

THE ANALYSIS OF MATCHED LAYERS

L. HALPERN* and S. PETIT-BERGEZ

LAGA, Institut Galilée, Université Paris XIII
93430 Villetaneuse, France
*halpern@math.univ-paris13.fr

J. RAUCH

Department of Mathematics,
University of Michigan,
Ann Arbor, MI 48109, USA

Received 3 January 2011

En écrivant ce papier, les auteurs ont toujours présente à l'esprit leur amitié pour Michelle Schatzman. Comment la garder vivante sinon en manifestant chaque jour la curiosité, l'exigence scientifique et le plaisir du partage qui étaient les siens.

A systematic analysis of matched layers is undertaken with special attention to better understand the remarkable method of Bérenger. We prove that the Bérenger and closely related layers define well-posed transmission problems in great generality. When the Bérenger method or one of its close relatives is well-posed, perfect matching is proved. The proofs use the energy method, Fourier–Laplace transform, and real coordinate changes for Laplace transformed equations. It is proved that the loss of derivatives associated with the Bérenger method does not occur for elliptic generators. More generally, an essentially necessary and sufficient condition for loss of derivatives in Bérenger's method is proved. The sufficiency relies on the energy method with pseudodifferential multiplier. Amplifying and nonamplifying layers are identified by a geometric optics computation. Among the various flavors of Bérenger's algorithm for Maxwell's equations, our favorite choice leads to a strongly well-posed augmented system and is both perfect and nonamplifying in great generality. We construct by an extrapolation argument an alternative matched layer method which preserves the strong hyperbolicity of the original problem and though not perfectly matched has *leading* reflection coefficient equal to zero at all angles of incidence. Open problems are indicated throughout.

Keywords: PML; WKB; hyperbolic operators; weak well-posedness; geometric optics; extrapolation; reflection; amplification

AMS Subject Classification: 65M12, 65M55, 30E10.

Contents

1. Introduction	160
2. Well-Posed First-Order Cauchy Problems	165

3.	Analysis of the Bérenger’s PML by Energy Methods	176
4.	Analysis of Layers with Only One Absorption by Fourier–Laplace Transform	189
5.	Plane Waves, Geometric Optics, and Amplifying Layers	213
6.	Harmoniously Matched Layers	219

1. Introduction

This paper analyses absorbing layer methods for calculating approximations to the solution, U , of first-order systems of hyperbolic partial differential equations,

$$L(\partial_t, \partial_x)U := \partial_t U + \sum_{l=1}^d A_l \partial_l U = F, \quad (t, x) \in \mathbb{R}^{1+d}, \quad U(t, x) \in \mathbb{C}^N. \quad (1.1)$$

Approximate values are sought on a finite domain. The source term F and/or initial condition is compactly supported in the domain. The absorbing layer strategy surrounds the domain with a layer of finite thickness intended to be absorbing and weakly reflective.

The simplest case is dimension $d = 1$ with computational domain $x_1 < 0$ and absorbing layer in $x_1 > 0$. For the first example consider inhomogeneous initial data and zero right-hand side. The simplest absorbing layers add a lower order term $\sigma \mathbf{1}_{x_1 > 0} C U$ where $\mathbf{1}$ denotes the characteristic function, for example,

$$\partial_t U + \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \partial_1 U + \sigma \mathbf{1}_{x_1 > 0} \begin{pmatrix} 1 & 0 \\ c & b \end{pmatrix} U = 0.$$

To get a feeling for the reflections, consider the solution $U(t, x_1)$ so that,

$$\text{for } t < 0, \quad U = (\delta(x_1 - t), 0).$$

Then

$$U_1 = \delta(x_1 - t)e^{-\sigma x_1}, \quad (\partial_t - \partial_1 + b\sigma \mathbf{1}_{x_1 > 0})U_2 = -\sigma \mathbf{1}_{x_1 > 0} c U_1.$$

If $c \neq 0$, then $\nabla_{t,x_1} U_2$ is discontinuous across the ray $\{x_1 = -t\}$. From the perspective of a numerical method, such a reflected singularity is undesirable.

The reflected singularity from a discontinuous lower order term is weaker than the singularity of the incident wave. For the equation

$$\partial_t U + A_1 \partial_1 U + \sigma \mathbf{1}_{x_1 > 0} C U = 0,$$

if C is diagonal in a basis diagonalizing A_1 , the reflections are avoided. The ease of eliminating reflections for this problem with $d = 1$ is deceptive. No such simple remedy exists in dimensions $d > 1$. For symmetric hyperbolic systems $A_1 = A_1^*$, it is wise to choose $C = C^* \geq 0$ so that the absorption term is dissipative in the $L^2(\mathbb{R}^d)$ norm.

Consider next the wave equation with friction $\partial_{tt}v - \partial_{11}v + 2\sigma\mathbf{1}_{x_1>0}\partial_tv = 0$ written in characteristic coordinates $(U_1, U_2) = (\partial_tv - \partial_1v, \partial_tv + \partial_1v)$ with absorption $B = \sigma C$:

$$\partial_tU + \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \partial_1U + \sigma\mathbf{1}_{x_1>0}CU = 0, \quad C = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}.$$

The absorption matrix C is symmetric and non-negative but does not commute with A_1 . It produces unacceptably strong reflections. The absorption from Israeli and Orszag [16], $\partial_{tt}v - \partial_{11}v + \sigma(\partial_tv + \partial_1v) = 0$, absorbs only rightward waves and corresponds to

$$C = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} = \pi_+(A_1),$$

introducing the notation $\pi_+(A_1)$ for the spectral projector on the eigenspace corresponding to strictly positive eigenvalues of A_1 . The general non-negative symmetric choice commuting with A_1 is a positive multiple of

$$C = \pi_+(A_1) + \nu\pi_-(A_1), \quad \nu \geq 0. \tag{1.2}$$

We call these *smart layers*. They dissipate the L^2 norm. As observed by Israeli and Orszag, the numerical performance of the smart layers is not as good as one would hope. One reduces reflections by choosing $\sigma(x) \geq 0$ vanishing to order $k \geq 0$ at the origin. That reduces the rate of absorption and thereby increases the width of the layer required. The leading reflection by such smart layers of incoming wave packets of amplitude $O(1)$ and wavelength ε is $O(\varepsilon^{k+1})$. The leading reflection is linear in σ . In Sec. 6, we introduce the method of Harmoniously Matched Layers which remove the leading order reflections (at all angles of incidence) by an extrapolation.

Open problem. *Repeated extrapolation further reduces the order of reflection. It is easy to program and it is possible that an optimization could pay dividends.*

Elaborate absorbing layer strategies, like Bérenger’s PML introduce operators related to but often more complicated than the original operator L . The operators in the absorbing layer and in the domain of interest may not be the same. For the case of a layer in $\{x_1 > 0\}$, absorbing layer algorithms solve a transmission problem for an unknown (V, W) where V is a \mathbb{C}^N -valued function on $x_1 < 0$ and W is a function on $x_1 > 0$. The equations in $x_1 > 0$ are chosen to be absorbing and the transmission problem weakly reflective. The ingenious innovation of Bérenger was to realize that the operator R in the layer can differ substantially from L . He increased the number of unknown functions in the layer. So W is \mathbb{C}^M -valued with $M > N$.

The pair (V, W) is determined by a well-posed transmission problem,

$$LV = F \quad \text{on } \mathbb{R}_-^{1+d} := \{(t, x) : x_1 < 0\} \quad RW = 0 \quad \text{on } \mathbb{R}_+^{1+d}, \tag{1.3}$$

with the homogeneous transmission condition

$$(V, W) \in \mathcal{N} \quad \text{on } \{x_1 = 0\}. \tag{1.4}$$

Here $\mathcal{N} \subset \mathbb{C}^N \times \mathbb{C}^M$ is a linear subspace.^a The choice of the hyperbolic operator R and transmission condition \mathcal{N} is made with three goals:

- The transmission problem is well-posed, and not hard to approximate numerically.
- Waves from the left are at most weakly reflected at $x_1 = 0$.
- Waves moving rightward decay rapidly in $x_1 > 0$ so that the layer can be chosen thin.

The criterion for perfection that we adopt is that of Appelo, Hagström and Kreiss [3]. In the case of one absorption, it is formulated as follows.

Definition 1.1. A well-posed transmission problem is *perfectly matched* when for all F supported in $x_1 < 0, t \geq 0$, the solution supported in $t \geq 0$ satisfies $V = U|_{x_1 < 0}$.

We prove in Sec. 4.1.4 that Bérenger’s method with one discontinuous absorption σ_1 is perfect in this sense.

In practice one does not absorb in only one direction and the computational domain is rectangular. We give in Sec. 3.5 a definition with absorptions in more than one direction and a proof of perfection.

The strategy of Bérenger is quite ingenious. For an artificial boundary in two dimensions at $\{x_1 = r\}$ and domain of interest $\{x_1 < r\}$ it consists of two steps. The first is a doubling of the system and the second is insertion of an absorption term in $\{x_1 > r\}$. The doubled system involves the unknown $\tilde{U} := (U^1, U^2) \in \mathbb{C}^N \times \mathbb{C}^N$. When $F = 0$, the doubled equation without dissipation is

$$\partial_t U^j + A_j \partial_j (U^1 + U^2) = 0, \quad j = 1, 2.$$

The system with damping in x_1 changes the $j = 1$ equation to

$$\partial_t U^1 + A_1 \partial_1 (U^1 + U^2) + \sigma(x_1) U^1 = 0, \quad \text{supp } \sigma \subset \{x_1 \geq r\}.$$

Then $U := \sum_j U^j$ satisfies $L(\partial)U = 0$ in $x_1 < r$. In practice it is the restriction of U to $x_1 < r$ that is of interest. There are three distinct ways to view this. One can think of the unknowns as U defined in $x_1 < r$ and \tilde{U} in $x_1 \geq r$ with the transmission condition that $A_1 U = A_1 (U^1 + U^2)$ on $x_1 = r$. One is given initial values of U and takes initial values of \tilde{U} vanishing. This is the most natural choice and the one presented by Bérenger.

From the computational point of view, it is simpler to have the same unknowns throughout. The simplification is greater when one passes from the half space case

^aTransmission conditions which involve derivatives can also be treated. The algorithms of Bérenger and our HML do not require that generality.

to a computational domain equal to a rectangular domain in \mathbb{R}^d . One introduces \tilde{U} everywhere with transmission condition $[A_1(U^1 + U^2)] = 0$ where $[*]$ denotes the jump at $x_1 = r$. The transmission condition is then equivalent to the validity of the differential equation satisfied by \tilde{U} in all of \mathbb{R}^d . When one uses \tilde{U} everywhere, the initial values of \tilde{U} are taken equal to zero outside the computational domain. The initial values are constrained to satisfy $U = \sum_j U^j$ within the computational domain. The choice is otherwise arbitrary. For the case of the doubling above, the choice $U^j(0, x) = U(0, x)/2$ for $j = 1, 2$ is common.

If the domain of interest is $|x_1| \leq r$ one would choose $\sigma > 0$ on $|x_1| > r$ and vanishing for $|x_1| < r$. The transmission condition is $[A_1(U^1 + U^2)] = 0$ with the jump at $x_1 = r$ and also at $x_1 = -r$.

In a rectangular geometry in \mathbb{R}^d introduce $\tilde{U} := (U^1, \dots, U^d)$, where $U^l \in \mathbb{C}^N$ for $1 \leq l \leq d$. Then \tilde{U} with values in \mathbb{C}^{Nd} is required to satisfy (in the case $F = 0$),

$$(\tilde{L}(\partial_t, \partial_x)\tilde{U})_l := \partial_t U^l + A_l \partial_t \left(\sum_{j=1}^d U^j \right) + \sigma_l(x_l) U^l = 0, \quad 1 \leq l \leq d. \quad (1.5)$$

Each absorption coefficient $\sigma_l(x_l) \geq 0$ depends on only one variable. It is strictly positive between the inside rectangle, and a larger outside rectangle. In the layer between the rectangles, the solution is expected to decay. If \tilde{U} solves (1.5), then $U = \sum_{j=1}^d U^j$ solves (1.1) on the set $\{x : \sigma_l(x_l) = 0 \text{ for } 1 \leq l \leq d\}$ including the inner rectangle. In the case considered by Bérenger, the σ were discontinuous and the equations (1.5) are equivalent to transmission problems where on the discontinuity surface of σ_j , one imposes the transmission condition of continuity of $A_j \sum_\ell U^\ell$.

Our first technique is the energy method. In Sec. 3.2 we show that if $(\xi_1, \dots, \xi_d) = 0$ does not meet the characteristic variety of L , then the Bérenger method is well-posed *without loss of derivatives*. This applies in particular to linearized elasticity and suggests that in some ways the Bérenger method is better adapted to that situation than the Maxwell equations for which it was intended. In Sec. 3.3 we give a nontrivial extension of the method of Métrol and Vacus to show that Bérenger’s method for the Maxwell equations in dimension $d = 2$ (respectively $d = 3$) is well-posed provided that $\sigma_j(x_j) \in W^{1,\infty}(\mathbb{R}_{x_j})$ (respectively $\sigma_j(x_j) \in W^{2,\infty}(\mathbb{R}_{x_j})$). The method introduces a norm that is the sum of $L^2(\mathbb{R}_{t,x}^d)$ norms of suitable differential operators $P_\alpha(D)$ applied to U . It has the property that the norm at time t_1 is estimated in terms of the norm at time t_2 . If one introduces the vector of unknowns $P_\alpha(D)U$, this shows that the Bérenger problem becomes strongly well-posed without loss of derivatives. Such transformations are typical of weakly well-posed problems. (See the Dominics’ proof of Theorem 1.1 in §IV.1 of [30].)

When such an estimate is known, we prove sharp finite speed in Sec. 3.4 and perfection in Secs. 3.5 and 3.6, the latter concerned with several variants of the Bérenger strategy. The perfection proof passes by a study of the Laplace transform

on $\{\text{Im } \tau = 0\}$. The transformed problem is conjugated to the problem without absorption by a τ -dependent change of independent variable x , an idea inspired by [10].

Our second method is the Fourier–Laplace method. Bérenger introduced his PML for Maxwell’s equations with σ piecewise constant. Using a computation which resembles plane wave analysis of reflections for problems without lower order terms, Bérenger argued that the layers were perfectly matched for all wave numbers and all angles of incidence. Using variants of the same approach, other closely related PML were constructed afterward. Performance is observed to be enhanced using σ which are not discontinuous. Twice differentiable cubic functions are the most common. The Bérenger method is a very good method for Maxwell’s equations. The Fourier–Laplace method gives a framework for understanding the computations of Bérenger. In addition, it is the only method we know for proving well-posedness of Bérenger’s PML with discontinuous σ for Maxwell’s equations.

Plane wave analysis is sufficient to study reflection and transmission for linear constant coefficient operators without lower order terms. Problems with lower order terms require other tools as it is no longer true that the plane waves generate all solutions. The first level of generalization is to use the Fourier–Laplace transform for problems where an absorbing layer occupies $x_1 \geq 0$ and both L and R have constant coefficients. Hersh [11] found necessary and sufficient conditions for (weak) well-posedness of transmission problems. We recall those ideas in Sec. 4.1.1 including the modifications needed for characteristic interfaces, and verify in Sec. 4.1.3 that the condition is satisfied for the Bérenger splitting of general systems with one discontinuous absorption coefficient. To our knowledge this is the first proof that the Bérenger split transmission problem with discontinuous $\sigma(x_1)$ is well-posed.

We give necessary and sufficient conditions for perfection at a planar interface. In Sec. 4.1.4 we verify that the condition is satisfied for the Bérenger splitting.

In Sec. 4.1.5 we prove that in the case of Maxwell’s equations (and not in general) the perfection criterion follows by analytic continuation from the plane wave identities established by Bérenger.

In Sec. 4.2 we prove using the Fourier–Laplace method that Bérenger’s method with one coefficient $\sigma_1(x_1) \in \text{Lip}(\mathbb{R}_{x_1})$ is well-posed and perfectly matched. In our use of the Fourier–Laplace method, including this one, a central role is played by the Seidenberg–Tarski theorem estimating the asymptotic behavior of functions defined by real polynomial equations and inequalities. The Fourier–Laplace method is limited to coefficients that depend only on x_1 .

Our third method of analysis is to study the behavior of short wavelength asymptotic solutions. For such solutions we examine in Sec. 5 the decay in the absorbing layers, and reflections at discontinuities of $\sigma_j(x_j)$ or its derivatives when smoother transitions are used. For problems other than Maxwell, Hu [14] and Bécache, Fauqueux, and Joly [5] have already shown that the supposedly absorbing layers may in fact lead to growth. The study of short wavelength solutions in the

layer yields precise and clear criteria, also valid for variable coefficients, explaining the phenomenon.

The analysis of the reflection of short wavelength wave packets at the interface with the layer also leads us to propose in Sec. 6, a new absorbing layer strategy which we call Harmoniously Matched Layers. The method starts with a smart layer for a symmetric hyperbolic system. Then for wavelength ε , asymptotic solutions of amplitude $O(1)$ and discontinuous σ , the leading order reflected wave at non-normal incidence typically has amplitude proportional to $\sigma\varepsilon$. Therefore an extrapolation using computations with two values of σ eliminates the reflections proportional to σ . This yields a method with leading order reflection $O(\varepsilon^2)$ at all angles of incidence. The resulting method inherits the simple L^2 estimates of the symmetric systems. More generally if the first discontinuous derivative of the absorption coefficient is the k th, then the reflection is $O([D^k\sigma]\varepsilon^{k+1})$ and the same extrapolation removes the leading order reflection. In Sec. 6.4 we investigate several implementations of this idea and show that the method with cubic σ is competitive with that of Bérenger with the same σ . On short wavelengths or random data, it performs better than the Bérenger method. On long wavelengths, Bérenger performs better.

Though we provide satisfactory answers to a wide range of questions about absorbing layers, there is a notable gap.

Open problem. *For the original strategy of Bérenger for Maxwell’s equations with discontinuous absorptions in more than one direction we do not know if the resulting problem is well-posed.*

Discussion. (1) In Sec. 3.6 we prove well-posedness and perfection for a closely related method. (2) In practice discontinuous σ have been abandoned, but it is striking that this problem remains open. (3) Once well-posedness is proved, perfection follows by the proof in Sec. 3.5.

2. Well-Posed First-Order Cauchy Problems

2.1. Basic definitions

Consider a first-order system of partial differential equations for \mathbb{C}^N valued functions in \mathbb{R}^{1+d} ,

$$\mathcal{L}(x, \partial_t, \partial_x)U := \partial_t U + \sum_{l=1}^d \mathcal{A}_l \partial_l U + \mathcal{B}(x)U = 0. \tag{2.1}$$

The principal part of \mathcal{L} , denoted \mathcal{L}_1 ,

$$\mathcal{L}_1(\partial_t, \partial_x) := \partial_t + \sum_{l=1}^d \mathcal{A}_l \partial_l,$$

has constant matrix coefficients \mathcal{A}_l . In the Bérenger strategy, the operators \tilde{L} are the centerpieces and they differ from L . It is for this reason that we introduce \mathcal{L} that can be L or \tilde{L} .

Definition 2.1. The *characteristic variety* $\text{Char}(\mathcal{L}) \subset \mathbb{C}^{1+d} \setminus \{0\}$ of \mathcal{L} is the set of (τ, ξ) such that $\det \mathcal{L}_1(\tau, \xi) = 0$.

Definition 2.2. The *smooth variety hypothesis* is satisfied at $(\underline{\tau}, \underline{\xi}) \in \text{Char}(\mathcal{L})$ if there is a conic neighborhood Ω of $\underline{\xi} \in \mathbb{R}^d \setminus \{0\}$ and a C^∞ function $\xi \mapsto \tau(\xi)$ on Ω so that on a neighborhood of $(\underline{\tau}, \underline{\xi})$, the characteristic variety has equation $\tau = \tau(\xi)$. At such a point the associated *group velocity* is defined to be $\mathbf{v} := -\nabla_\xi \tau(\xi)$.

Example 2.1. This hypothesis holds if and only if for ξ near $\underline{\xi}$ the spectrum of $\mathcal{L}(0, \xi)$ near $-\underline{\tau}$ consists of a single point with multiplicity independent of ξ . For the polynomial $(\tau + \xi_1)(\tau^2 - |\xi|^2)$ with $d > 1$ the hypothesis fails at and only at $\tau + \xi_1 = 0$ where two sheets of the variety are tangent. Replacing the first factor by $\tau + c\xi_1$ with $c > 1$ the hypothesis fails where the two sheets cross transversally. For $0 \leq c < 1$, the hypothesis holds everywhere.

The Cauchy problem for \mathcal{L} is to find a solution U defined on $[0, \infty[\times \mathbb{R}^d$ satisfying (2.1) with prescribed initial data $U(0, \cdot)$.

Definition 2.3. The Cauchy problem for \mathcal{L} is *weakly well-posed* if there exist $q > 0$, $K > 0$, and $\alpha \in \mathbb{R}$ so that for any initial values in $H^q(\mathbb{R}^d)$, there is a unique solution $U \in \mathcal{C}^0([0, +\infty[; L^2(\mathbb{R}^d))$ with

$$\forall t \geq 0, \quad \|U(t, \cdot)\|_{L^2(\mathbb{R}^d)} \leq Ke^{\alpha t} \|U(0, \cdot)\|_{H^q(\mathbb{R}^d)}. \tag{2.2}$$

When the conclusion holds with $q = 0$, the Cauchy problem is called *strongly well-posed*.

Theorem 2.1. (i) *The Cauchy problem for \mathcal{L}_1 is weakly well-posed if and only if for each $\xi \in \mathbb{R}^d$, the eigenvalues of $\mathcal{L}_1(0, \xi)$ are real.*

(ii) *The Cauchy problem for \mathcal{L}_1 is strongly well-posed if and only if for each $\xi \in \mathbb{R}^d$, the eigenvalues of $\mathcal{L}_1(0, \xi)$ are real and $\mathcal{L}_1(0, \xi)$ is uniformly diagonalisable, there is an invertible $S(\xi)$ satisfying,*

$$S(\xi)^{-1} \mathcal{L}_1(0, \xi) S(\xi) = \text{diagonal}, \quad S, S^{-1} \in L^\infty(\mathbb{R}_\xi^d).$$

(iii) *If \mathcal{B} has constant coefficients, then the Cauchy problem for \mathcal{L} is weakly well-posed if and only if there exists $M \geq 0$ such that for any $\xi \in \mathbb{R}^d$, $\det \mathcal{L}(\tau, \xi) = 0 \Rightarrow |\text{Im } \tau| \leq M$.*

Remark 2.1. (1) The algebraic conditions in (i) and (iii) express *weak hyperbolicity*, in the sense of Gårding. The necessity of uniform diagonalizability in (ii) expressing *strong hyperbolicity* is due to Kreiss [8, 19].

(2) An application of Grönwall’s inequality shows that if \mathcal{L}_1 satisfies the condition of Theorem 2.1(ii), then for all $B(x) \in L^\infty(\mathbb{R}^d; \text{Hom}(\mathbb{C}^N))$, the Cauchy problem for $\mathcal{L}_1 + B$ is strongly well-posed.

(3) By property (ii), if \mathcal{L} is strongly hyperbolic, then every eigenvalue $-\tau$ of $\mathcal{L}_1(0, \xi)$ is semi-simple. Equivalently, for any $(\tau, \xi) \in \text{Char}(\mathcal{L})$ the eigenvalue 0 of $\mathcal{L}_1(\tau, \xi)$ is semi-simple, i.e. its geometric multiplicity is equal to its algebraic multiplicity. It is equivalent to saying that $\text{Ker } \mathcal{L}_1(\tau, \xi) = \text{Ker}(\mathcal{L}_1(\tau, \xi))^2$, or that $\mathbb{C}^N = \text{Ker } \mathcal{L}_1(\tau, \xi) \oplus \text{Range } \mathcal{L}_1(\tau, \xi)$.

2.2. Characteristic variety and projectors for Bérenger’s \tilde{L}

To study the Cauchy problem for Bérenger’s split operators \tilde{L} , one starts with a study of the characteristic variety. The coefficients of Bérenger’s operator \tilde{L} are the $dN \times dN$ matrices,

$$\tilde{A}_l := \begin{pmatrix} 0 & \cdots & \cdots & \cdots & 0 \\ \vdots & & & & \vdots \\ A_l & \cdots & \cdots & \cdots & A_l \\ \vdots & & & & \vdots \\ 0 & \cdots & \cdots & \cdots & 0 \end{pmatrix}, \quad B(x) := \begin{pmatrix} \sigma_1(x_1)I_N & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_d(x_d)I_N \end{pmatrix}. \quad (2.3)$$

The principal symbol of \tilde{L} is

$$\tilde{L}_1(\tau, \xi) = \begin{pmatrix} \xi_1 A_1 + \tau I_N & \xi_1 A_1 & \cdots & \xi_1 A_1 \\ \xi_2 A_2 & \xi_2 A_2 + \tau I_N & \cdots & \xi_2 A_2 \\ \vdots & \vdots & \ddots & \vdots \\ \xi_d A_d & \xi_d A_d & \cdots & \xi_d A_d + \tau I_N \end{pmatrix}.$$

Theorem 2.2. (i) *The characteristic polynomial of \tilde{L} is*

$$\det \tilde{L}_1(\tau, \xi) = \tau^{N(d-1)} \det L(\tau, \xi). \quad (2.4)$$

The polynomial associated to the full symbol including the absorption is

$$\det \tilde{L}(\tau, \xi) = \det L \left(\prod_{j=1}^d (\tau + \sigma_j), \xi_1 \prod_{j \neq 1} (\tau + \sigma_j), \xi_2 \prod_{j \neq 2} (\tau + \sigma_j), \dots, \xi_d \prod_{j \neq d} (\tau + \sigma_j) \right). \quad (2.5)$$

If $(\tau, \xi) \in \text{Char } L$ with $\tau \neq 0$, the following properties hold.

(ii) *The mapping*

$$\mathcal{S} : \tilde{\Phi} = (\Phi_1, \dots, \Phi_d) \mapsto - \sum_{j=1}^d \Phi_j$$

is a linear bijection from $\text{Ker } \tilde{L}_1(\tau, \xi)$ onto $\text{Ker } L_1(\tau, \xi)$. Its inverse is given by

$$\Phi \mapsto \left(\frac{\xi_1}{\tau} A_1 \Phi, \dots, \frac{\xi_d}{\tau} A_d \Phi \right).$$

(iii) The kernel of the adjoint $\tilde{L}_1(\tau, \xi)^*$ is equal to the set of vectors $\tilde{\Phi} = (\Phi, \dots, \Phi)$ such that $\Phi \in \text{Ker } L_1(\tau, \xi)^*$. The range of $\tilde{L}_1(\tau, \xi)$ is equal to the set of vectors $\tilde{\Psi} = (\Psi_1, \dots, \Psi_d)$ such that $(\sum_{j=1}^d \Psi_j, \Phi) = 0$ for all $\Phi \in \text{Ker } L_1(\tau, \xi)^*$.

(iv) If moreover the eigenvalue 0 of $L_1(\tau, \xi)$ with $\tau \neq 0$ is semi-simple, the eigenvalue 0 of $\tilde{L}_1(\tau, \xi)$ is semi-simple. Equivalently,

$$\text{Ker } \tilde{L}_1(\tau, \xi) \oplus \text{Range } \tilde{L}_1(\tau, \xi) = \mathbb{C}^{dN}.$$

Proof. (i) Adding the sum of the other rows to the first row in the determinant of $\tilde{L}_1(\tau, \xi)$ yields,

$$\det \tilde{L}_1(\tau, \xi) = \begin{vmatrix} L(\tau, \xi) & \cdots & \cdots & L(\tau, \xi) \\ \xi_2 A_2 & \xi_2 A_2 + \tau I_N & \cdots & \xi_2 A_2 \\ \vdots & & \ddots & \vdots \\ \xi_d A_d & \cdots & \cdots & \xi_d A_d + \tau I_N \end{vmatrix}.$$

Subtracting the first column from the others yields,

$$\det \tilde{L}_1(\tau, \xi) = \begin{vmatrix} L(\tau, \xi) & 0 & \cdots & 0 \\ \xi_2 A_2 & \tau I_N & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ \xi_d A_d & 0 & \cdots & \tau I_N \end{vmatrix}.$$

The first result follows. For the second, write

$$\det \tilde{L}(\tau, \xi) = \begin{vmatrix} \xi_1 A_1 + (\tau + \sigma_1) I_N & \xi_1 A_1 & \cdots & \xi_1 A_1 \\ \xi_2 A_2 & \xi_2 A_2 + (\tau + \sigma_2) I_N & \cdots & \xi_2 A_2 \\ \vdots & & \ddots & \vdots \\ \xi_d A_d & \cdots & \cdots & \xi_d A_d + (\tau + \sigma_d) I_N \end{vmatrix}.$$

For each i , divide the i th row by $\tau + \sigma_i$ to find,

$$\det \tilde{L}(\tau, \xi) = \prod_{j=1}^d (\tau + \sigma_j)^N \det \tilde{L}_1 \left(1, \frac{\xi_1}{\tau + \xi_1}, \dots, \frac{\xi_d}{\tau + \xi_d} \right).$$

By formula (2.4) this implies

$$\det \tilde{L}(\tau, \xi) = \prod_{j=1}^d (\tau + \sigma_j)^N \det L \left(1, \frac{\xi_1}{\tau + \xi_1}, \dots, \frac{\xi_d}{\tau + \xi_d} \right),$$

which is equivalent to (2.5).

(ii) Suppose that $0 \neq \tilde{\Phi} = (\Phi_1, \dots, \Phi_d) \in \text{Ker } \tilde{L}_1(\tau, \xi)$. Then, for any l ,

$$\tau\Phi_l + \xi_l A_l \sum_{j=1}^d \Phi_j = 0. \tag{2.6}$$

Add to find

$$L_1(\tau, \xi) \sum_{j=1}^d \Phi_j = 0.$$

Therefore the map $\tilde{\Phi} \mapsto -\sum_j \Phi_j$ maps $\text{Ker } \tilde{L}_1(\tau, \xi)$ to $\text{Ker } L_1(\tau, \xi)$.

If $\sum_{j=1}^d \Phi_j = 0$, Eq. (2.6) implies that all the Φ_j vanish since $\tau \neq 0$. Therefore the mapping is injective.

Let $\Phi \in \text{Ker } L_1(0, \xi)$. Define

$$\Phi_j = \frac{\xi_j}{\tau} A_j \Phi. \tag{2.7}$$

This defines an element $\tilde{\Phi} = (\Phi_1, \dots, \Phi_d)$ in $\text{Ker } \tilde{L}_1(0, \xi)$ with $\mathcal{S}\tilde{\Phi} = \Phi$, so the mapping is surjective with inverse given by (2.7).

(iii) Since $\tilde{L}_1(\tau, \xi)^* \tilde{\Phi} = (L_1(\tau, \xi)^* \Phi, \dots, L_1(\tau, \xi)^* \Phi)$ it follows that the set of $\tilde{\Phi}$ is included in the kernel. Since the matrices are square, $\text{Ker } L_1(\tau, \xi)$ and $\text{Ker } L_1(\tau, \xi)^*$ have the same dimension. The set of $\tilde{\Phi}$ has dimension equal to this common dimension which by (ii) is equal to the dimension of $\text{Ker } \tilde{L}_1(\tau, \xi)$ proving that they exhaust the kernel. The last property follows directly from the fact that $\text{Range } \tilde{L}_1(\tau, \xi)$ is the orthogonal of $\text{Ker } L_1(\tau, \xi)^*$.

