

## On Energy Stable Runge-Kutta Methods for the Water Wave Equation and its Simplified Non-Local Hyperbolic Model

Lei Li<sup>1</sup>, Jian-Guo Liu<sup>2</sup>, Zibu Liu<sup>3,\*</sup>, Yi Yang<sup>4</sup> and Zhennan Zhou<sup>5</sup>

<sup>1</sup> School of Mathematical Sciences, Institute of Natural Sciences, MOE-LSC, Shanghai Jiao Tong University, Shanghai, 200240, P.R. China.

<sup>2</sup> Department of Mathematics and Department of Physics, Duke University, Durham, NC 27708, USA.

<sup>3</sup> Department of Mathematics, Duke University, Durham, NC 27708, USA.

<sup>4</sup> Nanjing Research Institute of Electronics Technology, Nanjing, P.R. China.

<sup>5</sup> Beijing International Center for Mathematical Research, Peking University, P.R. China.

Received 9 March 2021; Accepted (in revised version) 14 April 2022

---

**Abstract.** Although interest in numerical approximations of the water wave equation grows in recent years, the lack of rigorous analysis of its time discretization inhibits the design of more efficient algorithms. In practice of water wave simulations, the trade-off between efficiency and stability has been a challenging problem. Thus to shed light on the stability condition for simulations of water waves, we focus on a model simplified from the water wave equation of infinite depth. This model preserves two main properties of the water wave equation: non-locality and hyperbolicity. For the constant coefficient case, we conduct systematic stability studies of the fully discrete approximation of such systems with the Fourier spectral approximation in space and general Runge-Kutta methods in time. As a result, an optimal time discretization strategy is provided in the form of a modified CFL condition, i.e.  $\Delta t = \mathcal{O}(\sqrt{\Delta x})$ . Meanwhile, the energy stable property is established for certain explicit Runge-Kutta methods. This CFL condition solves the problem of efficiency and stability: it allows numerical schemes to stay stable while resolves oscillations at the lowest requirement, which only produces acceptable computational load. In the variable coefficient case, the convergence of the semi-discrete approximation of it is presented, which naturally connects to the water wave equation. Analogue of these results for the water wave equation of finite depth is also discussed. To validate these theoretic observation, extensive numerical tests have been performed to verify the stability conditions. Simulations of the simplified hyperbolic model in the high frequency regime and the water wave equation are also provided.

---

\*Corresponding author. *Email addresses:* leili2010@sjtu.edu.cn (L. Li), jliu@phy.duke.edu (J.-G. Liu), zibu.liu@duke.edu (Z. Liu), sailors2008@sina.cn (Y. Yang), zhennan@bicmr.pku.edu.cn (Z. Zhou)

**AMS subject classifications:** 65-XX

**Key words:** Runge-Kutta methods, non-locality, hyperbolicity.

## 1 Introduction

Simulation of water wave equations has been a challenging problem due to the bad well-posedness of the equation. To understand this difficulty, we will first review previous remarkable theoretic work on water wave equations. In both Wu's work [28] and Beale and Hou's work [5], the authors used Riemann mappings to find the right variables and rewrote the water wave equation. Both versions of the equation in [28] and [5] have the common leading order structure: a non-local hyperbolic equation, say (1.6). By analyzing this simplified model, we derive an optimal discretization strategy in the form of a CFL condition. This condition is rigorously proved for system (1.6) and numerically verified for water wave equations. Following this road map, we will first derive the simplified model.

The unsteady system of incompressible free surface flow in two-dimension has attracted much theoretic and numerical attention [6, 10, 11, 28]. Governed by the irrotational Euler equation, this free surface flow problem is also referred to as the water wave problem which dates back to the early 20th century [23, 25]. By observing the equation in both Eulerian coordinate and Lagrangian coordinate, insightful analytical results were derived since then. One can refer to [11] for a review of recent related results. In [5], Beale, Hou, and Lowengrub formulated the water wave equation in Lagrangian coordinates and considered the linearization of it which is a nonlocal system. Later in [28], by reducing the system to a nonlocal hyperbolic equation in Eulerian coordinate, Wu derived impressive results on the well-posedness of the water wave equation. These two works directed us to focus on a simplified system that inherits the common dominating structure shared by both works: a nonlocal hyperbolic model, which played a critical role in the proof of well-posedness in [28].

