

A FAST SECOND-ORDER DISCRETIZATION SCHEME FOR THE LINEARIZED GREEN-NAGHDI SYSTEM WITH ABSORBING BOUNDARY CONDITIONS

GANG PANG^{1,*}, SONGSONG JI² AND XAVIER ANTOINE³

Abstract. In this paper, we present a fully discrete second-order finite-difference scheme with fast evaluation of the convolution involved in the absorbing boundary conditions to solve the one-dimensional linearized Green-Naghdi system. The Padé expansion of the square-root function in the complex plane is used to implement the fast convolution. By introducing a constant damping parameter into the governing equations, the convergence analysis is developed when the damping term fulfills some conditions. In addition, the scheme is stable and leads to a highly reduced computational cost and low memory storage. A numerical example is provided to support the theoretical analysis and to illustrate the performance of the fast numerical scheme.

Mathematics Subject Classification. 76B15, 65M06, 65R10.

Received August 5, 2021. Accepted May 19, 2022.

1. INTRODUCTION

Under the effect of the gravity, the motion of an irrotational and incompressible fluid is described by the free-surface Euler equations. Because of the complexity of this system, asymptotic models for the water wave problem [19] were developed over the years. In particular, the Green-Naghdi model [13] includes the dispersive effects and writes in the two-dimensional space as

$$\begin{cases} H_t + \operatorname{div}(H\vec{U}) = 0, \\ \left(H\vec{U}\right)_t + \operatorname{div}(H\vec{U} \otimes \vec{U} + pI) = 0, \\ p = \frac{gH^2}{2} + \frac{1}{3}H^2\ddot{H}, \end{cases} \quad \mathbf{x} \in \mathbb{R}^2, t > 0. \quad (1.1)$$

In the above notations, we designate by H the fluid depth and by \vec{U} the depth-averaged horizontal velocity, $(v)_t$ stands for the derivation of v with respect to the time variable t and \dot{v} is the material derivative. System

Keywords and phrases. Linearized Green-Naghdi system, absorbing boundary conditions, convolution quadrature, Padé approximation; fast algorithm, convergence analysis.

¹ School of Mathematical Science, Beihang University, Beijing 100083, P.R. China.

² HEDPS, CAPT, and LTCS, College of Engineering, Peking University, Beijing 100871, P.R. China.

³ Université de Lorraine, CNRS, Inria, IECL, F-54000 Nancy, France.

*Corresponding author: gangpang@buaa.edu.cn

(1.1) describes the bidirectional propagation of dispersive water waves in the shallow water regime [19]. A one-dimensional simplification of (1.1), linearized around the steady state $(H, u) := (H_0, 0) + (u_1, u_2)$, with $|(u_1, u_2)| \ll 1$, can be derived as the one-dimensional linearized Green-Naghdi (GN) system [18]

$$\begin{cases} (u_1)_t + (u_2)_x = 0, \\ (u_2)_t + (u_1)_x = \kappa(u_2)_{xxt}, \\ u_1(x, 0) = u_1(x), u_2(x, 0) = u_2(x), \\ \lim_{|x| \rightarrow +\infty} u_1(x, t) = 0, \lim_{|x| \rightarrow +\infty} u_2(x, t) = 0, \end{cases} \quad \begin{array}{l} x \in \mathbb{R}, t > 0, \\ x \in \mathbb{R}, \\ t > 0. \end{array} \quad (1.2)$$

Since (1.2) is set in an unbounded domain, then suitable boundary conditions need to be introduced to get a finite spatial computational domain in view of a discretization. This problem is well-known in the literature and fits into the framework of designing Absorbing Boundary Conditions (ABCs) and Perfectly Matched Layers (when a fictitious layer is added) for systems of PDEs. We refer to [4, 8, 14, 26] for overviews of the various approaches that can be used, their pros/cons and the related discretization aspects and computational difficulties that are met. For the GN system (1.2), a major contribution has been recently achieved by Kazakova and Noble in [18]. In this work, exact ABCs for the fully discrete system of linearized Green-Naghdi equations (1.2) was proposed on staggered grids and the stability of the exact ABCs was proved. Nevertheless, and similarly to the one-dimensional linear Schrödinger equation, these absorbing boundary conditions require the expensive evaluation of nonlocal time convolution-type operators at the fictitious boundary points that lead to prohibitive computational costs and memory storage, in particular for long time computations. In addition, instabilities may arise if the evaluation of the ABCs through the \mathcal{Z} -transform is not carefully implemented (see *e.g.* [1, 9, 17, 23, 27] for some examples).

Probably the most emblematic linear dispersive PDE analyzed in the literature for constructing ABCs is the Schrödinger equation. The first scheme with ABCs was introduced by Baskakov and Popov [11] for computational acoustics based on the parabolic equation. It is now well-known that the resulting scheme with ABCs suffered from stability issues [4]. Since then, many developments allowed to solve numerous problems related to the one-dimensional case [2, 4, 5, 8, 9, 12, 16, 17], and extensions of some of the methods were proposed for higher-dimensional and nonlinear problems [3, 6, 7, 10, 15, 20, 23]. In addition, some of these contributions also include the numerical analysis of the schemes and the derivation of fast and stable evaluation schemes of the nonlocal half-order time derivative operator arising in the definition of the ABCs. The case of the GN system remains much less studied [18] and therefore requires still more understanding. The aim of the present paper is to contribute to the development of ABCs for (1.2) by deriving alternative efficient formulations, and to carefully analyze the convergence of the fully discrete scheme.

To this end, we consider the recent approach introduced by Li *et al.* [21] for the one-dimensional Schrödinger equation and extend it to the GN system (1.2). More precisely, to overcome the numerical instability, a convergent numerical method, integrating a fast evaluation of the exact ABC, is proposed for solving the Cauchy problem (1.2). To this end, the GN system is first reformulated in Section 2 under an equivalent form by introducing a constant damping term, and then a modified Crank–Nicolson scheme is built in Section 3 to discretize the equivalent problem according to the time variable. More specifically, a semi-discrete ABC for the temporally discretized problem is derived for the Crank–Nicolson scheme based on the \mathcal{Z} -transform. Then, a second-order finite difference-scheme for the full spatial discretization is considered. A fast algorithm is introduced in Section 4 to approximate the discrete convolution kernel involved in the exact semi-discrete ABC by using the Padé rational expansion of the square-root function [22]. The damping parameter and the number of Padé terms are chosen to maintain the convergence order of the resulting discrete scheme [21]. In Section 5, a convergence analysis for the proposed numerical method is developed, showing that the scheme is second-order both in space and time. A numerical example illustrates the properties of the scheme in Section 6. Section 7 finally concludes the paper.

2. EXACT ABCS FOR THE ONE-DIMENSIONAL LINEARIZED GREEN-NAGHDI SYSTEM

We consider the initial-value problem for the linearized GN system set on the whole space

$$\begin{aligned} \partial_t u_1(x, t) + \partial_x u_2(x, t) &= 0, \\ \partial_t u_2(x, t) + \partial_x u_1(x, t) &= \kappa \partial_{xx} u_2(x, t), & \forall x \in \mathbb{R}, \forall t > 0, \\ u_1(x, 0) &= u_1(x), u_2(x, 0) = u_2(x), & \forall x \in \mathbb{R}, \\ \lim_{|x| \rightarrow +\infty} u_1(x, t) &= 0, \lim_{|x| \rightarrow +\infty} u_2(x, t) = 0, & \forall t > 0. \end{aligned} \quad (2.1)$$

Let us introduce the new unknown functions: $v_i(x, t) = e^{-\sigma t} u_i(x, t)$, $i = 1, 2$, where $\sigma > 0$ is a parameter used to later control the stability of the fast algorithm. It is straightforward to check that the function $v_i(x, t)$ solves the following initial-value problem:

$$\begin{aligned} \partial_t v_1(x, t) + \sigma v_1(x, t) + \partial_x v_2(x, t) &= 0, \\ \partial_t v_2(x, t) + \sigma v_2(x, t) + \partial_x v_1(x, t) &= \kappa \partial_{xx} (\partial_t v_2(x, t) + \sigma v_2(x, t)), & \forall x \in \mathbb{R}, \forall t > 0, \\ v_1(x, 0) &= u_1(x), v_2(x, 0) = u_2(x), & \forall x \in \mathbb{R}, \\ \lim_{|x| \rightarrow +\infty} v_1(x, t) &= 0, \lim_{|x| \rightarrow +\infty} v_2(x, t) = 0, & \forall t > 0. \end{aligned} \quad (2.2)$$

To obtain exact ABCs for (2.2), we first consider the following exterior problem on the semi-infinite interval $[x_+, +\infty)$:

$$\partial_t v_1(x, t) + \sigma v_1(x, t) + \partial_x v_2(x, t) = 0, \quad (2.3a)$$

$$\partial_t v_2(x, t) + \sigma v_2(x, t) + \partial_x v_1(x, t) = \kappa \partial_{xx} (\partial_t v_2(x, t) + \sigma v_2(x, t)), \quad \forall x \in [x_+, +\infty), \forall t > 0, \quad (2.3b)$$

$$v_1(x, 0) = 0, v_2(x, 0) = 0, \quad \forall x \in [x_+, +\infty), \quad (2.3c)$$

$$\lim_{x \rightarrow +\infty} v_1(x, t) = 0, \lim_{x \rightarrow +\infty} v_2(x, t) = 0, \quad \forall t > 0. \quad (2.3d)$$

The Laplace transform in time (denoted by $\widehat{v}(s)$ for a function $v(t)$) of system (2.3) yields

$$(s + \sigma) \widehat{v}_1(x, s) + \partial_x \widehat{v}_2(x, s) = 0, \quad (2.4)$$

$$(s + \sigma) \widehat{v}_2(x, s) + \partial_x \widehat{v}_1(x, s) = \kappa (s + \sigma) \partial_{xx} \widehat{v}_2(x, s), \quad \forall x \in [x_+, \infty), \forall s \in \mathbb{C}_+, \quad (2.5)$$

$$\lim_{x \rightarrow +\infty} \widehat{v}_1(x, s) = 0, \lim_{x \rightarrow +\infty} \widehat{v}_2(x, s) = 0, \quad \forall s \in \mathbb{C}_+, \quad (2.6)$$

where \mathbb{C}_+ stands for the right half-space of complex numbers with positive real part. After some manipulations, one has: $(s + \sigma)^2 \widehat{v}_2(x, s) = (\kappa (s + \sigma)^2 + 1) \partial_{xx} \widehat{v}_2(x, s)$, $i = 1, 2$, which means that the general solution writes

$$\widehat{v}_2(x, s) = c_1(s) \exp(-x \sqrt{S(s)}) + c_2(s) \exp(x \sqrt{S(s)}), \quad (2.7)$$

where $\sqrt{\cdot}$ denotes the square-root with nonnegative real part and

$$S(s) := \frac{(s + \sigma)^2}{1 + \kappa (s + \sigma)^2}. \quad (2.8)$$

The behavior of the wave at infinity (2.6) implies that $c_2(s) = 0$. By differentiating (2.7), we obtain

$$\partial_x \widehat{v}_2(x, s) = -\sqrt{S(s)} \widehat{v}_2(x, s), \quad \forall x \in [x_+, +\infty), \forall s \in \mathbb{C}_+, \quad (2.9)$$

whose inverse Laplace transform yields an exact absorbing boundary condition for v_2 at x_+

$$(\mathcal{T} * v_2)(x_+, t) := -\mathcal{L}^{-1} \left[\sqrt{S(s)} \widehat{v}_2(x_+, s) \right] (t) = \partial_x v_2(x_+, t), \quad \forall t > 0. \quad (2.10)$$

In (2.10), \mathcal{L}^{-1} denotes the inverse Laplace transform with respect to the variable s and we have

$$(\mathbb{T} * v_2(x_+, \cdot))(t) = \int_0^t \mathbb{T}(t-s)v_2(x_+, s) \, ds.$$

A similar boundary condition can be derived at a left point $x_- < x_+$:

$$(\mathbb{T} * v_2)(x_-, t) + \partial_x v_2(x_-, t) = 0, \quad \forall t > 0. \quad (2.11)$$

These boundary conditions degenerate to the boundary conditions of the wave equation when $\kappa \rightarrow 0$.

In view of (2.10) and (2.11), the solution of (2.2) is the same as the solution of the following problem in the bounded domain (x_-, x_+)

$$\begin{aligned} \partial_t v_1(x, t) + \sigma v_1(x, t) + \partial_x v_2(x, t) &= 0, \\ \partial_t v_2(x, t) + \sigma v_2(x, t) + \partial_x v_1(x, t) &= \kappa \partial_{xx}(\partial_t v_2(x, t) + \sigma v_2(x, t)), & \forall x \in (x_-, x_+), \forall t > 0, \\ (\mathbb{T} * v_2)(x_{\pm}, t) &= \partial_{\nu} v_2(x_{\pm}, t), & \forall t > 0, \\ \partial_t v_1(x_{\pm}, t) + \sigma v_1(x_{\pm}, t) \pm \partial_{\nu} v_2(x_{\pm}, t) &= 0, \\ v_1(x, 0) = u_1(x), v_2(x, 0) &= u_2(x), & \forall x \in [x_-, x_+], \end{aligned} \quad (2.12)$$

where ∂_{ν} denotes the outward normal derivative at the boundary points x_{\pm} .

3. DISCRETIZATION OF THE ONE-DIMENSIONAL GN SYSTEM WITH EXACT SEMI-DISCRETE ABC

In this section, we discretize the one-dimensional linearized Green-Naghdi system in time by using the Crank-Nicolson scheme and derive an associated exact semi-discrete ABC. Then, we propose a second-order finite-difference scheme for the spatial discretization. To this end, we first introduce the notations related to the \mathcal{Z} -transform in the following subsection.

3.1. The \mathcal{Z} -transform of a sequence of functions

Let us consider a Hilbert space \mathcal{H} equipped with an inner product $(\cdot, \cdot)_{\mathcal{H}}$ and induced norm $\|\cdot\|_{\mathcal{H}}$. We introduce the semi-infinite sequence spaces:

$$\ell^2(\mathcal{H}) = \left\{ u = \{u^n\}_{n=0}^{\infty} : u^n \in \mathcal{H}, \|u\|_{\ell^2(\mathcal{H})} \equiv \left(\sum_{n=0}^{\infty} |u^n|^2 \right)^{\frac{1}{2}} < \infty \right\}, \quad (3.1)$$

and

$$\ell_0^2(\mathcal{H}) = \{u = \{u^n\}_{n=0}^{\infty} \in \ell^2(\mathcal{H}) : u^0 = 0\}, \quad (3.2)$$

with the inner product: $(u, v)_{\ell^2(\mathcal{H})} \equiv \sum_{n=0}^{\infty} (u^n, v^n)_{\mathcal{H}}, \forall u, v \in \ell^2(\mathcal{H})$. For any element $u = \{u^n\}_{n=0}^{\infty} \in \ell^2(\mathcal{H})$, we define its \mathcal{Z} -transform as $\tilde{u}(z) = \sum_{n=0}^{\infty} u^n z^n$. The following Parseval's identity holds:

$$(u, v)_{\ell^2(\mathcal{H})} = \int_{\partial \mathbb{D}} (\tilde{u}(z), \tilde{v}(z))_{\mathcal{H}} \nu(dz), \quad \forall u, v \in \ell^2(\mathcal{H}). \quad (3.3)$$

In the above, ν stands for the normalized Haar measure on the unit circle $\partial \mathbb{D}$ of the complex plane, and $\nu(dz) = \frac{1}{2\pi} d\theta$ through the change of variable $z = e^{i\theta}$, with $\theta \in [-\pi, \pi)$.

For a sequence $u = \{u^n\}_{n=0}^{\infty} \in \ell^2(\mathcal{H})$, we define the operator S by: $Su = \{u^{n+1}\}_{n=0}^{\infty}$. The average operator E and the forward difference quotient operator D_{τ} with step τ are given by $E = (S+I)/2$ and $D_{\tau} = (S-I)/\tau$, respectively. We also introduce the following notations: $Su^n = (Su)^n$, $Eu^n = (Eu)^n$ and $D_{\tau}u^n = (D_{\tau}u)^n$.