(iv) It suffices to show that the intersection of these spaces consists of the zero vector. Equivalently, it suffices to show that there is no $\Phi \neq 0$ in $\text{Ker } L_1(\tau, \xi)$ such that

$$\forall \Psi \in \text{Ker } L_1(\tau, \xi)^*, \quad \left(\sum_{j=1}^d \frac{\xi_j}{\tau} A_j \Phi, \Psi \right) = 0.$$

The quantity above is equal to $-(\Phi, \Psi)$, and Φ would belong to $(\text{Ker } L_1(\tau, \xi)^*)^\perp = \text{Range } L_1(\tau, \xi)$. Since $\tau \neq 0$ and $\text{Ker } L_1(\tau, \xi) \cap \text{Range } L_1(\tau, \xi) = 0$, this would imply that $\Phi = 0$, leading to a contradiction. □

Denote by $\Pi_L(\tau, \xi)$ (respectively $\Pi_{\tilde{L}}(\tau, \xi)$), the spectral projector onto the kernel of $L_1(\tau, \xi)$ (respectively $\tilde{L}_1(\tau, \xi)$) along its range. For L it is given by

$$\Pi_L(\tau, \xi) = \frac{1}{2\pi i} \oint_{|z|=\rho} (zI - L_1(\tau, \xi))^{-1} dz$$

with ρ small. Like the characteristic variety, Π_L depends only on the principal symbol L_1 . It is characterized by,

$$\Pi_L^2 = \Pi_L, \quad \Pi_L L_1(\tau, \xi) = 0, \quad L_1(\tau, \xi) \Pi_L = 0, \quad \text{rank } \Pi_L = \dim \text{Ker } L_1(\tau, \xi), \tag{2.8}$$

where the τ, ξ dependence of Π_L is suppressed for ease of reading. The first three conditions assert that $\Pi_L(\tau, \xi)$ is a projector annihilating $\text{Range } L_1(\tau, \xi)$ and projecting onto a subspace of $\text{Ker } L_1(\tau, \xi)$. That it maps onto the kernel is implied by the last equality.

Proposition 2.1. *The matrix $\Pi_{\tilde{L}}(\tau, \xi)$ is given by*

$$\Pi_{\tilde{L}}(\tau, \xi) = - \begin{pmatrix} \frac{\xi_1 A_1}{\tau} \Pi_L(\tau, \xi) & \cdots & \frac{\xi_1 A_1}{\tau} \Pi_L(\tau, \xi) \\ \vdots & & \vdots \\ \frac{\xi_d A_d}{\tau} \Pi_L(\tau, \xi) & \cdots & \frac{\xi_d A_d}{\tau} \Pi_L(\tau, \xi) \end{pmatrix}.$$

Proof. Call the matrix on the right $M(\tau, \xi)$. The properties of the projectors associated to L yield formulas for the (i, j) block of the products

$$\begin{aligned} (M(\tau, \xi) \tilde{L}_1(\tau, \xi))_{i,j} &= -\frac{\xi_i A_i}{\tau} \Pi_L L_1 = 0, \\ (\tilde{L}_1(\tau, \xi) M(\tau, \xi))_{i,j} &= -\frac{\xi_i A_i}{\tau} L_1 \Pi_L L_1 = 0 \quad \text{and} \\ (M(\tau, \xi) M(\tau, \xi))_{i,j} &= \frac{\xi_i A_i}{\tau^2} \Pi_L (L_1 - \tau I) \Pi_L = -\frac{\xi_i A_i}{\tau} \Pi_L^2 = (M(\tau, \xi))_{i,j}. \end{aligned}$$

This proves the first three equalities of (2.8). Since M projects onto a subspace of $\text{Ker } \tilde{L}_1$, $\text{rank } M \leq \dim \text{Ker } \tilde{L}_1$. Apply M to a vector $(\Psi, 0, \dots, 0)$ and compare with part (ii) of Theorem 2.2 to see that the range of M contains $\text{Ker } \tilde{L}_1(\tau, \xi)$, so $\text{rank } M \geq \dim \text{Ker } \tilde{L}_1$. This proves the last equality of (2.8). \square

Remark 2.2. (1) The characteristic varieties of L and \tilde{L} are identical in $\tau \neq 0$.
 (2) In particular, the smooth variety hypothesis is satisfied at (τ, ξ) with $\tau \neq 0$ for one if and only if it holds for both, and the varieties have the same equations and the same group velocities.
 (3) When the smooth variety hypothesis is satisfied, the spectral projection $\Pi_{\tilde{L}}(\tau(\underline{\xi}), \underline{\xi})$ is analytic in $\underline{\xi}$, hence of constant rank. It follows that 0 is a semi-simple eigenvalue of $\tilde{L}(\tau(\underline{\xi}), \underline{\xi})$ on a conic neighborhood of $\underline{\xi}$.

If the eigenvalue 0 of $L_1(\tau, \xi)$ is semi-simple, the kernel and the range of $L_1(\tau, \xi)$ are complementary subspaces as mentioned in Remark 2.1, (3), and the partial inverse $Q_L(\tau, \xi)$ of $L_1(\tau, \xi)$ is uniquely determined by

$$Q_L(\tau, \xi) \Pi_L(\tau, \xi) = 0, \quad Q_L(\tau, \xi) L_1(\tau, \xi) = I - \Pi_L(\tau, \xi). \tag{2.9}$$

The partial inverse $Q_{\tilde{L}}(\tau, \xi)$ is defined in the same way from $\tilde{L}_1(\tau, \xi)$.

2.3. The Cauchy problem for Bérenger’s split operators

Part (i) of Theorem 2.2 proves the following:

Corollary 2.1. *If the Cauchy problem for L is weakly well-posed, then so is the Cauchy problem for the principal part \tilde{L}_1 .*

An important observation is that though the Cauchy problem for \tilde{L}_1 is at least weakly well-posed, the root $\tau = 0$ is for all ξ a multiple root. When there are such multiple roots, it is possible that order zero perturbations of \tilde{L}_1 may lead to ill-posed Cauchy problems. The next example shows that this phenomenon occurs for the Bérenger split operators with constant absorption σ_j . Theorem 2.4 shows that when $\tau = 0$ is a root of constant multiplicity of $\det L_1(\tau, \xi) = 0$, the constant coefficient Bérenger operators have well-posed Cauchy problems. Cases where the problem are strongly well-posed are identified. In the latter cases, arbitrary bounded zeroth-order perturbations do not destroy the strong well-posedness.

Example 2.2. (1) For $L := \partial_t + \partial_1 + \partial_2$, $\det L(\tau, \xi) = \tau + \xi_1 + \xi_2$. Therefore $\tau = 0$ is a root if and only if $\xi_1 + \xi_2 = 0$. The doubled system with absorption $\sigma = 1$ in x_1 is

$$\begin{aligned} \tilde{L} &:= \partial_t + \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix} \partial_1 + \begin{pmatrix} 0 & 0 \\ 1 & 1 \end{pmatrix} \partial_2 + \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \quad \text{and} \\ \det \tilde{L}(\tau, \xi, -\xi) &= \det \begin{pmatrix} \tau + \xi + 1 & \xi \\ -\xi & \tau - \xi \end{pmatrix} \\ &= (\tau + \xi + 1)(\tau - \xi) + \xi^2 = \tau^2 + \tau - \xi. \end{aligned}$$

The roots of $\det \tilde{L}(\tau, \xi, -\xi) = 0$ are $\tau = (-1 \pm \sqrt{1 + 4\xi})/2$. Taking $\xi \rightarrow -\infty$ shows that the Cauchy problem for \tilde{L} is not weakly well-posed by Theorem 2.1(iii).

(2) More generally if $\tau + \xi_1 + \xi_2$ is a factor of $\det L_1(\tau, \xi)$, then for $\sigma \neq 0$ the operator \tilde{L} is not even weakly hyperbolic. In this case (2.5) implies that $\tau^2 + \tau - \xi$ is a factor of $\det \tilde{L}(\tau, \xi, -\xi)$.

(3) This generalizes to linear hyperbolic factors in arbitrary dimension.

A key tool is the following special case of Theorem A.2.5 in [13].

Theorem 2.3. (Seidenberg–Tarski theorem) *If $Q(\rho, \zeta)$, $R(\rho, \zeta)$, and $S(\rho, \zeta)$ are polynomials with real coefficients in the $n+1$ real variables $(\rho, \zeta_1, \dots, \zeta_n)$ and the set*

$$M(\rho) := \{ \zeta : R(\rho, \zeta) = 0, S(\rho, \zeta) \leq 0 \}$$

is nonempty when ρ is sufficiently large, define

$$\mu(\rho) := \sup_{\zeta \in M(\rho)} Q(\rho, \zeta).$$

Then either $\mu(\rho) = +\infty$ for ρ large, or there are $a \in \mathbb{Q}$ and $A \neq 0$ so that

$$\mu(\rho) = A\rho^a(1 + o(1)), \quad \rho \rightarrow \infty.$$

Theorem 2.4. *Suppose that $\tau = 0$ is an isolated root of constant multiplicity m of $\det L_1(\tau, \xi) = 0$.*

- (i) *If the Cauchy problem for L_1 is strongly well-posed, then for arbitrary constant absorptions $\sigma_j \in \mathbb{C}$, the Cauchy problem for $\tilde{L}_1 + B$ is weakly well-posed.*
- (ii) *If the Cauchy problem for L_1 is strongly well-posed, and if there is a $\xi \neq 0$ such that $\text{Ker } L(0, \xi) \neq \bigcap_{\xi_j \neq 0} \text{Ker } A_j$, then $\tilde{L}_1(0, \xi)$ is not diagonalizable. Therefore the Cauchy problem for \tilde{L} is not strongly well-posed.*
- (iii) *If the Cauchy problem for L is strongly well-posed and for all ξ , $\text{Ker } L_1(0, \xi) = \bigcap_{\xi_j \neq 0} \text{Ker } A_j$, then the Cauchy problem for \tilde{L} is strongly well-posed. This condition holds if $L_1(0, \partial_x)$ is elliptic, that is $\det L_1(0, \xi) \neq 0$ for all real ξ .*

Remark 2.3. (1) Part (i) is a generalization of results in [15] and Theorem 1 in [5]. In the latter paper, Bécache *et al.* treated the case $N = 2$ assuming that the nonzero eigenvalues of $L_1(0, \xi)$ are of multiplicity one. They conjectured that the result was true more generally. Like them we treat the roots near zero differently from those that are far from zero. The treatment of each of these cases is different from theirs. The tricky part is the roots near zero. We replace their use of Puiseux series by the related Seidenberg–Tarski Theorem 2.3.

(2) Abarbanel and Gottlieb [1] proved (ii) in the special case of Maxwell’s equations. The general argument below is simpler and yields a necessary and sufficient condition for loss of derivatives when the eigenvalue 0 of $L(0, \xi)$ is of constant multiplicity.

(3) Part (iii) is new, extending a result in the thesis of S. Petit-Bergez [25].

Proof. (i) For $\xi \in \mathbb{R}^d \setminus 0$, define for $\rho \in \mathbb{R}_+$,

$$E(\rho) := \max\{\text{Im } (\tau) : \det \tilde{L}(\tau, \xi) = 0, \xi \in \mathbb{R}^d, |\xi|^2 = \rho^2\}.$$

Apply the Seidenberg–Tarski Theorem 2.3 with real variables $\rho, \zeta = (\text{Re } \tau, \text{Im } \tau, \xi)$ and polynomials $R(\rho, \zeta) = |\det \tilde{L}(\tau, \xi)|^2 + (|\xi|^2 - \rho^2)^2$, $S = 0$ and $Q(\rho, \zeta) = \text{Im } \tau$. Conclude that there is an $\alpha \neq 0$ and a rational r so that

$$E(\rho) = \alpha\rho^r(1 + o(1)), \quad \rho \rightarrow \infty.$$

To prove the result it suffices to prove that $\text{Im } \tau$ is bounded, i.e. to show that $r \leq 0$. Suppose on the contrary that $r > 0$.

Given τ, ξ define $k \in S^{d-1}, \rho \in \mathbb{R}_+$ and θ by

$$k := \frac{\xi}{|\xi|}, \quad \xi = \rho k, \quad \theta := \frac{\tau}{\rho}.$$

Choose sequences $\tau(n)$ and $\xi(n)$ so that for $n \rightarrow \infty$,

$$\det \tilde{L}(\tau(n), \xi(n)) = 0, \quad \text{Im } \tau(n) = \alpha(\rho(n))^r(1 + o(1)). \tag{2.10}$$

Write

$$\tilde{L}(\tau, \xi) = \tilde{L}_1(\tau, \xi) + B = \rho \left(\tilde{L}_1(\theta, k) + \frac{1}{\rho} B \right) = \rho \left(\theta I_{2N \times 2N} + \tilde{L}_1(0, k) + \frac{1}{\rho} B \right).$$

The matrix $\tilde{L}(\tau, \xi)$ is singular if and only if $-\theta$ is an eigenvalue of $\tilde{L}_1(0, k) + \rho^{-1}B$.

For large ρ this is a small perturbation of $\tilde{L}_1(0, k)$. Choose $\mu > 0$ so that for $|k| = 1$, the only eigenvalue of $\tilde{L}_1(0, k)$ in the disk $|\theta| \leq 2\mu$ is $\theta = 0$.

Because of the strong well-posedness of L , there is a uniformly independent basis of unit eigenvectors for the eigenvalues of $L_1(0, k)$ in $|\theta| \geq \mu$. By Theorem 2.2(iv) there is a uniformly independent basis of unit eigenvectors for the eigenvalues of $\tilde{L}_1(0, k)$ in $|\theta| \geq \mu$.

It follows that there is a C_0 so that for $\rho > C_0$ the eigenvalues of $\tilde{L}_1(0, k) + \rho^{-1}B$ in $|\theta| > \mu$ differ from the corresponding eigenvalues of $\tilde{L}_1(0, k)$ by no more than C_0/ρ . In particular, their imaginary parts are not larger than C_0/ρ . Therefore, the corresponding eigenvalues $\tau = \rho\theta$ have bounded imaginary parts. Thus for n large, $E(\rho(n))$ can be reached only for the eigenvalues $-\theta(n)$ which are perturbations of the eigenvalue 0 of $\tilde{L}_1(0, k(n))$.

Perturbation by $O(1/\rho)$ of the uniformly bounded family of $dN \times dN$ matrices, $\tilde{L}_1(0, k)$, can move the eigenvalues by no more than $O(\rho^{-\frac{1}{dN}})$. Since the unperturbed eigenvalue is 0, $|\theta(n)| \leq C\rho(n)^{-1/dN}$, so

$$|\tau(n)| \leq C\rho(n)^{1-\frac{1}{dN}}, \quad \text{Im } \tau(n) = \alpha\rho(n)^r(1 + o(1)), \quad \alpha \neq 0.$$

Therefore $r \leq 1 - 1/dN < 1$ and

$$\prod_{j=1}^d (\tau(n) + \sigma_j) = \tau(n)^d(1 + o(1)), \quad \xi_\ell(n) \prod_{j \neq \ell} (\tau(n) + \sigma_j) = \xi_\ell(n)\tau(n)^{d-1}(1 + o(1)).$$

Insert in identity (2.5) to find

$$\det L_1(\tau(n)^d(1 + o(1)), \xi(n)\tau(n)^{d-1}(1 + o(1))) = 0.$$

Divide the argument by $\rho(n)\tau(n)^{d-1}$ and use homogeneity to find

$$\det L_1 \left(\frac{\tau(n)}{\rho(n)}(1 + o(1)), k(n)(1 + o(1)) \right) = 0.$$

The constant multiplicity hypothesis shows that

$$\det L_1(\tau, \xi) = \tau^m F_1(\tau, \xi) \quad \text{and} \quad \forall \xi \in \mathbb{R}^d, \quad F_1(0, \xi) \neq 0. \tag{2.11}$$

Since for n large $(\tau(n)/\rho(n))(1 + o(1)) \neq 0$ we have

$$F_1 \left(\frac{\tau(n)}{\rho(n)}(1 + o(1)), k(n)(1 + o(1)) \right) = 0.$$

Passing to a subsequence, we may assume that the bounded sequence $k(n) \rightarrow k$. In addition, $\tau(n)/\rho(n) \rightarrow 0$ so passing to the limit yields $F_1(0, k) = 0$ contradicting (2.11). This contradiction proves (i).

(ii) Theorem 2.2(i) shows that 0 is an eigenvalue of $\tilde{L}_1(0, \xi)$ with algebraic multiplicity equal to $N(d - 1) + m$. It remains to see that with the assumption, the dimension of $\text{Ker } \tilde{L}_1(0, \xi)$ is strictly smaller than $N(d - 1) + m$. By definition

$$\text{Ker } \tilde{L}_1(0, \xi) = \left\{ \tilde{\Phi} = (\Phi_1, \dots, \Phi_d) : \sum_{j=1}^d \Phi_j \in \cap_p \text{Ker } (\xi_p A_p) \right\}.$$

Define

$$\mathcal{E}_1 := \left\{ \tilde{\Phi} = (\Phi_1, \dots, \Phi_d) : \sum_{j=1}^d \Phi_j = 0 \right\}.$$

Then $\mathcal{E}_1 \subset \text{Ker } \tilde{L}_1(0, \xi)$ and $\dim \mathcal{E}_1 = N(d - 1)$.

Define

$$\mathcal{E}_2 := \text{Ker } L_1(0, \xi) \otimes O^{d-1}, \quad \dim \mathcal{E}_2 = m.$$

If $\tilde{\Phi} \in \text{Ker } \tilde{L}_1(0, \xi)$, $\sum_{j=1}^d \Phi_j \in \text{Ker } L_1(0, \xi)$, write

$$\tilde{\Phi} = \left(\sum \Phi_j, 0, \dots, 0 \right) - W, \quad \left(\sum \Phi_j, 0, \dots, 0 \right) \in \mathcal{E}_2, \quad W \in \mathcal{E}_1.$$

Thus, $\text{Ker } \tilde{L}_1(0, \xi) \subset \mathcal{E}_1 \oplus \mathcal{E}_2$.

Pick V in $\text{Ker } L_1(0, \xi)$, but not in $\cap_{\xi_j \neq 0} \text{Ker } A_j$. Then

$$\tilde{V} = (V, 0, \dots, 0) \in \mathcal{E}_2 \quad \text{and} \quad V \notin \text{Ker } \tilde{L}_1(0, \xi).$$

This proves that $\text{Ker } \tilde{L}_1(0, \xi)$ is a proper subset of $\mathcal{E}_1 \oplus \mathcal{E}_2$, so

$$\dim(\text{Ker } \tilde{L}_1(0, \xi)) < \dim \mathcal{E}_1 + \dim \mathcal{E}_2 = N(d - 1) + m.$$

Thus the geometric multiplicity of the eigenvalue 0 is strictly less than its algebraic multiplicity. Therefore, $\tilde{L}_1(0, \xi)$ is not diagonalizable. This proves (ii).

(iii) To prove that the split problem is strongly well-posed, it suffices to consider the principal part. Suppose $L(0, \xi)$ is uniformly diagonalizable on a conic neighborhood of $\underline{\xi} \in \mathbb{R}^N \setminus 0$. For $\tilde{U} = (U^1, \dots, U^d)$, introduce

$$\tilde{V} = (V^1, \dots, V^d) \quad \text{with} \quad V^1 := \sum_{j=1}^d U^j \quad \text{and} \quad V^l := U^l \quad \text{for } 2 \leq l \leq d. \quad (2.12)$$

Then

$$\tilde{L}_1(\partial_t, \partial_x)\tilde{U} = 0 \Leftrightarrow \partial_t\tilde{V} + \tilde{Q}(\partial_x)\tilde{V} = 0, \quad \text{with } \tilde{Q}(\xi) := \begin{pmatrix} L_1(0, \xi) & 0 & \cdots & 0 \\ \xi_2 A_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \xi_d A_d & 0 & \cdots & 0 \end{pmatrix}.$$

The eigenvalues of $\tilde{Q}(\xi)$ are those of $L_1(0, \xi)$, therefore real. It suffices to diagonalize uniformly $\tilde{Q}(\xi)$ on the conic neighborhood of $\underline{\xi}$. By homogeneity it suffices to consider ξ with $|\xi| = 1$. By hypothesis there exist a real diagonal matrix $D(\xi)$ and an invertible matrix $S(\xi)$ so that

$$L_1(0, \xi) = S(\xi)D(\xi)S^{-1}(\xi) \quad \text{and} \quad \exists K > 0, \quad \forall \xi \in \mathbb{R}^d, \quad \|S(\xi)\| + \|S^{-1}(\xi)\| \leq K.$$

Seek a diagonalization of $\tilde{Q}(\xi)$ on $|\xi| = 1$ in the form,

$$\tilde{S}(\xi) = \begin{pmatrix} S(\xi) & 0 & \cdots & 0 \\ \xi_2 A_2 Q(\xi) S(\xi) & \text{Id} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \xi_d A_d Q(\xi) S(\xi) & 0 & \cdots & \text{Id} \end{pmatrix} \text{ so } (\tilde{S}(\xi))^{-1} = \begin{pmatrix} (S(\xi))^{-1} & 0 & \cdots & 0 \\ -\xi_2 A_2 Q(\xi) & \text{Id} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -\xi_d A_d Q(\xi) & 0 & \cdots & \text{Id} \end{pmatrix} \tag{2.13}$$

with $Q(\xi)$ to be determined. Then,

$$\tilde{S}^{-1}(\xi)\tilde{Q}(\xi)\tilde{S}(\xi) = \begin{pmatrix} D(\xi) & 0 & \cdots & 0 \\ \xi_2 A_2 (I - Q(\xi)L_1(0, \xi))S(\xi) & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \xi_d A_d (I - Q(\xi)L_1(0, \xi))S(\xi) & 0 & \cdots & 0 \end{pmatrix}$$

and

$$\tilde{S}^{-1}(\xi)\tilde{Q}(\xi)\tilde{S}(\xi) \text{ is diagonal } \Leftrightarrow \xi_j A_j (I - Q(\xi)L_1(0, \xi)) = 0, \quad 2 \leq j \leq d. \tag{2.14}$$

From the strong well-posedness of L_1 , it follows that uniformly in ξ one has $\text{Ker } L_1(0, \xi) \oplus \text{Range } L_1(0, \xi) = \mathbb{C}^N$. Choose Q equal to the left inverse of $L_1(0, \xi)$ defined in (2.9). Since $\text{Ker } L_1(0, \xi) = \cap \text{Ker } \xi_j A_j$, the condition on the right in (2.14) holds so $\tilde{S}(\xi)$ diagonalizes $\tilde{Q}(\xi)$. Since $S(\xi)$ and $S(\xi)^{-1}$ are bounded on a conic neighborhood, it follows that $\tilde{S}(\xi)$ and $\tilde{S}(\xi)^{-1}$ are bounded on a neighborhood of $\underline{\xi}$ in $|\xi| = 1$. A finite cover of the sphere, completes the proof. \square

Remark 2.4. Denote by $\tilde{S}(\xi)$ the function homogeneous of degree zero given by (2.13) for $|\xi| = 1$ with Q constructed in the proof. Then $\|\tilde{S}(D)\tilde{V}(t)\|_{L^2(\mathbb{R}^d)}$ with \tilde{V} from (2.12) is a norm equivalent to $\|\tilde{U}(t)\|_{L^2(\mathbb{R}^d)}$ and is conserved for solutions of $\tilde{L}_1(\partial)\tilde{U} = 0$. Those solutions yield a unitary group with respect to the norm $\|\tilde{S}(D)\tilde{V}\|_{L^2(\mathbb{R}^d)}$.

3. Analysis of the Bérenger’s PML by Energy Methods

This section contains results proving that the initial value problems so defined are well-posed. We begin with the case of Gevrey absorptions, then $W^{2,\infty}$, and finally the case of the Heaviside function.

In Sec. 3.1, we prove that when \tilde{L} is only weakly well-posed, Gevrey regular σ_l lead to well-posed initial value problems in Gevrey classes. Commonly used σ are not this smooth.

The strongest result, from Sec. 3.2, applies when $L(0, \partial_x)$ is elliptic. Important cases are the wave equation and linearized elasticity. In these cases the operator \tilde{L}_1 is strongly hyperbolic so remains strongly hyperbolic even with general bounded zeroth-order perturbations. Thus for bounded $\sigma_l(x_l)$, the initial value problem is strongly well-posed.

In Sec. 3.3, we analyze the case of \tilde{L} associated to Maxwell’s equations with finitely smooth σ . We follow the lead of [25] and extend the analysis of [22] to several absorptions σ_l and to higher dimensions. Related estimates for the linearized Euler equation have been studied by Métivier [21].

The results of this section do not treat the case of Bérenger’s method for Maxwell’s equations with discontinuous σ_j . The case of one absorption is treated in Sec. 4. A closely related method is treated by an energy method in Sec. 3.6.

3.1. General operators and Gevrey absorption

The next result is implied by Bronstein’s theorem [6, 7, 23, 24]. It shows that when $L_1(0, \xi)$ has only real eigenvalues and the σ_j belong to the appropriate Gevrey class, then the Cauchy problem for \tilde{L} is solvable for Gevrey data.

Definition 3.1. For $1 \leq s < \infty, f \in \mathcal{S}'(\mathbb{R}^d)$ belongs to the Gevrey class $G^s(\mathbb{R}^d)$ when

$$\exists C, M, \quad \forall \alpha \in \mathbb{N}^d, \quad \|\partial^\alpha f\|_{L^2(\mathbb{R}^d)} \leq M\alpha!C^{|\alpha|}.$$

Then $G^s \subset \cap_\sigma H^\sigma(\mathbb{R}^d) \subset C^\infty(\mathbb{R}^d)$. For $s > 1$, the compactly supported elements of G^s are dense. If $|\hat{f}(\xi)| \leq Ce^{-|\xi|^a}$ with $0 < a < 1$, then $u \in G^{1/a}$.

Theorem 3.1. *If the principal part L_1 is weakly hyperbolic and $\sigma_j \in G^{N/N+1}(\mathbb{R}^d)$, then for arbitrary $f \in G^{N/(N+1)}(\mathbb{R}^d)$ there is one and only one solution $u \in C^\infty(\mathbb{R}^{1+d})$ to*

$$\tilde{L}u = 0, \quad u(0, \cdot) = f.$$

The solution depends continuously on f .

3.2. Strong hyperbolicity when $L(0, \partial)$ is elliptic

Theorem 3.2. *If L is strongly well-posed and $L(0, \partial)$ is elliptic, then \tilde{L} is strongly well-posed for any absorption $(\sigma_1(x_1), \dots, \sigma_d(x_d))$ in $(L^\infty(\mathbb{R}))^d$.*

Proof. Kreiss' Theorem 2.1 asserts that an operator with constant coefficient principal part is uniformly well-posed if and only if the principal part is uniformly diagonalizable on a conic neighborhood of each $\underline{\xi} \neq 0$. Therefore the corollary follows from the third part of Theorem 2.4. \square

Example 3.1. This result implies that the PML model for the elastodynamic system is strongly well-posed. The system is written in the velocity-stress (v, Σ) formulation,

$$\rho \partial_t v - \operatorname{div} \Sigma = 0, \quad \partial_t \Sigma - C \varepsilon(v) = 0, \quad \varepsilon_{ij}(v) := (\partial_i v_j + \partial_j v_i),$$

with positive definite elasticity tensor C and $\Sigma := C \varepsilon$. See [5], where the authors showed that such layers may be amplifying (see Sec. 5).

3.3. The method of Métral–Vacus extended to the 3D PML Maxwell system

Métral and Vacus proved in [22] a stability estimate for Bérenger's two-dimensional PML Maxwell system with one absorption $\sigma_1(x_1) \in W^{1,\infty}(\mathbb{R})$ and $x = (x_1, x_2) \in \mathbb{R}^2$. There are two crucial elements in their method. First following Bérenger, they do not split all variables in all directions. This section begins by showing that the partially split model is equivalent to the fully split model restricted to functions \tilde{U} some of whose components vanish. The \tilde{L} evolution leaves this space invariant and its evolution on that subspace determines its behavior everywhere.

The second element is that on the partially split subspace there is an *a priori* estimate bounding the norm at time t by the same norm at time 0. This looks inconsistent with the fact that the Cauchy problem is only weakly well-posed. However, the norm is not homogeneous. Certain linear combinations of components have more derivatives estimated than others. The observation of [22] is that the system satisfied by the fields and certain combinations of the fields and their derivatives, yields a large but symmetrizable first-order system. These estimates have been obtained, and extended in Sabrina Petit's thesis [25] in the 2D case with two coefficients, and in the 3D case for an absorption in only one direction.

In this section, motivated in part by the clarification of the role of symmetrizers in the work of Métivier [21] for the 2D variable coefficient Euler equations in geophysics, we construct analogous more elaborate functionals which suffice for the general case of three absorptions in three dimensions. They require $\sigma_j \in W^{2,\infty}(\mathbb{R})$.

Maxwell's equations for $\partial_t E_1$ and $\partial_t B_1$ contain only partial derivatives with respect to x_2, x_3 and not x_1 . In such a situation Bérenger splits the corresponding equations in directions x_2, x_3 but not in direction x_1 . To see why this is a special case of the general splitting algorithm (1.5) reason as follows. If the equation for $\partial_t U_j$ from L does not contain any terms in ∂_k , that is the j th row of A_k vanishes, then the equation for the j th component of the unknown U^k corresponding to the

splitting for the k th space variable is,

$$\partial_t U_j^k + \sigma_k(x_k)U_j^k = 0, \quad U_j^k = e^{-\sigma_k(x_k)t}U_j^k(0, x). \tag{3.1}$$

Substituting this into the other equations reduces the number of unknowns by one. The simplest strategy is to take initial data $U_j^k(0, x) = 0$ which yields the operator \tilde{L} restricted to the invariant subspace of functions so that $U_j^k = 0$. Conversely if one knows how to solve that restricted system then the full system can be reduced to the restricted system with an extra source term from (3.1).

Summary. To study the fully split system, it is sufficient to study the system restricted to $\{U_j^k = 0\}$. Performing this reduction for each missing spatial derivative, corresponds to splitting equations only along directions containing the corresponding spatial derivatives.

An extreme case of this reduction occurs if an equation contains no spatial derivatives, that unknown is eliminated entirely. For the Maxwell system which is the subject of this section this does not occur. The use of unsplit variables

- reduces the size of \tilde{U} reducing computational cost,
- corresponds to Bérenger’s original algorithm,
- is important for the method of Métral–Vacus which takes advantage of the vanishing components U_k^k .

Consider the 3D Maxwell equations,

$$\partial_t E - \nabla \times H = 0, \quad \partial_t H + \nabla \times E = 0.$$

Defining $U = E + iH$, they take the symmetric hyperbolic form (1.1) with hermitian matrices

$$A_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & -i \\ 0 & i & 0 \end{pmatrix}, \quad A_2 = \begin{pmatrix} 0 & 0 & i \\ 0 & 0 & 0 \\ -i & 0 & 0 \end{pmatrix} \quad \text{and} \quad A_3 = \begin{pmatrix} 0 & -i & 0 \\ i & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}. \tag{3.2}$$

Introduce the splitting (1.5) with some components unsplit. Define the subspace \mathcal{H} with vanishing components corresponding to the unsplit components

$$\mathcal{H} := \{\tilde{U} = (U^1, U^2, U^3) \in H^2(\mathbb{R}^3; \mathbb{C}^3)^3 : U_1^1 = 0, U_2^2 = 0, U_3^3 = 0\}.$$

For $\tilde{U} = (U^1, U^2, U^3)$ in \mathcal{H} , define

$$U := U^1 + U^2 + U^3, \quad V^j := \partial_j U, \quad V^{i,j} := \partial_{ij} U, \quad W := \sum_k \sigma_k(x_k)U^k, \\ W^j := \partial_j W, \quad Z := \sum_k \partial_k (W_k + \sigma_k(x_k)U_k), \quad Z^j := \partial_j Z, \tag{3.3}$$

$$\mathbb{V} := (U, V^i, V^{i,j}, W^j, U^j, W, Z^j) \in \mathbb{C}^{54}.$$

The function Z and therefore Z^j are \mathbb{C} -valued. The other slots in \mathbb{V} are \mathbb{C}^3 -valued. The second derivatives $V^{i,j}$ of U are ordered as $V^{1,1}, V^{2,1}, V^{3,1}, V^{2,2}, V^{3,2}, V^{3,3}$.

This convention is important when the equations for \mathbb{V} are written in matrix form. Computing in turn W^j, Z, Z^j requires two derivatives of σ_j .

