  $x$  obtained via  $r \mapsto x := r - a, x \geq 0$  and denote, for convenience, the shifted function  $u(x + a)$  again by  $u(x)$ . This way, we obtain our basic differential equation

$$(3.3.9) \quad u''(x) + \frac{1}{x+a}u'(x) + \left(1 - \frac{\nu^2}{(x+a)^2}\right)u = 0, \quad x > 0.$$

Now our concept is the following:

- (1) Laplace transform of (3.3.9), which yields again a differential equation, this time in the dual domain.
- (2) Solve the new differential equation by means of the variation of constant formula.
- (3) Identify the singularities which correspond to incoming waves.
- (4) Remove them by a proper choice of the initial conditions at  $x = 0$ .

In contrast to the classical approach, we can carry out the whole procedure without using properties of Hankel's functions.

**Step 1.** We multiply (3.3.9) with  $x+a$ , introduce the new variable  $v(x) = u(x)/(x+a)$ , and take into account the rule  $L(xu(x)) = -\hat{u}'(s)$ , see A.1.2, where the prime denotes the complex derivative with respect to the argument  $s$ , to obtain

$$\frac{d}{ds} \begin{bmatrix} s^2 + 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} = \begin{pmatrix} a(s^2 + 1) + s & -\nu^2 \\ -1 & a \end{pmatrix} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} - a \begin{pmatrix} su(0) + u'(0) \\ 0 \end{pmatrix}.$$

Equivalently, we have

$$(3.3.10) \quad \frac{d}{ds} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} = \begin{pmatrix} a - \frac{s}{s^2+1} & -\frac{\nu^2}{a} \\ -1 & a \end{pmatrix} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} - a \begin{pmatrix} \frac{su(0)+u'(0)}{s^2+1} \\ 0 \end{pmatrix}.$$

In view of the solution procedure, it is convenient, to extract the translation part and to denote the remaining matrix with  $A$

$$(3.3.11) \quad \frac{d}{ds} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} = a \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} + A \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} - a \begin{pmatrix} \frac{su(0)+u'(0)}{s^2+1} \\ 0 \end{pmatrix}$$

where  $I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$  and  $A = \begin{pmatrix} -\frac{s}{s^2+1} & -\frac{\nu^2}{s^2+1} \\ -1 & 0 \end{pmatrix}$

**Step 2.** Our goal is a solution formula of the type

$$(3.3.12) \quad \begin{pmatrix} \hat{u}(s) \\ \hat{v}(s) \end{pmatrix} = \Phi \int_s^\infty ds_1 \Phi^{-1} \mathbf{r}(u(0), u'(0)).$$

$\Phi$  is the matrix of fundamental solutions and must obey the matrix differential equation

$$(3.3.13) \quad \Phi'(s) = (x_0 I + \mathbf{A}(s)) \Phi(s),$$

and  $\mathbf{r}(u(0), u'(0))$  is the source vector of (3.3.11) which depends on the initial data. The determination of  $\Phi$  evolves in three sub-steps: (i) remove the translation part, (ii) factorize the solution to drop the part which corresponds to  $\nu = 0$ , (iii)

transform the variable  $s \rightarrow \tau$ , to arrive at a problem with constant coefficients.

*Step (i).* We introduce

$$\Phi = \Phi^{(1)} e^{as},$$

which results in

$$\Phi^{(1)'}(s) = \mathbf{A}(s)\Phi^{(1)}(s).$$

This way, a column  $(u, v)^T$  of  $\Phi$ , where we have dropped here the hat used so far to mark a Laplace transformed quantity, is transformed into a corresponding column of  $\Phi^{(1)}$

$$\begin{pmatrix} u \\ v \end{pmatrix} \mapsto \begin{pmatrix} u^{(1)} \\ v^{(1)} \end{pmatrix}.$$

*Step (ii).* We factorize the first components of every column of  $\Phi^{(1)}$

$$\begin{aligned} u^{(1)} &= u^{(2)}w, & w &= \frac{1}{\sqrt{s^2+1}} \\ v^{(1)} &= v^{(2)}, \end{aligned}$$

which results in

$$\Phi^{(2)'}(s) = -\frac{1}{\sqrt{s^2+1}} \begin{pmatrix} 0 & \nu^2 \\ 1 & 0 \end{pmatrix} \Phi^{(2)}(s)$$

*Step (iii).* Finally, we transform  $\Phi^{(2)}(s(\tau)) = \Phi^{(3)}(\tau)$  with

$$\frac{d\tau}{ds} = \frac{1}{\sqrt{s^2+1}}, \quad \tau = -\log(\sqrt{s^2+1} - s)$$

to obtain

$$\Phi^{(3)'}(s) = -\begin{pmatrix} 0 & \nu^2 \\ 1 & 0 \end{pmatrix} \Phi^{(3)}(s),$$

which in turn has a solution

$$\Phi^{(3)}(s) = \begin{pmatrix} e^{\nu\tau} & e^{-\nu\tau} \\ -\frac{1}{\nu}e^{\nu\tau} & \frac{1}{\nu}e^{-\nu\tau} \end{pmatrix}.$$

Tracing the whole procedure backwards, we find

$$(3.3.14) \quad \Phi^{(3)} \mapsto \Phi^{(2)} \mapsto \Phi^{(1)} \mapsto \Phi = \begin{pmatrix} \frac{(\sqrt{s^2+1}-s)^{-\nu}}{\sqrt{s^2+1}} & \frac{(\sqrt{s^2+1}-s)^\nu}{\sqrt{s^2+1}} \\ -\frac{1}{\nu}(\sqrt{s^2+1}-s)^{-\nu} & \frac{1}{\nu}(\sqrt{s^2+1}-s)^\nu \end{pmatrix} e^{as}.$$

Applying (3.3.12), we obtain finally the desired representation of  $\hat{u}(s)$ .

$$(3.3.15) \quad \hat{u}(s) = \frac{a}{2} \frac{(\sqrt{s^2+1}-s)^{-\nu}}{\sqrt{s^2+1}} e^{as} \int_s^\infty e^{-as_1} \frac{(\sqrt{s_1^2+1}-s_1)^\nu}{\sqrt{s_1^2+1}} (s_1 u(0) + u'(0)) ds_1 + \frac{a}{2} \frac{(\sqrt{s^2+1}-s)^\nu}{\sqrt{s^2+1}} e^{as} \int_s^\infty e^{-as_1} \frac{(\sqrt{s_1^2+1}-s_1)^{-\nu}}{\sqrt{s_1^2+1}} (s_1 u(0) + u'(0)) ds_1$$

**Step 3.** According to (3.3.15), we can represent  $\hat{u}(s)$  in factorized form

$$\hat{u}(s) = \frac{1}{\sqrt{s^2+1}} g(s),$$

with

$$(3.3.16) \quad g(s) = \frac{a}{2}(\sqrt{s^2+1}-s)^{-\nu}e^{as} \int_s^\infty e^{-as_1} \frac{(\sqrt{s_1^2+1}-s_1)^\nu}{\sqrt{s_1^2+1}} (s_1 u(0) + u'(0)) ds_1 + \frac{a}{2}(\sqrt{s^2+1}-s)^\nu e^{as} \int_s^\infty e^{-as_1} \frac{(\sqrt{s_1^2+1}-s_1)^{-\nu}}{\sqrt{s_1^2+1}} (s_1 u(0) + u'(0)) ds_1.$$

Consequently,  $\hat{u}(s)$  may have singularities at  $s = \pm i$ . In order to drop the singularity at  $s = -i$ , the following necessary condition must be fulfilled.

Pole condition for the Bessel equation

$$g(i) = 0.$$

**Step 4.** This, in turn, supplies the necessary condition

$$\begin{aligned} u'(0) &= \text{DtN}_\nu u(0) \\ \text{DtN}_\nu &= - \frac{\int_{-i}^\infty ds e^{-as} s \frac{i^{-\nu}(\sqrt{s^2+1}-s)^\nu + i^\nu(\sqrt{s^2+1}-s)^{-\nu}}{\sqrt{s^2+1}}}{\int_{-i}^\infty ds e^{-as} \frac{i^{-\nu}(\sqrt{s^2+1}-s)^\nu + i^\nu(\sqrt{s^2+1}-s)^{-\nu}}{\sqrt{s^2+1}}}. \end{aligned}$$

In [86] we showed by estimating the integrals that the DtN-number has the following properties

- (1)  $\text{DtN}_\nu \xrightarrow{\nu \rightarrow \infty} -\nu/a$ .
- (2)  $|\text{DtN}_\nu/\nu| \leq K < \infty$  for all  $\nu$ .
- (3)  $\text{Im}(\text{DtN}_\nu) > 0$  for all  $0 \leq \nu < \infty$ .
- (4)  $|\text{DtN}_\nu| \geq K > 0$  for all  $\nu \geq 0$ .

Similar results were shown by Harari and Hughes [55] and Ihlenburg [60, Chapter 3.2.2] dealing with properties of Hankel's functions. In Chapter 4 we shall again discuss these spectral properties from the viewpoint of our general theory.

### 3.4. Generalizations

In the following we show, how to extend the above constructions to higher dimensions. Starting with a separable problem, which consists of a sequence of 1D problems, we generalize the concept to general situations.

**Fields Outside a Sphere.** In order to formulate the pole condition for dimensions  $d > 1$ , we consider first fields outside a given sphere with radius  $a$ . There are two ways to generalize the concept. The first one uses a *separation* into functions depending on angular- and distance coordinates, the second uses a more general *parameterization* of the exterior domain. The first way is the more simpler one, the second way allows for a more general class of problems.

In the one-dimensional case we considered exterior functions  $u(x)$ ,  $x \geq a$ . Now we study functions  $u : \mathbb{R}^d \setminus \Omega_a \rightarrow \mathbb{C}$  with  $\Omega_a$  the interior of a sphere with radius  $a$ . Let  $Q : (a + R_+) \times S^{d-1} \rightarrow \mathbb{R}^d \setminus \Omega_a$  be an homeomorphism mapping generalized spherical coordinates to Euclidean coordinates outside the sphere. Let us write  $(uQ)(\mathbf{x}) =: \tilde{u}(\rho, \mathbf{x}_{S^{d-1}})$  and drop the tilde. Motivated by the diagonalization of

the Laplace-Beltrami operator on  $S^{d-1}$  based on the complete orthonormal system  $\{(\phi_j, \lambda_j) : j = 0, \dots, \infty\}$  of eigenvectors and eigenvalues, we consider first the case that  $u$  has a representation

$$(3.4.1) \quad u(\rho, \mathbf{x}_{S^{d-1}}) = \sum_{j=0}^{\infty} \phi_j(\mathbf{x}_{S^{d-1}}) v_j(\rho), \quad \langle \phi_i, \phi_j \rangle_{L^2(S^{d-1})} = \delta_{ij}, \quad \rho \geq a.$$

For  $d = 2$  the functions  $\phi_j$  are trigonometric monomials, for  $d = 3$  spherical harmonics. This will be the starting point of our analysis in Section 4.2. Here we want to use it only for our qualitative discussion. Now, a natural generalization of the one-dimensional pole conditions to higher dimensions is

**DEFINITION 3.4.1.** (Pole condition for generalized Fourier modes on  $\partial\Omega$ .) Functions with a separable representation of type (3.4.1) satisfy the pole condition if each of the functions  $v_j$ ,  $j = 0, \dots, \infty$ , satisfies the one-dimensional pole condition with respect to the distance variable  $(\rho - a)$ .

However, we want to have a characterization of  $u$  independent of the underlying partial differential equation. Therefore, in the general case, it does not make sense to use operator properties like the diagonalization above to define the pole condition. Therefore, we extend this to the case where no explicit series representation of  $u$  is given.

**DEFINITION 3.4.2.** (Pole condition for functions exterior to a sphere.) A function  $u(\cdot, \mathbf{x}_{S^{d-1}})$  satisfies the pole condition if it satisfies the one-dimensional pole condition for all  $\mathbf{x}_{S^{d-1}} \in S^{d-1}$ .

This definition is motivated by the following observation. Let  $d = 2$  and let us choose the  $\phi_j$  from (3.4.2) with a very small support  $\text{meassupp } v_j \rightarrow 0$ . Then, the support of the corresponding function in  $\mathbb{R}^2 \setminus \Omega_a$  is an angle segment with an interior angle  $\rightarrow 0$ . In the limiting case, the segment degenerates to a semi-infinite line. Our practical realization of the algorithms, however, follow a generalization of the first concept (3.4.1). Let us consider the case that  $u$  has a representation

$$(3.4.2) \quad u(\rho, \mathbf{x}_{S^{d-1}}) = \sum_{j=0}^{\infty} \psi_j(\mathbf{x}_{S^{d-1}}) v_j(\rho), \quad \rho \geq a$$

where the functions  $\psi_j : S^{d-1} \rightarrow \mathbb{C}$  form a set of linear independent functions. Clearly, (3.4.2) is a generalization of our first approach (3.4.1). In the numerical algorithms, the functions  $\psi_j$  will play a role as trial functions resulting automatically from a finite-element discretization, whereas the functions  $v_j$  has to be computed according to the pole condition. We call this representation a local separable representation, since we assume a factorization with respect to distance and angular coordinates, but we don't require a global orthogonality of any kind of the functions  $\psi_j$ . Since the pole condition for representations of type (3.4.2) plays a dominant role for all of our the numerical concepts, we want to give it an extra definition

**DEFINITION 3.4.3.** (Pole condition for locally separable representations.) Functions with a separable representation of type (3.4.2) satisfy the pole condition if each of the functions  $v_j$ ,  $j = 0, \dots, \infty$ , satisfies the one-dimensional pole condition with respect to the distance variable  $(\rho - a)$ .

Our numerical algorithms for the general Helmholtz scattering problem with potentials depending both on the angular and the distance coordinates will be based on the representation (3.4.2) and pole condition 3.4.3.

**General Case.** We generalize the above given special forms of the pole condition and give it in a general form. To this end we assume the following

CONDITION 3.4.4. Let  $\Omega \subset \mathbb{R}^d$  be the interior domain and  $\partial\Omega$  its boundary. We assume that there exists a homeomorphism  $Q$  such that

- (1)  $Q : \mathbb{R}_+ \times \partial\Omega \rightarrow \mathbb{R}^d \setminus \Omega$ .
- (2)  $\lim_{\xi \rightarrow \infty} |Q(\xi, \mathbf{x}_{\partial\Omega})| / \xi = C > 0$  for all  $\mathbf{x}_{\partial\Omega} \in \partial\Omega$ .

Item (1) supplies a generalized distance variable  $\xi$ , which plays the same role as the radius  $\rho - a$  above. Since the interpretation of the distance variable with an radius is not always useful, we decided to use the symbol  $\xi$  (motivated by  $|\mathbf{x}|$ ). Item (2) ensures that asymptotically the metric induced by  $Q$  is an Euclidean metric. By means of  $Q$  we can generalize the pole condition in a direct manner. Moreover, the general pole condition implies our understanding of outgoing functions.

DEFINITION 3.4.5. *General pole condition.* A function  $u(\mathbf{x})$  satisfies the pole condition if the Laplace transform of  $(uQ)(\cdot, \mathbf{x}_{\partial\Omega})$  has a holomorphic extension to the lower half-plane  $\{s \in \mathbb{C} : \text{Im } s < 0\}$  for all  $\mathbf{x}_{\partial\Omega} \in \partial\Omega$ .

DEFINITION 3.4.6. A function  $u(\mathbf{x})$  satisfying the pole condition Definition 3.4.5 is called outgoing.

Our definition of the general pole condition is based on a given parameterization  $Q$ . However, whether a function satisfies the pole condition or not does not depend on the special choice of  $Q$ . Using additional assumptions, we will show that if a function  $u$  satisfies the pole condition along a given path, it satisfies the pole condition in the neighborhood of this path, too. To discuss a particular example, we introduce to following notation:

- $\gamma_0 := \{(\xi, \eta_0), \xi \geq 0\}$  is a given path.
- $\hat{u}_0(s)$  is the Laplace transform of  $u_0(\xi, \eta_0) := (uQ)(\xi, \eta_0)$  along this path.
- $\gamma_1 := \{(\xi, \eta_0 + \delta\eta(\xi)), \xi \geq 0\}$  is a perturbation of  $\gamma_0$ , where  $\delta\eta(\xi)$  denotes the perturbation in the angular variable in dependence of the distance variable  $\xi$ .
- $\hat{u}_1(s)$  is the Laplace transform of  $u_1(\xi, \eta_0 + \delta\eta) := (uQ)(\xi, \eta_0 + \delta\eta)$  along this path.

Our basic additional assumption is that the function  $u(\xi, \eta)$  in the neighborhood of the path  $\gamma_0$  has a representation  $u(\xi, \eta) = f(\eta - \eta_0)u_0(\xi, \eta_0)$  and  $f$  has a convergent Taylor representation around 0.

LEMMA 3.4.7. *Let a parameterization  $Q(\xi, \eta)$  be given. Let  $\hat{u}_0(s)$  satisfy the pole condition along the path  $\gamma_0$ . Let there be  $\delta\eta$  and  $u(\xi, \eta)$  such that*

- (1)  $\delta\eta := p\left(\frac{1}{\xi+a}\right)$ , where  $p(t) = \sum_{m=1}^{\infty} p_m t^m$  with a convergence radius  $\rho > 1/\bar{a}$ ,  $0 < \bar{a} < a$  and
- (2)  $u(\xi, \eta) = \left(\sum_{j=0}^{\infty} c_j (\delta\eta)^j\right) u(\xi, \eta_0)$  absolutely and uniformly for  $0 < |\eta - \eta_0| \leq \rho_\eta$  and  $\xi \geq 0$ .

*Then,  $\hat{u}_1(s)$  computed along the perturbed path  $\gamma_1$  has a holomorphic extension to the lower half of the complex plane.*

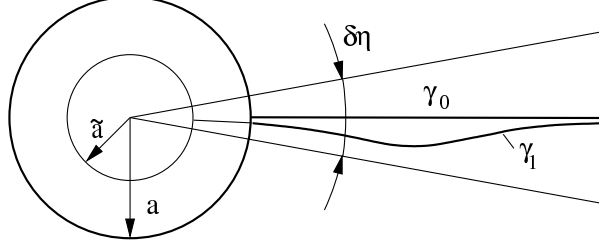


FIGURE 3.4.1. The path  $\gamma_0$ , on which the pole condition is satisfied and the disturbed path  $\gamma_1$ .

PROOF. We refer to Fig. 3.4.1. It holds

$$(\delta\eta)^j = \sum_{m=j}^{\infty} p_{m,j} \left( \frac{1}{\xi+a} \right)^m \quad \text{for } \xi > \tilde{a} - a, j \geq 1$$

with real coefficients  $p_{m,j}$  since the convergence radius of the product of two power series with convergence radius  $\rho$  is again  $\rho$ . Thus we get

$$\begin{aligned} u(\xi, \eta) &= u(\xi, \eta_0) \left( 1 + \sum_{j=1}^{\infty} c_j \left( \sum_{m=j}^{\infty} p_{m,j} \left( \frac{1}{\xi+a} \right)^m \right) \right) \\ &= u(\xi, \eta_0) + u(\xi, \eta_0) \sum_{j=1}^{\infty} \left( \frac{1}{\xi+a} \right)^j \left( \sum_{m=1}^j c_j p_{j,m} \right) \\ &= u(\xi, \eta_0) + u(\xi, \eta_0) \sum_{j=1}^{\infty} \left( \frac{1}{\xi+a} \right)^j q_j. \quad \text{for all } \xi \geq 0, 0 < |\eta - \eta_0| \leq \rho_\eta, \end{aligned}$$

where the coefficients  $q_j$  are defined by  $q_j := \sum_{m=1}^j c_j p_{j,m}$ . Since both series are absolute convergent, the reordering is legitimate and supplies a convergent power series. Thus we have a decomposition of function on the perturbed path into the function on the unperturbed path, which satisfies the pole condition and a remainder. The series is absolutely convergent for  $\xi \geq \tilde{a} - a$ , hence  $\lim_{j \rightarrow \infty} |q_j| / \tilde{a}^j \rightarrow 0$  and therefore  $|q_j| < \tilde{a}^j$  for  $j \rightarrow \infty$ . We apply the Laplace transform and find

$$\hat{u}(s) = \hat{u}_0(s) + \int_0^\infty d\xi e^{-s\xi} u(\xi, \eta_0) \left( \sum_{j=1}^{\infty} \left( \frac{1}{\xi+a} \right)^j q_j \right), \quad \text{Re } s > 0.$$

By the convolution lemma A.1.8 of the appendix this yields

$$(3.4.3) \quad \hat{u}(s) = \hat{u}_0(s) + \int_s^\infty ds_1 f(s_1 - s) \hat{u}_0(s_1)$$

$$(3.4.4) \quad f(s) = e^{-as} \sum_{j=1}^{\infty} \frac{q_j}{(j-1)!} s^{j-1}.$$

We have to show that the integral expression in (3.4.3) is holomorphic in the lower half of the complex plane. To this end we study first the convergence of  $g(s) := \sum_{j=1}^{\infty} q_j / (j-1)! s^{j-1}$ . By the root criterion

$$\begin{aligned}
\lim_{j \rightarrow \infty} \left( \frac{|q_j|}{(j-1)!} |s|^{j-1} \right)^{\frac{1}{j}} &\leq \lim_{j \rightarrow \infty} \left( \frac{\tilde{a}^j}{(j-1)!} |s|^{j-1} \right)^{\frac{1}{j}} \\
&\leq a \lim_{j \rightarrow \infty} \frac{|s|^{-\frac{1}{j}}}{((j-1)!)^{\frac{1}{j}}} \\
&= 0.
\end{aligned}$$

Thus the series converges on the entire complex plane. Next we have to show that the complex derivative of  $\exp(-as) g(s)$ , hence the complex derivative of  $g(s)$  exists. We can compute  $g'(s)$  term-wise, since

- (i)  $g(s)$  converges uniformly on the entire complex plane,
- (ii) the derivatives of the individual terms of the series are continuous, and
- (iii) the sum of the differentiated terms converges uniformly on the entire complex plane, which can be shown by the same technique as before.

Hence, the function  $f(s)$  is holomorphic on the entire complex plane. Further,  $g(s)$  is exponentially bounded, which follows from

$$\begin{aligned}
|g(s)| &= \left| \sum_{j=1}^{\infty} \frac{q_j}{(j-1)!} s^{j-1} \right| \\
&\leq C \left| \sum_{j=1}^{\infty} \frac{1}{(j-1)!} |\tilde{a}s|^{j-1} \right| \\
&= C e^{|\tilde{a}s|}
\end{aligned}$$

with  $C = \sup_{j \geq 1} |q_j / \tilde{a}^{j-1}|$ . The same holds for  $|g'(s)|$  with  $C = \sup_{j \geq 1} |q_{j+1} / \tilde{a}^{j-1}|$ . Thus it holds

$$\begin{aligned}
|f(s)| &\leq C e^{-a \operatorname{Re} s + |\tilde{a}s|} \\
|f'(s)| &\leq C e^{-a \operatorname{Re} s + |\tilde{a}s|}.
\end{aligned}$$

Using this we can show that (3.4.3) together with (3.4.4) is holomorphic in the lower half of the complex plane. A change of variables yields

$$(3.4.5) \quad \hat{u}(s) = \hat{u}_0(s) + \int_0^{\infty} ds_1 f(s_1) \hat{u}_0(s_1 + s).$$

We prove that the complex derivative with respect to  $s$  exists in the lower half of the complex plane. The following holds true:

- (i) for each fixed  $s_1$  the complex derivative of  $f(s_1) \hat{u}_0(s_1 + s)$  with respect to  $s$  exists,
- (ii) for each fixed  $s$  the integral exists, due to the exponential bounds,
- (iii) since  $|\partial_s (f(s_1) \hat{u}_0(s_1 + s))| = |f(s_1) \partial_s \hat{u}_0(s_1 + s)| \leq C |f(s_1)|$  there exists a function  $F(s_1) := C e^{-a \operatorname{Re} s_1 + |\tilde{a}s_1|}$  integrable on the interval  $\operatorname{Re}(s_1) \in [0, \infty)$ ,  $\operatorname{Im} s_1 < 0$  such that  $F(\operatorname{Re} s_1) \geq |(f(s_1) \hat{u}_0(s_1 + s))'|$  for all  $s$  in the lower half of the complex plane.

Hence we can apply Lebesgue's Dominated Convergence Theorem (cf. e. g. Zeidler [95, Appendix, p. 440]) which shows that the integral term can be differentiated

with respect to the parameter  $s$  and the result remains bounded. Thus the complex derivative of  $\widehat{u}(s)$  exists in the entire lower half of the complex plane.  $\square$

**Time-Dependent Schrödinger Equation.** We want to extend the definition of outgoing waves to time-dependent problems, where we concentrate only to so-called abstract Schrödinger type problems

$$\begin{aligned}\partial_t u(\mathbf{x}, t) &= -iAu(\mathbf{x}, t), & t > 0 \\ u(\mathbf{x}, 0) &= u_0(\mathbf{x})\end{aligned}$$

where the operator  $A : D(A) \subset X \rightarrow X$  is a self-adjoint operator on the complex Hilbert space  $X$ . Prototypes of this equation are the Schrödinger equation and the Fresnel equation discussed in Section 1.1. In these examples we have  $A : D(A) \subset L^2(\mathbb{R}^d) \rightarrow L^2(\mathbb{R}^d)$ , and together with the symmetry of  $A$  we find immediately the conservation property  $\partial_t(u, u)_{L^2(\mathbb{R}^d)} = 0$ . In what follows, we will not study existence and uniqueness of the abstract evolution problem for operators  $A$  depending both on space and time, but we want to extend the pole condition to functions  $(\mathbf{x}, t) \mapsto u(\mathbf{x}, t)$  and to apply this condition to derive discrete transparent boundary conditions satisfying discrete analogies to the continuous conservation property. To this end we give the obvious extension of the pole condition Definition 3.4.5 which refers only to spatial problems to Schrödinger-type evolution problems.

DEFINITION 3.4.8. A continuous, time-dependent function  $u(x, t)$  satisfies the pole condition if it satisfies the time-independent pole condition Definition 3.4.5 for each fixed time  $t \geq 0$ .

REMARK. The Schrödinger equation with a self-adjoint operator  $A$  is time-reversible, hence it could be natural to extend this definition to the time interval  $-\infty < t < \infty$ . On the other hand, practical computations always start with initial data given at a fixed initial time. Thus it seems to be sufficient to restrict the definition to the half-interval.



## Existence and Uniqueness Statements for Separable Problems

We restate our basic Helmholtz scattering problem (1.0.1) in a slightly modified form. Since our focus is on the solution of the exterior problem and the coupling to the interior problem, we do not consider interior sources and boundaries. This is done to simplify the presentation and does not mean any restriction with respect to the theory. Our starting point is the system

$$(4.0.6) \quad \begin{aligned} \Delta u(\mathbf{x}) + k^2(\mathbf{x}) u(\mathbf{x}) &= 0 && \text{in } \Omega_{\text{int}} \\ \Delta u_{\text{out}}(\mathbf{x}) + k^2(\mathbf{x}) u_{\text{out}}(\mathbf{x}) &= 0 && \text{outside } \overline{\Omega_{\text{int}}} \\ u(\mathbf{x}) &= u_{\text{src}}(\mathbf{x}) + u_{\text{out}}(\mathbf{x}) && \text{on } \partial\Omega \\ \partial u(\mathbf{x}) &= \partial u_{\text{src}}(\mathbf{x}) + \partial u_{\text{out}}(\mathbf{x}) && \text{on } \partial\Omega \\ \partial_n u_{\text{out}}(\mathbf{x}) &= B u_{\text{out}}(\mathbf{x}) && \text{on } \partial\Omega. \end{aligned}$$

Introducing the sesquilinear forms

$$(4.0.7) \quad a : H^1(\Omega_{\text{int}}) \times H^1(\Omega_{\text{int}}) \rightarrow \mathbb{C}$$

$$(4.0.8) \quad b : H^{1/2}(\Omega_{\text{int}}) \times H^{1/2}(\Omega_{\text{int}}) \rightarrow \mathbb{C}$$

$$\text{by} \quad \begin{aligned} a(v, u) &= (\nabla v, \nabla v) - (v, k^2(\mathbf{x})u) \\ b(v, u) &= \langle v|_{\partial\Omega_{\text{int}}}, B u|_{\partial\Omega_{\text{int}}} \rangle, \end{aligned}$$

which are derived by a formal application of Green's theorem, we obtain the *interior* problem in weak form: Find  $u \in H^1(\Omega_{\text{int}})$  such that for all  $v \in H^1(\Omega_{\text{int}})$

$$(4.0.9) \quad a(v, u) - b(v, u) = \langle v|_{\partial\Omega_{\text{int}}}, \partial_n u_{\text{src}} \rangle - b(v, u_{\text{src}}).$$

Here the source data  $u_{\text{src}}|_{\partial\Omega_{\text{int}}}$ ,  $\partial_n u_{\text{src}}|_{\partial\Omega_{\text{int}}}$  must be given such that the corresponding functionals are continuous. In the sequel of this chapter we consider mainly the *exterior* problem and its coupling to the interior problem.

### 4.1. Separable Coordinates

Since we deal in the following only with the exterior domain  $\Omega_{\text{ext}}$ , we suppress the subscript and write simply  $\Omega$  instead of  $\Omega_{\text{ext}}$ . Accordingly, we suppress the subscript *out* and write simply  $u$  instead of  $u_{\text{out}}$ . Further we write (4.0.6) explicitly in Cartesian coordinates

$$(4.1.1) \quad \begin{aligned} \Delta_{xy} u(x, y) + k^2 u(x, y) &= 0 && \text{in } \Omega_{xy} \\ \partial_n u(x, y) &= -B(x, y)u(x, y) && \text{on } \partial\Omega_{xy}, \end{aligned}$$

where  $\Delta_{xy} = \nabla_{xy}^T \nabla_{xy}$  with  $\nabla_{xy}$  the column vector  $(\partial_x, \partial_y)^T$  and  $\Omega_{xy}$  and  $\partial\Omega_{xy}$  are parameterizations of the exterior domain  $\Omega$  and the the boundary  $\partial\Omega$  with respect to Cartesian coordinates, that is

$$\begin{aligned}\Omega_{xy} &= \{(x, y) \in \mathbb{R}^2 : (x, y) \in \Omega\} \\ \partial\Omega_{xy} &= \{(x, y) \in \mathbb{R}^2 : (x, y) \in \partial\Omega\}.\end{aligned}$$

We want to transform the description based on Cartesian  $x, y$ -coordinates to other coordinates, denoted by  $\xi$  and  $\eta$ . To this end we introduce the  $C^1$ -diffeomorphism  $Q$  from the new domain

$$\Omega_{\xi\eta} = (\xi_0, \infty) \times [\eta_{\min}, \eta_{\max})$$

onto the original domain  $\Omega_{xy}$  by

$$\begin{aligned}Q : \Omega_{\xi\eta} &\rightarrow \Omega_{xy} \\ Q(\xi, \eta) &= \begin{pmatrix} x \\ y \end{pmatrix}.\end{aligned}$$

We consider the independent variable  $\xi$  with  $0 < \xi_0 \leq \xi$  as generalized distance variable and  $\eta$  with  $\eta_{\min} \leq \eta \leq \eta_{\max}$  as generalized angular variable. Let  $Q$  satisfy Condition 3.4.4. We denote the Jacobi matrix  $Q'$  by

$$J := Q' = \begin{pmatrix} \partial_\xi x & \partial_\eta x \\ \partial_\xi y & \partial_\eta y \end{pmatrix}$$

and its determinant by  $|J| = \partial(x, y)/\partial(\xi, \eta)$ . Given a differentiable function  $f$ , it holds, by the definition of the gradient,

$$\nabla_{xy} f = J^{-T} \nabla_{\xi\eta} f,$$

where  $f$  on the left-hand side has to be read as  $f(x, y)$  and  $f$  on the right-hand side has to be considered as function of  $\xi$  and  $\eta$  via  $f(x(\xi, \eta), y(\xi, \eta))$ . Considering the integral  $\int_{\Omega_{xy}} g \Delta_{xy} f \, dx \, dy$  with an arbitrary twice differentiable function  $f$  and a test function  $g \in C_0^\infty(\Omega_{xy})$  we find the transform of the Laplacian in terms of the the transformed gradient

$$\begin{aligned}\int_{\Omega_{xy}} g \Delta_{xy} f \, dx \, dy &= - \int_{\Omega_{xy}} (\nabla_{xy} g)^T \nabla_{xy} f \, dx \, dy \\ &= - \int_{\Omega_{\xi\eta}} (J^{-T} \nabla_{\xi\eta} g)^T J^{-T} \nabla_{\xi\eta} f |J| \, d\xi \, d\eta \\ &= \int_{\Omega_{xy}} g |J|^{-1} \nabla_{\xi\eta}^T (J^{-1} J^{-T} |J| \nabla_{\xi\eta}) f \, dx \, dy.\end{aligned}$$

The transformed Laplacian follows

$$(4.1.2) \quad \Delta_{xy} f = |J|^{-1} \nabla_{\xi\eta}^T (J^{-1} J^{-T} |J| \nabla_{\xi\eta}) f.$$

Next we restrict our consideration only to those transforms which can be derived from a conformal mapping.

**REMARK.** This does not mean that our consideration is based essentially on a conformal mapping and therefore mainly restricted to two-dimensional problems. We use this only as a convenient way to derive separated equations. For an example in higher space dimensions see Section 4.1.

Let us consider the coordinate transform  $Q$  as complex mapping  $(x + iy) = Q(\xi + i\eta)$  with an analytic function  $Q$ . Consequently we have

$$d(x + iy) = Q' d(\xi + i\eta)$$

which supplies

$$J = \begin{pmatrix} \operatorname{Re} Q' & -\operatorname{Im} Q' \\ \operatorname{Im} Q' & \operatorname{Re} Q' \end{pmatrix},$$

from which we find the well-known property of conformal mapping

$$|J| = Q' \overline{Q'} \quad \text{and} \quad J^{-1} J^{-T} |J| = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

Consequently the Helmholtz equation in transformed coordinates based on conformal mapping reads

$$(4.1.3) \quad (\partial_\xi^2 + \partial_\eta^2) u + |J| k^2 u = 0.$$

**DEFINITION 4.1.1.** We call the Helmholtz equation separable with respect to  $\xi, \eta$ -coordinates, if the Helmholtz term  $|J| k^2$  possesses a representation

$$|J(\xi, \eta)| k^2(\xi, \eta) = f_\xi + f_\eta$$

with functions  $f_\xi$  and  $f_\eta$  depending only on  $\xi$  and  $\eta$ , respectively.

If the transformed Helmholtz equation (4.1.3) is separable, we try in the usual way to find solutions of the form  $u(\xi, \eta) = u_1(\xi)u_2(\eta)$  by a solution of the separated equations

$$(4.1.4) \quad \partial_\xi^2 u_1 + f_\xi u_1 - K u_1 = 0$$

$$(4.1.5) \quad \partial_\eta^2 u_2 + f_\eta u_2 = -K u_2.$$

Since we will use in both cases under consideration, namely cylindric and elliptic-hyperbolic coordinates, the scaling  $\xi \mapsto \tilde{\xi} = \exp \xi$  we rewrite (4.1.4) using this transform. Depressing the tilde, we obtain

$$\xi^2 \partial_\xi^2 u_1 + \xi \partial_\xi u_1 + f_\xi u_1 - K u_1 = 0.$$

The exponential re-scaling of the distance variable  $\xi$  introduced by the technique of the conformal mapping ensures – as we will see – the required asymptotic equivalence between a distance measured in Cartesian coordinates and a distance measured by the generalized distance variable  $\xi$ .

**Cylindric Coordinates.** We use the mapping  $x + iy = Q(\xi + i\eta) := \exp(\xi + i\eta)$  with the angular variable in the interval  $-\pi < \eta \leq \pi$ . We obtain  $|J| = \exp(2\xi)$ , and after the scaling  $\xi \mapsto \tilde{\xi} = \exp \xi$  and dropping the tilde, we have the separated equations

$$(4.1.6) \quad \xi^2 \partial_\xi^2 u_1 + \xi \partial_\xi u_1 + (\xi^2 k^2(\xi) - K) u_1 = 0$$

$$(4.1.7) \quad \partial_\eta^2 u_2 = -K u_2.$$

Obviously, the Helmholtz equation in cylindric coordinates separates for  $\xi$ -dependent wavenumbers. For our further consideration, especially for the derivation of the boundary conditions, we allow position-dependent wavenumbers of the type

$$(4.1.8) \quad k^2(\xi) = k_0^2 + \frac{k_{-1}}{\xi} + \dots + \frac{k_{-n}}{\xi^n}.$$

**Elliptic-Hyperbolic Coordinates.** We use the mapping

$$x + iy = Q(\xi + i\eta) := \frac{q}{2} \left( e^{\xi+i\eta} + e^{-(\xi+i\eta)} \right)$$

with the angular variable in the interval  $-\pi < \eta \leq \pi$ . We compute

$$|J| = \frac{q^2}{4} (e^{2\xi} + e^{-2\xi}) - \frac{q^2}{2} \cos 2\eta.$$

Allowing only for an position-independent wavenumber  $k$ , and carrying out the exponential scaling as above, we obtain the separated equations

$$(4.1.9) \quad \xi^4 \partial_\xi^2 u_1 + \xi^3 \partial_\xi u_1 + \left( \frac{k^2 q^2}{4} (\xi^4 + 1) - K \xi^2 \right) u_1 = 0$$

$$(4.1.10) \quad \partial_\eta^2 u_2 - \frac{k^2 q^2}{2} \cos 2\eta u_2 = -K u_2.$$

**Generalization of the Potential.** Both of the separated equations in cylindrical and elliptic-hyperbolic coordinates, (4.1.6) and (4.1.9), possesses the common structure

$$(4.1.11) \quad \partial_\xi^2 u_1 + \xi^{-1} \partial_\xi u_1 + \left( \sum_{j=0}^{\infty} a_j \xi^{-j} \right) u_1 = 0$$

$$(4.1.12) \quad \partial_\eta^2 u_2 + f(\eta) u_2 = -K u_2,$$

where the function  $f(\eta)$  is a smooth, bounded and periodic function in the angle-like variable  $\eta$ .

Thus, these examples show exemplarily that the transform between curvilinear coordinates results merely in a change of the potential, where terms of the form  $a_j/\xi^j$ ,  $j \geq 0$  are added. A similar statement has been obtained in Section 3.4, Lemma 3.4.7, where we considered perturbed paths. Therefore we will base our theory on potentials of the above form.

**Generalization to Higher Space Dimensions.** In Section 4.2 we will need the Laplacian in  $\mathbb{R}^d$  exterior to a sphere. Therefore we add here a generalization of the conformal mapping used above to arbitrary space dimensions. We consider the map

$$Q : \mathbb{R}_+ \times S^{d-1} \rightarrow \mathbb{R}^d \\ Q(\xi, \mathbf{x}^0) = \mathbf{y},$$

with  $\mathbf{y} = (y_1, \dots, y_d)$  the Cartesian coordinates in  $\mathbb{R}^d$ . In generalization of the conformal mapping we assume the following:

- (1) Two orthogonal intersecting curves remain orthogonal intersecting under the map  $Q$ .
- (2) It holds the scaling  $e^{2\xi} = \sum_{j=1}^d y_j^2$ .
- (3)  $Q$  has a representation  $Q = e^\xi \tilde{Q}$ , where  $\tilde{Q} : S^{d-1} \rightarrow \mathbb{R}^{d \times d}$  depends only on the angular coordinates.

As in Section 4.1 we use the Jacobian  $J = Q'$ . Let  $d\mathbf{s}_{1,2} = (d\xi, d\mathbf{x}^0)_{1,2}^T$  be two orthogonal vectors. With  $d\mathbf{y}_{1,2} = J \cdot d\mathbf{s}_{1,2}$  and the orthogonality assumption:  $\langle d\mathbf{y}_1, d\mathbf{y}_2 \rangle = 0$  implies  $\langle d\mathbf{s}_1, d\mathbf{s}_2 \rangle = 0$ , it follows that the matrix

$$J^T J = e^{2\xi} \begin{pmatrix} f & \\ & \Phi \end{pmatrix}, \quad f : S^{d-1} \rightarrow \mathbb{R}, \quad \Phi : S^{d-1} \rightarrow \mathbb{R}^{(d-1) \times (d-1)}$$

is a diagonal matrix. Using further

$$|J| = e^{d\xi} g, \quad \text{with } g : S^{d-1} \rightarrow \mathbb{R},$$

we compute the transformed Laplacian:

$$\begin{aligned}\Delta_{\xi, S^{d-1}} u &= |J|^{-1} \nabla_{\xi, S^{d-1}}^T (J^{-1} J^{-T} |J| \nabla_{\xi, S^{d-1}}) u \\ &= \frac{1}{g e^{d\xi}} \left( \partial_{\xi} \quad \nabla_{S^{d-1}} \right) g e^{(d-2)\xi} \begin{pmatrix} f^{-1}(\mathbf{x}^0) \partial_{\xi} \\ \Phi^{-1}(\mathbf{x}^0) \nabla_{S^{d-1}} \end{pmatrix} u.\end{aligned}$$

As in the conformal mapping approach we let  $r := e^{\xi}$ , which results in

$$\begin{aligned}\Delta_{r, S^{d-1}} u &= \frac{f^{-1}(\mathbf{x}^0)}{r^{d-1}} \partial_r r^{d-1} \partial_r u + \frac{1}{r^2} \Delta_{S^{d-1}} u \\ \Delta_{S^{d-1}} &:= \frac{1}{g(\mathbf{x}^0)} \nabla_{S^{d-1}}^T g(\mathbf{x}^0) \Phi^{-1}(\mathbf{x}^0) \nabla_{S^{d-1}}, \quad \mathbf{x}^0 \in S^{d-1},\end{aligned}$$

where  $\Delta_{S^{d-1}}$  denotes the Laplace-Beltrami operator on  $S^{d-1}$ . Using the second assumption above, the special choice  $u = r^2$  shows  $f = 1$ , and we obtain finally the representation required in the following section

$$(4.1.13) \quad \Delta_{r, S^{d-1}} u = \frac{1}{r^{d-1}} \partial_r r^{d-1} \partial_r u + \frac{1}{r^2} \Delta_{S^{d-1}} u.$$

## 4.2. Preparation of Main Theory

Having outlined the structure of several example problems with potentials depending on the distance, we now want to analyze the scattering theory based on the Laplace transform and the pole condition on a model problem in arbitrary space dimensions. We shall supply the following key results:

- (1) The pole condition is equivalent to Sommerfeld's radiation condition, at least for a finite number of modes. This will be proofed in Theorem 4.4.4.
- (2) It generalizes Sommerfeld's radiation condition to problems with position dependent potentials (Theorem 4.4.4).
- (3) It ensures existence and uniqueness of the interior scattering problem (Section 4.4.4 and 4.4.5).
- (4) It supplies a new representation formula for the exterior solution of the scattering problem, which is not based on a Green's function (Theorem 4.4.1).
- (5) The new representation formula supplies directly the far-field.
- (6) It generalizes parts of the important theorems of Wilcox and Karp concerning the series representation of solutions in the exterior domain (Theorem 4.4.4).
- (7) The pole condition supplies a constructive tool to design numerical algorithms. Three different forms of the pole condition are given in Section 4.4.2.

In a joint work with Hohage and Zschiedrich [58], we present proofs of the two missing points:

- (1) The pole condition is equivalent to Sommerfeld's radiation condition even in the infinite dimensional case.
- (2) The boundary condition derived from the pole condition ensures existence and uniqueness of the exterior scattering problem.

To contain the following sections self-consistent, we start again from the Helmholtz equation for a model problem—simple enough to derive the desired properties—and introduce the corresponding variational formulation, regardless of the presentation given in the previous sections.

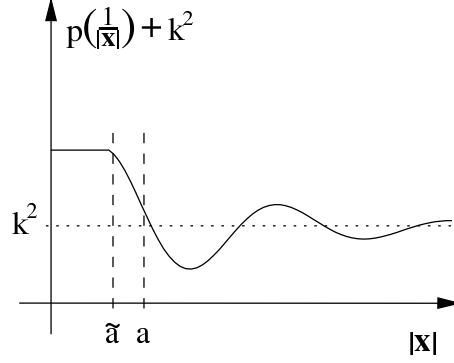


FIGURE 4.2.1. Potential of the model problem

**Model Problem.** We consider the Helmholtz problem (4.0.6) inside and outside the artificial domain  $\Omega_a = \{\mathbf{x} \in \mathbb{R}^d : |\mathbf{x}| < a\}^1$ ,  $d \geq 1$ , with boundary  $\Gamma_a = \{\mathbf{x} \in \mathbb{R}^d : |\mathbf{x}| = a\}$  and a radial symmetric potential. Further we set  $\Omega_{\text{int}} := \Omega_a$  and  $\Omega_{\text{ext}} = \mathbb{R}^d \setminus \overline{\Omega_a}$ . However, we underline that we can easily accommodate the setting to the more general situations described in the example section, e.g. to interior domains contained in  $\Omega_a$  and to problems with interior boundaries and sources and arbitrary piecewise continuous potentials. We investigate the model problem

$$(4.2.1) \quad \Delta u + p\left(\frac{1}{|\mathbf{x}|}\right)u + k^2u = 0 \quad \text{in } \mathbb{R}^d .$$

For the potential  $p : [0, \frac{1}{a}) \rightarrow \mathbb{R}$  we assume a power series representation

$$p(t) = \sum_{j=2}^{\infty} p_j t^j \quad \text{with convergence radius } \frac{1}{\tilde{a}} > \frac{1}{a} > 0,$$

where  $a$  is the above introduced radius of the artificial domain  $\Omega_a$ , see Fig. 4.2.1. The auxiliary radius  $\tilde{a} < a$  is needed later to obtain convenient bounds with respect to Laplace transformed functions. Additionally we assume that  $u$  exterior to  $\Omega_a$  can be composed from incoming and outgoing functions, that is  $u(\mathbf{x}) = u_{\text{in}}(\mathbf{x}) + u_{\text{out}}(\mathbf{x})$  for  $\mathbf{x} \in \Omega_{\text{ext}}$ . Further we assume that both the incoming and the outgoing functions are solutions of the exterior Helmholtz equation. We consider the incoming field  $u_{\text{in}}(\mathbf{x})$ ,  $\mathbf{x} \in \Omega_{\text{ext}}$ , as given and wish to compute the field  $u(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{R}^d$ , such that the field  $u - u_{\text{in}}$  is an outgoing field in the exterior domain.

Let us consider – as part of the full problem – the Helmholtz equation (4.2.1) for  $u_{\text{out}}$  on the exterior domain, where e.g. Neumann data on  $\partial\Omega_a$  are imposed. This problem is *not* uniquely solvable. This can already be seen at the simplest situation  $d = 1$  and the uniqueness of linear ODE's of second order. To obtain uniqueness of the exterior boundary value problem we need a second condition. As mentioned several times before, the established standard is to require the Sommerfeld radiation condition

$$(4.2.2) \quad \lim_{r \rightarrow \infty} r^{\frac{d-1}{2}} (\partial_r u_{\text{out}} - iku_{\text{out}}) = 0, \quad r = |\mathbf{x}|$$

<sup>1</sup>Different to the previous sections we write in the following  $a$  instead of  $r_0$  or  $\xi_0$  since we will use it several times as index of different quantities.

uniformly for all directions  $\frac{\mathbf{x}}{|\mathbf{x}|}$ . It can be shown (cf. [21]) that for a non-vanishing solution  $u_{\text{out}}$  to (4.2.1) and (4.2.2) the energy flux is positive

$$\text{Im} \int_{\partial\Omega_{\text{int}}} \overline{\partial_n u_{\text{out}}} u_{\text{out}} > 0.$$

This implies the uniqueness of the exterior solution, cf. [60]. In the following section we will repeat, in brief, how Sommerfeld's radiation condition, together with the positivity of the energy flux and properly defined Dirichlet-to-Neumann maps are used to prove existence and uniqueness of the interior scattering problem. This will direct our attention to the essential properties which we must derive from the pole condition.

**Dirichlet-to-Neumann Map.** We will study the interior problem by means of a variational formulation. To this end we multiply (4.2.1) by a test function  $-\bar{v}$  and integrate over the interior of the domain  $\Omega_a$ . A formal application of Green's theorem yields

$$(4.2.3) \quad \int_{\Omega_a} dx \left( \nabla \bar{v} \nabla u - \left( p \left( \frac{1}{|\mathbf{x}|} \right) + k^2 \right) \bar{v} u \right) - \int_{\Gamma_a} ds \bar{v} \partial_n u = 0$$

where  $\partial_n$  is the unit normal derivative with a normal vector  $\mathbf{n}$  pointing to the exterior of  $\Omega_a$ . Next we introduce the so-called Dirichlet-to-Neumann map  $B : H^{1/2}(\Gamma_a) \rightarrow H^{-1/2}(\Gamma_a)$  which maps Dirichlet data  $u_{\text{out}}|_{\Gamma_a}$  of the outgoing solution satisfying (4.2.1) and (4.2.2) to its Neumann data  $\partial_n u_{\text{out}}|_{\Gamma_a}$ . Later we will consider solutions in the exterior of  $\Omega_a$  based on the pole condition. For the moment, we concentrate on the solution of the problem interior to  $\Omega_a$ . According to (4.2.3), we introduce the sesquilinear form

$$(4.2.4) \quad \begin{aligned} a & : H^1(\Omega_a) \times H^1(\Omega_a) \rightarrow \mathbb{C} \\ a(v, u) & := \int_{\Omega_a} dx \left( \nabla \bar{v} \nabla u - \left( p \left( \frac{1}{|\mathbf{x}|} \right) + k^2 \right) \bar{v} u \right) - \int_{\Gamma_a} ds \bar{v} B u \end{aligned}$$

and the anti-linear functional

$$(4.2.5) \quad \begin{aligned} F & : H^1(\Omega_a) \rightarrow \mathbb{C} \\ F(v) & := \int_{\partial\Omega_a} ds \bar{v} (\partial_n u_{\text{src}} - B u_{\text{src}}), \end{aligned}$$

with given source data  $\partial_n u_{\text{src}} \in H^{-1/2}(\Gamma_a)$  and  $u_{\text{src}} \in H^{1/2}(\Gamma_a)$  so that the functional (4.2.5) is continuous. Now, the variational problem reads: Find  $u \in H^1(\Omega_a)$  such that

$$(4.2.6) \quad a(v, u) = F(v) \quad \text{for all } v \in H^1(\Omega_a).$$

The following theorem states the general existence and uniqueness result for the variational problem (4.2.6). Our further elaboration will not concentrate on the existence and uniqueness aspects, but it is important to know which properties the boundary operator has to satisfy.

**THEOREM 4.2.1.** *Let  $B$  be an operator with the following properties:*

- (i)  $B : H^{1/2}(\Gamma_a) \rightarrow H^{-1/2}(\Gamma_a)$  is linear and bounded.
- (ii) There exists a compact operator  $\tilde{B} : H^{1/2}(\Gamma_a) \rightarrow H^{-1/2}(\Gamma_a)$  such that
 
$$\text{Re} \int_{\Gamma_a} ds \left( -B + \tilde{B} \right) \bar{\phi} \phi \geq 0 \text{ for all } \phi \in H^{1/2}(\Gamma_a).$$

(iii)  $\text{Im} \int_{\Gamma_a} ds B \bar{\phi} \phi > 0$  for all  $\phi \in H^{1/2}(\Gamma_a)$ .

Then, the variational problem given by (4.2.4), (4.2.5) and (4.2.6) has a unique solution for  $u$  for all right-hand sides  $F$ , and depends continuously on  $F$ .

PROOF. Following [21, Theorem 5.7], we repeat the standard arguments in brief and leave out many details.

**Condition (i).** Since  $v \in H^1(\Omega_a)$ , the restriction to the boundary  $\Gamma_a$  results in  $v|_{\Gamma_a} \in H^{1/2}(\Gamma_a)$ . Let  $w := Bu$ . Hence  $w \in H^{-1/2}(\Gamma_a)$  and the integral  $\int_{\Gamma_a} ds \bar{v} w$  exists. Thus the condition guarantees that the sesquilinear form (4.2.4) is well defined.

**Condition (ii).** We use the operator  $\tilde{B}$  to decompose the sesquilinear form (4.2.4) into two parts  $a(v, u) = a_0(v, u) + a_1(v, u)$ :

$$\begin{aligned} a_{0,1} &: H^1(\Omega_a) \times H^1(\Omega_a) \rightarrow \mathbb{C} \\ a_0(v, u) &:= \int_{\Omega_a} dx \nabla \bar{v} \nabla u + \int_{\Gamma_a} ds \bar{v} (\tilde{B} - B) u + \int_{\Omega_a} dx \bar{v} u \\ a_1(v, u) &:= - \int_{\Omega_a} dx \left( p \left( \frac{1}{|\mathbf{x}|} \right) + k^2 \right) \bar{v} u - \int_{\Gamma_a} ds \bar{v} \tilde{B} u - \int_{\Omega_a} dx \bar{v} u \end{aligned}$$

Clearly,  $\text{Re} a_0(u, u) \geq c_1 \|u\|_{H^1(\Omega_a)}^2$  holds true for some constant  $c_1 > 0$ . Further we see

$$\begin{aligned} \int_{\Omega_a} dx \left( p \left( \frac{1}{|\mathbf{x}|} \right) + k^2 \right) \bar{u} u &\leq \sup_{\mathbf{x} \in \Omega_a} \left| p \left( \frac{1}{|\mathbf{x}|} \right) + k^2 \right| \int_{\Omega_a} dx \bar{u} u \\ \text{and} \quad \int_{\Gamma_a} ds \bar{u} \tilde{B} u &= \left\langle \text{Tr} u, \tilde{B} \text{Tr} u \right\rangle_{L^2(\Gamma_a)}. \end{aligned}$$

Together this yields the Gårding inequality

$$\text{Re} a(u, u) \geq c_1 \|u\|_{H^1(\Omega_a)}^2 - c_2 \|u\|_{L^2(\Omega_a)}^2 - \text{Re} \left\langle \text{Tr} u, \tilde{B} \text{Tr} u \right\rangle_{L^2(\Gamma_a)}$$

for all  $u \in H^1(\Omega_a)$  with positive constants  $c_1, c_2$  and  $\text{Tr} : H^1(\Omega_a) \rightarrow H^{1/2}(\Gamma_a)$  the trace operator. By Riesz' representation theorem, we can associate with the sesquilinear form  $a(\cdot, \cdot)$  an operator  $A : H^1(\Omega_a) \rightarrow H^1(\Omega_a)$ . This operator can be split into the sum  $A_0 + A_1$ , where  $A_0$  corresponds to  $a_0(\cdot, \cdot)$  and  $A_1$  to  $a_1(\cdot, \cdot)$ . The operator  $A_0$  is linear, continuous and boundedly invertible on  $H^1(\Omega_a)$ . The latter fact follows from the Lax-Milgram theorem. The operator  $A_1$  consists of two terms. The term corresponding to  $\left\langle \text{Tr} u, \tilde{B} \text{Tr} u \right\rangle_{L^2(\Gamma_a)}$  is compact, since  $\text{Tr}' \tilde{B} \text{Tr} : H^1(\Omega_a) \rightarrow$

$H^1(\Omega_a)$  is compact. Further, the term corresponding to  $\int_{\Omega_a} dx \left( p \left( \frac{1}{|\mathbf{x}|} \right) + k^2 \right) \bar{v} u$  is compact, since the embedding operator  $H^1(\Omega_a) \hookrightarrow L^2(\Omega_a)$  is compact. Thus the operator  $A_0 + A_1$  is of the form boundedly invertible plus compact operator, hence it is Fredholm of index zero.

**Condition (iii).** Let us consider the solution of the homogeneous operator equation  $(A_0 + A_1)u = 0$ , or equivalently, of the variational equation  $a(v, u) = 0$  for all  $v \in H^1(\Omega_a)$ . Taking the imaginary part of (4.2.6) and using condition (iii) shows that  $u$  has vanishing Cauchy data on  $\Gamma_a$ . Hence, by virtue of the Cauchy-Kowalewskaja Theorem and elliptic regularity results,  $u$  must vanish everywhere. Since the homogeneous equation has only the trivial solution, the original equation has a unique solution. Moreover, the solution depends continuously on the right-hand side, that is the inverse operator  $(A_0 + A_1)^{-1} : H^1(\Omega_a) \rightarrow H^1(\Omega_a)$  is continuous.  $\square$

Usually, the conditions (i)-(iii) are proved based on properties of the Hankel functions, see e.g. [21, 60]. This, however, works only for cases where a solution



representation based on Hankel functions is available. Our alternative approach presented below supplies these properties completely without using special functions. Moreover, it covers also cases with  $p \neq 0$ .

**Laplace Transformed Helmholtz Equation.** In order to simplify our analysis, we scale and shift the solution to obtain a new function

$$(4.2.7) \quad U(r, \mathbf{x}^0) := (r+a)^{\frac{d-1}{2}} u((r+a)\mathbf{x}^0)$$

with  $r \geq 0$  and  $\mathbf{x}^0 \in S^{d-1}$ . This scaling realizes  $\|U(r, \cdot)\|_{L^2(S^{d-1})} = \|u\|_{L^2((r+a)S^{d-1})}$ . It is well-known for the case of a vanishing potential  $p$ , (see again e. g. [21]), that the Sommerfeld condition (4.2.2) implies the asymptotic formula

$$U(r, \mathbf{x}^0) = U_\infty(\mathbf{x}^0) \left( e^{ikr} + \mathcal{O}\left(\frac{1}{r}\right) \right)$$

We will show that the same asymptotic formula holds true for our case of a non-vanishing potential  $p\left(\frac{1}{r}\right) \rightarrow 0$  as  $r \rightarrow \infty$ . Moreover, it is one of the most interesting aspects of our approach that it yields the far-field  $U_\infty$  directly both in theory and in numerical algorithms.

Let  $\{(\phi_j, \lambda_j) : j \in \mathbb{N}\}$  be a complete orthonormal system of eigenfunctions and eigenvalues of the Laplace-Beltrami operator on  $S^{d-1}$ . In two dimensions the eigenfunctions are trigonometric functions, in three dimensions spherical harmonics. We denote by  $U_j(r) := \int_{S^{d-1}} ds \bar{\phi}_j U(r, \cdot)$  the Fourier coefficients with respect to the given orthonormal system. We want to study first the case that the solution  $U$  possesses the factorized representation consisting only of *one* Fourier mode

$$U(r, \mathbf{x}^0) = U_j(r) \phi_j(\mathbf{x}^0), \quad \mathbf{x}^0 \in S^{d-1}.$$

We decompose the Laplacian in  $\mathbb{R}^d$  into a radial and a corresponding angular part using the representation formula (4.1.13)

$$\Delta = \frac{\partial^2}{\partial r^2} + \frac{d-1}{r+a} \frac{\partial}{\partial r} + \frac{1}{(r+a)^2} \Delta_{S^{d-1}}.$$

Based on this decomposition, we reformulate the Helmholtz equation (4.2.1) in terms of the Fourier coefficients:

$$(4.2.8) \quad U_j''(r) + \left[ k^2 + \frac{\frac{1}{4}(d-1)(3-d) + \lambda_j}{(r+a)^2} + p\left(\frac{1}{r+a}\right) \right] U_j(r) = 0.$$

Compared with Bessel's equation, there is no first order derivative, due to our scaling with the factor  $(r+a)^{\frac{d-1}{2}}$ .

Now we are ready to apply the Laplace transform (4.2.8). A direct application of the transformation formulas proved in the appendix supplies the key equations

$$(4.2.9) \quad (k^2 + s^2) \hat{U}_j + \int_s^\infty ds_1 P_j(s_1 - s) \hat{U}_j(s_1) = s U_j(0) + U_j'(0)$$

with

$$(4.2.10) \quad P_j(s) := p_*(s) + e^{-as} s \left( \frac{1}{4}(d-1)(3-d) + \lambda_j \right)$$

$$(4.2.11) \quad p_*(s) := e^{-as} \sum_{m=2}^{\infty} \frac{p_m}{(m-1)!} s^{m-1}.$$

The integral equation (4.2.9) is a linear Volterra integral equation of the second kind. It provides the basis of all our following considerations. Note, that the restriction of (4.2.11) to terms with index larger than two is of technical nature and not a restriction of the applicability of the method. In the following we will frequently use the quantity  $P_j(s)$  without referring to a given mode. To this end we introduce

$$(4.2.12) \quad P(s) := p_*(s) + ce^{-as}, \quad P(0) = 0$$

with a real constant  $c$ .