3.2. Exact ABCs for the semi-discretized one-dimensional linearized GN system

Let $\tau > 0$ be the uniform time step such that $N\tau = T = t_N$, with T the maximal time of computation. Let us set the discrete times as: $t_n = n\tau$, $0 \leq n \leq N$. System (2.2) is semi-discretized in time following

$$\begin{aligned} (D_\tau + \sigma E)v_1^n(x) + \partial_x E v_2^n(x) &= 0, \\ (D_\tau + \sigma E)v_2^n(x) + \partial_x E v_1^n(x) &= \kappa \partial_{xx}(D_\tau + \sigma E)v_2^n(x), & \forall x \in \mathbb{R}, \forall n \geq 0, \\ v_1^0(x) = u_1(x), v_2^0(x) &= u_2(x), & \forall x \in \mathbb{R}, \\ \lim_{|x| \rightarrow +\infty} v_1^n(x) = 0, \lim_{|x| \rightarrow +\infty} v_2^n(x) &= 0, & \forall n \geq 1, \end{aligned} \quad (3.4)$$

where $v_i^n(x) \approx v_i(x, t_n)$, for $i = 1, 2$. We now assume that the initial data u_1 and u_2 are compactly supported in the interval $[x_-, x_+]$. On $[x_+, +\infty)$, the semi-discrete problem (3.4) reduces to

$$\begin{aligned} (D_\tau + \sigma E)v_1^n(x) + \partial_x E v_2^n(x) &= 0, \\ (D_\tau + \sigma E)v_2^n(x) + \partial_x E v_1^n(x) &= \kappa \partial_{xx}(D_\tau + \sigma E)v_2^n(x), & \forall x \in [x_+, +\infty), \forall n \geq 0, \\ v_1^0(x) = 0, v_2^0(x) &= 0, & \forall x \in [x_+, +\infty), \\ \lim_{x \rightarrow +\infty} v_1^n(x) = 0, \lim_{x \rightarrow +\infty} v_2^n(x) &= 0, & \forall n \geq 1. \end{aligned} \quad (3.5)$$

Let us denote by $\tilde{u}(x, z)$ the \mathcal{Z} -transform of the sequence $\{u^n(x)\}_{n=0}^\infty$. Applying the \mathcal{Z} -transform to (3.5), we obtain

$$\begin{aligned} \frac{(2 - 2z + \sigma\tau(1 + z))^2}{\tau^2(1 + z)^2 + \kappa(2 - 2z + \sigma\tau(1 + z))^2} \tilde{v}_2(x, z) - \partial_{xx} \tilde{v}_2(x, z) &= 0, & \forall x \in [x_+, +\infty), \\ \lim_{x \rightarrow +\infty} \tilde{v}_2(x, z) &= 0, \end{aligned}$$

with general solution $\tilde{v}_2(x, z) = \tilde{c}_1 \exp(-x\sqrt{s(z)}) + \tilde{c}_2 \exp(x\sqrt{s(z)})$, setting, for all $\tau > 0$ and $\sigma > 0$,

$$s(z) = \frac{(2 - 2z + \sigma\tau(1 + z))^2}{\tau^2(1 + z)^2 + \kappa(2 - 2z + \sigma\tau(1 + z))^2} = S\left(\frac{2(1 - z)}{\tau(1 + z)}\right). \quad (3.6)$$

The condition at infinity, i.e. $\lim_{x \rightarrow +\infty} \tilde{v}_2(x, z) = 0$, implies that $\tilde{c}_2 = 0$, leading to

$$\partial_x \tilde{v}_2(x_+, z) = -\sqrt{s(z)} \tilde{v}_2(x_+, z), \quad \forall z \in \mathbb{D}, \quad (3.7)$$

which corresponds to the semi-discretization of (2.9) at $x = x_+$. Note that the function

$$\tilde{T}(z) = -\sqrt{s(z)} \quad (3.8)$$

is analytic in the unit disk \mathbb{D} . Thus, it admits a power series expansion

$$\tilde{T}(z) = \sum_{j=0}^{\infty} \mathcal{T}_j z^j, \quad \forall z \in \mathbb{D}. \quad (3.9)$$

Substituting (3.9) and $\tilde{v}_2(x, z) = \sum_{n=0}^{\infty} v_2^n(x) z^n$ into (3.7) yields an exact absorbing boundary condition for (3.4) at the right fictitious boundary point $x = x_+$:

$$(\mathcal{T} * v_2)^n(x_+) - \partial_x v_2^n(x_+) = 0, \quad \forall n \geq 0, \quad (3.10)$$

where $\mathcal{T} *$ is the convolution quadrature operator corresponding to the symbol $\tilde{T}(z)$, namely,

$$(\mathcal{T} * v_2)^n = \sum_{j=0}^n \mathcal{T}_j v_2^{n-j}. \quad (3.11)$$

To simplify the notations, for a function $v(x, t)$, we set: $(\mathcal{T} * v)(x, t_n) = \sum_{j=0}^n \mathcal{T}_j v(x, t_{n-j})$. The boundary condition (3.10) is the semi-discretization of (2.10).

Analogously, by analyzing (3.4) on $(-\infty, x_-]$, we derive an exact semi-discrete absorbing boundary condition at the left point $x = x_-$

$$(\mathcal{T} * v_2)^n(x_-) + \partial_x v_2^n(x_-) = 0, \quad \forall n \geq 1.$$

Consequently, the semi-discrete problem (3.4), originally defined on the whole space, can be reduced to the following semi-discrete problem on a bounded domain:

$$\begin{aligned} (D_\tau + \sigma E)v_1^n(x) + \partial_x E v_2^n(x) &= 0, \\ (D_\tau + \sigma E)v_2^n(x) + \partial_x E v_1^n(x) &= \kappa \partial_{xx}(D_\tau + \sigma E)v_2^n(x), & \forall x \in (x_-, x_+), \quad \forall n \geq 0, \\ (\mathcal{T} * v_2)^n(x_\pm) &= \partial_\nu v_2^n(x_\pm), & \forall n \geq 0, \\ (D_\tau + \sigma E)v_1^n(x_\pm) &= \mp \partial_\nu v_2^n(x_\pm), & \forall n \geq 0, \\ v_1^0(x) &= u_1(x), v_2^0(x) = u_2(x), & \forall x \in [x_-, x_+]. \end{aligned} \quad (3.12)$$

3.3. Spatial discretization

Let M be a positive integer and $h = (x_+ - x_-)/M$ the uniform mesh size. We define the mesh points: $x_k = x_- + (k - 1/2)h$, for $k = 0, 1, \dots, M + 1$, and $x_{k+1/2} = x_0 + (k + 1/2)h$, for $k = 0, 1, \dots, M$, where x_0 and x_{M+1} are two ghost points. In the time-stepping scheme (3.12), we use $(v_2)_k^n$ to denote the numerical approximation of $v_2^n(x_k)$, with $0 \leq k \leq M + 1$, and $(v_1)_k^n$ to define that of $v_1^n(x_{k-1/2})$, with $1 \leq k \leq M + 1$. Let $(v_2)^n = ((v_2)_0^n, \dots, (v_2)_{M+1}^n)$ and $(v_1)^n = ((v_1)_1^n, \dots, (v_1)_{M+1}^n)$. Being given a vector $\chi = (\chi_1, \dots, \chi_{M+1}) \in \mathbb{R}^{M+1}$ or $\omega = (\omega_0, \dots, \omega_{M+1}) \in \mathbb{R}^{M+2}$, we introduce the discrete gradients $\nabla_h \chi$ and $\nabla_h \omega$ such that

$$\begin{aligned} \nabla_h \chi &= \left(\frac{\chi_2 - \chi_1}{h}, \frac{\chi_3 - \chi_2}{h}, \dots, \frac{\chi_{M+1} - \chi_M}{h} \right), \\ \nabla_h \omega &= \left(\frac{\omega_1 - \omega_0}{h}, \frac{\omega_2 - \omega_1}{h}, \dots, \frac{\omega_{M+1} - \omega_M}{h} \right), \end{aligned}$$

respectively. The linear operator which maps the $(M + 2)$ -dimensional vector $\omega = (\omega_0, \dots, \omega_{M+1})$ to the M -dimensional vector $(\omega_1, \dots, \omega_M)$ will be denoted by \mathcal{P} . In addition, we introduce the Neumann and Dirichlet data associated with the $(M + 2)$ -dimensional vector ω as

$$\partial_\nu^- \omega = \frac{\omega_0 - \omega_1}{h}, \quad \partial_\nu^+ \omega = \frac{\omega_{M+1} - \omega_M}{h}, \quad \gamma^- \omega = \frac{\omega_0 + \omega_1}{2}, \quad \gamma^+ \omega = \frac{\omega_{M+1} + \omega_M}{2}.$$

Let us define the inner product for two M -dimensional vectors $\phi_1 = ((\phi_1)_1, \dots, (\phi_1)_M)$ and $\phi_2 = ((\phi_2)_1, \dots, (\phi_2)_M)$ by $(\phi_1, \phi_2)_h = h \sum_{k=1}^M \overline{(\phi_1)_k} (\phi_2)_k$, the inner product for two $(M + 2)$ -dimensional vectors $\omega_1 = ((\omega_1)_0, \dots, (\omega_1)_{M+1})$ and $\omega_2 = ((\omega_2)_0, \dots, (\omega_2)_{M+1})$ by

$$\langle \omega_1, \omega_2 \rangle_h = \frac{h}{2} \overline{(\omega_1)_0} (\omega_2)_0 + h \sum_{k=1}^M \overline{(\omega_1)_k} (\omega_2)_k + \frac{h}{2} \overline{(\omega_1)_{M+1}} (\omega_2)_{M+1},$$

and finally the inner product for two $(M + 1)$ -dimensional vectors $\chi_1 = ((\chi_1)_1, \dots, (\chi_1)_{M+1})$ and $\chi_2 = ((\chi_2)_1, \dots, (\chi_2)_{M+1})$ according to: $\{\chi_1, \chi_2\}_h = h \sum_{k=1}^{M+1} \overline{(\chi_1)_k} (\chi_2)_k$. In the above expressions, \bar{z} denotes the complex conjugate of a complex number z . The induced norms are such that: $\|\phi\|_h = \sqrt{(\phi, \phi)_h}$, $|\omega|_h = \sqrt{\langle \omega, \omega \rangle_h}$, and $\|\chi\|_h = \sqrt{\{\chi, \chi\}_h}$.

Let us now introduce a second-order spatial discretization Δ_h , which maps the $(M + 2)$ -dimensional vector ω to the M -dimensional vector space as

$$\Delta_h \omega = \left(\frac{\omega_0 - 2\omega_1 + \omega_2}{h^2}, \dots, \frac{\omega_{M-1} - 2\omega_M + \omega_{M+1}}{h^2} \right).$$

Thus, we have the discrete integration by parts formula

$$(\mathcal{P}\omega_2, \Delta_h \omega_1)_h = -\langle \nabla_h \omega_2, \nabla_h \omega_1 \rangle_h + \overline{\gamma^+ \omega_2} \partial_\nu^+ \omega_1 + \overline{\gamma^- \omega_2} \partial_\nu^- \omega_1. \quad (3.13)$$

Now, we define the vector $\nabla_h \mathcal{H}v_1^n$ by

$$\nabla_h \mathcal{H}v_1^n = \left(\frac{\mathcal{H}(v_1)_2^n - \mathcal{H}(v_1)_1^n}{h}, \dots, \frac{\mathcal{H}(v_1)_{M+1}^n - \mathcal{H}(v_1)_M^n}{h} \right),$$

where \mathcal{H} is any operator that only applies in the time direction. Similarly, the vector $\nabla_h \mathcal{H}v_2^n$ is such that

$$\nabla_h \mathcal{H}v_2^n = \left(\frac{\mathcal{H}(v_2)_1^n - \mathcal{H}(v_2)_0^n}{h}, \dots, \frac{\mathcal{H}(v_2)_{M+1}^n - \mathcal{H}(v_2)_M^n}{h} \right).$$

We can also introduce a vector $\Delta_h \mathcal{H}v_2^n$ by

$$\Delta_h \mathcal{H}v_2^n = \left(\frac{\mathcal{H}(v_2)_0^n - 2\mathcal{H}(v_2)_1^n + \mathcal{H}(v_2)_2^n}{h^2}, \dots, \frac{\mathcal{H}(v_2)_{M-1}^n - 2\mathcal{H}(v_2)_M^n + \mathcal{H}(v_2)_{M+1}^n}{h^2} \right).$$

Then, it is easy to see that the following identities hold: $\nabla_h \mathcal{H}v_1^n = \mathcal{H}\nabla_h v_1^n$, $\nabla_h \mathcal{H}v_2^n = \mathcal{H}\nabla_h v_2^n$, and $\Delta_h \mathcal{H}v_2^n = \mathcal{H}\Delta_h v_2^n$.

Now, in (3.12), replacing the function $v_1^n(x)$ by the vector $v_1^n = ((v_1)_1^n, \dots, (v_1)_{M+1}^n)$, $v_2^n(x)$ by $v_2^n = ((v_2)_0^n, \dots, (v_2)_{M+1}^n)$ and changing the continuous operator ∂_{xx} with its discrete analogue Δ_h , we obtain the following fully discrete finite-difference scheme

$$\begin{aligned} (D_\tau + \sigma E)v_1^n + \nabla_h^n E v_2^n &= 0, \\ (D_\tau + \sigma E)\mathcal{P}v_2^n + \nabla_h E v_1^n &= \kappa \Delta_h (D_\tau + \sigma E)v_2^n, \\ (\mathcal{T} * \gamma^\pm v_2)^n - \partial_\nu^\pm v_2^n &= 0, \\ v_1^0 &= (u_1(x_{1/2}), \dots, u_1(x_{M+1/2})), \quad v_2^0 = (u_2(x_0), \dots, u_2(x_{M+1})). \end{aligned} \quad \forall n \geq 0, \quad (3.14)$$

In fact, at x_\pm , we have $\partial_x v_2(x_\pm, t) = \pm \partial_\nu v_2(x_\pm, t)$. Therefore, we can write that

$$\partial_t v_1(x_\pm, t) + \sigma v_1(x_\pm, t) \pm \partial_\nu v_2(x_\pm, t) = 0.$$

For (3.14), the expression of v_1 at x_1 can be written as

$$(D_\tau + \sigma E)(v_1)_1^n + E \frac{(v_2)_1^n - (v_2)_0^n}{h} = (D_\tau + \sigma E)(v_1)_1^n - E \partial_\nu^- v_2^n = (D_\tau + \sigma E)(v_1)_1^n - E(\mathcal{T} * \gamma^- v_2)^n.$$

Thus, the boundary condition for v_1 can be supplied by v_2 due to the staggered grid.

4. FAST EVALUATION OF THE BOUNDARY DISCRETE CONVOLUTION $(\mathcal{T} * \gamma^\pm v_2)^n$

In this section, we introduce a fast algorithm for approximating the discrete convolution product $(\mathcal{T} * \gamma^\pm v_2)^n$ arising in (3.14). The stability of the proposed fast algorithm will be analyzed in the next section.

