

Domain Truncation, Absorbing Boundary Conditions, Schur Complements, and Padé Approximation

Martin J. Gander, Lukáš Jakabčín and Michal Outrata

September 30, 2023

Abstract

We show for a model problem that the truncation of an unbounded domain by an artificial Dirichlet boundary condition placed far away from the domain of interest is equivalent to a specific absorbing boundary condition placed closer to the domain of interest. This specific absorbing boundary condition can thus be implemented as a truncation layer terminated by a Dirichlet condition. We prove that the absorbing boundary condition thus obtained is a spectral Padé approximation about infinity of the transparent boundary condition. We also study numerically two improvements for this boundary condition – the truncation with an artificial Robin condition placed at the end of the truncation layer, and a Padé approximation about a different point than infinity. Both of these give new and substantially better results compared to using the artificial Dirichlet boundary condition at the end of the truncation layer. We prove our results in the context of linear algebra, using spectral analysis of finite and infinite Schur complements, which we relate to continued fractions. We illustrate our results with numerical experiments.

1 Introduction

The solution process of problems on unbounded domains usually requires a domain truncation, and hence artificial boundary conditions, leading to techniques such as *perfectly matched layers* (PML) or *absorbing boundary conditions* (ABC), see [5, 3]. At the discrete level, these closely relate to the problem of approximating the Schur complement in some sense, which inspired a number of iterative solvers, see, e.g., [11, 14] and the references therein. Our approach builds upon the eigendecomposition of the Schur complement, which for our model problem is very closely linked with the Fourier analysis of the Schur complement or, equivalently, the frequency domain analysis.

Domain truncation is also important in domain decomposition where a given computational domain is decomposed into many smaller subdomains, and then subdomain solutions are computed independently in parallel, see [14]. The solutions on the smaller subdomains can naturally be interpreted as solutions on truncated domains, and thus it is of interest to

use ABC or PML techniques at the interfaces between the subdomains, see also [9, 10, 11]. The classical Schwarz method [22] uses Dirichlet transmission conditions between subdomains and an overlap to achieve convergence [25]. In what follows the goal is to interpret the overlap as a specific ABC once the unknowns of the overlap are folded onto the interface (similarly to [20, 15]). Although the Schwarz method is not explicitly mentioned in what follows, it is one of the main applications for our results, see [14] for more information and corresponding numerical experiments. Note also that the Patch Substructuring Method [20, 15] is precisely such a method where the overlap was folded in.

Notably, the question of the *optimal PML* for problems with finite difference grids has been discussed in [17, 1] for the Laplace equation and then also extended to the Helmholtz equation in [8]. Our interest here is however different: we want to get a mathematical understanding of what object is obtained when one truncates an unbounded domain with a Dirichlet boundary condition after a finite layer of given length in which one still solves the partial differential equation, and how precisely the quality depends on this length. This is often done by people in applications for diffusive problems (e.g. finance), and is also done in the classical Schwarz method and all its variants like Additive and Multiplicative Schwarz. We only have as a second goal to try to improve this, by optimizing a Robin truncation at the end of the layer, or by modifying the equation in the entire layer in a simple way, linking this approach to PML. We also make substantial efforts to do that in a way that is both reasonably self-contained and easy to follow for readers from different mathematical communities, where such simple Dirichlet truncations are used for diffusive problems. This also includes introducing terminology for continued fractions and their types and properties in some detail, and also the Schur complement.

We start in Section 2 with some notation and definitions and continue in Section 3 by showing that there exists a limit of the Schur complement as the width of the truncation layer goes to infinity, and that the Schur complement of a finite width truncation with a Dirichlet condition is a spectral Padé approximation around infinity of the unbounded width limit. Next, we explore numerically how the spectral approximation changes when the Dirichlet condition at the end of the truncation layer is replaced by a Robin condition in Section 4, present an optimized choice for the Robin parameter and propose a new type of truncation layer in Section 5. We end with concluding remarks and possible extensions in Section 6.

2 Model Problem

We use as our model problem the partial differential equation (PDE)

$$\begin{aligned} (\eta - \Delta)u &= f & \text{in } \Omega &:= (0, +\infty) \times (0, 1), \quad \eta > 0, \\ u &= 0 & \text{on } \partial\Omega. \end{aligned} \tag{1}$$

We assume that the support of the right-hand side function f is *localized* in $\Omega^a := (0, a) \times (0, 1)$ and having $b \geq a$ we set $\Omega^b := (0, b) \times (0, 1) \subset \Omega$ as the artificially truncated region containing Ω_a , see Figure 1.

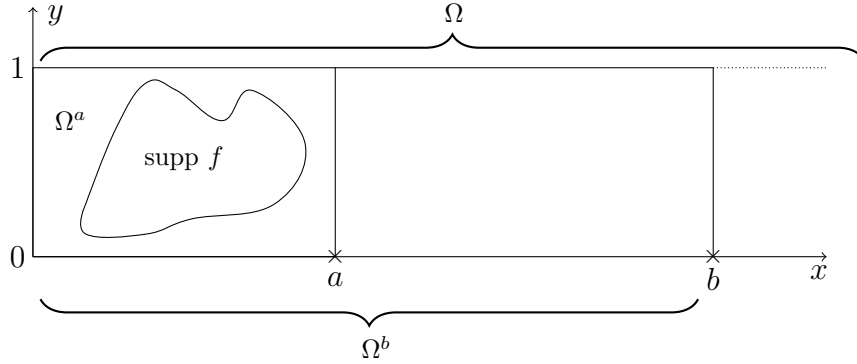


Figure 1: The unbounded strip domain in \mathbb{R}^2 with $\Omega = (0, +\infty) \times (0, 1)$.

Discretizing with a standard finite difference scheme, we denote by N the number of interior grid columns, and obtain the mesh size $h := 1/(N + 1)$. Assuming we have

$$a = (N^a + 1)h \quad \text{and} \quad b = (N^b + 1)h, \quad (2)$$

we obtain the discretized problems

$$A\mathbf{u} = \mathbf{f}, \quad A^b\mathbf{u}^b = \mathbf{f}^b, \quad A^a\mathbf{u}^a = \mathbf{f}^a, \quad (3)$$

with the right-hand side vectors $\mathbf{f}^a := [\mathbf{f}_1^T, \dots, \mathbf{f}_{N^a}^T]^T$, $\mathbf{f}^b := [(\mathbf{f}^a)^T, \mathbf{0}^T, \dots, \mathbf{0}^T]^T$ and $\mathbf{f} := [(\mathbf{f}^b)^T, \mathbf{0}^T, \dots]^T$, and the matrices

$$A^\star := \frac{1}{h^2} \begin{pmatrix} D_1 & -I & & & \\ -I & \ddots & \ddots & & \\ & \ddots & D_{N^\star-1} & -I & \\ & & -I & D_{N^\star} & \end{pmatrix}, \quad A := \frac{1}{h^2} \begin{pmatrix} h^2 A^b & & & & \\ & -I & & & \\ & -I & D_{N^b+1} & \ddots & \\ & & \ddots & \ddots & \end{pmatrix}, \quad (4)$$

where \star stands for either a or b (and thus changes the number of block rows and block columns) and each block has dimension N (vectors) or $N \times N$ (matrices) related to a particular set of grid column variables. The matrix I is the $N \times N$ identity and the diagonal blocks D_j are given by

$$D_j := D = \begin{pmatrix} \eta h^2 + 4 & -1 & & \\ -1 & \ddots & -1 & \\ & -1 & \eta h^2 + 4 & \end{pmatrix} \in \mathbb{R}^{N \times N}. \quad (5)$$

Here, it is enough to understand the infinite-dimensional system in (3) as the limit of the finite-dimensional one as $b \rightarrow +\infty$; for more details on infinite matrices, see, e.g., the concise review [23] or the historical overview [6].

Thanks to the localization of \mathbf{f} we can formulate a problem *only*¹ on Ω^a such that its solution coincides with $\mathbf{u}^b|_{\Omega^a}$, simply by eliminating the unknowns from the truncation domain

¹This is of particular interest for the domain decomposition methods, see Section 1 and also [14].

$\Omega^b \setminus \Omega^a$. This solution is then an approximation of $\mathbf{u}|_{\Omega^a}$. The continuous level formulation requires the *Dirichlet-to-Neumann* operator (see, e.g., [11] in the context of domain decomposition) and its approximation on finite difference grids in this context has been studied in [17, 1]. We carry out this elimination by “folding in” the variables $(\mathbf{u}_{N^b}^b, \dots, \mathbf{u}_{N^{a+1}}^b)$, starting with $\mathbf{u}_{N^b}^b$ and working our way from the right to the left on the grid. Recalling (3), for $b < +\infty$ these satisfy the equations

$$-\frac{\mathbf{u}_{N^{b-1-i}}^b}{h^2} + \frac{D_{N^b-i}\mathbf{u}_i^b}{h^2} - \frac{\mathbf{u}_{N^{b+1-i}}^b}{h^2} = 0, \quad -\frac{\mathbf{u}_{N^{b-1}}^b}{h^2} + \frac{D_{N^b}\mathbf{u}_{N^b}^b}{h^2} = 0, \quad (6)$$

with $i \in \{1, \dots, N^b - N^a\}$, where the index i counts the progress “from right to left” in the domain $\Omega^b \setminus \Omega^a$. The elimination process corresponds to the block Gaussian elimination (block size N) that eventually calculates the *Schur complement* of the unknowns $\mathbf{u}^b|_{\Omega^a}$ in A^b (see, [16, p. 103]). We summarize this in the definition below.