### 4.3. Main Theoretical Results

The integral equation (4.2.9) is the central object of our analysis. We shall derive from it all desired results. To simplify our investigation we set  $k = 1$ , which can always be obtained by a rescaling  $r \mapsto kr$  and define the new function

$$(4.3.1) \quad w(s) := \widehat{U}(s)(s^2 + 1)$$

where  $\widehat{U}$  stands for one of the modes  $\widehat{U}_j$ . With these notations (4.2.9) can be written as

$$(4.3.2) \quad w(s) + (Jw)(s) = sU(0) + U'(0)$$

with the integral operator

$$(4.3.3) \quad (Jw)(s) = \int_s^\infty ds_1 P(s - s_1) \frac{w(s_1)}{s_1^2 + 1}$$

and  $P$  is of the type (4.2.12). It is our first goal, to establish the solvability of the integral equation (4.3.2). To this end we introduce two cuts in the complex plane and define a metric appropriate for our purposes. Let the cuts be  $S_{\pm i} := \{\pm i - t : t > 0\}$ , see Fig. 4.3.1 and the domain  $V := \mathbb{C} \setminus (S_+ \cup S_-)$ . We want to analyze functions on  $V$  which are not continuous across the cuts. To deal with functions whose difference quotient, or Hölder difference quotient remains bounded for points with a small distance from each other, we define the following metric

$$d(s_1, s_2) = \sqrt{|s_1 - s_2|^2 + |\phi(s_1) - \phi(s_2)|^2}$$

with an auxiliary function  $\phi : V \rightarrow \mathbb{R}$

$$\phi(s) := \begin{cases} -\operatorname{Re} s & \text{for } |\operatorname{Im} s| < 1, \operatorname{Re} s < 0 \\ 0 & \text{else.} \end{cases}$$

Fig. 4.3.1 explains the construction of the metric. If two points lie outside the strip between the cuts, their distance is the usual one. If one lies inside and one outside the area, there distance is the usual one plus the “distance” needed to surround the cut.

Let  $\overline{V}$  denote the completion of  $V$  with respect to  $d$ . Then  $\overline{V}$  is the union of all points of  $V$  with the set of points

$$s_{\pm} := \lim_{\epsilon \downarrow 0} s \pm i\epsilon, \text{ where } s \in S_{-i} \cup S_{+i}.$$

The curves  $s_+$  and  $s_-$  are the “upper” and the “lower” side of the cuts. Now we can define a “jump” function  $[w]$  which corresponds to a given continuous (with respect to  $d$ ) function  $w : \overline{V} \rightarrow \mathbb{C}$  by

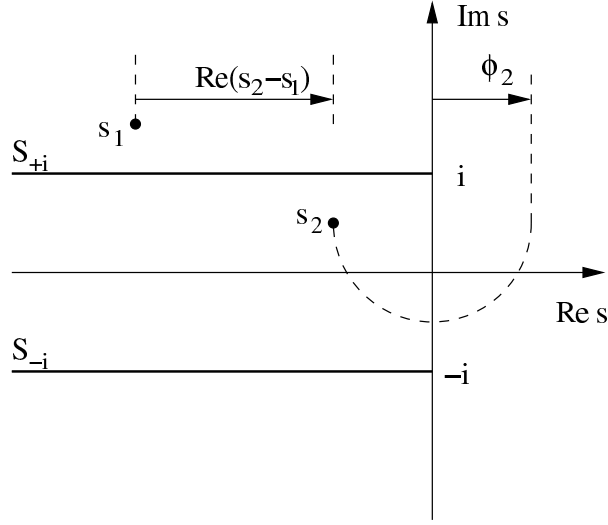


FIGURE 4.3.1. Definition of the cuts and construction of the metric  $d$

$$(4.3.4) \quad \begin{aligned} [w] & : S_{-i} \cup S_{+i} \rightarrow \mathbb{C} \\ [w](s) & := w(s_+) - w(s_-). \end{aligned}$$

If  $w$  is continuous on  $\mathbb{C}$  (with respect to the usual distance in  $\mathbb{C}$ ) so is  $[w]$  on  $S_{-i} \cup S_{+i}$  and  $\lim_{s \rightarrow \pm i} [w] = 0$ . We introduce the Banach space  $X$  with norm  $\|\cdot\|_X$  by

$$(4.3.5) \quad X = \left\{ w \in C(\overline{V}) : w \text{ holomorphic,} \right. \\ \left. w = o(|s|^2) \text{ uniformly for } |s| \rightarrow \infty \right\}$$

$$(4.3.6) \quad \|w\|_X = \sup_{s \in \overline{V}} \frac{|w(s)|}{1 + |s|^2}.$$

Further, in order to analyze the solvability of the integral equation (4.3.3) we need to isolate the singularities at  $s = \pm i$ . To this end we denote by  $|s|_1 := |\operatorname{Re} s| + |\operatorname{Im} s|$  and introduce the diamond shaped regions

$$D_{\pm} = \left\{ s \in \overline{V} : |s \pm i|_1 < \frac{1}{2} \right\}.$$

Now we can formulate the essential regularity results of the integral operator (4.3.3):

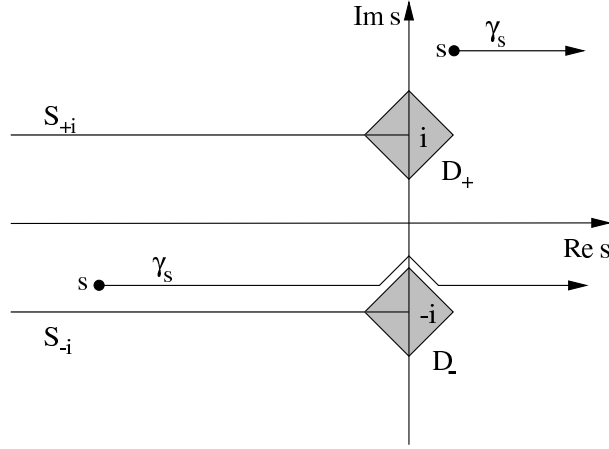
LEMMA 4.3.1. *Let  $0 < \alpha < 1$ . Then there exists a constant  $c$  such that for all  $w \in X$*

$$(4.3.7) \quad |(Jw)(s)| \leq c \left( \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s} \frac{|w(s_1)|}{1 + |s_1|^2} \right), \quad s \in \overline{V} \setminus (D_+ \cup D_-)$$

$$(4.3.8) \quad |(Jw)'(s)| \leq c \left( \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s} \frac{|w(s_1)|}{1 + |s_1|^2} \right), \quad s \in \overline{V} \setminus (D_+ \cup D_-)$$

$$(4.3.9) \quad |(Jw)(s) - (Jw)(\sigma)| \leq c d(s, \sigma)^\alpha \|w\|_X, \quad s, \sigma \in (D_+ \cup D_-)$$

$$(4.3.10) \quad |(Jw)(s)| \leq c \|w\|_X, \quad s \in \overline{V}.$$

FIGURE 4.3.2. Paths  $\gamma_s$  outside the domains  $D_{\pm}$ 

PROOF. We prove first the bound (4.3.7). The operator is an integral over a holomorphic function, hence we may deform the contour without changing the result. First we consider the rather simple case that  $\operatorname{Re} s > 1/2$ , that is  $s$  is on the right-hand side of  $D_{\pm}$ . With  $s_1 = t + s$ ,  $t \in \mathbb{R}_+$ , we get from (4.3.3)

$$\begin{aligned} |(Jw)(s)| &= \left| \int_0^{\infty} dt P(t) \frac{w(t+s)}{(t+s)^2 + 1} \right| \\ &\leq \left| \int_0^{\infty} dt P(t) \right| \left( \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s} \frac{|w(s_1)|}{1 + |s_1|^2} \right). \end{aligned}$$

Now we have from (4.2.12) and (4.2.11) and the convolution Lemma A.1.8

$$|P(t)| \leq c_1 e^{-at + (a-\tilde{a})t} + c_2 t e^{-at}.$$

Clearly, there exists a positive constant  $c$  such that

$$|P(t)| \leq c e^{-\tilde{a}t}, \quad t \geq 0.$$

Thus  $|\int_0^{\infty} dt P(t)| \leq c/\tilde{a}$  and we have proved (4.3.7) for this particular choice of the starting point  $s$ . Now we extend this to the general situation  $s \in \overline{V} \setminus (D_+ \cup D_-)$ . For  $\operatorname{Im} s > 0$  we use the contour parameterized by

$$\begin{aligned} \gamma_s(t) &:= s + t + i\psi_s(t) \\ \psi_s(t) &:= \begin{cases} \frac{1}{2} - |s + t - i|_1 & \text{if } |s + t - i|_1 \leq \frac{1}{2} \quad \text{and} \quad \operatorname{Im}(s + t - i) > 0 \\ -\frac{1}{2} + |s + t - i|_1 & \text{if } |s + t - i|_1 \leq \frac{1}{2} \quad \text{and} \quad \operatorname{Im}(s + t - i) \leq 0 \\ 0 & \text{else,} \end{cases} \end{aligned}$$

see Fig. (4.3.2). For  $\operatorname{Im} s \leq 0$  the contour follows accordingly replacing  $+i$  by  $-i$ .

This path has the following properties:

- (1) It neither intersect the cuts  $S_{\pm i}$  nor the domains  $D_{\pm}$ ,
- (2)  $\operatorname{meas}(\operatorname{supp} \psi_s) \leq 1$ ,
- (3)  $|\psi'_s| \leq 1$ ,
- (4)  $|\psi_s| \leq 1/2$ ,
- (5)  $\lim_{t \rightarrow \infty} \Psi_s(t) = 0$ , and  $\Psi_s(0) = 0$ .

Using this path, which is an extension of the foregoing one, we write now

$$\begin{aligned} |(Jw)(s)| &= \left| \int_0^\infty dt \frac{d\gamma_s}{dt} P(t + i\psi_s(t)) \frac{w(t + i\psi_s(t) + s)}{(t + i\psi_s(t) + s)^2 + 1} \right| \\ &= \left| \int_0^\infty dt \frac{d\gamma_s}{dt} P(t + i\psi_s(t)) \frac{|t + i\psi_s(t) + s|^2 + 1}{(t + i\psi_s(t) + s)^2 + 1} \frac{w(t + i\psi_s(t) + s)}{|t + i\psi_s(t) + s|^2 + 1} \right|. \end{aligned}$$

We want to analyze the influence of  $i\psi_s(t)$  in the argument of  $P(t + i\psi_s(t))$ . Our goal is to show that the deformation of the contour caused by  $\Psi_s$  is small enough to maintain an exponential bound with respect to  $|P(s)|$ . By the convolution lemma A.1.8 of the appendix we know  $|P(\tau + i\sigma)| \leq ce^{-a\tau + (a-\tilde{a})|\tau + i\sigma|}$ ,  $\sigma, \tau \in \mathbb{R}$ . Let  $\tau$  be the solution of

$$(4.3.11) \quad -a\tau + (a - \tilde{a}) \left| \tau + \frac{1}{2}i \right| = -\frac{\tilde{a}}{2}\tau,$$

that is

$$\tau = \frac{|a - \tilde{a}|}{2\sqrt{(a - \frac{\tilde{a}}{2})^2 - (a - \tilde{a})^2}}.$$

Since  $|\Psi_s| \leq 1/2$ , equation (4.3.11) extends to the inequality

$$-at + (a - \tilde{a})|t + i\Psi_s(t)| \leq -\frac{\tilde{a}}{2}t \quad \text{for all } t \geq \tau.$$

Thus we obtain

$$|P(t + i\Psi_s(t))| \leq ce^{-at + (a-\tilde{a})|t + i\Psi_s(t)|} \leq ce^{-\tilde{a}/2t} \quad \text{for } t \geq \tau.$$

Since  $P$  is continuous, this holds also true for  $0 \leq t \leq \tau$ , possibly with a larger constant  $c$ :

$$|P(t + i\Psi_s(t))| \leq ce^{-\tilde{a}/2t} \quad \text{for } t \geq 0.$$

Further, the quotients

$$\begin{aligned} \sup_{s \in V \setminus (D_+ \cup D_-)} \sup_{t \geq 0} \left| \frac{|t + i\psi_s(t) + s|^2 + 1}{(t + i\psi_s(t) + s)^2 + 1} \right| &< \infty \\ \sup_{s \in V \setminus (D_+ \cup D_-)} \sup_{t \geq 0} \left| \frac{(t + s)^2 + 1}{|t + i\psi_s(t) + s|^2 + 1} \right| &< \infty \end{aligned}$$

are bounded. Hence

$$\begin{aligned} |(Jw)(s)| &\leq c \left( \int_0^\infty dt e^{-\tilde{a}/2t} \right) \left( \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s} \frac{|w(s_1)|}{1 + |s_1|^2} \right) \\ &= \frac{2c}{\tilde{a}} \left( \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s} \frac{|w(s_1)|}{1 + |s_1|^2} \right). \end{aligned}$$

This is (4.3.7).

Next we want to prove (4.3.8). First we observe that

$$(Jw)'(s) = \int_s^\infty ds_1 P'(s_1 - s) \frac{w(s_1)}{s_1^2 + 1}$$

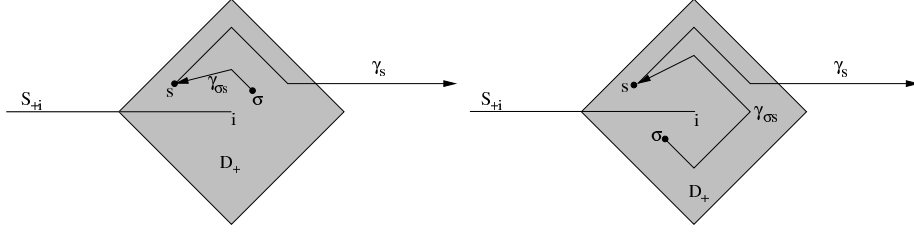


FIGURE 4.3.3. The two different choices of the paths inside the domains  $D_{\pm}$

since by assumption  $P(0) = 0$ . Further,  $P'$  possesses the same exponential bound as  $P$ , by virtue of the convolution lemma A.1.8. Hence we can repeat the above procedure to obtain (4.3.8).

Now to the harder part. We want to show the Hölder continuity (4.3.9). Without loss of generality we may assume that  $s, \sigma \in D_+$  and that  $|s - i|_1 \geq |\sigma - i|_1$ . We say that  $s$  and  $\sigma$  are on *opposite* sides (with respect to  $S_i$ ) if both  $\operatorname{Re} s < 0$ ,  $\operatorname{Re} \sigma < 0$  and  $(\operatorname{Im}(s - i))(\operatorname{Im}(\sigma - i)) < 0$  or if  $s$  and  $\sigma$  are limits of such points. First we assume that  $s$  and  $\sigma$  are not on opposite sides of  $S_i$ . We consider a path  $\gamma_{\sigma,s}$  connecting the points  $\sigma$  and  $s$ . Let  $\gamma_{\sigma,s}$  be the shortest path such that

$$|s_1 - i|_1 \geq |\sigma - i|_1, \quad \text{for all } s_1 \in \gamma_{\sigma,s}.$$

Fig. 4.3.3 illustrates this situation.

The length of the path can be estimated by  $l(\gamma_{\sigma,s}) \leq 3\delta$ , where  $\delta$  is the distance between the points  $s$  and  $\sigma$ ,  $\delta = |s - \sigma|$ . We extend the path  $\gamma_{\sigma,s}$  from  $s$  to  $+\infty$  as shown in Fig. 4.3.3 by a path  $\gamma_s$  such that

$$|s_1 - i|_1 \geq |s - i|_1, \quad \text{for all } s_1 \in \gamma_s.$$

Next, we compute the difference

$$\begin{aligned} (Jw)(s) - (Jw)(\sigma) &= \int_s^\infty ds_1 P(s_1 - s) \frac{w(s_1)}{s_1^2 + 1} - \int_\sigma^\infty ds_1 P(s_1 - \sigma) \frac{w(s_1)}{s_1^2 + 1} \\ &= - \int_\sigma^s ds_1 P(s_1 - \sigma) \frac{w(s_1)}{s_1^2 + 1} \\ &\quad + \int_s^\infty ds_1 (P(s_1 - s) - P(s_1 - \sigma)) \frac{w(s_1)}{s_1^2 + 1} \\ &= \int_{\gamma_{\sigma,s}} ds_1 P(s_1 - \sigma) \frac{w(s_1)}{s_1^2 + 1} \\ &\quad + \int_{\gamma_s} ds_1 (P(s_1 - s) - P(s_1 - \sigma)) \frac{w(s_1)}{s_1^2 + 1}. \end{aligned}$$

Using  $|s| \leq |s|_1 \leq \sqrt{2}|s|$  we obtain for  $s_1 \in \gamma_{\sigma,s}$

$$\begin{aligned} \left| \frac{P(s_1 - \sigma)}{s_1^2 + 1} \right| &= \left| \frac{1}{s_1 + i} \right| \left| \frac{P(s_1 - \sigma)}{s_1 - i} \right| \\ &\leq \left| \frac{\sqrt{2}}{s_1 + i} \right| \frac{|P(s_1 - \sigma)|}{|s_1 - i|_1}. \end{aligned}$$

We show that the second factor in the last inequality is bounded. First,  $|s_1 - \sigma|$  goes faster to zero than  $|s_1 - i|$  as  $s_1 \rightarrow \sigma$ , since by construction  $|s_1 - i|_1 \geq |s_1 - \sigma|_1$ .

Second,  $|P(t)| = \mathcal{O}(|t|)$  as  $|t| \rightarrow 0$ , due to the properties (4.2.10) and (4.2.11). Hence we find

$$\left| \frac{P(s_1 - \sigma)}{s_1^2 + 1} \right| \leq c, \quad \text{for all } s_1 \in \gamma_{\sigma, s}.$$

Together with the bound on  $l(\gamma_{\sigma s})$  this yields for the first integral

$$\left| \int_{\gamma_{\sigma s}} ds_1 P(s_1 - \sigma) \frac{w(s_1)}{s_1^2 + 1} \right| \leq c\delta \|w\|_X.$$

We estimate the second integral. We apply the mean value theorem

$$\begin{aligned} & \left| \int_{\gamma_s} ds_1 (P(s_1 - s) - P(s_1 - \sigma)) \frac{w(s_1)}{s_1^2 + 1} \right| \\ & \leq \delta \int_{\gamma_s} |ds_1| \sup_{z \in [s, \sigma]} |P'(s_1 - z)| \left| \frac{w(s_1)}{s_1^2 + 1} \right|. \end{aligned}$$

The derivative of  $P(\cdot)$  is uncritical. We consider the term  $1/(s_1^2 + 1)$ . To bound the integrand inside of  $D_+$  we take into account that due to the choice of  $\gamma_s$  we have

$$\begin{aligned} |s_1 - i|_1 & \geq |s_1 - s|_1 - |s - i|_1 && \text{inverse triangle inequality} \\ & \geq |s_1 - s|_1 - |s_1 - i|_1 \\ \text{and } 2|s_1 - i|_1 & \geq |s - i|_1 + |\sigma - i|_1 \\ & \geq |s - \sigma|_1 \\ & \geq \delta. \end{aligned}$$

The sum of these inequalities supplies

$$4|s_1 - i|_1 \geq |s_1 - s|_1 + \delta.$$

This implies that if  $s_1 \rightarrow s$  there is a residual distance between  $s_1$  and  $i$  larger than  $\delta/4$ . Together with  $|s_1 - s|_1 \geq \operatorname{Re}(s_1 - s)$  and  $|s_1 + i|_1 \geq 1$  (we are in  $D_+$ ) we get

$$\begin{aligned} |s_1^2 + 1| & \geq \frac{1}{2} |s_1 + i|_1 |s_1 - i|_1 \\ & \geq \frac{1}{8} (\operatorname{Re}(s_1 - s) + \delta). \end{aligned}$$

This is the result which we need to integrate the critical term  $1/(s_1^2 + 1)$  inside  $D_+$ . Outside of  $D_+$  the bound

$$\left| \frac{w(s_1)}{s_1^2 + 1} \right| \leq c \|w\|_X$$

holds true. With the bound on  $P'(\cdot)$  from Lemma A.1.8 and  $t^* := \sup\{t \geq 0 : s + t \in D_+\}$  we obtain

$$\begin{aligned}
& \left| \int_{\gamma_s} ds_1 (P(s_1 - s) - P(s_1 - \sigma)) \frac{w(s_1)}{s_1^2 + 1} \right| \\
& \leq c\delta \|w\|_X \left( \int_0^{t^*} dt \frac{1}{t + \delta} + \int_{t^*}^{\infty} dt e^{-\tilde{a}t/2} \right) \\
& \leq c\delta \|w\|_X \left( \ln \frac{t^* + \delta}{\delta} + \frac{\tilde{a}}{2} \right) \\
& \leq c\delta^\alpha \|w\|_X.
\end{aligned}$$

Since  $\delta \leq d(s, \sigma)$ , we have proved (4.3.9) if the points  $s$  and  $\sigma$  not on opposite sides of  $S_i$ . The other case can be traced back to the above one using the triangle inequality and our last estimates

$$\begin{aligned}
|(Jw)(\sigma) - (Jw)(s)| & \leq |(Jw)(\sigma) - (Jw)(i)| + |(Jw)(i) - (Jw)(s)| \\
& \leq c(|\sigma - i|^\alpha + |i - s|^\alpha).
\end{aligned}$$

Without loss of generality we may assume that  $\text{Im}(\sigma - i) < 0$  and  $\text{Im}(s - i) > 0$ , i. e.  $d(\sigma, s) = |\text{Re}\sigma| + |s - \sigma|$ . Further, inserting the following inequalities in the last expression

$$\begin{aligned}
|\text{Im}(\sigma - i)| + |\text{Im}(i - s)| & = |\text{Im}(\sigma - s)| \\
& \leq d(\sigma, s), \\
|\text{Re}(s - i)| & = |\text{Re}s| \\
& \leq |\text{Re}\sigma| + |\text{Re}(s - \sigma)| \\
& \leq d(s, \sigma), \\
|\text{Re}(\sigma - i)| & = \text{Re}\sigma \\
& \leq |d(s, \sigma)|,
\end{aligned}$$

we obtain the desired result (4.3.9). The last equation (4.3.10) follows directly from (4.3.9) and (4.3.7).  $\square$

The following is the essential prerequisite to prove the unique solvability of our integral equation (4.3.2).

**LEMMA 4.3.2.** *The integral operator  $J$  defined in (4.3.3) is a compact operator from  $X$  to  $X$ .*

**PROOF.** It follows from the definition of  $J$  in (4.3.3) that  $Jw$  is holomorphic for  $w \in V$ . Further, by the above lemma, equation (4.3.7),  $Jw$  remains bounded for  $w \in \overline{V}$ . Together with the global bound (4.3.10) from Lemma 4.3.1 we see that  $J$  is an operator  $J : X \rightarrow X$ . To prove compactness we use an iterated application of the Arzelá-Ascoli theorem. Let  $(w_n)_{n \in \mathbb{N}}$  be a bounded sequence in  $X$ ,  $\|w_n\|_X \leq 1$  for all  $n \in \mathbb{N}$ . We have to show that the sequence  $v_n = Jw_n$  has a convergent subsequence in  $X$ . First we consider the restriction of  $v_n$  to some compact subset  $K \subset \overline{V}$ . By Theorem 4.3.1, equations (4.3.9) and (4.3.10), the set of restricted functions  $\{v_n|_K\}$  is equi-continuous on  $K$  with respect to the metric  $d$ . Hence, by the Arzelá-Ascoli theorem, there exists a subsequence of  $v_n$  which converges with respect to the maximum norm on  $K$ ,  $\|\phi\|_{\infty, K} = \sup_{s \in K} |\phi(s)|$ . Next we introduce a sequence of sets  $K_j$  which become larger and larger. We define  $K_j = \{s \in \overline{V} : |s| \leq j, j \in \mathbb{N}\}$ . Our induction argument is the following. By our conclusion above, there exists a subsequence  $(v_{n_1(l)})_l$  which converges with respect to  $\|\cdot\|_{\infty, K_1}$ . Next we consider



$K_2$ . Applying the same argument, we obtain a subsequence  $(v_{n_2(l)})_l$  which converges with respect to  $\|\cdot\|_{\infty, K_2}$ . Repeating this process of selecting subsequences, we arrive at an array  $v_{n_j(l)}$  with the property that each row is a subsequence of the forgoing row. The diagonal subsequence  $v_n(l) := v_{n_l(l)}$  converges to some function  $v$  with respect to the supremum norm on each  $K_j$ . In particular,  $\lim_{l \rightarrow \infty} v_n(l)(s) = v(s)$  for all  $s \in \bar{V}$ . It remains to show that  $\|v_n(l) - v\|_X \rightarrow 0$ . Let  $\epsilon > 0$ . By virtue of the global bound Lemma 4.3.1 and equation (4.3.10) there exists a constant  $C > 0$  such that  $|v_n(l)(s)| \leq C$  for all  $s \in \bar{V}$  and  $l \in \mathbb{N}$ . Therefore we find

$$\frac{|v(s) - v_n(l)(s)|}{1 + |s|^2} \leq \epsilon, \quad \text{for all } l \in \mathbb{N} \quad \text{and } |s| \geq \sqrt{\frac{2C}{\epsilon}}.$$

Now let  $J \geq \sqrt{2C/\epsilon}, J \in \mathbb{N}$ . Since  $v_n(l)$  converges to  $v$  with respect to  $\|\cdot\|_{\infty, K_J}$ , there exists  $L \in \mathbb{N}$  such that

$$\sup_{s \in K_J} \frac{|v(s) - v_n(l)(s)|}{1 + |s|^2} \leq \|v(s) - v_n(l)(s)\|_{\infty, K_J} \leq \epsilon$$

for  $l \geq L$ . Putting these inequalities together, we find

$$\|v(s) - v_n(l)(s)\|_X \leq \epsilon \quad \text{for } l \geq L.$$

□

Now we are ready to state the first essential result.

**THEOREM 4.3.3.** *The integral equation (4.3.2) has a unique solution in  $X$ .*

**PROOF.** Since  $J$  is a compact operator from  $X$  onto  $X$ , the assertion follows from Riesz theory if we can show the uniqueness of the solution. Therefore we consider the homogeneous equation

$$w_0 + Jw_0 = 0$$

and want to show that necessarily  $w_0 = 0$ .

First we deduce that  $w_0$  must be bounded. We find

$$\begin{aligned} |w_0(s)| &= |Jw_0(s)|, \quad s \in \bar{V} \\ &\leq \frac{c}{a} \|w_0\|_X \\ &= \frac{c}{a} \sup_{s \in \bar{V}} \frac{|w_0(s)|}{1 + |s|^2} \end{aligned}$$

by Lemma 4.3.1, equation (4.3.10), the norm definition (??) and the fact that  $w_0 = o(|s|^2)$ . It follows from Lemma 4.3.1, equation (4.3.7) that there exists a  $s^* \in \bar{V} \setminus (D_+ \cup D_-)$  which specializes

$$|Jw_0(s)| \leq c \left( \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s} \frac{|w_0(s_1)|}{1 + |s_1|^2} \right), \quad s \in \bar{V} \setminus (D_+ \cup D_-)$$

to

$$|Jw_0(s^*)| \leq \frac{1}{4} \sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s^*} |w_0(s_1)|, \quad s^* \in \bar{V} \setminus (D_+ \cup D_-).$$

To compute the supremum we consider only the case  $\operatorname{Re} s^* > 0$ . We make  $\operatorname{Re} s^*$  even larger, until the bound

$$\sup_{\operatorname{Re} s_1 \geq \operatorname{Re} s^*} |w(s_1)| \leq 2 |w(s^*)|$$

holds. This must be possible, otherwise  $w$  cannot be a bounded function. Thus we obtain

$$\begin{aligned} |w_0(s^*)| &= |Jw_0(s^*)| \\ &\leq \frac{1}{2} |w(s^*)|. \end{aligned}$$

It follows that  $w(s^*) = 0$ . This implies  $w(s) = 0$  for all  $s$  with  $\operatorname{Re} s \geq \operatorname{Re} s^*$ . Since  $w$  is holomorphic and continuous, it follows that  $w(s) = 0$  for all  $s \in \overline{V}$ .  $\square$

**Cut Functions.** Now we want to study the above mentioned cut functions, using the definition of the jump functions (4.3.4) at p. 67,

$$(4.3.12) \quad \psi_{\pm}(t) := \frac{[\widehat{U}](\pm i - t)}{w(\pm i)}, \quad t > 0.$$

The normalization with respect to the values  $w(\pm i)$  allows an intuitive formulation of the theory. For the moment we assume  $w(\pm i) \neq 0$ , later we will see how to deal with the case  $w(+i) = 0$  and the case  $w(-i) = 0$ . It is a remarkable result that we can compute the cut functions without knowing the values  $U(0)$  and  $U'(0)$ , which define the right-hand side of the underlying integral equation (4.3.2). The introduction of the cut functions is of basic interest both for the theory and for the numerical algorithms for the following reasons:

- (1) The cut functions show that we can always eliminate one singularity by a proper choice of the boundary conditions.
- (2) It can be shown that, once one of the singularities is removed in one complex half-plane, the  $\widehat{U}$  becomes a holomorphic function in the corresponding half plane.
- (3) The cut functions allow the construction of an *explicit representation formula* for the spatial function in the exterior domain.
- (4) They can be computed efficiently.
- (5) Their computation supplies *automatically* the far field, without any other numerical effort.

LEMMA 4.3.4. *If  $w(\pm i) \neq 0$ , the cut functions  $\psi_{\pm}$  defined by (4.3.12) for  $t > 0$  satisfy the integral equations*

$$(4.3.13) \quad \psi_+(t) + \int_0^t dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1) = \pi \frac{P(t)}{t(t-2i)},$$

$$(4.3.14) \quad \psi_-(t) + \int_0^t dt_1 \frac{P(t-t_1)}{t(t+2i)} \psi_-(t_1) = \pi \frac{P(t)}{t(t+2i)}.$$

Here  $P(\cdot)$  denotes the infinite series defined in (4.2.10), with dropped index  $j$ .

PROOF. We prove only the first of these equations since the proof of the second is similar. First we write our basic integral equation (4.3.2) with  $w$  replaced according to (4.3.1)

$$(s_{\pm}^2 + 1) \widehat{U}(s_{\pm}) + \int_{\gamma_{\pm}^{\epsilon}} ds_1 P(s_1 - s_{\pm}) \widehat{U}(s_1) = s_{\pm} U(0) + U'(0)$$

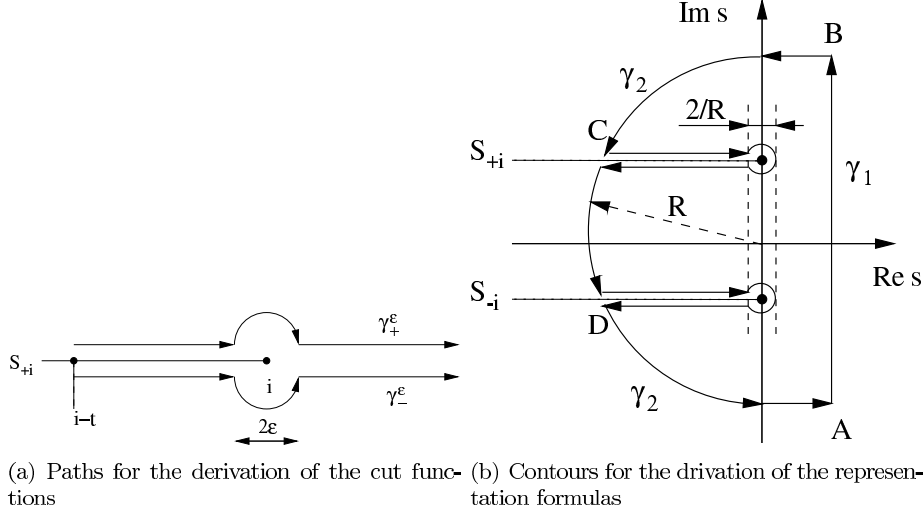


FIGURE 4.3.4.

for  $s = i - t \in S_i$  and  $\epsilon > 0$ . The paths  $\gamma_{\pm}^{\epsilon}$  are shown in Fig. 4.3.4. Subtracting the equation with the minus sign from the one with the plus sign yields

$$t(t-2i) \left[ \widehat{U} \right] (i-t) + \int_{-t}^{-\epsilon} dt P(t+t_1) \left[ \widehat{U} \right] (i+t_1) \\ + \int_0^{2\pi} d\phi (i\epsilon) e^{-i\phi} P(t - \epsilon e^{-i\phi}) \frac{w(i - \epsilon e^{-i\phi})}{(2i - \epsilon e^{-i\phi})(-\epsilon e^{-i\phi})} = 0.$$

Carrying out the limit  $\epsilon \rightarrow 0$  and dividing by  $w(i)t(t-2i)$  yields (4.3.13).  $\square$

Referring to the book of Kreß [66, Theorem 10.15] we know that the Volterra integral equations for the cut functions (4.3.13), (4.3.14) have a unique solution, since both the kernel and the right-hand sides are bounded. This is, of course what we expect, since our basic integral equation (4.3.2) is uniquely solvable. Next we will derive essential properties of the cut functions  $\psi_{\pm}$  near 0 and  $\infty$ . We will need these properties for the characterization of the near and the far-field of the corresponding spatial solution. Since  $\psi_{-}(t) = \overline{\psi_{+}(t)}$ , we give the first lemma only for  $\psi_{+}$ .

LEMMA 4.3.5. *The cut function  $\psi_{+}$  defined in (4.3.12) belongs to  $C^{\infty}(\mathbb{R}_{+})$ . The derivatives of  $\psi_{+}$  at 0 can be computed recursively as follows:*

$$(4.3.15) \quad \psi_{+}(0) = \pi \lim_{t \rightarrow 0} \frac{P(t)}{t(t-2i)}$$

$$(4.3.16) \quad \psi_{+}^{(k+1)}(0) = \pi \lim_{t \rightarrow 0} \left( \frac{P(t)}{t(t-2i)} \right)^{(k+1)} \\ + \frac{(k+1)!}{2i} \sum_{j=1}^{k+1} \frac{1}{(2i)^{k+1-j}(j+1)!} \sum_{n=1}^j P^{(n)}(0) \psi_{+}^{(j-n)}(0).$$

PROOF. Introducing the action of the integral operator by  $(K\psi_{+})(t)$ , we write the Volterra integral equation (4.3.13)

$$(4.3.17) \quad \begin{aligned} \psi_+(t) + (K\psi_+)(t) &= \pi \frac{P(t)}{t(t-2i)}, \quad t > 0 \\ (K\psi_+)(t) &:= \int_0^t dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1). \end{aligned}$$

By repeated partial integration we obtain

$$\begin{aligned} \int_0^t dt_1 P(t-t_1) \frac{t^j}{j!} &= \sum_{l=1}^{\infty} P^{(l)}(0) \int_0^t dt_1 \frac{(t-t_1)^l}{l!} \frac{t^j}{j!} \\ &= \sum_{l=1}^{\infty} P^{(l)}(0) \frac{t^{l+j+1}}{(l+j+1)!}. \end{aligned}$$

The change of the order of integration and summation is justified, since the Taylor expansion of  $P(\cdot)$  converges uniformly around 0. The last term is an analytic function in  $t$ . Next, suppose that a given function  $v(t)$  is a polynomial

$$v(t) := \sum_{j=0}^{\infty} \frac{v^{(j)}(0)}{j!} t^j.$$

Then, we compute

$$(4.3.18) \quad \begin{aligned} (Kv)(t) &= \frac{1}{t(t-2i)} \int_0^t dt_1 P(t-t_1) \sum_{j=0}^{\infty} \frac{v^{(j)}(0)}{j!} t^j \\ &= \frac{1}{t-2i} \sum_{l=1}^{\infty} \sum_{j=0}^{\infty} P^{(l)}(0) v^{(j)}(0) \frac{t^{l+j}}{(l+j+1)!} \\ &= -\frac{1}{2i} \sum_{r=0}^{\infty} \left(\frac{t}{2i}\right)^r \sum_{m=1}^{\infty} \frac{t^m}{(m+1)!} \sum_{n=1}^m P^{(n)}(0) v^{(m-n)}(0) \\ &= -\frac{1}{2i} \sum_{l=0}^{\infty} t^l \sum_{j=1}^l \left(\frac{1}{2i}\right)^{l-j} \frac{1}{(j+1)!} \sum_{n=1}^j P^{(n)}(0) v^{(j-n)}(0). \end{aligned}$$

In particular,  $(Kv)(t)$  is analytic at  $t = 0$ .

We prove by induction on  $n$  that  $\psi_+ \in C^n([0, \infty))$  for  $n \in \mathbb{N}_0$ . We start with  $n = 0$ . Note that the right-hand side of (??) can be continuously extended to a function in  $C^\infty([0, \infty))$ , due to the fact  $P(0) = 0$ . First we consider  $\psi_+$  on the interval  $[0, t^*]$  and show that there exists  $t^* > 0$  such that  $\|K\| < 1$ , where the integral operator  $K$  is restricted accordingly to  $K : C([0, t^*]) \rightarrow C([0, t^*])$ . The definition of  $P(\cdot)$  implies

$$P(t) = \mathcal{O}(t) \quad (t \rightarrow 0).$$

Based on the definition of  $K$  in (??) it follows

$$\begin{aligned}
|(K\psi_+)(t^*)| &= \left| \int_0^{t^*} dt_1 \frac{P(t^* - t_1)}{t^*(t^* - 2i)} \psi_+(t_1) \right| \\
&\leq C \left| \int_0^{t^*} dt_1 \frac{t^* - t_1}{t^*(t^* - 2i)} \psi_+(t_1) \right| \\
&\leq C \left( \max_{0 \leq t \leq t^*} |\psi_+(t)| \right) \int_0^{t^*} dt_1 \left| \frac{t^*}{t^*2} \right| \\
&= C \left( \max_{0 \leq t \leq t^*} |\psi_+(t)| \right) \frac{t^*}{2}.
\end{aligned}$$

Hence for  $t^* < 2/C$  it follows  $\|K\| < 1$  and  $\lim_{t \rightarrow 0} (K\psi_+)(t^*) = 0$ . Thus the integral operator vanishes in comparison to  $\psi_+$  as  $t \rightarrow 0$ , which shows the first statement of the lemma, equation (??). Moreover, by Banach's fixed-point theorem there exists a unique continuous solution on the interval  $[0, t^*]$ . Since the kernel of the integral operator  $K$  is bounded for all  $t \geq t^*$  we obtain simultaneously  $\psi_+ \in C([0, \infty))$ .

Induction step. Now let  $\psi_+ \in C^k([0, \infty))$ , for some  $k \geq 0$ . Then there exists a residual function  $R_k$  such that

$$\psi_+(t) = \sum_{j=0}^k \psi_+^{(j)}(0) \frac{t^j}{j!} + R_k(t) \quad \text{and } R_k(t) = o(t^k) \quad (t \rightarrow 0).$$

We compute the action of the integral operator on the residual function

$$\begin{aligned}
|(KR_k)(t)| &\leq \frac{1}{t|t-2i|} \int_0^t dt_1 |P(t-t_1)| \sup_{0 \leq t_1 \leq t} |R_k(t_1)| \\
&= o(t^{k+1}) \quad (t \rightarrow 0)
\end{aligned}$$

since  $P(t) = \mathcal{O}(t)$  ( $t \rightarrow 0$ ). Hence, the application of the integral operator increases the order by one,  $KR_k \in C^{k+1}([0, \infty))$ , also we see  $(KR_k)^{(k+1)}(0) = 0$ . Now it follows from the integral equation (??), the analyticity of its right-hand side, and our series representation (4.3.18), by a component-wise comparison, that  $\psi_+ \in C^{k+1}([0, \infty))$  and that  $\psi_+^{(k+1)}(0)$  satisfies (??).  $\square$

LEMMA 4.3.6. *Let positive numbers  $\epsilon_1, \tilde{a}$  and  $c$  be given such that  $a - \tilde{a} - \epsilon_1 > 0$ , and let*

$$|P(t)| \leq ce^{-(a-\tilde{a}-\epsilon_1)t}, \quad t > 0.$$

*Let  $\epsilon_2$  such that  $0 < \epsilon_2 < a - \tilde{a} - \epsilon_1$ . Then, there exists  $0 < \zeta < t^* := \sqrt{c/(a - \tilde{a} - \epsilon_1 - \epsilon_2)}$  such that for all  $t^* \leq t$*

$$|\psi_{\pm}(t)| \leq (a - \tilde{a} - \epsilon_1) (\pi + t^* |\psi_+(\zeta)|) e^{-\epsilon_2(t-t^*)}.$$

PROOF. The idea of the proof is based on Gronwall's lemma, cf. [25, Chapt. 3.1, p. 63]. As in Lemma 4.3.5, we start from the integral equation for the cut function  $\psi$ , (??).

$$\psi_+(t) + \int_0^t dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1) = \pi \frac{P(t)}{t(t-2i)}, \quad t > 0.$$

Next we take into account the structure of  $P(\cdot)$ , given by its definition (4.2.12) together with (4.2.11), and our estimate Lemma A.1.8, that is

$$\begin{aligned}
|P(t)| &\leq c_1 e^{(\tilde{a}-a)t} + c_2 t e^{-at}, & t > 0, & \quad \text{with } \tilde{a} < a \\
&\leq ct e^{(\tilde{a}-a)t} \\
&\leq ce^{(\tilde{a}-a+\epsilon_1)t}, & \text{with } \epsilon_1 > 0 & \text{ such that } \tilde{a} - a + \epsilon_1 < 0.
\end{aligned}$$

Hence, there are always positive constants  $\epsilon_1$  and  $c$  satisfying the assumptions of the lemma. We introduce the notation

$$|P(t)| \leq ce^{-\delta_a t}, \quad \delta_a := a - \tilde{a} - \epsilon_1 > 0.$$

As first step, we split the integration interval into two intervals  $[0, t^*]$  and  $[t^*, t]$ . Accordingly, the integral equation is given by

$$\begin{aligned}
\psi_+(t) &= \pi \frac{P(t)}{t(t-2i)} - \int_0^{t^*} dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1) \\
&\quad - \int_{t^*}^t dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1).
\end{aligned}$$

Next we multiply by  $\exp(\delta_a t)$  and take the absolute values

$$\begin{aligned}
(4.3.19) \quad |e^{\delta_a t} \psi_+(t)| &\leq \pi \left| e^{\delta_a t} \frac{P(t)}{t(t-2i)} \right| + \int_0^{t^*} dt_1 e^{\delta_a t} \left| \frac{P(t-t_1)}{t(t-2i)} \right| |\psi_+(t_1)| \\
&\quad \int_{t^*}^t dt_1 e^{\delta_a t} \left| \frac{P(t-t_1)}{t(t-2i)} \right| |\psi_+(t_1)|, \quad t \geq t^*.
\end{aligned}$$

We consider the integral terms separately. First we have, by the mean value theorem,

$$\begin{aligned}
\int_0^{t^*} dt_1 e^{\delta_a t} \left| \frac{P(t-t_1)}{t(t-2i)} \right| |\psi_+(t_1)| &= \frac{1}{|t(t-2i)|} \int_0^{t^*} dt_1 e^{\delta_a t} P(t-t_1) |\psi_+(t_1)| \\
&= \frac{1}{|t(t-2i)|} t^* e^{\delta_a t} P(t-\zeta) |\psi_+(\zeta)|, \quad 0 < \zeta < t^* \\
&\leq \frac{1}{|(t^*-2i)|} e^{\delta_a t} ce^{-\delta_a(t-\zeta)} |\psi_+(\zeta)| \\
&= \frac{ce^{\delta_a \zeta}}{|(t^*-2i)|} |\psi_+(\zeta)|.
\end{aligned}$$

Next, let us choose real positive  $t^*, \epsilon_2$  such that

$$|t^*(t^*-2i)| > \frac{c}{\delta_a - \epsilon_2}, \quad \delta_a - \epsilon_2 > 0,$$

which is satisfied for

$$t^* = \sqrt{\frac{c}{\delta_a - \epsilon_2}}.$$

This allows to bound the second integral as

$$\begin{aligned}
\int_{t^*}^t dt_1 e^{\delta_a t} \left| \frac{P(t-t_1)}{t(t-2i)} \right| |\psi_+(t_1)| &\leq \int_{t^*}^t dt_1 e^{\delta_a t} c e^{-\delta_a(t-t_1)} \left| \frac{1}{t(t-2i)} \right| |\psi_+(t_1)| \\
&= \int_{t^*}^t dt_1 e^{\delta_a t_1} c \left| \frac{1}{t(t-2i)} \right| |\psi_+(t_1)| \\
&\leq \int_{t^*}^t dt_1 e^{\delta_a t_1} (\delta_a - \epsilon_2) |\psi_+(t_1)|.
\end{aligned}$$

We bound the first term of the right-hand side of (4.3.19) by the same technique

$$\begin{aligned}
\pi \left| e^{\delta_a t} \frac{P(t)}{t(t-2i)} \right| &\leq \pi \left| e^{\delta_a t} \frac{c e^{-\delta_a t}}{t(t-2i)} \right| \\
&= \pi \left| \frac{c}{t(t-2i)} \right| \\
&\leq \pi (\delta_a - \epsilon_2).
\end{aligned}$$

Hence, all together we find the inequality

$$\begin{aligned}
\phi(t) &\leq \rho + \int_{t^*}^t dt_1 (\delta_a - \epsilon_2) |\phi(t_1)|, \quad t \geq t^* \\
\phi(t) &:= |e^{\delta_a t} \psi_+(t)| \\
\rho &:= \pi (\delta_a - \epsilon_2) + \frac{c e^{\delta_a \zeta}}{|(t^* - 2i)|} |\psi_+(\zeta)|.
\end{aligned}$$

It holds, by Gronwall's lemma

$$\begin{aligned}
\phi(t) &\leq \rho e^{\int_{t^*}^t dt_1 (\delta_a - \epsilon_2)} \\
&= \rho e^{(\delta_a - \epsilon_2)(t - t^*)}.
\end{aligned}$$

This yields

$$|\psi_+(t)| \leq \rho e^{-(\delta_a - \epsilon_2)t^* - \epsilon_2 t}$$

which results via

$$\begin{aligned}
|\psi_+(t)| &\leq (\delta_a - \epsilon_2) \left( \pi e^{-\delta_a t^*} + t^* |\psi_+(\zeta)| \right) e^{-\epsilon_2(t-t^*)} \\
0 &< \zeta < t^* \\
0 &< \epsilon_2 < \delta_a.
\end{aligned}$$

finally in

$$|\psi_+(t)| \leq \delta_a (\pi + t^* |\psi_+(\zeta)|) e^{-\epsilon_2(t-t^*)}, \quad t^* \leq t, 0 < \epsilon_2 < \delta_a, 0 < \zeta < t^*$$

Since  $\psi_+$  and  $\psi_-$  are complex conjugates, the same statement holds for  $\psi_-$ .  $\square$

The statements of Lemma 4.3.5 and Lemma 4.3.6 are – besides their value for the theoretical characterization – of vital interest for the numerical implementation. Since  $\psi_{\pm} \in C^\infty(\mathbb{R}_+)$  we may apply high order discretization methods with good effect, and since  $\psi_{\pm}$  decays exponentially for large distances on the cut, we may restrict the size of the support of approximating functions to a moderate size. In fact, this rapid decay is one of the characteristic differences between the numerical algorithms based on the cut function formulation and the direct computation of  $\widehat{U}$  with the pole condition realized along the real axis.

#### 4.4. Important Consequences

In the following we will use the theory of Section 4.2 to prove the essential results needed both for the theoretical justification of the pole condition approach as well as for the numerical implementation.

**4.4.1. Representation formula.** First we want to prove one of the most practical results of our analysis – a representation formula of the spatial function  $U$  in terms of the cut functions  $\psi_{\pm}$ . The derived formulas give not only the exterior solution, we will also use them as linking conditions between the interior problem and the exterior problem in our algorithms. Moreover, they form the basis of our understanding of the perfectly matched layer method.

**THEOREM 4.4.1.** *The function  $U(r)$  and its derivatives  $U^{(k)}(r)$ , ( $k \geq 1$ ), have holomorphic extensions to the right half of the complex plane and satisfy the representation formulas*

$$(4.4.1) \quad U(r) = w(i)e^{ir} \left( \frac{1}{2i} - \frac{1}{2\pi i} \int_0^{\infty} dt e^{-tr} \psi_+(t) \right) \\ + w(-i)e^{-ir} \left( -\frac{1}{2i} + \frac{1}{2\pi i} \int_0^{\infty} dt e^{-tr} \psi_-(t) \right)$$

$$(4.4.2) \quad U^{(k)}(r) = w(i)e^{ir} \left( \frac{i^k}{2i} - \frac{1}{2\pi i} \int_0^{\infty} dt (i-t)^k e^{-tr} \psi_+(t) \right) \\ + w(-i)e^{-ir} \left( -\frac{(-i)^k}{2i} - \frac{1}{2\pi i} \int_0^{\infty} dt (-i-t)^k e^{-tr} \psi_-(t) \right)$$

for  $\text{Re}(r) \geq 0$ .

**PROOF.** The idea is to compute the inverse of the Laplace transformed spatial function  $\widehat{U}$  by means of the inversion formula utilizing Cauchy's contour integral, contour deformation and our results on the decay behavior of  $\widehat{U}$  and  $\psi_{\pm}$  for large arguments. Let the contour parallel to the imaginary axis be denoted by  $\gamma_1(t) := 1 + it$ ,  $-R \leq t \leq R$ , see Fig. 4.3.4. The inversion formula is

$$U(r) = \lim_{R \rightarrow \infty} \frac{1}{2\pi i} \int_{\gamma_1} ds e^{rs} \widehat{U}(s).$$

In standard manner we use Cauchy's integral theorem

$$\frac{1}{2\pi i} \left( \int_{\gamma_1} ds e^{rs} \widehat{U}(s) + \int_{\gamma_2} ds e^{rs} \widehat{U}(s) \right) = 0$$

to compute the inversion integral based on the contributions from the path  $\gamma_2$ . The integrals from  $B$  to  $C$ , from  $C$  to  $D$  and from  $D$  to  $A$  vanish as  $R \rightarrow \infty$ , since  $|\widehat{U}(s)| = \mathcal{O}(|s|^{-1})$  as  $|s| \rightarrow \infty$  and since the integrand decays exponentially as  $\text{Re}s \rightarrow -\infty$ . The integrals along  $S_{\pm}$  converge to  $w(i) \int_0^{\infty} dt \exp(r[\pm i - t]) \psi_{\pm}(t)$  as  $R \rightarrow \infty$ . A computation similar to that in the proof of Lemma 4.3.4 shows that the integral around  $i$  converges to  $\pi w(i)e^{ir}$ , and the integral around  $-i$  to  $\pi w(-i)e^{-ir}$  as  $R \rightarrow \infty$ . Summing the parts, we obtain (4.4.1). Differentiating this equation and changing the order of differentiating and integration, which is possible by Lebesgue's Dominated Convergence Theorem, we obtain the second equation, (4.4.2).  $\square$

Thus we have generalized the partial fraction decomposition technique used in the introductory examples to the general case. Eq. (4.4.1) supplies in a very natural way the desired decomposition of a function into its incoming and its outgoing part.



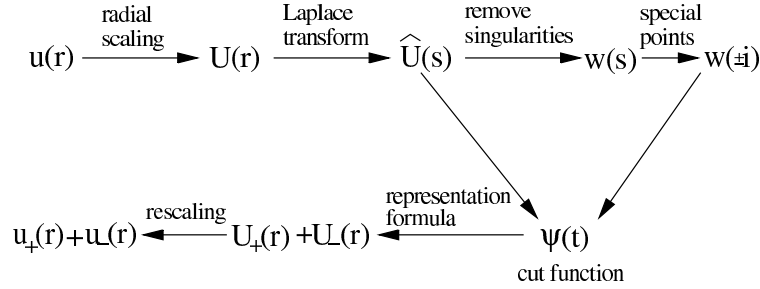


FIGURE 4.4.1. Decomposition of  $u$  into an incoming part  $u_-$  and an outgoing part  $u_+$

Further we note that the incoming part is proportional to  $w(-i)$ , the outgoing part to  $w(+i)$ . Since with  $w(i)$  the pole of  $\widehat{U}(s)$  vanishes, too, we refer to  $w_{\pm}(i)$  as pole strengths. Now the proof of the following corollary is very easy. Nevertheless it contains the central property that the pole condition can always be satisfied by a proper choice of the initial conditions.

COROLLARY 4.4.2. *The matrix  $L(r)$  defined by*

$$\begin{pmatrix} U(r) \\ U'(r) \end{pmatrix} = L(r) \begin{pmatrix} w(i) \\ w(-i) \end{pmatrix}$$

using the representation formulas (4.4.1) and (4.4.2) is regular for all  $r \geq 0$ . Hence there exist complex numbers  $U(r), U'(r)$  such that  $w(-i) = 0$ . In this case, the spatial function  $U$  satisfies the pole condition, i.e.  $\widehat{U}$  is holomorphic in the lower half of the complex plane.

PROOF. Let

$$L(r) \begin{pmatrix} w(i) \\ w(-i) \end{pmatrix} = 0.$$

Since  $U$  solves a linear second order differential equation,  $U(r) = U'(r) = 0$  implies  $U = 0$ , thus  $w(i) = w(-i) = 0$ . If the pole strength  $w(-i)$  vanishes, the cut function  $\psi_-$  vanishes, too, since  $\psi_- = w(-i)[U](-i-t)$ . Therefore,  $[w](-i-t) = 0$  for all  $t > 0$ , which tells us that  $w$  is continuous in the lower half of the complex plane. Using Morera's Theorem and a contour deformation around the cut  $S_{-i}$  it can be shown that  $w$  is even holomorphic in the lower half plane. As  $w(-i) = 0$ ,  $\widehat{U}$  is holomorphic in the lower half-plane.  $\square$

**4.4.2. Equivalent Formulations of the Pole Condition.** We want to summarize the statements obtained so far concerning the pole condition. Before doing so, we summarize the interdependence of the different quantities introduced so far. Fig. 4.4.1 gives the corresponding graphical representation. After the radial scaling we obtain our basic spatial function  $U(r)$ , which is then Laplace transformed to get  $\widehat{U}(s)$ . The function  $\widehat{U}(s)$  is multiplied with  $s^2 + 1$  which yields a function  $w(s)$  without any singularity. Based on  $\widehat{U}(s)$  a cut function  $\psi(t)$  is defined living only on the cuts  $S_{\pm}$ , and together with the values of  $w$  at  $s = \pm i$ , the spatial function  $U(r)$  can be recovered, decomposed into an outgoing and an incoming part.

The representation formulas of Theorem 4.4.1 tell us in which way the function  $\widehat{U}$  is expressed by the cut functions  $\psi_{\pm}$  together with the pole strengths  $w(\pm i)$ . If  $w(+i) = w(-i) = 0$ ,  $\widehat{U}$  vanishes identically, consequently  $U(r) \equiv 0$ . If  $w(-i) = 0$ , then  $\widehat{U} = \widehat{U}_+$ , hence the corresponding spatial function  $U(r)$  is a pure outgoing

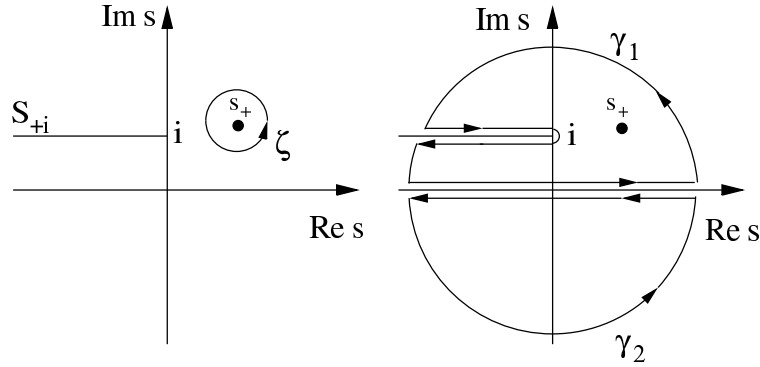


FIGURE 4.4.2. Deformation of the contour from a small circle (left) to the final contour

function and  $\widehat{U}$  satisfies the pole condition, i.e. it is holomorphic in the lower half of the complex plane. Thus we may write the current form of the pole condition simply as

Pole condition in terms of the pole strength  
 $w(-i) = 0.$

This is the most simple form of the pole condition. However, this form is not constructive because it does not show how to realize the pole condition. On the other hand, the vanishing pole strength  $w(-i) = 0$  means that  $U(r)$  satisfies the pole condition if it has a representation

Pole condition in terms of the representation formulas  
 –cut function approach–

$$U(r) = w(i)e^{ir} \left( \frac{1}{2i} - \frac{1}{2\pi i} \int_0^\infty dt e^{-tr} \psi_+(t) \right)$$

with some  $w(i) \in \mathbb{C}$  and  $\psi_+$  a solution of

$$\psi_+(t) + \int_0^t dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1) = \pi \frac{P(t)}{t(t-2i)}.$$

Practically this requires the computation of the cut function and the pole strength  $w(i)$ . This, in turn, requires to compute the representation formulas (4.4.1) and (4.4.2) at  $r = 0$  for  $k = 0, 1$ . We shall discuss this in detail in the application chapter, where we will use these formulas to construct a numerical algorithm, which we call in accordance with this theoretical framework *cut-function approach*.

There is a third form of the pole condition, formulated for the first time in [81].

LEMMA 4.4.3. *Let  $\widehat{U}$  satisfy the pole condition. Then,*

$$0 = \int_{-\infty}^{\infty} d\tau \frac{\widehat{U}(\tau)}{\tau - s_+} \quad \text{for all } s_+ \in V \text{ with } \text{Im}(s_+) > 0.$$

PROOF. We apply Cauchy's integral theorem. Let  $s_+$  given as defined above. First, we choose a small circle around  $s_+$ , see Fig. 4.4.2. Since  $s_+$  is the only singularity inside the small circle, it holds

$$\widehat{U}(s_+) = \frac{1}{2\pi i} \oint_{\gamma} d\zeta \frac{U(\zeta)}{\zeta - s_+}.$$

We deform the contour as shown in Fig. 4.4.2 and split the path  $\gamma$  into the two paths  $\gamma_1$  and  $\gamma_2$ . This defines two functions  $\widehat{U}_1(s_+)$  and  $\widehat{U}_2(s_+)$  both depending on the radius  $R$ , and we have  $\widehat{U}(s_+) = \widehat{U}_1(s_+) + \widehat{U}_2(s_+)$ . We let  $R \rightarrow \infty$  and denote the limiting functions again by  $\widehat{U}_1(s_+)$  and  $\widehat{U}_2(s_+)$ . Now, since  $U(\zeta)$  is holomorphic, we obtain

$$0 = \frac{1}{2\pi i} \oint_{\gamma_2} d\zeta \frac{\widehat{U}(\zeta)}{\zeta - s_+} \quad \text{for all } s_+ \in V \text{ with } \text{Im}(s_+) > 0.$$

Additionally, it holds by construction,  $\widehat{U}(s) = \mathcal{O}(|s|^{-1})$  as  $|s| \rightarrow \infty$ . Thus we have  $|\widehat{U}(s)/(s - s_+)| = \mathcal{O}(|s|^{-2})$  and the contour integral on the large semi-circle vanishes. This is the desired result.  $\square$

REMARK. The decomposition  $\widehat{U}(s_+) = \widehat{U}_1(s_+) + \widehat{U}_2(s_+)$  via the Cauchy integral corresponds to the partial fraction decomposition discussed in connection with our introductory 1D example, which has motivated us to study the properties of the singularities in more general cases.

Now we are ready to give the third form of the pole condition.

Pole condition in terms of the axis integral  
–real axis approach–

$$0 = \int_{-\infty}^{\infty} d\tau \frac{\widehat{U}(\tau)}{\tau - s_+}, \quad \tau \in \mathbb{R}, \text{ for fixed } s_+ \in V \text{ with } \text{Im}(s_+) > 0.$$

with  $\widehat{U}$  a solution of

$$\widehat{U}(\tau) + \int_{\tau}^{\infty} d\tau_1 \frac{P(\tau_1 - \tau)\widehat{U}(\tau_1)}{(k^2 + \tau^2)} = \frac{\tau U(0) + U'(0)}{(k^2 + \tau^2)}, \quad \tau \in \mathbb{R}.$$

Observe the remarkable symmetry of these formulas in comparison to those related to the cut function approach. Observe further, that the formulas of the real axis approach does not possess any singularity on the real axis, which neither makes it necessary to carry out the transform (4.2.7) nor to exclude the term  $1/r$  in the definition of the potentials  $p_*(1/r)$ , (4.2.11).

#### 4.4.3. Asymptotic Expansion of the Far-Field.

THEOREM 4.4.4. *For all  $m \geq 0$ ,  $U$  and  $U'$  of an outgoing function satisfy the asymptotic formulas*

$$(4.4.3) \quad \frac{U(r)}{w(i)} = e^{ir} \left( \frac{1}{2i} - \frac{1}{2\pi i} \sum_{j=1}^m \frac{\psi_+^{(j-1)}(0)}{r^j} + \mathcal{O}\left(\frac{1}{r^{m+1}}\right) \right)$$

$$(4.4.4) \quad \frac{U'(r)}{w(i)} = e^{ir} \left( \frac{1}{2} - \frac{1}{2\pi i} \sum_{j=1}^m \frac{i\psi_+^{(j-1)}(0) - (j-1)\psi_+^{(j-2)}}{r^j} + \mathcal{O}\left(\frac{1}{r^{m+1}}\right) \right)$$

for  $r \rightarrow \infty$  such that  $|\arg r| \leq \phi < \pi/2$ . Here  $0 \cdot \psi_+^{(-1)}(0) := 0$ .