4.1. Rational approximation of the convolution quadrature

In [22], for a nonnegative integer $m > 0$ and $\text{Re}(s) \geq -1$, the Padé approximation of the function $\sqrt{1+s}$ can be expressed as

$$\sqrt{1+s} \approx 1 + \sum_{j=1}^m \frac{\alpha_j s}{1 + \beta_j s},$$

where the coefficients are given by

$$\alpha_j = \frac{2}{2m+1} \sin^2\left(\frac{j\pi}{2m+1}\right), \quad \beta_j = \cos^2\left(\frac{j\pi}{2m+1}\right), \quad j = 1, \dots, m.$$

Based on this Padé approximation, a rational approximation of the square-root function \sqrt{s} on the closed right half complex plane can be written as

$$\sqrt{s} = \sqrt{1+s-1} \approx 1 + \sum_{j=1}^m \frac{\alpha_j(s-1)}{1+\beta_j(s-1)} \equiv R_m(s), \quad \operatorname{Re}(s) \geq 0.$$

Thus, we deduce

$$\begin{aligned} R_m(s) &= \lambda - \sum_{j=1}^m \frac{1}{g_j s + h_j}, \quad \lambda = 1 + \sum_{j=1}^m \alpha_j \beta_j^{-1}, \\ h_j &= \alpha_j^{-1} \beta_j (1 - \beta_j), \quad g_j = \alpha_j^{-1} \beta_j^2, \quad j = 1, \dots, m. \end{aligned} \quad (4.1)$$

For all $\tau > 0$, $\sigma > 0$, and for $s(z)$ defined by (3.6), we can introduce the rational approximation $\tilde{T}^{(m)}(z)$ of the symbol $\tilde{T}(z)$ as

$$\tilde{T}^{(m)}(z) := -R_m(s(z)), \quad \forall m \geq 0. \quad (4.2)$$

We denote by $\mathcal{T}^{(m)*}$ the convolution operator analogously defined as (3.11) by replacing the convolution coefficients with the series expansion coefficients of the function $\tilde{T}^{(m)}(z)$. By considering the rational approximation $\mathcal{T}^{(m)*}$ of \mathcal{T}^* in (3.14), we obtain the following fully discrete scheme:

$$(D_\tau + \sigma E)v_1^n + \nabla_h^n E v_2^n = 0, \quad (4.3)$$

$$(D_\tau + \sigma E)\mathcal{P}v_2^n + \nabla_h E v_1^n = \kappa \triangle_h (D_\tau + \sigma E)v_2^n, \quad \forall n \geq 0, \quad (4.4)$$

$$(\mathcal{T}^{(m)} * \gamma^\pm v_2)^n - \partial_\nu^\pm v_2^n = 0, \quad \forall n \geq 0, \quad (4.5)$$

$$v_1^0 = (u_1(x_{1/2}), \dots, u_1(x_{M+1/2})), \quad v_2^0 = (u_2(x_0), \dots, u_2(x_{M+1})). \quad (4.6)$$

In fact, equation (4.5) can be solved by the fast algorithm described in the next subsection.

4.2. Fast evaluation of $(\mathcal{T}^{(m)} * \gamma^\pm v_2)^n$

By applying (4.1) to (4.2), we obtain the sequence of equalities

$$\begin{aligned} \tilde{T}^{(m)}(z) &= -\lambda + \sum_{j=1}^m \frac{1}{g_j s(z) + h_j} = -\lambda + \sum_{j=1}^m \frac{\kappa(2 + \sigma\tau + (\sigma\tau - 2)z)^2 + \tau^2(1+z)^2}{(g_j + \kappa h_j)(2 + \sigma\tau + (\sigma\tau - 2)z)^2 + h_j \tau^2(1+z)^2} \\ &= -\lambda + \sum_{j=1}^m \left(\lambda_j + \frac{e_j z + f_j}{(a_j z + b_j)^2 - (c_j z + d_j)^2} \right) \\ &= -\lambda + \sum_{j=1}^m \lambda_j + \sum_{j=1}^m \left(\frac{A_j}{(a_j + c_j)z + b_j + d_j} + \frac{B_j}{(a_j - c_j)z + b_j - d_j} \right), \end{aligned} \quad (4.7)$$

setting

$$\begin{aligned} \lambda_j &= \frac{\kappa(\tau\sigma - 2)^2 + \tau^2}{(g_j + \kappa h_j)(\tau\sigma - 2)^2 + h_j \tau^2}, \quad e_j = 2(\kappa - g_j \lambda_j - \kappa h_j \lambda_j)(\sigma^2 \tau^2 - 4) + 2\tau^2(1 - \lambda_j h_j), \\ f_j &= (\kappa - g_j \lambda_j - \kappa h_j \lambda_j)(2 + \sigma\tau)^2 + \tau^2(1 - \lambda_j h_j), \end{aligned}$$

$$a_j = \sqrt{\kappa h_j + g_j}(\sigma\tau - 2), \quad b_j = \sqrt{\kappa h_j + g_j}(\sigma\tau + 2), \quad c_j = i\sqrt{h_j}\tau, \quad d_j = i\sqrt{h_j}\tau,$$

$$A_j = \frac{-e_j c_j + b_j f_j - e_j a_j + f_j d_j}{a_j^2 + b_j^2 - c_j^2 - d_j^2}, \quad B_j = \frac{e_j c_j + b_j f_j - e_j a_j - f_j d_j}{a_j^2 + b_j^2 - c_j^2 - d_j^2}.$$

Therefore, we have

$$\tilde{T}^{(m)}(z) = \left(-\lambda + \sum_{j=1}^m \lambda_j\right) + \sum_{j=1}^m \sum_{n=0}^{\infty} \left(\frac{A_j}{b_j + d_j} \left(-\frac{a_j + c_j}{b_j + d_j}\right)^n z^n + \frac{B_j}{b_j - d_j} \left(\frac{c_j - a_j}{b_j - d_j}\right)^n z^n \right).$$

Thus, $\tilde{T}^{(m)}(z)$ can be uniformly rewritten as

$$\tilde{T}^{(m)}(z) = \sum_{k=1}^{2m+1} \sum_{n=0}^{\infty} C_k (\gamma_k)^n z^n,$$

which implies that

$$\mathcal{T}_j^{(m)} = \sum_{k=1}^{2m+1} C_k (\gamma_k)^j,$$

with

$$C_{2k-1} = \frac{A_k}{b_k + d_k}, \quad C_{2k} = \frac{B_k}{b_k - d_k}, \quad \gamma_{2k-1} = -\frac{a_k + c_k}{b_k + d_k} \quad \text{and} \quad \gamma_{2k} = -\frac{-a_k + c_k}{b_k - d_k},$$

for $1 \leq k \leq m$. We can take $C_{2m+1} = (-\lambda + \sum_{k=1}^m \lambda_k)$ and $\gamma_{2m+1} = 0$. Therefore, $(\mathcal{T}^{(m)} * \gamma^{\pm} v_2)^n = \sum_{j=0}^n \mathcal{T}_j^{(m)} (\gamma^{\pm} v_2)^{n-j}$ can be implemented by a fast convolution in (4.5).

For fixed k , by defining

$$\mathcal{G}_k^n[v] := C_k \sum_{j=0}^n (\gamma_k)^{n-j} v^j,$$

with $1 \leq k \leq 2m+1$, we derive that,

$$\begin{aligned} \gamma_k \mathcal{G}_k^{n-1}[\gamma^{\pm} v_2] + C_k (\gamma^{\pm} v_2)^n &= \gamma_k C_k \sum_{j=0}^{n-1} (\gamma_k)^{n-1-j} (\gamma^{\pm} v_2)^j + C_k (\gamma^{\pm} v_2)^n \\ &= C_k \sum_{j=0}^{n-1} (\gamma_k)^{n-j} (\gamma^{\pm} v_2)^j + C_k (\gamma^{\pm} v_2)^n = C_k \sum_{j=0}^n (\gamma_k)^{n-j} (\gamma^{\pm} v_2)^j = \mathcal{G}_k^n[\gamma^{\pm} v_2]. \end{aligned} \quad (4.8)$$

Thus, the boundary term $(\mathcal{T}^{(m)} * \gamma^{\pm} v_2)^n = \sum_{j=0}^n \mathcal{T}_j^{(m)} (\gamma^{\pm} v_2)^{n-j}$ can be written in the form of fast convolution by (4.8), namely,

$$(\mathcal{T}^{(m)} * \gamma^{\pm} v_2)^n = \sum_{k=1}^{2m+1} \mathcal{G}_k^n[\gamma^{\pm} v_2] = \sum_{k=1}^{2m+1} \gamma_k \mathcal{G}_k^{n-1}[\gamma^{\pm} v_2] + (\gamma^{\pm} v_2)^n \sum_{k=1}^{2m+1} C_k. \quad (4.9)$$

From (4.9), we can see that the computational cost at the n -th time step is $\mathcal{O}(2m+1)$, rather than $\mathcal{O}(n)$ which is the computational cost of the convolution. Thus, the overall computational cost for all the n time steps is $\mathcal{O}(n(2m+1))$ instead of $\mathcal{O}(n^2)$. Therefore, for large time steps n , the total computational cost is greatly reduced. Such a cost is much lower than the one proposed in [18] for the linearized Green-Naghdi equation. It is similar to the techniques introduced for example in [21, 27] to optimize both the memory storage and computational cost when evaluating convolutions.

4.3. Properties of the rational approximation $\tilde{T}^{(m)}(z)$

Let us now prove some properties of $\mathcal{T}_j^{(m)}$ used to prove the error estimates.

Proposition 4.1. *Let us assume that the condition $\sigma \geq \frac{1}{\sqrt{2\kappa}}$ is satisfied, the time step τ is small enough and m is sufficiently large, i.e. it fulfills*

$$2m + 1 \geq \frac{\ln \epsilon}{\ln(1 - \delta)}, \quad \text{for some } \epsilon \in \left(0, \frac{\mu\sqrt{\kappa}\tau^3}{8}\right], \quad (4.10)$$

with

$$\mu(\kappa, \sigma) = \frac{\sqrt{2}\sigma}{2\sqrt{1 + \kappa\sigma^2}}, \quad \delta(\kappa, \sigma) = \frac{\sqrt{2}\sigma\sqrt{1 + \kappa\sigma^2}}{\sigma^2 + \sqrt{2}\sigma\sqrt{1 + \kappa\sigma^2} + 1 + \kappa\sigma^2}. \quad (4.11)$$

Then, the following inequalities hold

$$\max_{z \in \partial\mathbb{D}} |\tilde{T}^{(m)}(z) - \tilde{T}(z)| \leq \mu \frac{\tau^3}{2}, \quad (4.12)$$

$$\operatorname{Re} \sum_{k=0}^n \overline{(\tau D_\tau + \sigma E) u^k} (D_\tau + \sigma E) \left(\mathcal{T}^{(m)} * u \right)^k \leq 0, \quad (4.13)$$

$$\operatorname{Re} \sum_{k=0}^n \overline{(D_\tau + \sigma E) u^k} (D_\tau + \sigma E) \left(\mathcal{T}^{(m)} * u \right)^k \leq 0, \quad \forall n \geq 0, \quad (4.14)$$

for any complex-valued sequence $u = \{u^n\}_{n=0}^\infty$ such that $u^0 = 0$.

Before proving Proposition 4.1, we need to state a few lemmas. Let us introduce

$$r(s) := \frac{\sqrt{s} - 1}{\sqrt{s} + 1}. \quad (4.15)$$

Then, by using (3.8), one can prove that the symbol $\tilde{T}(z)$ satisfies the following inequalities.

Lemma 4.2. *We have the following inequalities*

$$\max_{z \in \partial\mathbb{D}} |\tilde{T}(z)| \leq \frac{1}{\sqrt{\kappa}}, \quad \min_{z \in \partial\mathbb{D}} |\tilde{T}(z)| \geq \frac{\sigma}{\sqrt{1 + \kappa\sigma^2}}. \quad (4.16)$$

Furthermore, under the condition $\sigma \geq \frac{1}{\sqrt{2\kappa}}$ and for $s(z)$ defined by (3.6), we can prove that

$$\max_{z \in \partial\mathbb{D}} \operatorname{Re} \tilde{T}(z) \leq -\mu(\kappa, \sigma) = -\frac{\sqrt{2}\sigma}{2\sqrt{1 + \kappa\sigma^2}}, \quad (4.17)$$

$$\max_{z \in \partial\mathbb{D}} |r(s(z))| \leq 1 - \delta(\kappa, \sigma), \quad (4.18)$$

$$-\arctan\left(\frac{1}{\sqrt{\tau}}\right) \leq \arg_{z \in \partial\mathbb{D}} \left[-\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} \tilde{T}(z) \right] \leq \arctan\left(\frac{1}{\sqrt{\tau}}\right). \quad (4.19)$$

Proof of Lemma 4.2. From (2.8), (3.6), (3.8), and setting $\rho = \frac{2(1-z)}{i\tau(1+z)}$, then $\rho \in \mathbb{R}$ for $z \in \partial\mathbb{D}$, and

$$\tilde{T}(z) = -\sqrt{S(i\rho)} = -\left(\frac{(\rho^2 + \sigma^2)^2}{\kappa^2(\rho^2 + \sigma^2)^2 + 1 + 2\kappa\sigma^2 - 2\rho^2\kappa} \right)^{\frac{1}{4}} \exp(i\theta/2), \quad (4.20)$$

where

$$\theta = \arg\left(\frac{2\sigma\rho + i(\rho^2 - \sigma^2)}{2\sigma\rho + i(\rho^2 - \sigma^2 - 1/\kappa)}\right). \quad (4.21)$$

It is straightforward to verify that

$$|\tilde{T}(z)| = \left(\frac{1}{\kappa^2 + (1 + 2\kappa\sigma^2 - 2\rho^2\kappa)/(\rho^2 + \sigma^2)^2}\right)^{\frac{1}{4}}, \quad (4.22)$$

which means that $|\tilde{T}(z)|$ increases with respect to $|\rho|$. Thus, we conclude that

$$\frac{\sigma}{\sqrt{1 + \kappa\sigma^2}} \leq |\tilde{T}(z)| \leq \frac{1}{\sqrt{\kappa}}, \quad \text{if } \rho \in \mathbb{R}, \quad (4.23)$$

hence proving (4.16).

From (4.21) we have

$$\theta = \arctan\left(\frac{2\sigma\rho/\kappa}{\rho^4 + (2\sigma^2 - 1/\kappa)\rho^2 + \sigma^4 + \sigma^2/\kappa}\right) := \arctan\left(\frac{2\sigma\rho/\kappa}{\Theta(\rho)}\right). \quad (4.24)$$

First, we discuss the case $\rho \in [0, +\infty)$. It is straightforward to show that $\Theta(\rho) \geq 0$, for $\sigma \geq \frac{1}{\sqrt{2\kappa}}$. Therefore, from (4.24), we derive that $0 \leq \theta \leq \pi/2$ for $\rho \in [0, +\infty)$, which means $0 \leq \theta/2 \leq \pi/4$. Similarly, we have $-\pi/2 \leq \theta \leq 0$ for $\rho \in (-\infty, 0]$, and then $-\pi/4 \leq \theta/2 \leq 0$. Thus, this yields

$$\operatorname{Re} \tilde{T}(z) \leq -\frac{\sqrt{2}\sigma}{2\sqrt{1 + \kappa\sigma^2}},$$

which proves (4.17).

Recalling that $s(z) = \frac{(2-2z+\sigma\tau(1+z))^2}{\tau^2(1+z)^2+(2-2z+\sigma\tau(1+z))^2}$, we deduce

$$\sqrt{s(z)} = -\tilde{T}(z) = |\tilde{T}(z)| \exp\left(i\frac{\theta}{2}\right), \quad (4.25)$$

with $\cos(\frac{\theta}{2}) \geq \frac{\sqrt{2}}{2}$ for $\rho \in (-\infty, +\infty)$. Using the above expression of $\sqrt{s(z)}$, we have

$$\begin{aligned} |r(s(z))| &= \left|\frac{\sqrt{s(z)} - 1}{\sqrt{s(z)} + 1}\right| = \sqrt{1 - \frac{4|\tilde{T}(z)|\cos(\frac{\theta}{2})}{|\tilde{T}(z)|^2 + 2|\tilde{T}(z)|\cos(\frac{\theta}{2}) + 1}} \\ &\leq \sqrt{1 - \frac{2\sqrt{2}|\tilde{T}(z)|}{|\tilde{T}(z)|^2 + \sqrt{2}|\tilde{T}(z)| + 1}} \leq 1 - \frac{\sqrt{2}|\tilde{T}(z)|}{|\tilde{T}(z)|^2 + \sqrt{2}|\tilde{T}(z)| + 1}, \end{aligned}$$

where the last inequality is a consequence of: $(1-x)^{\frac{1}{2}} = 1 - \frac{1}{2}x - \frac{1}{8}x^2 + \cdots \leq 1 - \frac{1}{2}x$. By considering

$$\frac{\sqrt{2}r}{r^2 + \sqrt{2}r + 1} = \frac{\sqrt{2}}{(r + \sqrt{2} + 1/r)},$$

we see that the minimum value of $\frac{\sqrt{2}|\tilde{T}(z)|}{|\tilde{T}(z)|^2 + \sqrt{2}|\tilde{T}(z)| + 1}$ is attained at $|\tilde{T}(z)| = \frac{\sigma}{\sqrt{1+\kappa\sigma^2}}$, and

$$\frac{\sqrt{2}|\tilde{T}(z)|}{|\tilde{T}(z)|^2 + \sqrt{2}|\tilde{T}(z)| + 1} \geq \frac{\sqrt{2}\sigma\sqrt{1+\kappa\sigma^2}}{\sigma^2 + \sqrt{2}\sigma\sqrt{1+\kappa\sigma^2} + 1 + \kappa\sigma^2} = \delta,$$

leading to (4.18).