The unknown in (1.5) is $\tilde{U} = (U^1, U^2, U^3)$. The U^j appear in the fifth slot of \mathbb{V} . Therefore,

$$\|\mathbb{V}(t, \cdot)\|_{(L^2(\mathbb{R}^3))^{54}} \geq \|\tilde{U}(t, \cdot)\|_{(L^2(\mathbb{R}^3))^9}.$$

For the Cauchy problem the initial data is $\tilde{U}_0 = (U_0^1, U_0^2, U_0^3)$, from which \mathbb{V}_0 is deduced by the derivations above, and

$$\|\mathbb{V}_0\|_{(L^2(\mathbb{R}^3))^{54}} \leq C\|\tilde{U}_0\|_{(H^2(\mathbb{R}^3))^9}.$$

Theorem 3.3. *If σ_j , for $j = 1, 2, 3$, belong to $W^{2,\infty}(\mathbb{R})$, then for any $\tilde{U}_0 = (U_0^1, U_0^2, U_0^3)$ in \mathcal{H} there is a unique solution \tilde{U} in $L^2(0, T; \mathcal{H})$ of the split Cauchy problem (1.5) with initial value \tilde{U}_0 . Furthermore, there is a $C_1 > 0$ independent of \tilde{U}_0 so that for all positive time t ,*

$$\|\tilde{U}(t, \cdot)\|_{(L^2(\mathbb{R}^3))^9} \leq C_1 e^{C_1 t} \|\tilde{U}_0\|_{(H^2(\mathbb{R}^3))^9}. \tag{3.4}$$

Proof. The main step is to derive a system of equations satisfied by $\mathbb{V}(t, x)$ together with a symmetrizer $S(D)$. These imply an estimate for $t \geq 0$,

$$\|\mathbb{V}(t)\|_{L^2(\mathbb{R}^3)} \leq C_2 e^{C_2 t} \|\mathbb{V}(0)\|_{L^2(\mathbb{R}^3)}. \tag{3.5}$$

From this estimate it easily follows that the Cauchy problem for the \mathbb{V} -equations is uniquely solvable. It is true but not immediate that if the initial values of \mathbb{V} are computed from those of \tilde{U} then the solution \mathbb{V} comes from a solution \tilde{U} of the Bérenger system. The strategy has three steps:

- Discretize the Bérenger system in x only.
- Derive an estimate analogous to (3.5) for the semidiscrete problem. The estimate is uniform as the discretization parameter tends to zero. The proof is a semidiscrete analogue of (3.4).
- Solve the semidiscrete problem and pass to the limit to prove the existence.

This is done for the case $d = 2$ in Ref. 25 to which we refer for details. Uniqueness of the solutions to the \mathbb{V} -system and therefore \tilde{U} is simpler and classical and is also in Ref. 25.

Equation (1.5) yields,

$$\partial_t U^j + \sum_k A_k V^j + \sigma_j U^j = 0. \tag{3.6}$$

Summing on j yields

$$\partial_t U + L(0, \partial)U + W = 0. \tag{3.7}$$

Differentiate in direction x_j to find,

$$\partial_t V^j + L(0, \partial)V^j + W^j = 0. \tag{3.8}$$

Differentiate once more to get

$$\partial_t V^{i,j} + L(0, \partial) V^{i,j} + \partial_i W^j = 0. \quad (3.9)$$

The quantity $\partial_i W^j$ on the left is replaced using the next lemma.

Lemma 3.1.

$$\partial_j W = L(0, \partial) A_j W + Z e_j - \sum_k E_{jk} (\sigma'_k U + \sigma_k V^k), \quad (3.10)$$

where e_j is the j th vector of the standard basis, and E_{ij} is the 3×3 matrix all of whose entries vanish except the (i, j) element that is equal to 1.

Proof. First evaluate $L(0, \partial) A_j W$ to find

$$L(0, \partial) A_j W = \sum_k A_k A_j \partial_k W.$$

The matrices in Maxwell's equations satisfy $A_j A_k = -E_{jk}$ for $j \neq k$, and $A_j^2 = \sum_{k \neq j} E_{kk} = I - E_{jj}$. This yields

$$\begin{aligned} L(0, \partial) A_j W &= - \sum_{k \neq j} E_{jk} \partial_k W + (I - E_{jj}) \partial_j W \\ &= \partial_j W - \sum_k E_{jk} \partial_k W = \partial_j W - \operatorname{div}(W) e_j. \end{aligned}$$

Introduce the definition of Z to find

$$\begin{aligned} W^j &:= \partial_j W = L(0, \partial) A_j W + \operatorname{div}(W) e_j \\ &= L(0, \partial) A_j W + Z e_j - \left(\sum_k \partial_k (\sigma_k U_k) \right) e_j. \end{aligned}$$

Compute

$$\begin{aligned} \left(\sum_k \partial_k (\sigma_k U_k) \right) e_j &= \sum_k E_{jk} \partial_k (\sigma_k U) \\ &= \sum_k E_{jk} \sigma'_k U + \sum_k E_{jk} \sigma_k V^k, \end{aligned}$$

which proves (3.10). The proof of the lemma is complete. \square

Differentiate (3.10) in space to obtain

$$\partial_i W^j = L(0, \partial) A_j W^i + \partial_i Z e_j - \sum_k E_{jk} \partial_i (\sigma_k U + \sigma'_k V^k).$$

Inserting into (3.9) yields

$$\partial_t V^{i,j} + L(0, \partial) V^{i,j} + L(0, \partial) A_j W^i + Z^{ij} - \sum_k E_{jk} \partial_i (\sigma_k U + \sigma'_k V^k) = 0.$$

This is equivalent to

$$\begin{aligned} & \partial_t V^{i,j} + L(0, \partial) V^{i,j} + L(0, \partial) A_j W^i + Z^i e_j - \sigma'_i E_{ji} U \\ & - \left(\sigma''_i E_{ji} + \sum_k \sigma_k E_{jk} \right) V^i - \sum_k \sigma'_k E_{jk} V^{i,k} = 0. \end{aligned} \tag{3.11}$$

To close the system it remains to evaluate the time derivatives of W, W^j and $\partial_j Z$.

$$\partial_t W = \sum_k \sigma_k \partial_t U^k = - \sum_k \sigma_k (A_k U + \sigma_k U^k) = - \sum_k \sigma_k A_k V^k - \sum_k \sigma_k^2 U^k.$$

Using the particular form of the equations yields

$$\begin{aligned} \sum_k \sigma_k^2 U^k &= \left(\sum_k \sigma_k \right) \left(\sum_k \sigma_k U^k \right) - \sum_k \sigma_k \left(\sum_{l \neq k} \sigma_l U^l \right) \\ &= \left(\sum_k \sigma_k \right) W - \text{diag}(\sigma_2 \sigma_3, \sigma_1 \sigma_3, \sigma_1 \sigma_2) U - \text{diag}(\sigma_1, \sigma_2, \sigma_3) W. \end{aligned}$$

Therefore

$$\begin{aligned} \partial_t W + \sum_k \sigma_k A_k V^k + \left(\sum_k \sigma_k \right) W - \text{diag}(\sigma_2 \sigma_3, \sigma_1 \sigma_3, \sigma_1 \sigma_2) U \\ - \text{diag}(\sigma_1, \sigma_2, \sigma_3) W = 0. \end{aligned} \tag{3.12}$$

Differentiate in x_i to find

$$\begin{aligned} \partial_t W^i + \sum_k (\partial_i(\sigma_k) A_k V^k + \sigma_k A_k V^{ki}) + \partial_i \left(\sum_k \sigma_k \right) W + \sum_k \sigma_k W^i \\ - \partial_i (\text{diag}(\sigma_2 \sigma_3, \sigma_1 \sigma_3, \sigma_1 \sigma_2)) U - \text{diag}(\sigma_2 \sigma_3, \sigma_1 \sigma_3, \sigma_1 \sigma_2) V^i \\ - \partial_i (\text{diag}(\sigma_1, \sigma_2, \sigma_3)) W - \text{diag}(\sigma_1, \sigma_2, \sigma_3) W^i = 0. \end{aligned} \tag{3.13}$$

Next compute

$$\partial_t Z = \sum_i \partial_i \partial_t (W_i + \sigma_i U_i).$$

Consider the pair of equations

$$\partial_t U_1^2 + i \partial_2 U_3 + \sigma_2 U_1^2 = 0 \quad \text{and} \quad \partial_t U_1^3 - i \partial_3 U_2 + \sigma_3 U_1^3 = 0$$

and add the two equations. Also add σ_2 times the first to σ_3 times the second. This yields two equations,

$$\partial_t U_1 + i(\partial_2 U_3 - \partial_3 U_2) + W_1 = 0 \quad \text{and} \quad \partial_t W_1 + i(\sigma_2 \partial_2 U_3 - \sigma_3 \partial_3 U_2) + \sigma_2^2 U_1^2 + \sigma_3^2 U_1^3 = 0.$$

Rewrite the last term as $\sigma_2^2 U_1^2 + \sigma_3^2 U_1^3 = (\sigma_2 + \sigma_3) W_1 - \sigma_2 \sigma_3 U_1$, to find

$$\partial_t U_1 + i(\partial_2 U_3 - \partial_3 U_2) + W_1 = 0 \quad \text{and}$$

$$\partial_t W_1 + i(\sigma_2 \partial_2 U_3 - \sigma_3 \partial_3 U_2) + (\sigma_2 + \sigma_3) W_1 - \sigma_2 \sigma_3 U_1 = 0.$$

Multiply the first equation by σ_1 and add the second to obtain

$$\partial_t(W_1 + \sigma_1 U_1) + i((\sigma_1 + \sigma_2)\partial_2 U_3 - (\sigma_1 + \sigma_3)\partial_3 U_2) + \left(\sum_k \sigma_k\right)W_1 - \sigma_2\sigma_3 U_1 = 0.$$

The other indices follow by permutation. Differentiate in x_k and add to find

$$\begin{aligned} \sum_k \partial_k \partial_t(W_k + \sigma_k U_k) + i \sum_k \partial_k((\sigma_k + \sigma_{k+1})\partial_{k+1} U_{k+2} - (\sigma_k + \sigma_{k+2})\partial_{k+2} U_{k+1}) \\ + \sum_i \partial_i \left(\left(\sum_k \sigma_k\right)W_i \right) - \sum_k \partial_k(\sigma_{k+1}\sigma_{k+2}U_k) = 0. \end{aligned}$$

The terms with two spatial derivatives cancel. This leaves

$$\begin{aligned} \partial_t Z + i \sum_k \sigma'_k(\partial_{k+1} U_{k+2} - \partial_{k+2} U_{k+1}) \\ + \left(\sum_k \sigma_k\right)\left(\sum_k W_k^k\right) + \sum_k \sigma'_k W_k - \sum_k \sigma_{k+1}\sigma_{k+2}V_k^k = 0. \end{aligned}$$

Since

$$Z = \sum_k \partial_k(W_k + \sigma_k(x_k)U_k) = \sum_k (W_k^k + \sigma'_k U_k + \sigma_k V_k^k),$$

we can replace $(\sum_k \sigma_k)(\sum_k W_k^k)$ in the previous equation by

$$\left(\sum_k \sigma_k\right)\left(Z - \sum_k (\sigma'_k U_k + \sigma_k V_k^k)\right)$$

so

$$\begin{aligned} \partial_t Z + \left(\sum_k \sigma_k\right)Z + i \sum_k \sigma'_k(V_{k+2}^{k+1} - V_{k+1}^{k+2}) - \left(\sum_k \sigma_k\right)\left(\sum_k (\sigma'_k U_k + \sigma_k V_k^k)\right) \\ + \sum_k \sigma'_k W_k - \sum_k \sigma_{k+1}\sigma_{k+2}V_k^k = 0. \end{aligned}$$

Differentiating in x_j yields

$$\begin{aligned} \partial_t Z^j + \left(\sum_k \sigma_k\right)Z^j + \sigma'_j Z + i\sigma''_j(V_{j+2}^{j+1} - V_{j+1}^{j+2}) + i \sum_k \sigma'_k(V_{k+2}^{k+1,j} - V_{k+1}^{k+2,j}) \\ - \sigma'_j \left(\sum_k (\sigma'_k U_k + \sigma_k V_k^k)\right) - \left(\sum_k \sigma_k\right)(\sigma''_j U_j + \sigma'_j V_j^j) \\ - \left(\sum_k \sigma_k\right)\left(\sum_k (\sigma'_k V_k^j + \sigma_k V_{j,k}^k)\right) + \sigma'_j W_j + \sum_k \sigma'_k W_k^j \\ - \sum_k \partial_j(\sigma_{k+1}\sigma_{k+2})V_k^k - \sum_k \partial_j(\sigma_{k+1}\sigma_{k+2})V_k^{j,k} = 0. \end{aligned}$$

Replace Z by $\sum_k (W_k^k + \sigma_k V_k^k)$ to end up with

$$\begin{aligned} & \partial_t Z^j + \left(\sum_k \sigma_k \right) Z^j + \sigma'_j \sum_k (W_k^k + \sigma_k V_k^k) + i \sigma''_j (V_{j+2}^{j+1} - V_{j+1}^{j+2}) \\ & + i \sum_k \sigma'_k (V_{k+2}^{k+1,j} - V_{k+1}^{k+2,j}) - \sigma'_j \left(\sum_k (\sigma'_k U_k + \sigma_k V_k^k) \right) \\ & - \left(\sum_k \sigma_k \right) (\sigma''_j U_j + \sigma'_j V_j^j) - \left(\sum_k \sigma_k \right) \left(\sum_k (\sigma'_k V_k^j + \sigma_k V_{j,k}^k) \right) + \sigma'_j W_j \\ & + \sum_k \sigma'_k W_k^j - \sum_k \partial_j (\sigma_{k+1} \sigma_{k+2}) V_k^k - \sum_k \partial_j (\sigma_{k+1} \sigma_{k+2}) V_k^{j,k} = 0. \end{aligned} \quad (3.14)$$

Summarizing, \mathbb{V} is solution of a first-order system, $\partial_t \mathbb{V} + P(\partial_x) \mathbb{V} + B(x) \mathbb{V} = 0$, whose principal symbol is given by

$$P(\partial) = \begin{pmatrix} I_4 \otimes L(0, \partial) & 0_{4,6} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} & 0_{4,4} \otimes 0_{3,3} \\ 0_{6,4} \otimes 0_{3,3} & I_6 \otimes L(0, \partial) & (I_6 \otimes L(0, \partial))M & 0_{6,3} \otimes 0_{3,3} & 0_{6,4} \otimes 0_{3,3} \\ 0_{3,4} \otimes 0_{3,3} & 0_{3,6} \otimes 0_{3,3} & 0_{3,3} \otimes 0_{3,3} & 0_{3,3} \otimes 0_{3,3} & 0_{3,4} \otimes 0_{3,3} \\ 0_{3,4} \otimes 0_{3,3} & 0_{3,6} \otimes 0_{3,3} & 0_{3,3} \otimes 0_{3,3} & 0_{3,3} \otimes 0_{3,3} & 0_{3,4} \otimes 0_{3,3} \\ 0_{4,4} \otimes 0_{3,3} & 0_{4,6} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} & 0_{4,4} \otimes 0_{3,3} \end{pmatrix}.$$

Here the $V^{i,j}$ are ordered as indicated before the theorem and

$$M := \begin{pmatrix} A_1 & 0 & 0 \\ 0 & A_1 & 0 \\ 0 & 0 & A_1 \\ 0 & A_2 & 0 \\ 0 & 0 & A_2 \\ 0 & 0 & A_3 \end{pmatrix}.$$

To symmetrize it suffices to construct a symmetrizer for the upper left-hand block

$$Q(\partial) := \begin{pmatrix} I_4 \otimes L(0, \partial) & 0_{4,6} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} \\ 0_{6,4} \otimes 0_{3,3} & I_6 \otimes L(0, \partial) & (I_6 \otimes L(0, \partial))M \\ 0_{3,4} \otimes 0_{3,3} & 0_{3,6} \otimes 0_{3,3} & 0_{3,3} \otimes 0_{3,3} \end{pmatrix}.$$

We verify that

$$\begin{aligned} \tilde{S} & := \begin{pmatrix} I_4 \otimes I_3 & 0_{4,6} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} \\ 0_{6,4} \otimes 0_{3,3} & I_6 \otimes I_3 & (I_6 \otimes I_3)M \\ 0_{3,4} \otimes 0_{3,3} & 0_{3,6} \otimes 0_{3,3} & I_3 \otimes I_3 \end{pmatrix} \quad \text{with} \\ \tilde{S}^{-1} & = \begin{pmatrix} I_4 \otimes I_3 & 0_{4,6} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} \\ 0_{6,4} \otimes 0_{3,3} & I_6 \otimes I_3 & -M(I_3 \otimes I_3) \\ 0_{3,4} \otimes 0_{3,3} & 0_{3,6} \otimes 0_{3,3} & I_3 \otimes I_3 \end{pmatrix} \end{aligned}$$

is a symmetrizer for $Q(i\xi)$. Compute

$$\tilde{S}Q\tilde{S}^{-1} = \begin{pmatrix} I_4 \otimes L(0, \cdot) & 0_{4,6} \otimes 0_{3,3} & 0_{4,3} \otimes 0_{3,3} \\ 0_{6,4} \otimes 0_{3,3} & I_6 \otimes L(0, \cdot) & 0_{6,4} \otimes 0_{3,3} \\ 0_{3,4} \otimes 0_{3,3} & 0_{3,6} \otimes 0_{3,3} & 0_{3,3} \otimes 0_{3,3} \end{pmatrix}$$

which is symmetric since $L(0, \cdot)$ is.

Therefore $P(\xi)$ is symmetrizable by a matrix independent of ξ . Hence, the Cauchy problem for (3.7), (3.8), (3.11)–(3.14) is strongly well-posed. The norm of the zeroth-order terms depends on the coefficients σ_j and their derivatives up to order 2. The estimate of the theorem follows.

Remark 3.1. We recall a computation from [25], showing that when there are only 2 coefficients σ_1 and σ_2 only one derivative of σ_j is needed. This is always the case in dimension $d = 2$. When $\sigma_3 \equiv 0$, split W as

$$W = E_{33}W + \begin{pmatrix} W_1 \\ W_2 \\ 0 \end{pmatrix}.$$

Then

$$\begin{pmatrix} W_1 \\ W_2 \\ 0 \end{pmatrix} = \begin{pmatrix} \sigma_2 U_1^2 \\ \sigma_1 U_2^1 \\ 0 \end{pmatrix} = \text{diag}(\sigma_2, \sigma_1, 0)U.$$

Rewrite (3.7) as

$$\partial_t U + L(0, \partial)U + E_{33}W + \text{diag}(\sigma_2, \sigma_1, 0)U = 0. \tag{3.15}$$

Differentiate with respect to x_1 and x_2 to obtain

$$\partial_t V^j + L(0, \partial)V^j + E_{33}\partial_j W + \partial_j(\text{diag}(\sigma_2, \sigma_1, 0))U + \text{diag}(\sigma_2, \sigma_1, 0)W = 0. \tag{3.16}$$

To find an equation on W , proceed as in the 3D proof to get,

$$\partial_t W + \sum_k \sigma_k A_k V^k + \left(\sum_k \sigma_k \right) W - \sigma_1 \sigma_2 U = 0. \tag{3.17}$$

Therefore \mathbb{V} is a solution of a first-order system, whose principal symbol is given by

$$P(\partial) = \begin{pmatrix} L(0, \partial) & 0 & 0 & 0 \\ 0 & L(0, \partial) & 0 & E_{33}\partial_1 \\ 0 & 0 & L(0, \partial) & E_{33}\partial_2 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

A symmetrizer is given by

$$\tilde{S} = \begin{pmatrix} I & 0 & 0 & 0 \\ 0 & I & 0 & iE_{23} \\ 0 & 0 & I & iE_{13} \\ 0 & 0 & 0 & I \end{pmatrix}, \quad \text{with } \tilde{S}^{-1} = \begin{pmatrix} I & 0 & 0 & 0 \\ 0 & I & 0 & -iE_{23} \\ 0 & 0 & I & -iE_{13} \\ 0 & 0 & 0 & I_3 \end{pmatrix}.$$

3.4. Sharp finite speed for Bérenger’s PML

Recall some notions associated with estimates on the domains of influence and determinacy for a hyperbolic operator L (see [17]). The *timelike cones* are the connected components of $(1, 0, \dots, 0)$ in the complement of the characteristic variety. The *forward propagation cone* is dual to the timelike cone. A Lipschitzian curve $[a, b] \ni t \mapsto (t, \gamma(t))$ is an influence curve when γ' belongs to the propagation cone for Lebesgue almost all $t \in [a, b]$.

Theorem 3.4. *Suppose that L defines a strongly well-posed Cauchy problem and that the multiplicity of $\tau = 0$ as a root of $\det L_1(\tau, \xi) = 0$ is independent of $\xi \in \mathbb{R}^d \setminus 0$. The support of the solution of the Bérenger transmission problem is contained in the union of the propagation curves of L starting in the support of the source terms when either of the following conditions is satisfied.*

- (i) $\forall \xi \in \mathbb{R}^d \setminus 0, \text{Ker } L_1(0, \xi) = \bigcap_j \text{Ker } \xi_j A_j$, and $\forall j, \sigma_j \in L^\infty$.
- (ii) L_1 is Maxwell’s equation and $\forall j, \sigma_j \in W^{2, \infty}$.

Proof. The characteristic varieties of L and \tilde{L} satisfy $\text{Char } \tilde{L} = \text{Char } L \cup \{\tau = 0\}$. When $\{\tau = 0\}$ has multiplicity as a root of $\det L_1(\tau, \xi) = 0$ independent of $\xi \in \mathbb{R}^d \setminus 0$, the timelike cones of L and \tilde{L} coincide. Therefore the propagation cones and influence curves coincide too.

Case (i). Part (ii) of Theorem 2.4 proves that \tilde{L}_1 defines a strongly well-posed Cauchy problem. It follows that the sharp propagation conclusion of the theorem is valid for $\tilde{L}_1 + B(t, x)$ for any bounded $B(t, x)$. This follows on remarking that the solution of $(\tilde{L} + B)\tilde{U} = 0$ with initial data \tilde{U}_0 is the limit at $\nu \rightarrow \infty$ of Picard iterates \tilde{U}^ν . The first, \tilde{U}^1 , is defined as the solution of the Cauchy problem without B . For $\nu > 1$ the iterates are defined by,

$$\tilde{L}_1 \tilde{U}^{\nu+1} + B(t, x) \tilde{U}^\nu = 0, \quad \tilde{U}^{\nu+1}(0, \cdot) = \tilde{U}_0.$$

Since \tilde{L}_1 has constant coefficients, sharp finite speed is classical for that operator. An induction proves that each iterate is supported in the union of influence curves starting in the support of \tilde{U}_0 .

Case (ii). Reason as above constructing by Picard iteration approximations \mathbb{V}^ν converging to the solution \mathbb{V} from (3.3). Since the equation satisfied by \mathbb{V} is strongly well-posed the iterates converge. An induction shows that they are supported in the set of influence curves starting in the support of U_0 . □

3.5. Proof of perfection for Bérenger’s PML by a change of variables

This section continues the analysis of Bérenger’s method when the hypotheses of Theorem 3.4 are satisfied. In those cases well-posedness is proved by an energy method. In addition, suppose that

$$\forall j, \quad \exists L_j > 0, \quad \sigma_j = 0 \quad \text{when } |x_j| \leq L_j. \tag{3.18}$$

Denote by $R := \Pi_j]-L_j, L_j[$.

Definition 3.2. In this setting the method is perfectly matched when for arbitrary $F \in C^\infty_0(]0, \infty[\times R)$ the unique solutions \tilde{V} and \tilde{U} of

$$\tilde{L}\tilde{V} = F, \quad \tilde{V}|_{t \leq 0} = 0, \quad \tilde{L}_1\tilde{U} = F, \quad \tilde{U}|_{t \leq 0} = 0, \tag{3.19}$$

with \tilde{L} as in (1.5) satisfies

$$\tilde{V}|_{\mathbb{R} \times R} = \tilde{U}|_{\mathbb{R} \times R}. \tag{3.20}$$

Theorem 3.5. *With the assumptions of Theorem 3.4 and σ_j as above, Bérenger’s method is perfectly matched.*

Proof. Taking the Laplace transform of the \tilde{V} equation in (3.19) yields a transform holomorphic in $\text{Re } \tau > \tau_0$ with values in $L^2(\mathbb{R}^d)$ satisfying for $1 \leq j \leq d$,

$$\widehat{V}^j + (\tau + \sigma_j(x_j))^{-1} A_j \partial_j \widehat{V} = \widehat{F}^j, \quad \text{with } V := \sum_j V^j, \quad F := \sum_j F^j. \tag{3.21}$$

Multiply by τ and sum on j to obtain

$$\tau \widehat{V} + \sum_j \frac{\tau}{\tau + \sigma_j(x_j)} A_j \partial_j \widehat{V} = \tau \widehat{F}. \tag{3.22}$$

When τ is *fixed real and positive* this equation can be transformed to the corresponding equation without the σ_j by a change of variables. The change of variables depends on τ . The resulting equation is exactly that determining $\widehat{U} := \sum_j \widehat{U}^j$. In this way we find that \widehat{V} is obtained from \widehat{U} by this change of variables. This idea is inspired by Diaz and Joly in [10].

For real $\tau > 0$ define d bi-Lipschitzian homeomorphisms $X_j(x_j)$ of \mathbb{R} to itself by

$$\frac{dX_j(x_j)}{dx_j} = \frac{\tau + \sigma_j(x_j)}{\tau}, \quad X_j(0) = 0.$$

Then,

$$\frac{\partial}{\partial x_j} = \frac{\partial X_j}{\partial x_j} \frac{\partial}{\partial X_j} = \frac{\tau + \sigma_j(x_j)}{\tau} \frac{\partial}{\partial X_j}, \quad \frac{\tau}{\tau + \sigma_j(x_j)} \frac{\partial}{\partial x_j} = \frac{\partial}{\partial X_j}.$$

Therefore if $\widehat{U}(X)$ is the solution of

$$\tau \widehat{U}(X) + \sum_j A_j \frac{\partial}{\partial X_j} \widehat{U} = \widehat{F}(X), \tag{3.23}$$

then the solution \widehat{V} of (3.22) is given by $\widehat{V}(x) := \widehat{U}(X(x))$ since the latter function of x satisfies the equation determining \widehat{V} .

Since $X(x) = x$ for $x \in R$, this proves that the transforms of \widetilde{U} and \widetilde{V} satisfy for real $\tau > \tau_0$

$$\sum_j \widehat{V}^j(\tau, x) = \sum_j \widehat{U}^j(\tau, x), \quad x \in R. \tag{3.24}$$

Since both sides of the identity in (3.24) are holomorphic in $\text{Re } \tau > \tau_0$, it follows that the identity extends to that domain by analytic continuation.

Equation (3.21) and its analogue for \widetilde{U} then imply that for all j , $\widehat{V}^j|_R = \widehat{U}^j|_R$. Uniqueness of the Laplace transform implies $V^j|_R = U^j|_R$ for all t proving perfection. □

Remark 3.2. The proof is very general. It shows that once the initial value problem defined by \widetilde{L} is well-posed, there is perfect matching. The proof works more generally for at least weakly well-posed methods for which the Laplace transform can be reduced to (3.22) for real τ . Our favorite version of the Bérenger algorithm is analyzed in this way in Sec. 3.6.

3.6. Perfection for methods related to Bérenger’s PML

Consider (3.22) with $F = 0$. This equation is the starting point for many authors to construct well-posed PML. It has been viewed as a complex stretching of coordinates (see [29, 9, 26, 12]). This idea, for τ real, becomes an honest change of variables as in [10], that is at the heart of the proof in Sec. 3.5. In the case of Maxwell system, it can be viewed as a system with modified constitutive equations (a lossy medium [27, 2]), or recovered as above from the Bérenger’s system. The system (3.22) is not differential because of the division by $\tau + \sigma_j(x_j)$. In order to recover a hyperbolic system, a change of unknowns is performed. We adopt the approach in [20] for the Maxwell system.

Lemma 3.2. *With matrices given in (3.2), define $S_j := (\tau + \sigma_j(x_j))/\tau$. There exists a pair of invertible matrices M, N , unique up to a multiplication by the same constant, such that*

$$S_j^{-1} N A_j = A_j M, \quad j = 1, 2, 3. \tag{3.25}$$

They are given by

$$M = \gamma \begin{pmatrix} S_1 & 0 & 0 \\ 0 & S_2 & 0 \\ 0 & 0 & S_3 \end{pmatrix}, \quad N = \gamma \begin{pmatrix} S_2 S_3 & 0 & 0 \\ 0 & S_1 S_3 & 0 \\ 0 & 0 & S_1 S_2 \end{pmatrix}, \quad \gamma \in \mathbb{C} \setminus \{0\}. \tag{3.26}$$

Proof. Since

$$A_j e_j = 0, \quad A_j e_{j+1} = -ie_{j+2}, \quad A_j e_{j+2} = ie_{j+1},$$

it is easy to see by applying (3.25) to e_j that M is necessarily diagonal, $M = \text{diag}(m_1, m_2, m_3)$. Applying (3.25) to e_{j+1} and e_{j+2} shows that for any j ,

$$N e_{j+1} = m_{j+2} S_j e_{j+1}, \quad N e_{j+2} = m_{j+1} S_j e_{j+2}.$$

This implies that N is also diagonal, equal to $\text{diag}(m_2 S_3, m_3 S_1, m_1 S_2)$, and

$$m_1 S_3 = m_3 S_1, \quad m_2 S_1 = m_1 S_2, \quad m_3 S_2 = m_2 S_3.$$

This leaves no choice but to choose (3.26). □

In the rest of the analysis take $\gamma = 1$. In (3.22) with $F = 0$ replace \widehat{V} by U . Insert (3.25) to obtain

$$\tau N U + \sum_j A_j M \partial_j U = 0. \tag{3.27}$$

The fact that σ_j depends only on x_j and the form of the matrices guarantees $A_j \partial_j M = 0$. This yields

$$A_j M \partial_j U \equiv A_j \partial_j (MU).$$

Define a new unknown $V := MU$ to find

$$N M^{-1} \tau V + \sum_j A_j \partial_j V = 0. \tag{3.28}$$

$N M^{-1} = \text{diag}(S_1^{-1} S_2 S_3, S_2^{-1} S_3 S_1, S_3^{-1} S_1 S_2)$. Next compute a rational fraction expansion of $\tau S_1^{-1} S_2 S_3$ as

$$\frac{(\tau + \sigma_2)(\tau + \sigma_3)}{\tau + \sigma_1} = \tau + (\sigma_2 + \sigma_3 - \sigma_1) + \frac{\sigma_1^2 + \sigma_2 \sigma_3 - \sigma_1(\sigma_2 + \sigma_3)}{\tau + \sigma_1}.$$

Introduce a new unknown W by

$$\tau N M^{-1} V = \tau V + \Sigma^1 V + \Sigma^2 W, \quad \text{equivalently } W_j = \frac{1}{\tau + \sigma_j(x_j)} V_j = \frac{1}{\tau} U_j,$$

with

$$\begin{aligned} \Sigma &:= \text{diag}(\sigma_1, \sigma_2, \sigma_3), \\ \Sigma^{(1)} &:= \text{diag}(\sigma_2 + \sigma_3 - \sigma_1, \sigma_3 + \sigma_1 - \sigma_2, \sigma_1 + \sigma_2 - \sigma_3), \\ \Sigma^{(2)} &:= \text{diag}((\sigma_1 - \sigma_2)(\sigma_1 - \sigma_3), (\sigma_2 - \sigma_1)(\sigma_2 - \sigma_3), (\sigma_3 - \sigma_1)(\sigma_3 - \sigma_2)). \end{aligned}$$

This leads to a system in the unknowns V and W

$$L(\partial_t, \partial_x) V + \Sigma^{(1)} V + \Sigma^{(2)} W = 0, \quad \partial_t W + \Sigma W - V = 0. \tag{3.29}$$

Finally, U is recovered from

$$U = \partial_t W = V - \Sigma W.$$

The system of equations for V, W is strongly well-posed since L is symmetric hyperbolic. In the case of a single layer in the x_1 direction, there is only one coefficient σ and therefore a single complex-valued supplementary variable. The equations for the magnetic and electric fields are

$$\begin{aligned} \partial_t E_1 - (\nabla \wedge H)_1 - \sigma E_1 + \sigma^2 W_1 &= 0, & \partial_t H_1 + (\nabla \wedge E)_1 - \sigma H_1 + \sigma^2 W_2 &= 0, \\ \partial_t E_2 - (\nabla \wedge H)_2 + \sigma E_2 &= 0, & \partial_t H_2 + (\nabla \wedge E)_2 + \sigma H_2 &= 0, \\ \partial_t E_3 - (\nabla \wedge H)_3 + \sigma E_3 &= 0, & \partial_t H_3 + (\nabla \wedge E)_3 + \sigma H_3 &= 0, \\ \partial_t W_1 + \sigma W_1 - E_1 &= 0, & \partial_t W_2 + \sigma W_2 - H_1 &= 0. \end{aligned}$$

In 2D this is identical to the layers in [27] and equivalent to those in [2].