The water wave equation is formulated as follows (see [5]) in Lagrangian coordinates. Consider a  $2\pi$ -periodic two-dimensional fluid with infinite depth whose surface is described by  $z: \mathbb{R} \times [0, \infty) \rightarrow \mathbb{C}$ :

$$z(\alpha, t) = x(\alpha, t) + iy(\alpha, t). \quad (1.1)$$

Here  $\alpha \in \mathbb{R}$  is a material coordinate that parametrizes the undisturbed surface. Periodicity of the fluid wave implies that  $s(\alpha, t) := z(\alpha, t) - \alpha$  is a  $2\pi$ -periodic function in  $\alpha$ . Because the fluid is inviscid and irrotational, the velocity can be written as  $\nabla\Phi$  where  $\Phi(x, y, t)$  is the velocity potential. Let

$$\phi: \mathbb{R} \times [0, \infty) \rightarrow \mathbb{R}, \quad (\alpha, t) \mapsto \phi(\alpha, t) := \Phi(x(\alpha, t), y(\alpha, t), t)$$

be the evaluation of the velocity potential at the surface so that  $\phi(\alpha, t)$  is a periodic function in  $\alpha$  with period  $2\pi$ . In [5], starting from the irrotational Euler equation, Beale, Hou, and Lowengrub derived the equation for the interface of fluid with infinite depth which is parametrized by  $z(\alpha, t)$  in Lagrangian coordinate, i.e.,  $z(\alpha, t)$  and  $\phi(\alpha, t)$  satisfy the following system of equations:

$$\begin{cases} \bar{z}_t = \frac{1}{4\pi i} \int_{-\pi}^{\pi} \gamma(\alpha') \cot\left(\frac{z(\alpha) - z(\alpha')}{2}\right) d\alpha' + \frac{\gamma(\alpha)}{2z_\alpha(\alpha)} =: w(\alpha, t), \\ \phi_t = \frac{1}{2}|w|^2 - gy, \\ \phi_\alpha = \frac{\gamma}{2} + \operatorname{Re} \left[ \frac{z_\alpha}{4\pi i} \int_{-\pi}^{\pi} \gamma(\alpha') \cot\left(\frac{z(\alpha) - z(\alpha')}{2}\right) d\alpha' \right], \end{cases} \quad (1.2)$$

where  $\bar{z}$  means the complex conjugate of  $z$ ,  $\gamma(\alpha, t)$  quantifies the derivative of the dipole strength and  $g$  is the gravity. This derivation could also be found in [3, 6] (Equations (1)-(3) in [6]). Authors of [5] proved that the linearization of (1.2) leads to (1.3) (Equation (2.8) in [5]), namely for  $(\alpha, t) \in \mathbb{R} \times (0, \infty)$ :

$$\begin{cases} \partial_t \eta = \sigma(\alpha, t) \Lambda \zeta + g_1, \\ \partial_t \zeta = -c(\alpha, t) \eta, \\ \partial_t \delta = g_2. \end{cases} \quad (1.3)$$

Here  $\sigma$  and  $c$  are positive, which depend on the solution of the water wave equation, but independent of  $\eta$  and  $\zeta$ .  $\eta$  is the normal component of the perturbation of the position of the interface.  $\delta$  is a combination of the tangential and normal components of the perturbation of the position.  $\zeta$  describes the perturbation of the potential.  $g_1$  and  $g_2$  are some extra terms in the linearization which contain linear terms in  $\eta, \delta$ . The operator

$$\Lambda = (-\Delta)^{1/2} = H\partial_\alpha \quad (1.4)$$

is the 1/2-fractional Laplacian with Fourier symbol  $|k|$ , where  $H$  is the Hilbert transform whose Fourier symbol is  $-i\operatorname{sgn}(k)$ . On  $\mathbb{R}$ , the Hilbert transform  $H$  is given by

$$H(f)(x) = \frac{1}{\pi} \text{p.v.} \int_{-\infty}^{\infty} \frac{f(y)}{x-y} dy.$$

As we shall see, system (1.3) is  $L^2$  stable and dispersive.

Another system with the same structure of (1.3) is derived in [28] in the Eulerian coordinates. Wu achieved remarkable results in [28] in proving the well-posedness of the water wave problem in Sobolev spaces. Wu used a conformal mapping formulation and reduced the water wave system to a quasi-linear hyperbolic system (see (4.6) and (5.8 $\epsilon$ ) in [28] and let  $w = -v$ ) for  $(\beta, t) \in \mathbb{R} \times (0, \infty)$ :

$$\begin{cases} u_t = \sigma(\beta, t) \Lambda v + b(\beta, t) \partial_\beta u + g_1, \\ v_t = -c(\beta, t) u + b(\beta, t) \partial_\beta v + g_2, \end{cases} \quad (1.5)$$

where  $\sigma > 0$ ,  $c > 0$  and  $\beta$  is the Eulerian coordinate instead of a material coordinate. Let  $X$  be the  $x$ -coordinate of the interface, then  $u = X_{tt}$  and  $v = -X_t$  in (1.5). Due to the extra time derivative, this equation can also be viewed as a linearization of the water wave equations in [5]. The usage of conformal mappings in this work was successful. It transformed all the nonlocal terms on a time-dependent domain into the half Laplacian, i.e.  $\Lambda = (-\Delta)^{1/2} = H\partial_\beta$  on  $\mathbb{R}$  which is time-invariant. Wu referred to this system as the 'hyperbolic system' in [28]. We will also preserve this description in the present paper.