PROOF. We start from the outgoing part of the representation formula (4.4.1)

$$\frac{U(r)}{w(i)} = e^{ir} \left( \frac{1}{2i} - \frac{1}{2\pi i} \int_0^\infty dt e^{-tr} \psi_+(t) \right).$$

Note that the integral term is the Laplace transform of  $\psi_+$ , where the dual variable is  $r$ . Using the asymptotic formula

$$(Lf)(s) = \sum_{j=1}^m \frac{f^{(j-1)}(0)}{s^j} + \mathcal{O}(s^{-m-1}), \quad s \rightarrow \infty, |\arg s| \leq \phi < \pi/2,$$

which holds for any bounded function  $f \in C^m(\mathbb{R}_+)$ , (cf. [26, Vol. 2, p. 47] and Lemma 4.3.5 we obtain

$$\begin{aligned} \frac{U(r)}{w(i)} &= e^{ir} \left( \frac{1}{2i} - \frac{1}{2\pi i} (L\psi_+)(r) \right) \\ &= e^{ir} \left( \frac{1}{2i} - \frac{1}{2\pi i} \sum_{j=1}^m \frac{\psi_+^{(j-1)}(0)}{r^j} + \mathcal{O}(r^{-m-1}) \right). \end{aligned}$$

This is (4.4.3). The asymptotic formula for  $U'$  follows analogously from the representation formula for  $U'$ , equation (4.4.2) written as

$$\frac{U'(r)}{w(i)} = e^{ir} \left( \frac{1}{2} - \frac{1}{2\pi i} (L[(i-t)\psi_+])(r) \right)$$

and the identity

$$\frac{d^j}{dt^j} ((i-t)\psi_+(t)) \Big|_{t \rightarrow 0} = i\psi_+^j(0) - j\psi_+^{(j-1)}(0).$$

□

COROLLARY 4.4.5. *The Hankel functions  $H_\nu^{(1)}$  of the first kind of order  $\nu$  satisfy*

$$H_\nu^{(1)}(r) = \sqrt{\frac{2}{\pi r}} e^{i(r - \frac{\nu\pi}{2} - \frac{\pi}{4})} \left( \sum_{k=0}^m \left\{ \prod_{j=1}^k \frac{\nu^2 - (j - \frac{1}{2})^2}{-2ijr} \right\} + \mathcal{O}(r^{-m-1}) \right)$$

for  $r \rightarrow \infty$  such that  $|\arg r| \leq \phi < \pi/2$ .

PROOF. We take

$$\begin{aligned} p_*(t) &= 0 \\ P(t) &= e^{-at} t \left( \frac{1}{4} - \nu^2 \right) \end{aligned}$$

which results from (4.2.12) and (4.2.10) for a constant potential. With the normalization

$$w(i) = \sqrt{\frac{8}{\pi}} \exp\left(i\frac{\pi}{4} - \frac{\nu\pi}{2}\right)$$

we get  $H_\nu^{(1)}(r+a) = (r+a)^{1/2} U(r)$ . Carrying out the limit  $a \rightarrow 0$  and using (??), (??) from Lemma 4.3.5 and the identity

$$\frac{1}{t(t-2i)} = \frac{1}{-2it} \sum_{j=0}^{\infty} t^j (2i)^{-j}$$

supplies

$$\begin{aligned} \psi_+(0) &= 2\pi \frac{\nu^2 - \frac{1}{4} w(i)}{2i} \frac{1}{2} \\ \psi_+^{(k+1)}(0) &= \left( \frac{k+1}{2i} + \frac{1}{2i(k+2)} \left( \frac{1}{4} - \nu^2 \right) \right) \psi_+^{(k)}(0) \\ &= \frac{1}{-2i(k+2)} \left( \nu^2 - \left( k + \frac{3}{2} \right)^2 \right) \psi_+^{(k)}(0). \end{aligned}$$

□

**4.4.4. Spectral Properties of the Dirichlet-to-Neumann Map .** We define a complete orthonormal system of eigenfunctions of the Laplace-Beltrami operator on  $\Gamma_a$  by  $\phi_j^a(ax^0) := a^{-\frac{d-1}{2}} \phi_j(x^0)$ . Then, for  $f \in H^{1/2}(\Gamma_a)$ , we expect that the solution to our underlying Helmholtz equation (4.2.1) satisfying the pole condition and the boundary condition  $\text{Tr}_{\Gamma_a} u = f$  is given by

$$u((r+a)\mathbf{x}^0) := \sum_{j=1}^{\infty} \langle \phi_j^a, f \rangle_{\Gamma_a} (r+a)^{-\frac{d-1}{2}} U_j(r) \phi_j(\mathbf{x}^0)$$

for  $r \geq 0$  and  $\mathbf{x}^0 \in S^{d-1}$ , where the functions  $U_j(r)$  are solutions of (4.2.8). This will be shown below. Differentiating the last equation with respect to  $r$  we formally find the following formula for the Dirichlet-to-Neumann map

$$\frac{\partial}{\partial r} u(ax^0) = \sum_{j=1}^{\infty} \text{DtN}(\lambda_j, a) \langle \phi_j^a, f \rangle_{\Gamma_a} \phi_j(\mathbf{x})$$

with the eigenvalues

$$\begin{aligned} \text{DtN}(\lambda_j, a) &= \frac{\left( (r+a)^{-\frac{d-1}{2}} U_j(r) \right)' \Big|_{r=0}}{a^{-\frac{d-1}{2}} U_j(0)} \\ (4.4.5) \quad &= \frac{U_j'(0)}{U_j(0)} - \frac{d-1}{2a}. \end{aligned}$$

Since the Sobolev norm of index  $s \in \mathbb{R}$  is given by

$$\begin{aligned} \|f\|_{H^s(\Gamma_a)}^2 &= \left\| (I - \Delta_{\Gamma_a})^{s/2} f \right\|_{L^2(\Gamma_a)}^2 \\ &= \sum_{j=1}^{\infty} (1 - \lambda_j)^s |\langle \phi_j^a, f \rangle|^2 \end{aligned}$$

(cf. e.g. the book of Taylor [88]) the properties (i)-(iii) in Theorem 4.2.1 are equivalent to

$$(4.4.6) \quad |\text{DtN}(\lambda_j, a)| = \mathcal{O}\left(\sqrt{|\lambda_j|}\right), \quad j \rightarrow \infty$$

$$(4.4.7) \quad \text{Re}(\text{DtN}(\lambda_j, a) + l_j) \leq 0 \text{ for some sequence } |l_j| = o\left(\sqrt{|\lambda_j|}\right)$$

$$(4.4.8) \quad \text{Im DtN}(\lambda_j, a) > 0 \text{ for all } j.$$

The most difficult part is to establish the spectral property (4.4.6). Unfortunately, we cannot give a complete proof for the spectral property. In our earlier work we showed how to prove this property in the context of the pole condition for the case of a constant potential and two space dimensions [86]. Since the main tool was an explicit representation of the functions  $\tilde{U}$ , obtained via a sequence of transforms, this technique does not generalize to arbitrary potentials. In a joint work with Hohage and Zschiedrich [58] a complete proof is presented. However, the techniques applied there do not fit well in our presentation here, so we will restrict ourselves to an incomplete, but convincing motivation of the desired property.

We establish the property (4.4.8) first. The idea is to apply the conservation property of the energy flux, cf. [86].

LEMMA 4.4.6. *If the spatial function  $U_j$  satisfies the pole condition, then*

$$\text{Im}\left(U_j'(r)\overline{U_j(r)}\right) = \left|\frac{w(i)}{2}\right|^2$$

for all  $r \geq 0$ . In particular,  $U_j(r) \neq 0$  for all  $r \geq 0$  and

$$\text{Im}\left(\frac{U_j'(r)}{U_j(r)}\right) > 0.$$

PROOF. We start from our basic equation for  $U$  in the spatial domain (4.2.8)

$$U_j''(r) + K(r)U_j(r) = 0.$$

$$K(r) := \left[ k^2 + \frac{\frac{1}{4}(d-1)(3-d) + \lambda_j}{(r+a)^2} + p\left(\frac{1}{r+a}\right) \right]$$

with  $K$  a real function. Multiplying by  $\overline{U_j}$ , integrating the expression on the interval  $[0, r]$  yields

$$\begin{aligned} 0 &= \int_0^r dr_1 \overline{U_j}(U_j''(r_1) + K(r_1)U_j(r_1)) \\ &= \overline{U_j(r)}U_j'(r) - \overline{U_j(0)}U_j'(0) + \int_0^r dr_1 \left(\overline{U_j'(r_1)}U_j'(r_1) + K(r_1)\overline{U_j(r_1)}U_j(r_1)\right). \end{aligned}$$

Taking the imaginary part of this equation, we obtain  $\text{Im}\left(\overline{U_j(r)}U_j'(r)\right) = \text{Im}\left(\overline{U_j(0)}U_j'(0)\right) = \text{const.}$  The constant can be evaluated using Theorem 4.4.4 via

$$\lim_{r \rightarrow \infty} \text{Im}\left(\overline{U_j(r)}U_j'(r)\right) = \left(-\frac{1}{2i}\overline{w(i)}e^{-ir}\right) \left(\frac{1}{2}w(i)e^{ir}\right) = \left|\frac{w(i)}{2}\right|^2.$$

It follows that both  $U(r) \neq 0$  and  $U'(r) \neq 0$  for all  $r \geq 0$ . Additionally, we know by ODE-theory that  $U(r)$  is bounded on each compact interval. By Theorem 4.4.4, it remains bounded on the whole semi-infinite interval. Hence it follows that

$$\text{Im}\left(\frac{U_j'(r)}{U_j(r)}\right) = \left|\frac{w(i)}{2\overline{U_j(r)}U_j(r)}\right|^2 > 0.$$

□

Next we consider the spectral properties (4.4.6) and (4.4.7).

**4.4.5. High-Frequency Limit of the Cut Functions.** We look at the limit process  $\lambda_j \rightarrow -\infty$ . Since the eigenvalues of the Laplace-Beltrami operator are all non-positive, we may equivalently write  $\lambda_j = -\nu^2$  with  $\nu \geq 0$  and let  $\nu \rightarrow \infty$ . Let us consider our integral equation (74) with the polynomial  $P(\cdot)$  for the constant potential, see (4.2.12), chosen according to our high-frequency consideration

$$\begin{aligned} \psi_+(t) + \int_0^t dt_1 \frac{P(t-t_1)}{t(t-2i)} \psi_+(t_1) &= \pi \frac{P(t)}{t(t-2i)} \\ P(t) &:= -\nu^2 e^{-at} t. \end{aligned}$$

We introduce a new variable  $\tau$  by  $t := \nu\tau$  and normalize  $\psi_+$  to  $\tilde{\psi}_\nu := c_\nu \psi_+$  via

$$(4.4.9) \quad \int_0^\infty d\tau \tilde{\psi}_\nu(\tau) = 1.$$

By Lemma 4.3.6, this is always possible. Next, we reformulate the integral equation using the new quantities

$$(4.4.10) \quad \tilde{\psi}_\nu(\tau) + \int_0^\tau d\tau_1 \frac{\tilde{P}(\tau-\tau_1)}{\nu\tau(\tau-\frac{2i}{\nu})} \tilde{\psi}_\nu(\tau_1) = \pi \frac{c_\nu \tilde{P}(\tau)}{\nu^2 \tau (\tau - \frac{2i}{\nu})}$$

$$(4.4.11) \quad \tilde{P}(\tau) = -\nu^3 e^{-a\nu\tau} \tau.$$

Now the idea is to approximate the integral for functions  $\tilde{\psi}_\nu(\tau)$  which vary at each point  $\tau$  (with respect to the new scale) much slower than  $\exp(-a\nu\tau)$  by

$$\begin{aligned} \int_0^\tau d\tau_1 \frac{\tilde{P}(\tau-\tau_1)}{\nu\tau(\tau-\frac{2i}{\nu})} \tilde{\psi}_\nu(\tau_1) &= - \int_0^\tau d\tau_1 \frac{\nu^2 e^{-a\nu(\tau-\tau_1)} (\tau-\tau_1)}{\tau(\tau-\frac{2i}{\nu})} \tilde{\psi}_\nu(\tau_1) \\ &\approx - \frac{\tilde{\psi}_\nu(\tau)}{(a\tau)^2} \quad \text{for } \nu \text{ large.} \end{aligned}$$

Hence, the integral operator, which is compact, converges to a multiplication operator, which is not compact. In order to motivate our approach further, let us assume for the moment, that this integral approximation is exact and let us replace the integral operator in (4.4.10) by the multiplication operator and multiply the whole equation by  $(a\tau)^2$

$$(4.4.12) \quad \tilde{\psi}_\nu(\tau) \left( (a\tau)^2 - 1 \right) = -\pi a^2 c_\nu \nu e^{-a\nu\tau}.$$

The right-hand side tends to zero exponentially fast for any  $\tau \geq \epsilon > 0$  if  $c_\nu$  is bounded. Hence  $\tilde{\psi}_\nu$  tends to zero on the same interval except at the position  $\tau = 1/a$ . Consequently we expect  $\tilde{\psi}_\nu(\cdot) \rightarrow \delta(\cdot - 1/a)$  as  $\nu \rightarrow \infty$ . If we integrate both sides of (??) over the entire real axis, we find that the left-hand side tends to zero for  $\nu \rightarrow \infty$  due to the delta property. The integration of the right-hand side yields  $\pi a c_\nu$ . Hence we expect  $c_\nu \xrightarrow{\nu \rightarrow \infty} 0$ . By the foregoing it seems reasonable to conjecture the following:

- (1)  $\tilde{\psi}_\nu(\cdot) \rightarrow \delta(\cdot - 1/a)$  as  $\nu \rightarrow \infty$ , or equivalently,
- (2)  $\psi_\nu(\cdot) \rightarrow 1/c_\nu \delta(\cdot/\nu - 1/a)$  as  $\nu \rightarrow \infty$ , and
- (3)  $c_\nu \rightarrow 0$  as  $\nu \rightarrow \infty$ .

Using these assumptions we are able to derive the spectral properties (4.4.6) and (4.4.7). Based on the delta-property we obtain directly from (4.4.5) and our basic representation Theorem 4.4.1

$$\begin{aligned} \lim_{\nu \rightarrow \infty} \frac{\text{DtN}(-\nu^2, a)}{\nu} &= -\lim_{\nu \rightarrow \infty} \frac{d-1}{2a\nu} + \lim_{\nu \rightarrow \infty} \frac{\frac{1}{2} - \frac{\nu}{2\pi i c_\nu} \int_0^\infty d\tau (i - \nu\tau) \tilde{\psi}_\nu(\tau)}{\frac{\nu}{2i} - \frac{\nu^2}{2\pi i c_\nu} \int_0^\infty d\tau \tilde{\psi}_\nu(\tau)} \\ &= -\frac{1}{a}. \end{aligned}$$

This implies both (4.4.6) and (4.4.7). Thus we have shown, up to the yet incomplete part in the proof of the spectral property, the unique solvability of the scattering problem interior to the artificial boundary  $\Gamma_a$ . We remember again that a complete, but technical proof is given in [58].

**4.4.6. Relation to the PML Method.** There is a very close relation between the pole condition and the PML method, as indicated already in Section 2.5, p. 35. This is surprising, since the motivations behind both methods are very different. We want to discuss this relation in a brief, formal manner; nevertheless an extension of the above theory to the infinite dimensional case can be used to prove rigorously the convergence of the PML method even for problems with position dependent potentials. This work is part of the ongoing research and will be published together with Hohage and Zschiedrich in [57]. The proof, which we will present there, combines the technique of Lassas and Somersalo [67] with our representation formula (4.4.1).

To get the point, however, the more qualitative discussion given below is sufficient. First we specialize the modified Helmholtz equation (4.2.8) from page 65 to the most simple case, namely  $d = 1$ ,  $\lambda_j = 0$ ,  $k = 1$ , and  $p(\cdot) \equiv 0$ . This is the the normalized Helmholtz equation in 1D. The incoming field  $U(r) := \exp(-ir)$  is a solution of this equation and its Laplace transform  $\hat{U}(s) := 1/(s+i)$  has a pole at  $-i$ . If we shift the pole by a distance  $\sigma$ ,  $\sigma \in \mathbb{R}_+$ , to the right, cf. Fig. 4.4.3, we obtain the modified Laplace transform  $\hat{U}_B(s) := 1/(s - [-i + \sigma])$ . The corresponding solution in the spatial domain is  $U_B(r) := \exp([-i + \sigma]r) = \exp(-i[1 + i\sigma]r)$ , where the subscript B in  $\hat{U}_B$  and  $U_B$  indicates the modification due to Bérenger. We can give this a new interpretation as continuation of the real distance variable  $r$  into the upper half of the complex plane. The desired effect of the stretching into the complex plane is that undamped oscillating, incoming functions are mapped to exponentially increasing functions and undamped oscillating outgoing functions are mapped to exponentially decreasing functions.

We want to extend this to the full transformed Helmholtz equation (4.2.8) in arbitrary space dimensions. We define  $\alpha = 1 + i\sigma$ ,  $0 \leq \arg \alpha < \pi/2$  and discuss the complexified distance variable  $\alpha r$ , where  $r$  denotes as before the real, positive physical distance. To analyze the effect of the complexification, we study the representation formula (4.4.1) (p. 80), with the cut functions  $\psi_\pm$  replaced by their original definition (4.3.12)

$$\begin{aligned} U_B(\alpha r) &= e^{i\alpha r} \left( \frac{1}{2i} - \frac{1}{2\pi i} \int_0^\infty dt e^{-t\alpha r} [\hat{U}_+] (i-t) \right) \\ &\quad + e^{-i\alpha r} \left( -\frac{1}{2i} + \frac{1}{2\pi i} \int_0^\infty dt e^{-t\alpha r} [\hat{U}_-] (-i-t) \right). \end{aligned}$$

By Lemma 4.3.5, p. 75, the cut functions are  $C^\infty(\mathbb{R}_+)$ -functions, and by Lemma 4.3.6, p. 77, they decay exponentially for  $r$  sufficiently large. Hence the integrals in the representation formula tend to 0 for  $r$  sufficiently large. Consequently, the behavior



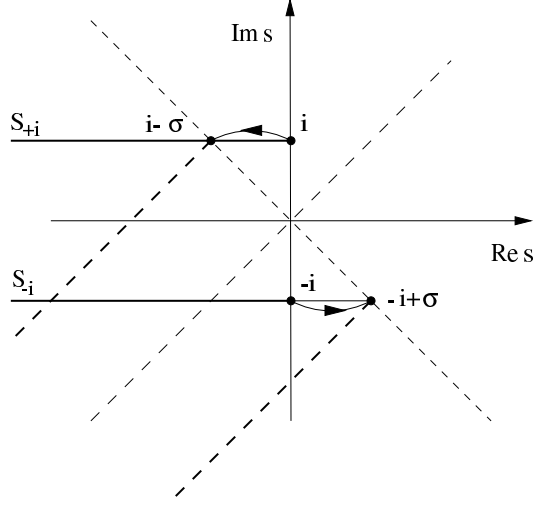


FIGURE 4.4.3. Original and rotated complex plane

of  $U_B(\alpha r)$  at sufficient large distances  $r$  is determined by the exponential factors before the parenthesis. The first term belongs to an exponentially decreasing, the second one to an exponentially increasing function. Exact transparent boundary conditions drop *incoming* functions. Now the central idea of PML is to drop the *incoming* functions, which at the same time are the exponential *increasing* functions by fixing zero-Dirichlet boundary conditions at the outer artificial boundary of the PML layer. Obviously, this forces  $\widehat{U}_-$  to get small. In fact, Lassas and Somersalo proved in [67] that the error of the interior domain decreases like  $\mathcal{O}(\exp(-Cr))$  with  $C$  a positive constant and  $r$  the thickness of the PML-layer.

There is an even closer relation to our analysis based on the pole condition. Since the singularities at  $\pm i$  of our original equations are mapped to  $\pm \alpha i$ , we may consider the complexified problem in a rotated complex plane  $\mathbb{C}_\alpha := \alpha\mathbb{C}$ , compare Fig. 4.4.3. This is the natural transform, since the complexification results in a transform

$$(4.4.13) \quad \alpha \widehat{U}_{B,\alpha}(\alpha s) = \widehat{U}_B(s),$$

see Lemma A.1.9 of the Appendix A.1.

We introduce rotated cuts, see Fig. 4.4.3, define cut functions along these cuts, and carry out the whole analysis as for the original system. Finally we obtain the representation formula

$$U_{B,\alpha}(r) = e^{i\alpha r} \left( \frac{1}{2i} - \frac{1}{2\pi i} \int_0^\infty dt e^{-t\alpha r} \left[ \widehat{U}_{B,\alpha,+} \right] (\alpha(i-t)) \right) \\ + e^{-i\alpha r} \left( -\frac{1}{2i} + \frac{1}{2\pi i} \int_0^\infty dt e^{-t\alpha r} \left[ \widehat{U}_{B,\alpha,-} \right] (\alpha(-i-t)) \right).$$

This coincides, of course, exactly with the representation formula derived for the original system, setting  $U_{B,\alpha}(r) = U_B(\alpha r)$ . Thus we may express the relation between the pole condition approach and the PML method as follows:

- (1) The complexification with a constant factor  $\alpha$  corresponds to a rotation and rescaling in the Laplace domain.

- (2) Singularities on the negative imaginary axis remain singularities in the lower half of the complex plane, hence the corresponding incoming functions remain incoming functions.
- (3) Incoming functions asymptotically increase exponentially.
- (4) The stretching into the complex plane offers the possibility to remove the singularities in the lower half of the complex plane by imposing zero-Dirichlet boundary conditions at infinity, which may be approximated by Dirichlet boundary conditions at sufficient large distances.

## Numerical Treatment of Helmholtz-Type Scattering Problems

We present two qualitative different approaches to solve scattering problems based on the pole condition:

- (1) The factorization method.
- (2) The Laplace domain method.

The factorization technique eliminates the poles in the lower half of the complex plane in a direct manner. That is, these poles are computed by a factorization of matrices and then eliminated by a proper choice of the Dirichlet and Neumann data on the boundary. In contrast, the Laplace domain method computes the Laplace transform either on the real axis (real axis approach) or on the cuts and eliminates the undesired poles by an additional integral condition (real axis approach) or simply by selecting of only one of the two possible cuts (cut function approach). In rather simple cases the factorization approach is preferable, since no additional discretization in the spectral domain is necessary. In more involved situations or non-separable coordinates, the Laplace domain method becomes advantageous, since no additional effort for matrix factorizations is needed.

### 5.1. Factorization Approach

We use the factorization approach to compute the reflection from two different types infinite obstacles: (1) The classical reflection from an infinite plane, (2) Reflection and diffraction from the half-plane already discussed in Section 2.2. This allows a direct comparison with the boundary integral method.

**5.1.1. Reflection at an Infinite Plane.** The reflection of a plane wave on an infinite plane with homogeneous Dirichlet or Neumann boundary conditions is the most simple example of a scattering problem, and analytical solutions are obtained immediately. However, the very closely related example of reflection and diffraction on a semi-infinite plane causes much more difficulties, and special tools from complex function theory like the Wiener-Hopf technique have to be applied. The numerical solution of the simple reflection on the infinite plane based on our Laplace technique with the pole condition in its core seems to be more difficult than the common technique based on an ansatz of plane waves. However, there is no much additional effort if we step over to the semi-infinite scattering problem, neither in the difficulty of the problem solution nor in the computational complexity.

**Continuous Case.** To prepare the semi-discrete case, which yields the desired algorithm, we look first at the continuous case. We wish to solve the Helmholtz equation

$$\begin{aligned} \partial_x^2 u + \partial_y^2 u + k^2 u &= 0 \quad \text{for } y \leq 0 \\ \text{with } \partial_y u(x, y) &= u'(x, 0) = 0 \quad \text{at } y = 0 \quad \text{for all } x \in \mathbb{R} \end{aligned}$$

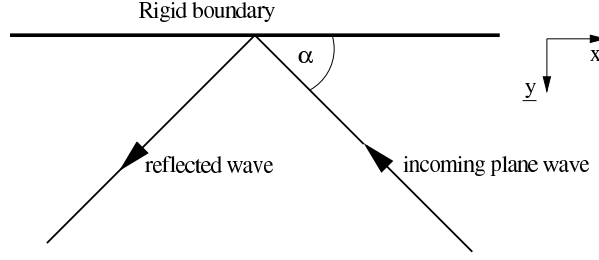


FIGURE 5.1.1. Reflection problem

in the space of twice differentiable, bounded functions on  $\{(x, y) \in \mathbb{R}^2 : y \leq 0\}$ , where an incident plane wave

$$(5.1.1) \quad u_{\text{in}}(x, y) = e^{ik_x x + ik_y y}, \quad \text{with } k_x^2 + k_y^2 = k^2, \quad k_x, k_y \in \mathbb{R}, |k_x| < k$$

excites the reflected field. Note that in this special case the general source field  $u_{\text{src}}$  of the scattering problem (4.0.6) is equal to an incoming field, i. e.  $u_{\text{in}} = u_{\text{src}}$ , in the sense of our definition of incoming and outgoing fields. Further the above bounds on  $k_x, k_y$  exclude explicitly the case  $|k_x| = k$ , that is we exclude phase fronts perpendicular to the reflecting plane. The situation is depicted in Fig. 5.1.1. We define  $\underline{y} := -y$  and use the  $\underline{y}$ -coordinate as unique distance variable; this prepares the reflection problem of the semi-infinite plane. The transverse variable  $x$  describes the position on the reflecting plane. We will solve the reflection problem as follows:

- (1) Fourier transform of the field in  $x$ -direction, viewing the distance coordinate  $\underline{y}$  as real parameter. Thus each Fourier component depends on  $\underline{y}$ .
- (2) Laplace transform of the result with respect to  $\underline{y}$ .
- (3) Imposing the pole condition to each Laplace transformed Fourier component. This supplies directly the field on the axis  $\underline{y} = 0$ .
- (4) Inverse Laplace and Fourier transform to obtain the complete field.

Let  $\mathcal{S}$  be the space of rapidly decreasing functions and  $\mathcal{S}'$  the space of tempered generalized functions (tempered distributions).  $\mathcal{S}'$  includes both the Dirac delta distribution and the undamped exponential function. By Fourier transform  $F : \mathcal{S}' \rightarrow \mathcal{S}'$  along the boundary  $\underline{y} = 0$

$$\widehat{v}(K_x) = \int_{-\infty}^{\infty} e^{-iK_x x} v(x) dx,$$

where the functions  $v$  are associated with the corresponding functionals  $v_f : \mathcal{S} \rightarrow \mathbb{C}$ , hence  $v_f \in \mathcal{S}'$ , via  $v_f(w) = \int_{-\infty}^{\infty} dx vw$  for all  $w \in \mathcal{S}$ , the Helmholtz equation transforms into an ODE with respect to the distance variable  $\underline{y}$

$$-K_x^2 \widehat{u}(K_x, \underline{y}) + \partial_{\underline{y}}^2 \widehat{u}(K_x, \underline{y}) + k^2 \widehat{u}(K_x, \underline{y}) = 0.$$

Thus each Fourier component is indicated by a Fourier frequency  $K_x$ . Note that we did not impose any boundary condition so far.

Both incident and reflected waves have to satisfy this equation. Next we perform the Laplace transform with respect to  $\underline{y}$ , taking  $s$  as dual variable,

$$(5.1.2) \quad (s^2 + k^2 - K_x^2) \widehat{\widehat{u}}(K_x, s) = s\widehat{u}(K_x, 0) + \widehat{u}'(K_x, 0),$$

which yields the corresponding solution

$$(5.1.3) \quad \widehat{u}(K_x, s) = \frac{s\widehat{u}(K_x, 0) + \widehat{u}'(K_x, 0)}{s^2 + k^2 - K_x^2}$$

in terms of the transversal wavenumber  $K_x$  and the Laplace variable  $s$ . We consider  $K_x$  as real parameter. The poles of this equation are given by

$$(5.1.4) \quad s_{\pm} = \pm i\sqrt{k^2 - K_x^2}.$$

The poles  $s_-$  belong to incoming or increasing (in positive  $y$ -direction) fields. Thus (5.1.3) satisfies the pole condition for each Fourier component if the numerator vanishes at the critical points  $s = s_-$ . Since we require that the outgoing solution  $\widehat{u}_{\text{out}}(K_x, s)$  with boundary data  $\widehat{u}_{\text{out}}(K_x, 0)$  and  $\widehat{u}'_{\text{out}}(K_x, 0)$  possesses no pole  $s_-$  the following condition has to be satisfied:

Pole condition for the simple reflection problem

$$(5.1.5) \quad s_- \widehat{u}_{\text{out}}(K_x, 0) + \widehat{u}'_{\text{out}}(K_x, 0) = 0.$$

Using the superposition  $\widehat{u} = \widehat{u}_{\text{in}} + \widehat{u}_{\text{out}}$ , we reformulate this condition in terms of the incoming and the complete field

$$s_- [\widehat{u}(K_x, 0) - \widehat{u}_{\text{in}}(K_x, 0)] + \widehat{u}'(K_x, 0) - \widehat{u}'_{\text{in}}(K_x, 0) = 0.$$

This yields, taking the imposed Neumann condition  $\widehat{u}'(K_x, 0) = 0$  into account,

$$(5.1.6) \quad \widehat{u}(K_x, 0) = \widehat{u}_{\text{in}}(K_x, 0) + \frac{\widehat{u}'_{\text{in}}(K_x, 0)}{s_-}, \quad s_- \neq 0.$$

By definition of the incoming field (5.1.1) we have  $\widehat{u}_{\text{in}}(K_x, 0) = \delta(K_x - k_x)$  and  $\widehat{u}'_{\text{in}}(K_x, 0) = -ik_y \delta(K_x - k_x)$ <sup>1</sup>, hence  $\widehat{u}_{\text{in}}(K_x, 0), \widehat{u}'_{\text{in}}(K_x, 0) \in \mathcal{S}'$ . From (5.1.4) we see that  $1/s_-$  becomes singular at the Fourier frequencies  $|K_x| = k$ . On the other hand, by definition of the incoming field and the bound  $|k_x| < k$  we find  $\lim_{|K_x| \rightarrow k} \widehat{u}'_{\text{in}}(K_x, 0)/s_-(K_x) = 0$ . Thus  $\widehat{u}(K_x, 0)$  computed using (5.1.6) is in  $\mathcal{S}'$  and the inverse Fourier transform exists. Computing the inverse Fourier transform of (5.1.6) yields

$$\begin{aligned} u(x, 0) &= u_{\text{in}}(x, 0) + F^{-1} \left( \frac{-ik_y \delta(K_x - k_x)}{-i\sqrt{k^2 - K_x^2}} \right) \\ &= u_{\text{in}}(x, 0) + \frac{k_y}{\sqrt{k^2 - k_x^2}} e^{ik_x x} \\ &= u_{\text{in}}(x, 0) + u_{\text{in}}(x, 0). \end{aligned}$$

Thus we obtain the boundary condition for the superposed field

$$u(x, 0) = 2u_{\text{in}}(x, 0), \quad \text{or equivalently,} \quad \widehat{u}(K_x, 0) = 2\widehat{u}_{\text{in}}(K_x, 0).$$

Once this boundary condition is known, the corresponding field can be computed directly. Using the solution representation (5.1.3), the imposed boundary condition  $\widehat{u}'(K_x, 0) = 0$  and the just obtained boundary condition we find

<sup>1</sup>The minus sign occurs since the derivative refers to  $y$ .

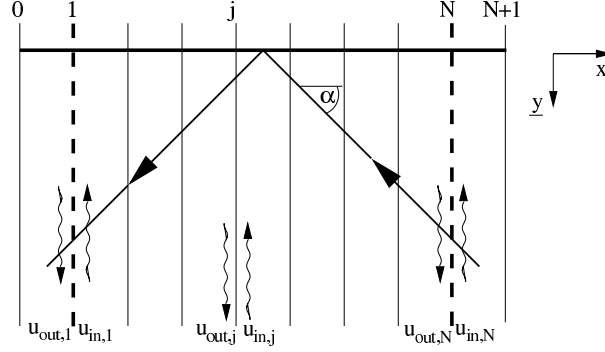


FIGURE 5.1.2. Reflection from an infinite plane: semi-discrete approximation

$$\widehat{u}(K_x, s) = \frac{2s\widehat{u}_{\text{in}}(K_x, 0)}{s^2 + k^2 - K_x^2}.$$

Partial fraction decomposition results in

$$\widehat{u}(K_x, s) = \left( \frac{1}{s - s_-} + \frac{1}{s + s_-} \right) \widehat{u}_{\text{in}}(K_x, 0)$$

which supplies, after inverse Laplace and Fourier transform,

$$(5.1.7) \quad u(x, y) = e^{i(k_x x + k_y y)} + e^{i(k_x x - k_y y)}.$$

**REMARK.** To derive (5.1.7) we have imposed an *angle condition* with respect to the incident field: the angle between the *wavevector* of the incident field and the reflecting *plane* must be larger than zero. In the limiting case  $\alpha \rightarrow 0$ , i. e.  $\tilde{k}_x \rightarrow k$  and  $k_y \rightarrow 0$ , we get from (5.1.7)  $u(x, y) = 2\exp(ikx)$ . On the other hand, the incident plane wave  $u_{\text{in}}(x, y) = \exp(ikx)$  itself is a solution of the Helmholtz equation. Further *any* function  $u(x, y) = \text{const} \cdot \exp(ikx)$  solves the reflection problem with the imposed boundary condition. Hence the solution does not depend continuously on the incident angle if  $\alpha \rightarrow 0$ , which justifies the angle condition introduced above. In the semi-discrete case discussed in the next section we will find corresponding conditions which ensures the existence of the discrete solution.

**Semi-discrete Case.** We repeat the procedure above *without* using a Fourier transform. Instead, we restrict the infinite plane to a finite (but long) interval  $[-a, a]$ . The situation is schematically depicted in Fig. 5.1.2.

Our approach is as follows:

- (1) Laplace transform of the semi-discrete problem with respect to  $y$ .
- (2) Discretization in  $x$ -direction.
- (3) The pole condition with respect to the solution on each ray is traced back to a pole condition with respect to eigenmodes of the discrete operator.

Note that the first two points might be interchanged. We consider a finite element (FE) as well as a finite difference (FD) implementation. Let us assume for simplicity that the interval is divided into  $N + 1$  uniform subintervals. Further let us assume that the fields on the vertical boundaries  $u_0(y)$  and  $u_{N+1}(y)$  in case of finite differences or  $\partial_x u_1(y)$  and  $\partial_x u_N(y)$  in case of finite elements and the incident fields  $u_{\text{in},j}(y)$ ,  $j = 1, \dots, N$ , are given. First we perform a Laplace transform of the Helmholtz equation in positive  $y$ -direction, using  $y$  as distance variable as before:

$$(5.1.8) \quad \partial_x^2 \widehat{u}(x, s) + s^2 \widehat{u}(x, s) + k^2 \widehat{u}(x, s) = s u(x, 0) + u'(x, 0).$$

Next we want to discretize this equation in  $x$ -direction by the standard finite difference and finite element methods. To this end we introduce the vector of unknown functions with elements  $\widehat{u}_j(s) = \widehat{u}(x_j, s)$ ,  $j = 0, \dots, N+1$ , and the vectors containing the Laplace transformed functions along the rays and the corresponding initial values at  $\underline{y} = 0$

$$\widehat{\mathbf{u}}(s) = \begin{pmatrix} \widehat{u}_1(s) \\ \vdots \\ \widehat{u}_N(s) \end{pmatrix}, \quad \mathbf{u}_0 = \begin{pmatrix} u_1(0) \\ \vdots \\ u_N(0) \end{pmatrix}, \quad \mathbf{u}'_0 = \begin{pmatrix} u'_1(0) \\ \vdots \\ u'_N(0) \end{pmatrix}.$$

A discretization results in the algebraic system<sup>2</sup>, depending on the complex parameter  $s$

$$(5.1.9) \quad (\mathbf{A} + s^2\mathbf{M}) \widehat{\mathbf{u}}(s) = \widehat{\mathbf{u}}_b(s) + s\mathbf{M}\mathbf{u}_0 + \mathbf{M}\mathbf{u}'_0.$$

In case of a finite element discretization the symmetric matrices  $\mathbf{A} \in \mathbb{R}^{N \times N}$  and  $\mathbf{M} \in \mathbb{R}^{N \times N}$  are computed based on the test and trial space  $V_h = \text{span}\{v_1, \dots, v_N\}$  with continuous, piecewise linear functions  $v_j$ , where we set  $v_1(-a) = 1$  and  $v_N(a) = 1$ . The finite element matrices  $\mathbf{A}$  and  $\mathbf{M}$  are computed by  $(\mathbf{A})_{ij} = -(\partial_x v_i, \partial_x v_j) + k^2(v_i, v_j)$  and  $(\mathbf{M})_{ij} = (v_i, v_j)$ , the corresponding standard finite difference matrices are considered in Example 5.1.3. Note that  $\mathbf{M}$  is symmetric positive definite and becomes the identity matrix in case of a finite difference discretization. The vertical boundary vector  $\widehat{\mathbf{u}}_b \in \mathbb{C}^N$  results from the boundary terms  $v(\pm a)\partial_x \widehat{u}(\pm a, s)$  of the integration by parts in case of the finite element formulation or from the leftmost and the rightmost contribution of the three-point finite difference stencil in case of the finite difference formulation

$$(5.1.10) \quad \text{FE} : \widehat{\mathbf{u}}_b(s) := \begin{pmatrix} -\partial_x \widehat{u}_1(s) \\ \mathbf{0} \\ \partial_x \widehat{u}_N(s) \end{pmatrix}, \quad \text{FD} : \widehat{\mathbf{u}}_b(s) := -\frac{1}{h^2} \begin{pmatrix} \widehat{u}_0(s) \\ \mathbf{0} \\ \widehat{u}_{N+1}(s) \end{pmatrix}.$$

Eq. (5.1.9) is the semi-discrete analogue to (5.1.2), completed with the influence of the prescribed vertical boundaries.

We will formulate the boundary condition based on the spectral properties of the discrete operator  $\mathbf{A} + s^2\mathbf{M}$  which appears on the left-hand side of (5.1.9). Therefore we consider, in preparation of the derivation of the boundary condition, the generalized, real symmetric eigenvalue problem

$$(5.1.11) \quad \mathbf{A}\mathbf{u}_j = \lambda_j^2 \mathbf{M}\mathbf{u}_j$$

with eigenvalues  $\lambda_j^2$  and eigenvectors  $\mathbf{u}_j$ . The eigenvalues are real and satisfy  $\lambda_j^2 \leq k^2$ . The eigenvectors  $\mathbf{u}_j$  can be normalized to vectors  $\mathbf{u}_j^0$  such that

$$\mathbf{U}^H \mathbf{M} \mathbf{U} = \mathbf{I}, \quad \text{with} \quad \mathbf{U} = (\mathbf{u}_1^0, \dots, \mathbf{u}_N^0).$$

By means of  $\mathbf{U}$  we can transform the matrix  $\mathbf{A}$  to diagonal form

$$(5.1.12) \quad \mathbf{U}^H \mathbf{A} \mathbf{U} = \Lambda^2, \quad \Lambda^2 = \text{diag}(\lambda_j^2).$$

Let the square root of  $\lambda_j^2$  be defined as usual

---

<sup>2</sup>In this section we denote vectors arising from the discretization of continuous functions by bold roman letters. This simplifies, e. g., to distinguish between  $u$  as a function and  $\mathbf{u}$  as a vector.

$$\lambda_j = \begin{cases} \sqrt{\lambda^2} & \text{if } \lambda^2 > 0 \\ i\sqrt{-\lambda^2} & \text{if } \lambda^2 < 0 \\ 0 & \text{if } \lambda^2 = 0. \end{cases}$$

Further, let  $S_{\pm} = \pm i\Lambda$  with  $\Lambda = \text{diag}(\lambda_j)$ . Accordingly, we define the sets  $\sigma_{\pm} = \{\pm i\lambda_1, \dots, \pm i\lambda_N\}$ . Note that  $\sigma_- \cap \sigma_+$  is either empty or contains the zero eigenvalue.

Now we are able to derive the desired boundary condition. By the transform (5.1.12) we rewrite the semi-discrete, Laplace transformed Helmholtz equation (5.1.9) to

$$(\mathbf{U}^H \mathbf{A} \mathbf{U} + s^2 \mathbf{U}^H \mathbf{M} \mathbf{U}) \mathbf{U}^{-1} \hat{\mathbf{u}}(s) = \mathbf{U}^H (\hat{\mathbf{u}}_b(s) + s \mathbf{M} \mathbf{u}_0 + \mathbf{M} \mathbf{u}'_0)$$

which yields

$$(5.1.13) \quad (\Lambda^2 + s^2 \mathbf{I}) \mathbf{U}^H \mathbf{M} \hat{\mathbf{u}}(s) = \mathbf{U}^H (\hat{\mathbf{u}}_b(s) + s \mathbf{M} \mathbf{u}_0 + \mathbf{M} \mathbf{u}'_0).$$

We factorize the operator on the left-hand side

$$\Lambda^2 + s^2 \mathbf{I} = (\mathbf{I}s - S_-)(\mathbf{I}s + S_-).$$

Next we consider, as in the continuous case, only the outgoing (reflected) field resulting from data  $\hat{\mathbf{u}}_{b,\text{out}}(s)$  of the vertical boundary and the corresponding initial data on the plane  $\mathbf{u}_{0,\text{out}}, \mathbf{u}'_{0,\text{out}}$ . The outgoing field has to satisfy the basic semi-discrete equation (5.1.9). To avoid a singularity of the expression  $(\mathbf{I}s + S_-) \mathbf{U}^H \mathbf{M} \hat{\mathbf{u}}_{\text{out}}(s)$  at one of the diagonal entries of  $S_-$ , i.e. for  $s \in \sigma_-$ , and equivalently of  $\hat{\mathbf{u}}_{\text{out}}(s)$  at these points, we require that the right-hand side of (5.1.13) lies in the orthogonal complement of the nullspace of the transposed operator  $\mathbf{I}s - S_-$ :

$$\mathbf{U}^H (\hat{\mathbf{u}}_{b,\text{out}}(s) + s \mathbf{M} \mathbf{u}_{0,\text{out}} + \mathbf{M} \mathbf{u}'_{0,\text{out}}) \in \mathcal{N}_E^{\perp}(\mathbf{I}s - S_-) \quad \text{for } s \in \sigma_-.$$

Here, and in the following we use the notation:  $N^{\perp}(A)$  denotes the orthogonal complement of the null space of a matrix  $A$  with respect to  $\langle \cdot, \mathbf{M} \cdot \rangle$  and  $\mathcal{N}_E^{\perp}(A)$  denotes the orthogonal complement of the null space of a matrix  $A$  with respect to the Euclidean metric  $\langle \cdot, \cdot \rangle$ . Since  $S_-$  is a diagonal matrix, the null space of  $(\mathbf{I}s_j - S_-)$ , with  $s_j := -i\lambda_j \in \sigma_-$ , is spanned by a single Euclidean unit vector  $(\mathbf{e}_j)_{i=1,\dots,N} = \delta_{ij}$ . Thus, with  $\mathbf{u}_j^0$  the  $j$ th normalized eigenvector from (5.1.11), the above pole condition may be given alternatively as

$$(5.1.14) \quad \langle \mathbf{u}_j^0, \hat{\mathbf{u}}_{b,\text{out}}(s_j) + s_j \mathbf{M} \mathbf{u}_{0,\text{out}} + \mathbf{M} \mathbf{u}'_{0,\text{out}} \rangle = 0, \quad s_j \in \sigma_- \quad \text{and} \quad j = 1, \dots, N.$$

This equation supplies  $N$  conditions. Defining the vector  $\mathbf{b}_{v,\text{out}}$  with elements

$$(5.1.15) \quad (\mathbf{b}_{v,\text{out}})_j := \langle \mathbf{u}_j^0, \hat{\mathbf{u}}_{b,\text{out}}(s_j) \rangle, \quad j = 1, \dots, N$$

and taking into account  $\mathbf{U} \mathbf{U}^H \mathbf{M} = \mathbf{I}$  we can rewrite this in compact form

Pole condition for the semi-discrete reflection problem

$$(5.1.16) \quad \mathbf{0} = \mathbf{U} \mathbf{b}_{v,\text{out}} + \mathbf{U} S_- \mathbf{U}^H \mathbf{M} \mathbf{u}_{0,\text{out}} + \mathbf{u}'_{0,\text{out}}.$$



Eq. (5.1.16) is the semi-discrete analogue to condition (5.1.5) for the continuous problem. It avoids that the outgoing (reflected) field with boundary data  $\mathbf{u}_{0,\text{out}}, \mathbf{u}'_{0,\text{out}}$  possesses poles which belong to incoming or increasing functions.

Next we compute the field caused by an incident field  $\mathbf{u}_{\text{in}}$ . Let  $\hat{u}_{\text{in},j}(s) = \hat{u}_{\text{in}}(x_j, s)$  be given. This means, we require the knowledge of the incident field along the vertical rays, cf. Fig. 5.1.2. As we will see, it does not play a role whether the incident field is a true incoming field in the sense of our definition of incoming and outgoing fields or a general source field containing both incoming and outgoing parts. We define the vector of Laplace transformed incoming functions along the rays, and the corresponding Dirichlet and Neumann data on the reflecting plane by

$$\hat{\mathbf{u}}_{\text{in}}(s) := \begin{pmatrix} \hat{u}_{\text{in},1}(s) \\ \vdots \\ \hat{u}_{\text{in},N}(s) \end{pmatrix}, \quad \mathbf{u}_{\text{in},0} := \begin{pmatrix} u_{\text{in},1}(0) \\ \vdots \\ u_{\text{in},N}(0) \end{pmatrix}, \quad \mathbf{u}'_{\text{in},0} := \begin{pmatrix} u'_{\text{in},1}(0) \\ \vdots \\ u'_{\text{in},N}(0) \end{pmatrix}.$$

Additionally we introduce the vectors  $\mathbf{b}_v, \mathbf{b}_{v,\text{in}}$  in analogy to (5.1.15) by

$$(5.1.17) \quad \mathbf{b}_v := \begin{pmatrix} \langle \mathbf{u}_1^0, \hat{\mathbf{u}}_b(s_1) \rangle \\ \vdots \\ \langle \mathbf{u}_N^0, \hat{\mathbf{u}}_b(s_N) \rangle \end{pmatrix}, \quad \mathbf{b}_{v,\text{in}} := \begin{pmatrix} \langle \mathbf{u}_1^0, \hat{\mathbf{u}}_{b,\text{in}}(s_1) \rangle \\ \vdots \\ \langle \mathbf{u}_N^0, \hat{\mathbf{u}}_{b,\text{in}}(s_N) \rangle \end{pmatrix}.$$

Based on the superposition of incoming and outgoing fields we rewrite (5.1.16) in terms of the complete and the incoming field

$$(5.1.18) \quad \mathbf{0} = \mathbf{U}(\mathbf{b}_v - \mathbf{b}_{v,\text{in}}) + \mathbf{U}\mathbf{S}_-\mathbf{U}^{\text{HM}}(\mathbf{u}_0 - \mathbf{u}_{0,\text{in}}) + \mathbf{U}\mathbf{U}^{\text{HM}}(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}).$$

We want to discuss the unique solvability of our reflection problem (5.1.9) together with condition (5.1.18), where we consider as given:

- (1) The data  $\mathbf{u}_{0,\text{in}}, \mathbf{u}'_{0,\text{in}}$  on the reflecting plane determined by the incoming field.
- (2) The Neumann boundary condition  $\mathbf{u}'_0(x)$  on the reflecting plane.
- (3) The incoming fields  $\hat{\mathbf{u}}_{b,\text{in}}(s)$  on the vertical boundaries, see (5.1.10).
- (4) The complete field  $\hat{\mathbf{u}}_b(s)$  (superposed from the incoming and the outgoing fields) on the vertical boundaries.

First of all note that the semi-discrete scattering problem (5.1.9) has a unique solution in the spatial domain for all given data  $\mathbf{u}_0, \mathbf{u}'_0$ , and  $\mathbf{u}_b$ , as long as their Laplace transform exists. Thus the remaining question is under which circumstances we can compute a unique function  $\mathbf{u}_0$  such that the reflected wave is an outgoing wave. We summarize these conditions in the following

LEMMA 5.1.1. *Let the given data have the following properties:*

- (1) *Horizontal boundary: Let  $(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}) \in \mathcal{N}^\perp(A)$ .*
- (2) *Vertical boundary: Let  $(\hat{\mathbf{u}}_b(0) - \hat{\mathbf{u}}_{b,\text{in}}(0)) \in \mathcal{N}_E^\perp(A)$*

*Then, the following holds true :*

*The boundary condition (5.1.18) has a unique solution vector  $\mathbf{u}_0$  and the reflected wave is an outgoing wave.*

PROOF. Let  $0 \in \sigma_-$ , otherwise the statement is trivial. Obviously, (5.1.18) has a unique solution if both

$$(a) \quad \mathbf{U}^{\text{HM}}(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}) \in \mathcal{N}_E^\perp(S_-) \quad \text{and} \quad (b) \quad (\mathbf{b}_v - \mathbf{b}_{v,\text{in}}) \in \mathcal{N}_E^\perp(S_-).$$

Let the subscript  $j_0$  indicate elements of nullspaces. Let  $\mathbf{e}_{j_0} \in \mathcal{N}(S_-)$ . Hence  $\mathbf{u}_{j_0} := U\mathbf{e}_{j_0} \in \mathcal{N}(A)$ . It follows directly from our assumption (1) written as

$$0 = \langle u_{j_0}, M(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}) \rangle$$

that condition (a) is met

$$\begin{aligned} 0 &= \langle U\mathbf{e}_{j_0}, M(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}) \rangle \\ &= \langle \mathbf{e}_{j_0}, U^H M(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}) \rangle. \end{aligned}$$

Recall the definition of  $\mathbf{b}_v$  and  $\mathbf{b}_{v,\text{in}}$  in (5.1.17):

$$(\mathbf{b}_v)_j := \langle \mathbf{u}_j^0, \hat{\mathbf{u}}_b(s_j) \rangle, \quad (\mathbf{b}_{v,\text{in}})_j := \langle \mathbf{u}_j^0, \hat{\mathbf{u}}_{b,\text{in}}(s_j) \rangle, \quad j = 1, \dots, N.$$

By  $s_{j_0} = -i\lambda_{j_0} = 0$  it follows

$$\begin{aligned} \langle \mathbf{u}_{j_0}^0, \hat{\mathbf{u}}_b(s_{j_0}) \rangle &= \langle \mathbf{u}_{j_0}^0, \hat{\mathbf{u}}_b(0) \rangle \\ &= (\mathbf{b}_v)_{j_0}. \end{aligned}$$

Obviously we have also

$$\langle \mathbf{u}_{j_0}^0, \hat{\mathbf{u}}_b(0) - \hat{\mathbf{u}}_{b,\text{in}}(0) \rangle = (\mathbf{b}_v)_{j_0} - (\mathbf{b}_{v,\text{in}})_{j_0}$$

Hence it follows from assumption (2),  $(\hat{\mathbf{u}}_b(0) - \hat{\mathbf{u}}_{b,\text{in}}(0)) \in \mathcal{N}_E^\perp(A)$ , that  $(\mathbf{b}_v)_{j_0} - (\mathbf{b}_{v,\text{in}})_{j_0} = 0$ . Thus  $\langle \mathbf{b}_v - \mathbf{b}_{v,\text{in}}, \mathbf{e}_{j_0} \rangle = 0$ . This is condition (b). The boundary condition (5.1.16) determines the properties of the reflected field. The condition excludes all modes with poles  $\text{Im } s_j < 0$ . Thus the reflected field must be an outgoing field in the sense of our definition.  $\square$

The lemma tells us that the outgoing field has to satisfy two independent conditions, one concerning the horizontal and the other concerning the vertical boundaries. However, based the orthogonality condition (5.1.14) we can put these into a single condition using the notation of Lemma 5.1.1:

**COROLLARY 5.1.2.** *The boundary condition (5.1.18) has unique solution vector  $\mathbf{u}_0$  if*

$$\langle \mathbf{u}_{j_0}^0, \hat{\mathbf{u}}_b(0) - \hat{\mathbf{u}}_{b,\text{in}}(0) + M(\mathbf{u}'_0 - \mathbf{u}'_{0,\text{in}}) \rangle = 0 \quad \text{for all } j_0 \text{ with } \lambda_{j_0}^2 = 0$$

and the corresponding reflected wave is an outgoing wave.

The boundary condition (5.1.16) is the general semi-discrete boundary condition for the reflection from the (restricted) infinite plane, which we want to keep in this form, in view of its application to the scattering from a semi-infinite plane. We get the semi-discrete analogue to the continuous condition (5.1.6) if we set  $\mathbf{u}' = 0$ , as required in the problem formulation.

**EXAMPLE 5.1.3.** Let us apply the standard uniform finite-difference method to derive (5.1.9) to the case of rigid scattering  $u'(x, 0) = 0$ . We take  $x_j$ ,  $j = 0, \dots, N+1$ , such that  $h = x_{k+1} - x_k$ ,  $i = 0, \dots, N$  and  $x_1 = -a$ ,  $x_N = a$  and obtain the real symmetric matrices  $A, M \in \mathbb{R}^{N \times N}$

$$(5.1.19) \quad A := \frac{1}{h^2} \begin{pmatrix} -2 & 1 & & & \\ 1 & -2 & 1 & & \\ & & \ddots & \ddots & \\ & & & 1 & -2 \end{pmatrix} + k^2 \mathbf{I}, \quad M = \mathbf{I}.$$

Let the incoming field be a plane wave

$$u_{\text{in}}(x_j, \underline{y}) = e^{ik_x x_j + ik_y \underline{y}},$$

with real wavenumbers  $k_x$  and  $k_y$ ,  $k_y^2 < k^2$ . We set  $z = \exp(ik_x h)$  such that the source field takes the form  $u_{\text{in}}(x_j, \underline{y}) = z^j \exp(ik_y \underline{y})$ . Since  $k_x$  is real we find  $|z| = 1$ . We require that the incoming field satisfies the semi-discrete Helmholtz equation on the half-plane  $\underline{y} \geq 0$  and  $j = 1, \dots, N$ :

$$(5.1.20) \quad \frac{1}{h^2} (u_{\text{in}}(x_{j-1}, \underline{y}) - 2u_{\text{in}}(x_j, \underline{y}) + u_{\text{in}}(x_{j+1}, \underline{y})) + \partial_y^2 u_{\text{in}}(x_j, \underline{y}) + k^2 u_{\text{in}}(x_j, \underline{y}) = 0.$$

Replacing  $u_{\text{in}}(x_j, \underline{y})$  by the above ansatz, the semi-discrete dispersion relation

$$(5.1.21) \quad z^2 + z(-2 + h^2(-k_y^2 + k^2)) + 1 = 0$$

follows. We assume that the discretization step width  $h$  and the wavenumber  $k_y$  of the incident plane wave are given. Then, the two roots of the semi-discrete dispersion relation supply the two possible values for the transversal wavenumber  $k_x$ . We discuss this in brief. For

$$h < \frac{2}{\sqrt{k^2 - k_y^2}}$$

the zeros  $z_{1,2}$  of (5.1.21) have an imaginary part. Since the polynomial is real, they must be complex conjugates. Then, by Vieta's theorem  $z_1 z_2 = 1$ , hence  $|z_1| = |z_2| = 1$  as required. The expressions  $z_1^j, z_2^j$  are discrete periodic functions. The field

$$(5.1.22) \quad u_{\text{in}}(x_j, \underline{y}) = e^{\pm i(\arg z_1)j + ik_y \underline{y}}$$

obeys the semi-discrete dispersion relation and we find  $k_x = \pm(\arg z_1)/h$ .

In order to follow the treatment of this section, namely the formulation of the boundary condition in terms of the eigenvalues and eigenvectors of the discrete operator, we compute first all eigenvectors  $\mathbf{u}_j$  and corresponding eigenvalues  $\lambda_j^2$  of the eigenvalue problem

$$(5.1.23) \quad \mathbf{A} \mathbf{u}_j = \lambda_j^2 \mathbf{u}_j, \quad \text{with } \mathbf{u}_j = e^{i\lambda_j \underline{y}} (u_j(x_1), \dots, u_j(x_N))^T$$

and set  $\sigma(\mathbf{A}) = \{\lambda_1^2, \dots, \lambda_N^2\}$ . To solve the eigenvalue problem we insert the ansatz  $(\mathbf{v}_j)_l = z^l \exp(i\lambda_j \underline{y})$ ,  $l = 1, \dots, N$ , into the semi-discrete Helmholtz equation (5.1.20) and compute the zeros  $z_{1,2}$  of the semi-discrete dispersion relation

$$(5.1.24) \quad z^2 + z(-2 + h^2(k^2 - \lambda_j^2)) + 1 = 0.$$

We use these zeros to compose a vector  $(\widetilde{\mathbf{u}}_j)_{l=0, \dots, N+1} = bz_1^l + cz_2^l$ , with complex coefficients  $b, c$ . The restricted vector  $(\widetilde{\mathbf{u}}_j)_{l=1, \dots, N}$  constructed this way solves the eigenvalue problem (5.1.23) with arbitrary coefficients  $b, c$  except for the first and the last row of the system. The additional conditions  $(\widetilde{\mathbf{u}}_j)_0 = 0$  and  $(\widetilde{\mathbf{u}}_j)_{N+1} = 0$  then guaranty that  $(\widetilde{\mathbf{u}}_j)_{l=1, \dots, N} = (\mathbf{u}_j)_{l=1, \dots, N}$  where  $\mathbf{u}_j$  is in fact the desired eigenvector since it satisfies the eigenvalue equation (5.1.23). These additional conditions together with the properties  $|z_1| = |z_2| = 1$  and  $z_1 = \bar{z}_2$  yield

$$\begin{aligned} b &= -c \\ z_{1,2}^{2(N+1)} &= 1. \end{aligned}$$

The last equation has  $N + 1$  different zeros:  $z_{1,2}^{(j)} = \exp(\pm\pi ij/(N + 1))$ ,  $j = 0, \dots, N$ . Ignoring the constant solution corresponding to  $j = 0$ , which only yields the trivial eigenvector due to the vanishing boundary data, we find  $N$  periodic eigenfunctions and  $N$  corresponding eigenvalues

$$(5.1.25) \quad \begin{aligned} (\mathbf{u}_j)_{l=1,\dots,N} &= \sin\left(j\frac{l\pi}{N+1}\right) \\ \lambda_j^2 &= k^2 + \frac{2}{h^2} \left( \cos\left(j\frac{\pi}{N+1}\right) - 1 \right), \quad j = 1, \dots, N. \end{aligned}$$

From the eigenvalue formula (5.1.25) it becomes clear that the nullspace of  $\mathbf{A}$  in general is empty. Then, the unique solvability of the boundary condition (5.1.18) is trivial and Lemma 5.1.1 tells us that the reflected wave is an outgoing wave.

However, we may construct a situation where the nullspace of the matrix  $\mathbf{A}$  is non-empty. Let us analyze such a situation in the following. First we choose the parameter  $h$  such that there exists a zero-eigenvalue. We fix an integer  $j_0$ ,  $1 \leq j_0 \leq N$ , and set in the eigenvalue equation (5.1.25)  $\lambda_{j_0}^2 = 0$ , which results in a condition with respect to  $h$

$$h = \frac{1}{k} \sqrt{2 \left( 1 - \cos \frac{j_0 \pi}{N+1} \right)}.$$

The corresponding eigenvector  $\mathbf{u}_{j_0}$  meets

$$\mathbf{A}\mathbf{u}_{j_0} = \mathbf{0}.$$

A comparison of the semi-discrete dispersion relation (5.1.21) which yields wavenumbers  $k_y^2$  such that the corresponding semi-discrete periodic functions satisfy the semi-discrete Helmholtz equation (5.1.20) and the dispersion relation (5.1.24) which yields the eigenvalues  $\lambda_j^2$  of the eigenvalue problem (5.1.23) shows that only wavenumbers  $k_y^2$  are admissible, which are in the spectrum of  $\mathbf{A}$ . Otherwise the incident wave  $u_{\text{in}}(x_j, \underline{y})$  is *not* a discrete periodic function in  $x$ . Note that the vector  $\mathbf{u}_{0,\text{in}}$  defined by  $(\mathbf{u}_{0,\text{in}})_j = u_{\text{in}}(x_j, 0)$ ,  $j = 1, \dots, N$ , is not an eigenfunction of the eigenvalue problem (5.1.23), despite of  $k_y^2 \in \sigma(\mathbf{A})$ . This follows from the fact that the incidence wave has a representation  $z_1^j$ , whereas the eigenvector has a representation  $z_1^j - z_1^{-j}$ . Further we see from the eigenvalue formula (5.1.25) that  $\lambda_j^2 < 0$  for  $j > j_0$ , hence  $\text{Im}(\lambda_j) \neq 0$  for  $j > j_0$ . Since  $k_y^2 \in \sigma(\mathbf{A})$  we find a necessary restriction for the possible choices of  $k_y^2$ , namely  $k_y^2 = \lambda_j^2$  for some  $j$  with  $1 \leq j < j_0$ . Otherwise the incident wave would be an exponential increasing or decreasing function.