Definition 2.1 (Schur complement) *Having $b < \infty$ we can reduce $A^b \mathbf{u}^b = \mathbf{f}^b$ to*

$$\tilde{A}^a \mathbf{u}^b|_{\Omega^a} = \mathbf{f}^a, \quad \text{with} \quad \tilde{A}^a = \frac{1}{h^2} \begin{pmatrix} D_1 & -I & & & \\ -I & \ddots & & \ddots & \\ & \ddots & D_{N^a-1} & -I & \\ & & -I & T_{N^a}^b & \end{pmatrix}, \quad (7)$$

where the block $T_{N^a}^b$ is called the Schur complement. It can be calculated recursively (see [21, Sections 1.3.2 and 1.4.3]) as

$$T_{N^b}^b := D_{N^b} = D \quad \text{and} \quad T_{N^{b-i}}^b := D_i - (T_{N^{b-i+1}}^b)^{-1} = D - (T_{N^{b-i+1}}^b)^{-1}, \quad (8)$$

for $i \in \{1, \dots, N^b - N^a\}$.

Comparing \tilde{A}^a and A^a , the only change is in the last block where the Dirichlet boundary condition block has been replaced by the Schur complement $T_{N^a}^b$, representing the truncation layer (or the “far-field” domain) unknowns in $\Omega^b \setminus \Omega^a$. Hence $\mathbf{u}^b|_{\Omega^a}$ approaches $\mathbf{u}|_{\Omega^a}$ in the limit as $b \rightarrow \infty$, but increasing b makes the defining recurrence in (8) longer. If b goes to infinity, the corresponding Schur complement matrix $T_{N^a}^\infty$ is still governed by (8), namely

$$T_{N^a}^\infty = D - (T_{N^a}^\infty)^{-1}, \quad \text{i.e.,} \quad (T^\infty)^2 - DT^\infty + I = 0. \quad (9)$$

Notably, this equation does not depend on N^a , and hence also its solution $T_{N^a}^\infty \equiv T^\infty$. To solve (9), we start the following section by changing the basis we work in to the eigenbasis of D , effectively applying a discrete Fourier transform in the y variable.

3 Spectral analysis

Writing D from (5) as $D = D_{yy} + 2I$, where D_{yy} is the 3-point finite difference stencil discretization of $\eta - \partial_{yy}$ multiplied by h^2 , we recall that $D_{yy} = Q^T \text{diag}(z_1, \dots, z_N)Q$ with

$$z_k := \eta h^2 + 4 \sin^2 \left(\frac{k\pi}{2(N+1)} \right) \quad \text{and} \quad \mathbf{q}_k := \left[\sqrt{\frac{2}{N+1}} \sin \left(\frac{k\pi}{N+1} j \right) \right]_{j=1}^N \in \mathbb{R}^N, \quad (10)$$

where Q is unitary and symmetric, with the eigenvectors \mathbf{q}_k in its columns. We can thus write $D = Q^T \Lambda Q$ with $\Lambda := \text{diag}(2 + z_1, \dots, 2 + z_N)$ as the eigendecomposition of D .

Remark 1 *Calculating in the eigenbasis of D is a necessity for our Schur complement analysis but in treating each eigenmode separately we would add yet another index to the already loaded notation. That is why, instead of referring to the particular eigenvalues $2 + z_k$ of D or z_k of D_{yy} we introduce new variable z and treat all quantities depending on $2 + z_k$ or z_k as functions of z . This way we avoid the index k whenever we can but in some places the reference to a particular eigenvalue or eigenmode is unavoidable and we keep the index k reserved for the eigenmode notation throughout the text.*

3.1 Diagonalization and convergence of the Schur Complement

Changing the basis for the Schur complement definition in (8) gives

$$\hat{T}_{N^b}^b = QDQ^T = \Lambda \quad \text{and} \quad \hat{T}_{N^b-i}^b = QDQ^T - Q(T_{N^b-i+1}^b)^{-1}Q^T = \Lambda - (\hat{T}_{N^b-i+1}^b)^{-1},$$

where $i = 1, \dots, N^b - N^a$ and all of the matrices $\hat{T}_{N^b-i}^b$ are diagonal. Working with the diagonal entries only, each of them becomes a function of z_k and also follows the recurrence. Recalling Remark 1, we write

$$\hat{t}_{N^b}^b(z) = (2 + z) \quad \text{and} \quad \hat{t}_{N^b-i}^b(z) = (2 + z) - \frac{1}{\hat{t}_{N^b-i+1}^b(z)} \quad \text{for } i = 1, \dots, N^b - N^a,$$

but in order to further simplify the notation, we label these scalar functions only by i rather than $N^b - i$ and without relabeling obtain²

$$\hat{t}_0^b(z) = (2 + z) \quad \text{and} \quad \hat{t}_i^b(z) = (2 + z) - \frac{1}{\hat{t}_{i+1}^b(z)} \quad \text{for } i = 1, \dots, N^b - N^a. \quad (11)$$

We obtain an analogous recurrence for the solution \mathbf{u}^b in (6). Setting $\hat{\mathbf{u}}_{N^b-i}^b := Q\mathbf{u}_{N^b-i}^b$ we get

$$\begin{aligned} -\frac{\hat{\mathbf{u}}_{N^b-1}^b}{h^2} + \frac{\Lambda \hat{\mathbf{u}}_{N^b}^b}{h^2} &= -\frac{\hat{\mathbf{u}}_{N^b-1}^b}{h^2} + \frac{\hat{T}_{N^b}^b}{h^2} \hat{\mathbf{u}}_{N^b}^b = 0 \\ -\frac{\hat{\mathbf{u}}_{N^b-1-i}^b}{h^2} + \frac{\Lambda \hat{\mathbf{u}}_{N^b-i}^b}{h^2} - \frac{\hat{\mathbf{u}}_{N^b+1-i}^b}{h^2} &= -\frac{\hat{\mathbf{u}}_{N^b-1-i}^b}{h^2} + \frac{\hat{T}_{N^b-i}^b}{h^2} \hat{\mathbf{u}}_{N^b-i}^b = 0, \end{aligned} \quad (12)$$

²This way, the scalar function labeling directly corresponds to the number of steps of the block Gaussian elimination we have already carried out. This notation becomes the natural one for the mathematical tools used later in this manuscript.

with $i = 1, \dots, N^b - N^a$. Turning to the limit case $b \rightarrow +\infty$ for $T_{N^a}^b$, we can now treat each mode separately, obtaining a scalar problem instead of (9). Setting

$$\lim_{b \rightarrow +\infty} \hat{t}_{N^b - N^a}^b(z) =: \hat{t}^\infty(z), \quad (13)$$

we observe that

$$(\hat{t}^\infty)^2(z) - (2+z)\hat{t}^\infty(z) + 1 = 0, \quad (14)$$

with the two solutions $\hat{\tau}^{\infty,1}(z) = \frac{2+z+\sqrt{(2+z)^2-4}}{2}$ and $\hat{\tau}^{\infty,2}(z) = \frac{2+z-\sqrt{(2+z)^2-4}}{2}$, and

$$(\hat{\tau}^{\infty,1}(z)) (\hat{\tau}^{\infty,2}(z)) = 1 \quad \text{and} \quad 0 < \hat{\tau}^{\infty,2}(z) < 1 < \hat{\tau}^{\infty,1}(z). \quad (15)$$

Next, we show that one of the solutions $\hat{\tau}^{\infty,1}(z), \hat{\tau}^{\infty,2}(z)$ acts as the limit Schur complement for our solution vector $\mathbf{u}^b|_{\Omega^a}$.

The key observation is that the characteristic polynomial of the recurrence relation in (12) is preserved through the limit process and thus the solutions $\hat{\tau}^{\infty,1}(z), \hat{\tau}^{\infty,2}(z)$ of the limit equation (14) coincide with the roots of the characteristic polynomial of the recurrence relation in (12) given by $p_z(r) = -r^2 + (2+z)r - 1$. This together with the explicit formula for the solution of the recurrence relation (12) is enough to solve the matrix equation defining T^∞ in (9). In order to do so, we will evaluate the functions of z at the particular points of interest z_k , i.e., at the eigenvalues of the matrix D .