The slight difference between system (1.5) and (1.3) is that transport terms appear in the former but disappear in the latter. This is because: in (1.5), variable  $\beta$  is not the material coordinate but a variable associated with the conformal mapping. Except for this difference, both equation (1.5) and (1.3) share the same dominating structure, i.e. the system

$$\begin{aligned} u_t &= \sigma(x,t)\Lambda v + \lambda_1(x,t)u + \lambda_2(x,t)v + f_1, \\ v_t &= -c(x,t)u + f_2 \end{aligned} \quad (1.6)$$

for  $(x,t) \in \mathbb{R} \times (0, \infty)$ . This system is intrinsic to the water wave equation since it is detected in both Eulerian and Lagrangian coordinates. Motivated by the preceding discussion, we will focus on this nonlocal hyperbolic system in the rest of the paper.

As in [6], we impose periodic boundary condition and study the system on the torus, i.e.

$$\begin{aligned} u_t &= \sigma(\theta,t)\Lambda v + \lambda_1(\theta,t)u + \lambda_2(\theta,t)v + f_1(\theta,t), \\ v_t &= -c(\theta,t)u + f_2(\theta,t) \end{aligned} \quad (1.7)$$

with  $\theta \in \mathbb{T} = \mathbb{R}/2\pi\mathbb{Z}$  and  $t \in (0, \infty)$ . The Hilbert transform  $H$  still has the symbol  $-i\text{sgn}(k)$  but the formula now is given by

$$Hf(\theta) = \text{p.v.} \int_{\mathbb{T}} f(\tau) \cot\left(\frac{\theta - \tau}{2}\right) \frac{d\tau}{2\pi}.$$

We will focus on equation (1.7) in the following sections.

Consider the special case of (1.7) where  $\sigma, c$  are constant and  $\lambda_1 = \lambda_2 = f_1 = f_2 = 0$ . We derived much insight into the numerical simulation of the water wave from this case. In this constant coefficient case, the system is reduced to the following second-order (in time) nonlocal hyperbolic equation

$$u_{tt} = -\mu\Lambda u, \quad (1.8)$$

where  $\mu = \sigma c$ . For heuristic purposes, we carry out some preliminary analysis and present the basic properties of (1.8) in Section 2.1. A more careful analysis for the constant-coefficient case is conducted in Section 3.

Numerical studies of water waves have been performed in many papers [6, 8, 12, 13, 19, 20, 27]. The numerical methods can roughly be divided into two classes. In the first

class, conformal mapping was not used. In [6], water wave problems were solved by the discretization of an integral formulation (see Section 5 for more information). However, the convergence was merely proved with time being kept continuous. The discussion of the fully discretized system seems challenging. In the second class [8, 12, 13, 27], conformal mappings are used for numerical simulations but no rigorous numerical analysis for conformal mapping formulation has been performed. Meanwhile, although the analytical properties of the nonlocal hyperbolic system (1.6) are relatively well understood in [5, 28], the numerical studies of such equations have not been thoroughly investigated.

Thus to shed light on the distinct properties of such hyperbolic systems and water wave simulations, we intend to focus on numerical analysis of the simplified model (1.6). Because (1.6) has nonlocal terms and the nonlocal terms have a simple Fourier symbol (say  $-i\text{sgn}(k)$ ), we select the pseudo-spectral approximation in the spatial discretization, which is often favored by wave equations (see [4, 26]).

The primary goal of this paper is to analyze Runge-Kutta methods when applied to such nonlocal wave equations. In particular, we emphasize two insightful properties and implications that our analysis of Runge-Kutta methods provides:

First, we explore the optimal time step sizes in terms of a CFL type condition, when certain explicit Runge-Kutta methods are used. As we shall show in Section 2.1, the hyperbolic system (1.6) is also dispersive and may exhibit multi-scale behavior. Therefore, the time step constraint is more severe in the high-frequency regime. Consequently, finding optimal time steps with respect to the wave numbers is naturally desired [4, 15, 16, 21].

In detail, we have systematically analyzed stability conditions of general Runge-Kutta methods for the hyperbolic system (1.3) with constant coefficients, including the high-frequency regime. In constant coefficient case, we have shown that, naive time discretization of system (1.7) results in the familiar hyperbolic CFL constraint  $\Delta t = \mathcal{O}(\Delta x)$ . If we use Runge-Kutta methods whose absolute stable region contains a part of the imaginary axis, this CFL constraint can be relaxed to  $\Delta t = \mathcal{O}(\sqrt{\Delta x})$  which is huge progress. See Theorem 3.1 for detail.