The principal symbol and lower terms are

$$R_1 = \begin{pmatrix} L & 0 \\ 0 & I_3 \partial_t \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} \Sigma^{(1)} & \Sigma^{(2)} \\ -I_3 & \Sigma \end{pmatrix}.$$

Reversing the computation shows that $(V, W) \in \text{Ker } R(\tau, \xi)$ if and only if $V = MU$, $W = 1/\tau U$, and $L_1(\tau, \frac{\xi_1 \tau}{\tau + \sigma_1}, \dots, \frac{\xi_3 \tau}{\tau + \sigma_3})U = 0$. The characteristic polynomial is therefore the same as for Bérenger’s layer. Thus, by Theorem 2.2

$$\det R(\tau, \xi) = \tau^2 - \sum \frac{\xi_j^2 \tau^2}{(\tau + \sigma_j)^2}.$$

Theorem 3.6. *If $\sigma_j(x_j) \in L^\infty(\mathbb{R})$ and vanish for $|x_j| \leq L_j$, then the system (3.29) for $V := MU$ and $W_j := V_j/(\tau + \sigma_j(x_j))$ is strongly well-posed in $L^2(\mathbb{R}^d)$ and perfectly matched in the sense that for sources supported in $R := \Pi_j\{|x_j| \leq L_j\}$ the function U computed from V, W agrees in $\mathbb{R}_t \times R$ with the solution of Maxwell’s equation with corresponding sources.*

Proof. The proof of Theorem 3.4 applies with only minor modifications. □

Remark 3.3. Note the ease with which strong well-posedness is established and the lack of regularity required of the functions σ_j .

4. Analysis of Layers with Only One Absorption by Fourier–Laplace Transform

There are cases where the energy method presented above does not prove well-posedness. This is the case for the Bérenger algorithm when the ellipticity assumption is not satisfied and the absorptions are not regular. Notably for the Maxwell system and discontinuous absorptions. In this section we present a systematic analysis by Fourier–Laplace transformation of transmission problems with absorption in only one direction.

4.1. Fourier analysis of piecewise constant coefficient transmission problems

Return to the situation of (1.3) with operators L and R on the left and right half spaces and transmission condition (1.4). Suppose that both L and R are weakly hyperbolic in the sense of Gårding. An example is the classical method of Béranger with one absorption. Among other things we will prove that the method is well-posed and perfect. Note the open problem at the end of the introduction emphasizing that we do not know if the classic algorithm with two discontinuous absorptions is well-posed. In addition, we show, by a nontrivial analytic continuation argument in Sec. 4.1.5, that the perfection of Béranger’s method can be verified using the modified plane wave solutions from his original paper. It is our hope that the analysis may help in the construction of new perfectly matched layers.

4.1.1. *Hersh’s condition for transmission problems*

This section takes up the analysis of mixed problems following Hersh in [11]. In the present context we treat transmission problems which are essentially equivalent. The analysis of Hersh assumed the interface is noncharacteristic which is never the case for Maxwell’s equations. We address the changes that are needed to treat problems with characteristic interfaces.

First analyze the solution of the constant coefficient pure initial value problem $\mathcal{L}U = F$ on \mathbb{R}^{1+d} by Laplace transform in time and Fourier transform in $x' = (x_2, \dots, x_d)$. The transform

$$\widehat{U}(\tau, x_1, \eta) := \iint_0^\infty e^{-\tau t} (2\pi)^{-d/2} e^{-ix' \cdot \eta} U(t, x') dt dx'$$

decays as $|x_1| \rightarrow \infty$ and satisfies

$$\mathcal{L}(\tau, d/dx_1, i\eta)\widehat{U} = \widehat{F} \quad \text{in } \mathbb{R}.$$

When A_1 is invertible, this is a standard ordinary differential equation in x_1 . When A_1 is singular, the analysis requires care. The homogeneous equation $\mathcal{L}(\tau, d/dx_1, i\eta)\widehat{U} = 0$ has purely exponential solutions $e^{\rho x_1}$ corresponding to the roots ρ of the equation

$$\det \mathcal{L}(\tau, \rho, i\eta) = 0. \tag{4.1}$$

Hyperbolicity of \mathcal{L} guarantees that for $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$, this equation has no purely imaginary roots.

The number of boundary conditions at $x_1 = 0$ for the boundary value problem in the right half space is chosen equal to the number of roots with negative real part (see also Remark 4.1). That integer must be independent of τ, η . Since roots cannot cross the imaginary axis, the only way the integer can change is if roots escape to infinity. That can happen when the coefficient of the highest power of ρ vanishes. The next hypothesis rules that out.

Definition 4.1. A hyperbolic operator $\mathcal{L}(\partial_t, \partial_x)$ is *nondegenerate with respect to x_1* when there is a $\tau_1 > 0$ so that the degree in ρ of the polynomial $\det \mathcal{L}(\tau, \rho, i\eta)$ is independent of (τ, η) for $\text{Re } \tau > \tau_1, \eta \in \mathbb{R}^{d-1}$.

Example 4.1. In the noncharacteristic case, $\det A_1 \neq 0$, the condition is satisfied and the degree with respect to ρ is equal to N .

(2) For Maxwell’s equations written in the real 6×6 form, the degree with respect to ρ is equal to 4. If written in the complex form (3.2), the degree is 2.

(3) The formula for the characteristic polynomial in Theorem 2.2 shows that if L is nondegenerate, then so is the Béranger doubled operator \tilde{L} with one absorption σ_1 in $x_1 > 0$. The degree in ρ is the same for L and \tilde{L} .

(4) If $\mathcal{L} = \mathcal{L}_1 + B$ is nondegenerate with respect to x_1 , then so is the operator $P := \mathcal{L}_1(\partial) + a^{-1}B = a^{-1}\mathcal{L}(a\partial)$ for any $a > 0$. If the degree for \mathcal{L} is constant in $\text{Re } \tau > \tau_1$, then the degree for P is constant for $\text{Re } \tau > a^{-1}\tau_1$.

For the lemmas to follow, it is useful to transform so that \mathcal{A}_1 has block form.

Lemma 4.1. *If \mathcal{L} in (2.1) is nondegenerate with respect to x_1 then for $\text{Re } \tau > \tau_1, \eta \in \mathbb{R}^{d-1}$,*

- (i) *the degree in ρ of the polynomial $\det \mathcal{L}(\tau, \rho, i\eta)$ is equal to $\text{rank } \mathcal{A}_1$,*
- (ii) *the number of roots ρ with positive real part is equal to the number of negative eigenvalues of \mathcal{A}_1 .*

Proof. Since \mathcal{L} is nondegenerate, it suffices to study the case $\eta = 0$.

(i) Choose invertible K so that

$$K^{-1}\mathcal{A}_1K = \begin{pmatrix} \mathcal{A} & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathcal{A} \text{ an invertible square matrix of size } \text{rank } \mathcal{A}_1.$$

Then

$$K^{-1}\mathcal{L}(\tau, \rho, 0)K = \begin{pmatrix} \tau I + \mathcal{A}\rho & 0 \\ 0 & \tau I \end{pmatrix} + \text{matrix independent of } \tau, \rho.$$

It follows that the degree in ρ is no larger than $\text{rank } \mathcal{A}_1$.

The coefficient of $\rho^{\text{rank } \mathcal{A}_1}$ in $\det \mathcal{L}(\tau, \rho, 0)$ is a polynomial in τ of degree $\leq N - \text{rank } \mathcal{A}_1$. For large τ the coefficient is equal to

$$(\det \mathcal{A}) (\tau^{N-\text{rank } \mathcal{A}_1}) + \text{lower order in } \tau.$$

Thus the degree in ρ is $\text{rank } \mathcal{A}_1$ for such τ proving the result.

(ii)

$$\det \mathcal{L}(\tau, \rho, 0) = \det \left(\begin{pmatrix} \tau I + \rho \mathcal{A} & 0 \\ 0 & \tau I \end{pmatrix} + B \right).$$

For fixed τ , sufficiently large, $\rho \neq 0$, and ρ/τ is a root of the polynomial $p(x, 1/\tau)$ of degree $\text{rank } \mathcal{A}$:

$$p(x, \varepsilon) = \begin{vmatrix} I + x\mathcal{A} + \varepsilon B_{11} & \varepsilon B_{12} \\ \varepsilon B_{21} & I + \varepsilon B_{22} \end{vmatrix}.$$

This polynomial has exactly $\text{rank } \mathcal{A}$ roots. By Rouché’s theorem,

$$\frac{\rho_j}{\tau} \sim -\frac{1}{\lambda_j}, \quad \tau \gg 1,$$

where the λ_j are the $\text{rank } \mathcal{A}$ eigenvalues of \mathcal{A} repeated according to their algebraic multiplicity. Since the eigenvalues of \mathcal{A} and the nonzero eigenvalues of \mathcal{A}_1 are the same, this completes the proof. \square

Remark 4.1. For the transformed one-dimensional hyperbolic operator $\mathcal{L}(\partial_t, \partial_1, i\eta)$, the number of incoming characteristics at the boundary $x_1 = 0$ in the right half space is equal to the number of strictly positive eigenvalues of \mathcal{A}_1 . The second part of the lemma shows that this is equal to the number of roots with negative real part. The two natural ways to compute the number of necessary boundary conditions yield the same answer.

The next lemma shows that for nondegenerate operators, the characteristic case can be transformed to a standard ordinary differential equation.

Lemma 4.2. *Suppose that $A, M \in \text{Hom}(\mathbb{C}^N)$ and the equation $\det(A\rho + M) = 0$ has degree in ρ equal to $\text{rank } A$ and no purely imaginary roots. Then,*

(i) *The matrix M is invertible and all solutions of the homogeneous equation*

$$A \frac{dU}{dx_1} + MU = 0 \tag{4.2}$$

take values in the space $\mathbb{G} := M^{-1}(\text{Range } A)$ satisfying $\dim \mathbb{G} = \text{rank } A$.

(ii) *There is a $\widetilde{M} \in \text{Hom } \mathbb{G}$ so that a function U satisfies (4.2) if and only if U is \mathbb{G} valued and satisfies*

$$\frac{dU}{dx_1} + \widetilde{M}U = 0. \tag{4.3}$$

(iii) *The vector space \mathbb{U} of solutions of (4.2) is a linear subspace of $C^\infty(\mathbb{R})$ with dimension equal to $\text{rank } A$. The Cauchy problem with data in \mathbb{G} is well-posed.*

Proof. (i) Since $\rho = 0$ is not a root, M is invertible. The equation $U = M^{-1}AdU/dx_1$ shows that continuously differentiable solutions U take values in \mathbb{G} . More generally, if U is a distribution solution and $\psi \in C_0^\infty(\mathbb{R})$ takes values in the annihilator, \mathbb{G}^\perp of \mathbb{G} , then

$$\langle U, \psi \rangle = \langle M^{-1}AdU/dx_1, \psi \rangle = \langle dU/dx_1, (M^{-1}A)^*\psi \rangle.$$

But $\mathbb{G} = \text{range } M^{-1}A$ so $\mathbb{G}^\perp = \ker(M^{-1}A)^*$. Therefore $(M^{-1}A)^*\psi = 0$ so $\langle U, \psi \rangle = 0$ which is the desired conclusion.

(ii) Multiplying the equation by an invertible P and making the change of variable $U = KV$ transforms the equation to the equivalent form

$$PAK \frac{dV}{dx_1} + PMKV = 0.$$

Choose invertible P, K so that PAK has block form

$$PAK = \begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix},$$

where I is the rank $A \times \text{rank } A$ identity matrix. With $V = (V_1, V_2)$, one has the block forms

$$PMK = \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix} \frac{dV}{dx_1} + \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} V = 0.$$

One has,

$$\det(A\rho + M) = \det P^{-1} \det \begin{pmatrix} \rho I + H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} \det K^{-1}.$$

The first part of the preceding lemma implies that the determinant on the left is a polynomial of degree $\text{rank } A$ in ρ . It follows that H_{22} is invertible.

The solutions V satisfy $H_{21}V_1 + H_{22}V_2 = 0$ so take values in $\mathbb{V} := \{V_2 = -H_{22}^{-1}H_{21}V_1\}$. The function V is a solution if and only if it takes values in \mathbb{V} and

$$\frac{dV_1}{dx_1} + RV_1 = 0, \quad R := H_{11} - H_{12}H_{22}^{-1}H_{21}.$$

If $N : \mathbb{V} \rightarrow \mathbb{V}$ is the map,

$$(V_1, V_2) \mapsto (RV_1, -H_{22}^{-1}H_{21}RV_1),$$

then V is a solution if and only if it is \mathbb{V} -valued and satisfies $dV/dx_1 = NV$. Writing $V = K^{-1}U$ and $\widetilde{M} = -KN$ implies (ii).

(iii) Follows from (ii). □

Lemma 4.3. *If \mathcal{L} is hyperbolic and nondegenerate with respect to x_1 , then its principal part $\mathcal{L}_1(\partial_t, \partial_x)$ is also nondegenerate with respect to x_1 . The degree in ρ of $\det \mathcal{L}_1(\tau, \rho, i\eta)$ is constant for $\text{Re } \tau > 0$ and $\eta \in \mathbb{R}^{d-1}$.*

Proof. With notation from the preceding proof,

$$\mathcal{L}(\partial_t, \partial_{x_1}, \partial_{x'}) = P^{-1} \left(\begin{pmatrix} I\partial_{x_1} & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} H_{11}(\partial_t, \partial_{x'}) & H_{12}(\partial_t, \partial_{x'}) \\ H_{21}(\partial_t, \partial_{x'}) & H_{22}(\partial_t, \partial_{x'}) \end{pmatrix} \right) K^{-1}. \quad (4.4)$$

The proof of the last lemma showed that for $\text{Re } \tau > \tau_1$ and $\eta \in \mathbb{R}^{d-1}$, $H_{22}(\tau, \eta)$ is invertible.

The computation in Lemma 4.1 shows that for $\eta = 0$ and $\mathbb{R} \ni \tau \rightarrow \infty$ the coefficient of $\rho^{\text{rank } \mathcal{A}_1}$ has modulus $\geq c\tau^{N-\text{rank } \mathcal{A}_1}$ with $c > 0$. This implies that $\eta = 0$ is noncharacteristic for H_{22} . Therefore $H_{22}(\partial_t, \partial_{x'})$ is hyperbolic.

Replacing \mathcal{L} by its principal part \mathcal{L}_1 has the effect of replacing each operator $H_{ij}(\partial)$ by its principal part. This yields identity (4.4) with \mathcal{L} and the H_{ij} replaced by their principal parts.

Since the principal part of a hyperbolic operator is hyperbolic, it follows that $(H_{22})_1(\partial_t, \partial_{x'})$ is a homogeneous hyperbolic operator. Therefore $(H_{22})_1(\tau, i\eta)$ is invertible for $\eta \in \mathbb{R}^{d-1}$ and $\text{Re } \tau \neq 0$. Thus, the coefficient of $\rho^{\text{rank } \mathcal{A}_1}$ in $\det \mathcal{L}_1(\tau, \rho, i\eta)$ is nonzero for $\eta \in \mathbb{R}^{d-1}$ and $\text{Re } \tau \neq 0$. □

Lemma 4.4. *Suppose that the ordinary differential equation (4.2) satisfies the hypotheses of Lemma 4.2. Denote by E^\pm the linear space of solutions which tend exponentially to zero as $x_1 \rightarrow \pm\infty$ and by \dot{E}^\pm their traces at $x_1 = 0$. Then*

- (i) $\dot{E}^\pm \cap \ker A = \{0\}$,
- (ii) $\dim AE^\pm = \dim E^\pm$,
- (iii) The map $U \mapsto U(0)$ is an isomorphism from E^\pm to \dot{E}^\pm ,
- (iv) $A\dot{E}^+ \cap A\dot{E}^- = \{0\}$,
- (v) $A\dot{E}^+ \oplus A\dot{E}^- = \text{Range } A$.

Proof. (i) The absence of purely imaginary roots shows that every solution is uniquely the sum of two solutions. One grows exponentially at $+\infty$ and decays exponentially at $x_1 = -\infty$. The second grows at $-\infty$ and decays at $+\infty$. In particular the only bounded solution is the zero solution.

If $e_+ \in \dot{E}^+ \cap \ker A$, denote by $U(x_1)$ the solution with this Cauchy data. The function that is equal to U on $x_1 > 0$ and equal to 0 in $x_1 \leq 0$ is a distribution solution of (4.2) on all of \mathbb{R} since $A[U]_{x_1=0} = 0$. This solution is bounded hence identically equal to zero. Therefore $e_+ = 0$. The case for $\dot{E}^- \cap \ker A$ is analogous.

(ii) Follows from (i).

(iii) It is surjective by definition. If it were not injective for E^+ , there would be a nontrivial solution $U(x)$ exponentially decaying as $x_1 \rightarrow +\infty$ with $U(0) = 0$ violating (i).

(iv) The set $A\dot{E}^+$ consists of the values $AU_+(0)$ with U_+ satisfying (4.2) and exponentially decreasing in $x_1 > 0$. If the intersection were nontrivial there would be a solutions U_- decaying as $x_1 \rightarrow -\infty$ so that $AU_+(0) = AU_-(0)$. The function V equal to U_+ in $x_1 > 0$ and U_- in $x_1 < 0$ is then a distribution solution for all x_1 exponentially decaying in both directions. Hyperbolicity implies that $V = 0$ contradicting the nontriviality.

(v) Using (ii) and (iv), one sees that the direct sum on the left is a subspace of $\text{Range } A$ of full dimension. □

The next lemma is needed in Sec. 4.1.2.

Lemma 4.5. *Assume that the hypothesis and notations of Lemma 4.2 are in force. Then for $K \in \dot{E}^+$ there is an $F \in C_0^\infty(]-\infty, 0[)$ so that the unique solution of*

$$A \frac{dU}{dx_1} + MU = F, \quad \lim_{|x_1| \rightarrow \infty} \|U(x_1)\| = 0, \tag{4.5}$$

satisfies $U(0) = K$.

Proof. Consider first the case of invertible A . A change of dependent variable yields the block form for the new variable still denoted U

$$\frac{dU}{dx_1} + \begin{pmatrix} M_+ & 0 \\ 0 & M_- \end{pmatrix} U = F, \quad U = (U_1, U_2), \quad F = (F_1, F_2),$$

$$\text{spec } M_\pm \subset \{\pm \text{Re } z > 0\}.$$

Then $\dot{E}^+ = \{U_2 = 0\}$ so $K = (K_1, 0)$. Choose $F = (F_1, 0)$. Then $U(0) = K$ if and only if,

$$K_1 = \int_{-\infty}^0 e^{M_+ s} F_1(s) ds.$$

This is achieved with,

$$F_1(s) = \chi(s) e^{-M_+ s} K_1, \quad \chi \in C_0^\infty(]-\infty, 0[), \quad \int \chi(s) ds = 1.$$

When A is not invertible change variable as in Lemma 4.2 to find the block form

$$\begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix} \frac{dU}{dx_1} + \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} U = F,$$

with invertible H_{22} .

Part (i) of (4.2) implies that the map $\mathbb{G} \ni G = (G_1, G_2) \mapsto G_1$ is an isomorphism. Write $\mathbb{G} \ni K = (K_1, K_2)$. Choose $F = (F_1, 0)$. Then choose a \mathbb{G} -valued solution U defined by

$$\frac{dU_1}{dx_1} + H_{11}U_1 = F_1, \quad U_2 = -H_{22}^{-1}H_{21}U_1.$$

One has $U(0) = K$ if and only if $U_1(0) = K_1$. The construction in the invertible case completes the proof. □

Suppose that

$$L = \partial_t + A_1 \partial_1 + \dots \quad \text{and} \quad R = \partial_t + \mathcal{A}_1 \partial_1 + \dots$$

are nondegenerate with respect to x_1 . For $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$, define $E_L^\pm(\tau, \eta)$ to be the set of solutions of

$$L(\tau, d/dx_1, i\eta)V = 0$$

tending to zero as $x_1 \rightarrow \pm\infty$. Denote by $\dot{E}_L^\pm(\tau, \eta) \subset \mathbb{C}^N$ the linear space of traces at $x_1 = 0$ of solutions in $E_L^\pm(\tau, \eta)$.

Similarly with a possibly larger value still called τ_0 , there are $\dot{E}_R^\pm(\tau, \eta) \subset \mathbb{C}^M$ so that the solutions of $R(\tau, d/dx_1, i\eta)Z = 0$ taking values in $\dot{E}_R^\pm(\tau, \eta)$ are exactly those tending to zero exponentially as $x_1 \rightarrow \pm\infty$. The subspaces $E_L^\pm(\tau, \eta)$ and $E_R^\pm(\tau, \eta)$ depend smoothly on τ, η for $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$.

Lemma 4.1 implies that

$$\begin{aligned} \dim E_L^-(\tau, \eta) &= \# \text{ positive eigenvalues of } A_1, \\ \dim E_R^+(\tau, \eta) &= \# \text{ negative eigenvalues of } \mathcal{A}_1. \end{aligned} \tag{4.6}$$

Consider the inhomogeneous transmission problem,

$$LV = 0 \quad \text{when } x_1 < 0, \quad RW = 0 \quad \text{when } x_1 > 0, \tag{4.7}$$

$$(V, W) - g \in \mathcal{N} \quad \text{when } x_1 = 0. \tag{4.8}$$

The problem with inhomogeneous term F can be reduced to this form by subtracting on the left a solution of the hyperbolic Cauchy problem $LU = F$ on \mathbb{R}^{1+d} with $U|_{t < 0} = 0$. Denote by $\widehat{V}(\tau, x_1, \eta), \widehat{W}(\tau, x_1, \eta), \widehat{g}(\tau, \eta)$ the Fourier–Laplace transforms. The transform \widehat{U} is defined for $x_1 \in \mathbb{R}$, while \widehat{V} (respectively \widehat{W}) is defined for $x_1 < 0$ (respectively $x_1 > 0$). The transforms \widehat{V}, \widehat{W} decay as $|x_1| \rightarrow \infty$. \widehat{V}, \widehat{W} satisfy the ordinary differential transmission problem

$$L(\tau, d/dx_1, i\eta)\widehat{V} = 0 \quad \text{in } x_1 < 0, \quad R(\tau, d/dx_1, i\eta)\widehat{W} = 0 \quad \text{in } x_1 > 0, \tag{4.9}$$

$$(\widehat{V}(\tau, 0, \eta), \widehat{W}(\tau, 0, \eta)) - \widehat{g}(s, \eta) \in \mathcal{N}. \tag{4.10}$$

Hersh’s necessary and sufficient condition for well-posedness of the transmission problem is derived as follows. Uniqueness of solutions of (4.9), (4.10) for $\text{Re } \tau > \tau_0, \eta \in \mathbb{R}^{d-1}$ is equivalent to the fact that there are no exponentially decaying solutions of the homogeneous transmission problem. That is,

$$\mathcal{N} \cap (\dot{E}_L^-(\tau, \eta) \times \dot{E}_R^+(\tau, \eta)) = \{0\}. \tag{4.11}$$

In order to guarantee existence, one imposes the maximality condition,

$$\mathcal{N} \oplus (\dot{E}_L^-(\tau, \eta) \times \dot{E}_R^+(\tau, \eta)) = \mathbb{C}^N \times \mathbb{C}^M. \tag{4.12}$$

Using (4.6), this determines the dimension of \mathcal{N} from the coefficients A_1 and \mathcal{A}_1 of L and R respectively.

Definition 4.2. If the transmission problem (4.7), (4.8) satisfies (4.12) for all $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$ it is said to satisfy *Hersh’s condition*.

Theorem 4.1. *Hersh’s condition is satisfied if and only if there is an r and a λ_0 so that for all $\lambda > \lambda_0$ and g supported in $t \geq 0$ with $e^{-\lambda t}g \in H^{s+r}(\mathbb{R}_{t,x}^d)$ with*

values in $\mathbb{C}^N \times \mathbb{C}^M$ there is a unique V, W supported in $t \geq 0$ with

$$e^{-\lambda t}V \in H^s(\text{]} - \infty, \infty[\times\{x_1 < 0\}) \quad \text{and} \quad e^{-\lambda t}W \in H^s(\text{]} - \infty, \infty[\times\{x_1 > 0\})$$

satisfying the transmission problem (4.7), (4.8).

Sketch of Proof. We have shown that the Hersh condition permits one to compute a candidate Fourier–Laplace transform. We outline how the condition implies the desired estimate. The method is to use the Seidenberg–Tarski Theorem 2.3 to derive a lower bound on the real parts of the roots ν together with a contour integral representation. The same elements form the heart of [18], and Sec. 12.9 of [13]. In the present context we treat a transmission problem rather than a boundary value problem. In addition, one needs to use the earlier lemmas to treat the case when $x_1 = 0$ is characteristic.

Choose $\Lambda > \max\{\tau_0(L), \tau_0(R)\}$. The equations

$$\det L(\tau, \nu, i\eta) = 0, \quad \det R(\tau, \nu, i\eta) = 0$$

with $\text{Re } \tau \geq \Lambda, \eta \in \mathbb{R}^{d-1}$ have no purely imaginary roots. Define

$$\zeta(R) := \min \{|\text{Re } \nu| : \eta \in \mathbb{R}^{d-1}, \text{Re } \tau \geq \Lambda, |\tau|^2 + |\eta|^2 \leq R^2, \{\det L(\tau, \nu, i\eta) = 0 \text{ or } \det R(\tau, \nu, i\eta) = 0\}\}.$$

The Seidenberg–Tarski Theorem 2.3 implies that there is a $\rho \in \mathbb{Q}$ and $b \neq 0$ so that

$$\zeta(R) = R^\rho(b + o(R)), \quad \text{as } R \rightarrow \infty.$$

Thus, there are C, N so that, for any τ, η with $\text{Re } \tau \geq \Lambda$,

$$|\text{Re } \nu| \geq \frac{C}{1 + |(\tau, \eta)|^N}. \tag{4.13}$$

The solutions in $E_R^+(\tau, \eta)$ are written using a contour integral representation of \widehat{W} in the block form of Lemma 4.2. Here the matrix H_{ij} depends of (τ, η) . Denote by $D = D(\tau, \eta)$ the finite union of squares with centers at the roots with $\text{Re } \nu < 0$. The side of each square is the smaller of 1 and half the distance of the root to the imaginary axis. Then

$$\widehat{W}_I = \frac{1}{2\pi i} \oint_{\partial D} e^{\tau x_1} (\tau + (H_{11} + H_{12}H_{22}^{-1}H_{21}))^{-1} d\tau \widehat{W}_I, \quad \widehat{W}_{II} = H_{22}^{-1}H_{21}\widehat{W}_I. \tag{4.14}$$

The Seidenberg–Tarski Theorem 2.3 applied to

$$\max\{|w|^2 : |z|^2 = 1, H_{22}w = z, \text{Re } \tau \geq \Lambda, |\tau|^2 + |\eta|^2 \leq R^2\}$$

proves that

$$\|H_{22}(\tau, \eta)^{-1}\| = R^\beta(a + o(1)), \quad a \neq 0, \quad \beta \in \mathbb{Q}.$$

This estimate together with (4.13) yields with new C, N ,

$$\int_0^\infty |\widehat{W}(\tau, x_1, \eta)|^2 dx_1 \leq C(1 + |(\tau, \eta)|^{2N})|\widehat{W}_I(0)|^2.$$

With the analogous expression for V the solution of (4.7) satisfies

$$\int_{-\infty}^\infty |\widehat{V}(\tau, x_1, \eta)|^2 dx_1 \leq C(1 + |(\tau, \eta)|^{2N})|\widehat{V}_I(0)|^2.$$

The Hersh condition asserts that for each (τ, η) , $\widehat{W}_I(0)$ and $\widehat{V}_I(0)$ are uniquely determined by $\widehat{g}(\tau, \eta)$. Seidenberg–Tarski Theorem 2.3 yields an estimate

$$\|\widehat{W}_I(0), \widehat{V}_I(0)\| \leq C(1 + |(\tau, \eta)|^a)\|\widehat{g}(\tau, \eta)\|^2.$$

The last three estimates together with Parseval’s identity proves the desired estimate,

$$\exists C, N, \quad \forall g, \quad \forall \lambda > \Lambda, \quad \|e^{-\lambda t}U\|_{L^2(\mathbb{R}^{1+d})}^2 \leq C \sum_{|\alpha| \leq N} \|e^{-\lambda t}\partial_{t,x}^\alpha g\|_{L^2(\mathbb{R}^{1+(d-1)})}.$$

This estimate proves the existence part of the theorem. □

4.1.2. Necessary and sufficient condition for perfection

The Fourier–Laplace method is used to derive a necessary and sufficient condition for perfection of an absorbing layer. Begin with a closer analysis of the transform, $\widehat{U}(\tau, x_1, \eta)$, of the solution of the basic Eq. (1.1).

When A_1 is invertible, \widehat{U} is analyzed as follows. Denote by $\Pi_\pm(\tau, \eta)$ the projectors associated with the direct sum decomposition $\dot{E}_L^+(\tau, \eta) \oplus \dot{E}_L^-(\tau, \eta) = \mathbb{C}^N$. Define $S_\pm(\tau, x_1, \eta)$ as the $\text{Hom}(\mathbb{C}^N)$ -valued solutions of

$$L(\tau, d/dx_1, i\eta)S_\pm = 0, \quad S_\pm|_{x_1=0} = A_1^{-1}\Pi_\pm.$$

Then S_\pm decays exponentially as $x_1 \rightarrow \pm\infty$ and

$$\chi_{] -\infty, 0[} S_- + \chi_{[0, \infty[} S_+$$

is the unique tempered fundamental solution of $L(\tau, d/dx_1, i\eta)$. Decompose $\widehat{F} = \widehat{F}_- + \widehat{F}_+$, $\widehat{U} = \widehat{U}_+ + \widehat{U}_-$ according to $E_L^+(\tau, \eta) \oplus E_L^-(\tau, \eta) = \mathbb{C}^N$. Then \widehat{U}_- is the convolution of \widehat{F}_- with $\chi_{] -\infty, 0[} S_-$ and \widehat{U}_+ is the convolution of \widehat{F}_+ with $\chi_{[0, \infty[} S_+$. In particular, $\widehat{U}_-(\tau, 0, \eta)$ vanishes on a neighborhood of $[0, \infty[$ so $\widehat{U}(\tau, 0, \eta) = \widehat{U}_+(\tau, 0, \eta) \in \dot{E}_L^+(\tau, \eta)$. The value of \widehat{U} in $x_1 \geq 0$ satisfy the homogeneous ordinary differential equation $L(\tau, d/dx_1, i\eta)\widehat{U} = 0$ with initial value $\widehat{U}(0) \in \dot{E}_L^+(\tau, \eta)$.

To reach the same conclusion when A_1 is singular, apply the lemmas of the preceding section to the equation $L(\tau, d/dx_1, i\eta)Z = 0$. Lemma 4.4 applied to $A = A_1$ and $M = \tau I + i \sum_{j=2}^d A_j \eta_j$ shows that both $\dot{E}_L^\pm(\tau, \eta)$ are subspaces of \mathbb{G} and that the space of solutions is a direct sum $E_L^+(\tau, \eta) \oplus E_L^-(\tau, \eta)$. It follows that

$$\dot{E}_L^-(\tau, \eta) \oplus \dot{E}_L^+(\tau, \eta) = \mathbb{G}(\tau, \eta).$$

Repeating the analysis in the nonsingular case applied to (4.3) shows that $\widehat{U}(\tau, 0, \eta) \in \dot{E}_L^+(\tau, \eta)$.

Definition 4.3. For a transmission problem (L, R, \mathcal{N}) satisfying Hersh’s condition (Definition 4.2), $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$, define the reflection operator, $H(\tau, \eta) : \dot{E}_L^+(\tau, \eta) \rightarrow \dot{E}_L^-(\tau, \eta)$ as follows. Hersh’s condition implies that for each $K \in \dot{E}_L^+(\tau, \eta)$ there is a unique $(\dot{V}, \dot{W}) \in \dot{E}_L^-(\tau, \eta) \times \dot{E}_R^+(\tau, \eta)$ so that $(K, 0) \equiv (\dot{V}, \dot{W}) \text{ mod } \mathcal{N}$. Define, $H(\tau, \eta)K := \dot{V}$.

Theorem 4.2. Suppose that the transmission problem (L, R, \mathcal{N}) satisfies the Hersh condition. The following are equivalent.