Next, we have the sequence of equalities

$$\begin{aligned} & \arg_{z \in \partial \mathbb{D}} \left[-\frac{(z^{-1}-1)/\tau + \sigma(z^{-1}+1)/2}{(z^{-1}-1) + \sigma(z^{-1}+1)/2} \tilde{T}(z) \right] \\ &= \frac{1}{2} \left[\arg_{z \in \partial \mathbb{D}} \frac{\left(\frac{2(z^{-1}-1)}{\tau(z^{-1}+1)} + \sigma \right)^2}{1 + \kappa \left(\frac{2(z^{-1}-1)}{\tau(z^{-1}+1)} + \sigma \right)^2} + 2 \arctan_{z \in \partial \mathbb{D}} \frac{(z^{-1}-1)/\tau + \sigma(z^{-1}+1)/2}{(z^{-1}-1) + \sigma(z^{-1}+1)/2} \right] \\ &= \frac{1}{2} \left[\arctan_{\rho \in \mathbb{R}} \frac{2\sigma\rho/\kappa}{\Theta(\rho)} + 2 \arg_{\rho \in \mathbb{R}} \frac{\rho i + \sigma}{\tau \rho i + \sigma} \right] = \frac{1}{2} \left[\theta(\rho) + 2 \arctan_{\rho \in \mathbb{R}} \frac{\sigma(1-\tau)\rho}{\sigma^2 + \tau\rho^2} \right] \\ &= \arctan_{\rho \in \mathbb{R}} \frac{2\sigma\rho/\kappa}{\sqrt{(\Theta)^2 + (2\sigma\rho/\kappa)^2} + \Theta} + \arctan_{\rho \in \mathbb{R}} \frac{\sigma(1-\tau)\rho}{\sigma^2 + \tau\rho^2} \\ &= \arctan_{\rho \in \mathbb{R}} F_1(\rho) + \arctan_{\rho \in \mathbb{R}} F_2(\rho) = \arctan_{\rho \in \mathbb{R}} \frac{F_1(\rho) + F_2(\rho)}{1 - F_1(\rho)F_2(\rho)}, \end{aligned} \quad (4.26)$$

with

$$F_1(\rho) = \frac{2\sigma\rho/\kappa}{\sqrt{(\Theta)^2 + (2\sigma\rho/\kappa)^2} + \Theta}, \quad F_2(\rho) = \frac{\sigma(1-\tau)\rho}{\sigma^2 + \tau\rho^2}.$$

For $\rho \in [0, +\infty)$, we obtain

$$F_1(\rho)F_2(\rho) \leq \frac{2\sigma^2\rho^2/\kappa}{2\Theta\sigma^2} = \frac{\rho^2/\kappa}{\Theta} \leq \frac{\rho^2/\kappa}{\rho^4 + (2\sigma^2 - 1/\kappa)\rho^2 + \sigma^4 + \sigma^2/\kappa} \leq \frac{1/\kappa}{(2\sigma^2 - 1/\kappa) + 2\sqrt{\sigma^4 + \sigma^2/\kappa}} \leq \frac{1}{\sqrt{3}}. \quad (4.27)$$

In addition, we have

$$\frac{d}{d\rho} F_2(\rho) = \frac{d}{d\rho} \left[\arctan_{\rho \in \mathbb{R}} \frac{\sigma(1-\tau)\rho}{\sigma^2 + \tau\rho^2} \right] = \frac{\sigma(1-\tau)(\sigma^2 - \tau\rho^2)}{(\sigma^2 + \tau\rho^2)^2 + \sigma^2(1-\tau)^2\rho^2}. \quad (4.28)$$

Therefore $F_2(\rho)$ increases with respect to $|\rho|$ in $[0, \frac{\sigma}{\sqrt{\tau}}]$, decreases in $[\frac{\sigma}{\sqrt{\tau}}, \infty)$, and $F_2(\rho)$ reaches its maximum at $\frac{\sigma}{\sqrt{\tau}}$. Then, (4.26) and (4.27) indicate that we have

$$\arctan_{\rho \in \mathbb{R}^+} \frac{F_1(\rho) + F_2(\rho)}{1 - F_1(\rho)F_2(\rho)} \leq \frac{1 + F_2\left(\frac{\sigma}{\sqrt{\tau}}\right)}{1 - \frac{1}{\sqrt{3}}} \leq \frac{1 + \frac{1-\tau}{2\sqrt{\tau}}}{1 - \frac{1}{2}} \leq \frac{1}{\sqrt{\tau}}. \quad (4.29)$$

Similarly, we can write that

$$\arctan_{\rho \in \mathbb{R}^-} \frac{F_1(\rho) + F_2(\rho)}{1 - F_1(\rho)F_2(\rho)} \geq -\frac{1}{\sqrt{\tau}}. \quad (4.30)$$

Combining (4.26), (4.29) and (4.30), we prove (4.19). \square

Let us recall that $R_m(s)$ is defined by (4.1). Then, the following result was proved in [22].

Lemma 4.3. *Let us define: $e_m(s) := \sqrt{s} - R_m(s)$, for $m = 0, 1, 2, \dots$. Then, the following identity holds:*

$$e_m(s) = 2\sqrt{s} \frac{r^{2m+1}(s)}{1 + r^{2m+1}(s)}, \quad \text{if } \operatorname{Re}(s) \geq 0 \text{ and } s \neq 0, \quad (4.31)$$

where $r(s)$ is defined by (4.15).

Let us now consider the following lemma.

Lemma 4.4. *Under the conditions $\sigma \geq \frac{1}{\sqrt{2\kappa}}$ and (4.10), we have the two inequalities*

$$\max_{z \in \partial\mathbb{D}} \operatorname{Re} \tilde{T}^{(m)}(z) \leq 0, \quad (4.32)$$

$$\max_{z \in \partial\mathbb{D}} \operatorname{Re} \left[\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} \tilde{T}^{(m)}(z) \right] \leq 0. \quad (4.33)$$

Proof of Lemma 4.4. From (4.18) in Lemma 4.2 we have: $\max_{z \in \partial\mathbb{D}} |r(s(z))| \leq 1 - \delta(\kappa, \sigma)$. If $\sigma \geq \frac{1}{\sqrt{2\kappa}}$ and m satisfies (4.10), then one gets $|r(s(z))|^{2m+1} \leq [1 - \delta]^{2m+1} \leq 1/2$. From (4.31), we obtain

$$\max_{z \in \partial\mathbb{D}} \left| \frac{\tilde{T}(z) - \tilde{T}^{(m)}(z)}{\tilde{T}(z)} \right| = \max_{z \in \partial\mathbb{D}} \left| \frac{2r^{2m+1}(s(z))}{1 + r^{2m+1}(s(z))} \right| \leq \max_{z \in \partial\mathbb{D}} \frac{2|r(s(z))|^{2m+1}}{1 - |r(s(z))|^{2m+1}} \leq 4 \max_{z \in \partial\mathbb{D}} |r(s(z))|^{2m+1}.$$

Then, equations (4.12), (4.16) and (4.17) imply

$$\max_{z \in \partial\mathbb{D}} \operatorname{Re} \tilde{T}^{(m)}(z) = \max_{z \in \partial\mathbb{D}} \left[\operatorname{Re} \tilde{T}(z) - \operatorname{Re} \left(\tilde{T}(z) - \tilde{T}^{(m)}(z) \right) \right] \leq -\mu + \frac{\mu\sqrt{\kappa}}{2} \max_{z \in \partial\mathbb{D}} |\tilde{T}(z)| \leq -\mu + \mu/2 \leq 0,$$

which proves (4.32).

In addition, for τ small enough, using (4.19), we have

$$\begin{aligned} & \arg_{z \in \partial\mathbb{D}} \left[-\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} \tilde{T}^{(m)}(z) \right] \\ &= \arg_{z \in \partial\mathbb{D}} \left[-\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} \tilde{T}(z) \left(1 + \frac{\tilde{T}^{(m)} - \tilde{T}}{\tilde{T}} \right) \right] \\ &= \arg_{z \in \partial\mathbb{D}} \left[-\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} \tilde{T}(z) \right] + \arg_{z \in \partial\mathbb{D}} \left(1 + \frac{\tilde{T}^{(m)} - \tilde{T}}{\tilde{T}} \right) \\ &\leq \arctan\left(\frac{1}{\sqrt{\tau}}\right) + \arg_{z \in \partial\mathbb{D}} (1 + i\mu\sqrt{\kappa}\tau^3/2) \leq \arctan\left(\frac{1}{\sqrt{\tau}}\right) + \arctan\left(\frac{\sqrt{\tau}}{2}\right) \leq \arctan\left(\frac{4}{\sqrt{\tau}}\right). \end{aligned}$$

Thus, for small enough τ , we deduce (4.33) since

$$-\arctan\left(\frac{4}{\sqrt{\tau}}\right) \leq \arg_{z \in \partial\mathbb{D}} \left(-\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} [\tilde{T}^{(m)}(z)] \right) \leq \arctan\left(\frac{4}{\sqrt{\tau}}\right).$$

□

Now, we can prove Proposition 4.1, as a consequence of (4.32) and (4.33).

Proof of Proposition 4.1. Firstly, for small τ , if $\sigma \leq \frac{1}{\sqrt{2\kappa}}$ and m satisfies (4.10), then we have: $[1 - \delta]^{2m+1} \leq \epsilon \leq 1/2$. Lemma 4.3 then implies

$$\max_{z \in \partial\mathbb{D}} \left| \frac{\tilde{T}(z) - \tilde{T}^{(m)}(z)}{\tilde{T}(z)} \right| \leq \max_{z \in \partial\mathbb{D}} \frac{2|r(s(z))|^{2m+1}}{1 - |r(s(z))|^{2m+1}} \leq 4\epsilon.$$

Consequently, by using (4.16) and (4.10), we have

$$\max_{z \in \partial\mathbb{D}} |\tilde{T}^{(m)}(z) - \tilde{T}(z)| \leq 4\epsilon \max_{z \in \partial\mathbb{D}} |\tilde{T}(z)| \leq \mu \frac{\tau^3}{2},$$

which proves (4.12).

We construct $\{u^k\}_{k=0}^\infty$ such that $(D_\tau + \sigma E)u^k = 0$, for $k \geq n+1$. Thus, one has:

$$u^{k+1} = \frac{1 - \sigma\tau/2}{1 + \sigma\tau/2} u^k,$$

for $k \geq n+1$, which shows that the sequence $\{u^k\}$ is such that $(D_\tau + \sigma E)u^k = 0$, for $k \geq n+1$. From (4.32), we have

$$\begin{aligned} \operatorname{Re} \sum_{k=0}^n \overline{(D_\tau + \sigma E)u^k} \left((D_\tau + \sigma E)T^{(m)} * u \right)^k &= \operatorname{Re} \sum_{k=0}^\infty \overline{(D_\tau + \sigma E)u^k} \left((D_\tau + \sigma E)T^{(m)} * u \right)^k \\ &= \operatorname{Re} \left((D_\tau + \sigma E)u, (D_\tau + \sigma E)T^{(m)} * u \right)_\ell^2(\mathbb{C}) \\ &= \operatorname{Re} \int_{\partial\mathbb{D}} |\tilde{u}(z)|^2 \overline{[z^{-1} - 1 + \sigma\tau(z^{-1} + 1)/2]} \tilde{T}^{(m)}(z) [z^{-1} - 1 + \sigma\tau(z^{-1} + 1)/2] \nu(dz) / \tau^2 \\ &= \operatorname{Re} \int_{\partial\mathbb{D}} |z|^{-2} |\tilde{u}(z)|^2 \overline{[2 - 2z + \sigma\tau(1 + z)]} \tilde{T}^{(m)}(z) [2 - 2z + \sigma\tau(1 + z)] \nu(dz) / (4\tau^2) \\ &= \operatorname{Re} \int_{\partial\mathbb{D}} |z|^{-2} |\tilde{u}(z)|^2 \tilde{T}^{(m)}(z) |[2 - 2z + \sigma\tau(1 + z)]|^2 \nu(dz) / (4\tau^2) \leq 0, \end{aligned} \quad (4.34)$$

which provides (4.13).

In the same way, let us assume that we have: $(\tau D_\tau + \sigma E)u^k = 0$, for $k \geq n+1$. Thus, one has

$$u^{k+1} = \frac{1 - \sigma/2}{1 + \sigma/2} u^k,$$

for $k \geq n+1$. As a consequence, we deduce that $\{u^k\}$ satisfies $Eu^k = 0$, for $k \geq n+1$, and

$$\begin{aligned} \operatorname{Re} \sum_{k=0}^n \overline{(\tau D_\tau + \sigma E)u^k} \left((D_\tau + \sigma E)T^{(m)} * u \right)^k &= \operatorname{Re} \sum_{k=0}^\infty \overline{(\tau D_\tau + \sigma E)u^k} \left((D_\tau + \sigma E)T^{(m)} * u \right)^k \\ &= \operatorname{Re} \left((\tau D_\tau + \sigma E)u, (D_\tau + \sigma E)T^{(m)} * u \right)_{\ell^2(\mathbb{C})} \\ &= \operatorname{Re} \int_{\partial\mathbb{D}} |\tilde{u}(z)|^2 \tilde{T}^{(m)}(z) \overline{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} [(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2] \nu(dz) \\ &= \operatorname{Re} \int_{\partial\mathbb{D}} |(z^{-1} - 1) + \sigma(z^{-1} + 1)/2|^2 \left(\frac{(z^{-1} - 1)/\tau + \sigma(z^{-1} + 1)/2}{(z^{-1} - 1) + \sigma(z^{-1} + 1)/2} \right) \tilde{T}^{(m)}(z) |\tilde{u}(z)|^2 \nu(dz) \leq 0, \end{aligned} \quad (4.35)$$

by using (4.33), ending hence the proof of (4.14). \square

5. ERROR ESTIMATES OF THE SCHEME

Let us define the two error vectors

$$\begin{aligned}\varepsilon_1^n &= ((v_1)_1^n - v_1(x_{1/2}, t_n), \dots, (v_1)_{M+1}^n - v_1(x_{M+1/2}, t_n)), \\ \varepsilon_2^n &= ((v_2)_0^n - v_2(x_0, t_n), \dots, (v_2)_{M+1}^n - v_2(x_{M+1}, t_n)),\end{aligned}\tag{5.1}$$

where (v_1, v_2) is solution to (2.12) and (v_1^n, v_2^n) is the solution to the discrete system (4.3)–(4.6). We first give the main result concerning the error estimate.