In high-frequency regime, this relaxation on CFL condition provides an optimal time discretization strategy which reduces computational load when simulating (1.7). In this regime, we consider the equation of  $u'(x, t) = u(\epsilon x, \epsilon t), v'(x, t) = v(\epsilon x, \epsilon t)$  which is the rescaling of  $(u, v)$ . Still in the constant coefficient case, the equation of  $(u', v')$  is rewritten as

$$\begin{cases} u_t = \sigma \Lambda v, & (x, t) \in \mathbb{T}_1 \times [0, \infty), \\ \epsilon v_t = -cu, & (x, t) \in \mathbb{T}_1 \times [0, \infty). \end{cases}$$

Due to the factor  $\epsilon$  before  $v_t$ , careful treatment of CFL conditions is necessary. In theorem 3.2, we conclude that Runge-Kutta schemes whose stability regions cover part of the imaginary axis are stable as long as both time step and spatial grid size resolve the wave oscillation, i.e.

$$\Delta x = \mathcal{O}(\epsilon), \quad \Delta t = \mathcal{O}(\sqrt{\epsilon \Delta x}) = \mathcal{O}(\epsilon).$$

This result is sharp in the view that one cannot capture the accurate wave function without resolving its oscillations.

Second, as a corollary of the stability analysis, we investigated the energy stable property of Runge-Kutta methods. The key point is that the nonlocal hyperbolic system (1.8) with constant coefficients is energy preserving, here the energy is given by  $E = \int_{\mathbb{T}} (|v_t|^2 + \frac{1}{2}\mu v \Lambda v) dx$ . When applying pseudo-spectral approximation in spatial discretization, this Hamiltonian is still conserved by Parseval's equality. We will prove that, when applying Runge-Kutta methods in time discretization, the Hamiltonian is non-increasing as long as the applied Runge-Kutta method has an absolute stable region that covers part of the imaginary axis. Available numerical discretization include explicit- $k$  Runge-Kutta methods for  $k \geq 3$ . See Corollary 3.1 for detail.

In the variable coefficient case, we discuss the extension of stability analysis from the constant-coefficient case and to the full water wave simulations. The proposed CFL conditions (say  $\Delta t = \mathcal{O}(\sqrt{\Delta x})$ ) are verified in numerical experiments.

Notice that all preceding discussion is established for water waves of infinite depth, we will also consider an analog in finite depth case.

The rest of the paper is organized as follows: in Section 2.2, we introduce basic notations and the setup for the numerical analysis. In Section 3, we discuss thoroughly the discretization of the nonlocal system with Runge-Kutta (both explicit and implicit) methods in time and the Fourier spectral method in space. We then study the discretization of the system with variable coefficients in Section 4. We prove the convergence for the semi-discrete schemes using the Fourier spectral method or the filtered Fourier spectral method and then discuss the time discretization using Runge-Kutta methods. We then connect the nonlocal hyperbolic system to water wave equations in Section 5. Analog in water waves of finite depth is also discussed. Lastly, in Section 6, we perform numerical experiments. The stability conditions for the nonlocal hyperbolic system with variable coefficients and water wave equations are confirmed numerically. Numerical experiments suggest possible caustics for the system in the high-frequency regime.

## 2 Preliminaries and basic notations

In this section, we discuss the special case (1.8) and then introduce necessary notations for numerical analysis later.

### 2.1 Basic properties of the nonlocal hyperbolic equation

In this section, we present a concise review of basic properties of (1.8), a special case of the hyperbolic system (1.3). Multiplying by  $u_t$  on both sides of (1.8), and integrating over  $x$ , we derive

$$\frac{d}{dt} \int_{\mathbb{R}} \left( |u_t|^2 + \frac{1}{2} \mu u \Lambda u \right) dx = 0.$$

This means that the energy

$$E = \int_{\mathbb{R}} \left( |u_t|^2 + \frac{1}{2} \mu u \Lambda u \right) dx \quad (2.1)$$

is conserved in time. Because  $v$  also satisfies equation (1.8), we know that

$$E' = \int_{\mathbb{R}} \left( |v_t|^2 + \frac{1}{2} \mu v \Lambda v \right) dx$$

is also conserved. To derive the dispersion relation, suppose that the plane wave  $u(x, t) = Ae^{i(2\pi kx - \omega t)}$  is the solution to the equation (1.8). On the Fourier side, (1.8) reads as

$$\hat{u}_{tt} = -\mu |\xi| \hat{u}. \quad (2.2)$$

Notice that the Fourier transform of plane wave  $u(x, t) = Ae^{i(2\pi kx - \omega t)}$  is  $\hat{u}(\xi, t) = A\delta(k - \xi)e^{-i\omega t}$ , substituting it into the equation yields the dispersive relation:

$$-\omega^2 = -\mu |\xi| \Rightarrow \omega = \pm \sqrt{\mu |\xi|}. \quad (2.3)$$

Because  $\omega$  is real and  $\omega \neq \text{const} \times \xi$ , the system is dispersive.