- (i) The transmission problem is perfectly matched in the sense of Definition 1.1.
- (ii) There is a $\tau_0 \in \mathbb{R}$ so that for all $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$, $H(\tau, \eta) = 0$.
- (iii) There is a $\tau_0 \in \mathbb{R}$ so that for all $\text{Re } \tau > \tau_0$ and $\eta \in \mathbb{R}^{d-1}$,

$$\forall K_L \in \dot{E}_L^+(\tau, \eta), \quad \exists! K_R \in \dot{E}_R^+(\tau, \eta), \quad \text{such that } (K_L, K_R) \in \mathcal{N}. \quad (4.15)$$

Proof. Conditions (ii) and (iii) are clearly equivalent.

For the equivalence with (i), compare the values of \widehat{U} and \widehat{V} in $\{x_1 < 0\}$. Since both satisfy $LZ = F$ and decay as $x_1 \rightarrow -\infty$ it follows that $L(\widehat{V} - \widehat{U}) = 0$, so, $\widehat{V} - \widehat{U} := \Gamma$ is an \dot{E}_L^- valued solution of $L\Gamma = 0$. Since $F = 0$ in $x_1 > 0$, $\widehat{W} \in E_R^+$. The transmission condition requires that

$$\mathcal{N} \ni (\widehat{V}(0), \widehat{W}(0)) = (\widehat{U}(0) + \Gamma(0), \widehat{W}(0)) = (\widehat{U}(0), 0) + (\Gamma(0), \widehat{W}(0)). \quad (4.16)$$

Since $(\Gamma(0), \widehat{W}(0)) \in \dot{E}_L^-(\tau, \eta) \times \dot{E}_R^+(\tau, \eta)$, (4.16) expresses $(\widehat{U}(0), 0)$ as a sum of an element in \mathcal{N} and an element of $\dot{E}_L^-(\tau, \eta) \times \dot{E}_R^+(\tau, \eta)$. The Hersh condition (4.12) asserts that such a decomposition is unique. Therefore $(\widehat{V}(0), \widehat{W}(0))$ is uniquely determined from $\widehat{U}(0)$.

The method is perfectly matched if and only if for all F supported in $x_1 < 0$, $t \geq 0$

$$V = U|_{x_1 < 0}.$$

This occurs if and only if Γ vanishes for $x_1 < 0$ which holds if and only if $\Gamma(0) = 0$.

If the method is perfectly matched, then in the decomposition (4.16) one has $\Gamma(0) = 0$. Then $(\widehat{U}(0), \widehat{W}(0)) \in \mathcal{N}$. Lemma 4.5 asserts that for any $K \in \dot{E}_L^+$ there is an F so that $\widehat{U}(0) = K$. This proves that (4.15) holds.

Conversely if (4.15) holds, then in the decomposition (4.16), $\Gamma(0) = 0$ so $\Gamma = 0$. It follows that $U|_{x_1 < 0} = V$. □

Remark 4.2. (1) When (4.15) holds, the decomposition of $(K, 0) \in \mathbb{C}^N \times \mathbb{C}^M$ in the direct sum (4.12) is,

$$(K, 0) = (K, W(K)) - (0, W(K)) \in \mathcal{N} \oplus (E_L^- \times E_R^+).$$

(2) With $K = U(0)$ as above, the solution (V, W) of the ordinary differential equation transmission problem is given by $V = U|_{x < 0}$ and W is the solution of $RZ = 0$ with $Z(0) = -W(K)$.

(3) In the important case where $N = M$, invertible A_1 and \mathcal{A}_1 and transmission condition $\mathcal{N} = \{V = W\}$, the perfection criterion (iii) asserts that $\dot{E}_L^+(\tau, \eta) = \dot{E}_R^+(\tau, \eta)$.

We present a typical example showing that the natural absorbing layers are virtually never perfectly matched in dimension $d \geq 2$.

Proposition 4.1. *Consider the dissipative symmetric hyperbolic example with $d = N = M = 2$,*

$$A_1 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, \quad A_2 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad R = L + P, \quad P = P^* \geq 0,$$

$$\mathcal{N} = \{(V, W) : V = W\}.$$

- (i) *The transmission problem is perfectly matched if and only if $P = 0$.*
- (ii) *The corresponding problem with $d = 1$ is perfectly matched if and only if P is diagonal.*

Proof. Define

$$M_L := A_1^{-1}[\tau + i\eta A_2], \quad M_R := A_1^{-1}[\tau + i\eta A_2] + A_1^{-1}P,$$

so that $A_1^{-1}L(\tau, \partial_1, i\eta) = \partial_1 + M_L(\tau, \eta)$ and similarly for $A_1^{-1}R(\tau, \partial_1, i\eta)$. For $\operatorname{Re} \tau > 0$ and $\eta \in \mathbb{R}$, the matrices M_L and M_R have one eigenvalue with positive real part and one with negative real part. The eigenspace corresponding to positive (respectively negative) real part eigenvectors is equal to $\dot{E}_L^+(\tau, \eta)$ (respectively $\dot{E}_R^+(\tau, \eta)$). Therefore the necessary and sufficient condition for perfection is that for $\operatorname{Re} \tau > \tau_0$ and any η , $\dot{E}_L^+(\tau, i\eta) = \dot{E}_R^+(\tau, i\eta)$.

Since

$$L(\tau, \rho, i\eta) = \begin{pmatrix} \tau + \rho & i\eta \\ i\eta & \tau - \rho \end{pmatrix} \quad \text{and} \quad \det L(\tau, \rho, i\eta) = \tau^2 - \rho^2 + \eta^2,$$

the eigenvalue of $M_L(\tau, \eta)$ with positive real part is $\rho = \sqrt{\tau^2 + \eta^2}$. The eigenspace is the kernel of $L(\tau, \rho, i\eta)$. Therefore

$$\dot{E}_L^+(\tau, \eta) = \mathbb{C}(-i\eta, \tau + \rho). \tag{4.17}$$

Since $M_R = M_L + A_1^{-1}P$, a necessary condition is that the family of vectors $v(\eta, \tau) := (-i\eta, \tau + \rho)$ be eigenvectors of the constant matrix $A_1^{-1}P$, which is possible only if $A_1^{-1}P$ is a constant multiple of the identity. Therefore $P = cA_1$. Since $P \geq 0$ and A_1 has eigenvalues of both signs, it follows that $c = 0$ proving (i).

In the one-dimensional case there is just one eigenvector $(0, 1)$ which must be an eigenvector of $A_1^{-1}P$. Since $(0, 1)$ is also an eigenvector of A_1 , it follows that $(0, 1)$ must be an eigenvector of P . Since $P = P^*$, the orthogonal vector $(1, 0)$ is

also an eigenvector and P is diagonal. Conversely, if P is diagonal the condition is satisfied. □

Remark 4.3. (1) Examples verifying perfection for a family of absorbing layers related to but not including those of Bérenger are presented in [3]. To our knowledge, Hersh’s criterion for Bérenger’s layers has not been verified before.

(2) The perfection criterion is related to the plane wave criterion of Bérenger. We examine the relation in Sec. 4.1.6.

4.1.3. Hersh’s condition for Bérenger’s PML with piecewise constant σ_1

Of our earlier results, only those of Sec. 3.2 apply to discontinuous absorptions. So, if the generator is not elliptic, (for example, the PML Maxwell system of Bérenger), the preceding results do not prove that the initial value problem is well-posed. In this section we prove that the doubled operators of Bérenger define a (weakly) well-posed initial value problem provided that

$$\sigma_j \equiv 0 \quad \text{for } j \geq 2 \quad \text{and} \quad \sigma_1(x_1) \equiv \sigma^\pm \quad \text{in } \mathbb{R}_\pm^d, \tag{4.18}$$

and, the constant coefficient operators \tilde{L} on \mathbb{R}_\pm^d are both (weakly) hyperbolic.

The unknown \tilde{U} satisfies (1.5). Denote by $\tilde{U}^\pm = \{U_1^\pm, \dots, U_d^\pm\}$ the restriction of the unknown \tilde{U} to \mathbb{R}_\pm^d . They satisfy differential equations in the half spaces \mathbb{R}_\pm^d .

Lemma 4.6. For \tilde{U} locally square integrable on a neighborhood of $(\underline{t}, \underline{x}) \in \{x_1 = 0\}$, the following are equivalent.

- (i) $\tilde{L}\tilde{U} \in L^2$ on a neighborhood of $(\underline{t}, \underline{x})$ in \mathbb{R}^{1+d} in the sense of distributions.
- (ii) There is a neighborhood \mathcal{O} of $(\underline{t}, \underline{x})$ so that $\tilde{L}\tilde{U}^\pm$ is square integrable on $\mathcal{O} \cap \mathbb{R}_\pm^d$ and $[\tilde{A}_1\tilde{U}] = 0$.

Remark 4.4. The first hypothesis is often verified by combining $\tilde{L}\tilde{U} + \tilde{B}(x)\tilde{U} \in L^2_{\text{loc}}$, $\tilde{U} \in L^2_{\text{loc}}$ and $\tilde{B} \in L^\infty_{\text{loc}}$.

(2) $[\tilde{A}_1\tilde{U}]$ makes sense since the differential equation implies

$$\partial_1(\tilde{A}_1\tilde{U}^+) \in L^2_{\text{loc}}(]0, \varepsilon[; H^{-1}_{\text{loc}}(\mathbb{R}^d_{t,x'})).$$

With $\tilde{U} \in L^2(]0, \varepsilon[; H^{-1}_{\text{loc}}(\mathbb{R}^d))$ this implies that $\tilde{A}_1\tilde{U}^+ \in C(]0, \varepsilon[; H^{-1/2}_{\text{loc}}(\mathbb{R}^d))$. An analogous result holds for $\tilde{A}_1\tilde{U}^-$. Therefore the traces from both sides and the jump are well-defined elements of $H^{-1/2}_{\text{loc}}$.

(3) is clear on a formal level since if $\tilde{A}_1\tilde{U}$ were discontinuous there would be a $\delta(x_1)$ term from the differential operator \tilde{L} applied to \tilde{U} .

(4) The standard proof based on these remarks is omitted.

We have assumed that the nonzero data are initial values $f^\pm(x)$. By the usual subtraction one can convert the problem to one with homogeneous initial values and right-hand side and inhomogeneous transmission condition. In this way, the determination of \widetilde{W}^\pm is reduced to finding \widetilde{W}^\pm satisfying the inhomogeneous transmission problem

$$\widetilde{L}_1(\partial_t, \partial_x)\widetilde{W}^\pm + \widetilde{B}^\pm\widetilde{W}^\pm = 0, \quad \widetilde{A}_1[\widetilde{W}] = \widetilde{g}, \tag{4.19}$$

where

$$\widetilde{B}^\pm := \begin{pmatrix} \sigma^\pm I_N & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{pmatrix}, \tag{4.20}$$

and $\widetilde{g}(t, x)$ take values in $\text{Range } \widetilde{A}_1$. The unknowns \widetilde{W} and source \widetilde{g} are vectors of length dN .

Theorem 4.3. *Suppose that $L(\partial)$ is a hyperbolic operator nondegenerate with respect to x_1 and that the Bérenger’s doubled operator \widetilde{L} is weakly hyperbolic for $\sigma = \sigma^\pm$. Then, the transmission problem (4.19) with absorption (4.18) satisfies Hersh’s condition.*

Proof. Drop the tildes on the Fourier–Laplace transforms of \widetilde{W} , \widetilde{g} for ease of reading. The transformed problem is

$$\left(\widetilde{A}_1 \frac{d}{dx_1} + \widetilde{L}(\tau, 0, i\eta) \right) \widehat{W}^\pm = 0, \quad \widetilde{A}_1[\widehat{W}] = \widehat{g}. \tag{4.21}$$

The condition of Hersh is that for an arbitrary right-hand side \widehat{g} in $\text{range } \widetilde{A}_1$ this transmission problem has one and only one solution.

Denote by $E_L^\pm(\tau, \eta, \sigma)$ the spaces associated to the Bérenger operator \widetilde{L} with absorption σ . The uniqueness of solutions of (4.21) is equivalent to

$$\widetilde{A}_1 \dot{E}_L^-(\tau, \eta, \sigma^-) \cap \widetilde{A}_1 \dot{E}_L^+(\tau, \eta, \sigma^+) = \{0\}. \tag{4.22}$$

Existence is equivalent to

$$\widetilde{A}_1 \dot{E}_L^-(\tau, \eta, \sigma^-) + \widetilde{A}_1 \dot{E}_L^+(\tau, \eta, \sigma^+) = \text{range } \widetilde{A}_1. \tag{4.23}$$

Part (ii) of Lemma 4.4 implies that

$$\dim(\widetilde{A}_1 \dot{E}_L^-(\tau, \eta, \sigma^-)) + \dim(\widetilde{A}_1 \dot{E}_L^+(\tau, \eta, \sigma^+)) = \dim(\text{range } \widetilde{A}_1),$$

so (4.22) implies (4.23). It remains to prove (4.22).

For $\sigma > 0$, the split Béranger operator \tilde{L} is hyperbolic so for $\text{Re } \tau > \tau_0(\sigma)$ the solutions of

$$\tilde{L}(\tau, d/dx_1, i\eta)\hat{U} = 0 \tag{4.24}$$

are generated by exponentially growing and exponentially decaying solutions. The next lemma identifies these solutions in terms of the corresponding solutions of

$$L(\tau, d/dx_1, i\eta)\hat{V} = 0. \tag{4.25}$$

The result shows that the traces at $x_1 = 0, \hat{E}_L^\pm$, are independent of σ .

Lemma 4.7. For $\sigma > 0, \text{Re } \tau > \tau_0(\sigma), \eta \in \mathbb{R}^{d-1}$,

(i) The map

$$\hat{V}(x_1) \mapsto \left(\frac{\eta_1}{\tau} A_1 \hat{V}((\tau + \sigma)x_1/\tau), \frac{\eta_2}{\tau} A_2 \hat{V}((\tau + \sigma)x_1/\tau), \dots, \frac{\eta_d}{\tau} A_d \hat{V}((\tau + \sigma)x_1/\tau) \right)$$

is an isomorphism from solutions of (4.25) onto the solutions of (4.24).

- (ii) μ is a root of $\det L(\tau, \cdot, i\eta) = 0$ if and only if $\nu = (\tau + \sigma)\mu/\tau$ is a root of $\det \tilde{L}(\tau, \cdot, \eta) = 0$.
- (iii) For the roots in (ii), the real parts of μ and ν have the same sign. In particular, the map in (i) is an isomorphism $E_L^\pm(\tau, \eta) \mapsto E_{\tilde{L}}^\pm(\tau, \eta, \sigma)$.
- (iv) The map $\widehat{W} = (\widehat{W}_1, \dots, \widehat{W}_d) \mapsto \sum_j \widehat{W}_j$ is an isomorphism $E_L^\pm(\tau, \eta, \sigma) \rightarrow E_L^\pm(\tau, \eta)$.

Remark 4.5. In (i) it is important to know that the solutions $\hat{V}(x_1)$ are entire analytic functions of x_1 so it makes sense to evaluate \hat{V} at points off the x_1 -axis. In the literature this is sometimes called a *complex change of variables*. It is only reasonable for analytic solutions. A related idea is used in the Fourier–Laplace analysis for general $\sigma_1(x_1)$ presented in Sec. 4.2.

Proof of Lemma 4.7. (i) If $\hat{U} = (\hat{U}_1, \dots, \hat{U}_d)$ satisfies (4.24), then with $\widehat{W} := \sum_j \hat{U}_j$,

$$A_1 \frac{d\widehat{W}}{dx_1} + (\tau + \sigma)\hat{U}_1 = 0, \quad \tau\hat{U}_j + i\eta_j A_j \widehat{W} = 0, \quad j = 2, \dots, d. \tag{4.26}$$

Multiply the first by τ and the last $d - 1$ by $(\tau + \sigma)$. Sum and then divide by τ to find,

$$A_1 \frac{d\widehat{W}}{dx_1} + \frac{\tau + \sigma}{\tau} L(\tau, 0, i\eta)\widehat{W}^\pm = 0. \tag{4.27}$$

Conversely if \widehat{W} satisfies (4.27) and \hat{U}_j for $j \geq 2$ is defined from \widehat{W} using the last equations in (4.26) and $\hat{U}_1 := \widehat{W} - \sum_{j \geq 2} \hat{U}_j$ then U satisfies (4.24).

The solutions \widehat{W} to (4.27) are exactly the $\widehat{V}((\tau + \sigma)x_1/\tau)$ with \widehat{V} satisfying (4.25). This proves that the mapping in (i) is surjective.

The set of solutions \widehat{V} of (4.25) has dimension $\text{rank } A_1$. The set of solutions of (4.24) has dimension $\text{rank } \widetilde{A}_1 = \text{rank } A_1$ (see (2.3)), so surjectivity implies injectivity.

(ii) and (iv) follow from (i).

(iii) Denote by K the mapping from (i). When $\tau > 0$, $(\tau + \sigma)/\tau$ is also positive and real. Therefore K maps decaying (respectively increasing) solutions to decaying (respectively increasing) solutions. Thus for $\tau > \tau_0$ and real,

$$K(E_L^+(\tau, \eta)) = E_L^+(\tau, \eta, \sigma). \tag{4.28}$$

For all $\text{Re } \tau > \tau_0$, $K(E_L^+(\tau, \eta))$ is a subspace of solutions of (4.24) with dimension equal to $\dim E_L^+(\tau, \eta)$. If (4.28) were violated, $K(E^+(\tau, \eta))$ would contain exponentially growing solutions. If this happened at $\underline{\tau}, \underline{\eta}$ with $\text{Re } \underline{\tau} > \tau_0$, consider the values $\tau(r) = \text{Re } \underline{\tau} + r \text{Im } \underline{\tau}$ for $0 \leq r \leq 1$. For $r = 0$, (4.28) is satisfied while for $r = 1$ it is violated. Let

$$f(r) := \max \left\{ \text{Re } \frac{\tau(r) + \sigma}{\tau(r)} \mu : \det L(\tau(r), \mu, i\eta) = 0, \text{Re } \mu < 0 \right\}.$$

Then $f(0) < 0, f(1) > 0$ and f is continuous, so there is a $0 < \underline{r} < 1$ so that $f(\underline{r}) = 0$. Then for $\tau = \text{Re } \underline{\tau} + \underline{r} \text{Im } \underline{\tau}$ there is a purely imaginary root. This violates the hyperbolicity of \widetilde{L} establishing (4.28). This proves (iii) completing the proof of the lemma. \square

We now finish the proof of Theorem 4.3 by proving (4.22). Lemma 4.7 implies that the spaces of Cauchy data $\dot{E}_{\widetilde{L}}^{\pm}$ are independent of σ . Therefore if (4.22) is violated, then also

$$\widetilde{A}_1 \dot{E}_{\widetilde{L}}^-(\tau, \eta, \sigma^+) \cap \widetilde{A}_1 \dot{E}_{\widetilde{L}}^+(\tau, \eta, \sigma^+) \neq \{0\}.$$

This contradicts part (iv) of Lemma 4.4 for the operator \widetilde{L} with absorption σ^+ . The proof of Hersh’s condition is complete.

In these problems with only one nonzero absorption coefficient σ_1 and $\sigma_1 = 0$ when $x_1 < 0$ one can consider a transmission problem which is only split in $x_1 > 0$. The next result shows that this partially split problem satisfies Hersh’s condition if and only if the fully split problem does.

Introduce the partially split problem (L, R, \mathcal{N}) where

$$\begin{aligned} L &= L_1(\partial), \quad R = \widetilde{L}_1 + \widetilde{B}^+, \quad \text{with } \sigma^+ > 0, \\ \mathcal{N} &:= \left\{ (V, W) : V - \sum_j W_j \in \ker A_1 \right\}, \end{aligned} \tag{4.29}$$

with \widetilde{B}^+ given by (4.20) and the split variable on the right is $\widetilde{W} = (W_1, \dots, W_d)$.

Corollary 4.1. *Suppose that $\sigma_j = 0$ for $j \geq 2$, and $\sigma_1^+ > 0$. Then the partially split Bérenger transmission problem $(L_1, \widetilde{L}_1 + \widetilde{B}^+, \mathcal{N})$ defined by (4.29) satisfies Hersh’s condition if and only if the fully split problem does.*

Proof. Denote by $(V, \widetilde{W}) = (V, W_1, \dots, W_d)$ the variables for the partially split problem and $(\widehat{U}, \widehat{W}) = ((U_1, \dots, U_d), (W_1, \dots, W_d))$ the split variables. If $\widehat{U}(\tau, x_1, \eta), \widehat{W}(\tau, x_1, \eta)$ is an exponentially decaying solution of the split Laplace–Fourier transformed homogeneous transmission problem, then $(\widehat{V}, \widehat{W}) = (\sum_j \widehat{U}_j, \widehat{W})$ is an exponentially decaying solution of the partially split homogeneous transmission problem.

Conversely, if $\widehat{V}(\tau, x_1, \eta), \widehat{W}(\tau, x_1, \eta)$ is a solution of the homogeneous partially split problem, the computation leading to (4.26) shows that

$$\widehat{U}_1 := -\tau^{-1} A_1 \partial_1 \widehat{V}, \quad \widehat{U}_j := -\tau^{-1} i \eta_j A_j \widehat{V}, \quad j \geq 2,$$

is an exponentially decaying solution of the fully split homogeneous transmission problem.

Therefore, if either problem has decaying solutions for η real and $\text{Re } \tau$ arbitrarily large, then so does the other. □

4.1.4. Perfection for Bérenger’s PML with piecewise constant σ_1

Theorem 4.4. *With the hypotheses of Theorem 4.3, the Bérenger transmission problem is perfectly matched. The Bérenger transmission problem that is only split on the right is also perfectly matched.*

Proof. Verify condition (ii) of Theorem 4.2. For $K \in \dot{E}_L^+(\tau, \eta, \sigma^+)$ consider the unique decomposition guaranteed by the Hersh’s condition,

$$(K, 0) = (W^-, W^+) + (F^-, F^+), \tag{4.30}$$

where

$$(W^-, W^+) = ((W_1^-, \dots, W_d^-), (W_1^+, \dots, W_d^+)) \in \mathcal{N},$$

$$(F^-, F^+) \in \dot{E}_L^-(\tau, \eta, \sigma^-) \times \dot{E}_L^+(\tau, \eta, \sigma^+).$$

Perfection is equivalent to $F^- = 0$.

By inspection, one such decomposition (4.30) is given by

$$(K, 0) = (K, K) + (0, -K),$$

where we use the fact from Lemma 4.7 that

$$\dot{E}_L^+(\tau, \eta, \sigma^-) = \dot{E}_L^+(\tau, \eta, \sigma^+).$$

As this decomposition satisfies $F^- = 0$, the proof of the first assertion is complete.

For the partially split case, $K \in \dot{E}_L^+(\tau, \eta)$ has a unique decomposition from the Hersh’s condition,

$$(K, 0) = (W^-, W^+) + (F^-, F^+),$$

with

$$(W^-, W^+) = (W^-, (W_1^+, \dots, W_d^+)) \in \mathcal{N}, \quad (F^-, F^+) \in \dot{E}_L^-(\tau, \eta) \times \dot{E}_L^+(\tau, \eta, \sigma^+).$$

Define

$$W_j^+ := \frac{\eta_j}{\tau} A_j K.$$

Part (i) of Lemma 4.7 implies that $W^+ \in \dot{E}_L^+(\tau, \eta, \sigma^+)$. In addition, $\sum_j W_j^+ = K$ so $(K, W^+) \in \mathcal{N}$.

By inspection

$$(K, 0) = (K, W^+) + (0, -W^+)$$

is the unique Hersh decomposition. Since F^- vanishes for this one, the proof is complete. □

4.1.5. Analytic continuation for Maxwell like systems and Bérenger's plane waves

In this section we investigate Bérenger's method for operators, including the Maxwell system, whose characteristic polynomial is $\tau^p(\tau^2 - |\xi|^2)^q$. For ease of exposition we treat the case $d = 2$ and the explicit operator,

$$L = \partial_t + \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \partial_1 + \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \partial_2. \tag{4.31}$$

Analogous results are valid for the Maxwell system with only slightly more complicated formulas.

For $\text{Re } \tau > 0$ and $\eta \in \mathbb{R}$ there is exactly one root of $\det L_1(\tau, \rho, i\eta) = 0$ with $\text{Re } \rho > 0$ given by

$$\rho = \sqrt{\tau^2 + \eta^2}, \quad \text{Re } \rho > 0. \tag{4.32}$$

The corresponding eigenspace $\dot{E}_L^+(\tau, \eta)$ from (4.17) is spanned by $\Phi(\tau, \eta) = (-i\eta, \tau + \rho)$. If \tilde{L}_1 is the Bérenger operator doubled in the x_1 direction, one has the same roots and $\dot{E}_{\tilde{L}_1}^+$ is spanned by $(\rho A_1 \Phi, i\eta A_2 \Phi)$.

Proposition 4.2. (i) For each η , $\rho(\tau, \eta)$, $\dot{E}_L^-(\tau, \eta)$, and $\dot{E}_L^+(\tau, \eta)$ are holomorphic in $\text{Re } \tau > 0$ with continuous extension to $\text{Re } \tau \geq 0$.

(ii) If $\sigma_1 > 0$, then for $\text{Re } \tau > 0$ and $\eta \in \mathbb{R}$ the equation $\det \tilde{L}(\tau, \nu, i\eta) = 0$ has exactly one root ν with positive real part. It is given by $\nu = (\tau + \sigma_1)\rho/\tau$.

(iii) For $\sigma_1 \geq 0$, the relation (4.12) with \tilde{L}_1 on the left and \tilde{L} on the right is satisfied on $\{\text{Re } \tau \geq 0, \rho \neq 0\}$.

(iv) The mapping $H(\tau, \eta)$ is for each η holomorphic in $\text{Re } \tau > 0$ with continuous extension to $\{\text{Re } \tau \geq 0, \rho \neq 0\}$.

Proof. (i) For $\text{Re } \tau > 0$ there are two roots $\pm\rho$ with ρ from (4.32). One has strictly positive real part and the other strictly negative. Each is holomorphic in $\text{Re } \tau > 0$. Holomorphy for $\Phi(\tau, \eta)$ follows from its expression in terms of ρ . As Φ is a basis for $\dot{E}_L^+(\tau, \eta)$ holomorphy of the latter follows.

So long as the eigenvalues $\pm\rho$ remain apart as $\text{Re } \tau \rightarrow 0$ they and their eigenspaces will be holomorphic. The delicate case is when $\tau^2 + \eta^2 \rightarrow 0$. The limiting points are $(\pm i\eta, \eta)$.

If $\eta = 0$, then $\rho = \tau$ and the eigenspace is $(0, 1)$. Both are continuous up to the boundary.

When $\eta \neq 0$ one has $\rho \rightarrow 0$ so ρ is continuous up to the boundary. Then Φ is continuous up to the boundary and nonvanishing from its expression in terms of ρ . Therefore \dot{E}_L^+ is continuous up to the boundary.

(ii) It suffices to remark that this is an eigenvalue and then to show that the real part is positive. For the latter, compute

$$\frac{\partial}{\partial \sigma} \left(\frac{\tau + \sigma}{\tau} \sqrt{\tau^2 + \eta^2} \right) = \frac{\sqrt{\tau^2 + \eta^2}}{\tau} = \sqrt{1 + \eta^2/\tau^2}.$$

For $\text{Re } \tau > 0$ this has positive real part so the real part of the eigenvalue is increasing as a function of σ so is positive for all $\sigma \geq 0$.

(iii) It suffices to show that (4.11) is valid for such s, η . Suppose $(v, w) = (v_1, v_2, w_1, w_2) \in \dot{E}_{L_1}^- \times \dot{E}_L^+$. We must show that $v_1 + v_2 \neq w_1 + w_2$. Since $(w_1, w_2) \in \dot{E}_L^+$ it follows that $w_1 + w_2 \in \dot{E}_L^+$. Similarly $v_1 + v_2 \in \dot{E}_{L_1}^-$. Thus it suffices to show that $\dot{E}_{L_1}^-$ and \dot{E}_L^+ are uniformly transverse as $\text{Re } \tau \rightarrow 0$. It suffices to show that $(i\eta, \tau + \rho)$ and $(i\eta, \tau - \rho)$ are uniformly independent. This follows from $\rho \neq 0$.

(iv) The holomorphy of H follows from (i). The continuous extension follows from (i) and (iii). □

Since the method is perfectly matched, $H = 0$ for $\text{Re } \tau > 0$. By continuity the map vanishes for purely imaginary $\tau \neq 0$. This shows that for $\{\text{Re } \tau \geq 0\} \setminus 0$, the function equal to

$$e^{i\tau t + \rho(\tau, \eta)x_1 + i\eta x_2} \tilde{\Phi} \quad \text{for } x_1 < 0, \quad \text{and} \quad e^{i\tau t + \rho(\tau, \eta)x_1 + i\eta x_2} e^{-\sigma \rho x_1 / \tau} \tilde{\Phi} \quad \text{for } x_1 > 0,$$

satisfies the Bérenger transmission problem. For $\text{Re } \tau > 0$ these solutions decay (respectively grow) exponentially as $x_1 \rightarrow \infty$ (respectively $x_1 \rightarrow -\infty$). Though such solutions serve to verify perfection they do not look very physical in isolation.

On the other hand, when τ is purely imaginary and not equal to zero, the solution is a bounded plane wave in $x_1 < 0$ and is a plane wave modulated by an exponentially decaying factor in $x_1 > 0$. These are the solutions which Bérenger constructed to show that the method was perfectly matched.

In the language of the analytic objects constructed in the preceding lemma, Bérenger's plane wave solutions show that $H(is, \eta) = 0$ when s is real-valued with

$s^2 > \eta^2$. For η fixed, the function $\tau \mapsto H(\tau, \eta)$ is holomorphic in the right half plane continuous up to the imaginary axis punctured at $\pm i|\eta|$, and vanishes on the boundary interval $\tau = is \in i\mathbb{R}$ with $s^2 > \eta^2$. By Schwarz reflection and analytic continuation, this implies that H vanishes in the right half plane.

In summary, *the computation of Bérenger is actually sufficient to prove perfection for Maxwell's system given the structures provided in this paper.*

Remark 4.6. The perfection argument based on plane waves is not valid in full generality where the objects like \hat{E} and H are analytic in $\text{Re } \tau > \tau_0$ with $\tau_0 > 0$. This is the case, for example, whenever the absorbing layer is amplifying.

4.2. Fourier–Laplace analysis with variable $\sigma_1(x_1)$

Consider the case of only one nonzero $\sigma_1(x_1)$. If \tilde{L} is hyperbolic for one constant value $\underline{\sigma}_1 \neq 0$ the scaling $(t, x) \mapsto (at, ax)$ shows that \tilde{L} is hyperbolic for $\underline{\sigma}_1/a$. Therefore \tilde{L} is hyperbolic for all constant values σ_1 .

The results of Sec. 4.1 will be extended to the case $\sigma_j = 0$ for $j \geq 2$ and variable coefficient $\sigma_1(x_1)$. The Fourier–Laplace transform $\hat{U}(\tau, x_1, \eta)$ of the Bérenger split operator satisfies

$$\tilde{L}(\tau, d/dx_1, \eta)\hat{U} = \hat{F}, \quad -\infty < x_1 < \infty,$$

with variable coefficient $\sigma_1(x_1)$.

The first line of the proof of Lemma 4.7 yields (4.26) with $\sigma = \sigma_1(x_1)$. As in the proof of that lemma one derives (4.27) now with $\sigma = \sigma_1(x_1)$. The important observation is that the x_1 dependence of the coefficient appears only as a scalar prefactor in (4.27). Such equations will be analyzed in the same way as the equations in Lemma 4.2.

4.2.1. *Well-posedness by Fourier–Laplace with variable $\sigma_1(x_1)$*

Theorem 4.5. *Suppose $\sigma_j = 0$ for $j \geq 2$ and $\sigma_1(x_1) \in L^\infty(\mathbb{R})$ is real-valued. Suppose in addition that L is nondegenerate with respect to x_1 , and for one value $\sigma_1 \neq 0, \tilde{L}$ is hyperbolic. Then there is a $\tau_0 > 0$ and m so that for all $\lambda > \tau_0$ and $F \in e^{\lambda t} L^2(\mathbb{R}_t : H^m(\mathbb{R}_{t,x'}^d))$ there is a unique solution $\tilde{U} \in e^{\lambda t} L^2(\mathbb{R}^{d+1})$ to the Bérenger split problem $\tilde{L}\tilde{U} = F$. In addition, there is a constant C independent of F, λ so that*

$$\|e^{-\lambda t}\tilde{U}\|_{L^2(\mathbb{R}^{1+d})} \leq C\|e^{-\lambda t}F\|_{L^2(\mathbb{R}_t; H^m(\mathbb{R}_{t,x'}^d))}. \tag{4.33}$$

Remark 4.7. (1) The condition $\tilde{U} \in e^{\lambda t} L^2$ implies that \tilde{U} tends to zero at $t \rightarrow -\infty$ as does F .