After this preparation, we can analyze the solvability of our example problem in analogy to the continuous problem. Let  $k_y \in \mathbb{R}$ ,  $k_y^2 \in \sigma(\mathbf{A})$ ,  $0 < |k_y| < k$ , let  $k_x$  be computed according to (5.1.22). Further let the vertical boundary functions are the exact functions known from the continuous case

$$\begin{aligned} u_{\text{in}}(\pm a, \underline{y}) &= e^{i(\pm ak_x - k_y \underline{y})} \\ u(\pm a, \underline{y}) &= e^{i(\pm ak_x + k_y \underline{y})} + e^{i(\pm ak_x - k_y \underline{y})}. \end{aligned}$$

with Laplace transforms

$$\begin{aligned}\widehat{u}_{\text{in}}(\pm a, s) &= \frac{1}{s + ik_y} e^{\pm ik_x a} \\ \widehat{u}(\pm a, s) &= \frac{1}{s - ik_y} e^{\pm ik_x a} + \frac{1}{s + ik_y} e^{\pm ik_x a}.\end{aligned}$$

Further we compute

$$(5.1.26) \quad \widehat{\mathbf{u}}_{\text{b,in}}(0) = \frac{1}{ik_y} \begin{pmatrix} e^{-ik_x a} \\ \mathbf{0} \\ e^{ik_x a} \end{pmatrix}, \quad \widehat{\mathbf{u}}_{\text{b}}(0) = \frac{1}{ik_y} \begin{pmatrix} e^{-ik_x a} \\ \mathbf{0} \\ e^{ik_x a} \end{pmatrix} - \frac{1}{ik_y} \begin{pmatrix} e^{-ik_x a} \\ \mathbf{0} \\ e^{ik_x a} \end{pmatrix} = \mathbf{0}.$$

We have to show that these function meet the condition of Corollary 5.1.2. With (5.1.26) we can express this condition as

$$(5.1.27) \quad \langle \mathbf{u}'_{j_0}, \widehat{\mathbf{u}}_{\text{b,in}}(0) + \mathbf{u}'_{0,\text{in}} \rangle = 0.$$

To prove the condition, we begin with the fact that the incoming field  $\mathbf{u}_{\text{in}}(\underline{y})$  must satisfy the spatial counterpart of the semi-discrete Laplace transformed Helmholtz equation (5.1.9) for all  $\underline{y} \geq 0$ , see also this equation in component-wise form (5.1.20). We restrict this to  $\underline{y} = 0$ , set again  $\mathbf{u}_{0,\text{in}} := \mathbf{u}_{\text{in}}(0)$  and obtain

$$(5.1.28) \quad \mathbf{A}\mathbf{u}_{0,\text{in}} - k_y^2 \mathbf{u}_{0,\text{in}} = \mathbf{u}_{\text{b,in}}.$$

Further, it follows from (5.1.26)

$$\widehat{\mathbf{u}}_{\text{b,in}}(0) = \frac{1}{ik_y} \mathbf{u}_{\text{b,in}}.$$

Replacing  $\mathbf{u}_{\text{b,in}}$  in (5.1.28) by  $\widehat{\mathbf{u}}_{\text{b,in}}(0)$  using the last equation and taking the inner product with  $\mathbf{u}_{j_0}$  it follows

$$\langle \mathbf{u}_{j_0}, \mathbf{A}\mathbf{u}_{0,\text{in}} \rangle - k_y^2 \langle \mathbf{u}_{j_0}, \mathbf{u}_{0,\text{in}} \rangle = ik_y \langle \mathbf{u}_{j_0}, \widehat{\mathbf{u}}_{\text{b,in}}(0) \rangle.$$

Since  $\mathbf{u}_{j_0} \in \mathcal{N}(\mathbf{A})$  and  $\mathbf{A} = \mathbf{A}^T$  we get

$$ik_y \langle \mathbf{u}_{j_0}, \widehat{\mathbf{u}}_{\text{b,in}}(0) \rangle + k_y^2 \langle \mathbf{u}_{j_0}, \mathbf{u}_{0,\text{in}} \rangle = 0.$$

Finally we use  $\mathbf{u}'_{0,\text{in}} = -ik_y \mathbf{u}_{0,\text{in}}$ , which holds true by the definition of the semi-discrete plane wave, to find

$$ik_y \langle \mathbf{u}_{j_0}, \widehat{\mathbf{u}}_{\text{b,in}}(0) \rangle + k_y^2 \left\langle \mathbf{u}_{j_0}, \frac{1}{-ik_y} \mathbf{u}'_{0,\text{in}} \right\rangle = 0.$$

Since  $k_y \neq 0$  this shows (5.1.27). Hence the discrete problem has exactly one solution even in the case when the nullspace of  $\mathbf{A}$  is non-empty, if the incidence wave is a discrete periodic wave and  $k_y \neq 0$ . This coincides exactly with the angle condition in the continuous case.

We show by a direct computation that  $\widehat{\mathbf{u}}(s) = \widehat{\mathbf{u}}_{\text{in}}(s) + \widehat{\mathbf{u}}_{\text{out}}(s)$  with

$$(5.1.29) \quad \begin{aligned}(\widehat{\mathbf{u}}_{\text{in}})_j &= \frac{1}{s + ik_y} e^{ik_x x_j}, & (\mathbf{u}_{0,\text{in}})_j &= e^{ik_x x_j}, & (\mathbf{u}'_{0,\text{in}})_j &= -k_y (\mathbf{u}_{\text{in},0})_j \\ (\widehat{\mathbf{u}}_{\text{out}})_j &= \frac{1}{s - ik_y} e^{ik_x x_j}, & (\mathbf{u}_{0,\text{out}})_j &= e^{ik_x x_j}, & (\mathbf{u}'_{0,\text{out}})_j &= k_y (\mathbf{u}_{\text{out},0})_j\end{aligned}$$

( $0 \leq j \leq N + 1$ ) satisfies in fact the semi-discrete, Laplace transformed Helmholtz equation

$$(5.1.30) \quad (\mathbf{A} + s^2) \hat{\mathbf{u}}(s) = \hat{\mathbf{u}}_b(s) + s\mathbf{u}_0 + \mathbf{u}'_0.$$

Let us consider the  $j$ th row of the system. Using  $z^j = \exp(ik_x x_j)$  it follows

$$\begin{aligned} (\hat{\mathbf{u}}_{\text{in}})_j + (\hat{\mathbf{u}}_{\text{out}})_j &= \left( \frac{1}{s + ik_y} + \frac{1}{s - ik_y} \right) z^j \\ &= \frac{2s}{s^2 + k_y^2} z^j. \end{aligned}$$

Further recall that by the definition of  $\hat{\mathbf{u}}_b(s)$  in (5.1.10) and  $\mathbf{A}$  in (5.1.19)

$$(\mathbf{A}\hat{\mathbf{u}}(s) - \hat{\mathbf{u}}(s))_j = \left[ \frac{1}{h^2} (z^{j-1} - 2z^j + z^{j+1}) + k^2 z^j \right] \frac{2s}{s^2 + k_y^2}, \quad 1 \leq j \leq N.$$

Thus the left-hand side of (5.1.30) minus  $\hat{\mathbf{u}}(s)$  reads

$$(\mathbf{A}\hat{\mathbf{u}}(s) - \hat{\mathbf{u}}(s) + s^2 \hat{\mathbf{u}}(s))_j = \left[ \frac{1}{h^2} (z^{j-1} - 2z^j + z^{j+1}) + k^2 z^j + s^2 \right] \frac{2s}{s^2 + k_y^2}.$$

Taking into account that  $\mathbf{u}'_0 = 0$  by the imposed boundary condition we rewrite (5.1.30) to

$$\begin{aligned} \left[ \frac{1}{h^2} (z^{j-1} - 2z^j + z^{j+1}) + k^2 z^j + s^2 z^j \right] \frac{2s}{s^2 + k_y^2} &= s(\mathbf{u}_0)_j \\ &= 2sz^j. \end{aligned}$$

This results in

$$\frac{1}{h^2} (z^{-1} - 2 + z) + k^2 + s^2 = s^2 + k_y^2$$

which is obviously identical to the semi-discrete dispersion relation (5.1.21), hence it is satisfied. Thus the functions (5.1.29) obey the semi-discrete scattering problem. An inverse Laplace transform supplies the corresponding spatial solutions:

$$u(x_j, \underline{y}) = e^{i(k_x x_j + k_y \underline{y})} + e^{i(k_x x_j - k_y \underline{y})}, \quad j = 1, \dots, N,$$

which in fact meets the boundary condition  $\mathbf{u}'_0 = \mathbf{u}'_{0,\text{in}} + \mathbf{u}'_{0,\text{out}} = \mathbf{0}$ . Apparently, the obtained functions are the direct discrete counterpart to the continuous solution (5.1.7).

The following facts are remarkable:

- (1) We obtained exactly the same solution as obtained for the continuous problem, simply sampled at discrete values on the  $x$ -axis.
- (2) The rays needed to discretize the exterior domain do not need to follow the phase front of the incident waves.
- (3) The incident waves may have both incoming and outgoing parts, nevertheless, the reflected wave is a pure outgoing wave.

**5.1.2. Reflection and Diffraction by a Semi-Infinite Plane.** We apply the results to the reflection/diffraction by a semi-infinite plane. We decompose the  $x$ -axis into a left and a right semi-infinite interval, see Fig. 5.1.3 and require

$$(5.1.31) \quad \lim_{\underline{y} \downarrow 0} \partial_{\underline{y}} u_{21}(x, \underline{y}) = 0 \quad \text{and} \quad \lim_{\underline{y} \downarrow 0} \partial_{\underline{y}} u_{11}(x, \underline{y}) = 0 \quad \text{for } x \leq 0,$$

which are the imposed boundary conditions. The Laplace transformed Helmholtz equation (5.1.8) should hold both in the upper and in the lower half-plane. We denote the distance variable in the upper half-plane by  $y$ , the distance variable in

the lower half-plane with  $\underline{y}$ , as in the previous case. We discretize the axis  $y = 0$  in the interval  $[-a, a]$  and obtain, corresponding to (5.1.9), the two algebraic equations

$$(5.1.32) \quad (\mathbf{A} + s^2\mathbf{I}) \hat{\mathbf{u}}_1(s) = \hat{\mathbf{u}}_{b,1}(s) + s\mathbf{u}_{10} + \mathbf{u}'_{10}$$

$$(5.1.33) \quad (\mathbf{A} + s^2\mathbf{I}) \hat{\mathbf{u}}_2(s) = \hat{\mathbf{u}}_{b,2}(s) + s\mathbf{u}_{20} + \mathbf{u}'_{20},$$

where we want to consider only the finite difference case with  $\mathbf{M} = \mathbf{I}$ . The finite element approach can be treated in complete analogy. Let the incident wave be as in Example 5.1.3

$$(5.1.34) \quad u_{\text{in}}(x_j, \underline{y}) = e^{ik_x x_j} e^{-ik_y \underline{y}}, \quad k_y^2 \in \sigma(\mathbf{A}), \quad 0 < |k_y| < k,$$

where as before  $(\mathbf{u}_{0,\text{in}})_j = u_{\text{in}}(x_j, 0)$ ,  $j = 0, \dots, N+1$ . We make the interval width  $2a$  so large that the following vertical boundary conditions become reasonable

$$(5.1.35) \quad \begin{aligned} u_{11}(-a, y) &= 0 \\ u_{12}(a, y) &= e^{ik_x a} e^{ik_y y} \\ u_{21}(-a, \underline{y}) &= e^{-ik_x a} e^{ik_y \underline{y}} + e^{-ik_x a} e^{-ik_y \underline{y}} \\ u_{22}(a, \underline{y}) &= e^{ik_x a} e^{-ik_y \underline{y}}. \end{aligned}$$

Consequently, we define the vertical boundary data needed to compose the vertical boundary functions  $\hat{\mathbf{u}}_b$ , as they has to be used by the finite difference (FD) method:

FD :

$$\begin{aligned} u_{11}(-a-h, y) &= 0 \\ u_{12}(a+h, y) &= e^{ik_x(a+h)} e^{ik_y y} \\ u_{21}(-a-h, \underline{y}) &= e^{ik_x(-a-h)} e^{ik_y \underline{y}} + e^{ik_x(-a-h)} e^{-ik_y \underline{y}} \\ u_{22}(a+h, \underline{y}) &= e^{ik_x(a+h)} e^{-ik_y \underline{y}} \end{aligned}$$

Then, the boundary functions  $\hat{\mathbf{u}}_{b,1}$  and  $\hat{\mathbf{u}}_{b,2}$  are computed based on the definition (5.1.10), where again the subscript 1 refers to the upper, and the subscript 2 to the lower half-plane:

$$(5.1.36) \quad \hat{\mathbf{u}}_{b,1}(s) := -\frac{1}{h^2} \begin{pmatrix} \hat{u}_{11}(s) \\ \mathbf{0} \\ \hat{u}_{12}(s) \end{pmatrix} = -\frac{1}{h^2} \begin{pmatrix} 0 \\ \mathbf{0} \\ \frac{1}{s-ik_y} e^{ik_x(a+h)} \end{pmatrix}$$

$$(5.1.37) \quad \hat{\mathbf{u}}_{b,2}(s) := -\frac{1}{h^2} \begin{pmatrix} \hat{u}_{21}(s) \\ \mathbf{0} \\ \hat{u}_{22}(s) \end{pmatrix} = -\frac{1}{h^2} \begin{pmatrix} \frac{1}{s-ik_y} e^{ik_x(-a-h)} + \frac{1}{s+ik_y} e^{ik_x(-a-h)} \\ \mathbf{0} \\ \frac{1}{s+ik_y} e^{ik_x(a+h)} \end{pmatrix}$$

The corresponding vectors  $\mathbf{b}_v$  are computed by (5.1.15).

Next, we write simply the general boundary condition (5.1.16) twice, each for the upper and the lower half plane. We set  $\mathbf{B} := \mathbf{US}_-\mathbf{U}^{\text{H}}$  and drop the index 0 at the Dirichlet and Neumann data on  $y = \underline{y} = 0$  and set  $\delta \mathbf{b}_{v,k} := \mathbf{b}_{v,k} - \mathbf{b}_{\text{in},v,k}$ ,  $k \in \{1, 2\}$ . To distinguish between the lower and the upper half-plane, and the

right and left half-plane, respectively, we use indices given in Fig. 5.1.3. Thus the boundary condition (5.1.16) can be written as

$$\begin{aligned} \mathbf{B} \begin{pmatrix} \mathbf{u}_{11} \\ \mathbf{u}_{12} \end{pmatrix} &= - \left( \mathbf{U} \delta \mathbf{b}_{v,1} + \begin{pmatrix} \mathbf{u}'_{11} - \mathbf{u}'_{in,11} \\ \mathbf{u}'_{12} - \mathbf{u}'_{in,12} \end{pmatrix} \right) + \mathbf{B} \mathbf{u}_{in,1} \\ \mathbf{B} \begin{pmatrix} \mathbf{u}_{21} \\ \mathbf{u}_{22} \end{pmatrix} &= - \left( \mathbf{U} \delta \mathbf{b}_{v,2} + \begin{pmatrix} \mathbf{u}'_{21} - \mathbf{u}'_{in,21} \\ \mathbf{u}'_{22} - \mathbf{u}'_{in,22} \end{pmatrix} \right) + \mathbf{B} \mathbf{u}_{in,2}. \end{aligned}$$

We specify the reflection problem by:  $\mathbf{u}_{in,1} = 0$ ,  $\mathbf{u}'_{in,11} = \mathbf{u}'_{in,12} = 0$ ,  $\mathbf{u}'_{11} = \mathbf{u}'_{21} = 0$ . Obviously, we have four unknown vectors  $\mathbf{u}_{11}$ ,  $\mathbf{u}_{12}$ ,  $\mathbf{u}_{22}$ ,  $\mathbf{u}'_{12}$ , since  $\mathbf{u}'_{12} = -\mathbf{u}'_{21}$  by definition of the coordinate systems and  $\mathbf{u}_{12} = \mathbf{u}_{22}$ . Taking this into account, the problem simplifies to

$$(5.1.38) \quad \mathbf{B} \begin{pmatrix} \mathbf{u}_{11} \\ \mathbf{u}_{12} \end{pmatrix} = - \left( \mathbf{U} \delta \mathbf{b}_{v,1} + \begin{pmatrix} \mathbf{0} \\ \mathbf{u}'_{12} \end{pmatrix} \right)$$

$$(5.1.39) \quad \mathbf{B} \begin{pmatrix} \mathbf{u}_{21} \\ \mathbf{u}_{12} \end{pmatrix} = - \left( \mathbf{U} \delta \mathbf{b}_{v,2} + \begin{pmatrix} -\mathbf{u}'_{in,21} \\ -\mathbf{u}'_{12} - \mathbf{u}'_{in,22} \end{pmatrix} \right) + \mathbf{B} \mathbf{u}_{in,2}.$$

We want to show that the discrete half-plane scattering problem has exactly one solution under the similar conditions as the more simple discrete reflection problem of Section 5.1.1, pp. 91.

LEMMA 5.1.4. *Let the discrete scattering problem of a semi-infinite plane defined by the boundary conditions (5.1.34), the discretizations of the upper and lower half-planes (5.1.32), (5.1.33), the incident plane wave (5.1.34) and the vertical boundary conditions (5.1.35).*

*Further assume that the nullspace of the matrix*

$$\mathbf{C} := \begin{bmatrix} (\mathbf{U}\mathbf{S} - \mathbf{U}^{\mathbf{H}}\mathbf{M})_{1,\dots,N_1;1,\dots,N_1} & \mathbf{0} \\ (\mathbf{U}\mathbf{S} - \mathbf{U}^{\mathbf{H}}\mathbf{M})_{N_1+1,\dots,N;1,\dots,N_1} & -\mathbf{I}_{N_2,N_2} \end{bmatrix}$$

*with  $N_1$  the number of points on the rigid boundary and  $N_1 + N_2 = N$  is empty.*

*Then, the discrete problem has a unique solution and the reflected wave is an outgoing wave.*

PROOF. We have to show that the two equations (5.1.38), (5.1.39) have a unique solution. First we take the sum and the difference of (5.1.32) and (5.1.39)

$$(5.1.40) \quad \mathbf{B} \begin{pmatrix} \mathbf{u}_{11} + \mathbf{u}_{21} \\ 2\mathbf{u}_{12} \end{pmatrix} = - \left( \mathbf{U} (\delta \mathbf{b}_{v,1} + \delta \mathbf{b}_{v,2}) + \begin{pmatrix} -\mathbf{u}'_{in,21} \\ -\mathbf{u}'_{in,22} \end{pmatrix} \right) + \mathbf{B} \mathbf{u}_{in,2}$$

$$(5.1.41) \quad \mathbf{B} \begin{pmatrix} \mathbf{u}_{11} - \mathbf{u}_{21} \\ 0 \end{pmatrix} = - \left( \mathbf{U} (\delta \mathbf{b}_{v,1} - \delta \mathbf{b}_{v,2}) + \begin{pmatrix} \mathbf{u}'_{in,21} \\ 2\mathbf{u}'_{12} + \mathbf{u}'_{in,22} \end{pmatrix} \right) - \mathbf{B} \mathbf{u}_{in,2}$$

We consider (5.1.40) first. We aim to establish an orthogonality condition similar to that in Corollary 5.1.2. Let  $j_0$  indicate elements of nullspaces, that is  $\mathbf{u}_{j_0} \in \mathcal{N}(\mathbf{A})$  and the corresponding eigenvalue  $\lambda_{j_0}^2 \in \sigma(\mathbf{A})$  is zero. Further let  $\mathbf{u}_{j_0} = \mathbf{U} \mathbf{e}_{j_0}$ . The system (5.1.40) has a unique solution if



$$\begin{aligned}
& \left\langle \mathbf{e}_{j_0}, \delta \mathbf{b}_{v,1} + \delta \mathbf{b}_{v,2} + \mathbf{U}^H \begin{pmatrix} -\mathbf{u}'_{in,21} \\ -\mathbf{u}'_{in,22} \end{pmatrix} \right\rangle \\
&= (\delta \mathbf{b}_{v,1} + \delta \mathbf{b}_{v,2})_{j_0} + \left\langle \mathbf{u}_{j_0}, \begin{pmatrix} -\mathbf{u}'_{in,21} \\ -\mathbf{u}'_{in,22} \end{pmatrix} \right\rangle \\
&= \left\langle \mathbf{u}_{j_0}, \widehat{\mathbf{u}}_{b,1}(0) - \widehat{\mathbf{u}}_{b,in,1}(0) + \widehat{\mathbf{u}}_{b,2}(0) - \widehat{\mathbf{u}}_{b,in,2}(0) \right\rangle - \left\langle \mathbf{u}_{j_0}, \begin{pmatrix} \mathbf{u}'_{in,21} \\ \mathbf{u}'_{in,22} \end{pmatrix} \right\rangle \\
&= \left\langle \mathbf{u}_{j_0}, \widehat{\mathbf{u}}_{b,1}(0) - \widehat{\mathbf{u}}_{b,in,1}(0) + \widehat{\mathbf{u}}_{b,2}(0) - \widehat{\mathbf{u}}_{b,in,2}(0) - \begin{pmatrix} \mathbf{u}'_{in,21} \\ \mathbf{u}'_{in,22} \end{pmatrix} \right\rangle \\
&= 0.
\end{aligned}$$

The evaluation of (5.1.36) and (5.1.37) shows  $\widehat{\mathbf{u}}_{b,1}(0) + \widehat{\mathbf{u}}_{b,2}(0) = \mathbf{0}$ . Hence the condition reduces to

$$\left\langle \mathbf{u}_{j_0}, \widehat{\mathbf{u}}_{b,in,1}(0) + \widehat{\mathbf{u}}_{b,in,2}(0) + \begin{pmatrix} \mathbf{u}'_{in,21} \\ \mathbf{u}'_{in,22} \end{pmatrix} \right\rangle = 0.$$

This, however coincides with condition (5.1.27), p. 101, whose validity has been proved. Thus (5.1.40) has a unique solution.

Now we consider (5.1.41), which we want to rewrite, setting  $\mathbf{u} := \mathbf{u}_{11} - \mathbf{u}_{21}$ ,  $\mathbf{v} = 2\mathbf{u}'_{12}$ ,

$$\mathbf{r} := - \left( \mathbf{U} (\delta \mathbf{b}_{v,1} - \delta \mathbf{b}_{v,2}) + \begin{pmatrix} \mathbf{u}'_{in,21} \\ \mathbf{u}'_{in,22} \end{pmatrix} \right) - \mathbf{U} \mathbf{S}_- \mathbf{U}^H \mathbf{u}_{in,2},$$

and using the definition of B above to

$$\mathbf{U} \mathbf{S}_- \mathbf{U}^H \begin{pmatrix} \mathbf{u} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbf{0} \\ \mathbf{v} \end{pmatrix} = \mathbf{r}.$$

Using the matrix C this linear problem takes the form.

$$(5.1.42) \quad \mathbf{C} \begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix} = \mathbf{r}.$$

By assumption, (5.1.42) has a unique solution. Consequently, the whole problem has a unique solution.  $\square$

## 5.2. Laplace domain method

The Laplace domain methods in their two forms supply a technique to satisfy the pole condition without any explicit factorization. The price to pay is an additional discretization along axes in the spectral domain. We based the real axis approach on a conventional spline-collocation technique and the cut-function approach on a Runge-Kutta-technique. Both realizations are the very first implementations and far-away from being optimal. For example, it seems to be quite natural to replace the spline technique by a *wavelet* realization. Most of all questions are open in this direction. Nevertheless, the numerical results obtained so far demonstrate clearly the high potential of the Laplace domain methods. We describe both realizations in great detail, such that they may repeated immediately.

**Model equation.** We consider the numerical solution of our basic 1D model problem, the Bessel equation

$$(5.2.1) \quad \partial_r^2 u + \frac{d-1}{r} \partial_r u + k^2 u - \frac{\nu^2}{r^2} u = 0,$$

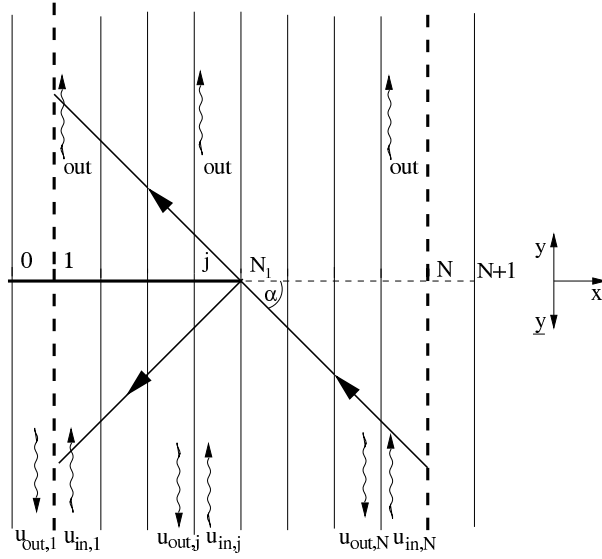


FIGURE 5.1.3. Diffraction problem: discretization by vertical rays

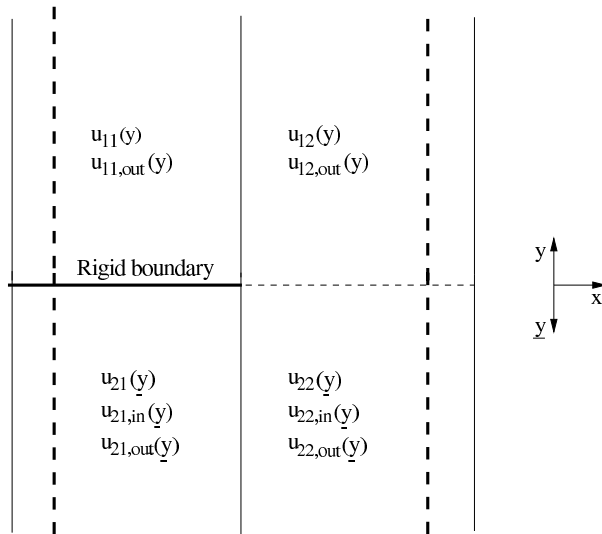


FIGURE 5.1.4. Notation of the solution on different sub-domains

where again  $d$  is the space dimension and  $\nu$  an arbitrary real separation constant. We define the interior and the exterior domain to  $\Omega_{\text{int}} = (r_1, r_0)$  and  $\Omega_{\text{ext}} = (r_0, \infty)$ ,  $r_0 > r_1 > 0$  and impose the Dirichlet condition  $u(r_1) = u_i$  on the interior boundary.

The map  $u \mapsto \tilde{u} = r^{-(d-1)/2}u$  transforms (5.2.1) into

$$(5.2.2) \quad \partial_r^2 u + \left[ k^2 - \frac{\frac{(d-1)^2}{4} - \frac{d-1}{2} + \nu^2}{r^2} \right] u = 0,$$

where we dropped the tilde on top of  $u$  for notational simplicity. Thus we arrive at the same equation on which our theoretical discussion in Section 4.2 is based.

**5.2.1. Real Axis Approach.** The solution of (5.2.2) includes three subproblems supplying the three parts of the final algebraic system: First, the variational

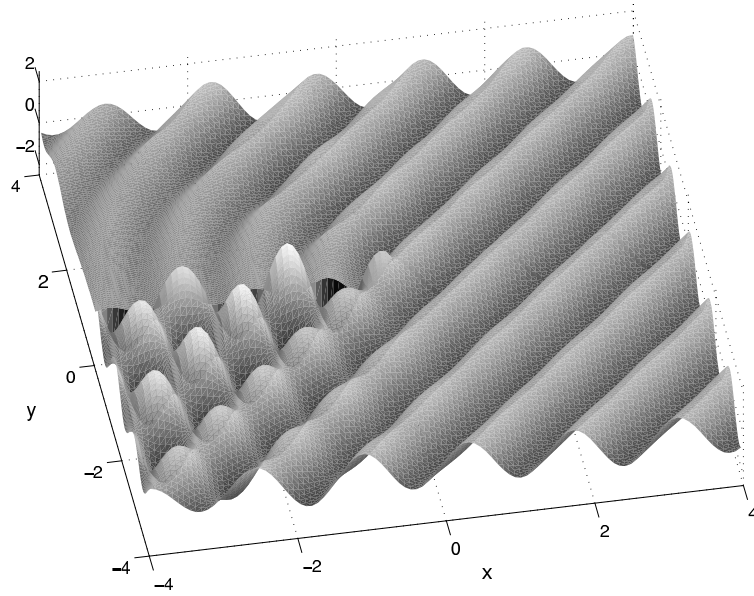


FIGURE 5.1.5. Diffraction of a plane wave with an incident angle of  $3/2\pi$  by a rigid semi-infinite plane sheet ( $y = 0, x \leq 0$ ).

formulation of (5.2.2) on  $\Omega_{\text{int}}$ , second, the pole condition in integral form on the real axis linking the exterior and the interior problem such that asymptotically only outgoing waves exist, and, third the Laplace transform of (5.2.2) on  $\Omega_{\text{ext}}$ .

Based on the function space

$$\begin{aligned} V &= \{v \in H^1(\Omega_{\text{int}}) : v(r_i) = u_i\} \\ V_0 &= H_0^1(\Omega_{\text{int}}) \end{aligned}$$

(note that  $H^1(\Omega_{\text{int}}) \subset C(\Omega_{\text{int}})$ ) and the sesquilinear form

$$\begin{aligned} a &: V_0 \times V \rightarrow \mathbb{C} \\ a(v, u) &= -(\nabla v, \nabla u) + \left( v, \left[ k^2 - \frac{(d-1)^2}{4} - \frac{d-1}{2} + \nu^2 \right] \frac{u}{r^2} + k^2 u \right), \end{aligned}$$

we rewrite (5.2.2) into variational form yielding the first part of our problem.

**Part I:** (Continuous interior problem): Find  $u \in V$  such that for all  $v \in V_0$

$$a(v, u) + (\bar{v}u')(r_0) = 0.$$

The interior and the exterior problem are linked together by the pole condition, which gives the second part of the problem.

**Part II:** (Continuous pole condition):

$$0 = \int_{\infty}^{-\infty} d\tau \frac{\hat{u}(\tau)}{\tau - s_+}, \quad \tau \in \mathbb{R}, \text{Im}(s_+) > 0.$$

On the exterior domain we consider the Laplace transform of (5.2.2), where the term  $u/r^2$  is given in its integral operator representation. This equation supplies the third part of the problem.

**Part III:** (Continuous Laplace transformed exterior problem):

$$(5.2.3) \quad \widehat{u}(s) - \frac{su(r_0) + u'(r_0)}{s^2 + k^2} - \frac{\frac{(d-1)^2}{4} - \frac{d-1}{2} + \nu^2}{s^2 + k^2} \int_s^\infty ds_1 e^{(s-s_1)r_0} (s-s_1) \widehat{u}(s_1) = 0.$$

The discretization of the sesquilinear form  $a(v, u)$  results in standard manner by the approximation of  $V$  by a finite-dimensional space  $V_h \subset V$ . Let  $\dim V_h = n$ . Thus we approximate  $a(v, u)$  by the symmetric matrix  $A_{\text{int}} \in \mathbb{R}^{n \times n}$  and arrive at the discrete interior problem .

**Part I: Discrete Interior Problem.**

Standard finite element technology supplies the system

$$A_{\text{int}} \mathbf{u} + \begin{bmatrix} \mathbf{0} \\ (\bar{v}u')(r_0) \end{bmatrix} = \mathbf{r}$$

where  $\mathbf{u} \in \mathbb{C}^n$  is the vector of interior unknowns and the right-hand side  $\mathbf{r} \in \mathbb{C}^n$  vanishes identically except for the entries arising from the given interior Dirichlet condition.

**Part II: Discrete Pole Condition.**

To approximate  $\widehat{u}$  on the real axis, we define a set of nodes  $\{x_i\}_{i=0}^m$  such that

$$-\infty < x_m < x_{m-1} < \dots < x_1 < x_0 < \infty.$$

On the interval  $[x_0, x_m]$  we use the cubic  $C^1$ -splines ( $i = 1, \dots, m-1$ )

$$\phi_i(x) = \begin{cases} \phi_{i,1}(x) = ((x-x_i)/h_i + 1)^2 (1 - 2(x-x_i)/h_i) & \text{for } x \in [x_{i-1}, x_i] \\ \phi_{i,2}(x) = ((x-x_i)/h_{i+1} - 1)^2 (1 + 2(x-x_i)/h_{i+1}) & \text{for } x \in [x_i, x_{i+1}] \\ 0 & \text{elsewhere,} \end{cases}$$

$$\psi_i(x) = \begin{cases} \psi_{i,1}(x) = ((x-x_i)/h_i + 1)^2 (x-x_i) & \text{for } x \in [x_{i-1}, x_i] \\ \psi_{i,2}(x) = ((x-x_i)/h_{i+1} - 1)^2 (x-x_i) & \text{for } x \in [x_i, x_{i+1}] \\ 0 & \text{elsewhere,} \end{cases}$$

with  $h_i = x_i - x_{i-1}$ , as basis functions, see Fig. 5.2.1 for a graphical representation.

The support of these basis functions is contained completely in the interval  $[x_0, x_m]$ . Because we need a representation on the full axis, we define additionally

$$\phi_0(x) = \begin{cases} \phi_{0,1}(x) & \text{for } x \geq x_0 \\ \phi_{0,2}(x) = ((x-x_0)/h_1 - 1)^2 (1 + 2(x-x_0)/h_1) & \text{for } x \in [x_0, x_1] \\ 0 & \text{elsewhere,} \end{cases}$$

$$\psi_0(x) = \begin{cases} \psi_{0,2}(x) = ((x-x_0)/h_1 - 1)^2 (x-x_0) & \text{for } x \in [x_0, x_1] \\ 0 & \text{elsewhere,} \end{cases}$$

for the basis functions crossing  $x_0$  and

$$\phi_m(x) = \begin{cases} \phi_{m,1}(x) = ((x-x_m)/h_m + 1)^2 (1 - 2(x-x_m)/h_m) & \text{for } x \in [x_{m-1}, x_m] \\ \phi_{m,2}(x) & \text{for } x \leq x_m \\ 0 & \text{elsewhere,} \end{cases}$$

$$\psi_m(x) = \begin{cases} \psi_{m,1}(x) = ((x-x_m)/h_m + 1)^2 (x-x_m) & \text{for } x \in [x_{m-1}, x_m] \\ 0 & \text{elsewhere,} \end{cases}$$

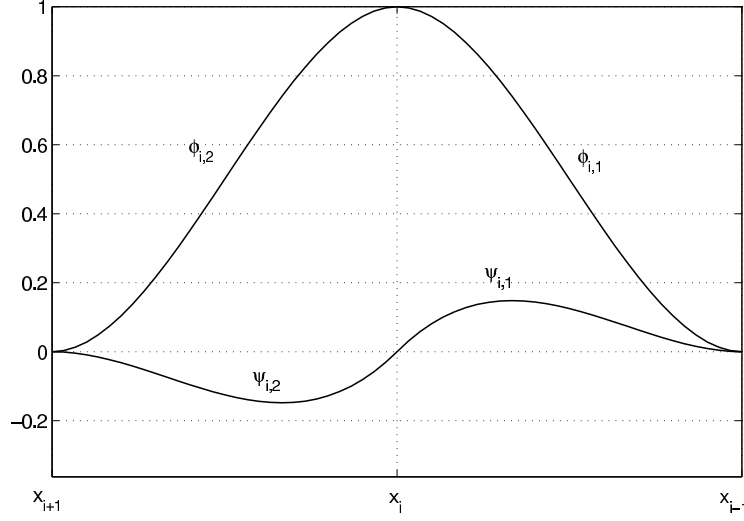


FIGURE 5.2.1. The function  $\phi_i$  is the Lagrange basis function with respect to  $\hat{u}(x_i)$ ,  $\psi_i$  is the derivative of the of the Lagrange basis function with respect to  $\hat{u}'(x_i)$ .

for the basis functions crossing  $x_m$ . The functions  $\phi_{0,1}(x)$  and  $\phi_{m,2}(x)$  are not given explicitly, here. They can be chosen from the asymptotic expressions of the solution of the basic integral equation or can be any other reasonable functions ensuring the continuity conditions  $\phi_{0,1}(x_0) = 1$  and  $\phi_{m,2}(x_m) = 1$ . For simplicity of the presentation we will use simply the constant functions  $\phi_{0,1}(x) = 1$  and  $\phi_{m,2}(x) = 1$  for  $x$  outside the interval  $[x_0, x_m]$ .

Thus we have *two* independent conditions per node or, equivalently, we may choose two collocation points per interval. This classic spline representation has the feature that the coefficients of  $\phi_i(x)$  have an immediate interpretation as nodal values of the function at  $x = x_i$  and the coefficients of  $\psi_i(x)$  as the derivatives at these positions. Note, that the interval lengths  $h_i$  in our definition are negative. We approximate a complex-valued function  $\hat{u}$  on the real axis by

$$(5.2.4) \quad \hat{u}_h(x) = \sum_{i=0}^m \hat{u}_i \phi_i(x) + \sum_{i=0}^m \hat{u}'_i \psi_i(x).$$

This requires  $2m+2$  coefficients for  $m$  intervals. Along with the set of nodes  $\{x_i\}_{i=0}^m$  we define the set of Gaussian collocation points  $\{x_{gi}\}_{i=1}^{2m}$  by

$$\begin{aligned} x_{g,2i-1} &= \frac{x_i + x_{i-1}}{2} - \frac{h_i}{2\sqrt{3}} \\ x_{g,2i} &= \frac{x_i + x_{i-1}}{2} + \frac{h_i}{2\sqrt{3}}, \quad h_i = x_i - x_{i-1}. \end{aligned}$$

Let  $\hat{\mathbf{u}}_g \in \mathbb{C}^{2m}$  be the vector of  $\hat{u}_h$  evaluated at the collocation points. Further, let  $\hat{\mathbf{u}}_c = (\hat{u}_0, \hat{u}_1, \hat{u}'_1, \hat{u}_2, \hat{u}'_2, \dots, \hat{u}_{m-1}, \hat{u}'_{m-1}, \hat{u}_m) \in \mathbb{C}^{2m}$  the vector of spline coefficients. Then, the evaluation of (5.2.4) at all collocation points  $x_{gi}, i = 1, \dots, 2m$  defines the transform

$$(5.2.5) \quad \hat{\mathbf{u}}_g = M_c \hat{\mathbf{u}}_c.$$

Next we approximate the integral condition by piecewise Gaussian quadratures. We replace the integral condition by a discrete approximation

$$0 = \frac{1}{2} \left( \sum_{i=1}^m \frac{h_i}{x_{g,2i-1} - s_+} \hat{u}_{g,2i-1} + \sum_{i=1}^m \frac{h_i}{x_{g,2i} - s_+} \hat{u}_{g,2i} \right)$$

with  $\hat{u}_{g,j}$ ,  $1 \leq j \leq 2m$  the  $j$ th component of  $\hat{\mathbf{u}}_g$ . It is convenient to rewrite this in form of an Euclidian inner product based on the  $m \times m$  identity matrix  $\mathbf{I}_{m,m}$  and the vector  $\mathbf{h} = (h_1, \dots, h_m)$

$$0 = \frac{1}{2} \left\langle \text{diag} \left( \frac{1}{x_{gi} - s_+} \right)_{i=1, \dots, 2m} \mathbf{I} \otimes \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mathbf{h}, \hat{\mathbf{u}}_g \right\rangle.$$

Based on (5.2.5), we obtain the discrete pole condition in terms of the spline coefficients

$$\frac{1}{2} \left\langle \mathbf{M}_c^T \text{diag} \left( \frac{1}{x_{gi} - s_+} \right)_{i=1, \dots, 2m} \mathbf{I}_{m,m} \otimes \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mathbf{h}, \hat{\mathbf{u}}_c \right\rangle = 0.$$

### Part III: Discrete Laplace Transformed Exterior Problem.

The spectral equation (5.2.3) is solved by collocation, i.e., it is rewritten in terms of the collocation coefficients and solved pointwise for any collocation point  $x_{gi}$ . The only non-standard problem here is the representation of the integral operator.

We define the complex number  $\hat{I}_i$  by

$$\hat{I}_i \hat{u}_h = \int_{x_{gi}}^{\infty} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \hat{u}_h(x_1).$$

The evaluation of the expression for each  $x_{gi}$  will supply one row of the full discrete integral operator. Using (5.2.4) we obtain

$$\begin{aligned} \hat{I}_i \hat{u}_h &= \int_{x_{gi}}^{\infty} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \sum_{j=0}^m \hat{u}_j \phi_j(x_1) \\ &\quad + \int_{x_{gi}}^{\infty} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \sum_{j=0}^m \hat{u}'_j \psi_j(x_1). \end{aligned}$$

This integral operator might be assembled in a FEM-like manner, i.e., the contributions of each segment is computed separately and added to the global operator. We describe this approach briefly. First we rewrite the representation (5.2.4) such that trial functions with support on the same interval are grouped together

$$\begin{aligned} \hat{u}_h(x) &= \sum_{i=1}^m \hat{u}_{i-1} \phi_{i-1,2}(x) + \hat{u}'_{i-1} \psi_{i-1,2}(x) + \hat{u}_i \phi_{i,1}(x) + \hat{u}'_i \psi_{i,1}(x) \\ &\quad + \hat{u}_0 \phi_{0,1}(x) + \hat{u}_m \phi_{m,2}(x) \\ &=: \sum_{i=1}^m [\hat{u}_{i-1}, \hat{u}'_{i-1}, \hat{u}_i, \hat{u}'_i]_{\Sigma} (x) + \hat{u}_0 \phi_{0,1}(x) + \hat{u}_m \phi_{m,2}(x) \end{aligned}$$

We use the expression in brackets as short-hand for the full expression and remark that  $[\hat{u}_{i-1}, \hat{u}'_{i-1}, \hat{u}_i, \hat{u}'_i]_{\Sigma}$  refers to the interval  $[x_i, x_{i-1}]$ . Note, that the boundary conditions  $\psi_{0,2}(x) = \psi_{m,1}(x) \equiv 0$  are understood implicitly within this notation. Let  $x_k$  be the smallest element of  $\{x_i\}_{i=0}^m$  larger than a given  $x_{gi}$ , from the set of Gaussian collocation points. Then we may write

$$\begin{aligned}
\widehat{I}_i \widehat{u}_h &= \int_{x_0}^{\infty} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \widehat{u}_0 \phi_{0,1}(x_1) \\
&+ \int_{x_k}^{x_0} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \sum_{j=1}^k [\widehat{u}_{j-1}, \widehat{u}'_{j-1}, \widehat{u}_j, \widehat{u}'_j]_{\Sigma}(x_1) \\
&+ \int_{x_{gi}}^{x_k} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) [\widehat{u}_k, \widehat{u}'_k, \widehat{u}_{k+1}, \widehat{u}'_{k+1}]_{\Sigma}(x_1).
\end{aligned}$$

The integral expression can be computed element-wise. To this end we introduce

$$\begin{aligned}
\widehat{I}_{i0} \widehat{u}_h &:= \int_{x_0}^{\infty} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \widehat{u}_0 \phi_{0,1}(x_1) \\
\widehat{I}_{ij} \widehat{u}_h &:= \int_{x_j}^{x_{j-1}} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) [\widehat{u}_{j-1}, \widehat{u}'_{j-1}, \widehat{u}_j, \widehat{u}'_j]_{\Sigma}(x_1) \\
&\quad \text{for } 1 \leq j \leq k \\
\widehat{I}_{i,k+1} \widehat{u}_h &:= \int_{x_{gi}}^{x_k} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) [\widehat{u}_k, \widehat{u}'_k, \widehat{u}_{k+1}, \widehat{u}'_{k+1}]_{\Sigma}(x_1),
\end{aligned}$$

which allows us to compute the contribution of the  $i$ th Gauss-point as

$$\widehat{I}_i \widehat{u}_h = \widehat{I}_{i0} \widehat{u}_h + \sum_{j=1}^k \widehat{I}_{ij} \widehat{u}_h + \widehat{I}_{i,k+1} \widehat{u}_h.$$

The indices of  $\widehat{I}_{ij}$  becomes easier readable if we remember that the first index,  $i$ , counts the number of the Gauss-point ( $1 \leq i \leq 2m$ ) and the second index,  $j$ , denotes the *end* point of the corresponding interval  $[x_{j-1}, x_j]$ ,  $1 \leq j \leq m$ . The element-wise contributions can be computed exactly or by numerical quadrature. On the interval  $x \geq x_0$  we continue  $\widehat{u}_h$  by a constant function – as discussed above – and set  $\phi_{0,1} = 1$ . Since  $\widehat{u} = \mathcal{O}(1/|x|)$  for  $x \rightarrow \infty$  on the real axis, this introduces an error of the order  $\mathcal{O}(1/|x_0|^2)$ . Of course this may be improved essentially by using the known asymptotic approximations of  $\widehat{u}$ . Further, as it is seen, the entries decay *exponentially* away from the  $i$ th segment. We associate to each number  $\widehat{I}_{i0} \widehat{u}_h$  a complex number  $\widehat{I}_{i0}$  and to each number  $\widehat{I}_{ij} \widehat{u}_h$  a local  $1 \times 4$  matrix  $\widehat{I}_{ij}$  through

$$\begin{aligned}
\widehat{I}_{i0} &:= \int_{x_0}^{\infty} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) \phi_{0,1}(x_1) \\
\widehat{I}_{ij} &:= \int_{x_j}^{x_{j-1}} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) [\phi_{j-1,2}(x_1), \psi_{j-1,2}(x_1), \phi_{j,1}(x_1), \psi_{j,1}(x_1)] \\
&\quad \text{for } 1 \leq j \leq k \\
\widehat{I}_{i,k+1} &:= \int_{x_{gi}}^{x_k} dx_1 e^{(x_{gi}-x_1)r_0} (x_{gi} - x_1) [\phi_{k,2}(x_1), \psi_{k,2}(x_1), \phi_{k+1,1}(x_1), \psi_{k+1,1}(x_1)],
\end{aligned}$$

The resulting discrete integral operator  $\widehat{\mathbf{I}}$  is composed from the local matrices  $\widehat{I}_{ij}$ . The corresponding procedure is given in Algorithm 1. This way we assemble a matrix  $\widehat{\mathbf{I}} \in \mathbb{C}^{2m \times (2m+2)}$ . We arrive at the desired matrix if we delete the columns corresponding to  $\psi_0$  (the column 2) and to  $\psi_m$  (the last column). Further we replace the vector of nodal values at the Gaussian points  $\widehat{\mathbf{u}}_g$  by the vector of unknown spline coefficients  $\widehat{\mathbf{u}}_c$  according to  $\widehat{\mathbf{u}}_g = \mathbf{M}_c \widehat{\mathbf{u}}_c$ , (5.2.5), and abbreviate  $c_\nu = -(d-1)^2/4 + (d-1)/2 - \nu^2$  to obtain

**Algorithm 1** FEM-like assembling of the integral operator

---

```

for  $k = 1$  to  $m$  do
  for  $j = k$  to  $m$  do
     $l = 2j - 1; n = 2j - 2k + 1$ 
     $\widehat{\mathbf{I}}(l, n : n + 3) = \widehat{\mathbf{I}}(l, n : n + 3) + \widehat{\mathbf{I}}_{l, n-1}$ 
     $l = l + 1$ 
     $\widehat{\mathbf{I}}(l, n : n + 3) = \widehat{\mathbf{I}}(l, n : n + 3) + \widehat{\mathbf{I}}_{l, n-1}$ 
  end for
end for
for  $l = 1$  to  $2m$  do
   $\widehat{\mathbf{I}}(l, 1) = \widehat{\mathbf{I}}(l, 1) + \widehat{\mathbf{I}}_{l, 0}$ 
end for

```

---

$$\begin{aligned}
\mathbf{M}_c \widehat{\mathbf{u}}_c - \left[ \frac{x_{gi}}{x_{gi}^2 + k^2} \right]_{i=1, \dots, 2m} u(r_0) - \left[ \frac{1}{x_{gi}^2 + k^2} \right]_{i=1, \dots, 2m} u'(r_0) \\
+ c_\nu \text{diag} \left( \frac{1}{x_{gi}^2 + k^2} \right) \widehat{\mathbf{I}} \widehat{\mathbf{u}}_c = \mathbf{0}.
\end{aligned}$$

All three contributions of the composite problem are assembled into an algebraic system of the form  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and solved simultaneously. This is summarized in Algorithm 2, and the structure of the composite system matrix is displayed in Fig. 5.2.2.

**Algorithm 2** Discrete operator with pole condition in spectral form

- 
1. FEM discretization of the interior problem
 
$$A_{11} = A_{\text{int}} + \begin{bmatrix} \mathbf{0} \\ (\overline{v}u')(r_0) \end{bmatrix}$$
  2. Discretization of the pole condition
 
$$A_{22} = \left[ \frac{1}{2} \quad \mathbf{M}_c^T \left( \text{diag} \left( \frac{1}{x_{gi} - s_+} \right)_{i=1, \dots, m} I_{m, m} \otimes \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mathbf{h} \right)^T \right]$$
  3. Discretization of the Laplace transformed exterior equation
 
$$A_{32} = \left[ -\mathbf{g}_0, -\mathbf{g}', \mathbf{M}_c + c_\nu \text{diag}(\mathbf{g}') \widehat{\mathbf{I}} \right]$$
 with  $\mathbf{g}_0 = \left[ \frac{x_{gi}}{x_{gi}^2 + k^2} \right]_{i=1, \dots, 2m}$  and  $\mathbf{g}' = \left[ \frac{1}{x_{gi}^2 + k^2} \right]_{i=1, \dots, 2m}$
  4. Assemble the system matrix  $A$ 

$$A = \begin{bmatrix} A_{11} & \mathbf{0} \\ \mathbf{0} & A_{22} \\ \mathbf{0} & A_{32} \end{bmatrix}$$
- 

**Numerical studies.** To test the numerical performance of the procedure, we define the following test problem:  $r_1 = 3$ ,  $r_0 = 3 + 2\pi$ , dimension  $d = 2$ ,  $k = 1$ . The interior problem is discretized by  $n$  linear finite elements, the real axis for the discretization of the Laplace-transformed equation by  $N$  collocation points, symmetrically placed around the origin. The real axis is discretized uniformly in the interval  $[-2, 2]$  by  $N/4$  of all points and with geometrically increasing segment lengths outside this interval. The full interval on the real axis were chosen to be  $[-5/8N, 5/8N]$ , hence its size increases if the number of collocation points increases. Fig. 5.2.3 shows the relative error of the computed DtN-number as a



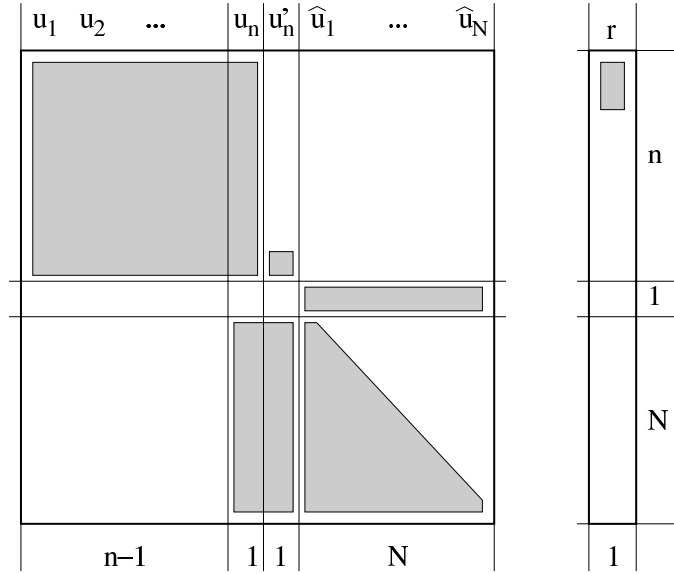


FIGURE 5.2.2. Occupation pattern of the composite matrix. The upper left block matrix is the sparse, real and symmetric FE-matrix, the lower right block matrix is empty in its  $(N-3) \times (N-3)$  upper triangular part. The number of interior degrees of freedom is  $n$ , the number of collocation points is  $N = 2m$ , with  $m$  the number of segments on the real axis.

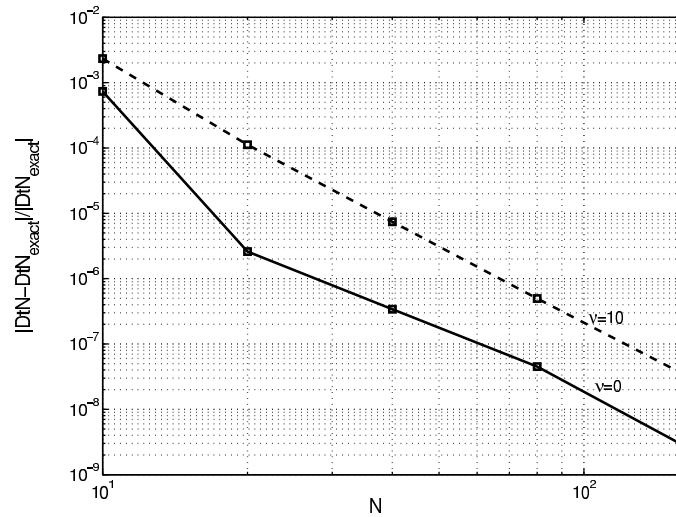


FIGURE 5.2.3. Convergence of the algorithm with respect to the relative error of the DtN-number as a function of the number of collocation points  $N$  on the real axis.

function of the number of collocation points  $N$  on the real axis for the frequencies  $\nu = 0$  and  $\nu = 10$ . The relative error has been computed by setting  $u(r_0) = 1$ , i.e. without to solve the interior problem, in order to obtain the result independent from the interior problem. The result shows that even with a relative small number of spectral collocation points ( $N \leq 160$ ) a good accuracy can be achieved. It

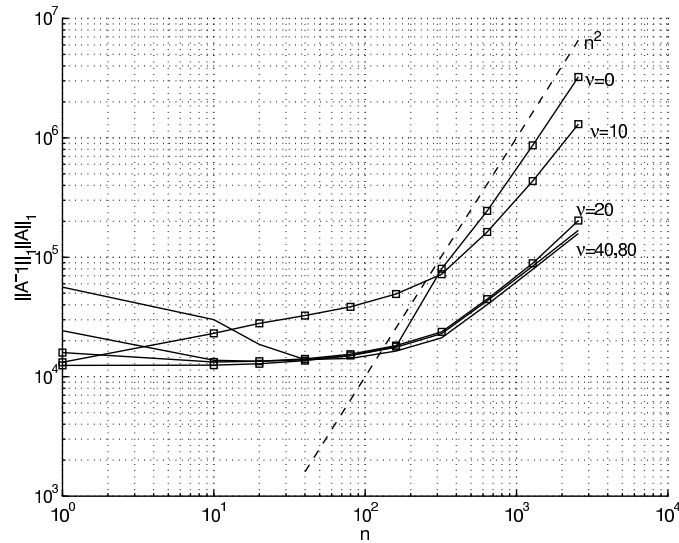


FIGURE 5.2.4. Condition number with respect to the 1-norm of the resulting matrix vs. the number of interior degrees of freedom  $n$ . The number of spectral points is kept fix with  $N = 160$ . The dashed line displays the function  $n^2$  indicating the condition number of the discrete Laplacian.

should be noted that this result was obtained with an explicit restriction of the triangular integral operator matrix to a bandwidth of  $N/4$ . A far better compression might be obtained using adaptive discretization techniques and more elaborated compression techniques. Fig. 5.2.4 shows the evolution of the relative condition number of the composite system matrix  $A$ . The important results are that, first, higher frequencies  $\nu$  do not cause an explosion of the condition number and, second, that the condition number is asymptotically dominated by the condition number of the corresponding discrete Laplacian of the interior problem. Hence the stability properties of the discrete Helmholtz scattering problem and the discrete Laplace problem are comparable.

This is underlined by the computation of the spectral norm of the inverse of the system matrix  $A$ . To ensure discrete solvability we want that the norm of the inverse remains bounded. Fig. 5.2.5 shows the computed results as a function of the number of interior nodes, where the number of spectral collocation points is kept to  $N = 160$ . The figure shows that the computed norms are even getting smaller as the number of unknown increases until it approaches a certain constant.

**Final remarks.** So far, we have not specified, how to carry out the integration of the spline elements. The numerical experiments demonstrated above used closed formulas obtained via a symbolic computation. This becomes possible, since all the involved integrals are of elementary nature (exponential functions times polynomials). For an efficient implementation, however, this should be replaced by numerical quadrature formulas. There is a variety of possible choices. The next section which deals with the implementation of the cut-function approach, presents one of the possibilities: the realization based on Runge-Kutta type integrators. This is a very efficient approach for the present case, too. Since we will discuss the details in the next section, we drop it here. However, in order to compare the structure of

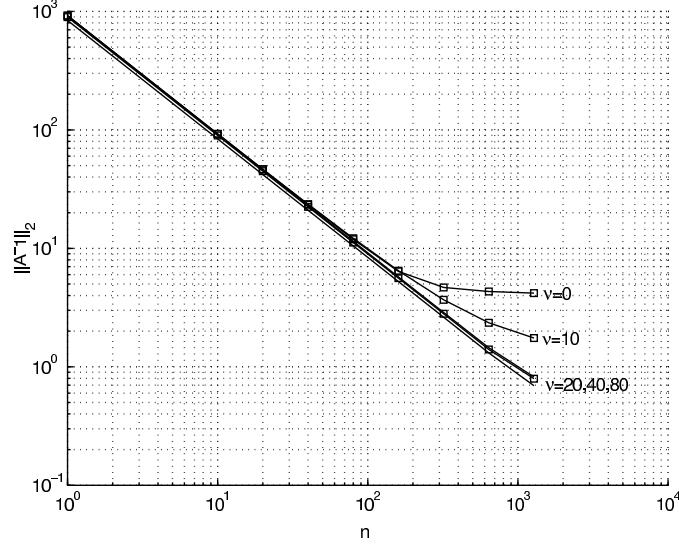


FIGURE 5.2.5. Spectral norm of the inverse of the resulting matrix vs. the number of interior degrees of freedom.

both methods, we give the pole condition and the integral equation using  $\tilde{\tau} = -s$ ,  $\hat{u}(s) = \hat{u}(-\tilde{\tau}) =: \tilde{u}(\tilde{\tau})$  in a slightly different form

$$\begin{aligned} \tilde{u}(\tilde{\tau}) + \frac{\tilde{\tau}u(r_0) - u'(r_0)}{\tilde{\tau}^2 + k^2} + c_\nu \int_{-\infty}^{\tilde{\tau}} d\tau_1 e^{(\tau_1 - \tilde{\tau})r_0} (\tau_1 - \tilde{\tau}) \tilde{u}(\tau_1) &= 0 \\ \int_{-\infty}^{\infty} d\tilde{\tau} \frac{\tilde{u}(\tilde{\tau})}{\tilde{\tau} + s_+} &= 0 \end{aligned}$$

We add that in general the Runge-Kutta type realization of the Volterra integral equation –unlike the ODE-case–is stable only for *explicit* schemes if large step sizes are used. This results from the fact that implicit methods have to compute  $\exp(a(\tau - \tilde{\tau}))$  as part of the kernel for  $\tau > \tilde{\tau}$ , which causes large entries in the discrete approximation and leads finally to extremely large condition numbers. Finally we remark that we implemented successfully the real axis approach to the Laplace domain method using the explicit Runge-Kutta methods from the explicit Euler scheme ( $p = 1$ ) to the Dormand-Prince RK5(4) scheme ( $p = 5$ ).

**5.2.2. Cut Function Approach.** We derive an algorithm, which applies the results of Chapter 4, especially of Sections 4.3, 4.4.2 (pp. 74) in a direct manner. This approach possesses the essential advantage to supply the far field without any extra numerical costs. We use the same model problem as in the foregoing section. For convenience, we repeat it here. We wish to solve Bessel's equation (5.2.1)

$$\partial_r^2 u + \frac{d-1}{r} \partial_r u + k^2 u - \frac{\nu^2}{r^2} u = 0$$

on the real half-axis  $\{r \in \mathbb{R}_+ : r > r_i\}$ . At  $r_i$  a Dirichlet condition is given. Bessel's equation is transformed to (5.2.2)

$$\partial_r^2 u + \left[ k^2 + \frac{c_\nu}{r^2} \right] u = 0$$

with  $c_\nu = -(d-1)^2/4 + (d-1)/2 - \nu^2$  and a new function  $u$ . We solve the transformed equation simultaneously on the domains  $\Omega_{\text{int}} = (r_i, r_0)$  and  $\Omega_{\text{ext}} = (r_0, \infty)$ ,

$r_0 > r_i > 0$ . Our strategy is the same as in Section 5.2, i. e., we decompose the problem into three subproblems: The interior problem (Part I), the linking conditions between the interior and the exterior problem (Part II), and the exterior problem in spectral form (Part III). First we give the three subproblems in continuous form.