Next, we derive the Green function of (1.8). The explicit expression of the Green function also suggest the dispersion relationship. In fact, the solution  $u(x, t)$  to the initial value problem

$$\begin{cases} u_{tt} + \mu \Lambda u = 0, & (x, t) \in \mathbb{R} \times \mathbb{R}^+, \\ u(x, 0) = f(x), \\ u_t(x, 0) = g(x) \end{cases} \quad (2.4)$$

can be written as

$$u(x, t) = f(x) * F(x, t) + g(x) * G(x, t). \quad (2.5)$$

Here  $*$  represents convolution in space, i.e.  $f(x) * F(x, t) = \int_{\mathbb{R}} f(x-y)F(y, t)dy$ . Function  $G(x, t)$  is the Green's function and  $F(x, t)$  is the time derivative of it. In fact,  $G(x, t)$  and  $F(x, t)$  satisfy following Cauchy problems respectively:

$$\begin{cases} u_{tt} + \mu \Lambda u = 0, & (x, t) \in \mathbb{R} \times \mathbb{R}^+, \\ u(x, 0) = 0, \\ u_t(x, 0) = \delta(x), \end{cases} \quad \begin{cases} u_{tt} + \mu \Lambda u = 0, & (x, t) \in \mathbb{R} \times \mathbb{R}^+, \\ u(x, 0) = \delta(x), \\ u_t(x, 0) = 0. \end{cases} \quad (2.6)$$

Remember that the nonlocal operator  $\Lambda$  has Fourier symbol  $|\xi|$ , thus  $G$  and  $F$  are respectively given by

$$G(x, t) = \mathcal{F}^{-1} \left( \frac{\sin \left( \sqrt{\mu |\xi|} t \right)}{\sqrt{\mu |\xi|}} \right), \quad F(x, t) = \mathcal{F}^{-1} \left( \cos \left( \sqrt{\mu |\xi|} t \right) \right). \quad (2.7)$$

Here,  $\mathcal{F}^{-1}$  is the inverse Fourier transform, i.e.  $(\mathcal{F}^{-1}f)(x) = \int_{\mathbb{R}} e^{2\pi i x \cdot k} f(k) dk$ . For the sake of completeness, we provide the details and more discussions of the Green's function in Section A.

Consider the special case where  $f(x) = \cos(kx), g(x) = 0$ . Then the solution of (2.4) is  $u(x,t) = \cos(kx)\cos(\sqrt{|\mu|k}|t|)$ . Thus, the angle velocity in space and time has the square-root relation in  $k$ :  $k$  and  $\sqrt{|\mu|k}$ . This can be explained by the dispersion relation (2.3). In fact, this also suggest a relaxed CFL condition as we shall explain later.

**Remark 2.1.** Eq. (1.8) is reminiscent of the surface quasi-geostrophic equations (SQG) studied in [7, 17]. However, the surface SQG equation is dissipative while (1.8) is dispersive.

### 2.2 Notations and setup for numerical analysis

In this work, we consider the one-dimensional nonlocal hyperbolic system (1.7) on  $\mathbb{T} = \mathbb{R}/2\pi\mathbb{Z}$ . The spatial discretization is selected as the Fourier pseudo-spectral method or the filtered Fourier pseudo-spectral method.

We discretize the spatial domain with grid size  $h = 2\pi/N$ , and we denote grid points by  $\theta_j = jh, j \in [N] = \{1, \dots, N\}$ , where  $N \in \mathbb{N}$  is even. We denote the time step size by  $\tau$ , and denote  $t^n = n\tau$ . The notation  $u_j^n$  represents the numerical value of  $u(\theta, t)$  at  $(\theta_j, t^n)$ , and  $\mathbf{u}^n$  represents the vector  $\mathbf{u}^n = (u_j^n)$ .

Given any  $N$ -vector  $\mathbf{f} = (f_j)$ , we expand each component as a sum of discrete Fourier modes via

$$f_j = \sum_{k \in [N]^*} \hat{f}_k e^{ik\theta_j}, \quad j \in [N],$$

where  $[N]^* := \{-\frac{1}{2}N+1, \dots, \frac{1}{2}N\}$  and the discrete Fourier transform  $\hat{f} = (\hat{f}_k)$  are given by

$$\hat{f}_k = \frac{1}{N} \sum_{j \in [N]} f_j e^{-ik\theta_j}, \quad k \in [N]^*.$$

Note that the Hilbert transform  $H$  and differentiation operators become certain multipliers when the Fourier transform is applied. When projected onto a uniform grid, those transforms between two functions reduce to corresponding relations between the discrete Fourier transforms of two functions confined on the grid.