Theorem 5.1. *Let us assume that the solutions $u_1(x, t)$ and $u_2(x, t)$ to system (2.1), or equivalently, the solutions $v_1(x, t)$ and $v_2(x, t)$ of (2.2), are sufficiently smooth. Let us suppose that $\sigma \geq \frac{1}{\sqrt{2\kappa}}$, τ is small enough and that m satisfies (4.10), with μ and δ given by (4.11). Then, we have the error bound*

$$\max_{1 \leq n \leq [T/\tau]} \left(\|\mathcal{P}\varepsilon_2^n\|_h^2 + \|\nabla_h \varepsilon_2^n\|_h^2 + \|\Delta_h \varepsilon_2^n\|_h^2 + \|\varepsilon_1^n\|_h^2 + \|\nabla_h \varepsilon_1^n\|_h^2 \right) \leq \mathcal{O}(\tau^2 + h^2).$$

It is straightforward to check that the error vector $(\varepsilon_1^n, \varepsilon_2^n)$ satisfies

$$(D_\tau + \sigma E)\varepsilon_1^n + \nabla_h E \varepsilon_2^n = f_1^n, \tag{5.2}$$

$$(D_\tau + \sigma E)\mathcal{P}\varepsilon_2^n + \nabla_h E \varepsilon_1^n = \kappa \Delta_h (D_\tau + \sigma E)\varepsilon_2^n + f_2^n, \quad \forall n \geq 0, \tag{5.3}$$

$$(\mathcal{T}^m * \gamma^\pm \varepsilon_2)^n - \partial_\nu^\pm \varepsilon_2^n = g_\pm^n, \quad \forall n \geq 0, \tag{5.4}$$

$$\varepsilon_1^0 = (0, \dots, 0) \quad \varepsilon_2^0 = (0, \dots, 0), \tag{5.5}$$

where $f_1^n = ((f_1)_1^n, \dots, (f_1)_{M+1}^n)$, $f_2^n = ((f_2)_1^n, \dots, (f_2)_M^n)$ and g_\pm^n are the interior/boundary truncation error vectors/numbers according to the time and space discretizations, *i.e.*

$$\begin{aligned}(f_1)_j^n &= \left[(D_\tau + \sigma E)v_1(x_{j-1/2}, t_n) - \left(\partial_t v_1(x_{j-1/2}, t_{n+\frac{1}{2}}) + \sigma v_1(x_{j-1/2}, t_{n+\frac{1}{2}}) \right) \right] \\ &\quad + \left[E(v_2(x_j, t_n) - v_2(x_{j-1}, t_n))/h - \partial_x v_2(x_{j-1/2}, t_{n+\frac{1}{2}}) \right], \quad 1 \leq j \leq M+1, \end{aligned} \tag{5.6}$$

$$\begin{aligned}(f_2)_j^n &= \left[(D_\tau + \sigma E)v_2(x_j, t_n) - \left(\partial_t v_2(x_j, t_{n+\frac{1}{2}}) + \sigma v_2(x_j, t_{n+\frac{1}{2}}) \right) \right] \\ &\quad + \left[E(v_1(x_{j+1/2}, t_n) - v_1(x_{j-1/2}, t_n))/h - \partial_x v_1(x_j, t_{n+\frac{1}{2}}) \right] \\ &\quad - \kappa \left[E(v_2(x_{j-1}, t_n) - 2v_2(x_j, t_n) + v_2(x_{j+1}, t_n))/h^2 - \partial_x^2 v_2(x_j, t_{n+\frac{1}{2}}) \right], \quad 1 \leq j \leq M, \end{aligned} \tag{5.7}$$

$$\begin{aligned}g_\pm^n &= \left(\mathcal{T}^{(m)} - \mathcal{T} \right) * \gamma^\pm v_2(t_n) + \left[(\mathcal{T} * \gamma^\pm v_2)(t_n) - (\mathcal{T} * \gamma^\pm v_2)(t_n) \right] \\ &\quad + \left[(\mathcal{T} * \gamma^\pm v_2)(t_n) - (\mathcal{T} * v_2)(x_\pm, t_n) \right] + \left[-\partial_\nu^\pm v_2(t_n) + \partial_\nu v_2(x_\pm, t_n) \right], \end{aligned} \tag{5.8}$$

with $v_2(t_n) = (v_2(x_1, t_n), \dots, v_2(x_M, t_n))$.

The proof of Theorem 5.1 is presented in the next two subsections as a consequence of Propositions 5.2 and 5.3.

5.1. Estimate for truncation errors

Let us first prove the estimate for the truncation errors of the boundary and interior schemes.

Proposition 5.2. *Under the conditions of Theorem 5.1, we have the following error estimate*

$$\|\nabla_h f_1^n\|_h + \|f_1^n\|_h + \|f_2^n\|_h + |g_\pm^n| + |D_\tau g_\pm^n| + |D_\tau^2 g_\pm^n| \leq C(\tau^2 + h^2), \tag{5.9}$$

with $D_\tau^2 g_\pm^n = \frac{g_\pm^{n+1} - 2g_\pm^n + g_\pm^{n-1}}{\tau^2}$, C being a strictly positive constant.

Proof of Proposition 5.2. The proof is separated into three estimates which are summed up at the end to prove the result.

Estimate of $|g_{\pm}^n|$. Here, we prove

$$g_{\pm}^n = \mathcal{O}(\tau^2 + h^2). \quad (5.10)$$

We divide the proof into two steps.

Step 1. Let us recall that $(\mathcal{T} * v_2)$ is defined by (3.11). We first derive the following bound

$$|(\mathcal{T} * v_2)(x_{\pm}, t) - (\mathbb{T} * v_2)(x_{\pm}, t)| \leq \mathcal{O}(\tau^2). \quad (5.11)$$

By using Taylor's expansion, it is straightforward to verify that

$$\left| \tilde{\mathcal{T}}(e^{-i\tau\xi}) - \left(-\sqrt{S(i\xi)} \right) \right| \leq C\tau^2|\xi|^3. \quad (5.12)$$

Since $v_2(x, 0) = v_1(x, 0) = 0$ for $x \in [x_+, \infty)$, thus $\partial_x v_1(x_+, 0)$, $\partial_x v_2(x_+, 0)$ and $\partial_{xx} v_2(x_+, 0)$ are also equal to zero. Then (2.3b) leads to the following equality

$$\partial_t v_2(x, 0) = \kappa \partial_{xx}(\partial_t v_2(x, 0)), \quad \forall x \in [x_+, +\infty),$$

implying that $\partial_t v_2(x, 0) = C e^{-\frac{x}{\sqrt{\kappa}}}$ for $x \in [x_+, +\infty)$ and

$$\partial_t \partial_x v_2(x_+, 0) = -\frac{C}{\sqrt{\kappa}} e^{-\frac{x_+}{\sqrt{\kappa}}}. \quad (5.13)$$

From (2.11), we obtain

$$\partial_x \partial_t v_2(x_+, 0) = \mathbb{T}(0) v_2(x_+, 0) + (\partial_t \mathbb{T} * v_2(x_+, \cdot))(0) = 0,$$

and then $C = 0$ in (5.13). Thus, we deduce: $\partial_t v_2(x_+, 0) = 0$. From (2.3a) and (2.3c), we also have $\partial_t v_1(x_+, 0) = 0$. Repeating the same procedure, it is easy to conclude that $v_2(x_{\pm}, t)$ and its time derivatives are zero at $t = 0$. Consequently, by extending $v_2(x_{\pm}, t)$ to zero on $t \in (-\infty, 0]$, we obtain a sufficiently smooth function $v_2(x_{\pm}, t)$ defined for $t \in \mathbb{R}$. We set

$$(\mathcal{T} * v_2)(x_{\pm}, t) := \sum_{j=0}^{\infty} \mathcal{T}_j v_2(x_{\pm}, t - j\tau), \quad \forall t \in \mathbb{R}, \quad (5.14)$$

which is consistent with definition (3.11) at $t = t_n$. The Fourier transform in time of (5.14) is

$$\begin{aligned} \mathcal{F}_t[(\mathcal{T} * v_2)(x_{\pm}, t)](\xi) &= \int_{\mathbb{R}} (\mathcal{T} * v_2)(x_{\pm}, t) e^{-it\xi} dt = \sum_{j=0}^{\infty} \int_{\mathbb{R}} \mathcal{T}_j v_2(x_{\pm}, t - j\tau) e^{-it\xi} dt \\ &= \tilde{\mathcal{T}}(e^{-i\tau\xi}) \mathcal{F}_t v_2(x_{\pm}, \xi) = -\sqrt{S(i\xi)} \mathcal{F}_t v_2(x_{\pm}, \xi) + \left(\tilde{\mathcal{T}}(e^{-i\tau\xi}) + \sqrt{S(i\xi)} \right) \mathcal{F}_t v_2(x_{\pm}, \xi) \\ &= \mathcal{F}_t[(\mathbb{T} * v_2)(x_{\pm}, t)](\xi) + \left(\tilde{\mathcal{T}}(e^{-i\tau\xi}) + \sqrt{S(i\xi)} \right) \mathcal{F}_t v_2(x_{\pm}, \xi). \end{aligned}$$

Let us recall that

$$\tilde{\mathcal{T}}(e^{-i\tau\xi}) = -\sqrt{S\left(\frac{2(1 - e^{-i\tau\xi})}{\tau(1 + e^{-i\tau\xi})}\right)}$$

with

$$S(s) = \frac{(\sigma + s)^2}{1 + \kappa(\sigma + s)^2}.$$

Therefore, we have

$$\begin{aligned}
\left| \tilde{\mathcal{T}}(e^{-i\tau\xi}) + \sqrt{S(i\xi)} \right| &= \left| \sqrt{S(i\xi)} - \sqrt{S\left(\frac{2(1-e^{-i\tau\xi})}{\tau(1+e^{-i\tau\xi})}\right)} \right| = \left| \sqrt{S(i\xi)} - \sqrt{S\left(i\frac{2\tan(\tau\xi/2)}{\tau}\right)} \right| \\
&= \left| \int_{i\xi}^{i\frac{\tan(\tau\xi/2)}{\tau/2}} \frac{d}{ds} \left(\sqrt{S(s)} \right) ds \right| = \left| \int_{i\xi}^{i\frac{\tan(\tau\xi/2)}{\tau/2}} \frac{1}{(1+\kappa(\sigma+s))^{3/2}} ds \right| \\
&\leq C \left| i\frac{\tan(\tau\xi/2)}{\tau/2} - i\xi \right| \leq \left| \int_0^\xi \left(\frac{1}{1+\frac{\tau^2\xi_1^2}{4}} - 1 \right) d\xi_1 \right| = \left| \int_0^\xi \frac{\frac{\tau^2\xi_1^2}{4}}{1+\frac{\tau^2\xi_1^2}{4}} d\xi_1 \right| \\
&\leq \frac{\tau^2}{4} \int_0^{|\xi|} \xi_1^2 d\xi_1 \leq C\tau^2|\xi|^3.
\end{aligned}$$

In addition, for $|\xi| \leq 1$, we obtain

$$|\xi|^3 \leq \frac{|\xi|^3 + |\xi|^4}{1 + |\xi|} \leq \frac{1 + |\xi|^4}{1 + |\xi|}.$$

For $|\xi| \geq 1$, we have

$$|\xi|^3 \leq \frac{|\xi|^3 + |\xi|^4}{1 + |\xi|} \leq \frac{2|\xi|^4}{1 + |\xi|} \leq 2\frac{1 + |\xi|^4}{1 + |\xi|}.$$

Thus, we deduce

$$|\xi|^3 \leq C\frac{1 + |\xi|^4}{1 + |\xi|}.$$

By the above two estimates, we have

$$\begin{aligned}
|\mathcal{T} * v_2(x_\pm, t) - (\mathbf{T} * v_2)(x_\pm, t)| &= \left| \mathcal{F}_\xi^{-1} \left[\left(\tilde{\mathcal{T}}(e^{-i\tau\xi}) + \sqrt{S(i\xi)} \right) \mathcal{F}_t v_2(x_\pm, \xi) \right] (t) \right| \\
&\leq \int_{\mathbb{R}} \left| \tilde{\mathcal{T}}(e^{-i\tau\xi}) + \sqrt{S(i\xi)} \right| |\mathcal{F}_t v_2(x_\pm, \xi)| d\xi \leq C\tau^2 \int_{\mathbb{R}} |\xi|^3 |\mathcal{F}_t v_2(x_\pm, \xi)| d\xi \\
&\leq C\tau^2 \int_{\mathbb{R}} \frac{1}{1 + |\xi|} (1 + |\xi|^4) |\mathcal{F}_t v_2(x_\pm, \xi)| d\xi \\
&\leq C\tau^2 \left(\int_{\mathbb{R}} \frac{1}{(1 + |\xi|)^2} d\xi \right)^{\frac{1}{2}} \left(\int_{\mathbb{R}} (1 + |\xi|^4)^2 |\mathcal{F}_t v_2(x_\pm, \xi)|^2 d\xi \right)^{\frac{1}{2}} \\
&= C\tau^2 \left(\int_0^\infty (|v_2(x_\pm, t)|^2 + |\partial_t^4 v_2(x_\pm, t)|^2) dt \right)^{\frac{1}{2}}.
\end{aligned} \tag{5.15}$$

We finally obtain (5.11).

Step 2. The inequality (4.12) of Proposition 4.1 implies that: $\left| \tilde{\mathcal{T}}^{(m)}(z) - \tilde{\mathcal{T}}(z) \right| \leq C\tau^3$ for $|z| = 1$. Since

$$\mathcal{T}_j^{(m)} = \int_{\partial\mathbb{D}} \tilde{\mathcal{T}}^{(m)}(z) z^{-j} \nu(dz) \quad \text{and} \quad \mathcal{T}_j = \int_{\partial\mathbb{D}} \tilde{\mathcal{T}}(z) z^{-j} \nu(dz),$$

this shows that

$$\left| \mathcal{T}_j^{(m)} - \mathcal{T}_j \right| \leq \int_{\partial\mathbb{D}} \left| \tilde{\mathcal{T}}^{(m)}(z) - \tilde{\mathcal{T}}(z) \right| \nu(dz) \leq C\tau^3.$$

Thus, the following inequalities hold

$$\left| \sum_{j=0}^n \mathcal{T}_j^{(m)} v_2^{n-j} - \sum_{j=0}^n \mathcal{T}_j v_2^{n-j} \right| \leq \sum_{j=0}^n \left| \mathcal{T}_j^{(m)} - \mathcal{T}_j \right| |v_2^{n-j}| \leq \sum_{j=0}^n C\tau^3 \leq C\tau^2,$$

which leads to

$$\left(\mathcal{T}^{(m)} - \mathcal{T}\right) * \gamma^\pm v_2(t_n) = \mathcal{O}(\tau^2). \quad (5.16)$$

Besides, equation (5.11) yields

$$(\mathcal{T} * \gamma^\pm v_2)(t_n) - (\mathcal{T} * \gamma^\pm v_2)(t_n) = \mathcal{O}(\tau^2). \quad (5.17)$$

Since $\gamma^+ v_2 = ((v_2)_{M+1} + (v_2)_M)/2$ and $x_+ = (x_{M+1} + x_M)/2 = x_{M+\frac{1}{2}}$, it follows that

$$|(\mathcal{T} * \gamma^\pm v_2)(t_n) - (\mathcal{T} * v_2)(x_\pm, t_n)| = \mathcal{O}(h^2),$$

and

$$\begin{aligned} \partial_\nu^+ v_2(t_n) - \partial_\nu v_2(x_+, t_n) &= \frac{v_2(x_{M+1}, t_n) - v_2(x_M, t_n)}{h} - \partial_x v_2\left(x_{M+\frac{1}{2}}, t_n\right) = \mathcal{O}(h^2), \\ \partial_\nu^- v_2(t_n) - \partial_\nu v_2(x_-, t_n) &= -\frac{v_2(x_1, t_n) - v_2(x_0, t_n)}{h} + \partial_x v_2\left(x_{\frac{1}{2}}, t_n\right) = \mathcal{O}(h^2). \end{aligned} \quad (5.18)$$

Substituting (5.16)–(5.18) into (5.8) leads to (5.10).