**Part I:** (Continuous interior problem): Find  $u \in V$  such that for all  $v \in V_0$

$$a(v, u) + (\bar{v}u')(r_0) = 0,$$

**Part II:** (Continuous linking conditions).

$$\begin{aligned} u(r_0) - \frac{w(i)}{2i} + \frac{1}{2\pi i} \int_0^\infty dt \hat{\psi} &= 0 \\ u'(r_0) - \frac{w(i)}{2} + \frac{1}{2\pi i} \int_0^\infty dt (i-t) \hat{\psi} &= 0. \end{aligned}$$

These are special forms of our representation formulas, evaluated at  $r_0$ . Note, that we used the map  $w(i)\psi \mapsto \hat{\psi}$  to define a new function  $\hat{\psi}$  in spectral domain. Further, these representations suggest to use the pole strength  $w(i)$  as additional unknown for the numerical implementation. Finally, the restriction  $\hat{\psi} = \hat{\psi}_+$ , which means the identification of the spectral solution representation with the pure outgoing one, forces the whole solution to be purely outgoing. This plays the same role as the direct implementation of the pole condition as extra condition in our real axis approach, compare Section 5.2.

**Part III:** (Continuous Laplace transformed exterior problem):

$$(5.2.6) \quad \begin{aligned} \hat{\psi}(t) - \frac{\pi c_\nu}{t-2i} e^{-r_0 t} w(i) \\ + \frac{c_\nu}{t(t-2i)} \int_0^t dt_1 e^{(t-t_1)r_0} (t-t_1) \hat{\psi}(t_1) = 0. \end{aligned}$$

In this form, the similarity between both approaches, the direct and the cut function approach, becomes apparent. Note the close similarity between the spectral equation (5.2.3) and (5.2.6).

Next we derive the corresponding discrete equations.

### Part I: Discrete Interior Problem.

The discrete interior problem is given again canonically :

$$A_{\text{int}} \mathbf{u} + \begin{bmatrix} \mathbf{0} \\ (\bar{v}u')(r_0) \end{bmatrix} = \mathbf{r}$$

where  $\mathbf{u} \in \mathbb{C}^n$  is the vector of interior unknowns and the right-hand side  $\mathbf{r} \in \mathbb{C}^n$  vanishes identically except for the entries arising from the given interior Dirichlet condition. Next we discuss the discretization of  $\hat{\psi}$  on the real axis  $t \geq 0$ .

### Part II: Discrete Linking Conditions.

Our implementation of the cut function formulation follows a completely different approach than used in the real axis method. Whereas we applied a collocation method based on cubic  $C^1$ -splines in a FEM-like manner in the latter one, we develop the cut function method based on quadrature formulas for ordinary differential equations (ODE's). This is motivated by the observed strong stability of the integral equation (5.2.6) when solved as ODE in forward ( $+\infty$ ) direction. We extend the ODE approach to the integral formulation based on explicit and implicit Runge-Kutta-type collocation schemes. Following the motivation given in [27] and the presentation of Runge-Kutta-type collocation schemes given in [25],

we derive in the following the Runge-Kutta schemes for our basic linear Volterra integral equation.

Let us denote (5.2.6) by

$$\widehat{\psi}(t) = f(t) + \int_0^t dt_1 K(t, t_1) \widehat{\psi}(t_1)$$

with the kernel  $K(t, t_1)$  and  $f(t)$  as given above

$$\begin{aligned} K(t, t_1) &= -\frac{c_\nu}{t(t-2i)} e^{(t-t_1)r_0} (t-t_1) \\ f(t) &= \frac{\pi c_\nu}{t-2i} e^{-r_0 t} w(i). \end{aligned}$$

We define the set of nodes  $\{t_i\}_{i=0}^m$

$$0 = t_0 < t_1 < \dots < t_m < \infty$$

which decompose the positive real axis into  $m$  intervals between  $t_0$  and  $t_m$ . Let us denote  $\mathbf{t} = (t_0, \dots, t_{m-1})^T$ . First, we consider a solution between the points  $t_l$  and  $t_{l+1}$ ,  $0 \leq l < m$  and write

$$\begin{aligned} (5.2.7) \quad \widehat{\psi}(t_l + \tau) &= f(t_l + \tau) + \int_0^{t_l + \tau} dt_1 K(t_l + \tau, t_1) \widehat{\psi}(t_1) \\ &= \sum_{n=0}^{l-1} \int_{t_n}^{t_{n+1}} dt_1 K(t_l + \tau, t_1) \widehat{\psi}(t_1) \\ &\quad + \int_{t_l}^{t_l + \tau} dt_1 K(t_l + \tau, t_1) \widehat{\psi}(t_1) + f(t_l + \tau). \end{aligned}$$

The relation to the integration of ODE's becomes apparent, if we consider the simplified sub-problem

$$\widetilde{\psi}(t_l + \tau) = \widetilde{\psi}(t_l) + \int_{t_l}^{t_l + \tau} dt_1 \widetilde{K}(t_1) \widetilde{\psi}(t_1).$$

This is equivalent to the ODE

$$\frac{d}{d\tau} \widetilde{\psi}(t + \tau) = \widetilde{K}(t + \tau) \widetilde{\psi}(t + \tau), \quad \widetilde{\psi}(t) = \psi_0.$$

Now, a general  $s_{RK}$ -stage Runge-Kutta scheme with consistency order  $p$  has the form (cf. e. g. [25, Sec. 6.2, Eq. 6.9])

$$\begin{aligned} g_i &= \widetilde{\psi}(t_l) + \tau_l \sum_{j=1}^{s_{RK}} a_{ij} K(t_l + c_j \tau_l) g_j, \quad i = 1, \dots, s \\ \widetilde{\psi}(t_l + \tau) &= \widetilde{\psi}(t_l) + \tau \sum_{j=1}^{s_{RK}} b_j K(t_l + c_j \tau) g_j + \mathcal{O}(\tau_l^{p+1}) \quad \tau_l = t_{l+1} - t_l. \end{aligned}$$

As it is common practice, we use the table (Butcher scheme)

$$\frac{\mathbf{c} \mid \mathbf{A}}{\mid \mathbf{b}}$$

with  $(\mathbf{A})_{ij} = a_{ij}$ ,  $\mathbf{c} = (c_1, \dots, c_{s_{RK}})$ ,  $\mathbf{b} = (b_1, \dots, b_{s_{RK}})$  to characterize a Runge-Kutta method. Following Engl [27], we adopt the method to solve our integral

equation. The form of the Runge-Kutta scheme in the form given above, suggests a direct generalization: taking the dependence of the kernel  $K$  from a second variable into account as well as the fact that  $f$  is given explicitly, we obtain the discretization of (5.2.7)

$$\begin{aligned}\widehat{\psi}(t_l + \tau_l c_i) &\approx \sum_{n=0}^{l-1} \tau_n \sum_{j=1}^{s_{RK}} b_j K(t_l + \tau_l c_i, t_n + \tau_n c_j) \widehat{\psi}(t_n + \tau_n c_j) \\ &\quad + \tau_l \sum_{j=1}^{s_{RK}} a_{ij} K(t_l + \tau_l c_i, t_l + \tau_l c_j) \widehat{\psi}(t_l + \tau_l c_j) \\ &\quad + f(t_l + \tau_l c_i).\end{aligned}$$

For the numerical implementation it is useful to define a vector  $\tau$  containing all intermediate time-steps

$$\tau = \mathbf{I}_{m,m} \otimes \mathbf{I}_{s_{RK},1} \mathbf{t} + \Delta t \otimes \mathbf{c}$$

with  $\Delta t$  the vector of interval lengths  $\Delta t = (t_1 - t_0, t_2 - t_1, \dots, t_m - t_{m-1})$  and the kernel-matrix

$$\begin{aligned}(\mathbf{K})_{pq} &= K(\tau_p, \tau_q), \quad \tau_p = (\tau)_p, \tau_q = (\tau)_q \\ K_{pq} &= 0 \quad \text{for } p + s_{RK} - 2 + p \bmod s_{RK} < q.\end{aligned}$$

Then, the integral operator component of the discrete system can be assembled by the sum of a sub-diagonal block matrix representing the quadrature formulas and a block-diagonal matrix corresponding to the Runge-Kutta scheme. Let  $\mathbf{T}$  be a lower triangular matrix, whose sub-diagonal part consists only of unit entries

$$\mathbf{T} = \begin{bmatrix} 0 & & & & \\ 1 & 0 & & & \\ \vdots & \ddots & \ddots & & \\ 1 & \dots & 1 & 0 & \end{bmatrix}, \quad \mathbf{T} \in \mathbb{R}^{m,m}.$$

The discretized integral operator can be realized in vectorized form:

$$\mathbf{M} = \text{diag } \Delta t \otimes \mathbf{A} \cdot * \mathbf{K} + (\mathbf{T} \text{diag } \Delta t) \otimes (\mathbf{I}_{s_{RK},1} \otimes \mathbf{b}) \cdot * \mathbf{K}.$$

The MatLab-type dot-multiplication  $\mathbf{A} \cdot * \mathbf{B}$  denotes element-wise multiplication of the elements of matrices of equal dimension:  $(\mathbf{A} \cdot * \mathbf{B})_{ij} = (\mathbf{A})_{ij} \cdot (\mathbf{B})_{ij}$ . This makes an implementation both simple and effective. For the study of the numerical performance we always used the implicit Runge-Kutta scheme of Gaussian type with  $p = 4$ , in order to have the possibility to compare the cut function approach with our real axis method. The corresponding Butcher scheme is given by

$$\begin{array}{c|cc} \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{4} & \frac{1}{4} - \frac{\sqrt{3}}{6} \\ \frac{1}{2} + \frac{\sqrt{3}}{6} & \frac{1}{4} + \frac{\sqrt{3}}{6} & \frac{1}{4} \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$

However, all the other common Gauss and Radau schemes from  $p = 1$  (implicit Euler method) to the implicit Gauss-method with  $p = 6$  have been implemented

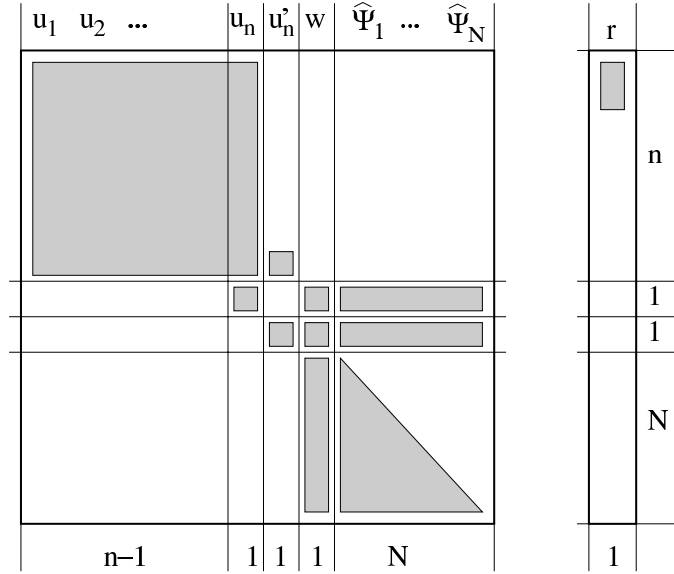


FIGURE 5.2.6. Structure of the algebraic system corresponding to the cut-function approach. Compare Fig. 5.2.2.

and tested successfully. Based on these components, we can express the linking conditions in discrete form

$$\begin{aligned}
 u(r_0) - \frac{w(i)}{2i} + \frac{1}{2\pi i}(\Delta \mathbf{t})^T \otimes \mathbf{b} \hat{\psi} &= 0 \\
 u'(r_0) - \frac{w(i)}{2} + \frac{1}{2\pi i}(i - \tau)^T \cdot * [(\Delta \mathbf{t})^T \otimes \mathbf{b}] \hat{\psi} &= 0.
 \end{aligned}$$

**Part III: Discrete Laplace Transformed Exterior Problem.**

Finally, we obtain the equation for the cut-function along the real axis

$$\hat{\psi} - \frac{\pi c_\nu}{\tau - 2i} e^{-r_0 \tau} w(i) - M \hat{\psi} = 0.$$

The second term in this equation is a vector whose entries has to be evaluated according to the entries of  $\tau$ .

**Numerical studies.** Fig. 5.2.6 shows the occupation pattern of the composite block-matrix system. Its structure is very similar to the one obtained for the real axis approach.

Fig. 5.2.7 presents the convergence result of the cut-function approach. The relative error of the DtN-number for the frequencies  $\nu = 0$  and  $\nu = 10$  are displayed versus the number of steps of the Runge-Kutta method in the spectral domain. This and all the following results refer to implicit Gauss-method with  $p = 4$ . It is seen that the accuracy of the high-frequency mode is lower than the of the zero-frequency mode. Compared with the real axis approach, Fig. 5.2.3, we note the accuracy of the zero-frequency mode in the cut-function approach at the same step numbers is better than in the real axis approach, whereas the situation is reversed for the high-frequency modes.

The essential advantage of the cut-function approach is that it supplies the far-field without additional costs. The corresponding convergence results are displayed

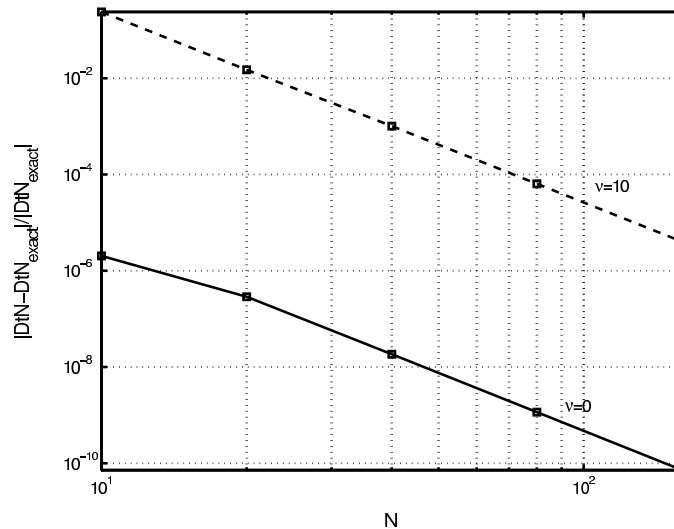


FIGURE 5.2.7. Convergence of the DtN-number in the cut-function approach vs. the number of intervals in the spectral domain.

in Fig. 5.2.8. As far-field we understand here the coefficient of first term of the asymptotic expansion of the Hankel function, normalized such that it matches the value  $u(r_i) = 1$ . Again we consider the modes  $\nu = 0$  and  $\nu = 10$ . Obviously, there is no essential improvement in the accuracy if we decrease the step length (increase the number of intervals). The reason is that the accuracy of the far-field is bounded by the accuracy of the interior problem. To show this, we refined the interior problem uniformly, ending up with a far-field accuracy improved by a factor 1/4. The accuracy of the far-field of the high-frequency mode is much worse than for the low-frequency mode, as expected. Nevertheless, it converges linearly in the log-log plot. The pole strength of the mode, hence the far field amplitude, depends exponentially on the frequency as  $\nu \rightarrow \infty$  if the amplitude on the inner radius of the interior problem is fixed. This result is verified numerically, as shown in Fig. 5.2.9.

The stability of the discrete problem is analyzed as in the real axis approach. First, Fig. 5.2.10 supplies the condition numbers of the discrete problem for different frequencies. Unlike the real axis approach, the condition number grows rapidly with the frequency, until a saturation is reached. If the number of interior nodes is large enough, the condition number of the interior problem dominates the condition number of the composite problem.

The norm of the inverse discrete operator remains bounded as the number of degrees of freedom in the spectral domain, as seen in Fig. 5.2.11. However, for a small number of discretization points, that is far away from the continuous solution, the norm is larger than for a finer discretization. One possibility to reduce the size of the condition numbers for high frequencies is to normalize the pole strength  $w(i)$  by the frequency. In Fig. 5.2.12 we used the normalization  $w(i) \mapsto \tilde{w}(i) := w(i)/(1 + \nu^2)$ . This results in a very effective reduction of the condition number and shows a way how to handle this property in higher dimensional cases.

### 5.3. Non-Separable, Discrete Problems in 2D

In the following we describe a canonical method to derive a finite set of equations approximating the exterior Helmholtz equation by a semi-discrete ansatz.



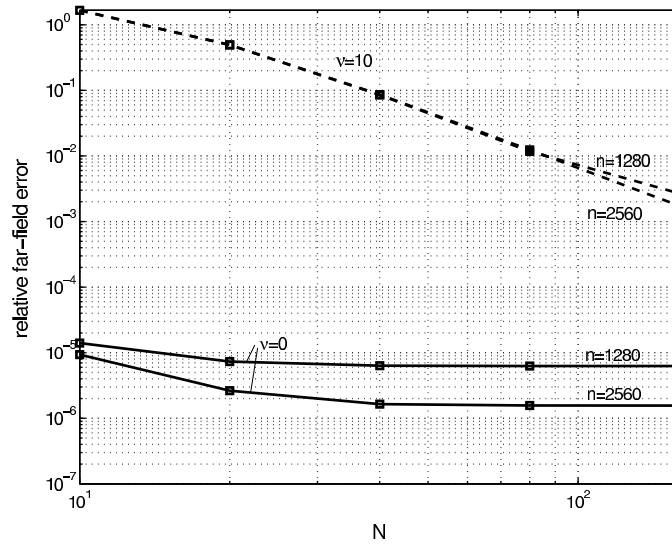


FIGURE 5.2.8. Convergence of the far-field in the cut-function approach vs. the number of intervals in the spectral domain for two different discretizations of the interior problem.

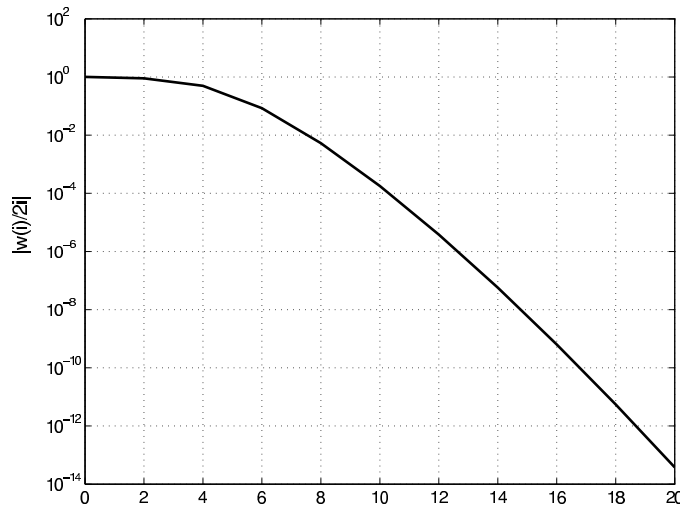


FIGURE 5.2.9. Pole strength as a function of the frequency  $\nu$

The technique is comparable to the method of (vertical) lines for the solution of parabolic problems, cf. the book of Johnson [62, Chapt. 8.3].

#### Basic Idea.

First, we decompose the exterior domain into a finite number of segments, as it is shown in the left-hand side of Fig. 5.3.1. The decomposition is based on straight rays connecting each vertex of the polygonal boundary  $\partial\Omega_{\text{int}}$  with infinity. The rays do not intersect each other. In each semi-infinite quadrilateral obtained this way we approximate the continuous function  $u(x, y)$  by a function  $u_h(x, y)$ , which is discrete in  $y$ , using a local ansatz

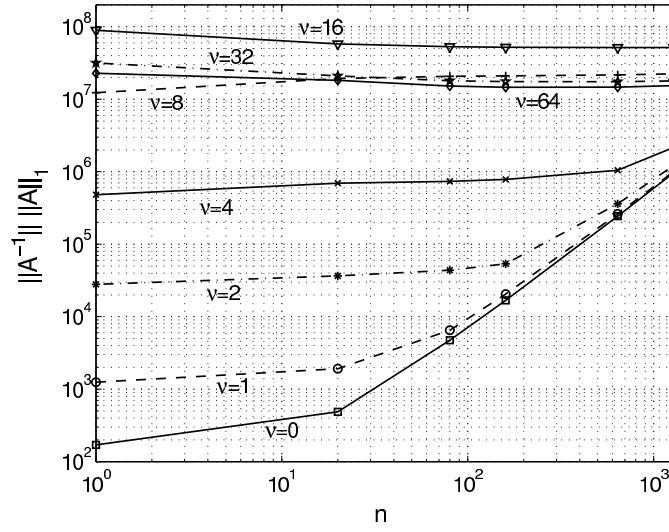


FIGURE 5.2.10. Condition of the composite discrete problem for different frequencies  $\nu$ .

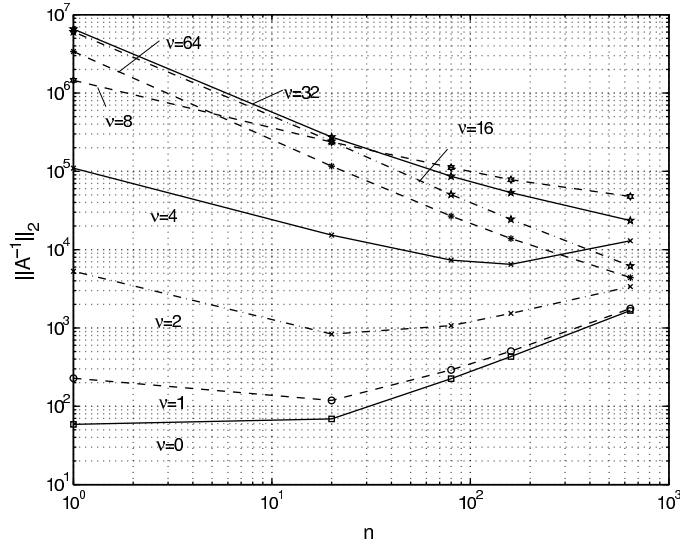


FIGURE 5.2.11. Spectral norm of the inverse discrete operator

$$u_h(x, y) = v_1(x, y)u_1(x) + v_2(x, y)u_2(x),$$

where  $u_1(x)$  and  $u_2(x)$  correspond to the rays  $y_1(x)$  and  $y_2(x)$ , see Fig. 5.3.2. This generalizes the *global* separability used in Sections 5.1–5.2.2 to a *local* separability inside the considered segment. The gray marked quadrilateral in the left-hand side of Fig. 5.3.1 shows that, in our example, the variable  $x$  plays the role of a distance variable and the variable  $y$  behaves as an angular-like variable. The functions  $v_1(x, y)$  and  $v_2(x, y)$  are given linear trial functions defined by

$$v_1(x, y) = \frac{y_2(x) - y}{h(x)} \quad v_2(x, y) = \frac{y - y_1(x)}{h(x)} \quad h(x) = y_2(x) - y_1(x).$$

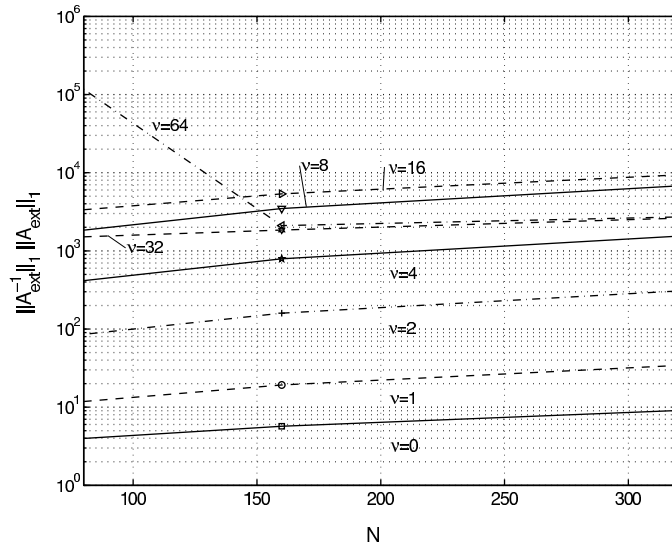


FIGURE 5.2.12. Condition number of the normalized exterior sub-problem.

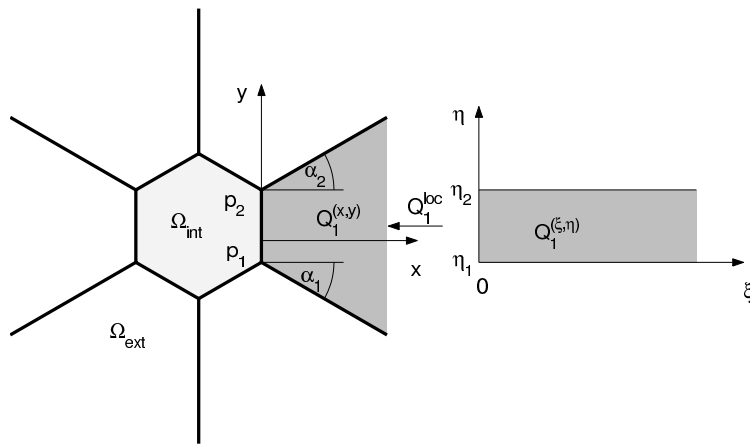


FIGURE 5.3.1. Discretization of the exterior domain by straight rays and mapping from the semi-infinite rectangular domain

This discretization can be extended to curved rays and higher order approximations in the angular-like variable in an obvious way. To obtain a global approximation we have to fit the *local* trial-functions defined on different quadrilaterals such that the resulting *global* trial function is continuous, hence exists at least in  $H_{loc}^1(\Omega_{ext})$ . The whole construction can be done in a FEM-like manner based on the transform of semi-infinite elements. In the following we make this idea more precise.

**Exterior Domain in Transformed Coordinates.** We proceed similarly as in Section 4.1. Let  $\Omega_{int}$  be a polygonally interior domain. Let  $\Omega_{ext} = ext\Omega_{int}$  the corresponding exterior domain and  $\partial\Omega_{int}$  its polygonal, simple and closed boundary

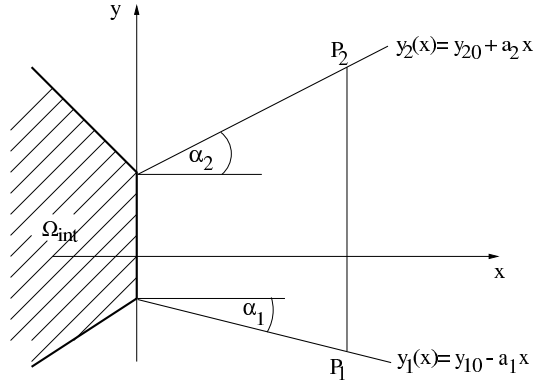


FIGURE 5.3.2. Linear rays belonging to nodes on the boundary ( $a_{1,2} = \tan \alpha_{1,2}$ ).

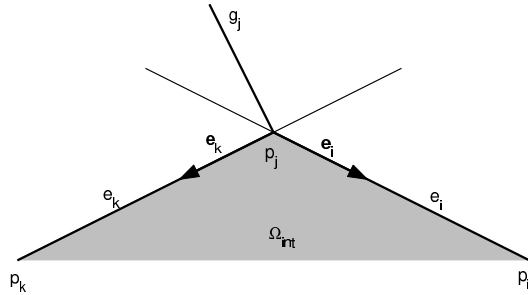


FIGURE 5.3.3. Construction of the ray  $g_j$  associated to vertex  $p_j$  if the interior angle  $\angle(e_i, e_k) < \pi$

with vertices  $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$ . Each vertex  $p_j$  is associated with a straight ray  $g_j$ .

DEFINITION 5.3.1. The set  $G = \{g_1, g_2, \dots, g_N\}$  is an admissible set of straight rays, if it meets the following conditions for any  $j \in 1, \dots, N$ :

- (1) (*Linearity.*) Let  $p_j$  be given. Let the edges  $e_i$  and  $e_k$ ,  $i, k \in 1, \dots, N$ ,  $i \neq k$ , have  $p_j$  as common vertex. Let the second vertex of  $e_i$  be  $p_i$  and the second vertex of  $e_k$  be  $p_k$ . Define two unit vectors  $\mathbf{e}_i = (p_i - p_j) / |p_i - p_j|$  and  $\mathbf{e}_k = (p_k - p_j) / |p_k - p_j|$  such that they are directed away from  $p_j$ .
  - (a) If the segment  $\overline{p_i p_k}$  intersects  $\Omega_{\text{int}}$  (convexity), then the ray  $g_j$  must have a representation  $g_j(\tau) = p_j + \tau(c_i \mathbf{e}_i + c_k \mathbf{e}_k)$  with  $\tau \in \mathbb{R}_+$  and both  $c_i$  and  $c_j$  strictly negative, compare Fig. 5.3.3.
  - (b) If the segment  $\overline{p_i p_k}$  intersects  $\text{ext } \Omega_{\text{int}}$  (concavity), then the same representation must hold with both  $c_i$  and  $c_j$  strictly positive, compare Fig. 5.3.4.
  - (c) If  $p_i, p_j, p_k$  are placed on a common straight line, then the ray  $g_j$  must have a representation  $g_j(\tau) = p_j + \tau(c_i \mathbf{e}_i + c_n \mathbf{n})$  with  $\mathbf{n}$  the unit normal vector directed outwards and  $c_n$  is strictly positive.
- (2) (*Intersection condition.*) The rays neither intersect each other in the exterior domain nor they have common points with  $\Omega_{\text{int}}$ . The only common points of the rays and  $\Omega_{\text{int}}$  are the vertices of the polygonal boundary.

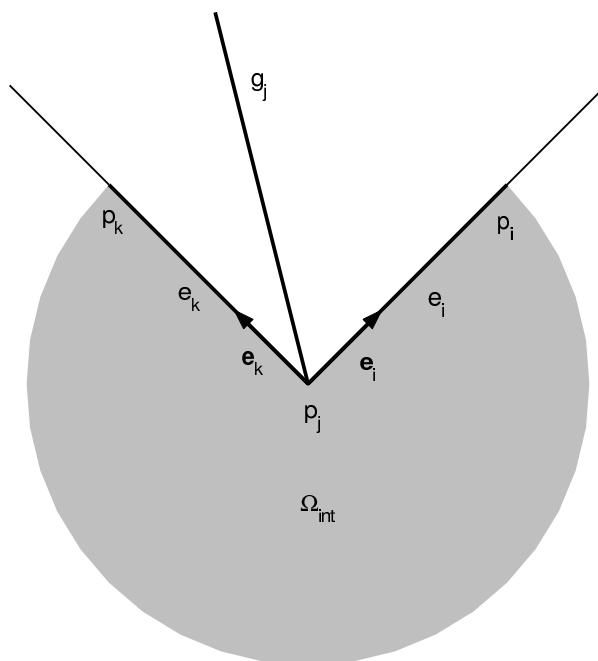


FIGURE 5.3.4. Construction of the ray  $g_j$  associated to vertex  $p_j$  if the interior angle  $\angle(e_i, e_k) > \pi$

- (3) (*Geometric similarity.*) There exist positive constants  $\gamma_1, \dots, \gamma_N$  such that for all  $\tau > 0$  the points  $p'_1 = g_1(\gamma_1\tau), \dots, p'_N = g_N(\gamma_N\tau)$  are the vertices of a new simple polygon, where the vertices are connected in the same order as the corresponding vertices of the original polygon, and the edges of the new polygon and the corresponding edges of the original polygon are *parallel* to each other, cf. Fig. 5.3.5.

Condition 3 is the crucial condition. It enables a direct application of the finite element technology to the exterior domain. We give two schemes which allow the construction of sets of admissible rays for convex and for certain types of non-convex domains. These non-convex domains are star-shaped domains which we want to define as follows:

DEFINITION 5.3.2. A domain  $\Omega$  is star-shaped with respect to a point  $p \in \Omega$ , if all rays starting at  $p$  hit the boundary  $\partial\Omega$  in exactly one point.

1. *Radial rays.* Let a nonempty convex domain be given, compare Fig. 5.3.5. Fix an arbitrary interior point. Connect this interior point by line segments with each of the vertices of the boundary. Extend these line segments to linear rays. This scheme can be extended to star-shaped non-convex domains as follows. Let a star-shaped non-convex domain be given, as it is shown in Fig. 5.3.6. Hence it exists an interior point such that any line segment which connects this interior point with a vertex of the boundary hits the boundary only at this vertex, and we can continue this line segment to the desired ray. By geometrical similarity, the radial ray constructions supplies an admissible set of rays.

2. *Generalized normal rays.* Let a nonempty convex domain be given, compare Fig. 5.3.7. Mark all vertices of the boundary whose aligned edges enclose an interior angle less than  $\pi$ . These are all given vertices  $p_1, \dots, p_4$  in Fig. 5.3.7. Construct the rays successively corresponding to all but the last marked vertex according to

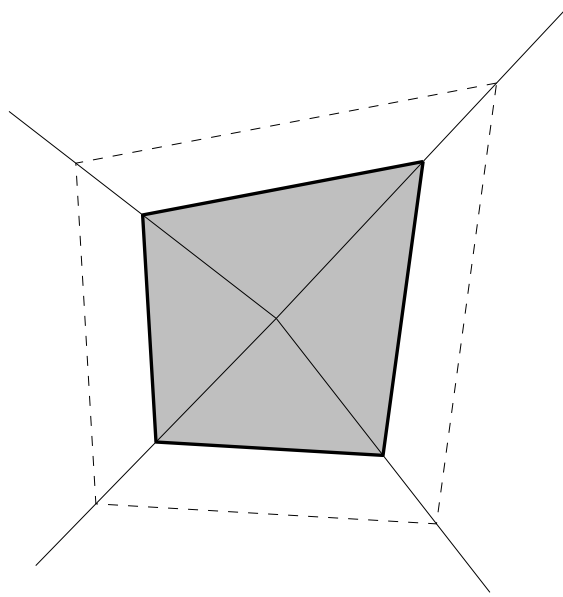


FIGURE 5.3.5. Radial ray construction for convex domains. The dashed line indicates the linearly scaled polygonal boundary. The scaling is based on a single scaling parameter, compare Definition 5.3.1, 3.

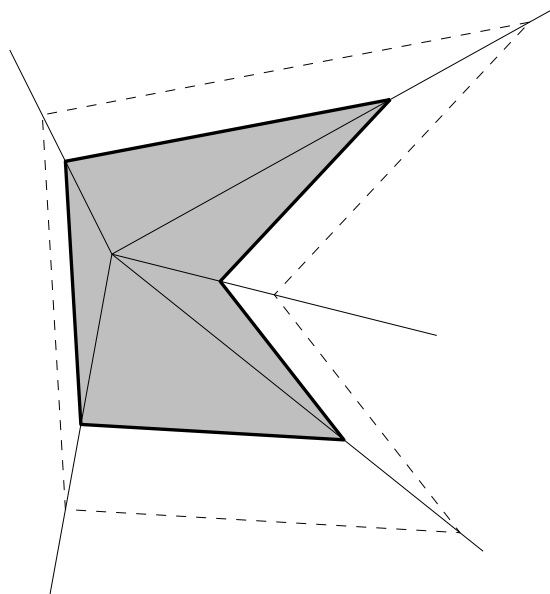


FIGURE 5.3.6. Radial ray construction for a star-shaped concave domain. Compare Fig. 5.3.5.

Definition 5.3.1, 1a). This corresponds to the rays  $g_1, \dots, g_3$  in our example. The last ray  $g_N$  (corresponding to  $g_4$  in Fig. 5.3.7) must be constructed in a special manner:

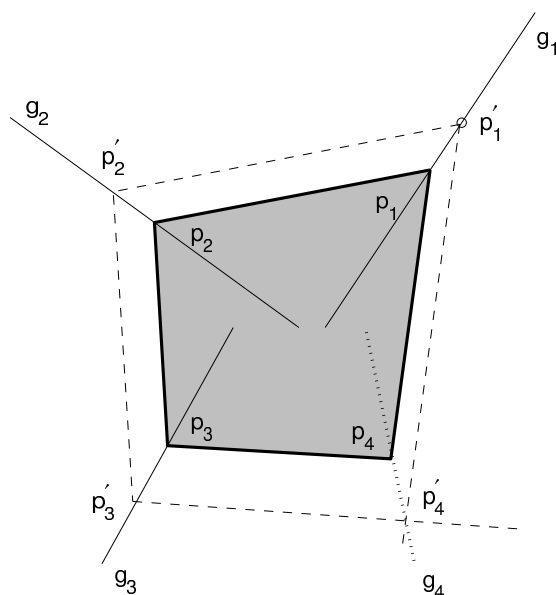


FIGURE 5.3.7. Normal ray construction for convex domains. They rays do not have a common intersection point. The last ray constructed is marked by the dotted line. It must ensure that the scaled polygonal boundary becomes a closed curve.

- (1) Fix an arbitrary point on the ray through the first marked vertex (in the example point  $p'_1$ ).
- (2) Construct the two lines which go through this point, and which are parallel to their corresponding boundary segments (the two dashed lines intersecting in  $p'_1$ ).
- (3) Determine the intersection point between this line and the next ray (the point  $p'_2$ ).
- (4) Continue the construction, starting from the intersection point and moving in positive direction.
- (5) The last ray must be constructed such that it goes through the last marked vertex and the intersection of the first and the last line constructed this way (the dotted ray in our example).

Construct the remaining rays according to Definition 5.3.1, 1c) and 2). The notion “normal” arises from the fact that the generalized normal rays tend to be scaled normal vectors of the boundary if the boundary tends to be smooth curve *and* the rays are constructed such that they halve the exterior angle between the aligned edges. Other schemes, e.g. a combination of the normal and radial ray scheme, are possible. A direct consequence from the geometric construction is the following

**LEMMA.** *Both normal and radial ray construction supply an admissible set of straight rays.*

A decomposition  $\mathcal{Q} = \{Q_1, Q_2, \dots, Q_N\}$  of  $\overline{\Omega_{\text{ext}}}$  based on  $N$  semi-infinite quadrilaterals  $Q_{1, \dots, N}$ , where each quadrilateral is bounded by two rays and one edge of  $\partial\Omega_{\text{int}}$ , provides the starting point for both the variational formulation and the finite element discretization. Fig. 5.3.1 illustrates the basic scheme.

**DEFINITION 5.3.3.** A decomposition  $\mathcal{Q}$  based on straight lines (rays) is called admissible, if the following properties are met:

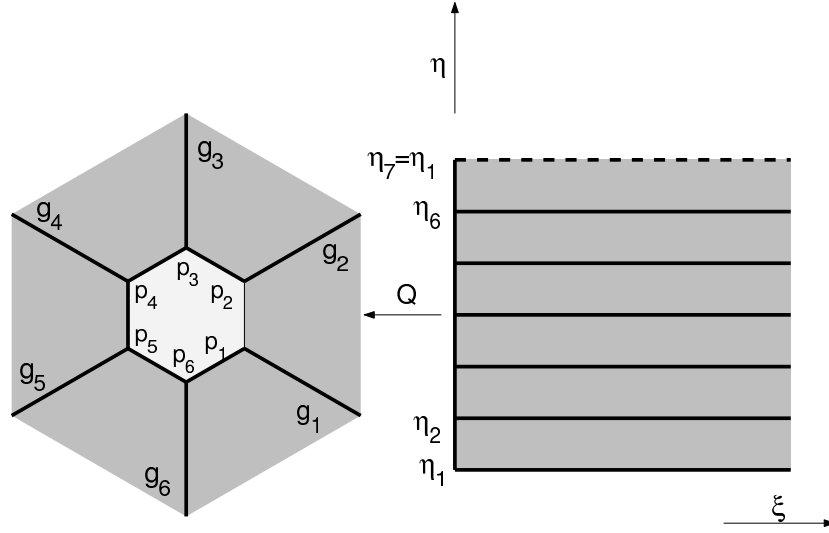


FIGURE 5.3.8. Mapping from the semi-infinite rectangular reference domain onto the exterior of the given problem. The dashed line indicates that the  $\eta$ -interval continuous periodically.

- (1) Two nodes of each  $Q_j$  are also nodes of  $\partial\Omega_{\text{int}}$ .
- (2)  $\bigcup_{j=1}^N Q_j = \overline{\Omega_{\text{ext}}}$ .
- (3) If  $i \neq j$ , then the intersection  $Q_i \cap Q_j$  is either empty or  $Q_i \cap Q_j = g_j(x, y)$ , where  $g_j(x, y)$  is a straight ray connecting the vertex  $p_j$  of  $\partial\Omega_{\text{int}}$  with infinity. This boundary vertex is a common point of  $Q_i$  and  $Q_j$ .

We consider each semi-infinite quadrilateral as image of a semi-infinite rectangle under a appropriate diffeomorphism. Let us denote the  $j$ th semi-infinite rectangle by  $Q_j^{(\xi, \eta)}$ , where  $\xi$  and  $\eta$  are the Cartesian coordinates of the pre-image –  $\xi$  playing the role of a generalized distance variable and  $\eta$  a role as a generalized angular variable – and by  $Q_j^{(x, y)}$  the  $j$ th semi-infinite quadrilateral defined above. For clarity, we added the superscript  $(x, y)$  to indicate that the image is defined in the original  $x, y$ -coordinates of the Cartesian reference systems. We denote the  $j$ th diffeomorphism by

$$Q_j^{(\text{loc})} : Q_j^{(\xi, \eta)} \rightarrow Q_j^{(x, y)}.$$

These local mappings are connected with each other such that a global homeomorphism

$$Q : \overline{\Omega_{\text{ext}}^{(\xi, \eta)}} \rightarrow \overline{\Omega_{\text{ext}}^{(x, y)}}$$

results. Moreover, we connect the local domains and mappings in a way that the global mapping  $Q$  meets the following crucial conditions

- (1)  $Q$  is at least twice continuously differentiable in  $\xi$ .
- (2)  $Q$  is continuous and piecewise continuously differentiable in  $\eta$ .

The latter fact is important since it allows a standard variational formulation in  $\eta$  with finite-dimensional test and trial spaces based on continuous  $C^0$ -elements. Fig. 5.3.8 displays this global mapping in a schematic diagram.



### Variational Formulation in Transformed Coordinates.

First we restate the variational formulation of the interior problem (4.0.9) in a slightly modified form. Without interior boundaries and sources, but with a given sufficiently regular normal derivative  $\mathbf{n}(s)\nabla u(s)$  on  $\partial\Omega_{\text{int}}$ , the variational form of the interior problem reads: Find  $u \in H^1(\Omega_{\text{int}})$  such that for all  $v \in H^1(\Omega_{\text{int}})$

$$(5.3.1) \quad \text{Interior problem :} \quad a(v, u) = \int_{\partial\Omega_{\text{int}}} ds \bar{\nu}(s) \mathbf{n}(s) \nabla u(s).$$

Again,  $a(v, u)$  denotes our standard sesquilinear form

$$\begin{aligned} a & : H^1(\Omega_{\text{int}}) \times H^1(\Omega_{\text{int}}) \rightarrow \mathbb{C} \\ a(v, u) & = (\nabla v, \nabla u) - (v, k^2(\mathbf{x})u). \end{aligned}$$

In the following, we derive a corresponding variational formulation of the exterior problem formulated in transformed coordinates. The essential difference as compared to the interior problem is that the variational (weak) form is used in  $\eta$ -direction only, whereas in  $\xi$ -direction the original (strong) form of the Helmholtz equation is used. In its core, this is a method of lines type formulation, which *formally* represents the exterior problem as elliptic initial value problem with initial data given on  $\partial\Omega_{\text{ext}}^{(\xi, \eta)}$ . Given three sesquilinear forms

$$a_i(v, u) : H_\pi^1[\eta_{\min}, \eta_{\max}] \times W^i(\Omega_{\text{ext}}^{(\xi, \eta)}) \rightarrow \mathbb{C}, \quad i \in \{0, 1, 2\},$$

which will be construct below, the variational form of the exterior problem reads: Find  $u \in W^2(\Omega_{\text{ext}}^{(\xi, \eta)})$  such that for all  $v \in H_\pi^1[\eta_{\min}, \eta_{\max}]$  and  $\xi \in \mathbb{R}_+$

$$(5.3.2) \quad \text{Exterior problem :} \quad \begin{aligned} a_0(v, u) + a_1(v, \partial_\xi u) + a_2(v, \partial_\xi^2 u) & = 0 \\ u(0) & = u_D \\ u'(0) & = u_N. \end{aligned}$$

The function spaces  $W^i(\Omega_{\text{ext}}^{(\xi, \eta)})$  might be defined as

$$W^i(\Omega_{\text{ext}}^{(\xi, \eta)}) := \left\{ \begin{array}{l} w(\xi_0, \eta) \in H_\pi^1[\eta_{\min}, \eta_{\max}] : \xi_0 \in \mathbb{R}, \text{ fixed and} \\ w(\xi, \eta_0) \in C^i(\mathbb{R}_+) : \eta_0 \in [\eta_{\min}, \eta_{\max}], \text{ fixed} \end{array} \right\}, \quad i \in \{0, 1, 2\}.$$

The construction of  $a_i(\cdot, \cdot)$  proceeds as follows:

- (1) The local mappings  $Q_j^{(\text{loc})}$ ,  $j = 1, \dots, N$  from the  $\xi, \eta$ -system to the  $x, y$ -system are derived.
- (2) The local mappings are glued together to obtain the global mapping  $Q : \overline{\Omega_{\text{ext}}^{(\xi, \eta)}} \rightarrow \overline{\Omega_{\text{ext}}^{(x, y)}}$ .
- (3) The Helmholtz equation is represented in  $\xi, \eta$ -coordinates.
- (4) The weak form with respect to  $\eta$  is derived.

The following geometrical consideration supplies the  $j$ th local mapping  $Q_j^{(\text{loc})}$ , compare Fig. 5.3.1. Let the set of vertices  $\mathcal{P}$  be assigned to the set  $\{\eta_1, \dots, \eta_N\}$  of generalized angular coordinates. The values of  $\eta_j$  might be derived from the arc length of the boundary (for 2D problems) or from the angle coordinate of a cylindrical coordinate system. Let us consider the semi-infinite quadrilateral  $Q_j^{(x, y)}$  as marked in Fig 5.3.1 by the gray area. It is bounded by the segment  $\overline{p_1 p_2}$ ,  $p_1, p_2 \in \mathcal{P}$ , and the rays  $g_1, g_2$ . Let us denote the unit vectors defining the direction of the rays  $g_1, g_2$  by  $\mathbf{g}_1$  and  $\mathbf{g}_2$ . Since  $g_1, g_2 \in \mathcal{G}$ , they possess a parameter representation

$$\begin{aligned} g_1(\tau) &= p_1 + \tau \mathbf{g}_1 \\ g_2(\tau) &= p_2 + \gamma \tau \mathbf{g}_2, \tau \in \mathbb{R}_+ \end{aligned}$$

and it must hold  $\overline{g_1(\tau)g_2(\tau)} \parallel \overline{p_1p_2}$ , compare the definition of the admissible rays Definition 5.3.1. Let the unit vector  $\mathbf{n}$  be normal to  $\overline{p_1p_2}$ . According to Definition 5.3.1, we have to choose  $\gamma$  such that  $\overline{g_1(\tau)g_2(\tau)} \perp \mathbf{n}$ , it follows  $\gamma = \cos \alpha_1 / \cos \alpha_2$ . So far,  $\tau$  is an arbitrary real parameter. We want to replace it by a distance variable  $\xi$  with a direct geometrical meaning. We define  $\xi$  as the distance between the straight line through  $p_1$  and  $p_2$  and the straight line through  $g_1(\tau)$  and  $g_2(\tau)$ . It can be seen from Fig. 5.3.1 that  $\xi = \tau \cos \alpha_1$ . This way we obtain a symmetric form of the parameter representation

$$\begin{aligned} g_1(\xi) &= p_1 + \frac{\xi}{\zeta \cos \alpha_1} \mathbf{g}_1 \\ g_2(\xi) &= p_2 + \frac{\xi}{\zeta \cos \alpha_2} \mathbf{g}_2, \xi \in \mathbb{R}_+, \end{aligned}$$

where the new scaling parameter is needed to fit the different local coordinate systems continuously to each other. We define the signs of the angles  $\alpha_{1,2}$  such that they are positive if the corresponding rays and the lengthened normal through the mid-point of boundary segment diverge in the semi-plane right from the lengthened boundary segment.

So far, a point given on  $g_1$  (or on  $g_2$ ) is uniquely related to the generalized distance  $\xi$ . Next we assign a given point *between* the rays  $g_1$  and  $g_2$  to a generalized distance  $\xi$  and a generalized angle  $\eta$ . To this end let  $\eta$  on  $g_1$  be equal to  $\eta_1$  and  $\eta$  on  $g_2$  equal to  $\eta_2$ . We perform a linear interpolation between these values inside  $Q_j^{(x,y)}$ . Consequently, a point  $p(x,y) \in Q_j^{(x,y)}$  can be expressed by

$$p(x,y) = g_1(\xi) + \frac{\eta - \eta_1}{\eta_2 - \eta_1} (g_2(\xi) - g_1(\xi)).$$

Further, with unit vectors  $\mathbf{g}_1$  and  $\mathbf{g}_2$  given in Cartesian coordinates

$$\mathbf{g}_1 = \begin{pmatrix} \cos \beta_1 \\ \sin \beta_1 \end{pmatrix} \quad \mathbf{g}_2 = \begin{pmatrix} \cos \beta_2 \\ \sin \beta_2 \end{pmatrix},$$

we can express  $p(x,y)$  in Cartesian coordinates too, which at the same time supplies the desired mapping

$$(5.3.3) \quad \begin{pmatrix} x \\ y \end{pmatrix} =: Q_j^{(\text{loc})}(\xi, \eta) \quad \text{with}$$

$$Q_j^{(\text{loc})}(\xi, \eta) = \left(1 - \frac{\eta - \eta_1}{\eta_2 - \eta_1}\right) \left[ \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} + \frac{\xi}{\zeta \cos \alpha_1} \begin{pmatrix} \cos \beta_1 \\ \sin \beta_1 \end{pmatrix} \right] + \frac{\eta - \eta_1}{\eta_2 - \eta_1} \left[ \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} + \frac{\xi}{\zeta \cos \alpha_2} \begin{pmatrix} \cos \beta_2 \\ \sin \beta_2 \end{pmatrix} \right].$$

LEMMA 5.3.4. *The mapping  $Q_j^{(\text{loc})} : Q_j^{(\xi, \eta)} \rightarrow Q_j^{(x, y)}$  is a local  $C^\infty$ -diffeomorphism.*

PROOF. We refer to the local inverse mapping theorem, cf. [96, Theorem 4.F, pp. 259]. We have to show

- (i)  $Q_j^{(\text{loc})}$  is a  $C^\infty$ -map on some neighborhood of all points  $(\xi_0, \eta_0) \in Q_j^{(\xi, \eta)}$ .
- (ii)  $\left(Q_j^{(\text{loc})}\right)'(\xi_0, \eta_0) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is bijective.

Ad (i). The map  $Q_j^{(\text{loc})}(\xi, \eta)$  is bilinear in  $\xi$  and  $\eta$ . Hence the property is obvious.

Ad (ii). The property is invariant with respect to translation and rotation. Hence, without loss of generalization, we may choose  $x_1 = x_2 = 0$ ,  $y_1 = 0$ ,  $y_2 = h$ ,  $\eta_1 = 0$ ,  $\eta_2 = h$ ,  $\beta_1 = -\alpha_1$  and  $\beta_2 = \alpha_2$ . Then the  $j$ th local mapping simplifies to

$$Q_j^{(\text{loc})}(\xi, \eta) = \left(1 - \frac{\eta}{h}\right) \frac{\xi}{\zeta} \begin{pmatrix} 1 \\ -\tan \alpha_1 \end{pmatrix} + \frac{\eta}{h} \left[ \begin{pmatrix} 0 \\ h \end{pmatrix} + \frac{\xi}{\zeta} \begin{pmatrix} 1 \\ \tan \alpha_2 \end{pmatrix} \right].$$

Introducing  $\mathbf{h} := (d\xi, d\eta)^T$  we compute

$$\left(Q_j^{(\text{loc})}\right)' \mathbf{h} = Q_j^{(\text{loc})}(d\xi, \eta_0) + Q_j^{(\text{loc})}(\xi_0, d\eta)$$

and it follows,

$$J_j := \left(Q_j^{(\text{loc})}\right)' = \begin{pmatrix} \frac{1}{\zeta} & 0 \\ -\frac{1-\eta_0}{\zeta} \tan \alpha_1 + \frac{\eta_0}{\zeta} \tan \alpha_2 & 1 + \frac{\xi_0}{h\zeta} (\tan \alpha_1 + \tan \alpha_2) \end{pmatrix}.$$

The determinant of the  $j$ th Jacobian  $J_j$  becomes

$$|J_j| = \frac{h\zeta + \xi_0 (\tan \alpha_1 + \tan \alpha_2)}{h\zeta^2}.$$

By construction,  $\alpha_1, \alpha_2 \geq 0$ , and both  $h$  and  $\zeta$  are strictly positive. Hence  $|J_j|$  is always positive, and the map is bijective.  $\square$

Next we glue two neighbored domains  $\overline{Q}_j^{(\xi, \eta)}$  and  $\overline{Q}_{j+1}^{(\xi, \eta)}$  together and set the free scaling parameters  $\zeta_j$  and  $\zeta_{j+1}$  such that the resulting mapping from the interior of the combined  $\xi, \eta$ -domains onto the corresponding combined  $x, y$ -domains remains  $C^\infty$  with respect to  $\xi$  and becomes globally  $C^0$  with respect to  $\eta$ . We denote the interior of the composite domains by  $Q_{j,j+1}^{(\xi, \eta)}$  and  $Q_{j,j+1}^{(x, y)}$ , respectively, and label the quantities  $\alpha, \beta$  and  $\zeta$  occurring in different subdomains by the index of the subdomain.

LEMMA 5.3.5. *Let  $\overline{Q}_j^{(\xi, \eta)}$  and  $\overline{Q}_{j+1}^{(\xi, \eta)}$  be two neighbored domains, and  $\overline{Q}_j^{(x, y)}$  and  $\overline{Q}_{j+1}^{(x, y)}$  their counterparts, separated by the ray  $g_{j+1}$ . The composite mapping*

$$Q_{j,j+1}^{(\text{loc})} : \text{int} \left( \overline{Q}_j^{(\xi, \eta)} \cup \overline{Q}_{j+1}^{(\xi, \eta)} \right) \rightarrow \text{int} \left( \overline{Q}_j^{(x, y)} \cup \overline{Q}_{j+1}^{(x, y)} \right)$$

*composed of the local mappings (5.3.3) for domains  $j$  and  $j+1$  is a homeomorphism if and only if*

$$(5.3.4) \quad \zeta_{j+1} \cos \alpha_{1,j+1} = \zeta_j \cos \alpha_{2,j}.$$

*Moreover  $Q_{j,j+1}^{(\text{loc})}$  is  $C^\infty$  with respect to  $\xi$ .*

PROOF. We consider a neighborhood of the ray  $g_{j+1}$ . The ray possesses the representations

$$g_{j+1}(\xi) = p_{j+1} + \frac{\xi}{\zeta_j \cos \alpha_{2,j}} \mathbf{e}_{j+1}, \quad \text{with } \xi \in \mathbb{R}_+$$

$$\text{and } g_{j+1}(\xi) = p_{j+1} + \frac{\xi}{\zeta_{j+1} \cos \alpha_{1,j+1}} \mathbf{e}_{j+1}.$$

Thus the condition  $\zeta_{j+1} \cos \alpha_{1,j+1} = \zeta_j \cos \alpha_{2,j}$  ensures the continuity

$$\lim_{\eta \uparrow \eta_{j+1}} Q_j^{(\text{loc})}(\xi, \eta) = \lim_{\eta \downarrow \eta_{j+1}} Q_{j+1}^{(\text{loc})}(\xi, \eta) \text{ for all } \xi \in \mathbb{R}_+.$$

Furthermore, it holds

$$\lim_{\eta \uparrow \eta_{j+1}} \partial_\xi Q_j^{(\text{loc})}(\xi, \eta) = \frac{1}{\zeta_j \cos \alpha_{2,j}} \begin{pmatrix} \cos \beta_{j+1} \\ \sin \beta_{j+1} \end{pmatrix}$$

as well as

$$\lim_{\eta \downarrow \eta_{j+1}} \partial_\xi Q_{j+1}^{(\text{loc})}(\xi, \eta_j) = \frac{1}{\zeta_{j+1} \cos \alpha_{1,j+1}} \begin{pmatrix} \cos \beta_{j+1} \\ \sin \beta_{j+1} \end{pmatrix}.$$

This shows that in the limit  $\eta \rightarrow \eta_{j+1}$  the derivatives of the local maps with respect to the generalized distance variable  $\xi$  coincide at both sides of the ray  $g_{j+1}$  to any order, if the angle condition (5.3.4) is met.  $\square$

Having established the local mapping properties between the domains  $Q_j^{(\xi, \eta)}$  and  $Q_j^{(x, y)}$  and the mapping properties of the composite mapping between two neighboring domains, we obtain the desired global mapping property.

**COROLLARY 5.3.6.** *Let  $\overline{Q}_{\text{ext}}^{(\xi, \eta)} = \cup_{j=1}^N \overline{Q}_j^{(\xi, \eta)}$  and  $\overline{Q}_{\text{ext}}^{(x, y)} = \cup_{j=1}^N \overline{Q}_j^{(x, y)}$ . The composite mapping*

$$Q : \overline{\Omega}_{\text{ext}}^{(\xi, \eta)} \rightarrow \overline{\Omega}_{\text{ext}}^{(x, y)}$$

*composed of the local mappings (5.3.3) for domains  $j = 1, \dots, N$  is a homeomorphism if one scaling parameter of the set  $\{\zeta_1, \dots, \zeta_N\}$  is given and the other satisfy the angle condition of Lemma 5.3.5.*

### Sesquilinear Forms in Transformed Coordinates.

The point of departure is the Laplacian in  $\xi, \eta$ -coordinates (4.1.2), which allows us to rewrite the Helmholtz equation in  $\xi, \eta$ -coordinates. To simplify notation, we do not introduce a new symbol for the solution  $u$  represented in  $\xi, \eta$ -coordinates. The transformed Helmholtz equation reads

$$\nabla_{\xi\eta}^T (J^{-1} J^{-T} |J| \nabla_{\xi\eta}) u + |J| k^2 u = 0.$$

Let us abbreviate  $F := J^{-1} J^{-T} |J|$ . Let  $V[\eta_{\min}, \eta_{\max}]$  denote the trial space of periodic, real functions  $v$  on the interval  $[\eta_{\min}, \eta_{\max}]$ . We multiply the transformed Helmholtz equation with  $v \in V$ , where  $V$  is a suitable space of test functions periodic on  $[\eta_{\min}, \eta_{\max}]$  and integrate over the interval  $\Gamma_\eta = [\eta_{\min}, \eta_{\max}]$

$$\int_{\Gamma_\eta} d\eta v(\eta) (\partial_\xi, \partial_\eta) F \nabla_{\xi\eta} u + \int_{\Gamma_\eta} d\eta |J| k^2 v(\eta) u = 0.$$

We rewrite this with a coordinate representation of the transposed gradient separating its  $\partial_\xi$ - and  $\partial_\eta$ -part

$$(5.3.5) \quad \int_{\Gamma_\eta} d\eta v(\eta) (\partial_\xi, 0) F \begin{pmatrix} \partial_\xi u \\ \partial_\eta u \end{pmatrix} \\ + \int_{\Gamma_\eta} d\eta v(\eta) (0, \partial_\eta) F \begin{pmatrix} \partial_\xi u \\ \partial_\eta u \end{pmatrix} + \int_{\Gamma_\eta} d\eta v(\eta) |J| k^2 u = 0.$$

Denoting the elements of the matrix  $F$  by  $F_{ij}$ ,  $i, j = 1, 2$  we compute, with respect to the first term of (5.3.5),

$$\begin{aligned}
& \int_{\Gamma_\eta} d\eta v(\eta) (\partial_\xi, 0) F (\partial_\xi u, \partial_\eta u)^T \\
&= \int_{\Gamma_\eta} d\eta v(\eta) [F_{11} \partial_{\xi\xi} u + F_{12} \partial_{\xi\eta} u + (\partial_\xi F_{11}) \partial_\xi u + (\partial_\xi F_{12}) \partial_\eta u] \\
&= \int_{\Gamma_\eta} d\eta v(\eta) [F_{11} \partial_{\xi\xi} u + (\partial_\xi F_{11}) \partial_\xi u + (\partial_\xi F_{12}) \partial_\eta u] - \int_{\Gamma_\eta} d\eta \partial_\eta (v F_{12}) \partial_\xi u.
\end{aligned}$$

Integration by parts supplies for the second term of (5.3.5)

$$\int_{\Gamma_\eta} d\eta v(\eta) (0, \partial_\eta) F (\partial_\xi u, \partial_\eta u)^T = - \int_{\Gamma_\eta} d\eta \partial_\eta v [F_{21} \partial_\xi u + F_{22} \partial_\eta u].$$

Here we used the fact that the elements of the space  $V$  are periodic functions, so no boundary terms occur. Collecting the terms which correspond to the same derivatives of  $u$ , we find the desired sesquilinear forms with  $v = v(\eta)$ ,  $u = u(\xi, \eta)$  as follows

$$(5.3.6) \quad a_2(v, \partial_{\xi\xi} u) = \int_{\Gamma_\eta} d\eta v F_{11} \partial_{\xi\xi} u$$

$$(5.3.7) \quad a_1(v, \partial_\xi u) = \int_{\Gamma_\eta} d\eta v \partial_\xi F_{11} \partial_\xi u - \int_{\Gamma_\eta} d\eta [\partial_\eta (v F_{12}) + (\partial_\eta v) F_{21}] \partial_\xi u$$

$$(5.3.8) \quad a_0(v, u) = \int_{\Gamma_\eta} d\eta v \partial_\xi F_{12} \partial_\eta u - \int_{\Gamma_\eta} d\eta \partial_\eta v F_{22} \partial_\eta u \\ + \int_{\Gamma_\eta} d\eta v |J| k^2 u.$$

Based on these forms, we derive a FEM-discretization along the  $\eta$ -coordinate in standard fashion.