We define the projected differential operator and the projected Hilbert transform  $H$  in the following. For two  $N$ - vector  $\mathbf{f}$  and  $\mathbf{g}$ , we write

$$\mathbf{g} = \mathcal{D}\mathbf{f} \quad \text{to mean} \quad \hat{g}_k = ik\hat{f}_k, \quad k \in [N]^*, \tag{2.8}$$

$$\mathbf{g} = \mathcal{H}\mathbf{f} \quad \text{to mean} \quad \hat{g}_k = -i\text{sgn}(k)\hat{f}_k, \quad k \in [N]^*. \tag{2.9}$$

We introduce the notation  $\mathcal{L} = \mathcal{D}\mathcal{H}$  as the projected  $\Lambda = \partial H$ , so that

$$\mathbf{g} = \mathcal{L}\mathbf{f} \quad \text{means} \quad \hat{g}_k = |k|\hat{f}_k, \quad k \in [N]^*.$$

Recall the discrete inner product between two  $N$ -vectors is defined as

$$\langle \mathbf{f}, \mathbf{g} \rangle = \sum_{j \in [N]} h f_j \bar{g}_j,$$

where  $\bar{g}$  means the complex conjugate. The discrete  $\ell^2$  and  $\ell^\infty$  norms are defined by

$$\|\mathbf{f}\|_2 = \sqrt{\langle \mathbf{f}, \mathbf{f} \rangle}, \quad \|\mathbf{f}\|_\infty = \max_{j \in [N]} |f_j|.$$

For discrete case, we still have Parseval's equality:

**Lemma 2.1.** *The discrete Parseval's equality holds*

$$\langle \mathbf{f}, \mathbf{g} \rangle = \sum_{j \in [N]} h f_j \bar{g}_j = 2\pi \sum_{k \in [N]^*} \hat{f}_k \hat{g}_k^*.$$

### 3 Discretization of the constant-coefficient equations

Consider the constant coefficient case of the simplified hyperbolic model, i.e.

$$\begin{cases} u_t = \sigma \Lambda v, \\ v_t = -c u, \end{cases} \quad (3.1)$$

where  $(\theta, t) \in \mathbb{T} \times (0, \infty)$ . On the Fourier side, this equation reads as

$$\partial_t \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} = \begin{pmatrix} 0 & \sigma |k| \\ -c & 0 \end{pmatrix} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} := \mathbf{A} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix}. \quad (3.2)$$

In this section, we will thoroughly analyze the stability of Runge-Kutta methods applied to (3.2). First, we will conduct the Von Neumann analysis [18] on Runge-Kutta methods. Both explicit and implicit ones will be systematically investigated. As a result, we will derive stability conditions in terms of CFL conditions. These CFL conditions provide necessary guidance for the simulation of water wave equations which is conducted in Section 6. Eq. (3.2) in the high-frequency regime is also considered. In the interest of avoiding aliasing error [24], an optimal discretization strategy is developed as a consequence of the analysis of Runge-Kutta methods. The strategy is optimal in the sense that it resolves oscillations at the lowest requirement. In what follows, we denote the Butcher tableau of a certain  $n$ -step Runge-Kutta method by

$$\begin{array}{c|c} \mathbf{p} & \mathbf{G} \\ \hline & \mathbf{w}^T \end{array}, \quad (3.3)$$

where  $\mathbf{G}$  is the Runge-Kutta matrix,  $\mathbf{w}$  are the weights and  $\mathbf{p}$  are the nodes.

### 3.1 Stability analysis

Direct calculation shows that  $A$  has 2 complex eigenvalues  $\lambda_{1,2} = \pm i\sqrt{c\sigma|k|}$ . Thus, matrix  $A$  is similar to the diagonal matrix  $D := \text{diag}\{\lambda_1, \lambda_2\}$ :

$$A = P^{-1}DP, \quad P = P_1P_2,$$

where

$$P_2 = \begin{pmatrix} 1 & 0 \\ 0 & \sqrt{\frac{\sigma|k|}{c}} \end{pmatrix}, \quad P_1 = \begin{pmatrix} 1 & 1 \\ i & -i \end{pmatrix}.$$