Estimate of $\|\nabla_h f_1^n\|_h^2 + \|f_1^n\|_h + \|f_2^n\|_h$. We now prove that

$$\|\nabla_h f_1^n\|_h + \|f_1^n\|_h + \|f_2^n\|_h \leq \mathcal{O}(\tau^2 + h^2). \quad (5.19)$$

Recalling (5.6), we estimate the three terms in the expression of $(f_1)_j^n$ separately. Firstly, we have

$$\begin{aligned} &(D_\tau + \sigma E)v_1(x_{j-1/2}, t_n) - \left(\partial_t v_1\left(x_{j-1/2}, t_{n+\frac{1}{2}}\right) + \sigma v_1\left(x_{j-1/2}, t_{n+\frac{1}{2}}\right)\right) \\ &= \left(\frac{v_1(x_{j-1/2}, t_{n+1}) - v_1(x_{j-1/2}, t_n)}{\tau} - \partial_t v_1\left(x_{j-1/2}, t_{n+\frac{1}{2}}\right)\right) \\ &\quad + \sigma \left(\frac{v_1(x_{j-1/2}, t_n) + v_1(x_{j-1/2}, t_{n+1})}{2} - v_1\left(x_{j-1/2}, t_{n+\frac{1}{2}}\right)\right) = \mathcal{O}(\tau^2). \end{aligned} \quad (5.20)$$

Secondly, the following estimate holds

$$\left[E(v_2(x_j, t_n) - v_2(x_{j-1}, t_n))/h - \partial_x v_2\left(x_{j-1/2}, t_{n+\frac{1}{2}}\right)\right] = \mathcal{O}(\tau^2 + h^2). \quad (5.21)$$

Thus, from (5.6), (5.20) and (5.21), we deduce

$$(f_1)_j^n = \mathcal{O}(\tau^2 + h^2), \quad 1 \leq j \leq M+1. \quad (5.22)$$

Similarly, recalling (5.7),

$$\frac{Ev_2(x_{j-1}, t_n) - 2Ev_2(x_j, t_n) + Ev_2(x_{j+1}, t_n)}{h^2} - \partial_x^2 v_2\left(x_j, t_{n+\frac{1}{2}}\right) = \mathcal{O}(\tau^2 + h^2), \quad (5.23)$$

we thus obtain

$$(f_2)_j^n = \mathcal{O}(\tau^2 + h^2), \quad 1 \leq j \leq M. \quad (5.24)$$

Finally, from (5.22) and (5.24) one gets

$$\|f_1^n\|_h = \mathcal{O}(\tau^2 + h^2), \quad \|f_2^n\|_h = \mathcal{O}(\tau^2 + h^2). \quad (5.25)$$

In a similar way (using a Taylor's expansion), one can prove that

$$\|\nabla_h f_1^n\|_h = \mathcal{O}(\tau^2 + h^2), \quad (5.26)$$

which leads to (5.19).

Estimate of $\|D_\tau g_\pm^n\|$. Since we have

$$\begin{aligned} D_\tau g_\pm^n &= \left(\mathcal{T}^{(m)} - \mathcal{T} \right) * \gamma^\pm D_\tau v_2(t_n) + [\mathcal{T} * \gamma^\pm D_\tau v_2(t_n) - (\mathcal{T} * \gamma^\pm D_\tau v_2)(t_n)] \\ &\quad + [(\mathcal{T} * \gamma^\pm D_\tau v_2)(t_n) - (\mathcal{T} * D_\tau v_2)(x_\pm, t_n)] + [-\partial_\nu^\pm D_\tau v_2(t_n) + \partial_\nu D_\tau v_2(x_\pm, t_n)], \end{aligned} \quad (5.27)$$

it follows that (5.27) can be estimated as (5.8) (replacing $v_2(x, t_n)$ by $D_\tau v_2(x, t_n)$), which provides

$$D_\tau g_\pm^n = \mathcal{O}(\tau^2 + h^2). \quad (5.28)$$

In a similar way, we can prove that

$$D_\tau^2 g_\pm^n = \mathcal{O}(\tau^2 + h^2). \quad (5.29)$$

Combing (5.10), (5.19), (5.28) and (5.29), we finally get (5.9). \square

5.2. Error estimates

Let us state the error estimate for system (4.3)–(4.6). Theorem 5.1 is then a consequence of Propositions 5.2 and 5.3.

Proposition 5.3. *If $\sigma \geq \frac{1}{\sqrt{2\kappa}}$, and the order m of the Padé approximation fulfills (4.10), the solution of (4.3)–(4.6) satisfies the following stability estimate:*

$$\begin{aligned} &\max_{1 \leq n \leq [T/\tau]} \left(\|\mathcal{P}\varepsilon_2^n\|_h^2 + \|\nabla_h \varepsilon_2^n\|_h^2 + \|\Delta_h \varepsilon_2^n\|_h^2 + \|\varepsilon_1^n\|_h^2 + \|\nabla_h \varepsilon_1^n\|_h^2 \right) \\ &\leq C_T \left[\max_{0 \leq k \leq n-1} \left(\|\nabla_h f_1^k\|_h^2 + \|f_2^k\|_h^2 + \|f_1^k\|_h^2 + |D_\tau g_\pm^k|^2 \right) + \max_{0 \leq k \leq n} |g_\pm^k|^2 + \max_{1 \leq k \leq n-1} |D_\tau^2 g_\pm^k|^2 \right], \end{aligned}$$

where C_T is a constant depending on T .

Proof of Proposition 5.3. Due to

$$D_\tau(\varepsilon_1)_m^n \cdot E(\varepsilon_1)_m^n = \frac{(\varepsilon_1)_m^{n+1} - (\varepsilon_1)_m^n}{\tau} \cdot \frac{(\varepsilon_1)_m^{n+1} + (\varepsilon_1)_m^n}{2} = \frac{|(\varepsilon_1)_m^{n+1}|^2 - |(\varepsilon_1)_m^n|^2}{2\tau},$$

taking the real part of the inner product between the left hand side of (5.2) and $E\varepsilon_1^n$ yields

$$\begin{aligned} \frac{1}{2} D_\tau (\|\varepsilon_1^n\|_h^2) + \sigma \|E\varepsilon_1^n\|_h^2 &= -\operatorname{Re} \{E\varepsilon_1^n, E\nabla_h \varepsilon_2^n\}_h + \operatorname{Re} (E\varepsilon_1^n, f_1^n)_h \\ &\leq \frac{\sigma}{4} \|E\varepsilon_1^n\|_h^2 + \frac{1}{\sigma} |E\nabla_h \varepsilon_2^n|_h^2 + \frac{\sigma}{4} \|\varepsilon_1^n\|_h^2 + \frac{1}{\sigma} \|f_1^n\|_h^2, \end{aligned} \quad (5.30)$$

using the inequality

$$|ab| = \left| \sqrt{\frac{\sigma}{2}} a \right| \left| \sqrt{\frac{2}{\sigma}} b \right| \leq \frac{1}{2} \left(\frac{\sigma}{2} a^2 + \frac{2}{\sigma} b^2 \right) = \frac{\sigma}{4} a^2 + \frac{1}{\sigma} b^2.$$

Summing up over n , we obtain

$$\|\varepsilon_1^n\|_h^2 \leq \mathcal{O}(\tau) \left(\sum_{k=0}^n |\nabla_h \varepsilon_2^k|_h^2 + \sum_{k=0}^{n-1} \|f_1^k\|_h^2 \right). \quad (5.31)$$

From (5.2), it is easy to see that: $(D_\tau + \sigma E)\nabla_h \varepsilon_1^n + \Delta_h E \varepsilon_2^n = \nabla_h f_1^n$. Next computing the inner product of the previous expression with $E \nabla_h \varepsilon_1^n$ and taking the real part, we deduce

$$\begin{aligned} \frac{1}{2} D_\tau \left(\|\nabla_h \varepsilon_1^n\|_h^2 \right) + \sigma \|E \nabla_h \varepsilon_1^n\|_h^2 &= -\operatorname{Re} (E \nabla_h \varepsilon_1^n, E \Delta_h \varepsilon_2^n)_h + \operatorname{Re} (E \nabla_h \varepsilon_1^n, \nabla_h f_1^n)_h \\ &\leq \frac{\sigma}{4} \|E \nabla_h \varepsilon_1^n\|_h^2 + \frac{1}{\sigma} \|E \Delta_h \varepsilon_2^n\|_h^2 + \frac{\sigma}{4} \|E \nabla_h \varepsilon_1^n\|_h^2 + \frac{1}{\sigma} \|\nabla_h f_1^n\|_h^2, \end{aligned}$$

which leads to

$$\|\nabla_h \varepsilon_1^n\|_h^2 \leq \mathcal{O}(\tau) \left(\sum_{k=0}^n \|\Delta_h \varepsilon_2^k\|_h^2 + \sum_{k=0}^{n-1} \|\nabla_h f_1^k\|_h^2 \right). \quad (5.32)$$

By taking the inner product of (5.3) with $(\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n$ and then the real part, one gets

$$\operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, (D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n + \nabla_h E \varepsilon_1^n)_h = \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, \kappa \Delta_h^n (D_\tau + \sigma E) \varepsilon_2 + f_2^n)_h. \quad (5.33)$$

The left hand side of (5.33) can be written as

$$\tau \|D_\tau \mathcal{P} \varepsilon_2^n\|_h^2 + \frac{\sigma(1+\tau)}{2} D_\tau (\|\mathcal{P} \varepsilon_2^n\|_h^2) + \sigma^2 \|E \mathcal{P} \varepsilon_2^n\|_h^2 + \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, \nabla_h E \varepsilon_1^n)_h. \quad (5.34)$$

By applying the discrete Green's formula (3.13), the boundary conditions (4.5), and (4.13), the right hand side of (5.33) can be rewritten following

$$\begin{aligned} & -\kappa \langle \nabla_h (\tau D_\tau + \sigma E) \varepsilon_2^n, \nabla_h (D_\tau + \sigma E) \varepsilon_2^n \rangle_h + \kappa \operatorname{Re} \left(\overline{\gamma^\pm (\tau D_\tau + \sigma E) \varepsilon_2^n} \partial_\nu^\pm (D_\tau + \sigma E) \varepsilon_2^n \right) \\ & + \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, f_2^n)_h \\ & = -\kappa \langle \nabla_h (\tau D_\tau + \sigma E) \varepsilon_2^n, \nabla_h (D_\tau + \sigma E) \varepsilon_2^n \rangle_h + \kappa \operatorname{Re} \left(\overline{(\tau D_\tau + \sigma E) \gamma^\pm \varepsilon_2^n} (D_\tau + \sigma E) \partial_\nu^\pm \varepsilon_2^n \right) \\ & + \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, f_2^n)_h \\ & = -\kappa \langle \nabla_h (\tau D_\tau + \sigma E) \varepsilon_2^n, \nabla_h (D_\tau + \sigma E) \varepsilon_2^n \rangle_h + \kappa \operatorname{Re} \left(\overline{(\tau D_\tau + \sigma E) \gamma^\pm \varepsilon_2^n} (D_\tau + \sigma E) \left(T^{(m)} * \gamma^\pm \varepsilon_2 \right)^n \right) \\ & - \kappa \operatorname{Re} \left(\overline{(\tau D_\tau + \sigma E) \gamma^\pm \varepsilon_2^n} (D_\tau + \sigma E) g_\pm^n \right) + \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, f_2^n)_h \\ & \leq -\kappa \langle \nabla_h (\tau D_\tau + \sigma E) \varepsilon_2^n, \nabla_h (D_\tau + \sigma E) \varepsilon_2^n \rangle_h \\ & - \kappa \operatorname{Re} \left(\overline{(\tau D_\tau + \sigma E) \gamma^\pm \varepsilon_2^n} (D_\tau + \sigma E) g_\pm^n \right) + \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, f_2^n)_h. \end{aligned} \quad (5.35)$$

Combining (5.34) and (5.35), we have

$$\begin{aligned} & \tau \|D_\tau \mathcal{P} \varepsilon_2^n\|_h^2 + \frac{\sigma(1+\tau)}{2} D_\tau (\|\mathcal{P} \varepsilon_2^n\|_h^2) + \sigma^2 \|E \mathcal{P} \varepsilon_2^n\|_h^2 + \kappa \langle \nabla_h (\tau D_\tau + \sigma E) \varepsilon_2^n, \nabla_h (D_\tau + \sigma E) \varepsilon_2^n \rangle_h \\ & = \tau \|D_\tau \mathcal{P} \varepsilon_2^n\|_h^2 + \frac{\sigma(1+\tau)}{2} D_\tau (\|\mathcal{P} \varepsilon_2^n\|_h^2) + \sigma^2 \|E \mathcal{P} \varepsilon_2^n\|_h^2 + \kappa \tau \|D_\tau \nabla_h \varepsilon_2^n\|_h^2 \\ & + \frac{\kappa \sigma(1+\tau)}{2} D_\tau (\|\nabla_h \varepsilon_2^n\|_h^2) + \kappa \sigma^2 \|E \nabla_h \varepsilon_2^n\|_h^2 \\ & \leq -\operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, \nabla_h E \varepsilon_1^n)_h + \operatorname{Re}((\tau D_\tau + \sigma E) \mathcal{P} \varepsilon_2^n, f_2^n)_h \\ & - \kappa \operatorname{Re} \left(\overline{(\tau D_\tau + \sigma E) \gamma^\pm \varepsilon_2^n} (D_\tau + \sigma E) g_\pm^n \right), \end{aligned}$$

from which we derive

$$D_\tau (\|\mathcal{P} \varepsilon_2^n\|_h^2) + D_\tau (\|\nabla_h \varepsilon_2^n\|_h^2) \leq \mathcal{O}(1) (E(\|\nabla_h \varepsilon_1^n\|_h^2) + E(|\gamma^\pm \varepsilon_2^n|^2)) + \mathcal{O}(1) (\|f_2^n\|_h^2 + |D_\tau g_\pm^n|^2 + E(|g_\pm^n|^2)). \quad (5.36)$$

This leads to

$$\|\mathcal{P}\varepsilon_2^n\|_h^2 + |\nabla_h \varepsilon_2^n|_h^2 \leq \mathcal{O}(\tau) \sum_{k=0}^n \left(\|\nabla_h \varepsilon_1^k\|_h^2 + |g_\pm^k|^2 + |\gamma^\pm \varepsilon_2^k|^2 \right) + \mathcal{O}(\tau) \sum_{k=0}^{n-1} \left(\|f_2^k\|_h^2 + |D_\tau g_\pm^k|^2 \right).$$

Again taking the inner product of (5.3) with $(D_\tau + \sigma E)\Delta_h \mathcal{P}\varepsilon_2^n$ and then the real part, we obtain

$$\operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, (D_\tau + \sigma E)\mathcal{P}\varepsilon_2^n + \nabla_h E\varepsilon_1)_h = \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, \kappa \Delta_h (D_\tau + \sigma E)\varepsilon_2^n + f_2^n)_h. \quad (5.37)$$

The right hand side of (5.37) can be written as

$$\sigma \kappa D_\tau (\|\Delta_h \varepsilon_2^n\|_h^2) + \kappa \|D_\tau \Delta_h \varepsilon_2^n\|_h^2 + \sigma^2 \kappa \|E \Delta_h \varepsilon_2^n\|_h^2 + \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, f_2^n)_h. \quad (5.38)$$

Now by using the integration by part (3.13), the boundary conditions (4.5), and (4.14), the left term of (5.37) reads

$$\begin{aligned} & -\langle \nabla_h (D_\tau + \sigma E)\varepsilon_2^n, \nabla_h (D_\tau + \sigma E)\varepsilon_2^n \rangle_h + \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, \nabla_h E\varepsilon_1^n)_h \\ & + \operatorname{Re}\left(\overline{\partial^\pm (D_\tau + \sigma E)\varepsilon_2^n} \gamma^\pm (D_\tau + \sigma E)\varepsilon_2^n\right) \\ & = -\langle \nabla_h (D_\tau + \sigma E)\varepsilon_2^n, \nabla_h (D_\tau + \sigma E)\varepsilon_2^n \rangle_h + \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, \nabla_h E\varepsilon_1^n)_h \\ & + \operatorname{Re}\left(\overline{(D_\tau + \sigma E)\partial^\pm \varepsilon_2^n} (D_\tau + \sigma E)\gamma^\pm \varepsilon_2^n\right) \\ & = -\langle \nabla_h (D_\tau + \sigma E)\varepsilon_2^n, \nabla_h (D_\tau + \sigma E)\varepsilon_2^n \rangle_h + \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, \nabla_h E\varepsilon_1^n)_h \\ & + \operatorname{Re}\left(\overline{(D_\tau + \sigma E)(\mathcal{T}^{(m)} * \gamma^\pm \varepsilon_2)^n} (D_\tau + \sigma E)\gamma^\pm \varepsilon_2^n\right) - \operatorname{Re}\left(\overline{(D_\tau + \sigma E)g_\pm^n} (D_\tau + \sigma E)\gamma^\pm \varepsilon_2^n\right) \\ & \leq -\langle \nabla_h (D_\tau + \sigma E)\varepsilon_2^n, \nabla_h (D_\tau + \sigma E)\varepsilon_2^n \rangle_h + \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, \nabla_h E\varepsilon_1^n)_h \\ & - \operatorname{Re}\left(\overline{(D_\tau + \sigma E)g_\pm^n} (D_\tau + \sigma E)\gamma^\pm \varepsilon_2^n\right). \end{aligned} \quad (5.39)$$