### Finite Element Discretization.

Let us denote the semi-discretized function  $u(\xi, \eta)$  by  $u_h(\xi, \eta)$ . We introduce the discrete, *local* separation ansatz

$$u_h(\xi, \eta) = \sum_{j=1}^N v_j(\eta) u_j(\xi)$$

with functions  $v_j \in V$ , thus employing a discretization of Galerkin-type. We consider the sesquilinear forms (5.3.6)-(5.3.8) one after the other. It follows from the local separation ansatz that

$$\begin{aligned}
A_{2,ij} \partial_{\xi\xi} u_j &:= a_2(v_i, v_j \partial_{\xi\xi} u_j) \\
&= \left( \int_{\Gamma_\eta} d\eta v_i v_j F_{11} \right) \partial_{\xi\xi} u_j.
\end{aligned}$$

Hence we obtain

$$A_{2,ij} = \int_{\Gamma_\eta} d\eta v_i v_j F_{11}.$$

The same way we get

$$\begin{aligned}
A_{1,ij} \partial_\xi u_j &:= a_1(v_i, v_j \partial_\xi u_j) \\
&= \left( \int_{\Gamma_\eta} d\eta v_i \partial_\xi F_{11} v_j - \int_{\Gamma_\eta} d\eta \partial_\eta (v_i F_{12} + v_i F_{21}) v_j \right) \partial_\xi u_j \\
&= \left( \int_{\Gamma_\eta} d\eta v_i [\partial_\xi F_{11} v_j + F_{12} \partial_\eta v_j] - \int_{\Gamma_\eta} d\eta \partial_\eta v_i F_{21} v_j \right) \partial_\xi u_j.
\end{aligned}$$

The last equation follows by an integration by parts with respect to the term containing the factor  $F_{12}$ . This supplies a convenient and easy to implement form of the FEM-matrix

$$A_{1,ij} = \int_{\Gamma_\eta} d\eta (v_i v_j \partial_\xi F_{11} + v_i \partial_\eta v_j F_{12} - \partial_\eta v_i v_j F_{21}).$$

Finally, we compute the FEM-matrix which corresponds to  $u_j(\xi)$

$$\begin{aligned}
A_{0,ij} u_j &:= a_0(v_i, v_j u_j) \\
&= \left( \int_{\Gamma_\eta} d\eta v_i \partial_\xi F_{12} \partial_\eta v_j - \int_{\Gamma_\eta} d\eta \partial_\eta v_i F_{22} \partial_\eta v_j + \int_{\Gamma_\eta} d\eta |J| k^2 v_i v_j \right) u_j
\end{aligned}$$

supplying

$$A_{0,ij} = \int_{\Gamma_\eta} d\eta (v_i \partial_\eta v_j \partial_\xi F_{12} - \partial_\eta v_i \partial_\eta v_j F_{22} + v_i v_j |J| k^2).$$

### Linear $C^0$ - Elements.

We compute the discrete scheme based on simple linear  $C^0$ - elements. The generalization to higher order elements follows as usual. For simplicity, we assume a strict ordering of the set of vertices  $\mathcal{P}$  obtained if one follows the boundary in positive direction. The linear nodal basis function  $v_i(\eta)$  associated with vertex  $p_i$  is given by

$$v_i(\eta) = \begin{cases} (\eta - \eta_{i-1}) / (\eta_i - \eta_{i-1}) & \text{for } \eta \in [\eta_{i-1}, \eta_i] \\ (\eta_{i+1} - \eta) / (\eta_{i+1} - \eta_i) & \text{for } \eta \in [\eta_i, \eta_{i+1}] \\ 0 & \text{elsewhere} \end{cases}$$

We derive local FEM-matrices from the unit element

$$v^{(0)}(\eta) = \begin{cases} v_1^{(0)} & \text{for } \eta \in [-1, 0] \\ v_2^{(0)} & \text{for } \eta \in [0, 1] \\ 0 & \text{elsewhere} \end{cases}, \quad v_1^{(0)} = \eta, \quad v_2^{(0)} = 1 - \eta.$$

On the segment  $\overline{p_i p_{i+1}}$  we introduce the normalization  $\tilde{\eta} = (\eta - \eta_i) / h$ ,  $h = \eta_{i+1} - \eta_i$ . Accordingly, we transform segment-wise  $\tilde{F}(\xi, \tilde{\eta}) = F(\xi, (\eta - \eta_i) / h)$  and  $\tilde{J}(\xi, \tilde{\eta}) = J(\xi, (\eta - \eta_i) / h)$ . Further, let us renumber the vertices of the segment:  $p_1 := p_i$ ,  $p_2 = p_{i+1}$  and drop the vertices-counting indices at the local matrices  $M_0, M_1, M_2$ . Thus we can reformulate these matrices, with  $(i, j) \in \{1, 2\}$  as

$$\begin{aligned}
A_2 &= h \int_0^1 d\tilde{\eta} v_i^{(0)} v_j^{(0)} \widetilde{F_{11}} \\
A_1 &= h \int_0^1 d\tilde{\eta} \left[ v_i^{(0)} v_j^{(0)} \partial_\xi \widetilde{F_{11}} + \frac{1}{h} v_i^{(0)} \partial_{\tilde{\eta}} v_j^{(0)} \widetilde{F_{12}} - \frac{1}{h} \partial_{\tilde{\eta}} v_i v_j \widetilde{F_{21}} \right] \\
A_0 &= h \int_0^1 d\tilde{\eta} \left[ \frac{1}{h} v_i^{(0)} \partial_{\tilde{\eta}} v_j^{(0)} \partial_\xi \widetilde{F_{12}} - \frac{1}{h^2} \partial_{\tilde{\eta}} v_i^{(0)} \partial_{\tilde{\eta}} v_j^{(0)} \widetilde{F_{22}} + v_i^{(0)} v_j^{(0)} |\tilde{J}| k^2 \right]
\end{aligned}$$

From now, we drop all the waves indicating the normalized argument the dependence of the normalized argument  $\tilde{\eta}$ . We introduce the  $2 \times 2$  block matrix  $W$

$$W = \begin{pmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{pmatrix}$$

whose elements are the following  $2 \times 2$  matrices

$$W_{11} = v^{(0)} v^{(0)T}, \quad W_{12} = v^{(0)} \partial_\eta v^{(0)T}, \quad W_{21} = \partial_\eta v^{(0)} v^{(0)T}, \quad W_{22} = \partial_\eta v^{(0)} \partial_\eta v^{(0)T}.$$

This allows a more compact notation of the elementary matrices:

$$\begin{aligned}
A_2 &= h \int_0^1 d\eta W_{11} F_{11} \\
A_1 &= \int_0^1 d\eta [W_{12} F_{12} + h W_{11} \partial_\xi F_{11} - W_{21} F_{21}] \\
A_0 &= \int_0^1 d\eta \left[ W_{12} \partial_\xi F_{12} - \frac{1}{h} W_{22} F_{22} + h W_{11} |J| k^2 \right]
\end{aligned}$$

To compute these matrices, we must compute the transformation matrix  $F$  and the Jacobian  $|J|$ . Since neither  $F$  nor  $|J|$  depend on the choice of the origins of the  $\xi, \eta$ - or the  $x, y$ -coordinate system, we set without loss of generality  $x_1 = x_2 = 0$ ,  $y_1 = 0$ ,  $y_2 = h$ ,  $\eta_1 = 0$ ,  $\eta_2 = h$ ,  $\beta_1 = -\alpha_1$  and  $\beta_2 = \alpha_2$ . Using the abbreviations  $a_1 = \tan \alpha_1$ ,  $a_2 = \tan \alpha_2$ , and  $a = \tan \alpha_1 + \tan \alpha_2$ , we obtain

$$J = \begin{pmatrix} \frac{1}{\zeta} & 0 \\ -\frac{1-\eta}{\zeta} a_1 + \frac{\eta}{\zeta} a_2 & 1 + \frac{\xi}{h\zeta} a \end{pmatrix} \quad \text{and} \quad J^{-1} = \begin{pmatrix} \zeta & -\frac{\zeta h(-a_1 + \eta a)}{h\zeta + a\xi} \\ 0 & \frac{\zeta h}{h\zeta + a\xi} \end{pmatrix}.$$

It follows

$$|J| = \frac{h\zeta + \xi a}{h\zeta^2} \quad \text{and} \quad F = \begin{pmatrix} \zeta + \frac{\xi}{h} a & a_1 - \eta a \\ a_1 - \eta a & h \frac{(a_1 - \eta a)^2 + 1}{h\zeta + \xi a} \end{pmatrix}.$$

Based on these results we can compute the elementary FEM-matrices. Now we label the matrices and the quantities  $h$ ,  $a_1$ ,  $a_2$ ,  $a$ , and  $\zeta$  with a subscript  $j$  and the matrices with a superscript  $(j)$  to indicate that these quantities belong to  $Q_j^{(\xi, \eta)}$  and may differ between different local domains.

$$\begin{aligned}
A_2^{(j)} &= (h_j \zeta_j + \xi a_j) \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix} \\
A_1^{(j)} &= \frac{1}{3} \begin{pmatrix} a & -a_2 + 2a_1 \\ 2a_2 - a_1 & a \end{pmatrix}^{(j)}
\end{aligned}$$

$$\begin{aligned} A_0^{(j)} &= \frac{1}{3} \frac{1}{h_j \zeta_j + \xi a_j} \begin{pmatrix} -a_1^2 + a_1 a_2 - a_2^2 - 3 & a_1^2 - a_1 a_2 + a_2^2 + 3 \\ a_1^2 - a_1 a_2 + a_2^2 + 3 & -a_1^2 + a_1 a_2 - a_2^2 - 3 \end{pmatrix}^{(j)} \\ &\quad + (h_j \zeta_j + \xi a_j) \frac{k_j^2}{\zeta_j^2} \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix}. \end{aligned}$$

The matrices depend on the quantity  $(h_j \zeta_j + \xi a_j)$ , which is the distance dependent width  $h(\xi)$  of each semi-infinite element. In view of the Laplace transform, it is convenient to regroup the matrices  $A_{0,1,2}^{(j)}$  according to this factor. The local matrices  $A_2^{(j)}$  are associated to  $a_2(v, \partial_\xi^2 u)$ , the matrices  $A_1^{(j)}$  to  $a_1(v, \partial_\xi u)$  and  $A_0^{(j)}$  to  $a(v, u)$ , according to the exterior variational problem (5.3.2). Thus we have, in a formal notation, the local approximation

$$(5.3.9) \quad a_0(v, u) + a_1(v, \partial_\xi u) + a_2(v, \partial_\xi^2 u) \xrightarrow{\text{local approx}} A_0^{(j)} \mathbf{u} + A_1^{(j)} \partial_\xi \mathbf{u} + A_2^{(j)} \partial_\xi^2 \mathbf{u}.$$

We aim to replace the expression on the right-hand side by sum over matrices which indicate besides the order of the associated derivative the power of  $(h_j \zeta_j + \xi a_j)$ . To this end we introduce the matrices  $M_{i,j}^{(k)} \in \mathbb{R}^{2 \times 2}$

$$(5.3.10) \quad \sum_{j=0}^2 \left( \sum_{i=-1}^1 (h_k \zeta_k + \xi a_k)^i M_{i,j}^{(k)} \right) \partial_\xi^j \mathbf{u} := A_0^{(k)} \mathbf{u} + A_1^{(k)} \partial_\xi \mathbf{u} + A_2^{(k)} \partial_\xi^2 \mathbf{u}$$

with

$$\begin{aligned} M_{-1,0}^{(k)} &= \frac{1}{3} \begin{pmatrix} -a_1^2 + a_1 a_2 - a_2^2 - 3 & a_1^2 - a_1 a_2 + a_2^2 + 3 \\ a_1^2 - a_1 a_2 + a_2^2 + 3 & -a_1^2 + a_1 a_2 - a_2^2 - 3 \end{pmatrix}^{(k)} \\ M_{1,0}^{(k)} &= \frac{k_k^2}{\zeta_k^2} \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix} \\ M_{0,1}^{(k)} &= \frac{1}{3} \begin{pmatrix} a & -a_2 + 2a_1 \\ 2a_2 - a_1 & a \end{pmatrix}^{(k)} \\ M_{1,2}^{(k)} &= \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix}, \end{aligned}$$

and all the other matrices  $M_{i,j}^{(k)} = 0$ . In this form it becomes transparent that the chosen method-of-lines type discretization yields a simple algebraic structure.

### Coupling between interior and exterior domain.

Based on these matrices we are able to assemble the whole discrete system in  $\xi, \eta$ -coordinates. Exterior and interior domain interact via the boundary integral, which appears on the right-hand side of the variational formulation of the interior problem (in  $x, y$ -coordinates) (5.3.1)

$$\int_{\partial\Omega_{\text{int}}} ds \bar{\mathbf{v}}(s) \mathbf{n}(s) \nabla_{xy} u(s).$$

We express this in  $\xi, \eta$ -coordinates on the interval  $\overline{p_i p_j}$ . As before, we set without loss of generality  $x_1 = x_2 = 0$ ,  $y_1 = 0$ ,  $y_2 = h$ ,  $\eta_1 = 0$ ,  $\eta_2 = h$  consequently  $\mathbf{n} = (1, 0)^T$ . We find with  $i, j \in \{1, 2\}$



$$\begin{aligned}
(5.3.11) \quad & \int_{p_1}^{p_2} ds \bar{\nabla}_i(s) \mathbf{n}(s) \nabla_{xy} u(s) \\
&= h \int_0^1 d\eta v_i(\eta) (1, 0) J^{-T} \nabla_{\xi\eta} (v_j(\eta) u_j(\xi)) \\
&= h \int_0^1 d\eta v_i(\eta) \begin{pmatrix} (J^{-1})_{11} & (J^{-1})_{12} \end{pmatrix} \begin{pmatrix} v_j \partial_\xi u_j \\ \partial_\eta v_j u_j \end{pmatrix} \\
&= \left( h \int_0^1 d\eta v_i(\eta) (J^{-1})_{11} v_j \right) \partial_\xi u_j + \left( h \int_0^1 d\eta v_i(\eta) (J^{-1})_{12} \partial_\eta v_j \right) u_j.
\end{aligned}$$

The last equation is of central importance for the whole concept. It supplies the link between the discretized interior and exterior problems. It becomes apparent that we do not need to introduce the normal derivative. Instead, we trace back the problem in a formulation in  $\xi, \eta$ -coordinates.

We want to express the terms in parenthesis by local matrices. For continuous, piece-wise linear elements we obtain the elementary matrices, using the notation of the previous section,

$$B_1 := \left( h \int_0^1 d\eta v_i(\eta) (J^{-1})_{11} v_j \right)$$

which results in the local matrix for the  $k$ th segment associated to the (first) derivative in  $\xi$ -direction

$$B_1^{(k)} = \frac{1}{3} h_k \zeta_k \begin{pmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix}.$$

In the same way we obtain

$$B_0 := \left( h \int_0^1 d\eta v_i(\eta) (J^{-1})_{12} \partial_\eta v_j \right)$$

with the result

$$B_0^{(k)} = \frac{1}{6} \begin{pmatrix} a - 3a_1 & -a + 3a_1 \\ 2a - 3a_1 & -2a + 3a_1 \end{pmatrix}^{(k)}.$$

Thus, we can extend the discretization scheme to the boundary terms, using the formal notation of (5.3.9),

$$(5.3.12) \quad \int_{\partial\Omega_{\text{int}}} ds \bar{\nabla}(s) \mathbf{n}(s) \nabla_{xy} u(s) \xrightarrow{\text{local approx}} \left( B_0^{(k)} \mathbf{u} + B_1^{(k)} \partial_\xi \mathbf{u} \right)_{\xi=0}.$$

### Algorithm.

Now we are able to assemble the whole discrete problem. Since most of all steps are completely parallel to the discussed 1D examples, we give only a brief overview and leave out the many obvious details.

**Step 1:** Discretization of the interior problem in standard manner, according to the left-hand side of (5.3.1).

**Step 2:** Discretization of the boundary terms of the interior problem, according to the right-hand side of (5.3.1), using the scheme (5.3.12). Dirichlet and Neumann data (with respect to  $\xi$ ) appear as unknowns, cf. Section 5.2 and Section 5.2.2.

- Step 3:** Assembling of the Laplace transformed exterior problem. To this end, the discretized exterior problem has to be assembled according to (5.3.9), followed by a Laplace transformation of the system. In the spatial domain, exactly the terms  $(h_k \zeta_k + \xi a_k)^i \partial_\xi^j u$ ,  $i \in \{-1, 0, 1\}$ ,  $j \in \{0, 1, 2\}$  appear, which already have appeared in the 1D situation.
- Step 4:** Formulate the linking condition exactly as in Section 5.2 as axis-integral condition, or as in Section 5.2.2 as cut-function formulation. Each segment supplies one pair of linking conditions.
- Step 5:** Discretize in the Laplace domain, either along the real axis (compare Section 5.2) or along the cuts (compare Section 5.2.2).
- Step 6:** Solve simultaneously the interior (spatial) and the exterior (spectral) problem.

### Asymptotic Approximation.

Following the example of the asymptotic (high frequency approximation) of the DtN-number of Bessel's equation Section 2.4, pp. 28, we can apply the same factorization technique to the system of ODE's

$$Lu(\xi) := \left[ \left( \sum_{j=0}^2 \xi^j M_{j2} \right) \partial_\xi^2 + \left( \sum_{j=0}^2 \xi^j M_{j1} \right) \partial_\xi + \left( \sum_{j=0}^2 \xi^j M_{j0} \right) \right] u(\xi) = 0,$$

which results from (5.3.10) by a multiplication with factors of the form  $h\zeta + \xi a$  and an assembling of all local matrices. Now we can repeat the procedure of Section 2.4 to derive a family of boundary conditions at the position  $\xi = 0$ . The symbol of  $L$  is given by

$$\sigma(L) = \left( \sum_{j=0}^2 \xi^j M_{j2} \right) (-\rho^2) + \left( \sum_{j=0}^2 \xi^j M_{j1} \right) (i\rho) + \left( \sum_{j=0}^2 \xi^j M_{j0} \right).$$

Here we used  $\xi$  as distance variable as before, and  $\rho$  as corresponding dual variable. Note that there is no dual pair  $\phi, \nu$  as in the continuous case. Introducing the matrices

$$K_2 = \left( \sum_{j=0}^2 \xi^j M_{j2} \right), \quad K_1 = K_2^{-1} \left( \sum_{j=0}^2 \xi^j M_{j1} \right), \quad K_0 = K_2^{-1} \left( \sum_{j=0}^2 \xi^j M_{j0} \right)$$

the symbol reads

$$\sigma(L) = K_2 (-I\rho^2 + K_1 i\rho + K_0).$$

Obviously, this is the discrete counterpart of the symbol of the continuous Helmholtz equation, cf. 2.4.13, p. 31. Accordingly, the factorization

$$(-I\rho^2 + K_0) = \left( Ii\rho + iK_0^{1/2} \right) \left( Ii\rho - iK_0^{1/2} \right)$$

defines the highest order pseudodifferential operator  $\Lambda_\pm^1 = \mp iK_0^{1/2}$ , and the same consideration as in Section 2.4 yields the two equations for the symbols  $\Lambda_-^0(x, \rho)$  and  $\Lambda_-^{-1}(x, \rho)$

$$\begin{aligned} 0 &= K_1 i\rho + \partial_x \Lambda_\pm^1 - \Lambda_\pm^0 (\Lambda_\pm^1 - \Lambda_\pm^1) \\ 0 &= \quad \quad + \partial_\xi \Lambda_\pm^0 - \Lambda_\pm^{-1} (\Lambda_\pm^1 - \Lambda_\pm^1) + \Lambda_\pm^0 \Lambda_\pm^0. \end{aligned}$$

We wrote  $\Lambda_\pm^j$  for the symbols of the pseudodifferential operators to indicate that these are parameter dependent matrices in our discrete case, where the parameters

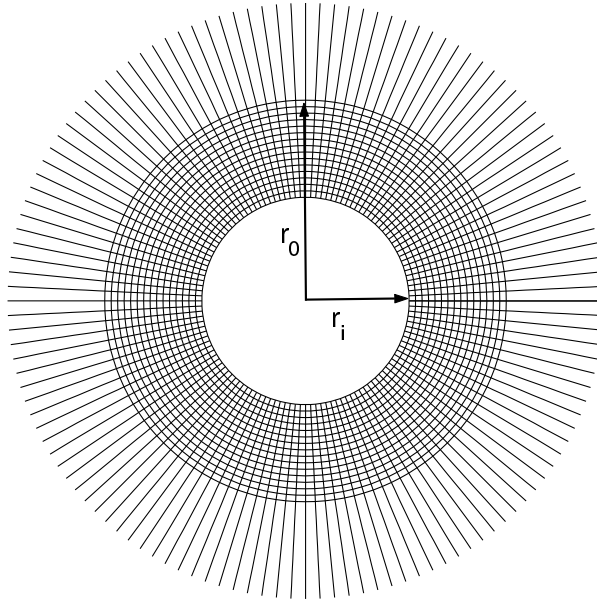


FIGURE 5.4.1. Regular mesh used for the radially symmetric scattering problem.

are  $\xi$  and  $\rho$ . If  $K_0^{1/2}$  can be computed in a numerically cheap way, this gives a practicable way to compute asymptotic boundary conditions directly from the given semi-discrete problem. Without following this way further here we remark that a cheap computation of a DtN-operator could supply an useful preconditioner for the algebraic system.

#### 5.4. Numerical Examples

The following 2D examples are the very first numerical results and are presented here to give an impression of the technique. Future work is needed to both analyze and improve the numerical sub-problems.

The first numerical experiments are meant to repeat the numerical studies of Section 5.2.2, pp. 115, in a 2D setting. We compute the solution of the Helmholtz equation with  $k = 1$  exterior to a disk with radius  $r_1 = 2\pi$ . On the boundary of the disk Dirichlet data  $u|_{\Gamma}(\phi) = \cos(\nu\phi)$ ,  $0 \leq \phi \leq 2\pi$  are prescribed. We carried out the experiments with  $\nu = 0$  and  $\nu = 10$  as characteristic angular frequencies. The interior computational domain is a circular domain between the radii  $r_1$  and  $r_0$ , where  $r_0 > r_1$  is the radius of the artificial circular boundary. In the following experiments we set  $r_0 = 4\pi$ , so that the radius of the artificial domain covers two wavelengths. The corresponding regular mesh is shown in Fig. 5.4.1.

Fig. 5.4.2 displays the convergence of the interior solution, if the initial mesh is subsequently uniformly refined, starting with a number of unknowns on the artificial boundary  $N_\phi = 32$  and ending up with  $N_\phi = 256$ . During the refinement of the interior domain we kept fix the discretization of the spectral domain. Here we used 14 intervals and the implicit Runge-Kutta collocation scheme of consistence order 4, p. 118, which results in 28 additional unknowns per node on the artificial boundary.

A natural measure for the quality of the discretization in the spectral domain is the far-field error shown in Fig. 5.4.3. Roughly, the far-field errors behave similarly like the interior errors. The errors corresponding to the higher angular frequency

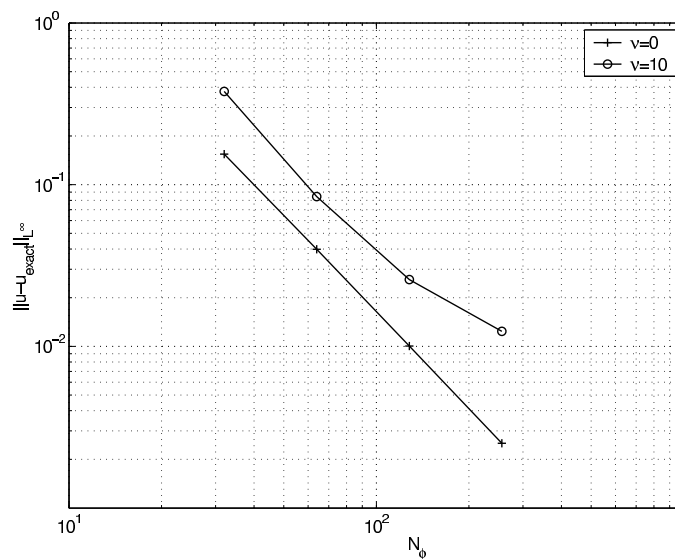


FIGURE 5.4.2. Convergence of the interior solution computed for a sequence of uniformly refined interior meshes. The discretization in the spectral domain is kept fix.  $N_\phi$  denotes the number of nodes on the artificial boundary.

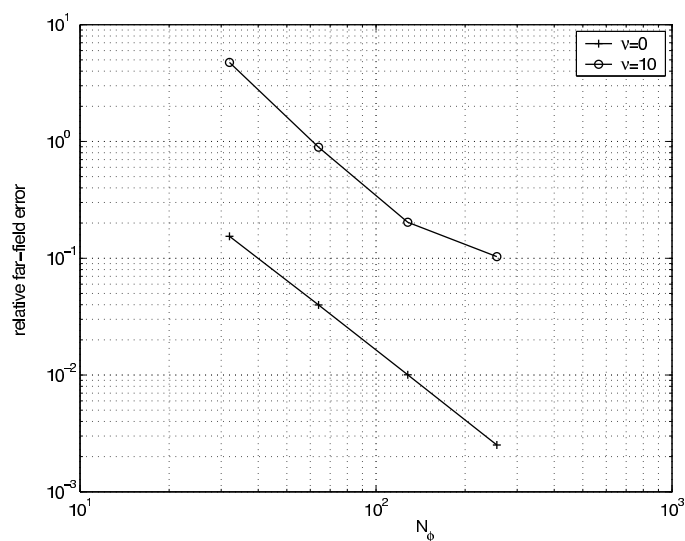


FIGURE 5.4.3. Far-field error corresponding to the radially symmetric scattering problem

$\nu = 10$  are larger than the errors corresponding to  $\nu = 0$ . This similar behavior shows that in the given situation the far-field error is dominated by the interior discretization error. The cut functions  $\psi$  corresponding to the frequencies  $\nu = 0$  and  $\nu = 10$  for the analyzed radially symmetric experiment are displayed in Fig. 5.4.4. The smooth behavior of these curves indicate that they can be efficiently treated by a number of numerical approximation methods so that future implementations of the cut function approach will result in fast methods.

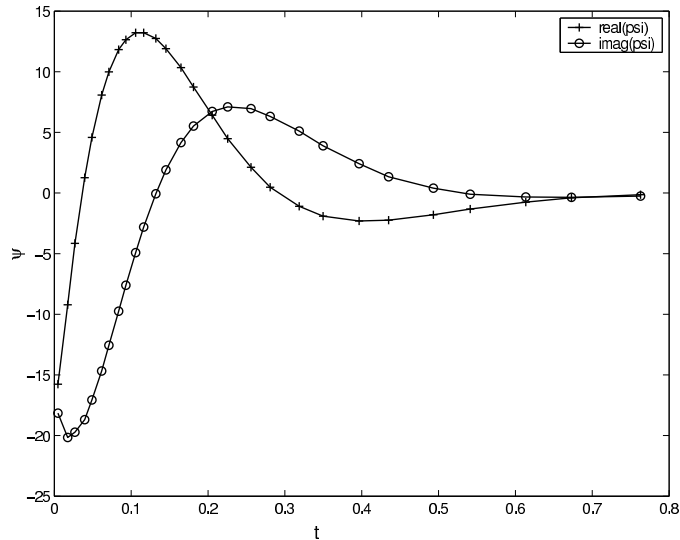


FIGURE 5.4.4. Cut function corresponding to  $\nu = 10$ . The + -signs indicate the collocation points.

Next we want to give briefly some other numerical experiments. Fig. 5.4.5 shows the triangulation of an elliptic computational domain and the matching ray-like discretization of the exterior domain. A plane wave propagates from left to right and causes the computed reflection from an interior Dirichlet boundary (circle).

Fig. 5.4.6 displays the the first convergence result obtained by our method. This result has a very preliminary character, since we have not implemented the very successful spectral approaches from Sections 5.2 and 5.2.2. Instead, we discretized the whole system as described above, but then we approximated the space of the Dirichlet functions on the boundary by a small number of Fourier modes. This approximation is the same as performed in the classical derivation of DtN-operators based on factorizations, see Section 2.1. In future implementations, this approximation will be replaced by the Laplace domain methods. Fig. 5.4.6 shows the convergence of the interior problem with respect to the dimension of the subspace, that is the number of Fourier modes used for the approximation of  $u$  on  $\partial\Omega$ , for three subsequently uniform refined meshes. The interesting result is that if the approximation quality of the subspace is sufficient, then the remaining error is the discretization error of the interior problem.

**Application of radial rays.** Fig. 5.4.7 shows an application to the same problem with an quadratic computational domain. The discretization of the exterior domain has been performed with *radial* rays. This allows in a natural way to consider corners, without to introduce any special technique.

**Application of normal rays.**

The example of Fig. 5.4.8 shows the reflection from a semi-transparent mirror. Here the exterior domain contains a waveguide, which extends from the left boundary to infinity. This waveguide inhomogeneity is easily taken into account by a *normal* ray discretization of the exterior domain.

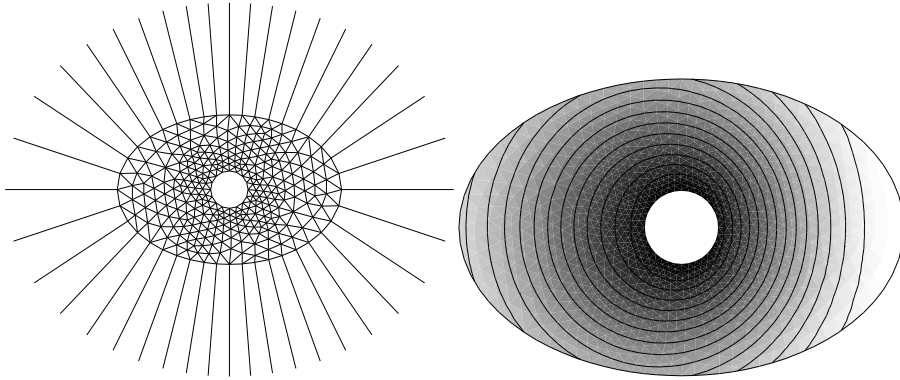


FIGURE 5.4.5. Inner and outer discretization of an elliptic computational domain and the corresponding numerical solution of the scattering problem.

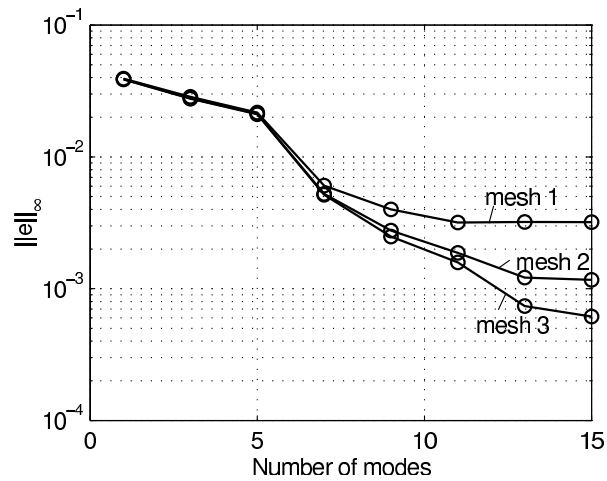


FIGURE 5.4.6. Convergence of the error  $e := u_{\text{exact}} - u_{\text{discrete}}$  in maximum norm vs. the number of modes for an initial (mesh 1) and two subsequently refined triangulations (meshes 2 and 3). The error decreases until the discretization error is reached.

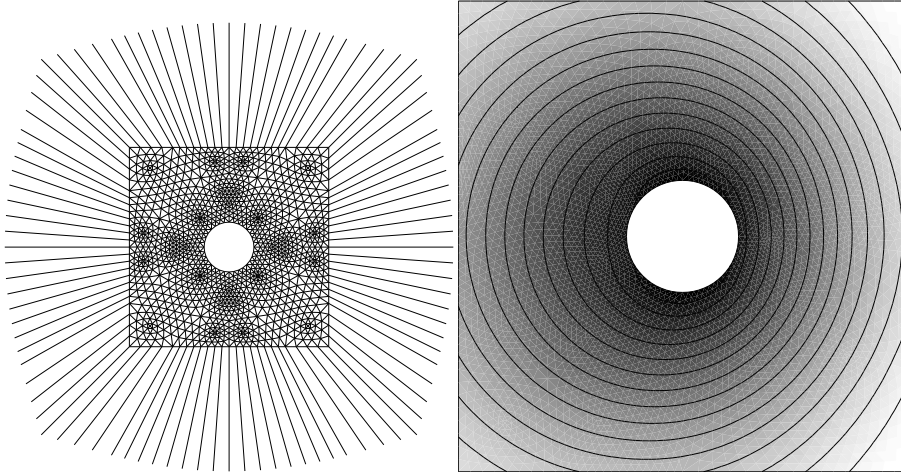


FIGURE 5.4.7. Inner and outer discretization of a quadratic computational domain and the corresponding numerical solution of the scattering problem.

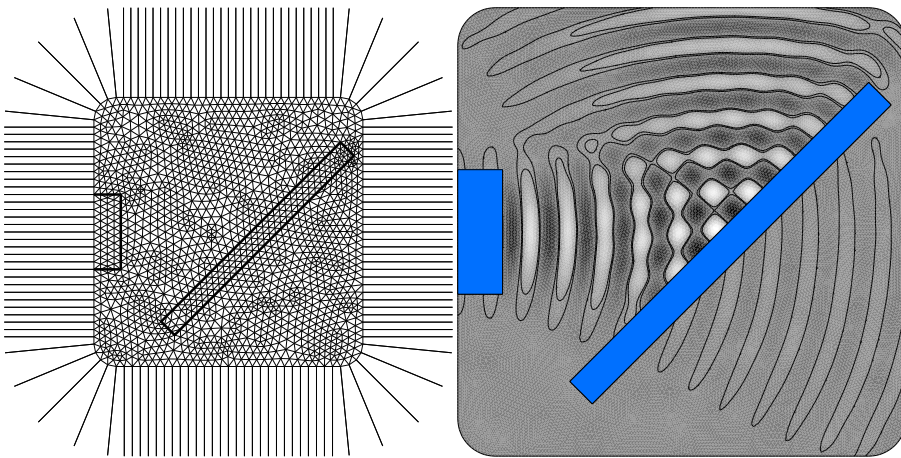


FIGURE 5.4.8. Inner and outer discretization of the reflection problem and the corresponding numerical solution. This is an example of an inhomogeneous scattering problem.

## Numerical Treatment of Schrödinger-Type Scattering Problems

This chapter applies the pole condition approach to time-dependent Schrödinger and to wide-angle one way equations. We want to discuss the problem: Given the (possibly variable coefficient) evolution problem and a (possibly non-uniform) time-discretization, how to derive corresponding transparent boundary conditions? We do not aim, in first place, to derive fast algorithms in the sense of e.g. fast convolution algorithms, but on the relation between the time-evolution and the pole condition.

Therefore we study evolution problems with only one spatial dimension. This, however, supplies the framework for the application of the methods developed in the previous sections, which, in turn, will lead to solvers for variable coefficient problems in higher dimensions. We give the principal structure of the corresponding algorithms; the analysis and realization of corresponding algorithms, however, is part of future research.

Our approach here is the same as in Section 5.1, namely a direct computation of the poles by a factorization technique and their elimination by a proper choice of the boundary conditions. In one space dimension this concept leads directly to the desired boundary conditions even in non-trivial cases like the one-way wide-angle equations. Based on these boundary conditions, we can prove the unconditional stability of the numerical propagation algorithms and of some discrete conservation properties.

### 6.1. Time-Dependent Schrödinger Equation

First we focus on the construction of transparent boundary conditions for discretized Schrödinger equations. The presentation follows the one given by the author in [80], and the joint work with Friese and Yevick published in [83].

Let the continuous evolution equation be given:

$$(6.1.1) \quad \begin{aligned} \partial_t u &= -\frac{i}{c} (\partial_x^2 u + V(x, t)u), & x \in \mathbb{R}, & t > 0 \\ u(x, 0) &= u_0, & \text{and} & \lim_{x \rightarrow \pm\infty} u(x) = 0, \end{aligned}$$

cf. Section 1.1, p. 12. Here  $c$  is a negative real constant and  $V(x, t)$  denotes a real and bounded potential. The operator  $A = 1/c(\partial_x^2 + V(x, t))$  with  $A : D(A) \subset L^2(\mathbb{R}) \rightarrow L^2(\mathbb{R})$  is symmetric. Hence, we have for all  $u_0 \in D(A)$  the conservation property  $\partial_t(u, u)_{L^2} = 0$ . We consider the evolution equation (6.1.1) in the infinite domain  $\Omega = \{x, t \in \mathbb{R} : t > 0\}$ . Further, we assume that the initial function  $u_0(x)$  has support only in a finite interval. We decompose the infinite domain  $\Omega$  into an interior domain  $\Omega_{\text{int}}$  containing the object of interest and its complement with respect to the whole domain  $\Omega$ . In our 1D-case we define  $\Omega_{\text{int}} = \{x, t \in \mathbb{R} : x_- < x < x_+, t > 0\}$ , and the exterior domain consists of the two domains  $\Omega_- = \{x, t \in \mathbb{R} : x < x_-, t > 0\}$  and  $\Omega_+ = \{x, t \in \mathbb{R} : x > x_+, t > 0\}$ , for the left and right semi-infinite sub-domains, respectively. Further, we assume for simplicity



that the potential function  $V$  depends in the exterior domains  $\Omega_{\pm}$  only on  $t$ , i. e. ,  $V = V(t)$ . A dependence of the potential on the transversal direction  $x$  can be taken into account, as well as higher space dimensions, by the same techniques as described in Sections 4 and 5. However, we want to simplify the problem to the extend possible such that we can concentrate on the treatment of the time evolution.

Besides the discrete pole condition based approach, which we present here, there is a large number of other methods tailored to evolution problems. In the following, we extend the discussion of some of the most important methods given in Section 1 to time-dependent methods.

**Artificially absorbing layers and PML-method.** From the physical point of view, it is a natural to modify the potential function in the exterior domain such that an artificial absorption is generated. The parameters describing the absorber-function have to be adjusted such that backward diffraction from the absorber is small over a prescribed spectral range (e. g. Kosloff and Kosloff [65] or Yevick [92]). The main advantage of such an approach is its simplicity for two and three-dimensional problems.

The PML-technique of Bérenger [11] has been constructed originally to solve time-dependent Maxwell's equations, hence this approach is a natural candidate for Schrödinger-type problems. The method turned out to be very efficient for many problems. Here we have the same simple formulation for 2D and 3D problems but also the same problem in finding the best choice of parameters for given task. A numerical study of the PML-technique with respect to Schrödinger-type equations in the context of optical beam propagation may be found in [94].

**Nonlocal boundary conditions.** As for Helmholtz-type scattering problems, a solution representation based on Green' functions provides the basis for the construction of transparent boundary conditions. The main steps to derive such boundary conditions are the following. At first, the partial differential equation is converted to a second-order differential equation in space by Laplace transforming in *time*. The transformed differential equation is then solved allowing only for decaying modes in the exterior domains. The back transform to the time-domain yields convolution formulas which are used to construct transparent boundary conditions. The derivation of such formulas for a variety of constant coefficient problems is well understood, see the review paper of Hagstrom [50]. The back transform contains convolution integrals of the form

$$(6.1.2) \quad \int_0^T d\tau f(t-\tau)g(\tau), \quad 0 \leq t \leq T,$$

which can be discretized directly, as suggested by Baskakov and Popov [8]. However, such approaches may lead to numerical instabilities. An analysis of a discretization scheme based on the implicit mid-point rule for the time-evolution may be found in the thesis of Mayfield [72]. She showed that numerical stability, even in the case of uniform discretization both in time (with time-step  $\Delta t$ ) and space (with space-step  $\Delta x$ ), is given only in disjointed intervals for  $\Delta t/\Delta x^2$ . Alternatively, the problem of finding boundary conditions for the discretized equation, can be consistently formulated for discrete time. This repeats the continuous procedure in the time-discrete situation, cf. e.g. [69]. In this manner, Arnold [4] composes a boundary condition which incorporates both a uniform space and a uniform time discretization. Besides the uniformity of discretization, the time-independency of the potential in the exterior domain is an essential prerequisite. The key point in his approach is that the uniform space discretization of the interior domain is continued to the exterior domain as well. In contrast, the approach [82] supposes a given, possible

nonuniform, time-discretization and an arbitrary inner space-discretization. The space-variable in the exterior domain is not discretized. Additionally, the outer potential may be a function of time. We will label this approach the semi-discrete method. In [85] we have shown that our semi-discrete method covers Arnold's fully discrete method.

The non-local (in time) boundary conditions are often considered to be inefficient, if a long evolution history has to be taken into account. Therefore a number of fast solution techniques has been proposed. One strategy is to approximate the convolution integrals by analytical means such that they can be evaluated fast. In its core, the function  $f$  contained in the convolution integral (6.1.2) is approximated by a finite sum of exponential functions, say

$$f(t) = \sum_{j=1}^m c_j e^{\lambda_j t}$$

which is always possible if the Laplace transform of  $f(t)$  is a rational function. For each exponential function  $\exp(\lambda_j t)$  the convolution integral (6.1.2) may be expressed as a function  $t \mapsto y$ :

$$y(t, \lambda_j) := \int_0^t d\tau e^{\lambda_j(t-\tau)} g(\tau).$$

Differentiation with respect to  $t$  shows that  $y$  satisfies the initial-value problem

$$y' = \lambda_j y + g, \quad y(0) = 0.$$

The remarkable fact is, that the computation of the convolution integral, which is a nonlocal operation, can be replaced by a solution of an initial value problem, where the computation of a new time-step requires only local operations. The described situation occurs exactly for the wave equation in 3D, when the solution of the wave equation is represented by separated spherical harmonics. The idea and algorithms go back to Grote and Keller [43, 42]. Unfortunately, most of all practical problems do not have this simple structure, as became evident even from our analysis of the Helmholtz equation in Sections 1 and 4, where the characteristic singularities were cut singularities. However, it is a convincing idea, to approximate a given function by a finite number of exponentials such that the above strategy may be applied to more general problems. This concept is discussed in detail in the review of Hagstrom [50]. Moreover, in [1] Alpert, Greengard and Hagstrom give explicit rigorous bounds on the degree  $m$  required to obtain a uniform approximation of  $f$  for the whole time-interval  $T$ .

A second strategy is to employ the fast convolution technique due to Hairer, Lubich and Schlichte [51] to obtain fast methods. Recently, such an approach has been combined with the initial-value formulation of the convolution problem. In [70] Lubich and Schädle have shown that nearly optimal algorithms can be constructed this way. As mentioned above, we will not consider the efficiency aspect of our algorithms here, since there is some reasonable hope that the pole condition approach together with the Laplace domain method will yield efficient algorithms, too.

**Local approximations of non-local boundary conditions.** In case of higher space-dimensions ( $\geq 2$ ) one obtains local approximations in space using the theory of pseudo-differential operators, exactly as described in Section 2.4. After transforming back into the time-domain the resulting transparent boundary conditions in general become local in space but nonlocal in time, as discussed above. To avoid the additional numerical costs due to the non-locality in time, the dispersion relation between time- and space-variables may be rational approximated in the

dual domain. This construction scheme has been successfully applied by a number of authors. Following the pioneering work of Engquist and Majda [28] on hyperbolic equations, advanced approximation techniques have been proposed for mixed parabolic-hyperbolic systems (Halpern, [52]) or parabolic equations (Hagstrom, [49]). Again, the review of Hagstrom [50] gives an excellent overview.

### Longitudinal Discretization.

Let us rewrite (6.1.1) as

$$\begin{aligned}\partial_t u &= f(u, t), & (x, t) \in \Omega \\ \lim_{x \rightarrow \pm\infty} u(x, t) &= 0.\end{aligned}$$

To solve this equation numerically, we apply the implicit one-step discretization

$$\begin{aligned}u_{i+1} - u_i &= \tau f(\theta u_{i+1} + (1 - \theta)u_i, t_i + \theta\tau) \\ \tau &= t_{i+1} - t_i, & i = 0, 1, \dots \\ 0 &< \theta \leq 1.\end{aligned}$$

Using the definition of  $f(u, t)$  from (6.1.1), we find

$$(6.1.3) \quad u_{i+1} - u_i = -i \frac{\tau}{c} ((\partial_x^2 + V)(\theta u_{i+1} + (1 - \theta)u_i)).$$

By rearranging the terms we obtain a sequence of inhomogeneous Helmholtz-type equations

$$(6.1.4) \quad \begin{aligned}\partial_x^2 u_{i+1} - \lambda^2 u_{i+1} &= -\Theta \partial_x^2 u_i + \kappa^2 u_i \\ \Theta &= \frac{1 - \theta}{\theta} \\ \lambda^2(x, t_i + \theta\tau) &= \frac{ic}{\tau\theta} - V(x, t_i + \theta\tau) \\ \kappa^2(x, t_i + \theta\tau) &= -\frac{ic}{\tau\theta} - \Theta V(x, t_i + \theta\tau).\end{aligned}$$

We now seek solutions  $u_i$ ,  $i \geq 1$  of (6.1.4) that vanish at infinity. We focus here on obtaining an exterior solution which enables the boundary conditions to be constructed, independently of the numerical method employed to solve the interior problem. For this purpose, we fix the right boundary at  $x_+ = 0$ ,  $t > 0$  and search for solutions  $u_i(x)$ ,  $i \geq 1$ ,  $x \geq 0$  in the right exterior domain. These exterior solutions have to obey the boundary condition at infinity

$$(6.1.5) \quad \lim_{x \rightarrow \infty} u_i(x) = 0, \quad i \geq 1.$$

**Continuous treatment of the space coordinate.** In the semi-discretized equation (6.1.3) the space-variable  $x$  appears as continuous variable. Therefore (6.1.3) forms a sequence of ordinary differential equations. In order to apply the pole condition, we perform a Laplace transformation with respect to the spatial coordinate  $x$ . In contrast to our previous notation, we denote the dual variable to  $x$  by  $p$ . Later we will introduce the quantity  $s$  as shift operator. With

$$\hat{u}_i(p) = Lu_i(x) = \int_0^\infty e^{-px} u_i(x) dx, \quad p \in \mathbb{C}, \quad \text{Re}(p) > 0,$$

we find the equivalent transformed system to

$$(6.1.6) \quad \begin{aligned}(p^2 - \lambda_{i+1}^2) \hat{u}_{i+1}(p) + (\Theta p^2 - \kappa_{i+1}^2) \hat{u}_i(p) \\ = pu_{i+1}(0) + \partial_x u_{i+1}|_{x=0} + \Theta(pu_i(0) + \partial_x u_i|_{x=0}).\end{aligned}$$

**Discrete treatment of the space coordinate.** The alternative idea is to consider a uniformly discretized space. That is, we associate the solution points at the  $i$ th propagation step with physical locations according to the formula

$$u_i^{(j)} = u_i(j \cdot h), \quad j \geq -1, \quad i \geq 0.$$

Here  $u_i^{(-1)}$  is the rightmost inner value in  $\Omega_i$  while  $u_i^0$  is located on the boundary between the internal and the right external region. The equation corresponding to 6.1.6 is obtained by introducing the sequences

$$(6.1.7) \quad \begin{aligned} U_i &= \{u_i^{(0)}, u_i^{(1)}, u_i^{(2)}, \dots\} \\ U_i^+ &= \{u_i^{(1)}, u_i^{(2)}, u_i^{(3)}, \dots\} \\ U_i^- &= \{u_i^{(-1)}, u_i^{(0)}, u_i^{(1)}, \dots\} \end{aligned}$$

with  $Z$ -transforms

$$\begin{aligned} ZU_i &= u_i^{(0)} + \frac{1}{z}u_i^{(1)} + \frac{1}{z^2}u_i^{(2)} + \dots \\ ZU_i^+ &= u_i^{(1)} + \frac{1}{z}u_i^{(2)} + \frac{1}{z^2}u_i^{(3)} + \dots \\ ZU_i^- &= u_i^{(-1)} + \frac{1}{z}u_i^{(0)} + \frac{1}{z^2}u_i^{(1)} + \dots \end{aligned}$$

By  $Z$ -transforming the finite-difference form of (6.1.4) and taking into account the relations between  $ZU_i$ ,  $ZU_i^+$ , and  $ZU_i^-$ , we obtain the discrete counterpart to (6.1.6)

$$(6.1.8) \quad \begin{aligned} &\left(z - (2 + h^2\lambda^2) + \frac{1}{z}\right) ZU_{i+1}(z) + \Theta \left(z - (2 + h^2\kappa^2/\Theta) + \frac{1}{z}\right) ZU_i(z) \\ &= u_{i+1}^{(-1)} - zu_{i+1}^{(0)} + \Theta(u_i^{(-1)} - zu_i^{(0)}) \end{aligned}$$

Comparing (6.1.6) and (6.1.8) we do not find any qualitative difference in the structure of these equations with respect to the time-evolution. The only difference is that the factors before  $ZU_{i+1}$  and  $ZU_i$  are slightly more involved. Based on this similarity, we restrict all of our following considerations to (6.1.6). The corresponding results for (6.1.8) are obtained in complete analogy.

### Construction of Discrete Transparent Boundary Conditions.

The main tool in deriving the transparent boundary conditions is the discrete analog of Mikusiński's operational calculus [74]. The difference is that Mikusiński considered an algebra based on continuous, complex valued functions over a semi-infinite time-interval, whereas our basic elements are semi-infinite sequences of complex numbers.

As we have seen, sequences like (6.1.7) play a central role. Let  $U$  and  $V$  such sequences. We need the operations addition (denoted by  $U + V$ ), multiplication (denoted by  $UV$ , which is in fact a convolution) and the map  $U \mapsto V = AU$ , where  $A$  is an infinite dimensional matrix. The sequence  $s = \{0, 1, 0, \dots\}$  is called the shift operator. In the Appendix A.2 the properties needed for our purposes are collected. In order to reformulate our problem based on the algebraic calculus, we rewrite the sequence (6.1.6) into a matrix equation with infinite dimensional matrices:

$$\begin{aligned}
(6.1.9) \quad & \left[ p^2 \begin{pmatrix} 1 & & & & & \\ \Theta & \ddots & & & & \\ & \ddots & 1 & & & \\ & & \Theta & 1 & & \\ & & & \ddots & \ddots & \\ & & & & \ddots & \ddots \end{pmatrix} - \begin{pmatrix} \lambda_1^2 & & & & & \\ \kappa_1^2 & \ddots & & & & \\ & \ddots & \lambda_i^2 & & & \\ & & \kappa_{i+1}^2 & \lambda_{i+1}^2 & & \\ & & & \ddots & \ddots & \\ & & & & \ddots & \ddots \end{pmatrix} \right] \begin{pmatrix} \hat{u}_1 \\ \vdots \\ \hat{u}_i \\ \hat{u}_{i+1} \\ \vdots \end{pmatrix} \\
& = \begin{pmatrix} 1 & & & & & \\ \Theta & \ddots & & & & \\ & \ddots & 1 & & & \\ & & \Theta & 1 & & \\ & & & \ddots & \ddots & \\ & & & & \ddots & \ddots \end{pmatrix} \left[ p \begin{pmatrix} u_1(0) \\ \vdots \\ u_i(0) \\ u_{i+1}(0) \\ \vdots \end{pmatrix} + \begin{pmatrix} u'_1(0) \\ \vdots \\ u'_i(0) \\ u'_{i+1}(0) \\ \vdots \end{pmatrix} \right].
\end{aligned}$$

Denoting the matrix containing the  $\Theta$ 's by  $Q$  and the matrix containing the  $\lambda_j^2$ 's and  $\kappa_j^2$ 's by  $T$  we write, instead of (6.1.9), equivalently

$$(6.1.10) \quad (p^2Q - T)U = Q(p\{u(0)\} + \{u'(0)\}).$$

**Uniform discretization.** Starting with uniform discretization and time-independent exterior potentials, i. e. ,  $\lambda_1^2 = \lambda_2^2 = \dots = \lambda^2$  and  $\kappa_1^2 = \kappa_2^2 = \dots = \kappa^2$ , we obtain from (6.1.9), using the shift-operator definition,

$$(6.1.11) \quad \left( p^2 - \frac{\lambda^2 + s\kappa^2}{1 + s\Theta} \right) U = p\{u(0)\} + \{u'(0)\}.$$

This way, we have expressed the operators  $T$  and  $Q$  completely in terms of the shift operator  $s$ . Both  $T$  and  $Q$  are in the space of sequences  $\mathcal{C}$ , therefore  $Q^{-1}T$  is in the ring of fractions  $\mathcal{C}_q$  (see the Appendix A.2), hence all rational operations with respect to  $Q^{-1}T$  are well defined. The term  $Q^{-1}T$  has the algebraic structure of a lower triangular matrix, since  $T$  and  $Q$  are lower triangular matrices. Because  $\text{Im}(\lambda^2) \neq 0$ , one proves easily that there is always a uniquely defined lower triangular matrix  $L$  such that

$$L \cdot L = L^2 = Q^{-1}T \quad \text{with} \quad \text{Re}(\text{diag}(L)) > 0.$$

Formally, we write

$$(6.1.12) \quad L = \sqrt{\frac{\lambda^2 + s\kappa^2}{1 + s\Theta}}.$$

It is clear that for the expression (6.1.12) a Taylor series expansion in  $s$  at  $s = 0$  exists. To any given order, we can take a sufficient number of terms of the series expansion such that the square of the series expansion equals to  $Q^{-1}T$  up to the given order. Therefore we have equivalently to (6.1.12)

$$L = \lambda + a_1s + a_2s^2 + \dots, \quad \text{with} \quad \lambda = \sqrt{\lambda^2}, \text{Re}(\lambda) > 0,$$

see Example A.2.6 of the Appendix A.2 for a more detailed discussion. It holds  $\text{Im}\lambda^2 < 0$ , by its definition in (6.1.4). Hence  $\lambda$  as defined above has a negative imaginary part. We obtain the exterior solutions  $\{U\}$ , as function of the boundary data,

$$U = \frac{p\{u(0)\} + \{u'(0)\}}{(p - L)(p + L)}.$$

This expression has poles at  $\pm L$ . The operator  $+L$  contains  $+\lambda$ , the operator  $-L$  contains  $-\lambda$ . The special choice

Pole condition for the uniformly discretized Schrödinger equation:

$$(6.1.13) \quad \{u'(0)\} = -L \{u(0)\}$$

supplies the desired transparent boundary conditions. This becomes clear, if we consider the related exterior solutions  $U$

$$U = \frac{\{u(0)\}}{p + L}.$$

Since  $L$  is a lower triangular matrix with real parts of the diagonal entries all greater than zero, it follows that  $U$  contains only decaying modes.

**REMARK.** We have motivated the pole condition by the decay behavior of the exterior solution. Therefore we excluded poles with positive real part. Since  $\text{Im } \lambda < 0$  for  $\text{Re } \lambda > 0$ , we could require as well that all poles with negative imaginary part must be excluded. This, however, is exactly our pole condition defined in 3.4.8, p. 56.

**Non-uniform discretization (algebraic approach).** The difficulty, which appears in generalizing the uniform approach is that the operator  $T$  cannot be expressed in terms of the shift-operator  $s$ , because all  $\lambda_i^2$ 's and  $\kappa_i^2$ 's may be different from each other. Nevertheless, the whole approach remains the same, only the simplification of the operator representation is lost. Equation (6.1.10) now reads

$$(6.1.14) \quad \left( p^2 + \frac{1}{1 + s\Theta} T \right) \{U\} = p \{u(0)\} + \{u'(0)\}.$$

As the argumentation concerning the factorization of operators  $Q^{-1}T$  is equally valid here, we solve the problem of the transparent boundary conditions in the nonuniform case, if we find a factorization

$$L \cdot L = L^2 = \frac{1}{1 + s\Theta} T \quad \text{with} \quad \text{Re}(\text{diag}(L)) > 0.$$

This, however, is a standard task in numerical linear algebra: find the square root of a lower triangular matrix with non-vanishing diagonal elements such that the lower triangular matrix contains only diagonal elements with positive real part. Among these methods are the direct Cholesky-like factorization, which amounts to  $\mathcal{O}(n^2)$  operations in the  $n$ -th propagation step, Krylov-subspace methods, and suitable basis-transformations, which amount to  $\mathcal{O}(n)$  operations. Once the matrix  $L$  is obtained, we finally arrive again at the boundary condition (6.1.13).

### Laplace domain method.

We refer to the algebraic representation of the time-discrete evolution problem in its general form (6.1.10). Instead to compute the poles explicitly, as we have done so far, we may add to each equation of the system the pole condition as an additional constraint. Let us take, for illustrative purposes, the pole condition in its direct (axis integral) form, as given in Lemma 4.4.3.

Following the representation given in Section 5.2, the full algebraic system consists of three parts:

**Part I:** Interior problem in discrete form given by (6.1.4).

**Part II:** Pole condition in discrete form. It consists of a discretization of  $U$  along the real axis

$$0 = \int_{-\infty}^{\infty} d\tau \frac{U(\tau)}{\tau - p_+}$$

for some fixed  $p_+$  with  $\text{Im}(p_+) > 0$ , and  $\tau \in \mathbb{R}$ .

**Part III:** Laplace transformed exterior problem in discrete form. It consists of a discretization of (6.1.10) along the real axis

$$(\tau^2 Q - T)U(\tau) = Q(\tau \{u(0)\} + \{u'(0)\}), \quad \tau \in \mathbb{R}.$$

We require the simultaneous solution of all three parts. The discretization of  $U$  in the spectral domain, and the linking to the interior problem can be realized exactly as discussed in detail in Sections 5.2 and 5.2.2.

Let  $n_p$  the number of propagation steps and  $n_U$  the number of discretization points of  $U$  on the real axis. Then the proposed approach has the obvious properties:

- (1) A fixed discretization of  $U$  along the real axis with  $n_U$  points requires  $O(n_p n_U)$  additional operations to compute the boundary condition, which is optimal.
- (2) Since we can solve the system equation by equation, the additional memory requirement is  $O(n_U)$ , i. e., it is independent of the number of steps.

The analysis and implementation of this algorithm is part of the current research.

## 6.2. Wide Angle One-Way Equations

The following is based on the representation given by the author, Friese and Yevick in [83] and extends the foregoing to a more general case. We consider a scalar wave propagation in 2D modeled by the Helmholtz equation

$$(6.2.1) \quad \partial_z^2 u + \partial_x^2 u + k^2 u = 0$$

in the entire  $\mathbb{R}^2$ , where the wavenumber  $k$  is a position dependent, bounded function  $(x, z) \mapsto k$ . In many physical situations, we can further distinguish a principal propagation direction, here taken to be the  $z$ -direction and a transverse,  $x$  direction. For the particular case of a position-independent wavenumber, the operator  $\partial_z^2 + \partial_x^2 + k^2$  can be explicitly factorized, leading to the exact one-way Helmholtz equation

$$(6.2.2) \quad \partial_z u = ik \sqrt{1 + \frac{1}{k^2} \partial_x^2} u.$$

In the above expression, the formal square root operator is a pseudo-differential operator, which can be given a precise meaning by the calculus of pseudo-differential operators, see Section 2.4. In the following, we will replace square root operator  $\sqrt{1 + X}$ , with  $X = 1/k^2 \partial_x^2$ , by an rational expression. Such approximation can be formally written as

$$(6.2.3) \quad \partial_z u = ik \frac{(1 - b'_0 \partial_x^2) (1 - b'_1 \partial_x^2) \cdots (1 - b'_m \partial_x^2)}{(1 - b_0 \partial_x^2) (1 - b_1 \partial_x^2) \cdots (1 - b_n \partial_x^2)} u$$

with complex coefficients  $b'_0, \dots, b'_m$  and  $b_0, \dots, b_n$ . The approximation quality and the well-posedness of this formal equation has been examined in e. g. [89] and [6]. Additionally, the non-commutativity of the factors in (6.2.3) in case of non-constant coefficients  $b'_0, \dots, b'_m$  and  $b_0, \dots, b_n$  on the computational domain is the source of theoretical investigations which are not discussed here. Rather, we adopt the usual technique of *frozen coefficients*, in which the functions are considered to be constant in deriving the rational approximation, cf. e. g. [6].

From a physical perspective, while originally developed in computational underwater acoustics propagation methods based on 6.2.3 have proved critically important in many other contexts, as evident from some representative references, e. g. [20, 19, 48, 93, 59]. A high-order Padé approximation to the Helmholtz propagator ensures that widely-divergent beams are described to great accuracy. Thus the Padé approximation to the Helmholtz equation supplies a model which ranges between the Schrödinger equation and the Helmholtz equation. It combines the simplicity of an evolution algorithm with the generality of the Helmholtz equation.