This similarity relation reduces the coupled system (3.2) into the following decoupled system:

$$\partial_t \begin{pmatrix} \hat{u}_1 \\ \hat{v}_1 \end{pmatrix} = D \begin{pmatrix} \hat{u}_1 \\ \hat{v}_1 \end{pmatrix}, \quad \begin{pmatrix} \hat{u}_1 \\ \hat{v}_1 \end{pmatrix} := P \begin{pmatrix} \hat{u} \\ \hat{w} \end{pmatrix}. \quad (3.4)$$

Hence, the Von Neumann analysis on the original coupled system 3.2 reduces to that on each scalar equation in (3.4). Notice that  $P_2$  has entry  $\sqrt{\sigma|k|/c}$  which contains Fourier variable  $k$ , so  $L_2$ -stability of (3.4) does not directly imply  $L_2$ -stability of (3.2). However, this fact does not affect stability of Runge-Kutta methods on (3.2). As we will explain later, one can still ensure stability of certain Runge-Kutta methods of (3.2) under a weighted norm, which actually implies a weaker stability in  $L_2$ -norm. Therefore, we first consider stability analysis of (3.4).

Notice that the 2 eigenvalues of matrix  $D$  have the same absolute value, so we can focus on the stability analysis for only the first equation which is exactly the ODE system  $y' = \lambda_1 y$ .

In the following lemma, an explicit formula of the growth index is derived when a certain Runge-Kutta method is employed. This facilitates the Von Neumann analysis.

**Lemma 3.1.** *Denote  $e$  as the vector of ones. For the linear equation  $y' = \lambda_1 y$ , the RK method applied to this equation reduces to  $y_{n+1} = f(\tau, \nu)y_n$  where*

$$|f(\tau, \nu)| = \psi(\tau, \nu) := \sqrt{\frac{|\det(\mathbf{I} + \tau^2 \nu (\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)|}{|\det(\mathbf{I} + \tau^2 \nu \mathbf{G}^2)|}}. \quad (3.5)$$

Here  $\tau$  is the time step size,  $\mathbf{G}$  the Runge-Kutta matrix,  $\mathbf{w}$  the weights in the Runge-Kutta method (see (3.3)) and  $\nu = c\sigma|k|$ .

The proof merely requires standard techniques in numerical analysis textbooks, hence attached to Appendix B for the sake of completeness.

An important case of this lemma is for explicit Runge-Kutta methods. For explicit Runge-Kutta methods, matrix  $\mathbf{G}$  is a lower triangular matrix, thus  $\det(\mathbf{I} + \tau^2 \nu \mathbf{G}^2) = 1$ . Then formula (3.5) reduces to

$$\psi(\tau, \nu) = \sqrt{|\det(\mathbf{I} + \tau^2 \nu (\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)|}. \quad (3.6)$$

This implies that  $\psi^2(\tau, \nu)$  is a constant coefficient polynomial of  $\tau^2\nu$  which depends on  $\mathbf{G}, \mathbf{w}$ . This fact can also be derived from properties of explicit methods. If one uses explicit RK- $p$  method to discretize the first component of system (3.4), the growth index in Lemma 3.1 has the following expression [18]:

$$f(\tau, \nu) = \sum_{m=0}^p \frac{1}{m!} (\tau\lambda_1)^m = \sum_{m=0}^p \frac{1}{m!} (i\tau\sqrt{\nu})^m. \quad (3.7)$$

Thus  $\psi^2(\tau, \nu) = |f(\tau, \nu)|^2$  is an constant coefficient polynomial of  $\tau^2\nu$ .

Moreover, the absolute stable region of RK- $p$  method contains a part of imaginary axis, i.e.  $[-iC_1(p), iC_1(p)]$  if  $p \geq 3$  [18]. Thus requiring

$$-C_1(p) \leq \pm\tau\sqrt{\nu} \leq C_1(p) \quad (3.8)$$

is sufficient to ensure the stability of explicit RK- $p$  methods when solving system (3.4).

Denote the growth matrix for a certain Runge Kutta method by  $B(\tau)$ . According to [18], we recall the following definitions.

**Definition 3.1.** A scheme is called weakly stable if at least for  $\tau$  sufficiently small, there exists  $\alpha > 0$  such that  $|B(\tau)| \leq 1 + \alpha\tau$  holds. A scheme is called strongly stable if  $|B(\tau)| \leq 1$  holds at least for  $\tau$  sufficiently small. Here  $|B(\tau)|$  is the norm of matrix induced by the norm  $|\cdot|$  defined on  $\mathbb{R}^n$ .