Theorem 3.1 *The Schur complement $T_{N^a}^b$ defined in (8) converges to $T^{\infty,1}$ solution of the formal limit equation (9) as $b \rightarrow +\infty$, i.e., the eigenvectors of those matrices are equal and the eigenvalues $\hat{t}_{N^b - N^a}^b(z_k)$ of the Schur complement converge to $\hat{\tau}^{\infty,1}(z_k)$ for all $k = 1, \dots, N$.*

Proof For any b large enough, we fix a particular grid-column index $j \in \{N^a, \dots, N^b\}$ and observe that the solution subvector $\hat{\mathbf{u}}_j^b = [\hat{u}_{j,1}^b, \dots, \hat{u}_{j,N}^b]^T \in \mathbb{R}^N$ follows the recurrence in (12). This recurrence has a closed form solution, namely there exist pairs of constants $(\nu_1^b, \mu_1^b), \dots, (\nu_N^b, \mu_N^b)$ independent of j such that

$$\hat{\mathbf{u}}_j^b = \begin{bmatrix} \mu_1^b (\hat{\tau}^{\infty,1}(z_1))^{j-N^a} + \nu_1^b (\hat{\tau}^{\infty,2}(z_1))^{j-N^a} \\ \vdots \\ \mu_N^b (\hat{\tau}^{\infty,1}(z_N))^{j-N^a} + \nu_N^b (\hat{\tau}^{\infty,2}(z_N))^{j-N^a} \end{bmatrix}.$$

Furthermore, recalling (15) it follows that

$$(\hat{\tau}^{\infty,1}(z_k))^{N^b - N^a} \rightarrow +\infty \quad \text{as } b \rightarrow +\infty, \quad \text{for any } k = 1, \dots, N.$$

As $\hat{\mathbf{u}}_{N^b}^b = 0$ for any $b > a$ we have $|\hat{\mathbf{u}}_{N^b}^b| \rightarrow +\infty$ as $b \rightarrow +\infty$, showing that for each k necessarily $\mu_k^b \rightarrow 0$ as $b \rightarrow +\infty$. Since $\hat{\mathbf{u}}_{N^a}^b$ converges as $b \rightarrow +\infty$ (see, e.g., [24, Sections 2 and 3]) we obtain also the limits $\nu_k^\infty := \lim_{b \rightarrow +\infty} \nu_k^b$ as $b \rightarrow +\infty$ and therefore

$$\hat{\mathbf{u}}_j^\infty \equiv \lim_{b \rightarrow \infty} \hat{\mathbf{u}}_j^b = \begin{bmatrix} \nu_1^\infty (\hat{\tau}^{\infty,2}(z_1))^{j-N^a} \\ \vdots \\ \nu_N^\infty (\hat{\tau}^{\infty,2}(z_N))^{j-N^a} \end{bmatrix}. \quad (16)$$

Taking $j = N^a + 1$ we can solve the k -th entry of the recurrence in (12) for $\hat{t}_{N^b - N^a - 1}^b(z_k)$ and using the finite difference stencil, we obtain

$$\hat{t}_{N^b - N^a - 1}^b(z_k) = \frac{(2 + z_k)\hat{u}_{N^a + 1, k}^b - \hat{u}_{N^a + 2, k}^b}{\hat{u}_{N^a + 1, k}^b} = \frac{\hat{u}_{N^a, k}^b}{\hat{u}_{N^a + 1, k}^b} \rightarrow \frac{1}{\hat{\tau}_k^{\infty, 2}(z_k)} = \hat{\tau}_k^{\infty, 1}(z_k),$$

where we used (16) and (15) before and after taking the limit respectively. Using the defining equation (14) we obtain

$$\hat{t}_{N^b - N^a}^b(z_k) = \frac{1}{(2 + z_k) - \hat{t}_{N^b - N^a - 1}^b(z_k)} \rightarrow \frac{1}{(2 + z_k) - \hat{\tau}_k^{\infty, 1}(z_k)} = \frac{1}{\hat{\tau}_k^{\infty, 2}(z_k)} = \hat{\tau}_k^{\infty, 1}(z_k).$$

□

Hence Theorem 3.1 implies

$$\hat{t}^\infty(z) = \frac{2 + z + \sqrt{(2 + z)^2 - 4}}{2} = 1 + \frac{z}{2} + \frac{\sqrt{z^2 + 4z}}{2} = \left(1 + \frac{z}{2} \left(1 + \sqrt{1 + \frac{4}{z}}\right)\right), \quad (17)$$

and for $b < \infty$ we recall (11) and in the same fashion obtain

$$\hat{t}_0^b(z) = 2 + z, \quad \hat{t}_1^b(z) = 2 + z - \frac{1}{2 + z}, \quad \hat{t}_2^b(z) = 2 + z - \frac{1}{\hat{t}_{N^b - 1}^b(z)} = 2 + z - \frac{1}{2 + z - \frac{1}{2 + z}},$$

and by the recursive definition in (11), the i -th term is given by

$$\hat{t}_i^b(z) = \frac{2 + z}{h^2} - \frac{\frac{1}{h^2}}{2 + z - \frac{1}{2 + z - \frac{1}{2 + z - \frac{1}{2 + z}}}}, \quad (18)$$

having i “levels” of the fraction. After some elementary calculations $\hat{t}_i^b(z)$ can be written as a rational function of degree $i + 1$. Notice that each level of $\hat{t}_i^b(z)$ in (18) corresponds to elimination of unknowns from one grid column, i.e., to one step of the block Gaussian elimination mentioned above. This is not surprising but it gives perhaps a more pleasant way of viewing and analyzing the matrix recurrence in Definition 3.3. We continue by a simple observation regarding the functions \hat{t}^∞ and \hat{t}_i^b .

Remark 2 *By a subsequent re-insertion we obtain*

$$\hat{t}^\infty(z) = 2 + z - \frac{1}{\hat{t}^\infty(z)}, \quad \hat{t}^\infty(z) = 2 + z - \frac{1}{2 + z - \frac{1}{\hat{t}^\infty(z)}}, \quad \dots$$

and so on. This suggests that the function $\hat{t}^\infty(z)$ is equal to the infinite continued fraction

$$\hat{t}^\infty(z) = 2 + z - \frac{1}{2 + z - \frac{1}{2 + z - \dots}},$$

and $\hat{t}_i^b(z)$ in (18) are approximations in the sense of a truncation after i levels, hence $\hat{t}_i^b(z)$ is called a truncated continued fraction.

The theory of *continued fractions* links various areas of mathematics, e.g., Padé approximations, orthogonal polynomials, Vorobyev’s moment matching problem, Riemann-Stieltjes integrals, Gauss quadrature and the method of conjugate gradients (see [19, 26, 7] and also [18, Section 3.3.2 - 3.3.6] for further references). In this manuscript we restrict ourselves to assume *no* knowledge of this field. As a result, the text is self-contained and easier to access for a wider audience. But this comes with a price – using the full strength of the continued fractions theory we could meaningfully refine the results as well as connect these with the above mentioned areas. We postpone such presentation to an upcoming manuscript, which will make a good use of the present one as a stepping stone. This also justifies the use of [2] as our main reference, where the author uses continued fractions only as one of the possible tools to arrive at *Padé approximants* – precisely the spirit in which we will use the continued fractions. We refer the interested reader to [19, 26, 7, 18] for more detailed expositions of the connected topics.

We continue in Section 3.2 with a concise summary of the continued fraction results and the connected simple algebraic calculations. We formulate these in terms of an auxiliary variable α , given by

$$\alpha := \frac{4}{z}. \quad (19)$$

This change of variables is unavoidable as we will need to expand about $+\infty$ and the standard way of defining this is to expand the same function but of a reciprocal argument about 0 – hence (19). This is why we do not consider (19) as a proper change of variables, which would otherwise necessitate tedious calculations of the derivatives of the function composition. In fact, the *true* change of variables consists only in multiplying by 4 and hence does not require a re-computation of the derivatives. Hence we rewrite $\hat{t}^\infty(z), \hat{t}_i^b(z)$ as functions of α instead of z and for the sake of simplicity we do not relabel, i.e., we abuse the notation to have

$$\hat{t}(z) \equiv \hat{t}(z(\alpha)) := \hat{t}(\alpha), \quad \text{for } \hat{t} = \hat{t}^\infty \text{ or } \hat{t} = \hat{t}_i^b. \quad (20)$$

3.2 Padé Approximation and Continued Fractions

We follow the notation from [2], i.e., the $[M/L]$ -Padé approximant of $f(z)$ is denoted by $[M/L]_f \equiv [M/L]_f(z)$. We start with Padé theory and proceed with continued fractions ([2, Chapter 4]).

Theorem 3.2 ([2, Theorem 1.5.3, 1.5.4, 1.5.1]) *Let $f(z)$ be a real function of a real variable. Then the following holds provided the Padé approximants exist:*

1. Let $\alpha, \beta \in \mathbb{R}$. Then $\alpha + \beta[M/L]_f = [M/L]_{\alpha+\beta f}$.

2. Let $m \geq 1$ and $f(z) = \sum_{j=0}^{+\infty} c_j z^j$ be a formal power series. Setting $g(z) = \frac{1}{z^m} \left(f(z) - \sum_{j=0}^{m-1} c_j z^j \right)$ and assuming $M - m \geq L - 1$ we have

$$[M - m/L]_g(z) = \frac{1}{z^m} \left([M/L]_f(z) - \sum_{j=0}^{m-1} c_j z^j \right).$$

3. Let $f(0) \neq 0$ and set $g(z) = 1/f(z)$. Then $[M/L]_g(z) = 1/[L/M]_f(z)$.

Definition 3.3 A continued fraction is given by sequences of real numbers $\{a_j\}_j, \{b_j\}_j$ – the numerator and the denominator sequence of the continued fraction – and has the general form

$$b_0 + \frac{a_1}{b_1 + \frac{a_2}{b_2 + \frac{a_3}{\ddots}}} =: b_0 + \sum_{j=1}^{+\infty} \frac{a_j}{b_j} \equiv b_0 + \frac{a_1}{b_1 + \frac{a_2}{b_2 + \cdots}},$$

where the sum is to be understood only formally. The continued fraction is called infinite as long as $a_j, b_j \neq 0$ for all j . The n -th truncation (or convergent) of a continued fraction is given by

$$\frac{A_n}{B_n} = b_0 + \sum_{j=1}^n \frac{a_j}{b_j} = b_0 + \frac{a_1}{b_1 + \frac{a_2}{b_{n-2} + \frac{a_3}{\ddots + \frac{a_n}{b_n}}}},$$

where A_n and B_n are the n -th truncation (or convergent) numerator and denominator.