In the following sections, we will accordingly investigate the discrete solution to the following problem:

- (1) Let the coefficients  $b'_0, \dots, b'_m$  and  $b_0, \dots, b_n$  be piecewise continuous real functions on the bounded computational domain  $\Omega := (x_-, x_+) \times (0, z_{\max})$  and real constants outside the computational domain.
- (2) Let  $u$  be the solution of (6.2.3) on the unbounded domain  $\mathbb{R} \times [0, z_{\max}]$  with initial conditions  $u(x, 0) = u_0(x)$  compactly supported on the interval  $[x_-, x_+] \subset \mathbb{R}$  together with the asymptotic boundary condition  $\lim_{|x| \rightarrow \infty} u(x, z) = 0$  holds for all  $0 \leq z \leq z_{\max}$ .
- (3) Determine the solution  $u(x)$  on the bounded computational domain  $\Omega$  that agrees exactly with the unbounded result, restricted to  $\Omega$ .

The goal of this section is to construct transparent boundary conditions for arbitrary higher order evolution equations of the type (6.2.3) that insure the realization of the third point above. Similar to the Schrödinger equation, various techniques have been developed for the solution of such boundary problems. Most of all the techniques discussed in Section 6.1 has been applied to the wide-angle one-way equation, too.

Recently, a specialization of the problem defined in this introduction based on a Padé (2,2) approximation to the Helmholtz operator has been successfully analyzed by Arnold and Erhardt [5]. In this reference, the authors follow the common way with forth- and back Laplace transforming in *time*. As outlined in Section 6.1, this implies to compute the *inverse* Laplace transform, which restricts the applicability of the method.

We extend the properties of our method described above as follows:

- We allow for an arbitrary high order Padé-approximation of the square root operator on the entire transversal axes.
- We do not need to introduce a Green's function representation for the chosen Padé approximation.
- We do not need to compute an explicite inverse Laplace transform.
- We prove unconditional stability of the algorithm.

As a result, we are able to develop numerical techniques even for high Padé orders. It should be noted that this presentation is restricted to uniform step-sizes in the propagation direction, so that we can apply the shift-operator technique already used in the previous sections. Similar formulas can however be derived for non-equidistant step sizes using the algebraic approach of Section 6.1, but the analysis as well as the resulting formulas would be far more complicated. So we drop this here.

### Standard Wide-Angle Approximations.

We first summarize the standard wide-angle approximation to the Helmholtz equation (6.2.1) that will form the basis for our subsequent considerations. The first step in this analysis is to define a new field variable  $\tilde{u}(x, z) := u(x, z) \exp(-ik_0 z)$ . The reference-wavenumber  $k_0$  is chosen to effectively minimize the mean phase



TABLE 1. Padé coefficients

	(2, 0)		(2, 2)		(4, 2)			(4, 4)		
$j$	0	2	0	2	0	2	4	0	2	4
$c'_j$	1	1/2	1	3/4	1	1	1/8	1	5/4	5/16
$c_j$	1		1	1/4	1	1/2		1	3/4	1/16

velocity of the wavevector components of  $\tilde{u}(x, z)$  in the  $z$ -direction. The resulting spectrally shifted Helmholtz equation, where we drop the tilde on  $u(x, z)$  for notational simplicity, is

$$(6.2.4) \quad (\partial_z^2 + 2ik_0\partial_z + \partial_x^2 + k^2 - k_0^2) u(x, z) = 0.$$

Formally factorizing the above equation in analogy with (6.2.2) yields the following one-way equation after the scaling  $z := z \cdot k_0$  and  $x := x \cdot k_0$

$$(6.2.5) \quad \partial_z u = -i \left(1 - \sqrt{1 + X^2}\right) u, \quad \text{with } X^2 := \frac{k^2 - k_0^2}{k_0^2} + \partial_x^2.$$

Since  $k = \text{const}$  in the exterior domain, rational approximations of the form

$$(6.2.6) \quad \sqrt{1 + X^2} \simeq \frac{c'_0 + c'_2 X^2 + \dots + c'_{2m} X^{2m}}{c_0 + c_2 X^2 + \dots + c_{2n} X^{2n}}.$$

can be applied to the square-root operator. The most accurate of these are generally obtained for  $m = n$  or  $m = n + 1$ , see [89]. We will derive a general procedure for deriving transparent boundary conditions that can be applied to all such approximations in succeeding sections. However, to simplify our numerical algorithm, we will consider only Padé-type approximations of order  $(2m, 2n)$ . This yields an interpolation error  $O(X^{2m+2n+2})$  for  $X \rightarrow 0$ .

For the simple Padé approximation of (6.2.6) the coefficients  $c'_{2i}, c_{2j}$ ,  $i = 0, \dots, m$  and  $j = 0, \dots, n$  can be obtained through either an explicit factorization ([7]) or by the Newman-procedure [53]. The latter technique, upon which our computer codes are based, can be extended to wide-angle equations other than those based on Padé-approximations. Some of the resulting coefficients can be found in Table 1.

### Longitudinal Discretization.

The implicit midpoint discretization of (6.2.5) results in

$$(6.2.7) \quad \frac{u_i(x) - u_{i-1}(x)}{\Delta z} = -i \left(1 - \sqrt{1 + X^2}\right) \frac{u_i(x) + u_{i-1}(x)}{2}.$$

Here  $u_i(x)$ ,  $0 < i \leq n$  denotes  $u(x, z_0 + i\Delta z)$  with  $z_0$  the initial value of the longitudinal distance and  $\Delta z$  the propagation step length. Note here that the subscripts  $i$  denote the longitudinal stepnumber. Superscript  $i$  values instead denote the  $i$ th integer power of the corresponding quantity while algebraic factors of  $i$  refer to the imaginary unit. From (6.2.7), we immediately obtain the discrete evolution equation

$$\left(1 + \frac{i\Delta z}{2} \left(1 - \sqrt{1 + X^2}\right)\right) u_i(x) = \left(1 - \frac{i\Delta z}{2} \left(1 - \sqrt{1 + X^2}\right)\right) u_{i-1}(x).$$

After replacing the square-root operator by its rational approximant we obtain an equation of the form

$$(6.2.8) \quad u_i(x) = \frac{P'(X^2)}{P(X^2)} u_{i-1}(x).$$

Here  $P(X^2)$  and  $P'(X^2)$  are the following polynomials of degree  $m$  in the variable  $X^2$

$$\begin{aligned} P'(X^2) &= \left(1 - i\frac{\Delta z}{2}\right) C(X^2) + i\frac{\Delta z}{2} C'(X^2) \quad \text{and} \\ P(X^2) &= \left(1 + i\frac{\Delta z}{2}\right) C(X^2) - i\frac{\Delta z}{2} C'(X^2). \end{aligned}$$

The quantities  $C'(X^2)$  and  $C(X^2)$  are themselves polynomials defined by  $C'(X^2) = c'_0 + c'_2 X^2 + \dots + c'_{2m} X^{2m}$  and  $C(X^2) = c_0 + c_2 X^2 + \dots + c_{2n} X^{2n}$ , see (6.2.6) and Tab. (1). Applying a complex root finder yields

$$(6.2.9) \quad P'(X^2) = c' \prod_{j=1}^k (1 - a'_j X^2) \quad \text{and} \quad P(X^2) = c \prod_{j=1}^k (1 - a_j X^2),$$

with  $c$  and  $c'$  are constants. Choosing  $X^2 = 0$  it follows  $P'(0)/P(0) = c'/c$  and comparing with (6.2.7) shows  $c'/c = 1$ . Finally, inserting  $X^2 := (k^2 - k_0^2)/k_0^2 + \partial_x^2$  leads to the desired factorization.

### Discrete Evolution Equation.

From the rational approximation (6.2.8) and the factorization (6.2.9) we obtain the longitudinally discretized form of the evolution equation

$$(6.2.10) \quad u_i(x) = \left(\frac{1 - a'_k \partial_x^2}{1 - a_k \partial_x^2}\right) \dots \left(\frac{1 - a'_2 \partial_x^2}{1 - a_2 \partial_x^2}\right) \left(\frac{1 - a'_1 \partial_x^2}{1 - a_1 \partial_x^2}\right) u_{i-1}(x),$$

which is the exact counterpart of the continuous evolution equation (6.2.3). In terms of the intermediate functions  $g_i^{(1)}(x), \dots, g_i^{(k-1)}(x)$  given by

$$\begin{aligned} g_i^{(1)}(x) &= \left(\frac{1 - a'_1 \partial_x^2}{1 - a_1 \partial_x^2}\right) u_{i-1}(x) \\ g_i^{(2)}(x) &= \left(\frac{1 - a'_2 \partial_x^2}{1 - a_2 \partial_x^2}\right) g_i^{(1)}(x) \\ &\vdots \\ g_i^{(k-1)}(x) &= \left(\frac{1 - a'_{k-1} \partial_x^2}{1 - a_{k-1} \partial_x^2}\right) g_i^{(k-2)}(x) \\ u_i(x) &= \left(\frac{1 - a'_k \partial_x^2}{1 - a_k \partial_x^2}\right) g_i^{(k-1)}(x), \end{aligned}$$

the factorized, order  $2k$  discrete evolution problem (6.2.10) can be recast into the following system of  $k$  second-order differential equations

$$(6.2.11) \quad \begin{aligned} (1 - a_1 \partial_x^2) g_i^{(1)}(x) - (1 - a'_1 \partial_x^2) u_{i-1}(x) &= 0 \\ (1 - a_2 \partial_x^2) g_i^{(2)}(x) - (1 - a'_2 \partial_x^2) g_i^{(1)}(x) &= 0 \\ &\vdots \\ (1 - a_{k-1} \partial_x^2) g_i^{(k-1)}(x) - (1 - a'_{k-1} \partial_x^2) g_i^{(k-2)}(x) &= 0 \\ (1 - a_k \partial_x^2) u_i(x) - (1 - a'_k \partial_x^2) g_i^{(k-1)}(x) &= 0. \end{aligned}$$

This sequence of equations may be rewritten conveniently into matrix form

$$(6.2.12) \quad \left[ \begin{pmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & -1 & 1 & & \\ & & & -1 & 1 & \\ & & & & \ddots & \ddots \end{pmatrix} + \begin{pmatrix} -a_1 \partial_x^2 & & & & & \\ & \ddots & & & & \\ & & a'_k \partial_x^2 & -a_k \partial_x^2 & & \\ & & & a'_1 \partial_x^2 & -a_1 \partial_x^2 & \\ & & & & \ddots & \ddots \end{pmatrix} \right] g(x) = 0,$$

where the unknown functions are assembled into the vector

$$g(x) = \left( g_1^{(1)}, \dots, g_1^{(k-1)}, u_1, \dots, \underbrace{g_i^{(1)}, \dots, g_i^{(k-1)}}_{i\text{th macro step}}, u_i, \dots \right)^T.$$

Here quantities distinguished by a subscript  $i$  are evaluated during the  $i$ th physical propagation step, while the  $k-1$  intermediate partial steps that comprise each physical step and that are associated with different terms in the Padé approximation are distinguished through their superscripts.

Given an initial condition with respect to  $x$ , at some fixed point, say  $x_+$ , we can solve the system (6.2.12) equation by equation. Thus at step  $i$  the function  $u_{i-1}(x)$  maps over  $k-1$  intermediate functions  $g_i^{(j)}(x)$  to the solution  $u_i(x)$ . Hence we consider a mapping

$$(6.2.13) \quad u_{i-1}(x) \mapsto \underbrace{g_i^{(1)}(x) \mapsto g_i^{(2)}(x) \mapsto \dots \mapsto g_i^{(k-1)}(x)}_{g_i(x) := (g_i^{(1)}(x), g_i^{(2)}(x), \dots, g_i^{(k-1)}(x), u_i(x))} \mapsto u_i(x)$$

$$(6.2.14) \quad u_{i-1}(x) \mapsto g_i(x),$$

which depends on the initial conditions (with respect to  $x$ ). In the lowest order case  $k=1$  the vector  $g_i(x)$  consists of the single function  $u_i(x)$ .

To write (6.2.11) or (6.2.12), respectively, as a simple matrix equation with finite dimension, we introduce the discrete shift-operator  $s_k$  with the property that  $s_k u_i(x) = u_{i-1}(x)$ . This operator shifts our sequence by one physical step, corresponding to  $k$  terms in  $g$ . That is, if  $s$  is the shift operator which shifts the sequence in  $g$  by one place, we have  $s_k = s^k$ . Since in the following the index shift occurs only in powers of  $s$ , we redefine  $s := s_k$  to yield the more intuitive formulation denoted by  $su_i(x) = u_{i-1}(x)$ .

Introducing the notation  $g'$  for  $\partial_x g$  and eliminating  $u_{i-1}(x)$  in (6.2.12) in terms of  $u_i(x)$  then yields

$$(6.2.15) \quad (E(s) + A(s)\partial_x^2) g_i(x) = 0 \text{ for all macro steps } i \geq 1,$$

which must be solved subject to boundary conditions

$$(6.2.16) \quad g'_{i,+} = B_+(s)g'_{i,+} \text{ and } g'_{i,-} = B_-(s)g_{i,-}.$$

A comparison with the infinite dimensional system (6.2.12) demonstrates that the  $k \times k$ -matrices  $E(s)$  and  $A(s)$  and the  $k$ -element vectors  $g_i(x)$ ,  $g_{i,\pm}$  and  $g'_{i,\pm}$  are given by

$$E(s) = \begin{pmatrix} 1 & & & -s \\ -1 & 1 & & \\ & & \ddots & \\ & & -1 & 1 \\ & & & -1 & 1 \end{pmatrix}, A(s) = \begin{pmatrix} -a_1 & & & & sa'_1 \\ a'_2 & -a_2 & & & \\ & & \ddots & & \\ & & & a'_{k-1} & -a_{k-1} \\ & & & & a'_k & -a_k \end{pmatrix},$$

$$g_i(x) = \begin{pmatrix} g_i^{(1)}(x) \\ g_i^{(2)}(x) \\ \vdots \\ g_i^{(k-1)}(x) \\ u_i(x) \end{pmatrix}, \quad g_{i,\pm} = \begin{pmatrix} g_i^{(1)} \\ g_i^{(2)} \\ \vdots \\ g_i^{(k-1)} \\ u_i \end{pmatrix}_{x=x_{\pm}}, \quad g'_{i,\pm} = \begin{pmatrix} \dot{g}_i^{(1)} \\ \dot{g}_i^{(2)} \\ \vdots \\ \dot{g}_i^{(k-1)} \\ \dot{u}_i \end{pmatrix}_{x=x_{\pm}}.$$

We refer to  $E(s)$  and  $A(s)$  as operators, as they depend explicitly on the shift-operator. Further, they may be represented as  $D(s) = \sum_{j=0}^i D_j s^j$ , where the  $k \times k$  matrices  $D_j$ ,  $j = 0, \dots, i$  are complex. We will label operators of this general type convolution operators or convolution matrices; the properties of such matrices that are required to derive our subsequent algorithms are collected in the Appendix. The Dirichlet-to-Neumann operators  $B_{\pm}(s)$  that implement the desired boundary conditions must be constructed such that the asymptotic boundary condition  $\lim_{|x| \rightarrow \infty} u_i(x) = 0$  is fulfilled for all propagation steps. We will show later by construction that these boundary operators are also of convolution type.

### Discrete Transparent Boundary Conditions.

We now present our derivation of transparent boundary conditions for general wide-angle methods. To simplify the discussion, we consider  $g_i(x)$  only in the right exterior domain,  $x \geq x_+$ , and further shift the position of the right boundary  $x_+$  to the origin according to  $x \mapsto x + x_0$  and omit the  $\pm$ -subscript. As well, we designate the boundary value  $g_{i,+}$  by  $g_{i,0}$ . All our results of course apply equally to the left exterior domain.

Consider (6.2.15) as an initial value problem in the right exterior domain with data given on the shifted boundary  $x_+ = 0$ . Our objective is to construct an operator with the property that the corresponding exterior solution decays asymptotically for any given Dirichlet data  $g_{i,0}$ . To do this, we construct the Laplace transform  $\widehat{g}_i$  of the exterior solution vector  $g_i(x)$ . In the exterior domain, the Laplace transform of the equation system (6.2.15) is

$$(E(s) + p^2 A(s)) \widehat{g}_i(p) = A(s) (p g_{i,0} + g'_{i,0}).$$

The operator  $A(s)$  is invertible, since  $a_j \neq 0$ ,  $j = 1, \dots, k$ , cf. Section 6.2 and Lemma A.2.4. Therefore we can write equivalently

$$(p^2 \mathbf{I} - C^2(s)) \widehat{g}_i(p) = p g_{i,0} + g'_{i,0},$$

where the  $k \times k$ -matrix  $C^2(s) := -A(s)^{-1} E(s)$  is of convolution type and  $\mathbf{I}$  is the  $k \times k$ -identity-matrix. If we assume for the moment that we can obtain the square roots,  $\pm C(s)$ , of the matrix  $C^2(s)$  we have  $(p^2 \mathbf{I} - C^2(s)) = (p \mathbf{I} + C(s))(p \mathbf{I} - C(s))$ , since  $C(s)$  and  $\mathbf{I}$  commute. Thus the ansatz  $g'_{i,0} = B g_{i,0}$ , see (6.2.15), with  $B(s) = -C(s)$  eliminates all poles which belong to asymptotically increasing solutions, hence we derived the

Pole condition for wide angle one-way equations

$$g'_{i,0} = B g_{i,0}, \quad \text{with } B(s) = -C(s)$$

Further, we obtain the *solution representation*,

$$(6.2.17) \quad \widehat{g}_i(p) = (p \mathbf{I} + C(s))^{-1} g_{i,0} \quad \text{subject to boundary conditions } g'_{i,0} = B(s) g_{i,0}.$$

To derive the nonlocal boundary conditions, which is equivalent to finding the boundary operator  $B(s)$ , we apply the same procedure as in [82] or [85]. That is, we construct the matrix  $C(s)$  such that all poles  $p_j$ ,  $j = 1, \dots, k$  of  $(pI + C(s))^{-1}$  are located in the right half of the complex plane, i.e.  $\text{Re } p_j > 0$ ,  $j = 1, \dots, k$ . This ensures that the exterior solution decays appropriately as  $x \rightarrow \pm\infty$ . The operator  $(pI + C(s))^{-1}$  possesses a pole in  $p$  if the matrix  $(pI + C(s)|_{s=0})^{-1}$  possesses a pole in  $p$ , cf. the properties of operators of convolution type given in the Appendix.

The square root of  $C^2(s) := -A^{-1}(s)E(s)$  is obtained by decomposing the matrices  $A(s)$  and  $E(s)$  into its components with respect to  $s$ . The first of these is independent of the shift operator  $s$ , namely  $A_0 := A|_{s=0}$  and  $E_0 := E|_{s=0}$  while the second is  $s$ -dependent. We thus have  $A = A_0 + sA_1$  and  $E = E_0 + sE_1$  with

$$A_1 := \begin{pmatrix} 0 & \cdots & 0 & a'_1 \\ & & & 0 \\ & \mathbf{0} & & \vdots \\ & & & 0 \end{pmatrix} \quad E_1 := \begin{pmatrix} 0 & \cdots & 0 & -1 \\ & & & 0 \\ & \mathbf{0} & & \vdots \\ & & & 0 \end{pmatrix}.$$

For convenience, we define new matrices

$$\tilde{A}_1 := A_0^{-1}A_1 \quad \tilde{E}_0 := A_0^{-1}E_0 \quad \tilde{E}_1 := A_0^{-1}E_1.$$

From the previous definitions together with the ansatz  $C(s) = \sum_{j=0}^{i-1} C_j s^j$ , where  $C_j \in \mathbb{C}^{k \times k}$ ,  $0 \leq j \leq i$ , we observe that we must find matrices  $C_j$  such that

$$(I + \tilde{A}_1 s) (C_0 + C_1 s + C_2 s^2 + \dots) (C_0 + C_1 s + C_2 s^2 + \dots) = -\tilde{E}_0 - \tilde{E}_1 s.$$

It is a property of matrices of convolution type that the expression (6.2.18) can be computed element-wise and further ordered with respect to equal powers of  $s$  such that the resulting operator is of convolution type. That is we are allowed to reorder, for example,  $C_1 s C_0$  to  $C_1 C_0 s$ . (This property, of course, is not generally true.) Now, as is evident by comparing coefficients, to solve (6.2.18) we must first find  $C_0$  such that  $C_0^2 = -\tilde{E}_0$ . Because the matrix  $\tilde{E}_0$  is a quotient of two lower triangular matrices, the matrix  $C_0$  is also lower triangular. Further, the diagonal entries  $C_0$  can be chosen such that  $\text{Re}(C_{jj}) > 0$ ,  $j = 1, \dots, k$ , through the Algorithm (3) below. Hence all poles corresponding to  $C_0$ , with  $C_0$  constructed according to Algorithm 3, are located in the right half of the complex plane.

---

**Algorithm 3** Calculate  $C_0 = \sqrt{-\tilde{E}_0}$

---

```

for  $i = 1$  to  $k$  do
   $c_{ii} = \sqrt{-\tilde{E}_{0,ii}}$ 
  if  $i > 1$  then
    for  $j = i - 1$  to  $1$  do
       $c_{ij} = (-\tilde{E}_{0,ij} - \sum_{m=j+1}^{i-1} c_{im}c_{mj}) / (c_{ii} + c_{jj})$ 
    end for
  end if
end for

```

---

Subsequently the sequence  $C_1, C_2, \dots, C_{n-1}$  for the following  $n$  propagation steps is obtained by comparing coefficients of equal powers of  $s$  in (6.2.18). The corresponding pseudo-code is given in Algorithm 4, and consists mainly of solutions of Sylvester equations. From the structure of the algorithm we observe that if  $C_0$  is computed, the entire sequence is uniquely determined.

---

**Algorithm 4** Recursive calculation of  $C_j$ ,  $j = 1, \dots, n-1$

---

$Z := -(E_1 + A_1 C_0^2)$   
 Compute  $C_1$  from  $C_1 C_0 + C_0 C_1 = Z$   
**for**  $k = 2$  to  $n-1$  **do**  
      $Z \leftarrow -A_1 Z$   
     Compute  $C_k$  from  $C_k C_0 + C_0 C_k = Z - \sum_{j=1}^{k-1} C_j C_{k-j}$   
**end for**

---

Finally, the boundary condition at every step  $0 < i \leq n$  is obtained from  $C_j$ ,  $j = 0, \dots, n-1$ , by employing the definitions  $B = -C$  and

$$\begin{aligned}
 (6.2.19) \quad g'_{i,0} &= B(s)g_{i,0} \\
 &= (B_0 + sB_1 + \dots + s^{i-1}B_{i-1})g_{i,0} \\
 &= B_0g_{i,0} + B_1g_{i-1,0} + \dots + B_{i-1}g_{1,0}.
 \end{aligned}$$

Equation (6.2.19) provides the algorithmic basis for constructing nonlocal boundary conditions for any wide-angle approximation and discrete propagation method, as different propagation methods can be distinguished simply through the values of the defining coefficients  $a'_1, \dots, a'_k$  and  $a_1, \dots, a_k$ . As the operator  $B(s)$  possesses a Taylor representation in  $s$ , its action can be represented by matrix-vector multiplications of the Taylor coefficients  $B_j$  with boundary values describing the history of the evolution process. In our numerical implementation we order the boundary vectors  $g_{i,0}$  from lower to larger step numbers and introduce a composite boundary vector  $g_0 = (g_{1,0}^T, \dots, g_{i-1,0}^T, g_{i,0}^T)^T$ . Similarly, we generate the composite boundary matrices

$$(6.2.20) \quad B = (B_{i-1}, B_{i-2}, \dots, B_0), \quad 1 \leq i \leq n$$

$$(6.2.21) \quad C = -B.$$

in place of the boundary operator  $B(s)$ , after which the normal derivative  $g'_{i,0}$  from system (6.2.19) is computed through a matrix-vector multiplication. This procedure for implementing the discrete boundary condition in terms of a composite matrix  $C$  is summarized in Algorithm 7.

### Finite-Element Discretization.

From the representation (6.2.11), we will now generate a finite-element discretization on the interior domain. For illustrative purposes, consider the first equation of the system (6.2.11) at step  $i$ ,  $0 < i \leq n$ . Multiplying this equation by a trial function  $v \in H^1(\Omega)$ ,  $\Omega = (x_-, x_+)$ , integrating over  $\Omega$ , and integrating by parts yields

$$\begin{aligned}
 (6.2.22) \quad & \left( v, g_i^{(1)} \right) + \left( \partial_x v, a_1 \partial_x g_i^{(1)} \right) - \left( \bar{v} a_1 \partial_x g_i^{(1)} \right) \Big|_{x_-}^{x_+} = \\
 & \left( v, u_{i-1} \right) + \left( \partial_x v, a'_1 \partial_x u_{i-1} \right) - \left( \bar{v} a'_1 \partial_x u_{i-1} \right) \Big|_{x_-}^{x_+}.
 \end{aligned}$$

The variational problem corresponding to this equation is therefore to find a function  $g_i^{(1)} \in H^1(\Omega)$  such that (6.2.22) holds for any  $v \in H^1(\Omega)$ . The other equations of the system (6.2.11) can be reformulated similarly. The resulting system is then discretized by replacing the infinite-dimensional function space  $H^1(\Omega)$  by a finite-dimensional space  $V_h = \text{span}\{v_1, v_2, \dots, v_N\}$  with  $V_h \subset H^1(\Omega)$ . Hence the corresponding discrete problem is to determine a discrete approximation  $g_{h,i}^{(1)}$  of  $g_i^{(1)}$ .

with  $g_{h,i}^{(1)} \in V_h$  such that for all  $v_h \in V_h$

$$(6.2.23) \quad \begin{aligned} & (v_h, g_{h,i}^{(1)}) + \left( \partial_x v_h, a_1 \partial_x g_{h,i}^{(1)} \right) - \left( \bar{v}_h a_1 \partial_x g_{h,i}^{(1)} \right) \Big|_{x_-}^{x_+} \\ & = (v_h, u_{h,i-1}) + \left( \partial_x v_h, a'_1 \partial_x u_{h,i-1} \right) - \left( \bar{v}_h a'_1 \partial_x u_{h,i-1} \right) \Big|_{x_-}^{x_+}. \end{aligned}$$

To solve the above system, we employ standard linear  $C^0$ -elements. Let  $v_1$  and  $v_N$  denote the leftmost and the rightmost finite elements. A compact notation for (6.2.23) results if we set  $v_1|_{x_-} = 1$ ,  $v_N|_{x_+} = 1$ , and define the vectors  $b_i^{(1)}, b_i'^{(1)} \in \mathcal{C}^N$  with  $N = \dim V_h$  by

$$b_i^{(1)} = \begin{pmatrix} \left( -a_1 \partial_x g_{h,i}^{(1)} \right) \Big|_{x=x_-} \\ 0 \\ \vdots \\ 0 \\ \left( a_1 \partial_x g_{h,i}^{(1)} \right) \Big|_{x=x_+} \end{pmatrix} \quad \text{and} \quad b_i'^{(1)} = \begin{pmatrix} \left( -a'_1 \partial_x u_{h,i-1} \right) \Big|_{x=x_-} \\ 0 \\ \vdots \\ 0 \\ \left( a'_1 \partial_x u_{h,i-1} \right) \Big|_{x=x_+} \end{pmatrix},$$

and introduce the mass and stiffness matrices  $M \in \mathbb{R}^{N \times N}$ ,  $A_1, A'_1 \in \mathbb{C}^{N \times N}$  in standard fashion as  $(M)_{i,j} = (v_{h,i}, v_{h,j})$  and  $(A)_{i,j} = (\partial_x v_{h,i}, a_1 \partial_x v_{h,j})$ . Defining as well the vectors  $g_i^{(1)} = (g_{h,i,1}^{(1)}, \dots, g_{h,i,N}^{(1)})^T \in \mathbb{C}^N$  and  $u_{i-1} = (u_{h,i-1,1}, \dots, u_{h,i-1,N})^T \in \mathbb{C}^N$ , that are the discrete counterparts of the continuous functions  $g_i^{(1)}(x)$ ,  $u_i(x)$  we have

$$(6.2.24) \quad (M + A_1) g_i^{(1)} - b_i^{(1)} = (M + A'_1) u_{i-1} - b_i'^{(1)}.$$

If we know the solution  $u_{i-1}$  in the interior domain together with its normal derivative on the boundary and the normal derivative of  $g_i^{(1)}$  we can obtain the unknown intermediate vector  $g_i^{(1)}$ . Repeating this procedure for each of the equations of (6.2.11), we generate the following block matrix equation in terms of the matrices and vectors introduced in the preceding paragraph,

$$(6.2.25) \quad \begin{aligned} & \begin{pmatrix} M + A_1 & & & \\ & M + A_2 & & \\ & & \ddots & \\ & & & M + A_k \end{pmatrix} \begin{pmatrix} g_i^{(1)} \\ g_i^{(2)} \\ \vdots \\ u_i \end{pmatrix} - \begin{pmatrix} b_i^{(1)} \\ b_i^{(2)} \\ \vdots \\ b_i^{(k)} \end{pmatrix} \\ & = \begin{pmatrix} M + A'_1 & & & \\ & M + A'_2 & & \\ & & \ddots & \\ & & & M + A'_k \end{pmatrix} \begin{pmatrix} u_{i-1} \\ g_i^{(1)} \\ \vdots \\ g_i^{(k-1)} \end{pmatrix} - \begin{pmatrix} b_i'^{(1)} \\ b_i'^{(2)} \\ \vdots \\ b_i'^{(k)} \end{pmatrix}. \end{aligned}$$

To solve the system (6.2.25), the vectors  $b_i^{(j)}, b_i'^{(j)}$ ,  $j = 1, \dots, k$  must be constructed in accordance with the boundary conditions. The relationship between the discretized evolution equation (6.2.25) and the boundary condition (6.2.19) is determined by first decomposing the boundary condition at each boundary according to

$b_i^{(j)} = b_{i,-}^{(j)} + b_{i,+}^{(j)}$  with

$$b_{i,-}^{(j)} = \begin{pmatrix} \left. (-a_j \partial_x g_{i,h}^{(j)}) \right|_{x=x_-} \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \text{ and } b_{i,+}^{(j)} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \left. (a_j \partial_x g_{i,h}^{(j)}) \right|_{x=x_+} \end{pmatrix}, \quad b_{i,\pm}^{(j)} \in \mathbb{C}^N.$$

Performing the same decomposition for each vector  $b_i^{(j)}$  and assembling all nonzero entries of the vectors  $b_{i,\pm}^{(j)}, b'_{i,\pm}^{(j)}, j = 1, \dots, k$  into the four vectors  $b_{i,\pm}, b'_{i,\pm} \in \mathbb{C}^k$  we arrive at

$$(6.2.26) \quad b_{i,\pm} = \begin{pmatrix} \left. (\pm a_1 \partial_x g_{i,h}^{(1)}) \right|_{x=x_\pm} \\ \left. (\pm a_2 \partial_x g_{i,h}^{(2)}) \right|_{x=x_\pm} \\ \vdots \\ \left. (\pm a_{k-1} \partial_x g_{i,h}^{(k-1)}) \right|_{x=x_\pm} \\ \left. (\pm a_k \partial_x u_i) \right|_{x=x_\pm} \end{pmatrix} \text{ and } b'_{i,\pm} = \begin{pmatrix} \pm \left. (a'_1 \partial_x u_{i-1}) \right|_{x=x_\pm} \\ \pm \left. (a'_2 \partial_x g_{i,h}^{(1)}) \right|_{x=x_\pm} \\ \vdots \\ \pm \left. (a'_{k-1} \partial_x g_{i,h}^{(k-2)}) \right|_{x=x_\pm} \\ \pm \left. (a'_k \partial_x g_{i,h}^{(k-1)}) \right|_{x=x_\pm} \end{pmatrix}.$$

We now derive an equation relating the vectors  $b_{i,\pm}$  and  $b'_{i,\pm}$  to the boundary condition (6.2.19). Regarding first  $b_{i-1,+}$ , we have from (6.2.26) and (6.2.19)

$$\begin{aligned} b_{i,+} &= \text{diag}(a_1, \dots, a_k) g'_{0,i} \\ &= \text{diag}(a_1, \dots, a_k) B(s) g_{0,i} \\ &= \text{diag}(a_1, \dots, a_k) (B_0 + B_1 s + B_2 s^2 + \dots + B_{i-1} s^{i-1}) g_{0,i}. \end{aligned}$$

All of the expressions  $B_j s^j g_n(0) = B_j g_{n-j}(0)$  above with  $j = 1, \dots, n-1$  can be immediately evaluated based on the observation that the shift operator  $s$  decreases the index of  $B_j$  by unity. Further, as the matrix  $B_0$  is a lower triangular matrix we can arrange the algorithm such that the boundary condition at the current step only depends on boundary values at previous steps. To this end we decompose  $b_\pm$  as

$$(6.2.27) \quad \begin{aligned} b_\pm &= B_{d,\pm} g_{0,i} + B_{r,\pm} g_i \\ \text{with } B_{d,\pm} &:= \pm \text{diag}(a_1, \dots, a_k) \Big|_{x=x_\pm} \text{diag}(B_{0,\pm}) \\ \text{and } B_{r,\pm}(s) &:= \pm \text{diag}(a_1, \dots, a_k) \Big|_{x=x_\pm} \\ (6.2.28) \quad &\cdot (B_0 - \text{diag}(B_0) + s B_1 + \dots + s^{i-1} B_{i-1}) \Big|_{x=x_\pm}, \end{aligned}$$

where the matrices  $\text{diag}(B_{0,\pm})$  contain only the main diagonals of  $B_{0,\pm}$ . The diagonal matrices  $B_{d,\pm}$  are then inserted into the matrix of the final system which yields updated matrices  $A_j$  satisfying

$$(6.2.29) \quad A_j = A_j - \begin{pmatrix} B_{d,-}(j,j) & & \\ & 0 & \\ & & B_{d,+}(j,j) \end{pmatrix}, \quad j = 1, \dots, k.$$

The reduced matrix  $B_{r,\pm}$ , which is a lower triangular matrix with a zero diagonal, only couples previously determined boundary values and is therefore placed on the right-hand side of the evolution equation, cf. Algorithm 7.

A corresponding expression for  $b'_+$  results from the observation that the vector  $g_i$  is ordered as  $(g_i^{(1)}, g_i^{(2)}, \dots, g_i^{(k-1)}, u_i)^T$ , reflecting the algebraic structure of the



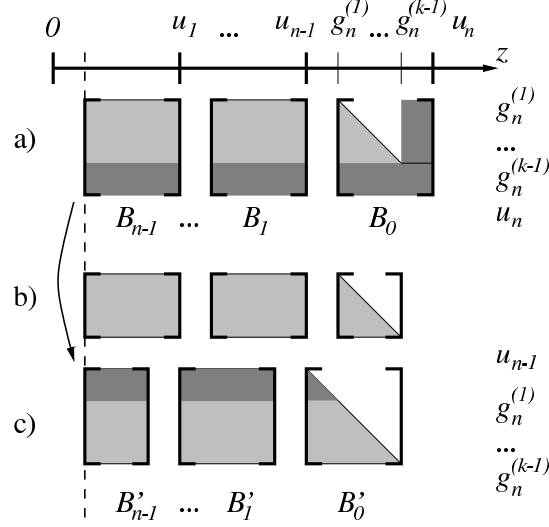


FIGURE 6.2.1. Construction of the boundary operator  $B'$  from the boundary operator  $B$

boundary condition (6.2.19). The discrete evolution system (6.2.25), however, requires the alternate ordering (6.2.26)  $(u_{i-1}, g_i^{(1)}, g_i^{(2)}, \dots, g_i^{(k-1)})^T$ . Thus we obtain a condition of the form

$$\begin{pmatrix} \partial u_{i-1} \\ \vdots \\ \partial g_i^{(k-1)} \end{pmatrix} = (B'_0 + B'_1 s + \dots + B'_{i-1} s^{i-1}) \begin{pmatrix} u_{i-1} \\ \vdots \\ g_i^{(k-1)} \end{pmatrix}$$

with an operator  $B'(s) := B'_0 + B'_1 s + \dots + B'_{i-1} s^{i-1}$ , if we rearrange the columns and rows of the operator  $B$  appropriately.

In the first step of this transformation, which is illustrated in Fig. 6.2.1 we remove the last row and the last column of the composite boundary matrix  $B = -C$  to obtain the reduced matrix shown in part b) of Fig. 6.2.1. We then place the former last row of  $B$  at the top of the reduced matrix, after shifting the row to the left and adjusting it to the dimension of the reduced matrix, as illustrated in part c) of Fig. 6.2.1. Finally, the resulting matrix is multiplied with the diagonal matrix  $\text{diag}(a'_1, \dots, a'_k)$  according to (6.2.26). The pseudo-code for these operations is given in Algorithm 5.

---

**Algorithm 5** Computation of the operator  $B'$  (see Fig. 6.2.1)

---

- 1: Compute  $B' := -C(1 : k - 1, 1 : nk - 1)$
  - 2: Compute  $b' := - \begin{bmatrix} C(k, k + 1 : nk), \underbrace{0, 0, \dots, 0}_{k-1} \end{bmatrix}$
  - 3: Compute  $B' \leftarrow \text{diag}(a'_1, \dots, a'_k) \begin{bmatrix} 0 & b' \\ 0 & B' \end{bmatrix}$
- 

To evolve the field over  $n$  propagation steps, we must first initialize the two boundary operators  $B_+$  and  $B_-$  with dimension  $k \times (kn)$ , acting on the right and left boundary, respectively. At the end of the simulation we will then possess two vectors  $g_+$  and  $g_-$  of dimension  $nk + 1$  that contain the required boundary values.

The initialization of the propagation algorithm, which includes the computation of the standard finite element matrices, the updating of these matrices (which corresponds to incorporating the boundary conditions) and the computation of the boundary matrices  $B_{\pm}$  are summarized in Algorithm 6.

---

**Algorithm 6** Computation of the finite element matrices  $A_j$ ,  $M$  and the boundary matrices  $B_{\pm}$

---

Compute $C_{0,\pm}$	{acc. to Algorithm 3}
Compute $C_{j,\pm}$ , $j = 1 : n - 1$	{acc. to Algorithm 4}
Compute $C_{\pm} := [C_{n-1}, C_{n-2}, \dots, C_0]_{\pm}$	{acc. to (6.2.21)}
Compute $A_j$ , $j = 1, \dots, k$ and $M$	{acc. to (6.2.23), (6.2.24)}
Compute $B'_{\pm}$	{acc. to Algorithm 5}
Compute $B_d$	{acc. to (6.2.27)}
Update $A_j$ , $j = 1 : k$ using $B_d$	{acc. to (6.2.29)}
Compute $B_{\pm} := -\text{diag}(a_1, \dots, a_k)_{\pm} [C_{n-1}, C_{n-2}, \dots, C_0 - \text{diag}(C_0)]_{\pm}$	{acc. to (6.2.28)}
$B_{\pm} \leftarrow [0, B_{\pm}(:, 1 : nk - 1)]$	
$B_{\pm} := [B_{n-1}, B_{n-2}, \dots, B_0] \leftarrow B_{\pm} - B'_{\pm}$	

---

The structure and the numerical details of the resulting propagation algorithm is finally described in Algorithm 7.

---

**Algorithm 7** Propagation algorithm

---

$g := u_0$	{set the initial values}	
$g_- := (g(1))$ , $g_+ := (g(N))$	{save the boundary values}	
<b>for</b> $i = 1$ to $n$ <b>do</b> {propagate $n$ steps}		
<b>for</b> $j = 1$ to $k$ <b>do</b> {solve $k$ intermediate problems}		
$b = (M + A'_j) g$		
$c_- = B_-(j, (n-i)k + 1 : (n-1)k + j)g_-$		
$c_+ = B_+(j, (n-i)k + 1 : (n-1)k + j)g_+$		
$b \leftarrow b + \begin{pmatrix} c_- \\ 0 \\ c_+ \end{pmatrix}$		
Compute $g$ from $(M + A_j) g = b$		
$g_- \leftarrow (g_-, g(1))^T$ , $g_+ \leftarrow (g_+, g(N))^T$		{save the boundary values}
<b>end for</b>		
$u_i := g$		{solution of the $i$ -th step}
<b>end for</b>		

---

### 6.3. Stability

The stability of the wide-angle propagation Algorithm 7 can be verified through a natural extension of our earlier analysis for the Schrödinger-type Padé-(2,0) approximant [85]. We summarize the result in

**THEOREM 6.3.1.** *Let the finite element space  $V_h$  used for the discretization (6.2.23) be independent of the step number. Let the coefficients  $a'$  and  $a$  of the discrete evolution equation (6.2.10) are complex conjugates with  $\text{Im}(a) \neq 0$ . Then the wide angle propagation Algorithm 7 together with the factorization Algorithms 3 and 4 is stable for each stepsize  $\Delta z$ . In particular, both  $L^2(\mathbb{R})$  and  $H^1(\mathbb{R})$  norms of the solution are conserved along the propagation.*

PROOF. We start from the representation of the exterior solution given by (6.2.17)

$$\widehat{g}_i(p) = (p\mathbf{I} + C(s))^{-1} g_{i,0}, \quad \text{with } C(s) = C_0 + C_1s + \dots + C_{i-1}s^{i-1},$$

in which by construction (Algorithm 3)  $\text{Re}(\text{diag}(C_0)) > 0$ . The operator  $p\mathbf{I} + C(s)$  is of convolution type. Its inverse exists and is again of convolution type  $(p\mathbf{I} + C(s))^{-1} = D_0 + D_1s + \dots$ , as long  $\det(p\mathbf{I} + C_0) \neq 0$ . Further it holds  $C_0^{-1} = D_0$ . Hence the zeros of  $\det(p\mathbf{I} + C_0)$  are exactly the poles of the solution vector  $\widehat{g}_i(p)$ . Consequently,  $g_i(x)$  consists only of exponentially decaying functions, in particular  $\lim_{x \rightarrow \infty} g_i(x) = 0$ .

Let  $\Omega_- = (-\infty, x_-)$ ,  $\Omega_+ = (x_+, \infty)$ , and  $\Omega = (x_-, x_+)$ . The unconditional stability of the propagation algorithm follows directly from this conservation law.

We consider first the discrete variational equation, (6.2.23), and abbreviate  $g_{h,i}^{(1)}$  by  $g$  and  $u_{h,i-1}$  by  $u$  and  $a_1, a'_1$  by  $a, a'$ . Accordingly, (6.2.23) reads

$$(v_h, g) + (\partial_x v_h, a \partial_x g) - (\bar{v}_h a \partial_x g)|_{x_-}^{x_+} = (v_h, u) + (\partial_x v_h, a' \partial_x u) - (\bar{v}_h a' \partial_x u)|_{x_-}^{x_+}.$$

We must compute the function  $g$  from its predecessor  $u$ . If we regard the special choice  $v_h = g$  and  $v_h = u$ , we obtain the equation system

$$\begin{aligned} (g, g) + (\partial_x g, a \partial_x g) - (\bar{g} a \partial_x g)|_{x_-}^{x_+} &= (g, u) + (\partial_x g, a' \partial_x u) - (\bar{g} a' \partial_x u)|_{x_-}^{x_+} \\ (u, g) + (\partial_x u, a \partial_x g) - (\bar{u} a \partial_x g)|_{x_-}^{x_+} &= (u, u) + (\partial_x u, a' \partial_x u) - (\bar{u} a' \partial_x u)|_{x_-}^{x_+}. \end{aligned}$$

We compute the conjugate complex of the second equation

$$(g, u) + (a \partial_x g, \partial_x u) - (u \overline{a \partial_x g})|_{x_-}^{x_+} = (u, u) + (a' \partial_x u, \partial_x u) - (u \overline{a' \partial_x u})|_{x_-}^{x_+}.$$

The sum of the last equation and the first of the above system yields

$$\begin{aligned} (g, g) + (\partial_x g, a \partial_x g) - (\bar{g} a \partial_x g + u \overline{a \partial_x g})|_{x_-}^{x_+} \\ = (u, u) + (a' \partial_x u, \partial_x u) - (\bar{g} a' \partial_x u + u \overline{a' \partial_x u})|_{x_-}^{x_+}. \end{aligned}$$

The same procedure applies to the two exterior domains. For the right exterior domain

$$\begin{aligned} (g, g)_+ + (\partial_x g, a \partial_x g)_+ - (\bar{g} a \partial_x g + u \overline{a \partial_x g})|_{x_+}^\infty \\ = (u, u)_+ + (a' \partial_x u, \partial_x u)_+ - (\bar{g} a' \partial_x u + u \overline{a' \partial_x u})|_{x_+}^\infty \end{aligned}$$

with an analogous result for the left exterior domain. Summing all three contributions and noting that the terms at infinity vanish while the normal derivatives at  $x_\pm$  cancel we obtain

$$\begin{aligned} (g, g)_- + (g, g) + (g, g)_+ + (\partial_x g, a \partial_x g)_- + (\partial_x g, a \partial_x g) + (\partial_x g, a \partial_x g)_+ \\ = (u, u)_- + (u, u) + (u, u)_+ + (a' \partial_x u, \partial_x u)_- + (a' \partial_x u, \partial_x u) + (a' \partial_x u, \partial_x u)_+ \end{aligned}$$

The imaginary part of the above equation yields

$$\begin{aligned} \text{Im}(a) [(\partial_x g, \partial_x g)_- + (\partial_x g, \partial_x g) + (\partial_x g, \partial_x g)_+] \\ = -\text{Im}(a') [(\partial_x u, \partial_x u)_- + (\partial_x u, \partial_x u) + (\partial_x u, \partial_x u)_+] \end{aligned}$$

Since  $a$  and  $a'$  are complex conjugates with  $\Im(a) \neq 0$ , we obtain the conservation law

$$\begin{aligned} (\partial_x g, \partial_x g)_- + (\partial_x g, \partial_x g) + (\partial_x g, \partial_x g)_+ \\ = (\partial_x u, \partial_x u)_- + (\partial_x u, \partial_x u) + (\partial_x u, \partial_x u)_+. \end{aligned}$$

This means that the  $H^1(\mathbb{R})$  semi-norm of two subsequent steps is conserved and implies the conservation of the  $L^2(\mathbb{R})$  norm

$$(g, g)_- + (g, g) + (g, g)_+ = (u, u)_- + (u, u) + (u, u)_+.$$

Since the  $L^2(\mathbb{R})$  norm is conserved for every intermediate step, it is conserved over the entire discrete evolution. Since the initial data  $u_0$  is required to be compactly supported on the interior domain, we obtain the bound

$$\|g_i\|_{L^2(\Omega), i \geq 0} \leq \|u_0\|_{L^2(\Omega)}$$

□

#### 6.4. Numerical Experiments

We now verify our theoretical considerations by computing the reflection from the computational window boundary of a Gaussian input beam described by

$$u(x, 0) = u_0(x) = \text{const} \exp\left(-\left(x/10\right)^2\right) \exp(ik_0x \sin \phi)$$

propagating in air. We set  $k(x) = k_0 = 2\pi/\lambda$  where the free space wavelength  $\lambda = 1.55\mu\text{m}$ , the propagation step size  $\Delta z = 0.4\mu\text{m}$  and  $\phi = \pi/4$ . In our calculations which are meant to duplicate the corresponding numerical experiments in a predecessor to this work [32], the computational domains are either  $\Omega = (-50, 50) \times (0, 400)\mu\text{m}^2$  or  $\Omega = (-50, 50) \times (0, 100)\mu\text{m}^2$  while the transverse step sizes  $\Delta x$  vary between  $0.01\mu\text{m} \leq \Delta x \leq 0.2\mu\text{m}$ . We consider first the intrinsic error associated with applying the implicit mid-point rule to plane wave solutions of the Helmholtz equation. Recall that for the exact one-way Helmholtz propagator

$$u(\Delta z) = u(0) \exp\left(-i\Delta z k_0 \sqrt{1 - \sin^2 \phi}\right),$$

for which the implicit mid-point rule yields

$$u_{\text{IMR}}(\Delta z) = u_{\text{IMR}}(0) \frac{1 - i\Delta z k_0/2 \left(1 - \sqrt{1 - \sin^2 \phi}\right)}{1 + i\Delta z k_0/2 \left(1 - \sqrt{1 - \sin^2 \phi}\right)}.$$

In this expression  $\phi$ ,  $-\pi/2 < \phi < \pi/2$  is the angle between the propagation direction and the  $z$ -axis. This yields a resulting phase error  $\log(u(\Delta z)) - \log(u_{\text{IMR}}(\Delta z))$  which we compare in Fig. 6.4.1 to that obtained by instead applying the Padé (2,0) (Schrödinger-type) approximation to the exact propagator. While at Padé order 2 the phase error of the Padé approximation is far greater, the opposite is true for a Padé (8,8) approximant as evident from Fig. 6.4.2.

Next we display in Fig. 6.4.3 the spectral norms of the matrices  $B_i$  for the boundary conditions associated with both the (2,0) and (8,8) Padé approximants as a function of the number of propagation steps. Here the matrices  $B_i$  are defined and computed as in Algorithm 6. Both approximations decay asymptotically as  $\|B_i\|_2 = \text{const } i^{-3/2}$ , independent of the order of the approximation. Note that every second coefficient of the Padé (2,0) approximation vanishes.

We now propagate the field from  $z = 0$  to  $z = 400\mu\text{m}$  with the (8,8) Padé procedure. The transverse grid spacing is here  $\Delta x = 0.2\mu\text{m}$  while the computational domain is  $\Omega = (-50, 50) \times (0, 400)\mu\text{m}^2$ . The contour lines for the logarithmic amplitude over the first  $100\mu\text{m}$  of propagation are shown in Fig. 6.4.4. While the incident field propagates as expected along the  $\theta = \pi/4$ -direction, residual reflections are generated by the finite transverse discretization error.

To underline the wide-angle property of the Padé (8,8) approximant, we first note that employing the (2,0) in place of the (8,8) Padé approximation for the square root operator leads to considerable phase errors, as evident from Fig. 6.4.5.

Performing the numerical experiment in the presence of second Gaussian beam that describes an angle of  $+\pi/4$  with respect to the  $z$ -axis. Fig. 6.4.6 demonstrates the wide-angle nature of the underlying propagation method.

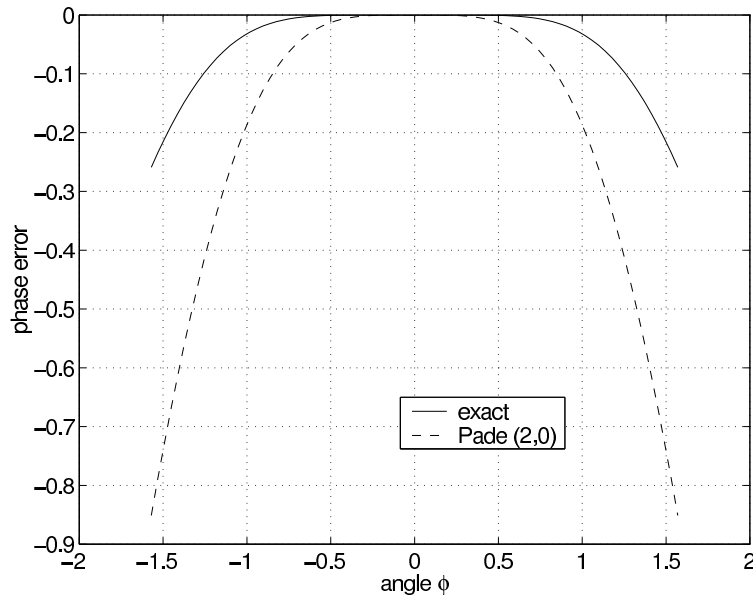


FIGURE 6.4.1. The phase error associated with the exact implicit mid-point discretization (solid line) compared to that of the corresponding Padé (2,0) approximant (dashed line)

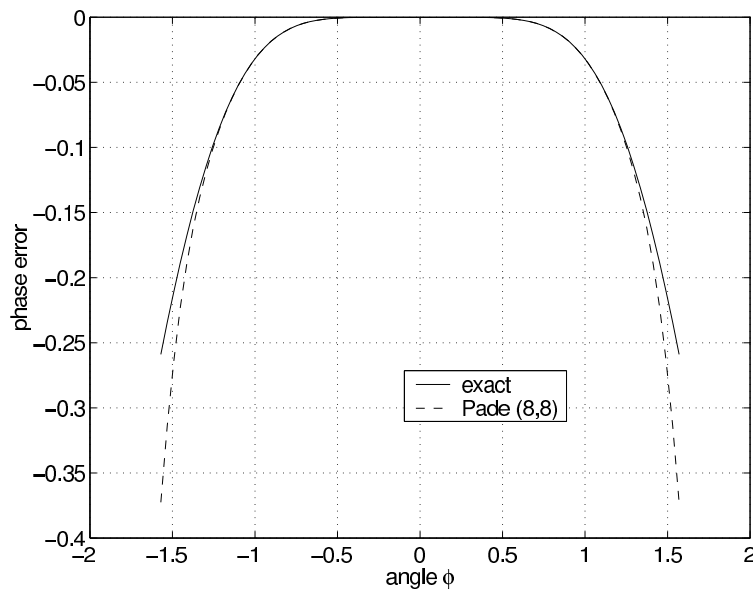


FIGURE 6.4.2. As in Fig. 6.4.1, but for a (8,8) Padé approximant.

The influence of the discretization error with respect to the transverse step length  $\Delta x$ , is evident if we repeat our previous (8,8) Padé simulation, Fig. 6.4.4 with  $\Delta x = 0.01\mu m$ , cf. Fig. 6.4.7. The boundary reflection, which vanishes in the  $\Delta x \rightarrow 0$  limit is indeed significantly reduced. The actual dependence of the reflection from the discretization error can be deduced from Fig. 6.4.8 which graphs the discrete  $L^2(-50, 50)$  norm as a function of propagation distance for transverse step sizes of  $0.2\mu m, 0.1\mu m, 0.05\mu m, 0.025\mu m$ , and  $0.01\mu m$ . The figure clearly shows

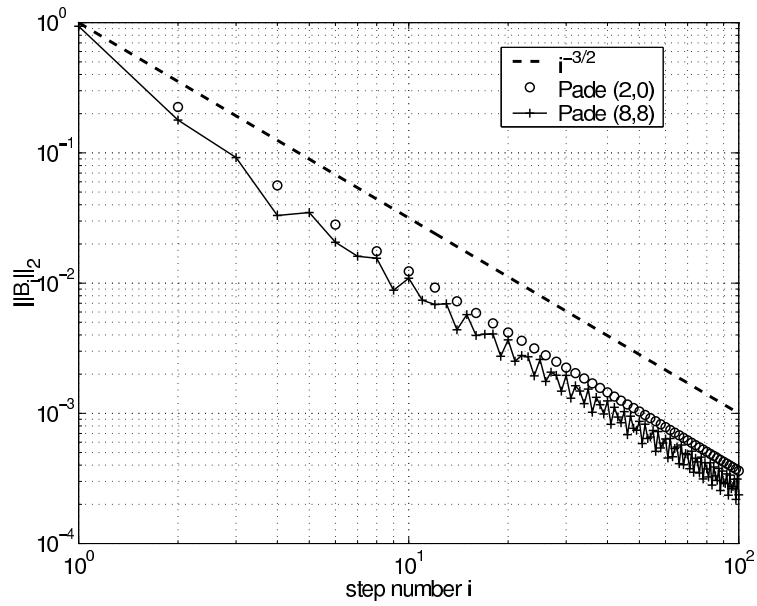


FIGURE 6.4.3. Spectral norm of the boundary matrices  $B_i$  as a function of the number of propagation steps.

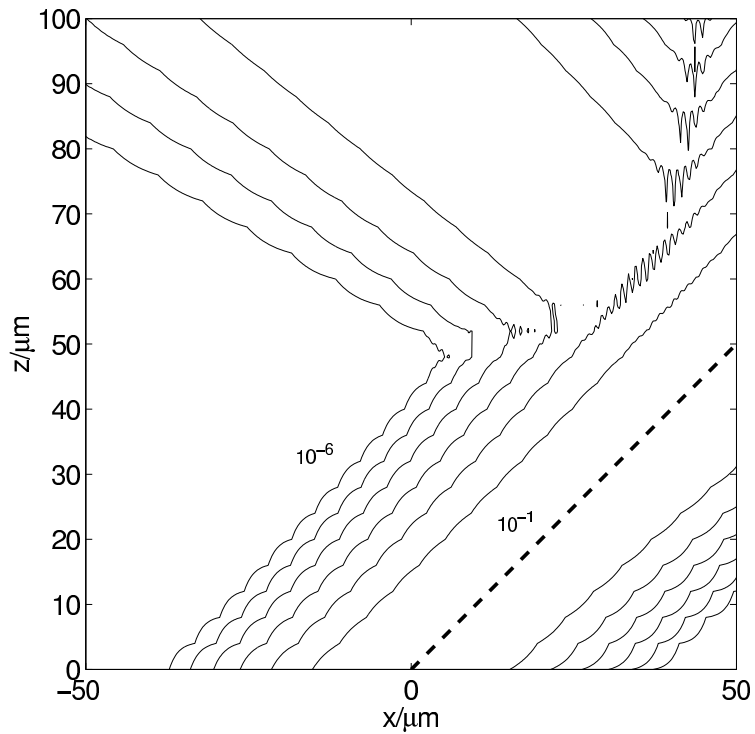


FIGURE 6.4.4. Gaussian beam propagation calculated with a (8,8) Padé propagator and  $\Delta x = 0.2\mu\text{m}$ . The dashed line represents the exact propagation angle  $\theta = \pi/4$ . The reflected field vanishes as  $\Delta x \rightarrow 0$ .

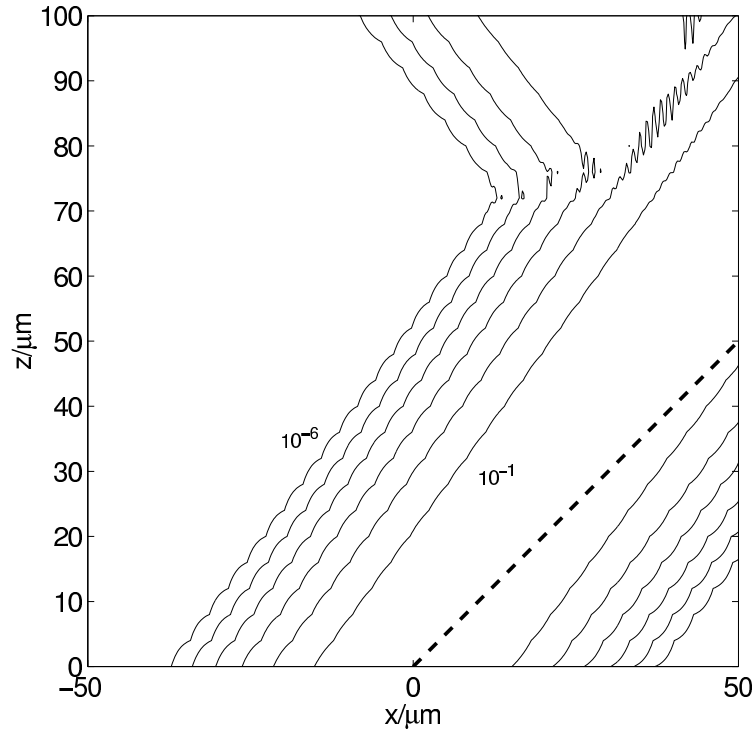


FIGURE 6.4.5. As in Fig. 6.4.4 except for a Padé (2,0) propagator. Note the error in the propagation angle.