Combining (5.38) and (5.39), we deduce

$$\begin{aligned} & \sigma \kappa D_\tau (\|\Delta_h \varepsilon_2^n\|_h^2) + \kappa \|D_\tau \Delta_h \varepsilon_2^n\|_h^2 + \sigma^2 \kappa \|E \Delta_h \varepsilon_2^n\|_h^2 + \langle \nabla_h^n (D_\tau + \sigma E)\varepsilon_2^n, \nabla_h^n (D_\tau + \sigma E)\varepsilon_2^n \rangle_h \\ & \leq \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, \nabla_h^n E\varepsilon_1)_h - \operatorname{Re}\left(\overline{(D_\tau + \sigma E)g_\pm^n} (D_\tau + \sigma E)\gamma^\pm \varepsilon_2^n\right) - \operatorname{Re}((D_\tau + \sigma E)\Delta_h \varepsilon_2^n, f_2^n)_h, \end{aligned}$$

from which we derive that

$$\begin{aligned} D_\tau \|\Delta_h \varepsilon_2^n\|_h^2 + D_\tau |\nabla_h \varepsilon_2^n|_h^2 & \leq \mathcal{O}(1) (E(\|\nabla_h \varepsilon_1^n\|_h^2) + E(|\gamma^\pm \varepsilon_2^n|^2)) \\ & + \mathcal{O}(1) (\|f_2^n\|_h^2 + |D_\tau g_\pm^n|^2 + E(|g_\pm^n|^2)) - \operatorname{Re}\left(\overline{(D_\tau + \sigma E)g_\pm^n} D_\tau \gamma^\pm \varepsilon_2^n\right). \end{aligned}$$

Summing up the index n and using a summation by parts in time for $\sum_{k=0}^{n-1} \operatorname{Re}\left(D_\tau \gamma^\pm \varepsilon_2^k \cdot \overline{(D_\tau + \sigma E)g_\pm^k}\right)$, we obtain

$$\begin{aligned} \tau^{-1} \left[\|\Delta_h \varepsilon_2^n\|_h^2 + |\nabla_h \varepsilon_2^n|_h^2 \right] & \leq \mathcal{O}(1) \tau^{-1} \left(|\gamma^\pm \varepsilon_2^n|^2 + |g_\pm^{n-1}|^2 + |g_\pm^n|^2 + |D_\tau g_\pm^{n-1}|^2 \right) \\ & + \mathcal{O}(1) \sum_{k=1}^n \left(\|\nabla_h \varepsilon_1^k\|_h^2 + |\gamma^\pm \varepsilon_2^k|^2 + |g_\pm^k|^2 \right) + \mathcal{O}(1) \sum_{k=0}^{n-1} \left(\|f_2^k\|_h^2 + |D_\tau g_\pm^k|^2 \right) \\ & + \mathcal{O}(1) \sum_{k=1}^{n-1} |D_\tau^2 g_\pm^k|^2. \end{aligned} \quad (5.40)$$

By the discrete Sobolev imbedding theorem

$$|\gamma^\pm \varepsilon_2^n|^2 \leq \mathcal{O}(\epsilon_1^{-1}) \|\mathcal{P}\varepsilon_2^n\|_h^2 + \epsilon_1 \|\nabla_h \varepsilon_2^n\|_h^2,$$

and choosing ϵ_1 small enough from (5.40), we have

$$\begin{aligned} \|\Delta_h \varepsilon_2^n\|_h^2 + \|\nabla_h \varepsilon_2^n\|_h^2 &\leq \mathcal{O}(1) \left(\|\mathcal{P}\varepsilon_2^n\|_h^2 + |g_\pm^{n-1}|^2 + |g_\pm^n|^2 + |D_\tau g_\pm^{n-1}|^2 \right) \\ &\quad + \mathcal{O}(\tau) \sum_{k=1}^n \left(\|\nabla_h \varepsilon_1^k\|_h^2 + |\gamma^\pm \varepsilon_2^k|^2 + |g_\pm^k|^2 \right) + \mathcal{O}(\tau) \sum_{k=0}^{n-1} \left(\|f_2^k\|_h^2 + |D_\tau g_\pm^k|^2 \right) \\ &\quad + \mathcal{O}(\tau) \sum_{k=1}^{n-1} |D_\tau^2 g_\pm^k|^2. \end{aligned} \quad (5.41)$$

Combining (5.31), (5.32), (5.37) and (5.41) together yields

$$\begin{aligned} &\|\Delta_h \varepsilon_2^n\|_h^2 + \|\nabla_h \varepsilon_2^n\|_h^2 + \|\mathcal{P}\varepsilon_2^n\|_h^2 + \|\nabla_h \varepsilon_1^n\|_h^2 + \|\varepsilon_1^n\|_h^2 \\ &\leq \mathcal{O}(1) \left(|g_\pm^{n-1}|^2 + |g_\pm^n|^2 + |D_\tau g_\pm^{n-1}|^2 \right) \\ &\quad + \mathcal{O}(\tau) \sum_{k=1}^n \left(\|\nabla_h \varepsilon_1^k\|_h^2 + \|\Delta_h \varepsilon_2^k\|_h^2 + |\gamma^\pm \varepsilon_2^k|^2 + |g_\pm^k|^2 \right) \\ &\quad + \mathcal{O}(\tau) \sum_{k=0}^{n-1} \left(\|f_2^k\|_h^2 + \|\nabla_h f_1^k\|_h^2 + \|f_1^k\|_h^2 + |D_\tau g_\pm^k|^2 \right) + \mathcal{O}(\tau) \sum_{k=1}^{n-1} |D_\tau^2 g_\pm^k|^2 \\ &\leq \mathcal{O}(\tau) \sum_{k=1}^n \left(\|\nabla_h \varepsilon_1^k\|_h^2 + \|\Delta_h \varepsilon_2^k\|_h^2 + \|\nabla_h \varepsilon_2^k\|_h^2 + \|\mathcal{P}\varepsilon_2^k\|_h^2 \right) \\ &\quad + \mathcal{O}(\tau) \sum_{k=0}^{n-1} \left(\|f_2^k\|_h^2 + \|\nabla_h f_1^k\|_h^2 + \|f_1^k\|_h^2 + |D_\tau g_\pm^k|^2 \right) \\ &\quad + \mathcal{O}(\tau) \sum_{k=0}^n |g_\pm^k|^2 + \mathcal{O}(\tau) \sum_{k=1}^{n-1} |D_\tau^2 g_\pm^k|^2 + \mathcal{O}(1) \left(|g_\pm^{n-1}|^2 + |g_\pm^n|^2 + |D_\tau g_\pm^{n-1}|^2 \right). \end{aligned} \quad (5.42)$$

Applying the discrete Gronwall's inequality [25], we derive (5.3). The proof of Proposition 5.3 is complete. \square

We remark that, From Proposition 4.1, for very small values of κ , we have: $m \sim \tau \kappa^{-1/2}$. From Theorem 5.1 and since $e^{\sigma t} v_i = u_i$, in order to achieve the accuracy of $\mathcal{O}(\tau^2 + h^2)$, we have $m \sim (|\log(\tau)| + |\log(\sigma)| + |\log(C)|) \kappa^{-1/2}$ and the time steps fulfill $N \sim C^{1/2} e^{\sigma T/2} T/\tau$, where C depends on κ . Thus, the resulting total computational complexity is $\mathcal{O}(mn) \sim (|\log(\tau)| + |\log(\kappa)| + |\log(C)|) C^{1/2} e^{\kappa^{-1/2} T/2} T/(\kappa^{1/2} \tau)$. Here, the constant C is not given by an explicit formula. From Proposition 4.1, we see that we can design a fast and stable approximation of the boundary condition due to the damping coefficient σ . Without any damping factor, then the approximate boundary condition can lead to some instabilities, similarly to [18].

6. NUMERICAL EXAMPLES

We now provide two illustrative numerical examples to validate the theoretical results derived in the preceding sections. In the calculations, we always take $\sigma = 1/\sqrt{2\kappa}$ and adapt the number of Padé expansion terms (see Thm. 5.1) following the rule

$$m = \frac{\ln \epsilon}{2 \ln(1 - \delta)}, \quad \epsilon = \frac{\mu \sqrt{\kappa} \tau^3}{8},$$

with μ and δ given by (4.11). Therefore for N fixed (with $N\tau = T$), the total computational cost to efficiently evaluate the convolution is $\mathcal{O}(mN) = \mathcal{O}(N \log(N))$ [21] by (4.9).

Example 6.1. To demonstrate the performance of our numerical scheme, we first consider a Gaussian initial distribution for the free-surface elevation and zero distribution for velocity, *i.e.*

$$u_1(x, 0) = \exp(-400(x - 0.5)), \quad (6.1)$$

and $u_2(x, 0) = 0$. The initial data is negligibly small outside the spatial domain of computation $[0, 1]$. In this numerical test, we chose $\kappa = 10^{-2}$. In addition, we set the maximal time at $T = 1$.

We report on Figure 1a the amplitude of the numerical surface elevation u_1 with ABCs on the computational domain. We set $N = M = 640$. The reference solution $(u_1^{\text{ref}}, u_2^{\text{ref}})$ is computed for $\tau = h = \frac{1}{2560}$ with a CN scheme, in a very large computational domain to avoid any influence of the boundary condition. We also draw on Figure 1b the amplitude of the velocity u_2 . Furthermore, on Figures 1c and 1d, we plot the error $\log_{10} |u_j^{\text{ref}} - u_j|$, $j = 1, 2$, in the domain of computation. As observed, the numerical and reference solutions are very similar, the error being related to the second-order accuracy of the scheme. There is no reflection related to the absorbing boundaries. Finally, we report on Figures 1e and 1f the L^∞ -error $\max_{x \in [0, 1]} |u_j(x, 1) - u_j^{\text{ref}}(x, T)|$, $j = 1, 2$, at final time $T = 1$, when recursively doubling the parameters of the discretization grid, *i.e.* $M = N$ from 80 to 640. A second-order convergence rate in L^∞ -norm is observed.

Example 6.2. To observe the dispersive behavior of the GN system, we consider now the following initial distribution for the free-surface elevation

$$u_1(x, 0) = \exp(-400(x - 0.5)) \sin(20\pi x), \quad (6.2)$$

and set to zero the initial velocity, *i.e.* $u_2(x, 0) = 0$. The initial data u_1 is small outside the spatial computation domain $[0, 1]$ so that it can be considered as numerically compactly supported. We fix $\kappa = 10^{-3}$ and the maximal evolution time $T = 1$.

We first plot on Figure 2a the modulus of the numerical surface elevation u_1 computed with ABCs on the computational domain, for $N = M = 1280$. The reference solution $(u_1^{\text{ref}}, u_2^{\text{ref}})$ is computed with $\tau = h = \frac{1}{6400}$ in a large enough domain with a CN scheme so that we do not see the effect of the boundary condition. Similarly, we draw on Figure 2b the amplitude of the velocity u_2 . In addition, on Figures 2c and 2d, we report the error $\log_{10} |u_j^{\text{ref}} - u_j|$, $j = 1, 2$, in the computational domain. As it can be observed, the numerical and reference solutions superpose, up to the second-order error of the scheme. No spurious reflection can be detected near the absorbing boundaries. For completeness, we plot on Figures 2e and 2f the L^∞ -error $\max_{x \in [0, 1]} |u_j(x, 1) - u_j^{\text{ref}}(x, T)|$, $j = 1, 2$, at final time $T = 1$, when recursively doubling the parameters of the discretization grid, *i.e.* $M = N$ from 160 to 1280. We observe a second-order convergence rate in L^∞ -norm.

To end, the CPU time (log scale, sec.) *vs.* N (log scale) is reported on Figure 3 for the fast evaluation of the convolution operator, fixing the number of spatial grid points to $M = 160$. The total number of time steps N increases from $N = 1.2 \times 10^5$ to $N = 7.2 \times 10^5$, with step 1.2×10^5 . We observe a slope equal to 1, showing that the cost is linear according to $\log(N)$, *i.e.* as $\mathcal{O}(N \log N)$ for the computational time.

We draw on the left of Figure 4 the amplitude of the velocity u_2 until $T = 2$ to display the stability for a longer simulation time. In addition, on the right of Figure 4, we report the error $\log_{10} |u_2^{\text{ref}} - u_2|$, in the computational domain. One can see the wave goes out of the computational domain. The error also increases due to the relation $e^{\sigma t} v_i(x, t) = u_i(x, t)$.

7. CONCLUSION

The one-dimensional linearized Green-Naghdi system in an unbounded domain was reformulated into an initial boundary-value problem in a bounded domain with transparent boundary conditions. A fully discrete

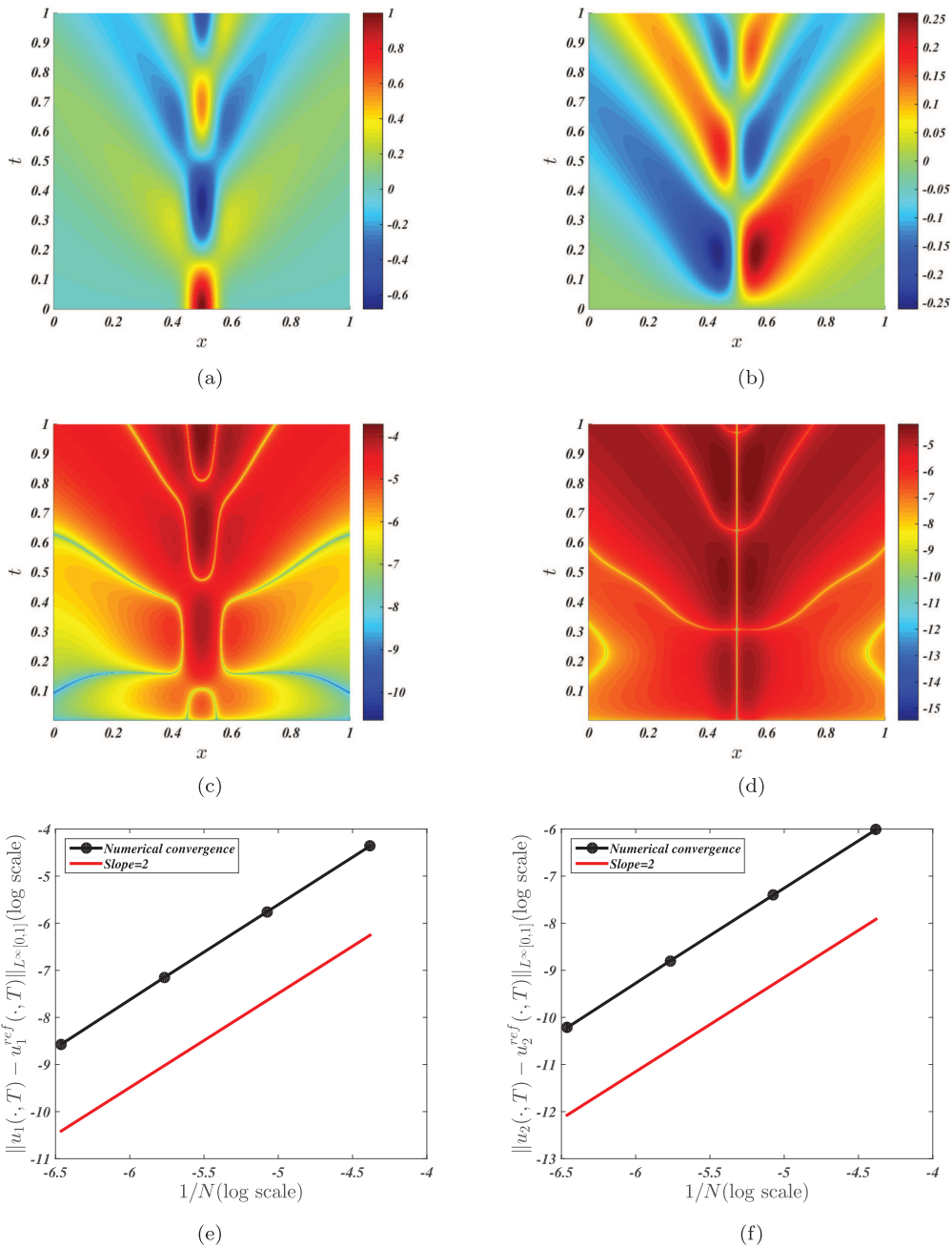


FIGURE 1. (Example 6.1) Illustration of the accuracy of the scheme for both the boundary reflection and the convergence rate. (a) Surface elevation u_1 . (b) Velocity u_2 . (c) Error $\log_{10}(|u_1^{\text{ref}} - u_1|)$. (d) Error $\log_{10}(|u_2^{\text{ref}} - u_2|)$. (e) $\|u_1(\cdot, T) - u_1^{\text{ref}}(\cdot, T)\|_{L^\infty[0,1]}$ vs. $1/N$ (log scale). (f) $\|u_2(\cdot, T) - u_2^{\text{ref}}(\cdot, T)\|_{L^\infty[0,1]}$ vs. $1/N$ (log scale).