Replacing the scalars a_j and/or b_j by linear (or affine) functions of a real variable z , A_n and B_n become polynomials in z and the n -th truncation of the continued fraction becomes a rational function in z . Different settings of this framework lead to different types of continued fractions. Most notably, a continued fraction is called regular C-fraction (short for regular classical continued fraction), provided it has the form

$$b_0 + \frac{a_1 z}{1 + \frac{a_2 z}{1 + \frac{a_3 z}{\ddots}}} \equiv b_0 + \frac{a_1 z}{1} + \frac{a_2 z}{1} + \cdots,$$

with $a_j \neq 0$ for all j . If, moreover, $a_j > 0$ for all j , then it is called S-fraction (short for Stieltjes continued fraction). If the continued fraction takes the form

$$b_0 + \frac{r_1}{z + s_1 - \frac{r_2}{z + s_2 - \frac{r_3}{\ddots}}} \equiv b_0 + \frac{r_1}{z + s_1} - \frac{r_2}{z + s_2} + \cdots,$$

with $r_j \neq 0$ for all j then it is called J-fraction (short for Jacobi continued fraction). For more details on the introduced types of continued fractions as well as other types of continued fractions (e.g., non-regular C-fraction, T-fraction, P-fraction, ...) we refer also to [19] and [26] and references therein.

First, we note that we have ignored the questions of convergence of infinite continued fractions and we refer the reader to [19] and [26]. Next, notice that one function can be represented by two seemingly different continued fractions (different in type and/or in the coefficient values) and one way to recognize their equality is via the *three-term recurrence relation* (see [2, Theorem 4.1.1, pp.106]). We have that

$$\begin{aligned} A_{-1} &= 1, & A_0 &= b_0, & A_n &= b_n A_{n-1} + a_n A_{n-2}, \\ B_{-1} &= 0, & B_0 &= 1, & B_n &= b_n B_{n-1} + a_n B_{n-2}, \end{aligned} \tag{21}$$

and assuming the n -th truncation (convergent) of two continued fractions are equal for any n , the infinite continued fractions are equal as well. Last but not least, we note that some authors will call a continued fraction an S -fraction even though the fraction itself does not meet the definition above but can be *transformed* into a continued fraction that does. We next recall a basic transformation rule of continued fractions.

Lemma 3.4 ([2, Section 4.1, pp. 105-106]) *Let $\{a_j\}_j, \{b_j\}_j$ be two real sequences of the numerators and denominators of a continued fraction as in Definition 3.3. Let $\{e_j\}_j$ be a sequence of real numbers different from zero. Then we have*

$$b_0 + \frac{a_1}{b_1 + \frac{a_2}{b_2 + \frac{a_3}{b_3 + \dots}}} = b_0 + \frac{e_1 a_1}{e_1 b_1 + \frac{e_1 e_2 a_2}{e_2 b_2 + \frac{e_2 e_3 a_3}{e_3 b_3 + \dots}}},$$

For purposes of this text we present immediately the continued fraction result for the square root function, which is of interest to us³.

Theorem 3.5 ([2, Section 4.6, Theorem 4.4.3 and formula (6.4) on pp. 139]) *For any $\alpha \in (-1, +\infty)$ ⁴ we have*

$$\sqrt{1 + \alpha} = 1 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{2 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{\ddots}}}}} = 1 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{2 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{\ddots}}}}} + \frac{a_n}{b_n + \dots}, \quad (22)$$

with $b_0 = 1$, $b_j = \frac{3+(-1)^j}{2}$ and $a_j = \frac{\alpha}{2}$, $j \geq 1$. Moreover, for any n the $[n, n]$ -Padé approximation of $\sqrt{1 + \alpha}$ expanded about $\alpha = 0$ is given by the $(2n)$ -th truncation of the continued fraction in (22) and the $[n + 1, n]$ -Padé approximation of $\sqrt{1 + \alpha}$ expanded about $\alpha = 0$ is given by the $(2n + 1)$ -st truncation of the continued fraction in (22).

Remark 3 *By a direct computation we see that*

$$\sqrt{1 + \alpha} = 1 + \frac{\alpha}{2 + \frac{\alpha}{2 + \frac{\alpha}{2} + \dots}},$$

i.e., the representation in (22) can be written as a cyclic S -fraction⁵ with $a_j = 1/2$ for all j .

³We refer to the book of Baker but the original result is due to Gauss, who showed a much more general result for the hypergeometric function ${}_2F_1$; for more details see [26, Chapter XVIII] or [19, Chapter VI].

⁴There is a misprint in [2, equation (6.4), page 139]. The authors state the convergence “for all z except $-\infty < z \leq 1$ ” but the result also holds for $z \in (-1, 1]$.

⁵Infinite continued fractions with periodic sequences $\{a_j\}, \{b_j\}$ are called cyclic continued fractions.

The rest of this section is devoted to auxiliary results, the first of which links a truncation of the S -fraction from Theorem 3.5 and a truncation of the J -fraction from Remark 2. Notice that the continued fractions are not identical but rather differ in the absolute term.

Lemma 3.6 *Let α be real and consider the two continued fractions*

$$\tau(\alpha) := \frac{\frac{\alpha}{2}}{2 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{2 + \frac{\frac{\alpha}{2}}{1 + \dots}}}} \quad \text{and} \quad \sigma(\alpha) := \frac{1}{1 + \frac{4}{\alpha} - \frac{1}{2 + \frac{4}{\alpha} - \frac{1}{2 + \frac{4}{\alpha} - \frac{1}{2 + \frac{4}{\alpha} - \dots}}}},$$

and denote their n -th truncations by $A_n(\alpha)/B_n(\alpha)$ and $C_n(\alpha)/D_n(\alpha)$. For any $n = 1, 2, \dots$ we have

$$A_{2n}(\alpha)/B_{2n}(\alpha) = C_n(\alpha)/D_n(\alpha).$$

Proof Using Lemma 3.4 we transform $\tau(\alpha)$ and without further relabeling we obtain

$$\tau(\alpha) := \frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \dots}}}}}}, \quad (23)$$

and by a direct computation confirm that equality holds for $n = 1$. Next, we notice that the continued fraction (23) can be written in cyclic form with *the core* R given by

$$R = \frac{4}{\alpha} + \frac{1}{1 + \frac{1}{R}}, \quad (24)$$

i.e., the continued fraction can be obtained by a successive re-insertion of the core equality (24) into itself, e.g.,

$$\underbrace{\frac{1}{\frac{4}{\alpha} + \frac{1}{1}}}_{= \frac{A_2(\alpha)}{B_2(\alpha)}}, \quad \underbrace{\frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \frac{1}{\frac{4}{\alpha} + \frac{1}{1}}}}}_{= \frac{A_4(\alpha)}{B_4(\alpha)}}, \quad \underbrace{\frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \frac{1}{\frac{4}{\alpha} + \frac{1}{1 + \dots}}}}}}}_{= \frac{A_6(\alpha)}{B_6(\alpha)}}, \quad \dots$$

In this way every re-insertion adds two elements of the numerator and denominator sequences, and using the algebraic identity

$$\frac{1}{1 + \frac{1}{R}} = 1 - \frac{1}{1 + R},$$

we reformulate the core equality (24) to obtain

$$1 + R = 2 + \frac{4}{\alpha} - \frac{1}{1 + R}, \quad (25)$$

and notice that the core equality in (25) is the one that generates the J -fraction $\sigma(\alpha)$.

Hence we have shown that for $n \geq 2$ the $2n$ re-insertions of the core R in the equality (24) is equal to n re-insertions of the core $1 + R$ in the equality (25), finishing the proof. \square

We now build upon Lemma 3.6 by contracting the S -fraction in (22) into a J -fraction.