(2) If F is supported in $t \geq t_0$ it follows from (4.33) on sending $\lambda \rightarrow \infty$ that \tilde{U} is supported in $t \geq t_0$.

Proof. The values of the Fourier–Laplace transform of $W = \sum U_j$ are computed from the ordinary differential equation

$$A_1 \frac{d\widehat{W}}{dx_1} + \frac{\tau + \sigma_1(x_1)}{\tau} L_1(\tau, 0, i\eta) \widehat{W} = \widehat{F}. \tag{4.34}$$

As in Lemmas 4.2 and 4.3, transform to the equivalent form,

$$\begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix} \frac{d\widehat{W}}{dx_1} + \frac{\tau + \sigma_1(x_1)}{\tau} \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} \widehat{W} = \widehat{F}, \quad H_{22} \text{ invertible.}$$

Denote the decomposition as $W = (W_I, W_{II})$ and similarly F . The invertibility of H_{22} from Lemma 4.3 yields,

$$\widehat{W}_{II} = H_{22}^{-1}(\widehat{F}_{II} - H_{21}\widehat{W}_I). \tag{4.35}$$

It suffices to find \widehat{W}_I which is determined from,

$$\begin{aligned} \frac{d\widehat{W}_I}{dx_1} + \frac{\tau + \sigma_1(x_1)}{\tau} M(\tau, \eta) \widehat{W}_I &= \widehat{G}, \quad M(\tau, \eta) := H_{11} - H_{12}H_{22}^{-1}H_{21}, \\ \widehat{G} &:= \widehat{F}_I + H_{22}^{-1}\widehat{F}_{II}. \end{aligned}$$

The hyperbolicity of \widetilde{L} implies that M has no purely imaginary eigenvalues. Correspondingly there is the decomposition, into the spectral parts with positive and negative imaginary parts,

$$W_I = W_I^+ + W_I^-, \quad G = G^+ + G^-, \quad M = M^+ \oplus M^-.$$

For σ constant, part (iii) of Lemma 4.7 (using the hyperbolicity of \widetilde{L}) implies that for $\text{Re } \tau$ sufficiently large (depending on σ), one has the spectral decomposition,

$$\frac{\tau + \sigma}{\tau} M(\tau, \eta) = \frac{\tau + \sigma}{\tau} M(\tau, \eta)^+ \oplus \frac{\tau + \sigma}{\tau} M(\tau, \eta)^-$$

corresponding to spectra with positive and negative real parts.

Lemma 4.8. *If $g(x_1)$ satisfies $dg(x_1)/dx_1 = \sigma_1(x_1)$. Then*

$$\left(\frac{d}{dx_1} + M \right) \left(e^{g(x_1)M/\tau} \widehat{U}_I \right) = e^{g(x_1)M/\tau} \left(\frac{d}{dx_1} + \frac{\tau + \sigma_1(x_1)}{\tau} M \right) \widehat{U}_I.$$

Proof of Lemma 4.8. Since $(de^{gM/\tau}\widehat{U}_I)/dx_1 = e^{gM/\tau}(g'M\widehat{U}_I/\tau + d\widehat{U}_I/dx_1)$ one has

$$\left(\frac{d}{dx_1} + M \right) \left(e^{g(x_1)M/\tau} \widehat{U}_I \right) = e^{gM/\tau} \left(\frac{d}{dx_1} \widehat{U}_I + \left(\frac{dg/dx_1 M}{\tau} + \frac{\tau M}{\tau} \right) \widehat{U}_I \right),$$

proving the desired identity. □

Therefore

$$\widehat{W}_I = e^{-g(x_1)M/\tau} \left(\frac{d}{dx_1} + M \right)^{-1} \left(e^{g(x_1)M/\tau} \widehat{G} \right).$$

The unique L^1 fundamental solution of $\partial_1 + M$ is equal to,

$$e^{-x_1 M^+} \chi_{[0, \infty[}(x_1) + e^{-x_1 M^-} \chi_{]-\infty, 0]}(x_1).$$

Therefore,

$$\begin{aligned} e^{gM/\tau} \widehat{W}_I^+ &= (e^{-x_1 M^+} \chi_{[0, \infty[}(x_1)) * (e^{gM/\tau} \widehat{G})^+, \\ e^{gM/\tau} \widehat{W}_I^- &= (e^{-x_1 M^-} \chi_{]-\infty, 0]}(x_1)) * (e^{gM/\tau} \widehat{G})^-. \end{aligned}$$

The kernel of the integral operator mapping \widehat{G}^+ to \widehat{W}_I^+ is equal to,

$$\exp\left(- (x_1 - y_1) \left[\frac{\tau + (g(x_1) - g(y_1))/(x_1 - y_1)}{\tau} M(\tau, \eta)^+ \right] \right) \chi_{x_1 \geq y_1}. \tag{4.36}$$

Lemma 4.9.

$$\begin{aligned} \exists \tau_0 = \tau_0(\mu), \quad \forall \operatorname{Re} \tau \geq \tau_0, \quad \forall \eta \in \mathbb{R}^{d-1}, \quad \forall \sigma \in [-\mu, \mu], \\ \operatorname{spec} \frac{\tau + \sigma}{\tau} M^+(\tau, \eta) \subset \{\operatorname{Re} z > 0\}. \end{aligned}$$

Proof of Lemma 4.9. Part (iii) of Lemma 4.7 allows one to choose τ_1 so that for $\sigma = \mu$ one has the desired conclusion for $\operatorname{Re} \tau > \tau_1$. Then for $\lambda \in \operatorname{spec} M(\tau, \eta)^+$ one has

$$\operatorname{Re} \lambda > 0, \quad \operatorname{Re} \left(1 + \frac{\mu}{\tau} \right) \lambda = \operatorname{Re} \frac{\tau + \mu}{\tau} > 0.$$

For $0 \leq \sigma \leq \mu$ write $\sigma = a + b\mu$ with non-negative a, b summing to 1. It follows that $\operatorname{Re}(1 + \mu/\tau)\lambda > 0$. This proves that τ_1 suffices to treat the non-negative values $0 \leq \sigma \leq \mu$.

Choosing τ_2 for $\sigma = -\mu$, that value suffices for $-\mu \leq \sigma \leq 0$. Set τ_0 equal to the maximum of τ_1 and τ_2 . □

The Seidenberg–Tarski Theorem 2.3 shows that the absolute values of the real parts of the eigenvalues of $M(\tau, \eta)$ are bounded below by $C(|\tau| + |\eta|)^{-N}$ for some N . And also that the spectral decomposition $V \mapsto (V^+, V^-)$ and its inverse are both bounded polynomially in $|\tau, \eta|$. More generally for τ, η, μ, σ as above,

$$\operatorname{spec} \frac{\tau + \sigma}{\tau} M^+(\tau, \eta) \subset \{\operatorname{Re} z > C(|\tau| + |\eta|)^{-N}\}.$$

Taking $\mu := \|f\|_{L^\infty}$ one finds that for all x_1, y_1 , the matrix

$$\frac{\tau + (g(x_1) - g(y_1))/(x_1 - y_1)}{\tau} M(\tau, \eta)^+ \tag{4.37}$$

has spectrum in

$$\{\operatorname{Re} z > C(|\tau| + |\eta|)^{-N}\}, \quad C > 0.$$

The elements of the matrix (4.37) are bounded above polynomially in $|\tau, \eta|$. Therefore the kernel (4.36) is bounded above by

$$|\tau, \eta|^N \exp(-c(x_1 - y_1)/|\tau, \eta|^N) \chi_{x_1 \geq y_1}, \quad c > 0. \tag{4.38}$$

This is proved using Schur’s theorem to reduce M^\pm to upper triangular form by orthogonal transformations of the spectral subspaces. Then solve the differential equation $X' + M^+ X = 0$ by back substitution to prove $\|\exp(\rho M^+)\| \leq C|\tau, \eta|^p e^{-c\rho/|\tau, \eta|^N}$.

The operator with kernel (4.38) is convolution by an element of $L^1(\mathbb{R})$ whose L^1 norm grows polynomially in $|\tau, \eta|$. By Young’s theorem one concludes that the operator with kernel (4.36) has norm in $\text{Hom}(L^2(\mathbb{R}))$ which grows at most polynomially in $|\tau, \eta|$.

There is an entirely analogous estimate for the expression for the spectrum with negative real part.

Therefore,

$$\begin{aligned} \|\widehat{W}_I(\tau, x_1, \eta)\|_{L^2(\mathbb{R})} &\leq C_1(1 + |\tau| + |\eta|)^N \|\widehat{G}(\tau, x_1, \eta)\|_{L^2(\mathbb{R})} \\ &\leq C_2(1 + |\tau| + |\eta|)^N \|\widehat{F}(\tau, x_1, \eta)\|_{L^2(\mathbb{R})}. \end{aligned}$$

A similar estimate for \widehat{W}_{II} follows from (4.35). Estimates for \widehat{U}_j follow from the second equation in (4.26). Plancherel’s theorem then implies (4.33), proving the existence part of well-posedness.

Uniqueness is proved by a duality argument of Hölmgren type using existence (backward in time) for the adjoint differential operator (details omitted).

4.2.2. Perfection for Bérenger’s PML with variable coefficient $\sigma_1(x_1)$

Lemma 4.10. *Suppose that A, M satisfy the hypothesis of Lemma 4.2 with \mathbb{G} and $\widetilde{M} \in \text{Hom } \mathbb{G}$ are from that lemma. Suppose in addition that $f \in L^\infty_{\text{loc}}(\mathbb{R}; \mathbb{C})$ and g is the unique solution of*

$$\frac{dg}{dx_1} = f, \quad g(0) = 0, \quad \text{so } g(x_1) = \int_0^{x_1} f(s) ds.$$

Then for $\gamma \in \mathbb{G}$ the unique solution of the equivalent initial value problems for the \mathbb{G} -valued function U ,

$$A \frac{dU}{dx_1} + f(x_1)MU = 0, \quad \text{equivalently,} \quad \frac{dU}{dx_1} + f(x_1)\widetilde{M}U = 0, \quad U(0) = \gamma,$$

is

$$U(x_1) = e^{-g(x_1)\widetilde{M}}\gamma.$$

Proof. Compute using the differential equation,

$$\frac{d}{dx_1}[e^{g(x_1)\widetilde{M}}U] = e^{g(x_1)\widetilde{M}} \left[\left(\frac{dg}{dx_1} \right) \widetilde{M} + \frac{dU}{dx_1} \right] = e^{g(x_1)\widetilde{M}} [f\widetilde{M} - f\widetilde{M}] = 0.$$

The lemma follows. □

The next result shows that when the Béranger split problem with absorption $\sigma(x_1)$ defines a stable time evolution, then the problem is perfectly matched. Either the split problem is ill posed, or it is well-posed and perfect.

Theorem 4.6. *Suppose that $\sigma_1(x) \in L^\infty(\mathbb{R})$ has support in $[0, \rho]$ for some $\rho > 0$, that $\sigma_j = 0$ for $j \neq 1$, and that the operator \tilde{L} with these absorptions is nondegenerate with respect to x_1 and defines a weakly well-posed time evolution. Then, the \tilde{L} evolution is perfectly matched in the sense that for $F \in C_0^\infty(\{t > 0\} \cap \{x_1 < 0\})$ the solutions \tilde{U} and \tilde{U}' with and without absorptions respectively,*

$$\tilde{L}_1 \tilde{U} = F, \quad \tilde{U}|_{t \leq 0} = 0 \quad \text{and} \quad \tilde{L} \tilde{U}' = F, \quad \tilde{U}'|_{t \leq 0} = 0$$

satisfy

$$\tilde{U}|_{x_1 < 0} = \tilde{U}'|_{x_1 < 0}.$$

Proof. Denote by \hat{U} and \hat{U}' the Fourier–Laplace transforms. The functions are characterized by

$$\tilde{L}_1(\tau, d/dx_1, \eta) \hat{U} = \hat{F} \quad \text{and} \quad \tilde{L}(\tau, d/dx_1, \eta) \hat{U}' = \hat{F},$$

both required to decay exponentially as $|x_1| \rightarrow \infty$. The strategy is to construct a solution of the problem defining \hat{U}' from the solution \hat{U} .

The equations for $\widehat{W} = \sum_j \hat{U}_j$ and $\widehat{W}' = \sum_j \hat{U}'_j$ in $x_1 \geq 0$ have the form

$$A_1 \frac{d\widehat{W}}{dx_1} + M\widehat{W} = 0, \quad A_1 \frac{d\widehat{W}'}{dx_1} + \frac{\tau + \sigma_1(x_1)}{\tau} M\widehat{W}' = 0.$$

Lemma 4.10 applies with $f(x) := (\tau + \sigma(x_1))/\tau$.

Define g as in that lemma. Set $\hat{V} = \widehat{W}$ in $x_1 \leq 0$. For $x_1 \geq 0$ define

$$\hat{V} := e^{-g(x_1)M} \widehat{W}(\tau, 0, \eta).$$

The resulting function satisfies the differential equation required of \widehat{W}' . In addition since $e^{-g(x_1)M}$ is independent of x_1 for $x_1 \geq \rho$, \hat{V} decays as rapidly as \widehat{W} . Therefore \hat{V} satisfies the conditions uniquely determining \widehat{W}' . Therefore $\hat{V} = \widehat{W}'$, and $\widehat{W}'|_{x_1 \leq 0} = \widehat{W}|_{x_1 < 0}$. Use (4.26) to recover \hat{U}, \hat{U}' from $\widehat{W}, \widehat{W}'$ shows that $\tilde{U}'|_{x_1 < 0} = \tilde{U}|_{x_1 < 0}$, proving perfection. \square

Example 4.2. (1) If $L_1(0, \partial_x)$ is elliptic then Corollary 3.2 shows that the evolution of \tilde{L} is strongly well-posed. This includes the case of anisotropic wave equations for which the layer is amplifying showing that perfection is not at all inconsistent with amplification.

(2) For the Maxwell equations and $\sigma_1(x_1) \in W^{2,\infty}(\mathbb{R})$ well-posedness is proved in the remark following Theorem 3.3 and we deduce perfection.

5. Plane Waves, Geometric Optics, and Amplifying Layers

This section includes a series of ideas all related to plane waves and short wavelength asymptotic solutions of WKB type. We first recall the derivation of such solutions from exact plane wave solutions by Fourier synthesis. Then we review the construction of short wavelength asymptotic expansions. These are then applied to examine the proposed absorption by the σ_j . In many common cases the supposedly absorbing layers lead to asymptotic solutions which grow in the layer. Related phenomena are studied by Hu, and Becache, Fauqueux, Joly [14, 5]. For the Maxwell equations for which the PML were designed, the layers are not amplifying. At the end of Sec. 5.3, situations where the amplification does not occur are identified.

5.1. Geometric optics by Fourier synthesis

When the coefficient σ vanishes identically, both L and \tilde{L} are homogeneous constant coefficient systems. When $(\underline{\tau}, \underline{\xi})$ is a smooth point of the characteristic variety, denote by $\tau = \tau(\xi)$ the smooth parametrization, and $\Pi_L(\tau, \xi)$ and $\Pi_{\tilde{L}}(\tau, \xi)$ the associated spectral projections for $\xi \approx \underline{\xi}$, see (2.8). The function $\tau(\xi)$ is homogeneous of degree 1, while the projectors are homogeneous of degree 0. The next argument works equally well for L and \tilde{L} .

For $G(\xi) \in C_0^\infty(\mathbb{R}^d)$ construct exact solutions for $0 < \varepsilon \ll 1$,

$$U^\varepsilon(t, x) := \int e^{i(\xi \cdot x + \tau(\xi)t)} \Pi_L(\tau, \xi) G(\xi - \underline{\xi}/\varepsilon) d\xi.$$

Make the change of variable

$$\xi - \underline{\xi}/\varepsilon := \zeta, \quad \xi = (\underline{\xi} + \varepsilon\zeta)/\varepsilon,$$

and extract the rapidly oscillating term $e^{i(\underline{\xi} \cdot x + \tau t)/\varepsilon}$ to find,

$$\begin{aligned} U^\varepsilon(t, x) &:= e^{i(\underline{\xi} \cdot x + \tau t)/\varepsilon} \int e^{i((\tau(\underline{\xi} + \varepsilon\zeta) - \tau(\underline{\xi}))t/\varepsilon + \zeta x)} \Pi_L(\tau(\underline{\xi} + \varepsilon\zeta), \underline{\xi} + \varepsilon\zeta) G(\zeta) d\zeta \\ &:= e^{i(\underline{\xi} \cdot x + \tau t)/\varepsilon} a(\varepsilon, t, x). \end{aligned} \tag{5.1}$$

Expanding in ε and keeping just the leading term yields the principal term in the geometric optics approximation

$$U^\varepsilon \approx e^{i(\underline{\xi} \cdot x + \tau t)/\varepsilon} \int e^{i(x \cdot \zeta - \mathbf{v}(\underline{\xi}) \cdot \zeta t)} \Pi_L(\underline{\tau}, \underline{\xi}) G(\zeta) d\zeta, \quad \mathbf{v}(\underline{\xi}) := -\partial_\xi \tau(\underline{\xi}).$$

One has

$$U^\varepsilon \approx e^{i(\underline{\xi} \cdot x + \tau t)/\varepsilon} a_0(x - \mathbf{v}(\underline{\xi})t), \quad a_0(x) := \int e^{ix \cdot \zeta} \Pi_L(\underline{\tau}, \underline{\xi}) G(\zeta) d\zeta.$$

A complete Taylor expansion yields the corrected approximations which satisfy the equation with an error $O(\varepsilon^N)$ for all N . We write $O(\varepsilon^\infty)$ for short. This yields infinitely accurate solutions,

$$U^\varepsilon(t, x) := e^{i(\underline{\xi}x + \underline{\tau}t)/\varepsilon} a(t, x, \varepsilon), \quad a(t, x, \varepsilon) \sim a_0(x - \mathbf{v}t) + \varepsilon a_1(t, x) + \dots \quad (5.2)$$

If 0 is a semisimple eigenvalue of $L(\tau, \xi)$, and $\Phi_0 \in \text{Ker } L(\tau, \xi) \setminus \{0\}$ of dimension 1, then the leading amplitude a_0 in the case of (1.1) (respectively (1.5)) is of the form

$$\alpha(t, x)\Phi_0, \quad \left(\text{respectively } \alpha(t, x) \left(\frac{\xi_1}{\tau} A_1 \Phi_0, \dots, \frac{\xi_d}{\tau} A_d \Phi_0 \right) \right)$$

with scalar-valued amplitude α satisfying,

$$(\partial_t + \mathbf{v} \cdot \partial_x)\alpha = 0.$$

This shows that α is constant on the rays which are lines with velocity equal to the group velocity $\mathbf{v}(\xi) := -\partial_\xi \tau(\xi)$.

For $g \in C_0^\infty \setminus 0$ the solutions do not have compact spatial support. This weakness is easily overcome. Choose $\chi \in C_0^\infty(\mathbb{R}^d)$ with $\chi = 1$ on a neighborhood of the origin. For $g \in \mathcal{S}(\mathbb{R}^d)$, define exact solutions by cutting off the integrand outside the domains of definition of $\tau(\xi)$ and $\Pi_L(\tau(\xi), \xi)$,

$$u^\varepsilon(t, x) := \int e^{i(\xi x + \tau(\xi)t)} \Pi_L(\tau(\xi), \xi) g(\xi - \underline{\xi}/\varepsilon) \chi(\sqrt{\varepsilon}(\xi - \underline{\xi}/\varepsilon)) d\xi. \quad (5.3)$$

The analysis above applies with the only change being the initial values. In the preceding case, these values were equal to the transform of $\Pi_L(\tau(\xi), \xi)g(\xi - \underline{\xi}/\varepsilon)$ and in the present case they are infinitely close to that quantity,

$$u^\varepsilon(0, x) = \int e^{ix \cdot \xi} \Pi_L(\tau(\xi), \xi) g(\xi - \underline{\xi}/\varepsilon) d\xi + O(\varepsilon^\infty).$$

This yields infinitely accurate approximate solutions (2.5) which have support in the tube of rays with feet in the support of $\int e^{ix \cdot \xi} \Pi_L(\tau(\xi), \xi) g(\xi) d\xi$.

5.2. Geometric optics with variable coefficients

The Fourier transform method of the preceding sections is limited to problems with constant coefficients. In this section the WKB method which works for variable coefficients is introduced. It will also serve for the analysis of reflected waves.

Let \mathcal{L} be the general operator in (2.1). Fix $(\tau, \xi) \in \text{Char } \mathcal{L}$ and seek asymptotic solutions

$$U^\varepsilon \sim e^{iS/\varepsilon} \sum_{j=0}^{+\infty} \varepsilon^j a_j(t, x), \quad \text{with the phase } S(t, x, \xi) = t\tau + x\xi. \quad (5.4)$$

More precisely we construct smooth functions $a_j(t, x)$ with $\text{supp } a_j \cap ([0, T] \times \mathbb{R}^d)$ compact so that if

$$a(t, x, \varepsilon) \sim \sum_{j=0}^{\infty} \varepsilon^j a_j(t, x),$$

in the sense of Taylor series at $\varepsilon = 0$, and $\text{supp } a \cap ([0, T] \times \mathbb{R}^d \times]0, \varepsilon])$ is compact, then

$$U^\varepsilon := e^{iS/\varepsilon} a(t, x, \varepsilon)$$

satisfies for all s, N ,

$$\|\mathcal{L}U^\varepsilon\|_{H^s([0, T] \times \mathbb{R}^d)} = O(\varepsilon^N).$$

In this case we say that (5.4) is an *infinitely accurate approximate solution*. The next result recalls some facts about such solutions.

Theorem 5.1. *Suppose Problem (2.1) is hyperbolic, $(\tau, \xi) \in \text{Char}(\mathcal{L})$ satisfies the smooth variety hypothesis, and 0 is a semisimple eigenvalue of $\mathcal{L}(\tau, \xi)$.*

(i) *If the coefficients a_j satisfy the recursion relation*

$$a_0(t, x) \in \text{Ker } \mathcal{L}_1(\tau, \xi), \tag{5.5a}$$

$$\forall j \geq 0, \quad i\mathcal{L}_1(\tau, \xi)a_{j+1}(t, x) + \mathcal{L}(\partial_t, \partial_x)a_j(t, x) = 0, \tag{5.5b}$$

then (5.4) is an infinitely accurate approximate solution of (2.1).

(ii) *If $g_j = \Pi_{\mathcal{L}}(\tau, \xi)g_j \in C_0^\infty(\mathbb{R}_x^d)$ are supported in a fixed compact K , then there is one and only one family of a_j satisfying (5.5) together with the initial conditions, $\Pi_{\mathcal{L}}(\tau, \xi)a_j(0, \cdot) = g_j$ and the polarization $\Pi_{\mathcal{L}}(\tau, \xi)a_0 = a_0$. They have support in the tube of rays with feet in K and speed of propagation $\mathbf{v}(\xi) = -\partial_\xi \tau(\xi)$.*

(iii) *The principal term a_0 is a solution of the transport equation*

$$\partial_t a_0 + \mathbf{v}(\xi) \cdot \partial_x a_0 + \Pi_{\mathcal{L}}(\tau, \xi)\mathcal{B}(x)\Pi_{\mathcal{L}}(\tau, \xi)a_0 = 0. \tag{5.6}$$

Proof. For simplicity note $\Pi_{\mathcal{L}} := \Pi_{\mathcal{L}}(\tau, \xi)$ when no ambiguity is to be feared.

Equations (5.5) are obtained by injecting U^ε in (2.1), to find an expression $\sim e^{iS/\varepsilon} \sum_{j \geq 0} \varepsilon^j w_j(t, x)$. In order that the w_j vanish it is necessary and sufficient that Eqs. (5.5) are satisfied.

Next examine the leading order terms to find the relations determining a_0 . Projecting the case $j = 0$ of (5.5) onto $\text{Ker } \mathcal{L}_1$ yields,

$$\Pi_{\mathcal{L}} \left(\partial_t + \sum_{l=1}^d \mathcal{A}_l \partial_l + \mathcal{B}(x) \right) a_0 = 0.$$

This yields a first-order system satisfied by $a_0 = \Pi_{\mathcal{L}} a_0$,

$$\partial_t a_0 + \sum_{l=1}^d \Pi_{\mathcal{L}} \mathcal{A}_l \Pi_{\mathcal{L}} \partial_l a_0 + \Pi_{\mathcal{L}} \mathcal{B}(x) \Pi_{\mathcal{L}} a_0 = 0. \tag{5.7}$$

The leading order part of this equation is a scalar transport operator. To see this differentiate $\mathcal{L}_1(\tau(\xi), \xi)\Pi_{\mathcal{L}}(\tau(\xi), \xi) = 0$ with respect to ξ_l to find

$$\left(\mathcal{A}_l + \frac{\partial\tau(\xi)}{\partial\xi_l}\text{Id}\right)\Pi_{\mathcal{L}}(\tau(\xi), \xi) + \mathcal{L}_1(\tau(\xi), \xi)\frac{\partial}{\partial\xi_l}(\Pi_{\mathcal{L}}(\tau(\xi), \xi)) = 0.$$

Multiplying on the left by $\Pi_{\mathcal{L}}(\tau(\xi), \xi)$ eliminates the second term yielding,

$$\Pi_{\mathcal{L}}\mathcal{A}_l\Pi_{\mathcal{L}} + \frac{\partial\tau(\xi)}{\partial\xi_l}\Pi_{\mathcal{L}} = 0.$$

Injecting this in (5.7) yields (5.6).

In order to compute the coefficients recursively, multiply (5.5b) on the left by the partial inverse $Q_{\mathcal{L}}(\tau, \xi)$, using the identity in (2.7), to obtain for $j \geq 1$,

$$(I - \Pi_{\mathcal{L}})a_j = iQ_{\mathcal{L}}\mathcal{L}(\partial_t, \partial_x)a_{j-1}. \tag{5.8}$$

Projecting (5.5b) on the kernel yields,

$$\Pi_{\mathcal{L}}\mathcal{L}(\partial_t, \partial_x)a_j = 0.$$

Writing a_j as

$$a_j = \Pi_{\mathcal{L}}a_j + (I - \Pi_{\mathcal{L}})a_j, \quad \text{yields } \Pi_{\mathcal{L}}\mathcal{L}(\partial_t, \partial_x)\Pi_{\mathcal{L}}a_j = -\Pi_{\mathcal{L}}\mathcal{L}(\partial_t, \partial_x)(I - \Pi_{\mathcal{L}})a_j.$$

This is again a transport equation, but with a right-hand side,

$$\partial_t\Pi_{\mathcal{L}}a_j + \mathbf{v} \cdot \partial_x\Pi_{\mathcal{L}}a_j + \Pi_{\mathcal{L}}B\Pi_{\mathcal{L}}a_j = -\Pi_{\mathcal{L}}\mathcal{L}(\partial_t, \partial_x)(I - \Pi_{\mathcal{L}})a_j. \tag{5.9}$$

(5.8) and (5.9) permit to calculate the coefficients recursively, knowing the initial values. □

Next apply the above algorithm to the PML operator \tilde{L} . Fix $(\tau, \xi) \in \text{Char } L$ and seek asymptotic solutions

$$\tilde{U}^\varepsilon \sim e^{iS/\varepsilon} \sum_{j=0}^{+\infty} \varepsilon^j \tilde{a}_j(t, x), \quad \text{with the phase } S(t, x) = t\tau + x \cdot \xi. \tag{5.10}$$

Corollary 5.1. *Suppose Problem (1.1) is strongly well-posed, $(\tau, \xi) \in \text{Char } L$ satisfies the smooth variety hypothesis, and 0 is a semisimple eigenvalue of $L(\tau, \xi)$.*

(i) *If the coefficients \tilde{a}_j satisfy the recursion relation*

$$\tilde{a}_0(t, x) \in \text{Ker } \tilde{L}_1(\tau, \xi), \tag{5.11a}$$

$$\forall j \geq 0, \quad (I - \Pi_{\tilde{L}})\tilde{a}_j(t, x) = iQ_{\tilde{L}}\tilde{L}(\partial_t, \partial_x)\tilde{a}_{j-1}(t, x), \tag{5.11b}$$

$$\partial_t\Pi_{\tilde{L}}\tilde{a}_j + \mathbf{v} \cdot \partial_x\Pi_{\tilde{L}}\tilde{a}_j + \beta(x)\Pi_{\tilde{L}}\tilde{a}_j = -\Pi_{\tilde{L}}\tilde{L}(\partial_t, \partial_x)(I - \Pi_{\tilde{L}})\tilde{a}_j, \tag{5.11c}$$

then (5.10) is an infinitely accurate approximate solution of (1.5).

- (ii) If $\tilde{g}_j(x) = \Pi_{\tilde{L}}\tilde{g}_j \in C_0^\infty(\mathbb{R}^d)$ are supported in a fixed compact K , then there is one and only one family of \tilde{a}_j satisfying (5.11) together with the initial conditions, $\Pi_L(\tau, \xi)\tilde{a}_j(0, x) = \tilde{g}_j$ and the polarization $\Pi_L(\tau, \xi)\tilde{a}_0 = \tilde{a}_0$. They have support in the tube of rays with feet in K and speed of propagation $\mathbf{v} = -\partial_\xi\tau(\xi)$.
- (iii) The principal term \tilde{a}_0 is a solution of the transport equation

$$\partial_t\tilde{a}_0 + \mathbf{v} \cdot \partial_x\tilde{a}_0 + \beta(x)\tilde{a}_0 = 0, \quad \text{with } \beta(x) = \sum_{l=1}^d \frac{\sigma_l(x_l)\xi_l}{\tau(\xi)} \frac{\partial\tau(\xi)}{\partial\xi_l}. \quad (5.12)$$

Proof. We need only identify the constant term in (5.6). Use the form of the projector given in Proposition 2.1, to obtain

$$\Pi_{\tilde{L}}B(x)\Pi_{\tilde{L}} = \beta(x)\Pi_{\tilde{L}}. \quad \square$$

5.3. Amplifying layers

The coefficient $\sigma_1(x_1) \geq 0$ is introduced with the idea that waves will be damped in the layer. In this section, we show that sometimes the anticipated decay is not achieved, and waves may be amplified. This was observed in [5]. The authors analyzed the phenomenon for σ constant in the layer. They showed that in an infinite layer solutions can in certain cases grow infinitely large. We present a related analysis using WKB solutions which has three advantages:

- (1) The analysis is valid for variable coefficients $\sigma_1(x_1)$ which corresponds to common practice.
- (2) The growth is seen immediately and not expressed in terms of large time asymptotics.
- (3) The analysis in [5] was in part restricted to $d = 2$ and eigenvectors of multiplicity one. We remove these restrictions.

It is because of (2) that we choose not to follow the authors of [5] in calling this phenomenon instability.

Theorem 5.2. *Suppose $(\tau, \xi) \in \text{Char}(L)$ satisfies the smooth variety hypothesis and $\beta(x)$ is as in (5.12). Suppose in addition that there is an interval on a ray*

$$\Gamma := \{(0, \underline{x}) + t(1, -\partial_\xi\tau(\xi)), 0 \leq t \leq t_0\}, \quad \text{so that } \int_0^{t_0} \beta(\underline{x} - t\partial_\xi\tau(\xi))dt < 0.$$

Then the corresponding WKB solution grows in the layer.