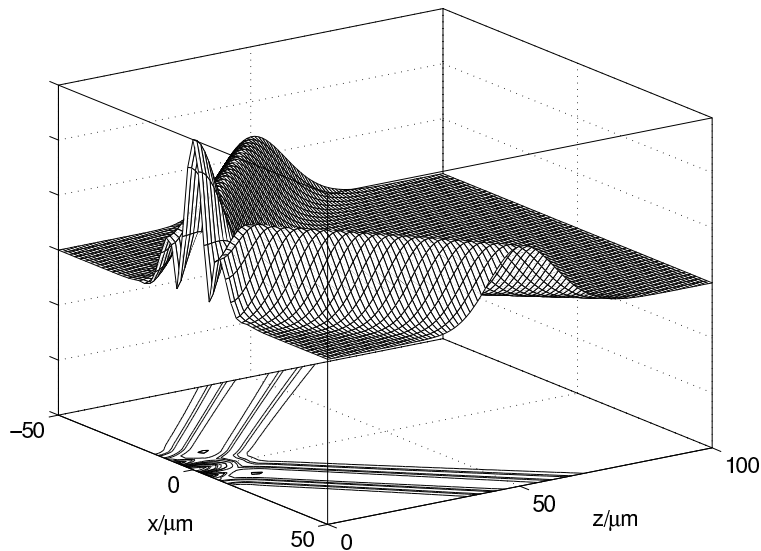


FIGURE 6.4.6. Propagation of two Gaussian beams at a relative angle of  $\pi/2$ .

that halving the transverse grid point spacing reduces the norm of the reflected field by a factor of 4. This behavior is entirely consistent with the  $O(\Delta x^2)$  discretization error of the underlying linear finite elements. Finally, to demonstrate the stability of our algorithm (subject to arithmetic error), we display in Fig. 6.4.9  $\|u\|$  computed

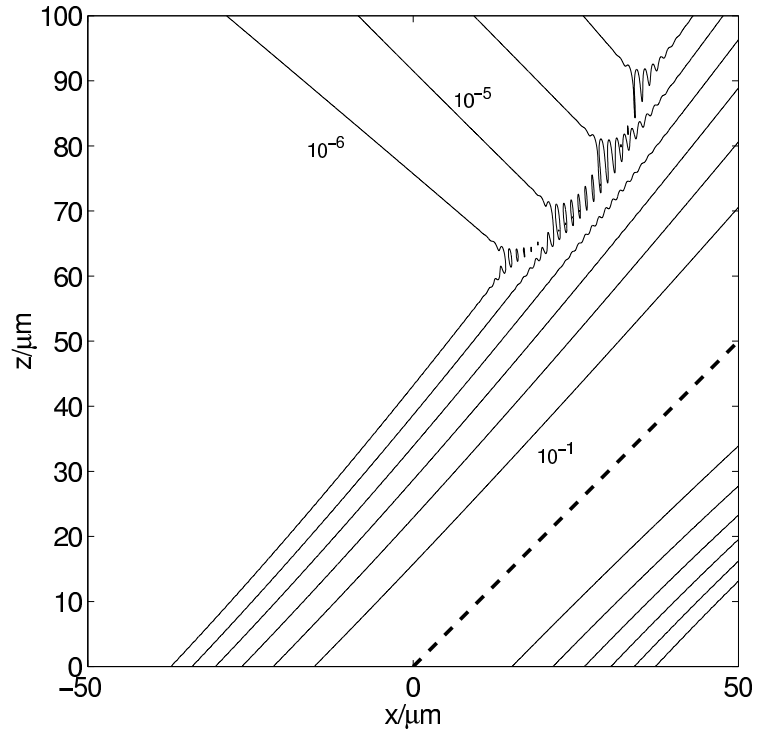


FIGURE 6.4.7. As in Fig. 6.4.4, but with  $\Delta x = 0.01 \mu m$ .

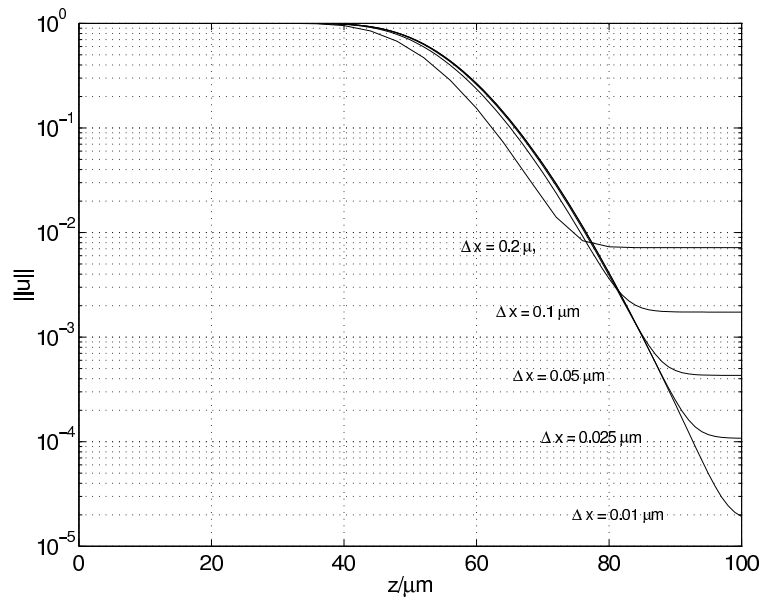


FIGURE 6.4.8. The  $L^2(-50, 50)$  norm  $\|u\|$  as a function of the propagation distance for varying step sizes  $\Delta x$ .



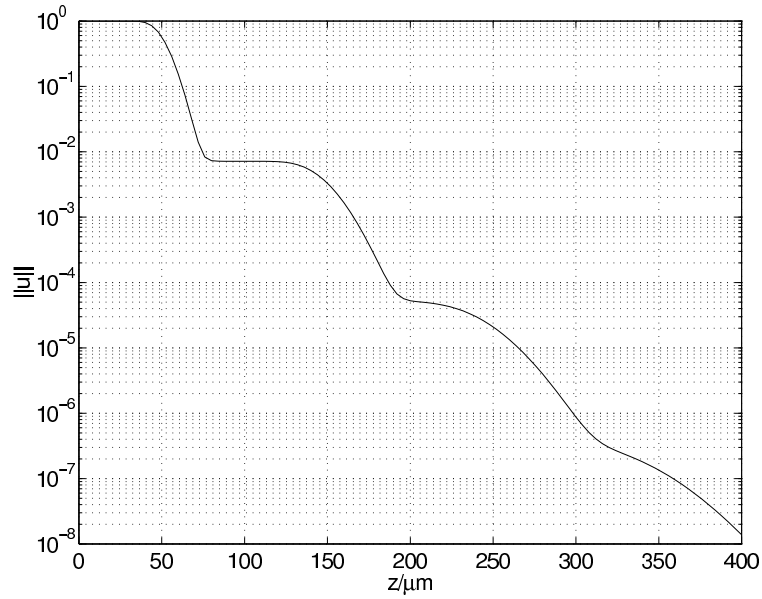


FIGURE 6.4.9. The discrete  $L_2$ -norm of the field within the computational window as a function of longitudinal distance for  $0 \leq z \leq 400 \mu m$  for the parameters of Fig. 6.4.4.

over a longer longitudinal interval  $0 \leq z \leq 400 \mu m$  for an (8,8) Padé propagator. The parameters are the same as in Fig. 6.4.4. Clearly the resulting curve, which displays successive plateaus corresponding to integer number of reflections of the Gaussian beam from the computational window boundaries indicates the stability of the algorithm.

## Conclusions

**Main Results.** The central result of this work is a new characterization of wave propagation on unbounded media. The most important notion is the concept of *pole condition*. All asymptotically outgoing waves have to satisfy the pole condition. The pole condition generalizes the famous Sommerfeld radiation condition to problems with inhomogeneous exterior domains and is equivalent to Sommerfeld's condition in the simpler case of homogeneous exterior domains. The pole condition is constructive. It allows both a theoretical analysis of wave propagation and a construction of numerical algorithms. As a side-effect of the approach, a new representation formula for the solution in the exterior domain has been established. There are different possibilities for numerical realizations of the theoretical concept. For problems in space dimensions  $d \geq 2$  we developed a new class of methods which we call *Laplace domain methods*. The Laplace domain methods solve the interior problem (given in the spatial domain) and the exterior problem (given in the *spectral* domain) simultaneously such that the pole condition is automatically satisfied.

The main focus of the numerical experiments is on the 2D Helmholtz equation. Nevertheless we applied the technique also to the time-dependent 1D Schrödinger equation and, in a generalization of this approach, to higher order approximations of the 1D one-way wide angle equations.

**Open Problems and Current Research.** Since our whole concept is new, many questions in theory, algorithms and applications are still open. We want to mention the most important ones.

Theory:

- In Section 4.2 we have shown the equivalence of the Sommerfeld radiation condition and the pole condition for the finite dimensional case. In a joint paper with Hohage and Zschiedrich [58], we extended this to the infinite dimensional case. The technique used there is a perturbation analysis of Riccati's equation, which is far-away from the concept of our algorithms. Is there a proof closer related to the numerical algorithms?
- In our analysis we excluded potentials of the form  $\mathcal{O}(1/|\mathbf{x}|)$  for  $|\mathbf{x}| \rightarrow \infty$ . The Laplace domain method in its axis integral form, cf. Section 5.2, however, is able to deal even with this type of potentials. How to extend the theory to include such potentials?
- Our theory is based throughout on separable problems which is a consequence of the applied modal analysis. In contrast, the numerical algorithms are able to deal with non-separable problems as well. Can we reformulate the theory to cover non-separable problems?

## Algorithms:

- The Laplace domain methods analyzed in Sections 5.2 and 5.2.2 for separable problems in 2D have to be implemented and analyzed for general non-separable problems in 2D and 3D.
- The Laplace domain method for time-dependent problems, cf. p. 150, has to be implemented and analyzed. The central question here is whether or not it is possible to obtain stability results as in Theorem 6.3.1 for the Laplace domain method. This would offer the possibility to develop a new class of propagation algorithms.
- A detailed analysis of the operation account and memory requirements has to be done.

## Applications:

- The discretization of the exterior domain is based on *linear* rays. To extend the applicability, curved rays might be considered also. Our consideration of hyperbolic-elliptic coordinates in Section 4, p. 59, can be seen as a first step in this direction.
- Further, it seems to be possible to extend the construction scheme to vectorial elements of Nédélec-type. This would allow to apply the pole condition idea to Maxwell problems. Such results will be contained in the PhD thesis of Zschiedrich [97].
- The concept can be applied in a natural way to eigenproblems on unbounded domains. A first successful study of 1D eigenproblems based on the pole condition approach will be reported in the diploma thesis of Geng [33]. The numerical results clearly demonstrate the superiority of the pole condition approach over the PML approach.

## Basic Properties of the Continuous and Discrete Transforms

### A.1. Basic Properties of the Laplace Transform

For convenience, we assemble in the following some basic properties of Laplace transform, as far as we need them for our further elaboration. Our representation follows mainly the one given by Dettman in [24]. We assume the following:

- (H1) The real-valued function  $u(\xi)$  is twice continuously differentiable in any closed finite interval  $[0, a]$ ,  $u(\xi) \in C^2[0, a]$ ,  $0 < a < \infty$ .
- (H2) The second derivative of  $u$  is exponentially bounded,  $|u''(\xi)| \leq C e^{d\xi}$ ,  $0 \leq \xi < \infty$ , for some real  $d$ .

An immediate consequence is that besides  $u''$  also  $u'$  and  $u$  are also exponentially bounded. It follows directly  $u, u' \leq C e^{b\xi}$  with  $b = \max(0, d)$ .

LEMMA A.1.1. *The Laplace transform  $L(u)$*

$$L(u) := \int_0^{\infty} e^{-s\xi} u(\xi) d\xi, \quad s \in \mathbb{C}$$

*is analytic in the half-plane  $\operatorname{Re}(s) > b$ .*

PROOF. The integral converges uniformly for  $\operatorname{Re}(s) > b$  since  $|u(\xi)| \leq C e^{b\xi}$ . Hence the derivative with respect to  $s$  can be computed under the integral sign, which shows the existence of the derivative.  $\square$

As usually, we denote the Laplace transform of  $u(\xi)$  by  $\hat{u}(s) := L(u(\xi))(s)$  and its derivative by  $\hat{u}'(s)$  or by  $\partial_s \hat{u}$ .

LEMMA A.1.2. *The Laplace transform of the expression  $\xi^n u(\xi)$  is analytic for  $\operatorname{Re}(s) > b$ . It holds*

$$L(\xi^n u) = (-1)^n \partial_s^n \hat{u}(s).$$

PROOF. Taking the first derivative in the half-plane  $\operatorname{Re}(s) > b$  supplies  $\hat{u}'(s) = \int_0^{\infty} \partial_s e^{-s\xi} u(\xi) d\xi$ . This is equivalent to  $L(-\xi u(\xi)) = \hat{u}'(s)$ . Repeating this  $n$ -times gives the result.  $\square$

LEMMA A.1.3. *Let  $u$  satisfy the hypothesis (H1) and (H2). Then it holds in the half-plane  $\operatorname{Re}(s) > b$*

$$\begin{aligned} L(u'(\xi)) &= s\hat{u}(s) - u(0) \\ L(u''(\xi)) &= s^2\hat{u}(s) - su(0) - u'(0) \end{aligned}$$

PROOF. The function  $u(\xi)$ ,  $u'(\xi)$ , and  $u''(\xi)$  are all exponentially bounded. Hence there Laplace transform exists in the half-plane  $\operatorname{Re}(s) > b$ . Integration by parts yields  $\int_0^{\infty} d\xi e^{-s\xi} u' = -u(0) + s \int_0^{\infty} d\xi e^{-s\xi} u = -u(0) + s\hat{u}(s)$ . Setting  $v'(\xi) = u''(\xi)$  and repeating the procedure for  $v'$  gives the second correspondence.  $\square$

COROLLARY A.1.4. *Let  $u$  satisfy the hypothesis (H1) and (H2). Then there is a constant  $M$  and a half-plane  $\operatorname{Re}(s) \geq \rho$  such that*

$$|\hat{u}(s)| \leq \frac{M}{|s|}.$$

PROOF. The exponential bound of  $\hat{u}'$  gives the estimate

$$|\hat{u}'| \leq \int_0^\infty d\xi C e^{b\xi} e^{-\operatorname{Re}(s)\xi} = \frac{C}{\operatorname{Re}(s) - b} \leq \frac{C}{\rho - b}$$

On the other hand, Lemma A.1.3 shows

$$\begin{aligned} |\hat{u}(s)| &\leq \frac{|\hat{u}'| + |u(0)|}{|s|} \\ &\leq \frac{\frac{C}{\rho - b} + |u(0)|}{|s|} \\ &= \frac{M}{|s|} \end{aligned}$$

where  $M = \max(C/(\rho - b), |u(0)|)$ . □

COROLLARY A.1.5. *Let  $u$  satisfy the hypothesis (H1) and (H2). It holds*

$$\lim_{\operatorname{Re}(s) \rightarrow \infty} s\hat{u}(s) = u(0).$$

This follows immediately from the fact that (H1) and (H2) it follows  $\lim_{\operatorname{Re}(s) \rightarrow \infty} |\hat{u}'(s)| = 0$ .

LEMMA A.1.6. *Let  $u$  satisfy the hypothesis (H1) and (H2). Let  $s_\infty$  be defined as a point in the complex plane with some finite imaginary part and an infinite, positive real part. It holds in a half-plane  $\operatorname{Re}(s) > b$*

$$\mathbb{L}\left(\frac{1}{\xi_0 + \xi} u(\xi)\right) = e^{s\xi_0} \int_s^{s_\infty} ds_1 e^{-s_1\xi_0} \hat{u}(s_1),$$

where the Laplace transformed expression remains analytic in this half-plane.

PROOF. Let

$$\hat{g}(s) = e^{s\xi_0} \int_s^{s_\infty} ds_1 e^{-s_1\xi_0} \hat{u}(s_1), \text{ in } \operatorname{Re}(s) > b$$

It follows

$$\begin{aligned} \hat{g}'(s) &= \xi_0 e^{s\xi_0} \int_s^{s_\infty} ds_1 e^{-s_1\xi_0} \hat{u}(s_1) - \hat{u}(s) \\ &= \xi_0 \hat{g}(s) - \hat{u}(s). \end{aligned}$$

By Lemma A.1.3,

$$u(\xi) = (\xi + \xi_0) g(x).$$

□

COROLLARY A.1.7. *It holds in a half-plane  $\operatorname{Re}(s) > b$*

$$\mathbb{L}\left(\frac{1}{(\xi_0 + \xi)^n} u(\xi)\right) = e^{s\xi_0} \int_s^{s_\infty} ds_n \dots \int_{s_2}^{s_\infty} ds_1 e^{-s_1\xi_0} \hat{u}(s_1),$$

where the Laplace transformed expression remains analytic in this half-plane.

The iterated integrals of Corollary A.1.7 can be generalized to the following (obviously not widely known) lemma, which supplies the Laplace transform of the product of two functions  $p\left(\frac{1}{\xi+\xi_0}\right)$  and  $u(\xi)$  via a convolution along the real axis.

LEMMA A.1.8. *Let  $p$  be a power series of the form  $p(t) = \sum_{m=1}^{\infty} p_m t^m$  with convergence radius  $\rho > \frac{1}{\xi_0}$ , further  $v(\xi) = p\left(\frac{1}{\xi}\right)$ , and let  $u$  be continuous and bounded on the positive real axis. Then*

$$(A.1.1) \quad \mathbf{L}(v(\xi + \xi_0)u(\xi))(s) = \int_s^{\infty} ds_1 e^{(s-s_1)\xi_0} (L^{-1}v)(s_1 - s) \widehat{u}(s_1)$$

in the half-plane  $\operatorname{Re}(s) > 0$ . If the convergence radius of  $p(t)$  even  $\rho > \frac{1}{\xi_0 - \tilde{\xi}_0} > \frac{1}{\xi_0}$  for some  $0 < \tilde{\xi}_0 < \xi_0$ , then there exists a constant  $C > 0$  such that

$$(A.1.2) \quad |e^{-s\xi_0} (L^{-1}v)(s)| \leq C e^{-\xi_0 \operatorname{Re} s + (\xi_0 - \tilde{\xi}_0)|s|}$$

$$(A.1.3) \quad \left| \frac{d}{ds} (e^{-s\xi_0} (L^{-1}v)(s)) \right| \leq C e^{-\xi_0 \operatorname{Re} s + (\xi_0 - \tilde{\xi}_0)|s|}$$

for all  $s \in \mathbb{C}$ .

PROOF. It holds

$$(L^{-1}v)(s) = \sum_{m=1}^{\infty} \frac{p_m}{(m-1)!} s^{m-1}.$$

For  $m = 1$  Eq. (A.1.1) recovers Lemma A.1.6. We want to show the result for  $m > 1$  by induction. Assume that (A.1.1) holds true for  $p = t^m$ . We show that it holds also for  $p = t^{m+1}$ . By Lemma A.1.6

$$\begin{aligned} & \mathbf{L}\left(\left(\frac{1}{\xi_0 + \xi}\right)^{m+1} u(\xi)\right)(s) \\ &= \mathbf{L}\left(\frac{1}{\xi_0 + \xi} \frac{u(\xi)}{(\xi_0 + \xi)^m}\right)(s) \\ &= \int_s^{\infty} ds_2 e^{(s-s_2)\xi_0} \left( \int_{s_2}^{\infty} ds_1 e^{(s_2-s_1)\xi_0} \mathbf{L}^{-1}\left(\left(\frac{1}{\xi}\right)^m\right)(s_1 - s_2) \widehat{u}(s_1) \right) \\ &= \int_s^{\infty} ds_2 \int_{s_2}^{\infty} ds_1 e^{(s-s_1)\xi_0} \frac{1}{(m-1)!} (s_1 - s_2)^{m-1} \widehat{u}(s_1) \\ &= \int_s^{\infty} ds_1 \int_s^{s_1} ds_2 e^{(s-s_1)\xi_0} \frac{1}{(m-1)!} (s_1 - s_2)^{m-1} \widehat{u}(s_1) \\ &= \int_s^{\infty} ds_1 e^{(s-s_1)\xi_0} \frac{1}{m!} (s_1 - s)^m \widehat{u}(s_1). \end{aligned}$$

So far the statement holds if  $p$  is a polynomial. We extend this to the case that  $p$  is an infinite series. From the definition of the convergence radius follows the exponential bound

$$\begin{aligned}
|(\mathbb{L}^{-1}v)(s)| &\leq \sum_{m=0}^{\infty} \frac{|p_{m+1}|}{m!} |s^m| \\
&= \sum_{m=0}^{\infty} \frac{|p_{m+1}|}{m!} \frac{(\xi_0 - \tilde{\xi}_0)^m}{(\xi_0 - \tilde{\xi}_0)^m} |s^m| \\
&\leq C \sum_{m=0}^{\infty} \frac{(\xi_0 - \tilde{\xi}_0)^m}{m!} |s^m| \\
&\leq C e^{(\xi_0 - \tilde{\xi}_0)|s|}
\end{aligned}$$

with  $C = |\xi_0 - \tilde{\xi}_0| \sup_{m \geq 0} |p_{m+1}| (\xi_0 - \tilde{\xi}_0)^{-(m+1)} < \infty$ , because all partial sums above are bounded by the right-hand side given by the exponential factor  $\exp(\xi_0 - \tilde{\xi}_0)|s|$  and because the power series converges uniformly for  $|t| < 1/\xi_0 < 1/(\xi_0 - \tilde{\xi}_0)$ . The latter fact makes it necessary that  $\lim_{m \rightarrow \infty} |p_{m+1}| (\xi_0 - \tilde{\xi}_0)^{-(m+1)} = 0$ , hence  $C < \infty$ . Obviously, we can repeat the same conclusion for

$$|(\mathbb{L}^{-1}v)'(s)| = \left| \sum_{m=0}^{\infty} \frac{|p_{m+1}|}{(m-1)!} |s^{m-1}| \right|$$

resulting in

$$|(\mathbb{L}^{-1}v)'(s)| \leq C e^{(\xi_0 - \tilde{\xi}_0)|s|}$$

with  $C = |\xi_0 - \tilde{\xi}_0|^2 \sup_{m \geq 0} |p_{m+1}| (\xi_0 - \tilde{\xi}_0)^{-(m+1)} < \infty$

It can be shown by Lebesgue's Dominated Convergence Theorem that

$$\begin{aligned}
\int_0^{\infty} d\xi e^{-s\xi} v(\xi + \xi_0) u(\xi) &= \lim_{M \rightarrow \infty} \int_0^{\infty} d\xi e^{-s\xi} \sum_{m=1}^M \frac{p_m}{(\xi + \xi_0)^m} u(\xi) \\
&= \lim_{M \rightarrow \infty} \int_s^{\infty} ds_1 e^{(s-s_1)\xi_0} \sum_{m=1}^M \frac{p_m}{(m-1)!} (s-s_1)^{m-1} \hat{u}(s_1) \\
&= \int_s^{\infty} ds_1 e^{(s-s_1)\xi_0} (\mathbb{L}^{-1}v)(s_1 - s) \hat{u}(s_1)
\end{aligned}$$

□

LEMMA A.1.9. *Let  $u : \mathbb{R}_+ \rightarrow \mathbb{C}$  be bounded. Further let  $u$  have a bounded analytic continuation such that  $u(r)$  is holomorphic for all complex  $r$  with  $\operatorname{Re} r \geq 0$  and  $0 \leq \arg r < \pi/2$ . Let  $\alpha \neq 0$  be a complex number with  $0 \leq \arg \alpha < \pi/2$ . Then it holds*

$$L(u(\alpha\xi))(s) = \frac{1}{\alpha} (Lu(\xi))\left(\frac{s}{\alpha}\right) \quad \text{for } \operatorname{Re} s > 0, \quad \xi \in \mathbb{R}_+.$$

PROOF. By definition, it holds

$$L(u(\alpha\xi))(s) = \int_0^{\infty} e^{-s\xi} u(\alpha\xi) d\xi, \quad \operatorname{Re} s > 0.$$

This is identical to

$$L(u(\alpha\xi))(s) = \int_0^{\infty} \frac{1}{\alpha} e^{-\frac{s}{\alpha}\alpha\xi} u(\alpha\xi) d\alpha\xi, \quad \operatorname{Re} s > 0.$$

Setting  $\zeta := \alpha\xi$  and applying Cauchy's Theorem, we obtain, with  $P_\infty := \lim_{\xi \rightarrow \infty} \alpha\xi$ ,

$$\begin{aligned} \int_0^{P_\infty} \frac{1}{\alpha} e^{-\frac{s}{\alpha}\zeta} u(\zeta) d\zeta &= \int_0^\infty \frac{1}{\alpha} e^{-\frac{s}{\alpha}\xi} u(\xi) d\xi \\ &= \frac{1}{\alpha} L(u(\xi)) \left( \frac{s}{\alpha} \right). \end{aligned}$$

□

## A.2. Discrete Operational Calculus

This is the discrete analog to Mikusiński's operational calculus [74]. Our goal is to define an algebra with infinite sequences as elements such that operations like quotients and square roots of sequences and are well defined. Let  $U$  denote the infinite sequence of complex numbers  $\{u_0, u_1, \dots, u_j, \dots\}$ ,  $|u_j| < \infty$  for all  $j \geq 0$ . Let  $\mathcal{C}$  be the set of all such sequences. In particular, if the  $u_j$  are functions  $u_j : \mathbb{C} \rightarrow \mathbb{C}$  of the complex argument  $p$ , we write  $U(p) := \{u_0(p), u_1(p), \dots, u_j(p), \dots\}$ .

DEFINITION A.2.1. We define the Null element and the unity element of  $\mathcal{C}$  by

$$[0] = \{0, 0, \dots\} \quad \text{and} \quad [1] = \{1, 0, 0, \dots\}.$$

For simplicity we write often 0 instead of  $[0]$ , and 1 instead of  $[1]$ , respectively, if no ambiguity occurs. We define the operations addition and multiplication between two elements  $U, V \in \mathcal{C}$  by

$$U + V = \{u_0 + v_0, u_1 + v_1, \dots\} \quad \text{and} \quad UV = \{g_0, g_1, \dots\} \quad \text{with} \quad g_j = \sum_{i=0}^j u_{j-i} v_i.$$

From these definitions we find for any  $U, V \in \mathcal{C}$  that  $U + V \in \mathcal{C}$  and  $UV \in \mathcal{C}$ , and further  $UV = VU$ ,  $U(VW) = (UV)W$ , and  $U(V + W) = UV + UW$ . Hence,  $\mathcal{C}$  is a commutative ring. It follows by induction that for  $U, V \in \mathcal{C}$  the equation  $UV = 0$  implies either  $U = 0$  or  $V = 0$  or both. Hence we can extend the ring  $\mathcal{C}$  to the ring of fractions  $\mathcal{C}_q = \{U/V \mid U, V \in \mathcal{C}, V \neq 0\}$ . To fractions  $U/V$  and  $U'/V'$  are defined to be equal if and only if  $UV' = VU'$ . The addition and multiplication of two fractions of sequences are defined by

$$\frac{U}{V} + \frac{U'}{V'} = \frac{UV' + VU'}{VV'} \quad \text{and} \quad \frac{U}{V} \frac{U'}{V'} = \frac{UU'}{VV'},$$

respectively. Some special elements of  $\mathcal{C}_q$  are labeled with own symbols. We define the multiplication operator by  $[a] := \{a, 0, 0, \dots\}$ ,  $a \in \mathbb{C}$ . The operation  $[a]U$  realizes a multiplication of each element of  $U$  by  $a$ , for simplicity we write  $aU$ .

DEFINITION A.2.2. The shift-operator  $s$  is defined by  $s := \{0, 1, 0, 0, \dots\}$ .

By definition, it holds  $sU = \{0, u_0, u_1, \dots\}$  and  $ss = s^2 = \{0, 0, 1, 0, \dots\}$ , etc. Two sequences  $U, V \in \mathcal{C}$  are equal, if their elements are equal,  $u_i = v_i$  for all  $i \geq 0$ . In particular, if we compare sequences elementwise, we use the following notation:  $(sU)_{i+1} = (U)_i = u_i$ . Even shorter we write  $su_{i+1} = u_i$  which has to be understood in the foregoing sense. The following property holds:  $V := s^k U$  implies  $v_j = 0$  for  $j = 0, \dots, k-1$ . Further, we say that  $U(p)$  has a complex pole  $\lambda$  if  $\lim_{p \rightarrow \lambda} |u_j(p)| = \infty$  for some  $j \geq 0$ .

DEFINITION A.2.3. A linear mapping  $A : \mathcal{C} \rightarrow \mathcal{C}$  is defined through infinite matrices  $A = (a_{ij})_{i,j=0,\dots,\infty}$ ,  $a_{ij} \in \mathbb{C}$  which map  $U$  into  $V$ ,  $U \mapsto V = AU$  by a matrix-vector multiplication  $V_i = \sum_{j=0}^i a_{ij} U_j$ . We denote the set of all of such matrices by  $\mathcal{L}$ .



In our application equations between sequences appear of the form

$$(A.2.1) \quad \left[ \left( \begin{array}{cccc} A & & & \\ & A & & \\ & & A & \\ & & & \ddots \end{array} \right) + \left( \begin{array}{cccc} B & & & \\ & B & & \\ & & B & \\ & & & \ddots \end{array} \right) s^k \right] U = V,$$

where  $A$  and  $B$  are lower triangular, complex  $k \times k$  matrices. Each row of this equation corresponds to one propagation step (we call this a macro-step), which itself consists of  $k$  intermediate steps, according to the order of the Padé approximation. In view of this separation of matrices, we separate the sequences  $U$  and  $V$  into subsequences of length  $k$ , and motivated by the matrix form of the basic equation we denote this subsequences column vectors, that is we compose  $U = \{u_0^T, u_1^T, \dots\}$  and  $V = \{v_0^T, v_1^T, \dots\}$  with  $u_j, v_j \in \mathbb{C}^k$ . Thus we can replace (A.2.1) by the shorter form

$$(A.2.2) \quad (A + Bs^k) u_j = v_j, \quad j \geq 0.$$

The second important type of matrices which occur are lower triangular block matrices of *convolution* type

$$(A.2.3) \quad \left( \begin{array}{cccc} B_0 & & & \\ B_1 & B_0 & & \\ B_2 & B_1 & B_0 & \\ \vdots & \vdots & \vdots & \ddots \end{array} \right)$$

where the  $B_j$ ,  $j \geq 0$ , are  $k \times k$  complex matrices. These matrices are special forms of infinite Töplitz matrices, cf. e.g. [14]. In the special case  $k = 1$  we can assign to each matrix  $B$  of convolution type a sequence  $B \in \mathcal{C}$  defined by  $B = \{B_0, B_1, B_2, \dots\}$  such that  $BU = BU$  for all  $U \in \mathcal{C}$ . A typical mapping  $U \mapsto V$  within our factorization procedure possesses the form

$$(A.2.4) \quad \left[ \left( \begin{array}{cccc} A & & & \\ & A & & \\ & & A & \\ & & & \ddots \end{array} \right) + \left( \begin{array}{cccc} B_0 & & & \\ B_1 & B_0 & & \\ B_2 & B_1 & B_0 & \\ \vdots & \vdots & \vdots & \ddots \end{array} \right) \right] U = V,$$

where  $A$  and  $B_0$  are lower triangular, complex  $k \times k$  matrices and  $B_j$ ,  $j \geq 1$  are arbitrary, complex  $k \times k$  matrices. By our notation introduced above, the  $j$ th row ( $j \geq 0$ ) reads

$$(A.2.5) \quad \left( A + \sum_{i=0}^j B_i s^{ik} \right) u_j = v_j, \quad j \geq 0,$$

which we further abbreviate by  $(A + B(s))u_j = v_j$  with an operator  $B(s) := \sum_{i=0}^j B_i s^{ik}$ . We denote equations of type A.2.5 as operator equations, because they are formulated using the shift operator. We say that the inverse of the operator appearing in brackets on the left-hand side of A.2.4 possesses a pole, if  $\det(A - B_0) = 0$ . Due to the equivalence of A.2.4 and the operator equation A.2.5 we say accordingly, that the operator equation A.2.2 possesses a pole if  $\det(A - B_0) = 0$ .

**PROPOSITION A.2.4.** *The product of two infinite matrices of convolution type is again of convolution type.*

PROOF. We consider the mapping  $U \mapsto V$

$$(A.2.6) \quad \left[ \left( \begin{array}{cccc} A_0 & & & \\ A_1 & A_0 & & \\ A_2 & A_1 & A_0 & \\ \vdots & \vdots & \vdots & \ddots \end{array} \right) \left( \begin{array}{cccc} B_0 & & & \\ B_1 & B_0 & & \\ B_2 & B_1 & B_0 & \\ \vdots & \vdots & \vdots & \ddots \end{array} \right) \right] U = V,$$

with matrices  $B_j$  as above and  $A_j$  of the same type as  $B_j$ . The  $j$ th row of (A.2.6) reads

$$(A.2.7) \quad \left[ \sum_{i=0}^j \underbrace{\left( \sum_{l=0}^i A_l B_{i-l} \right)}_{C_i} s^{ik} \right] u_j = v_j, \quad j \geq 0,$$

and the matrices  $C_i$  represent the resulting matrix.  $\square$

The expression in brackets on the left-hand side of (A.2.7) can be identified with the formal product

$$\left( \sum_{i=0}^j A_i s^{ik} \right) \left( \sum_{i=0}^j B_i s^{ik} \right) = A_0 B_0 + \sum_{l+m=1} A_l B_m s^k + \dots + \sum_{l+m=j} A_l B_m s^{jk},$$

$j, l, m \geq 0.$

Thus given an operator expression in polynomial form  $\sum_{i=0}^j C_i s^{ki}$  and matrices  $B_i$ ,  $i = 0, \dots, j$ , we can determine the coefficients  $A_i$ ,  $i = 0, \dots, j$ , of a factorization  $\left( \sum_{i=0}^j A_i s^{ik} \right) \left( \sum_{i=0}^j B_i s^{ik} \right) u_j = \sum_{i=0}^j C_i s^{ki} u_j$  by a comparison of coefficients. In the following we consider matrices of convolution type  $B(s) = \sum_{i=0}^j B_i s^{ik}$ ,  $B_i \in \mathbb{C}^{k \times k}$  with invertible matrices  $B_0$ . In particular, we consider the special choice  $B_0 = I$ . The general case is easily traced back to this case.

LEMMA A.2.5. *Given an operator  $B(s)$  of convolution type,*

$$B(s) = I + \sum_{i=1}^j B_i s^{ik}, \quad B_i \in \mathbb{C}^{k \times k}.$$

*It holds:*

- (i) *the inverse  $B^{-1}(s)$  of the operator  $B(s)$  exists and  $B^{-1}(s) \in \mathcal{L}$ .*
- (ii)  *$B^{-1}(s)$  possesses the convolution form (A.2.3). The operator  $C(s) := B^{-1}(s)$  has a representation*

$$(A.2.8) \quad C(s) = I + \sum_{i=1}^j C_i s^{ik}, \quad C_i \in \mathbb{C}^{k \times k},$$

*where the  $C_i$  are recursively determined by  $C_i = -\sum_{l=0}^{i-1} B_{i-l} C_l$ ,  $C_0 = I$ .*

PROOF. (i) Gaussian elimination.

(ii) Proof by induction. Given  $U \in \mathcal{C}$ , we set  $v_j := \left( I + \sum_{i=1}^j B_i s^{ik} \right) u_j$ . We want to show that there exists an operator  $C(s)$  with  $u_j = \left( I + \sum_{i=1}^j C_i s^{ik} \right) v_j$ . Let  $j = 0$  and  $j = 1$ . It holds, by definition of  $B(s)$ ,  $u_0 = v_0$  and  $B_1 u_0 + u_1 = v_1$ , hence

$$\begin{aligned} u_0 &= v_0 \\ u_1 &= -B_1 v_0 + v_1. \end{aligned}$$

It follows  $C_1 = -B_1$ . Assume that the representation formula (A.2.8) is true for all indices  $l$  up to some  $j \geq 1$ . Hence it holds

$$(A.2.9) \quad \begin{aligned} u_l &= \sum_{i=0}^l C_i v_{l-i}, \quad l = 0, \dots, j \text{ and } C_0 := I \\ u_j &= v_j - \sum_{i=1}^j B_i u_{j-i}, \end{aligned}$$

where the last equation is just the definition of  $B(s)$ . It follows, expressing  $u_l$  by the appropriate representation (A.2.9),

$$(A.2.10) \quad \begin{aligned} u_j &= v_j - \sum_{i=1}^j B_i u_{j-i} \\ &= v_j - \sum_{i=1}^j \sum_{m=0}^{j-i} B_i C_m v_{j-i-m} \end{aligned}$$

$$(A.2.11) \quad \begin{aligned} &= v_j - \sum_{k=j-1}^0 \sum_{m=0}^k B_{j-k} C_m v_{k-m} \end{aligned}$$

$$(A.2.12) \quad \begin{aligned} &= v_j - \sum_{k=0}^{j-1} \sum_{m=0}^k B_{j-k} C_m v_{k-m} \\ &= v_j - \sum_{k=0}^{j-1} \sum_{l=k}^0 B_{j-k} C_{k-l} v_l \end{aligned}$$

$$(A.2.12) \quad \begin{aligned} &= v_j - \sum_{k=0}^{j-1} \sum_{l=0}^k B_{j-k} C_{k-l} v_l \\ &= v_j - \sum_{l=0}^{j-1} \underbrace{\left( \sum_{k=l}^{j-1} B_{j-k} C_{k-l} \right)}_{-C_{j-l}} v_l \end{aligned}$$

In (A.2.10) we replaced  $j - i$  by  $k$  and eliminated  $i$ , in (A.2.11) we replaced  $k - m$  by  $l$  and eliminated  $m$ . In (A.2.12) we changed the order of summation. It follows  $C_i = -\sum_{l=0}^{i-1} B_{i-l} C_l$  for  $i \leq j$ . This gives in recursive manner  $C_1 = -B_1$ ,  $C_2 = -B_2 - B_1 C_1$ , etc. Note, that we used the induction assumption (A.2.9) only up to the index  $l = j - 1$ . By assumption, this representation holds true for  $l = j$ . In what follows, we show that (A.2.9) remains true for  $l = j + 1$ . By definition

$$u_{j+1} = v_{j+1} - \sum_{i=1}^{j+1} B_i u_{j+1-i}.$$

Hence we can repeat the procedure setting  $j := j + 1$  which yields

$$u_{j+1} = v_{j+1} - \sum_{l=0}^j \underbrace{\left( \sum_{k=l}^j B_{j+1-k} C_{k-l} \right)}_{-C_{j+1-l}} v_l$$

It follows  $C_i = -\sum_{l=0}^{i-1} B_{i-l} C_l$  at step  $j + 1$  for  $i \leq j + 1$ . That is exact the same recursion as for the  $j$ th step, hence the the matrices  $C_i$ ,  $0 \leq i \leq j$  corresponding to steps  $j$  and  $j + 1$ , respectively, are the same.  $\square$

EXAMPLE A.2.6. Let the operator equation

$$\left( \sum_{i=0}^j a_i s^i \right) \left( \sum_{i=0}^j a_i s^i \right) u_j = (1+s)u_j$$

be given. We to compute coefficients  $a_i$ ,  $i \geq 0$ , such that this equation holds for all given  $u_j$ ,  $j \geq 0$ . The operator on the right-hand side is of type (A.2.2) with  $k = 1$ ,  $A = 1$ ,  $B = 1$ . A comparison of coefficients supplies  $a_0^2 = 1$ ,  $2a_0a_1 = 1$ ,  $2a_0a_2 + 2a_1a_1 = 0, \dots$ . Once  $a_0$  is determined, the whole sequence follows uniquely. The operator equation should hold for any  $U \in \mathcal{C}$ . Hence it is convenient to write

$$\left( \sum_{i=0}^j a_i s^i \right)^2 = 1 + s$$

and to define the square root of  $1 + s$  by

$$\sqrt{1+s} = \left( \sum_{i=0}^j a_i s^i \right),$$

where the coefficients  $a_i$  are determined by the above procedure. Observe, that by construction the coefficients  $a_i$  are in fact the Taylor coefficients of  $\pm\sqrt{1+s}$ ,  $s$  considered as complex variable and the square root defined with the usual branch cut. Observe further, that the formal expression  $\sqrt{1+su_j}$  has to be read as  $\sum_{i=0}^j a_i s^i u_j = \sum_{i=0}^j a_i u_{j-i}$ . Hence the Taylor coefficients in the formal expression have set to zero for  $i > j$ .

## Bibliography

- [1] B. Alpert, L. Greengard, and T Hagstrom. Rapid evolution of non-reflecting boundary kernels for time-domain wave propagation. *SIAM J. Numer. Anal.*, 37:1138–1164, 2000.
- [2] S. Amini, P. J. Harris, and D. T. Wilton. Coupled Boundary and Finite Element Methods for the Solution of the Dynamic Fluid-Structure Interaction Problem. In C. A. Brebbia and S. A. Orszag, editors, *Lecture Notes in Engineering*, volume 77. Springer Verlag, New York, Berlin, 1992.
- [3] S. Amini and S. M. Kirkup. Solution of Helmholtz equation in the exterior domain by elementary boundary integral methods. *J. Comput. Phys.*, 118:208–221, 1995.
- [4] A. Arnold. Numerically absorbing boundary conditions for quantum evolution equations. In *Proceedings of the International Workshop on Computational Electronics, Tempe, USA, 1995*. VLSI-Design 63-5, 1995.
- [5] Anton Arnold and Matthias Erhardt. Discrete transparent boundary conditions for wide angle parabolic equations in underwater acoustics. *J. Comput. Phys.*, 145(2):611–638, 1998.
- [6] A. Bamberger, B. Engquist, L. Halpern, and P. Joly. Higher order paraxial wave equation approximations in heterogeneous media. *SIAM J. Appl. Math.*, 48(1):129–154, 1988.
- [7] A. Bamberger, B. Engquist, L. Halpern, and P. Joly. Parabolic wave equation approximations in heterogeneous media. *SIAM J. Appl. Math.*, 48(1):99–128, 1988.
- [8] V. A. Baskakov and A. V. Popov. Implementation of transparent boundaries for numerical solution of the schrödinger equation. *Wave Motion*, 14:123–128, 1991.
- [9] A. Bayliss, M. Gunzburger, and E. Turkel. Boundary conditions for the numerical solution of elliptic equations in exterior domains. *SIAM J. Appl. Math.*, 42:430–451, 1982.
- [10] A. Bayliss and E. Turkel. Radiation boundary conditions for wave-like equations. *Communic. on Pure and Applied Math.*, 33:707–725, 1980.
- [11] J.-P. Bérenger. A Perfectly Matched Layer for the Absorption of Electromagnetic Waves. *J. Comput. Phys.*, 114:185–200, 1994.
- [12] J.-P. Bérenger. Three Dimensional Perfectly Matched Layer for the Absorption of Electromagnetic Waves. *J. Comput. Phys.*, 127:363–379, 1996.
- [13] Max Born and Emil Wolf. *Principle of Optics*. Pergamon Press Oxford, New York, Seoul, Tokyo, 6 edition, 1980.
- [14] A. Böttcher and S. M Grudsky. *Töplitz Matrices, Asymptotic Linear Algebra, and Functional Analysis*. Birkhäuser, Basel, Boston, Berlin, 2000.
- [15] P. G. Carrier, L. Greengard, and V. Rokhlin. A fast adaptive multipole algorithm for particle simulations. *SIAM J. Sci. Stat. Comput.*, 9(4):159–184, 1988.
- [16] L. Collatz. *Differentialgleichungen*. Teubner, 7 edition, 1990.
- [17] Francis Collino and Peter Monk. The perfectly matched layer in curvilinear coordinates. *SIAM J. Sci. Comput.*, 19(6):2061–2090, 1998.
- [18] Francis Collino. Perfectly Matched Absorbing Layers for the Paraxial Equations. *J. Comp. Phys.*, 131(1):164–180, 1997.
- [19] M.D. Collins. Higher-order Padé approximations for accurate and stable elastic parabolic equations with application to interface wave propagation. *J. Acoust. Soc. Am.*, 89:1050–1057, 1991.
- [20] M.D. Collins. A split-step Padé solution for the parabolic equation method. *J. Acoust. Soc. Am.*, 93:1736–1742, 1993.
- [21] D. Colton and R. Kress. *Inverse Acoustic and Electromagnetic Scattering Theory*. Springer-Verlag Berlin, Heidelberg, New York, 1992.
- [22] L. Demkovicz and K. Gerdes. Convergence of the infinite element methods for the Helmholtz equation. *Numerische Mathematik*, 79(1):11–42, 1998.
- [23] John W. Dettman. *Mathematical methods in physics and engineering*. McGraw-Hill, 2 edition, 1969.
- [24] John W. Dettman. *Applied complex variables*. Dover Publications, Inc., New York, 1984.
- [25] Peter Deuffhard and Folkmar Bornemann. *Numerische Mathematik II. Integration gewöhnlicher Differentialgleichungen*. Walter de Gruyter, Berlin, New York, 1994.

- [26] Doetsch. *Handbuch der Laplace-Transformation*. Birkhäuser Verlag, Basel, Stuttgart, 1955.
- [27] H. W. Engl. *Integralgleichungen*. Springer, Vienna/New York, 1997.
- [28] B. Engquist and A. Majda. Absorbing Boundary Conditions for the Numerical Simulation of Waves. *Mathematics of Computation*, 31(139):629–651, July 1977.
- [29] Bjorn Engquist and Andrew Majda. Radiation Boundary Condition for Acoustic and Elastic Wave Propagation. *Communications on Pure and Applied Mathematics*, XXXII:313–357, 1979.
- [30] Oliver G. Ernst. *Fast Numerical Solution of Exterior Helmholtz Problems with Radiation Condition By Imbedding*. PhD thesis, Department of Scientific Computing and Computational Mathematics of Stanford University, 1994.
- [31] Richard P. Feynmann, Robert B. Leighton, and Matthew Sands. *The Feynmann lectures on physics*. Addison Wesley Publishing Company, Reading Massachusetts, 1963.
- [32] Tilmann Friese, Frank Schmidt, and David Yevick. Transparent boundary conditions for a wide-angle approximation of the one-way Helmholtz equation. *J. Comput. Phys.*, 165(2):645–659, 2000.
- [33] Roland Geng. Berechnung der Eigenlösung des Schrödinger-Operators auf unbeschränkten Gebieten in 1d. Master's thesis, Freie Universität Berlin, Fachbereich Mathematik und Informatik, submission: 2001. in preparation.
- [34] Klaus Giebermann. *Schnelle Summationsverfahren zur numerischen Lösung von Integralgleichungen für Streuprobleme im  $R^3$* . PhD thesis, Universität Karlsruhe, Fakultät für Mathematik, 1997.
- [35] D. Givoli. Nonreflecting boundary conditions. *J. Comput. Phys.*, 94:1–29, 1991.
- [36] D. Givoli. A spatially exact non-reflecting boundary condition for time dependent problems. *Comput. Methods Appl. Mech. Engrg.*, 95:97–113, 1992.
- [37] D. Givoli and J. B. Keller. Nonreflecting boundary conditions for elastic waves. *Wave Motion*, 12:261–279, 1990.
- [38] Dan Givoli and Joseph B. Keller. A finite element method for large domains. *Computer Methods in Applied Mechanics and Engineering*, 76(6):41–66, 1989.
- [39] Gene H. Golub and Charles F. Loan. *Matrix computations*. The John Hopkins University Press, third edition 1996 edition, 1983.
- [40] L. Greengard. *The rapid evaluation of potential fields in particle systems*. MIT Press, Cambridge, Massachusetts, 1987.
- [41] L. Greengard and V. Rokhlin. A fast algorithm for particle simulations. *J. Comp. Phys.*, 73(2):325–348, 1987.
- [42] M. J. Grote and J. B. Keller. Exact nonreflecting boundary conditions for time-dependent wave equation. *SIAM J. of Applied Mathematics*, 55(2):280–297, April 1995.
- [43] M. J. Grote and J. B. Keller. On nonreflecting boundary conditions. *J. Comput. Phys*, 122:231–243, 1995.
- [44] M. J. Grote and J. B. Keller. Nonreflecting boundary conditions for time-dependent scattering. *J. Comput. Phys.*, 127:52–65, 1996.
- [45] M. J. Grote and J. B. Keller. Nonreflecting boundary conditions for Maxwell's equations. *J. Comput. Phys.*, 1997. submitted.
- [46] W. Hackbusch, C. Lage, and S.A. Sauter. On the Efficient Realization of Sparse Matrix Techniques for Integral Equations with Focus on Panel Clustering, Cubature and Software Design Aspects. In Wolfgang E. Wendland, editor, *Boundary Element Topics: proceedings of the Final Conference of the Priority Research Programme Boundary Element Methods 1989-1995 of the German Research Foundation*, pages 51–75. Springer Verlag, 1997.
- [47] W. Hackbusch and Z. P. Nowak. On the fast matrix multiplication in the boundary element method by panel clustering. *Numerische Mathematik*, 54:463–491, 1989.
- [48] G.R. Hadley. Multistep method for wide-angle beam propagation. *Optics Letters*, 17:1743–1745, 1992.
- [49] T. M. Hagstrom. Asymptotic expansions and boundary conditions for time-dependent problems. *SIAM J. Num. Anal.*, 23:948–958, 1986.
- [50] Thomas Hagstrom. Radiation boundary condition for the numerical simulation of waves. *Acta Numerica*, 8:47–106, 1999.
- [51] E. Hairer, C. Lubich, and M. Schlichte. Fast numerical solution of weakly singular Volterra integral equations. *J. Comput. Appl. Math.*, 23(1):87–98, 1988.
- [52] L. Halpern. Artificial boundary conditions for incompletely parabolic perturbations of hyperbolic systems. *SIAM J. Math. Anal.*, 22:1256–1283, 1991.
- [53] L. Halpern and L. N. Trefethen. Wide-angle one-way wave equations. *J. Acoust. Soc. Amer.*, 84(4):1397–1404, 1988.
- [54] Laurence Halpern and Jeffrey Rauch. Error analysis for absorbing boundary conditions. *Numerische Mathematik*, 51:459–467, 1987.

- [55] Isaac Harari and Thomas. J. R. Hughes. Analysis of continuous formulations underlying the computation of time-harmonic acoustics in exterior domains. *Comput. Methods Appl. Mech. Engrg.*, 97:103–124, 1997.
- [56] Eriikki Heikkola, Yuri A. Kuznetsov, Pekka Neittaanmäki, and Jari Toivanen. Fictitious Domain Methods for the Numerical Solution of Two-Dimensional Scattering Problems. *Journal of Comput. Phys*, 145(1):89–109, 1998.
- [57] Thorsten Hohage, Frank Schmidt, and Lin Zschiedrich. Solving time-harmonic scattering problems based on the pole condition:Convergence of the PML method. Technical report, Konrad-Zuse-Zentrum (ZIB), 2001. ZIB-Report 23-01, in preparation.
- [58] Thorsten Hohage, Frank Schmidt, and Lin Zschiedrich. Solving time-harmonic scattering problems based on the pole condition:Theory. Preprint 01-01, Konrad-Zuse-Zentrum (ZIB), 2001.
- [59] W.P. Huang and C.L. Xu. A wide-angle vector beam propagation method. *Photon. Technol. Lett.*, 4:1118–1119, 1992.
- [60] F. Ihlenburg. *Finite Element Analysis of Acoustic Scattering*, volume 132 of *Applied Mathematical Sciences*. Springer-Verlag New York Berlin Heidelberg, 1998.
- [61] Frank Ihlenburg. On fundamental aspects of exterior approximations with infinite elements. *Journal of Computational Acoustics*, 1999. manuscript, submitted 1999.
- [62] Claes Johnson. *Numerical solution of partial differential equations by the finite element method*. Cambridge University Press Cambridge, New York, New Rochelle Melbourne Sydney, 1987.
- [63] S. N. Karp. A convergent "far-field" expansion for two dimensional radiation functions. *Comm. Pure Appl. Math.*, 14:427–434, 1961.
- [64] J. B. Keller and D. Givoli. Exact non-reflecting boundary conditions. *J. Comput. Phys.*, 82:172–192, 1989.
- [65] R. Kosloff and D. Kosloff. Absorbing boundaries for wave propagation problems. *J. Comput. Phys.*, 63:363–376, 1986.
- [66] R. Kress. *Linear Integral Equations*. Springer-Verlag Berlin, Heidelberg, New York, 1989.
- [67] M. Lassas and E. Somersalo. On the existence and convergence of the solution of pml equations. *Computing*, 60:229–241, 1998.
- [68] R. Leis. *Initial Boundary Value Problems in Mathematical Physics*. Teubner, 1986.
- [69] Ch. Lubich. Convolution quadrature and discretized operational calculus i, ii. *Numer. Math.*, 52:129–145 and 413–425, 1988.
- [70] Ch. Lubich and A. Schädle. Fast convolution for non-reflecting boundary conditions. Preprint, Mathematisches Institut, Universität Tübingen, April 2001.
- [71] Reinhard März. *Integrated optics: design and modeling*. Artec House, Inc., 1995.
- [72] B. Mayfield. *Nonlocal Boundary Conditions for the Schrödinger Equation*. PhD thesis, University of Rhode Island, Providence, RI, 1989.
- [73] Erhard Meister. Modern Wiener-Hopf Methods in Diffraction Theory. In D. Bainov and V. Covachev, editors, *5th Int. Coll. on Differential Equations*, pages 245–256, 1995.
- [74] J. Mikusiński. *Operational Calculus*. Pergamon Press, London, New York, 1959.
- [75] P. Petropoulos. Reflectionless sponge layers as absorbing boundary conditions for the numerical solution of Maxwell equations in rectangular, cylindrical, and spherical coordinates. *SIAM J. Appl. Math.*, 60:1037–1058, 2000.
- [76] V. Rokhlin. Rapid solution of integral equations of classical potential theory. *J. Comput. Phys.*, 60(2):187–207, 1985.
- [77] V. Rokhlin. Diagonal forms of translation operators for the Helmholtz equation in three dimensions. *Appl. Comput. Harmon. Anal.*, 1(1):82–93, 1993.
- [78] A. Ruhe. An algorithm for numerical determination of the structure of a general matrix. *BIT*, 10:196–216, 1970.
- [79] S. Sauter. *Über die effiziente Verwendung des Galerkinverfahrens zur Lösung Fredholmscher Integralgleichungen*. PhD thesis, Universität Kiel, 1992.
- [80] F. Schmidt. Construction of discrete transparent boundary conditions for Schrödinger-type equations. *Surv. Math. Ind.*, 9:87–100, 1999.
- [81] F. Schmidt. Discrete nonreflecting boundary conditions for the Helmholtz equation. In A. Bermúdez, D. Gómez, C. Hazard, P. Joly, and J. E. Roberts, editors, *Fifth International Conference on Mathematical and Numerical Aspects of Wave Propagation*, pages 921–925, Waves 2000, Santiago de Compostella, Spain, 2000.
- [82] F. Schmidt and P. Deuffhard. Discrete transparent boundary conditions for the numerical solution of Fresnel's equation. *Computers Math. Applic.*, 29:53–76, 1995.
- [83] F. Schmidt, T. Friese, and D. Yevick. Transparent boundary conditions for split-step Padé approximations of the one-way Helmholtz equation. Submitted to *J. Comp. Phys.*, Preprint SC 99-46, ZIB, 1999.

- [84] F. Schmidt and R. März. On the reference wave vector of paraxial Helmholtz equations. *IEEE Journal of Lightwave Technology*, 14:2395–2400, 1996.
- [85] F. Schmidt and D. Yevick. Discrete transparent boundary conditions for Schrödinger-type equations. *J. Comput. Phys.*, 134:96–107, 1997.
- [86] Frank Schmidt. An Alternate Derivation of the Exact DtN-Map on a Circle. Preprint SC 98-32, Konrad-Zuse-Zentrum Berlin, December 1998.
- [87] Michal E. Taylor. *Pseudodifferential Operators*. Princeton University Press, Princeton, New Jersey, 1981.
- [88] Michal E. Taylor. *Partial Differential Equations*. Springer Verlag, New York, 1996.
- [89] L. N. Trefethen and L. Halpern. Well-posedness of one-way wave equations and absorbing boundary conditions. *J. Acoust. Soc. Amer.*, 47:421–435, 1986.
- [90] Calvin H. Wilcox. A generalization of theorems of Rellich and Atkinson. *Proc. Amer. Math. Soc.*, pages 271–276, 1956.
- [91] Baolin Yang and Peter G. Petropoulos. Plane-wave analysis and comparison of split-field, bi-axial, and uniaxial PML methods as abcs for pseudospectral electromagnetic wave simulations in curvilinear coordinates. *J. Comp. Phys.*, 146(2):747–774, 1998.
- [92] D. Yevick. Optimal absorbing boundary conditions. *J. Opt. Soc. Am. A*, 12:107–110, 1995.
- [93] D. Yevick, J. Yu, W. Bardyszewski, and M. Glasner. Stability issues in vector electric field propagation. *Photon. Technol. Lett.*, 7:656–658, 1995.
- [94] D. Yevick, J. Yu, and F. Schmidt. Analytic studies of absorbing and impedance-matched boundary layers. *IEEE Photonics Technology Letters*, 9:73–75, 1997.
- [95] E. Zeidler. *Applied Functional Analysis: Applications to Mathematical Physics*, volume 108 of *Applied Mathematical Sciences*. Springer-Verlag New York Berlin Heidelberg, 1997.
- [96] E. Zeidler. *Applied Functional Analysis: Main Principles and Their Applications*, volume 109 of *Applied Mathematical Sciences*. Springer-Verlag New York Berlin Heidelberg, 1997.
- [97] Lin Zschiedrich. *Transparent boundary conditions for time-harmonic scattering problems and time-dependent Schrödinger equations*. PhD thesis, Fachbereich Mathematik und Informatik, FU Berlin, submission: 2002. in preparation.



# Index

- Acoustics, 11
- angular frequency, 11
- Arzelá-Ascoli theorem, 72
  
- Banach fixed-point theorem, 77
- Bessel
  - equation, 21, 32, 105, 115
  - function, 21, 47
- boundary
  - artificial, 17, 23, 27, 28, 35, 48, 88
  - rigid, 13
- boundary condition
  - absorbing, 16
  - Bayliss-Gunzburger-Turkel (BGT), 27
  - discrete transparent, 16
  - exact, 16
  - radiating, 16
  - transparent, 16
- boundary element method, 22
- Butcher scheme, 118
  
- Cauchy integral theorem, 80
- collocation points, 109
- collocation scheme, 116
- condition number, 114
- conformal mapping, 58
- conservation property, 38, 86
- consistency order, 117
- continuity equation, 11
- coordinates, separable, 57
- cut function, 74
  
- diffraction, 12
- Dirichlet condition, 115
- Dirichlet-to-Neumann map, 9, 63
- dispersion relation, 11
- divergence condition, 11
- domain
  - artificial, 35, 62
  - computational, 13
  - exterior, 13
  - inhomogeneous, 13
- Dormand-Prince integrator, 115
  
- energy flux, 86
- Engquist-Majda boundary condition, 28
- Euler
  - equation, 11
  - explicit discretization, 115
  - implicit discretization, 118
  
- exterior domain
  - semi-discretization, 120
  - transformed coordinates, 123
  
- factorization method, 91
- fictitious domain, 35
- field
  - electric, 10
  - magnetic, 10
  - outgoing, 9
  - source, 9
  - strength, 10
  - time harmonic, 11
- Floquet theorem, 40
- Fourier mode, 65
- Fourier transform, 92
- Fredholm integral equation, 24
- Fredholm index, 64
- Fresnel equation, 12
- fundamental solution, 49
  
- Galerkin method, 133
- Garding inequality, 64
- Gauss-Radau integrators, 118
- geometrical optics, 13, 22
- Givens matrix, 42
- Green
  - function, 22, 23
  - identity, 13, 24
  - theorem, 57, 63
- Gronwall lemma, 79
  
- Hölder continuity, 70
- Hankel function, 21
  - asymptotic series, 48
- Helmholtz equation, 9, 11, 57
  - 1D with constant potential, 37
  - cylindric coordinates, 59
  - elliptic-hyperbolic coordinates, 59
  - periodic coefficients, 37
  - radially symmetric, 37
  - separable, 59
  - transformed coordinates, 59
- high-frequency limit, 87
  
- implicit midpoint discretization, 153
- infinite element method, 25
- inhomogeneity
  - waveguide-type, 14
- integral equation

- Fredholm, 13
- Volterra, 66, 75
- Wiener-Hopf, 13
- Jacobi matrix, 58
- jump function, 66
- Karp expansion theorem, 26
- Laplace domain method, 105
  - cut function approach, 115
  - real axis approach, 106
- Laplace transform, 37
- Laplace-Beltrami operator, 52, 61
- Laplacian
  - transformed coordinates, 58, 132
- Lax-Milgram theorem, 64
- Maxwell equations, 10
- Mikusinski operational calculus, 148
- multipole method, 34
- Newton law, 11
- normal ray, 125
- Nyström method, 25
- one-way equation, 16
- Pade approximation, 153
- panel clustering, 34
- paraxial beam propagation, 12, 15
- partial fraction decomposition, 94
- perfectly matched layer (PML), 35, 88
- periodic potential, 38
- phase factor, 14
- phase front, 102
- Planck quantum action, 12
- plane
  - infinite, 91
  - semi-infinite, 12, 22
- PML method, 88
- pole condition
  - exterior to a sphere, 52
  - Fourier modes, 52
  - general, 53
  - in 1D, 38
  - time-dependent problems, 56
- potential, 9
- power flux, 38
- pseudodifferential operator, 29
- quasi-periodic solution, 40
- radial ray, 125
- reference index, 12
- reflection, 43
- representation formula, 80
- Riccati equation, 170
- Riesz representation theorem, 64
- Runge-Kutta integrator, 114, 117
- Schrödinger equation, 12, 56, 144
  - abstract, 56
- Schur decomposition, 41
- separable coordinates, 20
- shift operator, 149
- Sobolev space, weighted, 26
- Sommerfeld radiation condition, 9
- source, 57
- spectral property, 86
- splines, 116
- stopband, 43
- symbol class, 29
- Töplitz matrix, 177
- time-harmonic solution, 11
- transmission, 43
- Volterra integral equation, 117
- wave
  - acoustic, 11
  - angular frequency, 11
  - electromagnetic, 10
  - plane, 11
  - quantum mechanic, 12
  - scalar equation, 11
  - vector, 11
  - vectorial equation, 10
- waveguide, 13
- wavelet, 105
- wavenumber, 11
  - distance dependent, 59
  - position dependent, 9
- Weber function, 21
- wide-angle one way equations, 144
- Wiener-Hopf
  - integral equation, 13, 24
  - technique, 13, 24, 91
- Wilcox expansion theorem, 26
- z-transforms, 39