Proposition 3.7 *Let α be real and set the continued fractions $\tau(\alpha)$ and $\sigma(\alpha)$ as in Lemma 3.6. Moreover, we define the continued fractions*

$$\tilde{\tau}(\alpha) := \frac{1}{1 + \tau(\alpha)} \quad \text{and} \quad \phi(\alpha) := 1 - \frac{1}{2 + \frac{4}{\alpha} - \frac{1}{2 + \frac{4}{\alpha} - \frac{1}{2 + \frac{4}{\alpha} - \dots}}}$$

with n -th truncations $\tilde{A}_n(\alpha)/\tilde{B}_n(\alpha)$ and $E_n(\alpha)/F_n(\alpha)$ with $E_0 = F_0 = 1$. Then for $n \geq 0$

$$A_{2n+1}(\alpha)/B_{2n+1}(\alpha) = E_n(\alpha)/F_n(\alpha).$$

Proof The equality for $n = 0$ holds by inspection. Taking $n \geq 1$, we use Lemma 3.6 for the continued fraction $\tilde{\tau}(\alpha)$, and obtain

$$\tilde{A}_{2n+1}(\alpha)/\tilde{B}_{2n+1}(\alpha) = \frac{1}{1 + A_{2n}(\alpha)/B_{2n}(\alpha)} = \frac{1}{1 + C_n(\alpha)/D_n(\alpha)},$$

and having the truncations $C_n(\alpha), D_n(\alpha)$ of $\sigma(\alpha)$ from Lemma 3.6 it remains to show that

$$\frac{1}{1 + C_n(\alpha)/D_n(\alpha)} = 1 - E_n(\alpha)/F_n(\alpha). \quad (26)$$

The cyclic parts of both $\sigma(\alpha)$ and $\phi(\alpha)$ coincide and we denote them by $\tilde{\sigma}(\alpha)$,

$$\tilde{\sigma}(\alpha) := \frac{1}{2 + \frac{4}{\alpha} - \frac{1}{2 + \frac{4}{\alpha} - \dots}}. \quad (27)$$

Recalling (15), we have

$$\{\hat{t}^\infty\}^{-1}(\alpha) := \frac{1}{\hat{t}^\infty(\alpha)} = 1 + \frac{2}{\alpha} - \frac{2}{\alpha}\sqrt{1+\alpha}, \quad (28)$$

and using Theorem 3.2 part 3 we get

$$[i/i]_{\hat{t}^\infty}(\alpha) = \frac{1}{[i+1/i]_{\{\hat{t}^\infty\}^{-1}(\alpha)}}$$

for any $i \geq 1$. By a direct computation we obtain

$$\{\hat{t}^\infty\}^{-1}(\alpha) = 1 + \frac{2}{\alpha} - \frac{2}{\alpha}\sqrt{1+\alpha} = 1 - 2\frac{1}{\alpha}(\sqrt{1+\alpha} - 1),$$

and hence by the Padé approximant calculus (see Theorem 3.2 parts 1 and 2) we obtain

$$[i/i]_{\{\hat{t}^\infty\}^{-1}(\alpha)} = 1 - 2\frac{1}{\alpha}([i+1/i]_{\sqrt{1+\alpha}}(\alpha) - 1).$$

Using the continued fraction representation from Theorem 3.5, we obtain

$$[i/i]_{\{\hat{t}^\infty\}^{-1}(\alpha)} = 1 - \frac{2}{\alpha} \left(1 + 1 + \frac{A_{2i+1}(\alpha)}{B_{2i+1}(\alpha)} - 1 \right) = 1 - \frac{2}{\alpha} \frac{\frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{2 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{b_{2i-1} + \frac{\dots}{b_{2i} + \frac{a_{2i+1}}{b_{2i+1}}}}}}}}{\dots}}$$

where the sequences $\{a_j\}_j, \{b_j\}_j$ are given as in Theorem 3.5 and $A_{2i+1}(\alpha)/B_{2i+1}(\alpha)$ is the $(2i+1)$ -st truncation of the continued fraction $\tau(\alpha)$ from Lemma 3.6. Hence we have

$$[i/i]_{\{\hat{t}^\infty\}^{-1}(\alpha)} = 1 - \frac{1}{1 + \frac{\frac{\alpha}{2}}{2 + \frac{\frac{\alpha}{2}}{1 + \frac{\frac{\alpha}{2}}{b_{2i-1} + \frac{\dots}{b_{2i} + \frac{a_{2i+1}}{b_{2i+1}}}}}}}, \quad (29)$$

and by a straight-forward manipulation (see Proposition 3.7) we observe that the continued fraction on the right-hand side in (29) is the $(2i+1)$ -st truncation of the continued fraction

$$\tilde{\tau}(\alpha) := \frac{1}{1 + \tau(\alpha)}.$$

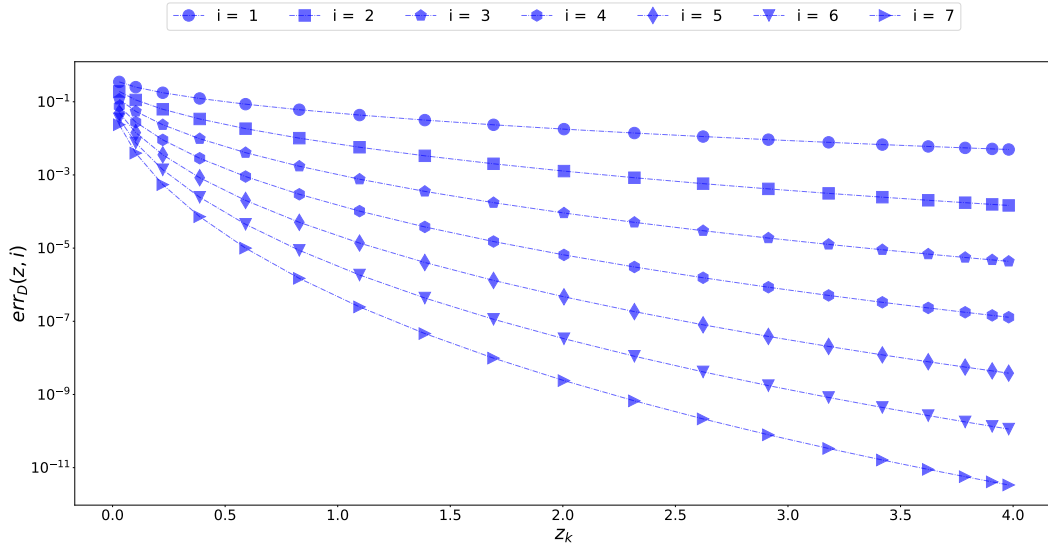


Figure 2: Plots of the function $err_D(z, i)$ at the points z_k with the parameters set to $N = 20$ and $\eta = 2$. The value i corresponds to the number of grid columns in $\Omega^b \setminus \Omega^a$.

We plot the error err_D for small i in and appropriate z in Figure 2. We see that $err_D(z, i)$ quickly decreases as z tends towards the right endpoint of the spectrum, which is to be understood as the expansion point. This becomes more pronounced for larger i , i.e., for higher order Padé approximations, i.e., when b increases. We see that the error is still large for z far away from the right endpoint, i.e., the ABC struggles with the low frequency mode approximation. We try improving this⁷ in the next section by considering a Robin boundary condition at the end of the truncation layer, $x = b$.

4 Robin boundary condition for truncation

We see that the Padé approximation error is far from optimal. Replacing the Dirichlet boundary condition at $x = b$ with a homogeneous Robin boundary condition⁸ with the Robin parameter $p \geq 0$ at b , i.e., with

$$\frac{\partial u}{\partial n} + pu = 0 \quad \text{at } x = b,$$

we hope to improve this. Using a centered finite difference approximation as before, the Robin condition can be discretized with the so-called *ghost point* trick. We use a centered

⁷In practice, the low frequency modes can be also solved by coupling our ABC with some effective low-frequency solver, e.g., some type of multigrid or multilevel DD scheme. However the focus here is to efficiently improve the ABC itself.

⁸A Robin boundary condition is a simple approximation to the transparent boundary condition and works in general substantially better; for subdomain truncation in domain decomposition see [9] and [14].

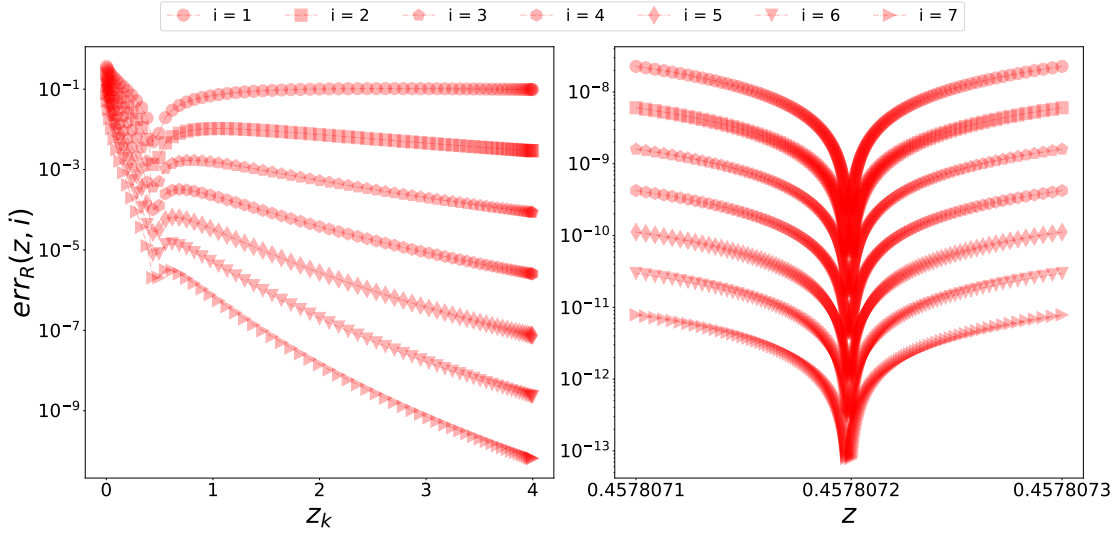


Figure 3: Left: plots of $err_R(z, i)$ at the points z_k (see (10)) evaluated for different number of grid columns i in $\Omega^b \setminus \Omega^a$, with $p = 50$, $N = 70$ and $\eta = 2$. Right: plots of the same functions under the same settings but zoomed in on the cusp (and thus plotted over artificial variables z rather than the eigenvalues z_k). In addition we also show the error err_P from Section 5.

and that $err_R(z, i)$ is smooth except at a finite number of points, equation (34) defines ζ_p as an implicit function of p and the other parameters of the problem. For $i = 1$ we get

$$err_R(z, i) = \left(1 + \frac{z}{2} \left(1 + \sqrt{1 + \frac{4}{z}} \right) \right) - 1 + ph + \frac{z}{2} = \frac{z}{2} \sqrt{1 + \frac{4}{z}} - ph,$$

which gives ζ_p as the positive root of the quadratic equation

$$\zeta_p^2 + 4\zeta_p - 4p^2h^2 = 0 \implies \zeta_p = -2 + 2\sqrt{1 + p^2h^2}. \quad (35)$$

Numerically, this formula worked for all different settings we have tried and, e.g., the numerical independence of ζ_p on i is already visible on the example in Figure 3.