Proof. The solution of the transport equation (5.12) evaluated on Γ is

$$\tilde{a}_0(t_0, x) = \exp\left(-\int_0^{t_0} \beta(\underline{x} - s\partial_\xi\tau(\xi))ds\right) \tilde{a}_0(0, \underline{x} + t_0\partial_\xi\tau(\xi)).$$

The exponential is strictly greater than 1, so

$$|\tilde{a}_0(t_0, \underline{x} + t_0 \partial_\xi \tau(\xi))| > |a_0(0, \underline{x})|. \quad \square$$

Example 5.1. (No amplification for Maxwell/D'Alembert) If the dispersion relation is $\tau^2 = |\xi|^2$ and $\sigma \geq 0$, then there is no amplification since

$$\beta = \sum_{j=1}^d \sigma_j(x_j) \frac{\xi_j^2}{\|\xi\|^2} \geq 0.$$

Example 5.2. (Amplification is common) For the dispersion relation $\tau^2 = q(\xi)$ where q is a positive definite quadratic form so that the ξ axes are not major and minor axes of the ellipse $q = 1$, there are always $\tau > 0, \xi$ so that x_1 layers with $\sigma_1 > 0$ are amplifying ([5]). There are two lines on $\{\tau = q(\xi)^{1/2}\}$ where $\partial q / \partial \xi_1 = 0$. The half cone on which $\partial q / \partial \xi_1 < 0$ corresponds to rays on which x_1 is increasing so they enter a layer $x_1 > 0$. The half cone $\{\partial q / \partial \xi_1 < 0\}$ is divided into two sectors by the plane $\xi_1 = 0$. The sector on which $\xi_1 > 0$ (respectively $\xi_1 < 0$) corresponds to growing (respectively decaying) solutions (see Fig. 1, left). This example shows that amplification is very common. Consequently for the dispersion relation $\tau^2 = q(\xi)$ it is wise to align coordinates along the major and minor axes of the ellipse to avoid amplification. However, if $(\tau^2 - q_1)(\tau^2 - q_2)$ divides the characteristic polynomial and the axes of q_1 and q_2 are distinct from each other, then no linear change of coordinates can avoid amplification in the layer.

A second example from [10] is the linearized compressible Euler equation with nonzero background velocity $(c, 0), c > 0$ for which amplified wave numbers at a right-hand boundary are indicated in bold in Fig. 1, right.

Summary. There is no amplification when the characteristic polynomial is a product of factors τ and $\tau^2 - q$ where q is a positive definite quadratic form with axes of inertia parallel to the coordinate axes. This includes the cases of Maxwell's equations in vacuum, for which the method was developed by Bérenger, the linearized Euler equations about the stationary state, and the linear isotrope elasticity equations. For these the quadratic forms q are multiples of $|\xi|^2$.

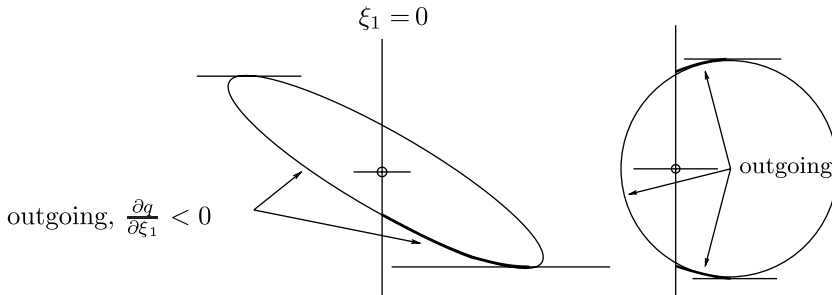


Fig. 1. Amplified outgoing wave numbers in bold.

Example 5.3. (Methods related to Bérenger, continued) For the model developed in Sec. 3.6 for the Maxwell equations, one can compute

$$\beta = 2 \frac{\xi_2^2 + \xi_3^2}{\xi_1^2} \sum_{j=1}^3 \sigma_j(x_j) \xi_j^2 > 0.$$

Thus, this model has exactly the same good properties as Bérenger’s, and is strongly well-posed. For Maxwell equations, it is therefore an attractive alternative. The advantage is twofold. The system with the auxiliary variable P is very compact. And it is strongly well-posed, even for discontinuous σ .

On the surface this result sounds almost too good to be true. However the Bérenger system in the case of Maxwell’s equations has almost exactly the same structure. The energy method proof when $\sigma_j'' \in L^\infty$ shows that there is a large vector \mathbb{V} consisting of the components of U together with differential operators $P_\alpha(D)$ applied to U and a strongly well-posed equation for \mathbb{V} . This means that if one were to introduce the additional variables in \mathbb{V} one obtains a system with some of the desirable properties of SPML (strong PML). However, the SPML reduction is much more compact, and, has a good energy estimate even when σ is discontinuous. The extension of this strategy to other equations is not straightforward. For elastodynamic models, see [28].

6. Harmoniously Matched Layers

This section introduces a new absorbing layer method. It is based on the following strategy. Start with an operator $L = L_1(\partial)$ on the left and consider a smart layer on the right

$$R(t, x, \partial) = L_1(\partial) + C(t, x), \quad C = \sigma(x_1)(\pi_+(A_1) + \nu\pi_-(A_1)), \quad \text{supp } \sigma \subset [0, \infty[, \tag{6.1}$$

generalizing (1.2). This method is embedded in a family of absorbing layers parametrized by $\mu \geq 0$,

$$R^\mu := L_1 + \mu C. \tag{6.2}$$

The method is nonreflective when $\mu = 0$ and is both reflective and dissipative for $\mu > 0$. When σ is discontinuous, the leading order reflection coefficient for wave packets of amplitude 1 oscillating as $e^{i(\tau t + x\xi)/\varepsilon}$ is of the form $\varepsilon\mu r(\tau, \xi)$. The leading order reflections can be removed by an extrapolation method using two values of μ . This simultaneously removes the leading reflections at all angles of incidence. We call the resulting method the *harmoniously matched layer*.

6.1. Reflection is linear in μ by scaling

In this section the linearity in μ of leading order reflections by the layer with R^μ from (6.2) is demonstrated by a scaling argument when $\sigma(x_1) = \mathbf{1}_{x_1 > 0}$. In the

next section the reflection is computed exactly for Maxwell’s equations yielding additional information.

If $(L_1 + C)U = 0$, then

$$\underline{U}(t, x) := U(\mu t, \mu x), \quad \text{satisfies } R^\mu \underline{U} = 0.$$

Suppose that U has an incoming wave of wavelength ε and reflected waves U_ℓ with amplitudes $\rho_\ell \varepsilon$. Then \underline{U} has an incoming wave with wavelength $\underline{\varepsilon} := \varepsilon/\mu$. The reflected waves have amplitudes

$$\rho_\ell \varepsilon = \rho_\ell \mu \frac{\varepsilon}{\mu} = \rho_\ell \mu \underline{\varepsilon}.$$

Denote by $\underline{\rho}_\ell$ the reflection coefficient of R^μ . The leading amplitude of the reflected ℓ wave is then $\underline{\rho}_\ell \underline{\varepsilon}$. The preceding identity shows that $\underline{\rho}_\ell = \rho_\ell \mu$ showing that the reflection coefficients are linear in μ .^b

6.2. Reflection for Maxwell with smart layers

In this section \mathcal{L} may denote one of two distinct operators. One option is the Maxwell operator L_1 from (3.2) for the \mathbb{C}^3 -valued field $E + iB$. The lower order term is $\mathcal{B} := \mu C$ from the smart layer (6.1). Alternatively \mathcal{L} may denote the Béranger operator \tilde{L} with lower order term $\mathcal{B} = \mu C$ with

$$C = \begin{pmatrix} \sigma(x_1)I_N & 0 & \cdots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \cdots & 0 & 0 \end{pmatrix}, \quad \text{supp } \sigma \subset [0, \infty[.$$

In both cases the absorption term is linear in μ . We compute the dependence of the reflection coefficient on μ .

Lemma 4.6 shows that the Cauchy problem is equivalent to homogeneous problems in each half-space with a transmission condition on $\Gamma := \{x_1 = 0\}$,

$$[\mathcal{A}_1 \mathcal{U}]_\Gamma = 0. \tag{6.3}$$

In order to cover both cases the operator, coefficients, and unknown are indicated with round letters.

^bThis argument can be made rigorous under the following conditions. The incoming wave is a wave packet with oscillatory part $e^{i(\tau t + x\xi)/\varepsilon}$ with $(\tau, \xi) \in \text{Char } L$. Denote (τ, ξ') the part determining the oscillations in $x_1 = 0$. Consider the roots ξ_1 of each of the equations, $\det L_1(\tau, \xi_1, \xi') = 0$. The nonreal roots are called *elliptic*. They lead to waves which have the structure of a boundary layer of thickness $\sim \varepsilon$. The real roots are called *hyperbolic*. The favorable situation is when all the hyperbolic roots are at smooth points of the characteristic variety and the group velocities are transverse to the boundary. In that case one can construct infinitely accurate asymptotic solutions of the transmission problem consisting of incoming, reflected, and transmitted wave packets corresponding to the hyperbolic roots, and, a finite number of boundary layers corresponding to elliptic roots. As this is a long story, we content ourselves with the Maxwell computation of the next subsection.

We study the reflection of high frequency waves in $x_1 \leq 0$ which approach the boundary $x_1 = 0$. The input is an incident wave with phase $S^I(t, x) := \tau t + \xi \cdot x$, where $\tau \neq 0$ and $\tau(\xi) = \pm|\xi|$. The phase is chosen so that the group velocity $\mathbf{v} = -\xi/\tau$ satisfies $v_1 > 0$. Denote $x' := (x_2, x_3)$, $\xi' = (\xi_2, \xi_3)$. Theorem 5.1 applies to the incident wave with $\mathcal{B} \equiv 0$. In $x_1 \leq 0$, the incident wave is

$$\mathcal{U}^\varepsilon := e^{iS^I(t,x)/\varepsilon} a^I(t, x, \varepsilon), \quad a^I(t, x, \varepsilon) \sim \sum_{j=0}^\infty \varepsilon^j a_j^I(t, x), \quad \mathcal{L}_1(\partial_t, \partial_x)\mathcal{U}^\varepsilon = O(\varepsilon^\infty). \tag{6.4}$$

Suppose that the amplitudes a_j^I are supported in a tube, \mathcal{T} , of rays with compact temporal crosssections $\mathcal{T} \cap \{t = 0\} \subset\subset \{x_1 < 0\}$.

We construct a transmitted wave with the same phase, and a reflected wave with phase $S^R(t, x) := \tau t + \xi^R x$, with $\xi^R := (-\xi_1, \xi')$. We first show that there are uniquely determined reflected and transmitted waves. Then we compute exactly the leading terms in their asymptotic expansions.

The reflected wave \mathcal{V}^ε is also supported in $x_1 \leq 0$. The group velocity for the reflected wave is equal to $\mathbf{v}^R := (-v_1, v')$, and in $x_1 \leq 0$,

$$\mathcal{V}^\varepsilon = e^{iS^R(t,x)/\varepsilon} a^R(t, x, \varepsilon), \quad a^R(t, x, \varepsilon) \sim \sum_{j=0}^\infty \varepsilon^j a_j^R(t, x), \quad \mathcal{L}_1(\partial_t, \partial_x)\mathcal{V}^\varepsilon = O(\varepsilon^\infty). \tag{6.5}$$

The transmitted wave is supported in $x_1 \geq 0$,

$$\mathcal{W}^\varepsilon := e^{iS^I(t,x)/\varepsilon} a^T(t, x, \varepsilon), \quad a^T(t, x, \varepsilon) \sim \sum_{j=0}^\infty \varepsilon^j a_j^T(t, x), \quad \mathcal{L}(\partial_t, \partial_x)\mathcal{W}^\varepsilon = O(\varepsilon^\infty). \tag{6.6}$$

Theorem 6.1. (i) *Given the incoming amplitudes a_j^I there are uniquely determined amplitudes a_j^T and a_j^R so that for any choice of the $a_j^{I,R,T}(t, x, \varepsilon) \sim \sum \varepsilon^j a_j^{I,R,T}(t, x)$, the $\mathcal{U}^\varepsilon, \mathcal{V}^\varepsilon$ and \mathcal{W}^ε are infinitely accurate solutions of the differential equations and the transmission condition is also satisfied to infinite order,*

$$\forall(t, x') \in \mathbb{R} \times \mathbb{R}^2, \quad \mathcal{A}_1(\mathcal{U}^\varepsilon + \mathcal{V}^\varepsilon)(t, 0_-, x') = \mathcal{A}_1(\mathcal{W}^\varepsilon)(t, 0_+, x') + O(\varepsilon^\infty). \tag{6.7}$$

(ii) *In the case of Bérenger’s PML, the coefficients \tilde{a}_j^R vanish identically for $j \geq 0$.*

(iii) *For the smart layer (6.1), (6.2) with $\sigma = \mathbf{1}_{x_1 > 0}$, the coefficient a_0^R vanishes identically. The reflection coefficient of the layer is equal to*

$$R(\tau, \xi) = i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{8\tau\xi_1^2} = i \frac{\mu(1 + \nu)}{8\tau} \frac{v_1^2 - 1}{v_1^2}.$$

That is, if $a_0^I(t, 0_-, x') = \alpha(t, x')\Phi(\tau, \xi) \in \text{Ker } L(\tau, \xi)$, then

$$a^R(t, 0_-, x') = \varepsilon R(\tau, \xi)\alpha(t, x')\Phi(\tau, \xi^R) + O(\varepsilon^2).$$

Furthermore, the amplitudes $a^{T,R}$ are such that on the interface Γ , we have for all $(i, j) \in \mathbb{N}^2, i = 0$ and $j \leq 1$, or $i \geq 1$ and $j \geq 0$,^c

$$\partial_1^j a_i^T - \partial_1^j a_i^I \in \mu(\mathbb{C}_{i+j-1}[\mu] \otimes \mathbb{C}^3), \quad \partial_1^j a_i^R \in \mu(\mathbb{C}_{i+j-1}[\mu] \otimes \mathbb{C}^3). \tag{6.8}$$

(iv) The smart layer with $\sigma(x_1)$ satisfying $\sigma(0) = \dots = \sigma^{(k-1)}(0) = 0, \sigma^{(k)}(0) \neq 0$ is nonreflecting at order k for any angle of incidence, i.e. if $a_0^I(t, 0_-, x') = \alpha(t, x')\Phi(\tau, \xi)$, there exists $R_k(\tau, \xi)$ such that

$$a^R(t, 0_-, x') = \varepsilon^k \sigma^{(k)}(0) R_k(\tau, \xi) \alpha(t, x') \Phi(\tau, \xi^R) + O(\varepsilon^{k+1}).$$

Furthermore the amplitudes $a^{T,R}$ are linear functions of μ on the interface Γ . That is denoting $c_i^R(\mu) = a_i^R|_\Gamma$ and $c_i^T(\mu) = a_i^T|_\Gamma - a_i^I|_\Gamma$, we have for all $i \geq k$ in \mathbb{N} ,

$$c_i^{R,T}(\mu) \in \mu(\mathbb{C}_{i-1}[\mu] \otimes \mathbb{C}^3).$$

Remark 6.1. (1) There exist choices of $a^{I,R,T}$ so that $\mathcal{U}^\varepsilon, \mathcal{V}^\varepsilon$, and \mathcal{W}^ε is an exact solution. Since the transmission problem is well posed, there is a uniquely determined corrector c^ε smooth and infinitely small on both sides so that adding c^ε yields an exact solution. Adding c^ε to the left corresponds to adding the infinitely small term $c^\varepsilon e^{iS^I/\varepsilon}$ to a^I with a similar remark on the right.

(2) Part (iv) of the theorem with $k = 0$ generalizes part (iii) to discontinuous and variable $\sigma(x_1)$.

(3) The basis elements, Φ^R for a_1^R and Φ^I for a_0^I are homogeneous of degree 2 in τ, ξ . Doubling τ, ξ and also ε leaves the incoming and reflected waves unchanged. Therefore $\varepsilon R(\tau, \xi)$ must be equal to $2\varepsilon R(2\tau, 2\xi)$. This explains why R is homogeneous of degree -1 .

(4) The reflection coefficient vanishes when $\xi' = 0$. Since it is an even function of $\xi, \nabla_\xi R = 0$ too.

Proof. The incoming solution is given.

(i) Seek the leading amplitudes a_0^T and a_0^R . We will show that $a_0^R = 0$ so it is actually a_1^R that is the leading amplitude of the reflected wave. A jump discontinuity in a lower order coefficient does not lead to reflection at leading order. Denote

$$\mathcal{L}_T := \partial_t + \mathcal{A}_2 \partial_2 + \mathcal{A}_3 \partial_3; \quad \mathcal{L}_1 := \mathcal{L}_T + \mathcal{A}_1 \partial_1; \quad \mathcal{L} := \mathcal{L}_1 + \mu C.$$

$\mathcal{T} := \partial_t + v_2 \partial_2 + v_3 \partial_3$ is the tangential transport operator. By Theorem 5.1, the amplitudes are polarized, i.e. $a_0^{I,R,T} = \Pi_{\mathcal{L}} a_0^{I,R,T}$, and a_0^T (respectively a_0^R) satisfies a forward transport equation in $x_1 \geq 0$ (respectively backward in $x_1 \leq 0$) with zero

^c $\mathbb{C}_j[\mu]$ denotes the space of polynomials of degree less than or equal to j with complex coefficients. $\mathbb{C}_j[\mu] \otimes \mathbb{C}^3$ is the corresponding space of polynomials with coefficients in \mathbb{C}^3 .

initial values in time,

$$\begin{aligned} (v_1 \partial_1 + \mathcal{T})a_0^I &= 0, & x_1 &\in \mathbb{R}, \\ (-v_1 \partial_1 + \mathcal{T})a_0^R &= 0, & x_1 &\in \mathbb{R}_-, \\ (v_1 \partial_1 + \mathcal{T} + \mu \Pi_{\mathcal{L}} C \Pi_{\mathcal{L}})a_0^T &= 0, & x_1 &\in \mathbb{R}_+. \end{aligned} \tag{6.9}$$

Therefore, to determine a_0^T and a_0^R everywhere, it suffices to know $a_0^T(t, 0_+, x')$ and $a_0^R(t, 0_-, x')$. These values are determined from the transmission condition (6.3):

$$\mathcal{A}_1(a_0^I(t, 0_-, x') + a_0^R(t, 0_-, x')) = \mathcal{A}_1 a_0^T(t, 0_+, x'). \tag{6.10}$$

The matrix \mathcal{A}_1 is singular. It is easy to see that $(\text{Ker } \mathcal{L}(\tau, \xi) \oplus \text{Ker } \mathcal{L}(\tau, \xi^R)) \cap \text{Ker } \mathcal{A}_1 = 0$. Therefore, $\mathcal{A}_1 \text{Ker } \mathcal{L}(\tau, \xi)$ and $\mathcal{A}_1 \text{Ker } \mathcal{L}(\tau, \xi^R)$ are complementary subspaces and generate $\text{Range } \mathcal{A}_1$. This proves that

$$a_0^R(t, 0_-, x') = 0, \quad a_0^T(t, 0_+, x') - a_0^I(t, 0_-, x') = 0. \tag{6.11}$$

By the transport equation, we conclude that

$$a_0^R \equiv 0, \quad \text{for } x_1 < 0. \tag{6.12}$$

The reflected zeroth-order term vanishes identically when $x \in \mathbb{R}_-^3$. We also deduce from the transport equation (6.9) that

$$v_1(\partial_1 a_0^T - \partial_1 a_0^I) + \mu \Pi_{\mathcal{L}} C \Pi_{\mathcal{L}} a_0^T = 0 \quad \text{on } \Gamma. \tag{6.13}$$

Next determine inductively the correctors. For simplicity, throughout the proof we note $\Pi_{\mathcal{L}} := \Pi_{\mathcal{L}}(\tau, \xi)$ and $\Pi_{\mathcal{L}}^R := \Pi_{\mathcal{L}}(\tau, \xi^R)$. Write the recursion relation (5.5) for $j \geq 1$ for the incident, reflected and transmitted waves. Split the amplitudes as

$$\begin{aligned} a_j^{I,T}(t, x) &= \Pi_{\mathcal{L}} a_j^{I,T}(t, x) + (I - \Pi_{\mathcal{L}}) a_j^{I,T}(t, x), \\ a_j^R(t, x) &= \Pi_{\mathcal{L}}^R a_j^R(t, x) + (I - \Pi_{\mathcal{L}}^R) a_j^R(t, x). \end{aligned}$$

$(I - \Pi_{\mathcal{L}}) a_j^T(t, x)$ and $(I - \Pi_{\mathcal{L}}^R) a_j^R(t, x)$ are determined directly by (5.8). To determine the projection on the kernel, split the transmission condition (6.3) and insert (5.8) on the interface to get,

$$\begin{aligned} &\mathcal{A}_1(\Pi_{\mathcal{L}} a_j^I(t, 0, x') - \Pi_{\mathcal{L}} a_j^T(t, 0_+, x') + \Pi_{\mathcal{L}}^R a_j^R(t, 0_-, x')) \\ &= -\mathcal{A}_1((I - \Pi_{\mathcal{L}}) a_j^I(t, 0, x') - (I - \Pi_{\mathcal{L}}) a_j^T(t, 0_+, x') + (I - \Pi_{\mathcal{L}}^R) a_j^R(t, 0_-, x')). \end{aligned} \tag{6.14}$$

As for the terms of order 0, this determines $\Pi_{\mathcal{L}} a_j^T(t, 0_+, x')$ and $\Pi_{\mathcal{L}}^R a_j^R(t, 0_-, x')$. By (5.9), the projections are solution of a transport equation, therefore uniquely determined by initial data and the values on the boundary. Borel's theorem allows one to construct

$$a^I(t, x, \varepsilon), \quad a^T(t, x, \varepsilon) \quad \text{and} \quad a^R(t, x, \varepsilon),$$

so that the transmission condition is exactly satisfied. With this choice, the approximate solution satisfies the transmission problem with infinitely small residual.

(ii) Theorem 4.4 implies that the exact solution in $x_1 \leq 0$ is equal to $\mathcal{U}^\varepsilon + O(\varepsilon^\infty)$. The error of the approximation is $O(\varepsilon^\infty)$ so the exact solution is equal to $\mathcal{U}^\varepsilon + \mathcal{V}^\varepsilon + O(\varepsilon^\infty)$. Therefore $\mathcal{V}^\varepsilon = (\mathcal{U}^\varepsilon + \mathcal{V}^\varepsilon) - \mathcal{U}^\varepsilon = O(\varepsilon^\infty)$ which is the desired conclusion.

(iii) For the smart layer (6.1), (6.2) with $\sigma = \mathbf{1}_{x_1 > 0}$, compute the first-order term by (5.11) with $j = 1$. First deduce from (5.8) that

$$\begin{aligned} (I - \Pi_L) a_1^I(t, x) &= iQ_L L_1(\partial_t, \partial_x) a_0^I(t, x), & x_1 \in \mathbb{R}, \\ (I - \Pi_L^R) a_1^R(t, x) &= 0, & x_1 \in \mathbb{R}_-, \\ (I - \Pi_L) a_1^T(t, x) &= iQ_L(L_1(\partial_t, \partial_x) + \mu C) a_0^T(t, x), & x_1 \in \mathbb{R}_+. \end{aligned} \tag{6.15}$$

Replace in (6.15) the x_1 derivatives using (6.9),

$$\begin{aligned} (I - \Pi_L) a_1^I(t, x) &= iQ_L(L_T + A_1 \partial_1) a_0^I(t, x) = iQ_L\left(L_T - \frac{1}{v_1} A_1 \mathcal{T}\right) a_0^I(t, x), \\ (I - \Pi_L) a_1^T(t, x) &= iQ_L(L_T + A_1 \partial_1 + \mu C) a_0^T(t, x) \\ &= iQ_L\left(L_T - \frac{1}{v_1} A_1 \mathcal{T} + \mu\left(-\frac{1}{v_1} A_1 \Pi_L C \Pi_L + C\right)\right) a_0^T(t, x), \end{aligned}$$

to obtain, with $C_1 := C - \frac{1}{v_1} A_1 \Pi_L C \Pi_L$,

$$(I - \Pi_L) a_1^T(t, 0_+, x') - (I - \Pi_L) a_1^I(t, 0, x') = i\mu Q_L C_1 a_0^I(t, 0, x'). \tag{6.16}$$

Using $(I - \Pi_L) a_1^R = 0$ in the transmission condition yields,

$$A_1(\Pi_L^R a_1^R + \Pi_L a_1^I - \Pi_L a_1^T) = i\mu A_1 Q_L C_1 a_0^I. \tag{6.17}$$

The eigenvalues of A_1 are 0 and ± 1 , with associated orthonormal set of eigenvectors $\Phi_0 = e_1$ and $\Phi_\pm = (0, 1, \pm i)/\sqrt{2}$. The projection operators on the positive and negative eigenspaces are $\pi_\pm(A_1) = \Phi_\pm \Phi_\pm^*$, and $C = \Phi_+ \Phi_+^* + \nu \Phi_- \Phi_-^*$. The kernel of $L(\tau, \xi)$ is one-dimensional, it is spanned by

$$\Phi(\tau, \xi) = \xi - \frac{\tau^2}{\xi_1} e_1 + i \frac{\tau}{\xi_1} \xi \wedge e_1 = \left(\xi_1 - \frac{\tau^2}{\xi_1}, i \frac{\tau}{\xi_1} \xi_3 + \xi_2, -i \frac{\tau}{\xi_1} \xi_2 + \xi_3 \right), \tag{6.18}$$

and the projection on $\text{Ker } L(\tau, \xi)$ is $\Pi_L = \frac{\Phi \Phi^*}{\Phi^* \Phi}$. Compute

$$\begin{aligned} \Pi_L C \Pi_L &= \frac{\Phi \Phi^*}{\Phi^* \Phi} (\Phi_+ \Phi_+^* + \nu \Phi_- \Phi_-^*) \frac{\Phi \Phi^*}{\Phi^* \Phi} \\ &= \frac{1}{(\Phi^* \Phi)^2} \Phi(\Phi^* \Phi_+)(\Phi_+^* \Phi) \Phi^* + \nu \Phi(\Phi^* \Phi_+)(\Phi_+^* \Phi) \Phi^* \\ &= \frac{|\Phi^* \Phi_+|^2 + \nu |\Phi^* \Phi_-|^2}{\Phi^* \Phi} \Pi_L. \end{aligned}$$

Define

$$\gamma := \frac{|\Phi^* \Phi_+|^2 + \nu |\Phi^* \Phi_-|^2}{\Phi^* \Phi}, \quad \text{so,} \quad \Pi_L C \Pi_L = \gamma \Pi_L. \tag{6.19}$$

Since a_0^I is polarized,

$$C_1 a_0^I = \left(C - \frac{\gamma}{v_1} A_1 \Pi_L \right) a_0^I = \left(C - \frac{\gamma}{v_1} A_1 \right) a_0^I := \tilde{C}_1 a_0^I.$$

To compute the right-hand side of (6.17), use

$$C = \begin{pmatrix} 0 & 0 & 0 \\ 0 & \frac{\nu + 1}{2} & i \frac{\nu - 1}{2} \\ 0 & -i \frac{\nu - 1}{2} & \frac{\nu + 1}{2} \end{pmatrix}, \quad \gamma = \frac{(\tau - \xi_1)^2 + \nu(\tau + \xi_1)^2}{4\tau^2}$$

$$= \frac{1}{4}((1 + v_1)^2 + \nu(1 - v_1)^2),$$

to find

$$\tilde{C}_1 = \frac{\nu + 1}{2} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & i \frac{v_1^2 + 1}{2v_1} \\ 0 & -i \frac{v_1^2 + 1}{2v_1} & 1 \end{pmatrix}.$$

Write $a_0^I = \alpha_0^I \Phi$ and compute

$$\tilde{C}_1 \Phi = (1 + \nu) \frac{\xi_1^2 - \tau^2}{4\tau \xi_1^2} \Psi, \quad \Psi = \begin{pmatrix} 0 \\ \xi_2 \tau - i \xi_1 \xi_3 \\ \xi_3 \tau + i \xi_2 \xi_3 \end{pmatrix},$$

to find a new version of (6.17),

$$A_1 (\Pi_L^R a_1^R + \Pi_L a_1^I - \Pi_L a_1^T) = i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{4\tau \xi_1^2} \alpha_0^I A_1 Q_L \Psi. \tag{6.20}$$

Next compute $Q_L \Psi$. First compute a basis of eigenvectors for $L(\tau, \xi)$. Φ_2 is such that $L(\tau, \xi)\Phi_2 = \tau\Phi_2$, Φ_3 is such that $L(\tau, \xi)\Phi_3 = 2\tau\Phi_3$. Choose

$$\Phi_2 = \xi, \quad \Phi_3 = \Phi(-\tau, \xi).$$

Note that

$$\Psi = \tau\xi - \xi_1(\tau e_1 + i\xi \wedge e_1)$$

and

$$\Phi = \xi - \frac{\tau}{\xi_1}(\tau e_1 - i\xi \wedge e_1), \quad \Phi_3 = \xi - \frac{\tau}{\xi_1}(\tau e_1 + i\xi \wedge e_1),$$

which gives

$$(\tau e_1 + i\xi \wedge e_1) = \frac{\xi_1}{\tau}(\xi - \Phi_3)$$

and

$$\Psi = \tau\xi - \frac{\xi_1^2}{\tau}(\xi - \Phi_2) = \frac{\tau^2 - \xi_1^2}{\tau}\xi + \frac{\xi_1^2}{\tau}\Phi_3.$$

Since Q_L is the left inverse of L , we have $Q_L\xi = \frac{1}{\tau}\xi$, and $Q_L\Phi_3 = \frac{1}{2\tau}\Phi_3$, which gives

$$Q_L\Psi = \frac{\tau^2 - \xi_1^2}{\tau^2}\xi + \frac{\xi_1^2}{2\tau^2}\Phi_3.$$

Write the coefficients on Γ as

$$\Pi_L a_1^{I,T} = \alpha_1^{I,T}\Phi, \quad \Pi_L^R a_1^R = \alpha_1^R\Phi^R$$

and inject into the transmission condition to obtain

$$\alpha_1^R A_1\Phi + (\alpha_1^I - \alpha_1^T)A_1\Phi^R = i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{4\tau\xi_1^2} \alpha_0^I A_1 \left(\frac{\tau^2 - \xi_1^2}{\tau^2}\xi + \frac{\xi_1^2}{2\tau^2}\Phi_3 \right).$$

Since the kernel of A_1 is e_1 , A_1 in the preceding identity may be replaced by the projection on (e_2, e_3) . The projection of Φ is $\varphi = \xi' - i\frac{\tau}{\xi_1}\xi \wedge e_1$, and note that Φ_3 and $\Phi(\tau, \xi^R)$ have the same projection, which is $\varphi_3 = \xi' + i\frac{\tau}{\xi_1}\xi \wedge e_1$. Write

$$\begin{aligned} \alpha_1^R\varphi_3 + (\alpha_1^I - \alpha_1^T)\varphi &= i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{4\tau\xi_1^2} \alpha_0^I \left(\frac{\tau^2 - \xi_1^2}{\tau^2}\xi' + \frac{\xi_1^2}{2\tau^2}\varphi_3 \right) \\ &= i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{4\tau\xi_1^2} \alpha_0^I \left(\frac{\tau^2 - \xi_1^2}{2\tau^2}(\varphi + \varphi_3) + \frac{\xi_1^2}{2\tau^2}\varphi_3 \right) \\ &= i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{4\tau\xi_1^2} \alpha_0^I \left(\frac{\tau^2 - \xi_1^2}{2\tau^2}\varphi + \frac{1}{2}\varphi_3 \right). \end{aligned}$$

The solutions are parametrized by α_0^I ,

$$\alpha_1^R = i\mu(1 + \nu) \frac{\xi_1^2 - \tau^2}{8\tau\xi_1^2} \alpha_0^I, \quad \alpha_1^I - \alpha_1^T = -i\mu(1 + \nu) \frac{(\xi_1^2 - \tau^2)^2}{8\tau^3\xi_1^2} \alpha_0^I.$$

Now use the results in (i) and prove (6.8) by induction on i . Equation (6.12) asserts that $a_0^R = 0$ for $x_1 < 0$. Equations (6.9) and (6.10), imply that at $x_1 = 0$, $a_0^{I,T} = \Pi_L a_0^{I,T}$, and

$$v_1(\partial_1 a_0^T - \partial_1 a_0^I) + \mathcal{T}(a_0^T - a_0^I) + \mu\gamma(a_0^T - a_0^I) + \mu\gamma a_0^I = 0, \quad x_1 \geq 0.$$

Differentiation in x_1 several times yields $\partial_1^j a_0^T - \partial_1^j a_0^I \in \mu\mathcal{C}_{j-1}[\mu] \otimes \mathbb{C}^3$, giving the results for $i = 0$ and $j \geq 1$.