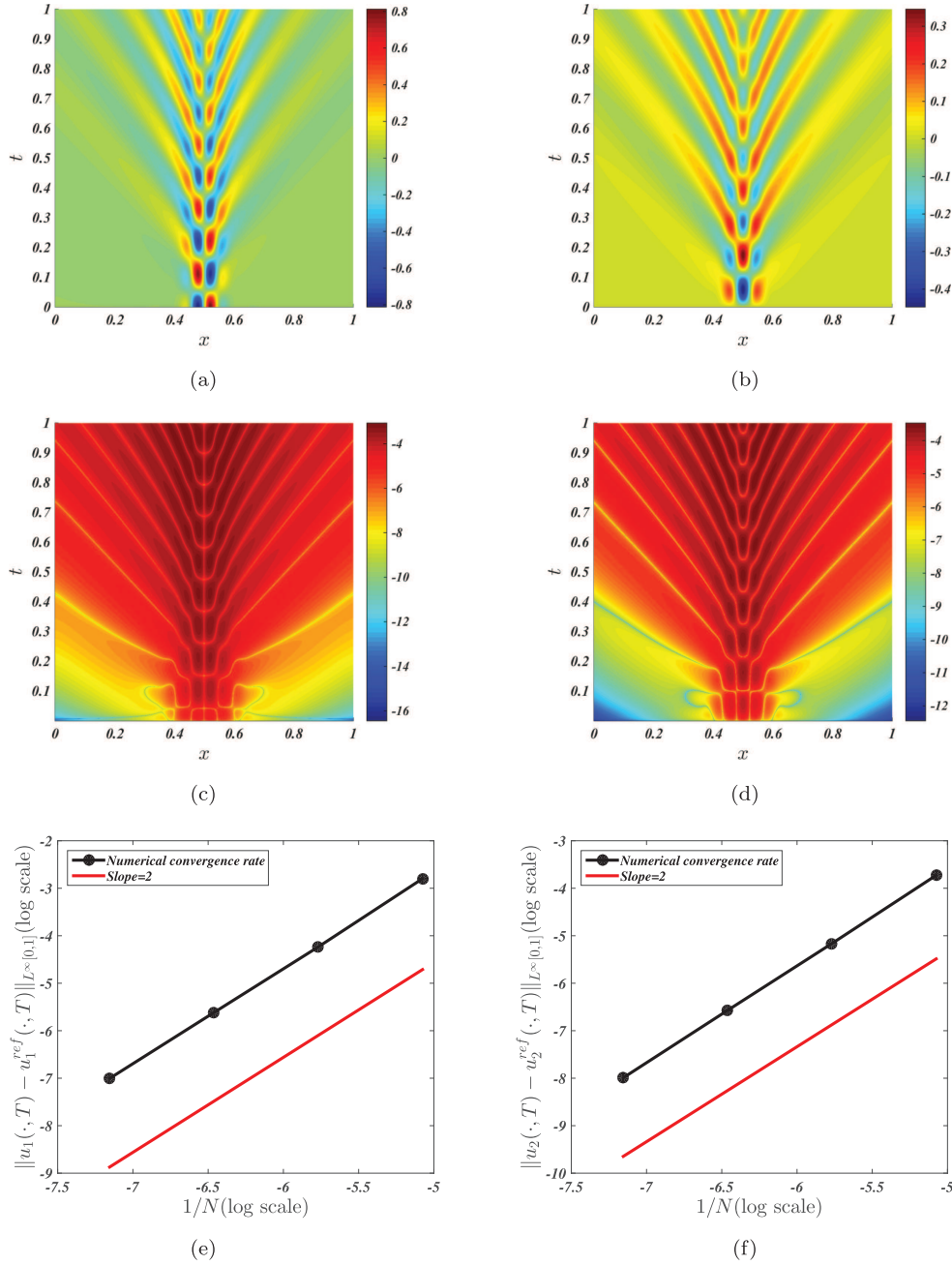


FIGURE 2. (Example 6.2) Illustration of the accuracy of the scheme for both the boundary reflection and the convergence rate. (a) Surface elevation u_1 . (b) Velocity u_2 . (c) Error $\log_{10}(|u_1^{\text{ref}} - u_1|)$. (d) Error $\log_{10}(|u_2^{\text{ref}} - u_2|)$. (e) $\|u_1(\cdot, T) - u_1^{\text{ref}}(\cdot, T)\|_{L^\infty[0,1]}$ vs. $1/N$ (log scale). (f) $\|u_2(\cdot, T) - u_2^{\text{ref}}(\cdot, T)\|_{L^\infty[0,1]}$ vs. $1/N$ (log scale).

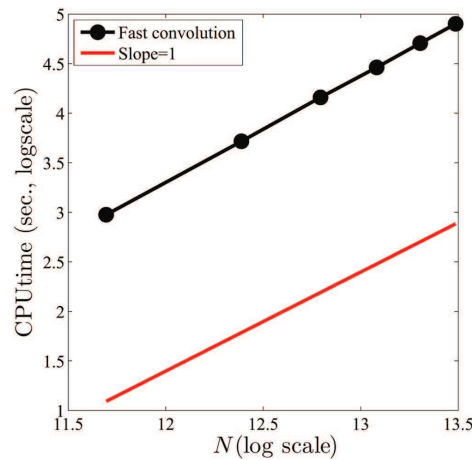


FIGURE 3. Computational time for the evaluation of the convolution by the fast algorithm *vs.* N (for $M = 160$).

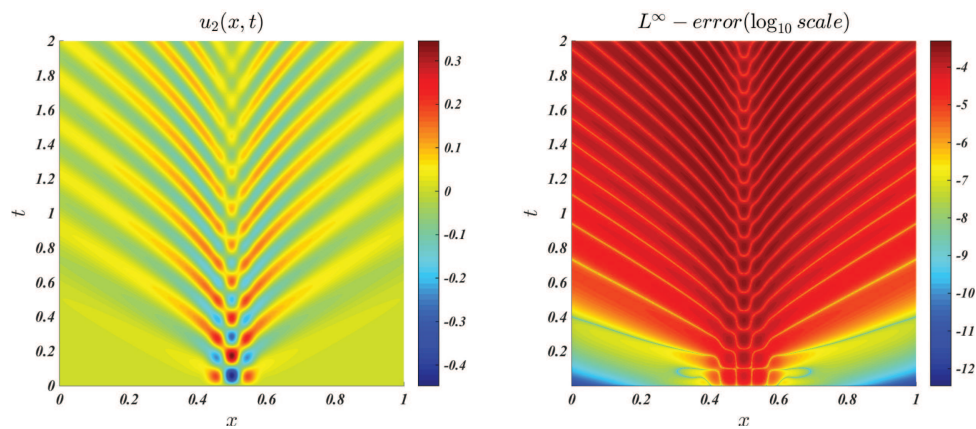


FIGURE 4. (Example 6.2) *Left:* numerical $u_2(x, t)$ for $(x, t) \in [0, 1] \times [0, 2]$. *Right:* error $\log_{10} |u_2^{ref} - u_2|$ for $(x, t) \in [0, 1] \times [0, 2]$.

Crank–Nicolson finite-difference method was proposed to solve the reformulated initial boundary-value problem but with an exact semi-discrete ABC. A fast convolution algorithm is introduced to deal with the convolutions for the exact semi-discrete ABC by using the Padé rational expansion. A criterion determining the damping term was proposed to guarantee the convergence. In this case, it was proved theoretically that the corresponding numerical scheme can achieve a second-order accuracy both in space and time. A numerical example validates the accuracy and efficiency of the proposed numerical method.

The problem that still needs to be solved is that the damping term $e^{-\sigma t}$ should satisfy the stability condition $\sigma \geq \frac{1}{\sqrt{2\kappa}}$. For a small dispersion parameter κ , the damping term $e^{-\sigma t}$ which decays too fast will bring some numerical errors due to the relation $e^{\sigma t} v_i(t) = u_i(t)$. We will deal with this problem in the forthcoming paper. Extensions to higher-dimensional problems still need further investigations. Finally, the variable coefficients and nonlinear cases of the Green-Naghdi system remain open problems as well as the case of the two-layer

Green-Naghdi system. These questions will be addressed in further works based on microlocal analysis techniques [3, 5, 7].

Acknowledgements. This research is partially supported by NSFC under grant Nos. 11502028. Research conducted within the context of the Sino-French International Associated Laboratory for Applied Mathematics-LIASFMA. X. Antoine thanks the LIASFMA funding support of the Université de Lorraine.

REFERENCES

- [1] B. Alpert, L. Greengard and T. Hagstrom, Rapid evaluation of nonreflecting boundary kernels for time-domain wave propagation. *SIAM J. Numer. Anal.* **37** (2000) 1138–1164.
- [2] X. Antoine and C. Besse, Unconditionally stable discretization schemes of non-reflecting boundary conditions for the one-dimensional Schrödinger equation. *J. Comput. Phys.* **188** (2003) 157–175.
- [3] X. Antoine, C. Besse and V. Mouysset, Numerical schemes for the simulation of the two-dimensional Schrödinger equation using non-reflecting boundary conditions. *Math. Comput.* **73** (2004) 1779–1799.
- [4] X. Antoine, A. Arnold, C. Besse, M. Ehrhardt and A. Schädle, A review of transparent and artificial boundary conditions techniques for linear and nonlinear Schrödinger equations. *Commun. Comput. Phys.* **4** (2008) 729–796.
- [5] X. Antoine, C. Besse and P. Klein, Absorbing boundary conditions for the one-dimensional Schrödinger equation with an exterior repulsive potential. *J. Comput. Phys.* **228** (2009) 312–335.
- [6] X. Antoine, C. Besse and P. Klein, Absorbing boundary conditions for general nonlinear Schrödinger equations. *SIAM J. Sci. Comput.* **33** (2011) 1008–1033.
- [7] X. Antoine, C. Besse and P. Klein, Absorbing boundary conditions for the two-dimensional Schrödinger equation with an exterior potential. Part II: discretization and numerical results. *Numer. Math.* **125** (2013) 191–223.
- [8] X. Antoine, E. Lorin and Q. Tang, A friendly review of absorbing boundary conditions and perfectly matched layers for classical and relativistic quantum waves equations. *Mol. Phys.* **115** (2017) 1861–1879.
- [9] A. Arnold, M. Ehrhardt and I. Sofronov, Discrete transparent boundary conditions for the Schrödinger equation: fast calculation, approximation, and stability. *Commun. Math. Sci.* **1** (2003) 501–556.
- [10] A. Arnold, M. Ehrhardt, M. Schulte and I. Sofronov, Discrete transparent boundary conditions for the Schrödinger equation on circular domains. *Commun. Math. Sci.* **10** (2012) 889–916.
- [11] V. Baskakov and A. Popov, Implementation of transparent boundaries for numerical solution of the Schrödinger equation. *Wave Motion* **14** (1991) 123–128.
- [12] T. Fevens and H. Jiang, Absorbing boundary conditions for the Schrödinger equation. *SIAM J. Sci. Comput.* **21** (1999) 255–282.
- [13] A. Green and P. Naghdi, A derivation of equations for wave propagation in water of variable depth. *J. Fluid Mech.* **78** (1976) 237–246.
- [14] T. Hagstrom, New results on absorbing layers and radiation boundary conditions. In: *Topics in Computational Wave Propagation. Lecture Notes in Computational Science and Engineering*, edited by M. Ainsworth, P. Davies, D. Duncan, B. Rynne and P. Martin. Vol. 31. Springer, Berlin, Heidelberg (2003).
- [15] H. Han and Z. Huang, Exact artificial boundary conditions for the Schrödinger equation in \mathbb{R}^2 . *Commun. Math. Sci.* **2** (2004) 79–94.
- [16] S. Ji, Y. Yang, G. Pang and X. Antoine, Accurate artificial boundary conditions for the semi-discretized linear Schrödinger and heat equations on rectangular domains. *Comput. Phys. Commun.* **222** (2018) 84–93.
- [17] S. Jiang and L. Greengard, Fast evaluation of nonreflecting boundary conditions for the Schrödinger equation in one dimension. *Comput. Math. App.* **47** (2004) 955–966.
- [18] M. Kazakova and P. Nobel, Discrete transparent boundary conditions for the linearized Green-Naghdi system of equations. *SIAM J. Numer. Anal.* **1** (2020) 657–683.
- [19] D. Lannes, *The Water Waves Problem: Mathematical Analysis and Asymptotics*. Providence, AMS (2013).
- [20] H. Li, X. Wu and J. Zhang, Local artificial boundary conditions for Schrödinger and heat equations by using high-order Azimuth derivatives on circular artificial boundary. *Comput. Phys. Commun.* **185** (2014) 1606–1615.
- [21] B. Li, J. Zhang and C. Zheng, An efficient second-order finite difference method for the one-dimensional Schrödinger equation with absorbing boundary conditions. *SIAM J. Numer. Anal.* **56** (2018) 766–791.
- [22] Y.Y. Lu, A Padé approximation method for square roots of symmetric positive definite matrices. *SIAM J. Matrix Anal. App.* **19** (1998) 833–845.
- [23] C. Lubich and A. Schädle, Fast convolution for nonreflecting boundary conditions. *SIAM J. Sci. Comput.* **24** (2002) 161–182.
- [24] G. Pang, L. Bian and S. Tang, ALmost EXact boundary condition for one-dimensional Schrödinger equation. *Phys. Rev. E* **86** (2012) 066709.
- [25] A. Quarteroni and A. Valli, *Numerical Approximation of Partial Differential Equations*. Springer, Berlin (1994).
- [26] S. Tsynkov, Numerical solution of problems on unbounded domains. A review. *Appl. Numer. Math.* **27** (1998) 465–532.

- [27] C. Zheng, Approximation, stability and fast evaluation of exact artificial boundary condition for one-dimensional heat equation. *J. Comput. Math.* **25** (2007) 730–745.

Subscribe to Open (S2O)

A fair and sustainable open access model



This journal is currently published in open access under a Subscribe-to-Open model (S2O). S2O is a transformative model that aims to move subscription journals to open access. Open access is the free, immediate, online availability of research articles combined with the rights to use these articles fully in the digital environment. We are thankful to our subscribers and sponsors for making it possible to publish this journal in open access, free of charge for authors.

Please help to maintain this journal in open access!

Check that your library subscribes to the journal, or make a personal donation to the S2O programme, by contacting subscribers@edpsciences.org

More information, including a list of sponsors and a financial transparency report, available at: <https://www.edpsciences.org/en/maths-s2o-programme>