Next, we try numerically to optimize p so that the infinity norm of $err_R(z)$ is minimized, i.e., we search for p that equioscillates the maximum of $err_R(z)$ on the left and on the right of ζ_p , and show the results in Figure 4. The relative improvement in the infinity norm of replacing the Dirichlet condition with the Robin one for that setting is roughly 5 fold.

Running the optimization while varying i , i.e., the number of grid columns from a to b , we obtain Table 1, again for $N = 200$ and $\eta = 2$. We see that the improvement over the Dirichlet truncation increases with increasing number of layers. The corresponding results over a larger range of i are shown graphically in Figure 5.

In Figure 5 we varied i as powers of 2 from $2^1 = 2$ to $2^8 = 256$ on the left and then up to 2^{15} on the right and observe a linear dependence in the log-log scale on the left, i.e., for

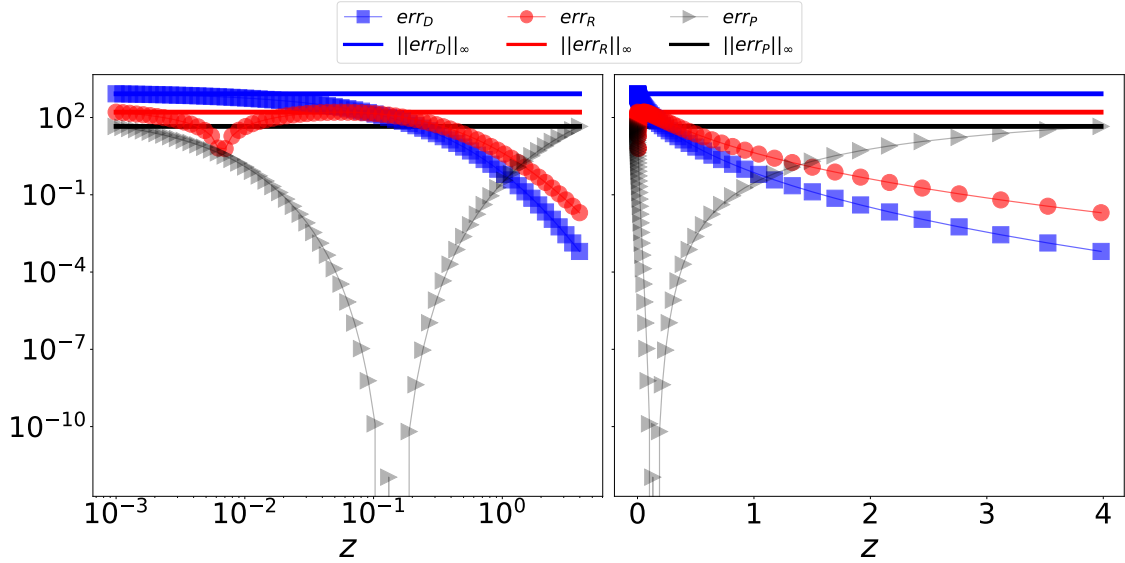


Figure 4: Left: minimization over p of the infinity norm of the Robin condition error, clearly showing the equioscillation. Right: optimized error compared with the corresponding Dirichlet condition error. We set $N = 200, i = 4$ and $\eta = 2$ and note that instead of z_k from (10) we take logarithmically equidistant z from the interval (30).

i	$p^*(i)$	$\frac{\ err_D\ _\infty}{\ err_R\ _\infty}$
1	27.4013	2.569
2	13.7783	3.924
4	8.2295	5.167
8	5.6016	6.598
16	4.3271	8.940

i	optimal z_0	$\frac{\ err_D\ _\infty}{\ err_P\ _\infty}$	$\frac{\ err_R\ _\infty}{\ err_P\ _\infty}$
1	0.4356	3.691	1.441
2	0.2101	10.091	2.572
4	0.1409	18.446	3.569
8	0.0932	86.163	13.058
16	0.0680	3595.822	402.186

Table 1: Evolution of the optimized Robin parameter $p^*(i)$ depending on expansion point z_0 depending on the the number of layers i and the improvement ratio from the Dirichlet condition error to the Robin condition error in the infinity norm.

Table 2: Evolution of the optimized Robin boundary condition error to the error of the approximation $\check{t}_{z_0}^i$.

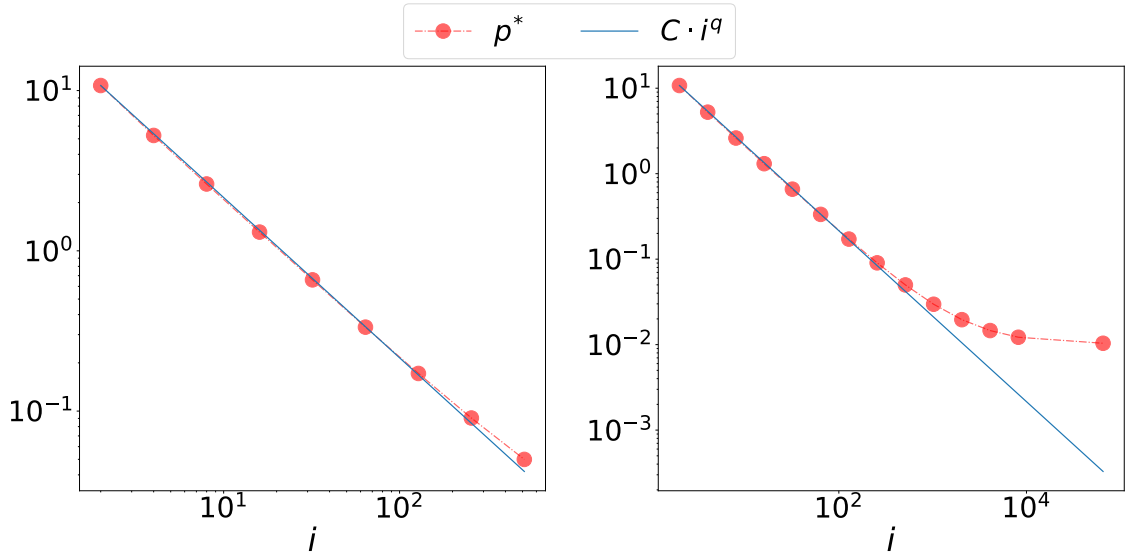


Figure 5: Dependence of the optimized Robin parameter $p^*(i)$ on the number of layers i added after a compared with the predicted behavior. The value i corresponds to the number of grid columns in $\Omega^b \setminus \Omega^a$.

values $i \leq 256$, and fitting the line gives the law

$$p^*(i) \sim C \cdot i^q, \quad \text{with } C \approx 11, q \approx -1. \quad (36)$$

The range $i \leq 256$ (and hence also the approximation (36)) in our eyes well covers the practically interesting values of i but it is clear that in general $p^*(i)$ does not follow the proposed relation (36).

Although the change and optimization of the Robin condition at $x = b$ offers a considerable improvement over the Dirichlet condition, we still observe for both of these the qualitatively identical behavior for z around the right endpoint of the spectrum. Naturally, we would like to shift this expansion, i.e., move the zero of the error from the right endpoint of the spectrum inside, and analogously to finding p^* and ζ_p we would like to get the *optimal* expansion point that minimizes the maximum of the approximation error. We explore this direction further in the following section.