Assuming the inductive hypothesis is true for i we prove it for $i + 1$. Write (5.11) in the form

$$\begin{aligned} (I - \Pi_L^R)a_{i+1}^R(t, x) &= iQ_L^R(L_T(\partial_t, \partial_{x'}) + A_1\partial_1)a_i^I(t, x), \\ (I - \Pi_L)a_{i+1}^I(t, x) &= iQ_L(L_T(\partial_t, \partial_{x'}) + A_1\partial_1)a_i^I(t, x), \\ (I - \Pi_L)a_{i+1}^T(t, x) &= iQ_L(L_T(\partial_t, \partial_{x'}) + A_1\partial_1 + \mu C)a_i^T(t, x). \end{aligned}$$

By induction, $\partial_1^j(I - \Pi_L^R)a_i^R(t, x) \in \mu \mathbb{C}_{i+j-1}[\sigma] \otimes \mathbb{C}^3$ on the interface. Write for $x_1 \geq 0$,

$$\begin{aligned} (I - \Pi_L)(a_{i+1}^T - a_{i+1}^I)(t, x) \\ = iQ_L(L_T + A_1\partial_1 + \mu C)(a_i^T - a_i^I)(t, x) + i\mu Q_L C a_i^I(t, x). \end{aligned} \quad (6.21)$$

The inductive hypothesis shows that on Γ , $\partial_1^j(L_T + A_1\partial_1 + \mu C)(a_i^T - a_i^I) \in \mu \mathbb{C}_{i+j}[\mu] \otimes \mathbb{C}^3$ and $\partial_1^j(\mu C a_i^I) \in \mu \mathbb{C}_0[\sigma] \otimes \mathbb{C}^3$. The result follows for $(I - \Pi_L)(a_{i+1}^T - a_{i+1}^I)$, and for $(I - \Pi_L^R)a_{i+1}^R$ in the same way. The transmission condition extends the assertion to the other parts $\Pi_L(a_{i+1}^T - a_{i+1}^I)$ and $\Pi_L^R a_{i+1}^R$.

(iv) Here σ vanishes to order k at $x_1 = 0$. (6.11) and (6.12) are still valid, and the transport equation (6.9) implies on the interface Γ that

$$\begin{aligned} \partial_1^j \Pi_L a_0^T - \partial_1^j \Pi_L a_0^I &= 0, \quad j = 0, \dots, k, \\ \partial_1^{k+1} \Pi_L a_0^T - \partial_1^{k+1} \Pi_L a_0^I &= -\mu \sigma^{(k)}(0) \frac{\gamma}{v_1} \Pi_L a_0^I. \end{aligned} \quad (6.22)$$

From (6.21) for $i = 0$, (6.15) and (6.22), derive

$$\begin{aligned} a_1^R &\equiv \Pi_L a_1^R \quad \text{everywhere,} \\ \partial_1^j (I - \Pi_L) a_1^T - \partial_1^j (I - \Pi_L) a_1^I &= 0, \quad j = 0, \dots, k - 1, \quad \text{on } \Gamma, \\ \partial_1^k (I - \Pi_L) a_1^T - \partial_1^k (I - \Pi_L) a_1^I &= i\mu \sigma^{(k)}(0) Q_L C_1 \Pi_L a_0^I \quad \text{on } \Gamma. \end{aligned} \quad (6.23)$$

Using the transmission conditions to obtain on the interface Γ ,

$$\Pi_L^R a_1^R = 0, \quad \Pi_L a_1^I = \Pi_L a_1^T.$$

Insert into the transport equations (5.9) to find $a_1^R = 0$ in \mathbb{R}_- , and

$$\begin{aligned} \partial_1^j \Pi_L a_1^T - \partial_1^j \Pi_L a_1^I &= 0, \quad j = 0, \dots, k - 1, \quad \text{on } \Gamma, \\ \partial_1^k \Pi_L a_1^T - \partial_1^k \Pi_L a_1^I &= -\frac{1}{v_1} \Pi_L A_1 (\partial_1^k (I - \Pi_L) a_1^T - \partial_1^k (I - \Pi_L) a_1^I) \\ &= -i \frac{\mu}{v_1} \sigma^{(k)}(0) \Pi_L A_1 Q_L C_1 \Pi_L a_0^I \quad \text{on } \Gamma. \end{aligned} \quad (6.24)$$

Recover

$$\begin{aligned} \partial_1^j a_1^T - \partial_1^j a_1^I &= 0, \quad j = 0, \dots, k - 1, \quad \text{on } \Gamma, \\ \partial_1^k a_1^T - \partial_1^k a_1^I &= i\mu \sigma^{(k)}(0) \left(I - \frac{1}{v_1} \Pi_L A_1 \right) \Pi_L A_1 Q_L C_1 \Pi_L a_0^I \quad \text{on } \Gamma. \end{aligned} \quad (6.25)$$

Now proceed iteratively, to see that, for $i < k + 1$,

$$\begin{aligned}
 a_i^R &\equiv \Pi_L a_i^R \equiv 0 \quad \text{everywhere,} \\
 \partial_1^j a_i^T - \partial_1^j a_i^I &= 0, \quad j = 0, \dots, k - i, \quad \text{on } \Gamma, \\
 \partial_1^{k-i+1} (I - \Pi_L) a_i^T - \partial_1^{k-i+1} (I - \Pi_L) a_i^I & \\
 &= i Q_L A_1 (\partial_1^{k-i+2} a_{i-1}^T - \partial_1^{k-i+2} a_{i-1}^I) \quad \text{on } \Gamma, \\
 \partial_1^{k-i+1} \Pi_L a_i^T - \partial_1^{k-i+1} \Pi_L a_i^I & \\
 &= -\frac{1}{v_1} \Pi_L A_1 (\partial_1^{k-i+1} (I - \Pi_L) a_i^T - \partial_1^{k-i+1} (I - \Pi_L) a_i^I).
 \end{aligned} \tag{6.26}$$

Denote by s_i the value of $\partial_1^{k-i+1} a_i^T - \partial_1^{k-i+1} a_i^I$ on Γ . Equation (6.26) yields the recursion relation

$$s_j = i \left(I - \frac{1}{v_1} \Pi_L A_1 \right) Q_L A_1 s_{j-1}, \quad s_1 = i \mu \sigma^{(k)}(0) \left(I - \frac{1}{v_1} \Pi_L A_1 \right) Q_L C_1 \Pi_L a_0^I,$$

which can be solved as

$$s_{k+1} = i \mu \sigma^{(k)}(0) M \Pi_L a_0^I, \quad \text{with } M := \left(i \left(I - \frac{1}{v_1} \Pi_L A_1 \right) Q_L A_1 \right)^k Q_L C_1.$$

The first nonzero reflected term is therefore $a_{k+1}^R = \Pi_L^R a_{k+1}^R$, and using the transmission condition yields

$$A_1 (\Pi_L a_{k+1}^R + \Pi_L a_{k+1}^I - \Pi_L a_{k+1}^T + s_{k+1}) = 0.$$

The incoming amplitude on Γ is $a_0^I = \alpha(t, x') \Phi$, the leading reflection is $a_{k+1}^R = \alpha_{k+1}^R \Phi^R$ and the leading transmission is $a_{k+1}^T = \alpha_{k+1}^T \Phi^T$. Using again the notation Φ' to denote the projection of a vector Φ on $\text{Vec}(e_2, e_3)$, this linear system is solved as

$$\begin{aligned}
 \alpha_{k+1}^R &= -\frac{\Phi' \wedge s'_{k+1}}{\Phi' \wedge (\Phi^R)'}, \quad \alpha_{k+1}^T - \alpha_{k+1}^I = \frac{(\Phi^R)' \wedge s'_{k+1}}{\Phi' \wedge (\Phi^R)'}. \\
 \alpha_{k+1}^R &= -i \mu \sigma^{(k)}(0) \alpha \frac{\Phi' \wedge (M \Phi)'}{\Phi' \wedge (\Phi^R)'}, \quad \alpha_{k+1}^T - \alpha_{k+1}^I = i \mu \sigma^{(k)}(0) \alpha \frac{(\Phi^R)' \wedge (M \Phi)'}{\Phi' \wedge (\Phi^R)'}.
 \end{aligned}$$

The proof of the linearity follows the same path as in (iii). □

6.3. *Harmoniously matched layers*

Based on Theorem 6.1 we construct an extrapolation method for symmetric hyperbolic operators with smart layers which eliminates the leading order reflection. The resulting method has desirable stability properties and is nearly as good as Bérenger’s algorithm for the Maxwell equations where his method is at its best. We think that the new method provides a good alternative in situations where Bérenger’s method is not so effective.

Consider the computational domain $x_1 \leq b_1$. The domain of interest is the interval $x_1 \leq a_1 < b_1$. The absorbing layer is located in $a_1 \leq x_1 \leq b_1$. The differential operator in the computational domain is symmetric hyperbolic L with smart layer

$$LU + \sigma_1(x_1)(\pi_+(A_1) + \nu\pi_-(A_1))U = 0, \quad \sigma_1 \geq 0, \quad \text{supp } \sigma_1 \subset \{x_1 \geq a_1\}.$$

At the outer boundary $x_1 = b_1$ of the absorbing layer impose the simplest weakly reflecting boundary condition

$$\pi_-(A_1)U = 0 \quad \text{when } x_1 = b_1.$$

This is a well-posed problem provided that A_1 has constant rank on $x_1 = b_1$. When $L = L_1(\partial)$ has constant coefficients, it generates a contraction group in $L^2(\{x_1 \leq b_1\})$.

The *harmoniously matched layer algorithms* compute a smart layer with coefficient σ_1 and also with coefficient $2\sigma_1$. In view of Theorem 6.1, subtracting the second from twice the first, $2U(\sigma_1) - U(2\sigma_1)$, yields a field with one more vanishing term in the reflected wave at the interface $x_1 = a_1$. This extrapolation removes the leading reflection.

The harmonious matched layer algorithms in a rectangular domain \mathcal{R} perform the same extrapolation with absorptions in all directions. With

$$LU + \sum_{j=1}^d \sigma_j(x_j)(\pi_+(A_j) + \nu\pi_-(A_j))U = 0, \quad \sigma_j \geq 0, \quad \text{supp } \sigma_j \subset \{|x_j| \geq a_j\}$$

with

$$\pi_{\mp}(A_j)U = 0 \quad \text{when } x_j = \pm b_j.$$

This initial boundary value problem on a rectangle has weak solutions.^d When $L = L_1(\partial)$, the $L^2(\mathcal{R})$ norm is nonincreasing in time. The extrapolation is $2U(\sigma_1, \dots, \sigma_d) - U(2\sigma_1, \dots, 2\sigma_d)$.

Open Problem. *For discontinuous σ_j , the uniqueness of solutions to the initial boundary value problem on the rectangular computational domain is not known because of the discontinuity of the boundary space $\ker A_j$ at the corner. Solutions are typically discontinuous. Uniqueness of strong solutions and existence of weak solutions is proved by the energy method. We do not know how to prove uniqueness of solutions with regularity not exceeding that of solutions known to exist. Similar problems plague virtually all methods on rectangular domains with absorbing boundary conditions imposed on the computational domain with corners. The present problem is one of the simplest of its kinds. The fact that algorithms designed to compute solutions encounter no difficulties is reason for optimism.*

^dThis can be proved by penalization. Denote by Ω the rectangular computational domain. Add $\Lambda \mathbf{1}_{\mathbb{R}^d \setminus \Omega}$ to L and solve on $\mathbb{R}_{t,x}^{1+d}$. The limit as $\Lambda \rightarrow \infty$ provides a solution in $L^\infty([0, T] : L^2(\mathcal{R}))$ [4].

6.4. Numerical experiments

Simulations are performed for the 2D transverse electric Maxwell system in the (x, y) coordinates,

$$\begin{aligned} \partial_t E_x - \partial_y H_z &= 0, \\ \partial_t E_y + \partial_x H_z &= 0, \\ \partial_t H_z + \partial_x E_y - \partial_y E_x &= 0, \end{aligned} \tag{6.27}$$

in a rectangle, with boundary conditions $n \wedge E = 0$ on the west, north and south boundaries. The layer will be imposed on the east boundary. Maxwell Béranger is given by

$$\begin{aligned} \partial_t E_x - \partial_y H_z &= 0, \\ \partial_t E_y + \partial_x H_z + \sigma(x) E_y &= 0, \\ \partial_t H_{zx} + \partial_x E_y + \sigma(x) H_{zx} &= 0, \\ \partial_t H_{zy} - \partial_y E_x &= 0, \\ H_z &= H_{zx} + H_{zy}. \end{aligned} \tag{6.28}$$

For the computation, these equations are used in the whole rectangle (see the discussion in the Introduction), with $\sigma = 0$ outside the layer. The boundary conditions are

$$E_y = H_z \quad \text{and} \quad E_x = 0 \quad \text{on the east,} \quad n \wedge E = 0 \quad \text{on the other boundaries.} \tag{6.29}$$

Since $\Pi_+(A_1) = \frac{E_y + H_z}{2}(0, 1, 1)$ and $\Pi_-(A_1) = \frac{E_y - H_z}{2}(0, 1, -1)$, the smart layers are:

$$\begin{aligned} \partial_t E_x - \partial_y H_z &= 0, \\ \partial_t E_y + \partial_x H_z + \frac{\sigma(x)}{2}(E_y + H_z + \nu(E_y - H_z)) &= 0, \\ \partial_t H_z - \partial_y E_x + \partial_x E_y + \frac{\sigma(x)}{2}(E_y + H_z - \nu(E_y - H_z)) &= 0. \end{aligned} \tag{6.30}$$

The boundary conditions (6.29) are imposed.

The Yee scheme for Maxwell is

$$\begin{aligned} \frac{(E_x)_{i+\frac{1}{2},j}^n - (E_x)_{i+\frac{1}{2},j}^{n-1}}{\Delta t} - \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} - (H_z)_{i+\frac{1}{2},j-\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta y} &= 0, \\ \frac{(E_y)_{i,j+\frac{1}{2}}^n - (E_y)_{i,j+\frac{1}{2}}^{n-1}}{\Delta t} + \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} - (H_z)_{i-\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta x} &= 0, \\ \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} - (H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta t} + \frac{(E_y)_{i+1,j+\frac{1}{2}}^n - (E_y)_{i,j+\frac{1}{2}}^n}{\Delta x} \\ - \frac{(E_x)_{i+\frac{1}{2},j+1}^n - (E_x)_{i+\frac{1}{2},j}^n}{\Delta y} &= 0. \end{aligned} \tag{6.31}$$

The Yee scheme for Maxwell Bérenger using the notations $\sigma_i = \sigma(x_i)$ and $\sigma_{i+\frac{1}{2}} = \sigma(x_{i+\frac{1}{2}})$ is,

$$\begin{aligned}
 & \frac{(E_x)_{i+\frac{1}{2},j}^n - (E_x)_{i+\frac{1}{2},j}^{n-1}}{\Delta t} - \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} - (H_z)_{i+\frac{1}{2},j-\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta y} = 0, \\
 & \frac{(E_y)_{i,j+\frac{1}{2}}^n - (E_y)_{i,j+\frac{1}{2}}^{n-1}}{\Delta t} + \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} - (H_z)_{i-\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta x} \\
 & \quad + \sigma_i \frac{(E_y)_{i,j+\frac{1}{2}}^n + (E_y)_{i,j+\frac{1}{2}}^{n-1}}{2} = 0, \\
 & \frac{(H_{zx})_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} - (H_{zx})_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta t} + \frac{(E_y)_{i+1,j+\frac{1}{2}}^n - (E_y)_{i,j+\frac{1}{2}}^n}{\Delta x} \quad (6.32) \\
 & \quad + \sigma_{i+\frac{1}{2}} \frac{(H_{zx})_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} + (H_{zx})_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{2} = 0, \\
 & \frac{(H_{zy})_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} - (H_{zy})_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta t} - \frac{(E_x)_{i+\frac{1}{2},j+1}^n - (E_x)_{i+\frac{1}{2},j}^n}{\Delta y} = 0, \\
 & \quad (H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} = (H_{zx})_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} + (H_{zy})_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}}.
 \end{aligned}$$

The Yee scheme for the smart layers is

$$\frac{(E_x)_{i+\frac{1}{2},j}^n - (E_x)_{i+\frac{1}{2},j}^{n-1}}{\Delta t} - \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} - (H_z)_{i+\frac{1}{2},j-\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta y} = 0, \quad (6.33a)$$

$$\begin{aligned}
 & \frac{(E_y)_{i,j+\frac{1}{2}}^n - (E_y)_{i,j+\frac{1}{2}}^{n-1}}{\Delta t} + \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} - (H_z)_{i-\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta x} \\
 & \quad + \frac{(1+\nu)\sigma_i}{2} \frac{(E_y)_{i,j+\frac{1}{2}}^n + (E_y)_{i,j+\frac{1}{2}}^{n-1}}{2} \\
 & \quad + \frac{(1-\nu)}{2} \frac{\sigma_{i+\frac{1}{2}}(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}} + \sigma_{i-\frac{1}{2}}(H_z)_{i-\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{2} = 0, \quad (6.33b)
 \end{aligned}$$

$$\begin{aligned}
 & \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} - (H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{\Delta t} + \frac{(E_y)_{i+1,j+\frac{1}{2}}^n - (E_y)_{i,j+\frac{1}{2}}^n}{\Delta x} \\
 & \quad - \frac{(E_x)_{i+\frac{1}{2},j+1}^n - (E_x)_{i+\frac{1}{2},j}^n}{\Delta y} + \frac{(1-\nu)}{2} \frac{\sigma_{i+1}(E_y)_{i+1,j+\frac{1}{2}}^n + \sigma_i(E_y)_{i,j+\frac{1}{2}}^n}{2} \\
 & \quad + \frac{(1+\nu)\sigma_{i+\frac{1}{2}}}{2} \frac{(H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n+\frac{1}{2}} + (H_z)_{i+\frac{1}{2},j+\frac{1}{2}}^{n-\frac{1}{2}}}{2} = 0. \quad (6.33c)
 \end{aligned}$$

The schemes are implemented using time windows to save memory.

The harmoniously matched layers can be implemented in several ways that we compare. The function $\sigma(x)$ is as above.

HML Version 1. Global extrapolation. Compute the solution of (6.33) with an absorption of σ , (E^1, H^1) and 2σ , (E^2, H^2) over the whole time window. Then $E_{x,y} = 2 * E_{x,y}^1 - E_{x,y}^2$ and $H_z = 2 * H_z^1 - H_z^2$.

HML Version 2. Local extrapolation. Compute *at each time step* the solution of (6.33) with an absorption of σ , (E^1, H^1) and 2σ , (E^2, H^2) over the whole time interval. Then $E_{x,y} = 2 * E_{x,y}^1 - E_{x,y}^2$ and $H_z = 2 * H_z^1 - H_z^2$. Save computation by taking advantage of the fact that the computation of E_x does not involve the absorption parameter. At each time step,

- (1) E_x is computed by (6.33a),
- (2) two values of E_y are computed by (6.33b): E_y^1 with an absorption parameter equal to σ , E_y^2 with an absorption parameter equal to 2σ .
- (3) two values of H_z are computed by (6.33c): H_z^1 with an absorption parameter equal to σ , H_z^2 with an absorption parameter equal to 2σ .

Then $E_y = 2 * E_y^1 - E_y^2$ and $H_z = 2 * H_z^1 - H_z^2$.

HML Version 3. Split field local extrapolation. At each time step,

- (1) E_x is computed by (6.33a),
- (2) two values of E_y are computed by (6.33b): E_y^1 with an absorption parameter equal to σ , E_y^2 with an absorption parameter equal to 2σ . Then $E_y = 2 * E_y^1 - E_y^2$.
- (3) two values of H_z are computed by (6.33c): H_z^1 with an absorption parameter equal to σ , H_z^2 with an absorption parameter equal to 2σ . Then $H_z = 2 * H_z^1 - H_z^2$.

We perform a series of experiments to illustrate the transmission properties of the layers. The coefficient ν is meant to achieve backward absorption and is taken equal to zero. The domain of interest is $(0, 6) \times (0, 10)$, the coefficient $\sigma(x)$ is supported in $6 \leq x \leq 10$. The time of computation is 4, the initial electric field is zero. The initial transverse magnetic field,

$$H_z^0 = \cos^2\left(\pi \frac{|\mathbf{x} - \mathbf{x}_c|}{r}\right) \cos\left(k\pi \mathbf{v} \cdot \frac{\mathbf{x} - \mathbf{x}_c}{r}\right) \chi_{|\mathbf{x} - \mathbf{x}_c| \leq r}$$

is compactly supported in the ball $B(\mathbf{x}_c, r)$, with $\mathbf{x}_c = (5, 5)$ and $r = 0.8$.

The time of computation is fixed such that there is no reflection on the exterior walls. The initial mesh is taken to be $\Delta x = \Delta y = 0.1$, $\Delta t = 0.0702$, and then divided by 2 twice.

In the first set of experiments, the absorption coefficient is constant in the layer, equal to 2. The initial magnetic field hits the layer at incidence $0^\circ(\mathbf{v} = (1, 0))$ or $45^\circ(\mathbf{v} = (1, 1))$.

Table 1. Comparison of the L^∞ errors for high frequency, discontinuous absorption.

Refinement	Normal incidence			45° incidence		
	0	1	2	0	1	2
Bérenger	9.4e-02	3.9e-02	7.9e-03	1.3e-01	2.9e-02	5.6e-03
Smart	5.2e-02	1.3e-02	5.1e-04	6.2e-02	1.1e-02	5.3e-03
HML V1	3.4e-02	3.1e-03	2.1e-05	4.5e-02	1.2e-03	5.5e-04
HML V2	2.5e-02	6.0e-03	1.2e-03	7.4e-02	1.1e-02	1.7e-03
HML V3	2.1e-02	4.2e-03	5.1e-04	4.5e-02	5.3e-03	5.7e-04

Table 2. Comparison of the L^∞ errors for low frequency, discontinuous absorption.

Refinement	Normal incidence			45° incidence		
	0	1	2	0	1	2
Bérenger	1.5e-02	7.1e-03	3.5e-03	1.3e-02	6.1e-03	3.0e-03
Smart	2.0e-02	2.0e-02	2.01e-02	4.3e-02	4.2e-02	4.2e-02
HML V1	1.7e-02	1.60e-02	1.6e-02	3.4e-02	3.3e-02	3.2e-02
HML V2	1.8e-02	1.1e-02	6.7e-03	3.1e-02	1.9e-02	1.1e-02
HML V3	4.3e-03	2.6e-03	1.4e-03	8.2e-03	4.8e-03	2.6e-03

Table 3. Comparison of the L^∞ errors for high frequency, continuous absorption.

Refinement	Normal incidence			45° incidence		
	0	1	2	0	1	2
Bérenger	3.8e-05	1.9e-07	2.1e-09	2.0e-04	9.1e-07	1.6e-09
Smart	2.7e-05	2.2e-07	1.7e-07	1.7e-04	9.0e-07	3.1e-08
HML V1	5.5e-07	6.0e-08	5.6e-08	5.6e-06	1.2e-08	4.7e-09
HML V2	6.8e-07	6.5e-08	3.1e-08	2.6e-06	8.1e-09	2.8e-09
HML V3	5.8e-08	2.4e-09	1.1e-09	1.5e-06	9.5e-10	9.0e-11

In Table 1 we compare the performances on a high frequency wave ($k = 10$), while in Table 2 we consider a low frequency wave ($k = 1$).

In Tables 3 and 4, we perform the same set of experiments, but the absorption coefficient is now a third degree polynomial in the layer, equal to $(x - 6)^3/8$.

The Bérenger layer performs well on every frequency and every angle of incidence. Among the three versions for the HML, the third version is the best, which should be analyzed thoroughly.

Next compare the method on a Gaussian initial value, supported in $(0, 6) \times (0, 10)$. Table 5 uses a constant absorption in the layer, while Table 6 uses the same smooth absorption as before.

Table 4. Comparison of the L^∞ errors for low frequency, continuous absorption.

Refinement	Normal incidence			45° incidence		
	0	1	2	0	1	2
Bérenger	6.2e-07	3.2e-08	7.8e-010	5.2e-07	2.9e-08	6.5e-010
Smart	5.3e-04	5.3e-04	5.2e-04	3.9e-04	3.8e-04	3.7e-04
HML V1	1.6e-04	1.6e-04	1.5e-04	8.6e-05	8.3e-05	8.2e-05
HML V2	4.1e-04	2.0e-04	9.6e-05	2.0e-04	9.8e-05	4.8e-05
HML V3	1.1e-05	5.4e-06	2.7e-06	5.9e-06	2.9e-06	1.4e-06

Table 5. Comparison of the L^∞ errors for a Gaussian initial magnetic field, constant absorption.

Refinement	0	1	2
Bérenger	1.5e-02	6.7e-03	3.3e-03
Smart	3.4e-02	3.4e-02	3.3e-02
HML V1	3.0e-02	2.9e-02	2.8e-02
HML V2	3.6e-02	2.5e-02	1.6e-02
HML V3	1.0e-02	6.6e-03	3.9e-03

Table 6. Comparison of the L^∞ errors for a Gaussian initial magnetic field, continuous absorption.

Refinement	0	1	2
Bérenger	7.5e-07	2.0e-08	8.3e-10
Smart	4.3e-04	4.2e-04	4.1e-04
HML V1	1.3e-04	1.2e-04	1.2e-04
HML V2	3.0e-04	1.5e-04	7.3e-05
HML V3	8.8e-06	4.3e-06	2.1e-06

Finally, take unstructured random initial value, supported in the ball centered at $(5, 5)$ and of radius 1. In Table 7, the absorption coefficient is constant in the layer, equal to 3.

In Table 8, the absorption coefficient is a function of x in the layer, equal to $(x - 6)^3/8$.

Summary. When comparing the reflection properties, the harmoniously matched layer, version 3, is competitive with the Bérenger layer. For very regular data, the Bérenger layers outperform everything. The performance of the HMLV3 gives hope the method with its stronger well-posedness, more robust absorption, and small

Table 7. Comparison of the L^∞ errors for a random initial magnetic field, constant absorption.

Refinement	0	1	2
Bérenger	5.7e-02	4.9e-02	4.4e-02
Smart	6.7e-02	6.3e-02	5.4e-02
HML V1	5.1e-02	4.5e-02	4.0e-02
HML V2	6.4e-02	3.0e-02	1.9e-02
HML V3	3.2e-02	1.5e-02	6.7e-03

Table 8. Comparison of the L^∞ errors for a random initial magnetic field, continuous absorption.

Refinement	0	1	2
Bérenger	1.1e-04	5.0e-05	4.4e-06
Smart	7.2e-04	6.9e-04	6.4e-04
HML V1	2.1e-04	2.2e-04	2.0e-04
HML V2	5.0e-04	2.7e-04	1.2e-04
HML V3	1.5e-05	7.9e-06	3.7e-06

reflection at all angles will be a good method where Bérenger has proven less good. For example, for nonconstant coefficients and nonlinear problems. We have taken pains to make the comparison where Bérenger is at its best. In 2D with a layer in a single direction the HML has an extra cost. Since there are five quantities to compute at each time step instead of four for Bérenger. This is no longer the case in three dimensions, since both strategies have to split six unknowns.

Open problems. (1) *Our analysis does not explain the much better behavior with continuous absorption, nor the advantages of HMLV3.* (2) *A comparison with other methods where only supplementary ordinary differential equations are added should be made.*

Acknowledgments

This research project has spanned many years. It owes a great deal to the support of the University Paris 13 where J.R. was often invited for one month visits. J.R. is partially supported by the National Science Foundation under grant NSF DMS 0405899.

References

1. S. Abarbanel and D. Gottlieb, A mathematical analysis of the PML method, *J. Comput. Phys.* **134** (1997) 357–363.
2. S. Abarbanel and D. Gottlieb, On the construction and analysis of absorbing layers in cem, *Appl. Numer. Math.* **27** (1998) 331–340.
3. D. Appelö, T. Hagström and G. Kreiss, Perfectly matched layers for hyperbolic systems: General formulation, well-posedness and stability, *SIAM J. Appl. Math.* **67** (2006) 1–23.
4. C. Bardos and J. Rauch, Maximal positive boundary value problems as limits of singular perturbation problems, *Trans. Amer. Math. Soc.* (1982) 377–408.
5. E. Bécache, S. Fauqueux and P. Joly, Stability of perfectly matched layers, group velocities and anisotropic waves, *J. Comput. Phys.* **188** (2003) 399–433.
6. M. D. Bronshtein, Smoothness of roots of polynomials depending on parameters, *Sib. Mat. Zh.* **20** (1979) 493–501, in Russian [English transl., *Siberian Math. J.* **20** (1980) 347–352].
7. M. D. Bronshtein, The Cauchy problem for hyperbolic operators with characteristics of variable multiplicity, *Trudy Moskov. Mat. Obshch.* **41** (1980) 83–99, in Russian [English transl., *Trans. Moscow Math. Soc.*, No. 1 (1982) 87–103].
8. J. Chazarain and A. Piriou, *Introduction à la Théorie des Équations aux Dérivées Partielles Linéaires* (Gauthier-Villars, 1981).
9. W. C. Chew and W. H. Weedon, A 3d perfectly matched medium from modified Maxwell's equations with stretched coordinates, *IEEE Microwave and Optical Tech. Lett.* **17** (1995) 599–604.
10. J. Diaz and P. Joly, A time domain analysis of PML models in acoustics, *Comput. Methods Appl. Mech. Engrg.* **195** (2006) 3820–3853.
11. R. Hersh, Mixed problems in several variables, *J. Math. Mech.* **12** (1963) 317–334.
12. J. S. Hesthaven, On the analysis and construction of perfectly matched layers for the linearized Euler equations, *J. Comput. Phys.* **142** (1998) 129–147.

13. L. Hörmander, *The Analysis of Linear Partial Differential Operators. II. Differential Operators with Constant Coefficients* (Springer-Verlag, 2005).
14. F. Q. Hu, On absorbing boundary conditions of linearized Euler equations by a perfectly matched layer, *J. Comput. Phys.* **129** (1996) 201–219.
15. F. Q. Hu, A stable, perfectly matched layer for linearized Euler equations in unsplit physical variables, *J. Comput. Phys.* **173** (2001) 455–480.
16. M. Israeli and S. A. Orszag, Approximation of radiation boundary conditions, *J. Comput. Phys.* **41** (1981) 115–135.
17. J. L. Joly, G. Métivier and J. Rauch, Hyperbolic domains of determinacy and Hamilton–Jacobi equations, *J. Hyp. Part. Differential Eqns.* **2** (2005) 713–744.
18. K. Kasahara, On weak well-posedness of mixed problems for hyperbolic systems, *Publ. Res. Inst. Math. Sci.* **6** (1970/71) 503–514.
19. H.-O. Kreiss and J. Lorenz, *Initial-Boundary Value Problems and the Navier–Stokes Equations* (Academic Press, 1989).
20. P. A. Mazet, S. Paintandre and A. Rahmouni, Interprétation dispersive du milieu PML de Bérenger, *C. R. Math. Acad. Sci. Paris* **327** (1998) 59–64.
21. L. Métivier, Utilisation des équations Euler–PML en milieu hétérogène borné pour la résolution d’un problème inverse en géophysique, in *ESAIM Proc, CANUM 2008*, (EDP Sciences, 2009), pp. 156–170.
22. J. Méttral and O. Vacus, Caractère bien posé du problème de Cauchy pour le système de Bérenger, *C. R. Math. Acad. Sci. Paris* **328** (1999) 847–852.
23. T. Nishitani, Energy inequality for non strictly hyperbolic operators in Gevrey class, *J. Math. Kyoto Univ.* **23** (1983) 739–773.
24. T. Nishitani, Sur les équations hyperboliques à coefficients hölderiens en t et de classe Gevrey en x , *Bull. Sci. Math.* **107** (1983) 113–138.
25. S. Petit-Bergez, Problèmes faiblement bien posés: Discrétisation et applications, Ph.D. thesis, Université Paris 13, 2006. <http://tel.archives-ouvertes.fr/tel-00545794/fr/>.
26. P. G. Petropoulos, Reflectionless sponge layers as absorbing boundary conditions for the numerical solution of Maxwell’s equations in rectangular, cylindrical and spherical coordinates, *SIAM J. Appl. Math.* **60** (2000) 1037–1058.
27. L. Zhao, P. G. Petropoulos and A. C. Cangellaris, A reflectionless sponge layer absorbing boundary condition for the solution of Maxwell’s equations with high-order staggered finite difference schemes, *J. Comput. Phys.* **139** (1998) 184–208.
28. A. Rahmouni, Un modèle PML bien posé pour l’élastodynamique anisotrope, *C. R. Math. Acad. Sci. Paris* **338** (2004) 963–968.
29. C. M. Rappaport, Interpreting and improving the PML absorbing boundary condition using anisotropic lossy mapping of space, *IEEE Trans. Magn.* **32** (1996) 968–974.
30. M. Taylor, *Pseudodifferential Operators* (Princeton Univ. Press, 1981).