5 Shifting the Padé expansion point

Taking some $\alpha_0 > 0$ and introducing the new variable

$$\tilde{\alpha} := \frac{\alpha - \alpha_0}{1 + \alpha_0} \quad \text{and hence} \quad \alpha(\tilde{\alpha}) = \tilde{\alpha} \cdot (1 + \alpha_0) + \alpha_0, \quad (37)$$

a direct computation gives

$$\sqrt{1 + \alpha} = \sqrt{1 + \alpha_0} \sqrt{1 + \tilde{\alpha}},$$

and expanding the right-hand side about 0 then corresponds to expanding the left-hand side about α_0 . Using Theorem 3.5, Lemma 3.6 and Proposition 3.7 we get

$$\sqrt{1+\alpha} = \sqrt{1+\alpha_0} \left(1 + \frac{\frac{\tilde{\alpha}}{2}}{1 + \frac{\frac{\tilde{\alpha}}{2}}{2 + \frac{\frac{\tilde{\alpha}}{2}}{1 + \frac{\frac{\tilde{\alpha}}{2}}{2+\dots}}}}} \right) = \sqrt{1+\alpha_0} \left(1 + \frac{\tilde{\alpha}}{2} \left(1 - \frac{1}{2 + \frac{4}{\tilde{\alpha}} - \frac{1}{2 + \frac{4}{\tilde{\alpha}} - \dots}} \right) \right).$$

Notice that the equality is valid only for the *formal, infinite* continued fraction and once we truncate, the correspondence follows from Proposition 3.7. Setting $\check{t}_{\alpha_0}^{\infty}(\tilde{\alpha}) := \hat{t}^{\infty}(z(\alpha(\tilde{\alpha})))$ we get

$$\check{t}_{\alpha_0}^{\infty}(\tilde{\alpha}) = 1 + \frac{2}{\tilde{\alpha}(1+\alpha_0)+\alpha_0} \left(1 + \left(1 + \frac{\tilde{\alpha}}{2} \right) \sqrt{1+\alpha_0} \right) - \frac{2}{\tilde{\alpha}(1+\alpha_0)+\alpha_0} \cdot \sqrt{1+\alpha_0} \cdot \frac{\tilde{\alpha}}{2} \cdot \frac{1}{2 + \frac{4}{\tilde{\alpha}} - \frac{1}{2 + \frac{4}{\tilde{\alpha}} - \dots}},$$

and based on Theorem 3.8 the truncation after i levels of $\check{t}_{\alpha_0}^{\infty}$ results in the $[i+1, i+1]$ -Padé approximant of \hat{t}^{∞} about α_0 . We define $\check{t}_{\alpha_0}^i(\tilde{\alpha})$ as

$$\check{t}_{\alpha_0}^i(\tilde{\alpha}) := 1 + \frac{2}{\tilde{\alpha}(1+\alpha_0)+\alpha_0} \left(1 + \left(1 + \frac{\tilde{\alpha}}{2} \right) \sqrt{1+\alpha_0} \right) - \frac{\tilde{\alpha}}{\tilde{\alpha}(1+\alpha_0)+\alpha_0} \cdot \sqrt{1+\alpha_0} \cdot \underbrace{\frac{1}{2 + \frac{4}{\tilde{\alpha}} - \frac{1}{2 + \frac{4}{\tilde{\alpha}} - \dots}}}_{i \text{ "levels"}},$$

and continue by focusing on the formulation of $\check{t}_{\alpha_0}^i$ as a function of z rather than $\tilde{\alpha}$. Recalling the definition of $\tilde{\alpha}$ in (37) we have¹⁰

$$z = \frac{4}{\alpha} = \frac{4}{\tilde{\alpha}(1 + \frac{4}{z_0}) + \frac{4}{z_0}},$$

which leads to

$$\tilde{\alpha} = \frac{4z_0 - 4}{4 + z_0} \quad \text{and hence} \quad \frac{4}{\tilde{\alpha}} = \frac{4 + z_0}{\frac{z_0}{z} - 1}.$$

Without relabeling¹¹ the function we can write

$$\check{t}_{z_0}^i(z) = 1 + \frac{z}{2} \left(1 + \left(1 + 2 \frac{\frac{z_0}{z} - 1}{4 + z_0} \right) \sqrt{1 + \frac{4}{z_0}} \right) - \frac{\left(\frac{1}{1 + \frac{4}{z_0}} - \frac{z}{4} \frac{\frac{z_0}{z_0}}{1 + \frac{4}{z_0}} \right) \sqrt{1 + \frac{4}{z_0}}}{2 + \frac{4+z_0}{\frac{z_0}{z} - 1} - \frac{1}{2 + \frac{4+z_0}{\frac{z_0}{z} - 1} - \dots}}, \quad (38)$$

$\underbrace{\hspace{15em}}_{i \text{ "levels"}}$

¹⁰Note that we used z_k above as the points of interest for the variable z for $k = 1, \dots, N$. In contrast to that, here we use $z_0 := 4/\alpha_0$ as the shifted expansion point of the Padé approximation, i.e., as the zero point of the corresponding approximation error of \hat{t}^{∞} . Thus the meaning of z_0 is *qualitatively* different compared to z_1, \dots, z_N .

¹¹Meant in the spirit of (20). We signal the variable by the expansion point in subscript from α_0 to z_0 .

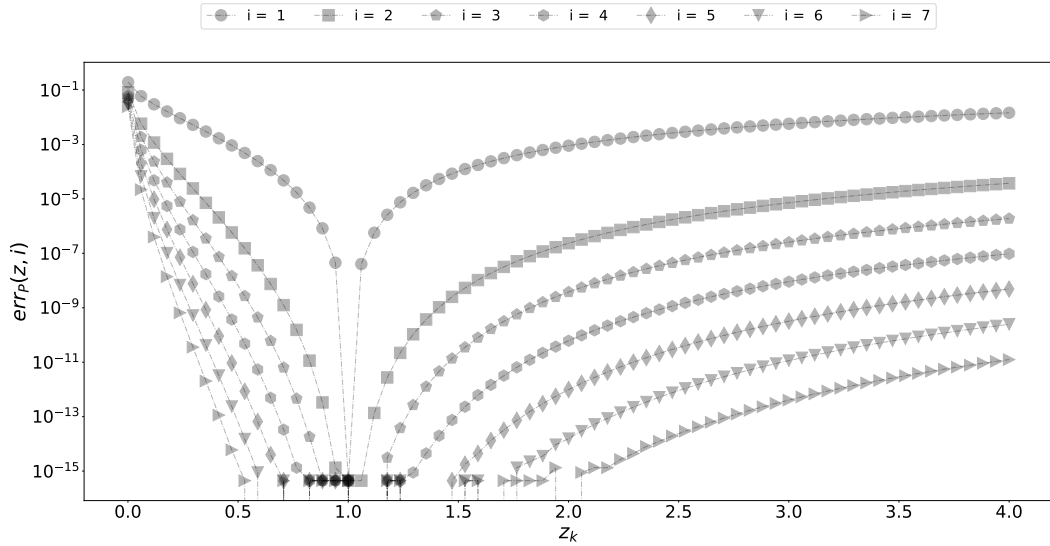


Figure 6: Plots of the function $err_P(z, i)$ at points equally spaced in the interval $[0, 4]$ evaluated for different values of i , for $\alpha_0 = 4$ (and thus $z_0 = 1$), $N = 70$ and $\eta = 2$. The value i corresponds to the number of grid columns in $\Omega^b \setminus \Omega^a$.

and thereby define the error function $err_P(z, i)$ (P for Padé) by

$$err_P(z, i) := |\hat{t}^\infty(z) - \check{t}_{z_0}^i(z)|.$$

The expectation is that the error function $err_P(z, i)$ should have one root at $z_0 = 4/\alpha_0$, which should get numerically more pronounced as i increases, in contrast to $err_R(z, i)$, and indeed, this is fully supported by the numerical results which we show in Figure 6.

Again, we turn our attention to finding the optimal z_0 , i.e., such that the error equioscillates on the left and on the right of z_0 and we present the results first in Figure 4, observing an 18-fold improvement over the Dirichlet case and hence roughly 3-fold improvement over the Robin case. The improvements become even more pronounced when increasing i as we show in Table 2. Finally, in Figure 7 we plot the evolution of the optimal z_0 as a function of i . We can see that for $i \leq 64$ there seems to be a trajectory for the optimal choice of z_0 , possibly convergent. But for i around 80 the error function becomes numerically equal to zero on the entire interval (10) and thus the optimization routine converges very close to or exactly at the initial guess, which was taken as 1.

We conclude this section by linking the above proposed approximation back to the physical problem and its solution methods by introducing a new PML technique that stems from the above approximation. Recalling the progress from (8) and (7) to (11), we need to move now in the opposite direction, starting with (38) and working up to the new block matrix we will call \check{A}^b .

Notice that $\check{t}_{z_0}^i(z)$ is structurally similar to $\bar{t}_i^b(z)$, containing structurally identical J -fractions. The first difference is in the absolute term added to the J -fractions and the second is the multiplicative factor in $\check{t}_{z_0}^i(z)$ in front of the J -fraction (which is not present

and thereby for any $i = 0, \dots, N^b - (N^a + 1)$ we have $\check{D}_{N^b-i} = \check{D} \neq \check{D}_{N^a}$ with

$$\begin{aligned} \check{D}_i &= Q^T \begin{pmatrix} 2 + z_1 \frac{4+z_0}{\mu_0-\mu_1} & & \\ & \ddots & \\ & & 2 + z_{N^r-1} \frac{4+z_0}{\mu_0-\mu_{N^r-1}} \end{pmatrix} Q \\ &= 2I + (4 + \eta h^2 + \mu_0)(D - 2I)(\mu_0 I - D_{yy})^{-1}, \end{aligned} \quad (40)$$

where Q is given in (10) and D is the diagonal block of the original problem, see (5). Focusing on the N^a -th block row, we obtain

$$1 + \frac{z}{2} \left(1 + \left(1 + 2 \frac{z_0 - 1}{4 + z_0} \right) \sqrt{1 + \frac{4}{z_0}} \right) = 1 + \frac{z}{2} + z \sqrt{\frac{1}{4} + \frac{1}{z_0}} + \frac{z_0 - z}{4 + z_0} \sqrt{1 + \frac{4}{z_0}}$$

for the absolute term, and the multiplicative terms reads

$$\left(\frac{1}{1 + \frac{4}{z_0}} - z \frac{1}{z_0 + 4} \right) \sqrt{1 + \frac{4}{z_0}}.$$

Hence, we set the diagonal block as

$$\begin{aligned} \check{D}_{N^a} &:= Q^T \begin{pmatrix} 1 + \frac{z_1}{2} + z_1 \sqrt{\frac{1}{4} + \frac{1}{z_0}} + \frac{\mu_0 - \mu_1}{z_0 + 4} \sqrt{1 + \frac{4}{z_0}} & & \\ & \ddots & \\ & & 1 + \frac{z_N}{2} + z_{N^r-1} \sqrt{\frac{1}{4} + \frac{1}{z_0}} + \frac{\mu_0 - \mu_{N^r-1}}{z_0 + 4} \sqrt{1 + \frac{4}{z_0}} \end{pmatrix} Q \\ &= \frac{1}{2} D + \sqrt{\frac{1}{4} + \frac{1}{z_0}} (D - 2I) + \frac{\sqrt{1 + \frac{4}{z_0}}}{z_0 + 4} (\mu_0 I - D_{yy}), \end{aligned}$$

and the off-diagonal block then reads

$$J = \frac{\sqrt{1 + \frac{4}{z_0}}}{1 + \frac{4}{z_0}} I - \frac{1}{z_0 + 4} (D - 2I).$$

We finish this section with the following remark.

Remark 5 *The formula (39) contains an explicit inverse, which is clearly unpractical but can be easily avoided by multiplying the block-rows $N^a + 1, \dots, N^b - 1, N^b$ in (39) by the matrix $M := \mu_0 I - D_{yy}$, which leads to*

$$\frac{1}{h^2} \begin{pmatrix} D_1 & -I & & & & \\ -I & \ddots & \ddots & & & \\ & \ddots & \check{D}_{N^a} & -J & & \\ & & -M & \check{D}_{N^a+1} M & \ddots & \\ & & & \ddots & \ddots & -M \\ & & & & -M & \check{D}_{N^b} M \end{pmatrix},$$

where no inverse of a matrix appears. An overall deeper understanding of \check{A}^b and its continuous counterpart are clearly of interest and will be discussed in future work.

6 Conclusion and future work

We proved for a model problem that truncation of the unbounded computational domain by a Dirichlet boundary condition at a certain distance from the domain of interest, is a spectral Padé approximation about infinity of the transparent boundary condition at the boundary of the domain of interest and that the degree of the Padé approximation increases with the distance. We then replaced the Dirichlet truncation condition by a Robin truncation condition at the end of the truncation layer, and showed that this considerably improves the behavior around a different point in the spectrum but loses the above Padé approximation property. We showed how to optimize the Robin parameter leading to an equioscillation property. In search to obtain a Padé approximation about a different point in the spectrum, we then proposed a new approximant in the eigenspace (leading to a new PML/ABC method for this problem), with much better truncation properties than the Robin truncation. In order to keep the exposition self-contained and of reasonable length we postponed the further results on the value of the optimal Robin parameter as well as the approximation properties of the new PML/ABC method (and the optimal choice of the expansion point) to a new upcoming manuscript, where we aim to lay out these in detail, making proper use of the theory of continued fractions.

Recognizing we worked with a very particular problem, there are some straightforward generalizations. First, none of the computations required the particular choice of D in (5). As long as D is symmetric and positive-definite, all of the computations still work and the only change is in the interval of interest for the minimization of the Robin parameter p and the shifted expansion point z_0 in Section 4 and Section 5. This even holds if D is only symmetric, non-singular and with eigenvalues outside the interval $(-\infty, -1]$. If the spectrum intersects the interval $(-\infty, -1]$, the square root becomes a complex number and the computations move to the complex plane – in fact this is true for any diagonalizable non-singular normal matrix D . The Helmholtz problem is the canonical example and in fact a very similar technique was used to establish a result related to Theorem 3.8 in [13]. If D is not normal, then the eigenvectors cannot be chosen to form an orthonormal basis of \mathbb{R}^N (or \mathbb{C}^N) but the formulas would follow (based on the spectrum) one of the above mentioned cases in the same way, but one could not use the results directly. For example, the improvement factor would not be of immediate interest as the condition number of the eigenbasis would play an important role in computing the optimized Robin parameter p . Alternatively, we can use the technique from [4], where the authors use integration over some contour enclosing the numerical range of the matrix. If the matrix is diagonalizable and singular, then the modes corresponding to the zero eigenvalues do not admit the formulation of the function $\hat{t}_i^b(z)$ as in (17) but the analysis would work for the rest of the modes, based on the normality and spectrum of the matrix. In the case that the matrix is not diagonalizable, it is not immediately clear how to generalize any of the results based on the available Jordan form.

7 Acknowledgement

We would like to thank Prof. Zdeněk Strakoš and Prof. Bernhard Beckermann for their very useful comments and references to the literature that helped significantly to improve the manuscript.

References

- [1] S. Asvadurov, V. Druskin, M. N. Guddati, and L. Knizhnerman. On optimal finite-difference approximation of PML. *SIAM Journal on Numerical Analysis*, 41(1):287–305, 2003.
- [2] G.A. Baker. *Padé Approximants Part I: Basic theory*. Addison-Wesley, Reading, 1981.
- [3] A. Bayliss and E. Turkel. Radiation boundary conditions for wave-like equations. *Communications in Pure and Applied Mathematics*, 33(6):707–725, 1980.
- [4] B. Beckermann and L. Reichel. Error estimates and evaluation of matrix functions via the Faber transform. *SIAM Journal on Numerical Analysis*, 47(5):3849–3883, 2009.
- [5] J. P. Bérenger. A perfectly matched layer for the absorption of electromagnetic waves. *Journal on Computational Physics*, 114(2):185–200, 1994.
- [6] M. Bernkopf. A history of infinite matrices. *Archive for History of Exact Sciences*, 4(4):308–358, 1968.
- [7] C. Brezinski. *History of Continued Fractions and Padé Approximants*, volume 12. Springer, 2012.
- [8] V. Druskin, S. Güttel, and L. Knizhnerman. Near-optimal perfectly matched layers for indefinite Helmholtz problems. *SIAM Review*, 58(1):90–116, 2016.
- [9] M. J. Gander. Optimized Schwarz methods. *SIAM Journal on Numerical Analysis*, 44(2):699–731, 2006.
- [10] M. J. Gander. Schwarz methods over the course of time. *Electronical Transactions on Numerical Analysis*, 31(5):228–255, 2008.
- [11] M. J. Gander and H. Zhang. A class of iterative solvers for the Helmholtz equation: factorizations, sweeping preconditioners, source transfer, single layer potentials, polarized traces, and optimized Schwarz methods. *SIAM Review*, 61(1):3–76, 2019.
- [12] M.J. Gander and A. Schädle. The Pole condition: A Padé approximation of the Dirichlet to Neumann operator. In *Domain Decomposition Methods in Science and Engineering XIX, Lecture Notes in Computational Science and Engineering*. Springer, 2010.

- [13] M.J. Gander and A. Schädle. On the relationship between the pole condition, absorbing boundary conditions and perfectly matched layers. In preparation, 2023.
- [14] M.J. Gander and H. Zhang. Schwarz methods by domain truncation. *Acta Numerica*, pages 1–134, 2022.
- [15] M.J. Gander, L. Halpern, F. Magoules, and F. Roux. Analysis of patch substructuring methods. *International Journal of Applied Mathematics and Computer Science*, 17(3): 395–402, 2007.
- [16] G. H. Golub and C. F. van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, 1996.
- [17] D. Ingerman, V. Druskin, and L. Knizhnerman. Optimal finite difference grids and rational approximations of the square root : I. Elliptic problems. *Communications on Pure and Applied Mathematics*, 53(8):1039–1066, 2000.
- [18] J. Liesen and Z. Strakoš. *Krylov Subspace Methods: Principles and Analysis*. Oxford University Press, Oxford, 2013.
- [19] L. Lorentzen and H. Waadeland. *Continued Fractions with Applications*. Elsevier, North Holland, 1992.
- [20] F. Magoulès, F.-X. Roux, and L. Series. Algebraic approximation of Dirichlet-to-Neumann maps for the equations of linear elasticity. *Computer Methods in Applied Mechanics and Engineering*, 195(29–32):3742–3759, 2006.
- [21] M. Outrata. *Schwarz methods, Schur complements, preconditioning and numerical linear algebra*. PhD thesis, University of Geneva, Math Department, University of Geneva, Math Department, 2022.
- [22] H. A. Schwarz. Über einen Grenzübergang durch alternierendes Verfahren. *Vierteljahrsschrift der Naturforschenden Gesellschaft in Zürich*, 15:272–286, 1870.
- [23] P. N. Shivakumar and K. C. Sivakumar. A review of infinite matrices and their applications. *Linear Algebra and its Applications*, 430(4):976–998, 2009.
- [24] P. N. Shivakumar and R. Wong. Linear equations in infinite matrices. *Linear Algebra and its Applications*, 7(1):53–62, 1973.
- [25] A. Toselli and O. Widlund. *Domain Decomposition Methods - Algorithms and Theory*. Springer Berlin, Heidelberg, 2004.
- [26] H. S. Wall. *Analytic Theory of Continued Fractions*. Chelsea Publishing Company, New York, 1967.