One can see that even with weak stability, the numerical scheme is convergent as  $\tau \rightarrow 0$  for any fixed time  $T > 0$ . Thus, we use (3.7) and Lemma 3.1 to prove the following theorem which provides the stability condition:

**Theorem 3.1.** *There are 2 cases in this theorem:*

1. *For any Runge-Kutta scheme, if there exists a positive real number  $C$  such that  $\tau \leq Ch$  holds, then the scheme is stable for system (3.4) or (3.2). More specifically,*
  - (a) *if  $\text{tr}(\mathbf{G}^2) > \text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)$ , then the scheme is strongly stable;*
  - (b) *if  $\text{tr}(\mathbf{G}^2) \leq \text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)$ , then the scheme is weakly stable.*

*Here  $\mathbf{e}$  is the vector of ones and  $\mathbf{G}, \mathbf{w}$  is defined in (3.3).*

2. *For any Runge-Kutta scheme whose absolute stable region contains a part of the imaginary axis  $[-C_1i, C_1i]$ , there exists a positive constant  $C_2$  such that if*

$$\tau \leq C_2\sqrt{h}, \quad (3.9)$$

*then the scheme is strongly stable.*

*Proof.* We first prove statement (a) and (b) in 1. Recall that  $\nu = c\sigma|k|$  and  $|k|h \leq \pi$ , thus we have

$$\begin{aligned} \tau\nu &= c\sigma|k|\tau \leq c\sigma\pi\tau/h \\ &\leq c\sigma C\pi. \end{aligned} \tag{3.10}$$

Therefore,  $\tau\nu = \mathcal{O}(1)$  and  $\tau^2\nu = \mathcal{O}(\tau)$  as  $\tau \rightarrow 0$ . Using Lemma 3.1 and the fact  $\det(\mathbf{I} + \tau\mathbf{X}) = 1 + \text{tr}(\mathbf{X})\tau + \mathcal{O}(\tau^2)$ , we have

$$\begin{aligned} \psi(\tau, \nu) &= \sqrt{\frac{|\det(\mathbf{I} + \tau^2\nu(\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)|}{|\det(\mathbf{I} + \tau^2\nu\mathbf{G}^2)|}} \\ &= \sqrt{\frac{|1 + \text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)\tau^2\nu + \mathcal{O}(\tau^4\nu^2)|}{|1 + \text{tr}(\mathbf{G}^2)\tau^2\nu + \mathcal{O}(\tau^4\nu^2)|}} \\ &= \sqrt{1 + [\text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2) - \text{tr}(\mathbf{G}^2)]\tau^2\nu + \mathcal{O}(\tau^4\nu^2)} \\ &= 1 + \frac{\text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2) - \text{tr}(\mathbf{G}^2)}{2}\tau^2\nu + \mathcal{O}(\tau^4\nu^2). \end{aligned} \tag{3.11}$$

Recall that  $\tau^2\nu = \mathcal{O}(\tau)$ , thus

$$\psi(\tau, \nu) = 1 + \frac{\text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2) - \text{tr}(\mathbf{G}^2)}{2}\tau^2\nu + \mathcal{O}(\tau^2). \tag{3.12}$$

Then we consider 2 cases respectively:

1. If  $\text{tr}(\mathbf{G}^2) > \text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)$ , then  $\psi(\tau, \nu) < 1$  when  $\tau$  goes to 0, which means strong stability. This proves statement (a).
2. If  $\text{tr}(\mathbf{G}^2) \leq \text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)$ , by the estimate (3.10), we have  $\psi(\tau, \nu) \leq 1 + L\tau + \mathcal{O}(\tau^2)$  when  $\tau$  goes to 0, where

$$L = \frac{\text{tr}((\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2) - \text{tr}(\mathbf{G}^2)}{2} \nu C\pi. \tag{3.13}$$

This implies weak stability which proves statement (b).

For the third statement, by estimate (3.10) again, we have  $\tau\sqrt{\nu} \leq \tau\sqrt{c\sigma\pi/h}$ . If it satisfies  $\tau\sqrt{c\sigma\pi/h} \leq C_1$ , namely

$$\tau \leq \frac{C_1}{\sqrt{c\sigma\pi}}\sqrt{h} := C_2\sqrt{h}, \tag{3.14}$$

then  $\tau\sqrt{\nu}$  stays in the absolutely stable region of the numerical scheme which indicates strong stability.  $\square$

**Remark 3.1.** In the proof of Theorem 3.1, we proved stability in the weighted norm  $\|\mathbf{P}_2 \cdot\|$ , which actually does not directly imply stability in the  $L^2$ - norm. In fact, in this case, we still have stability in the  $L^2$ - norm by Theorem 3.1: notice that  $|k| \leq \sqrt{\pi}/\sqrt{h}$ , so there exists constants  $C_1$  and  $C_2$  such that for all  $k \geq 1$ ,

$$\|(u,v)\| \leq \frac{C\|\mathbf{P}_2(u,v)^T\|}{\sqrt{h}}, \quad (3.15)$$

for the mode  $k=0$ , any numerical scheme is always stable since (3.1) reduces to  $\partial_t u = 0$ ,  $\partial_t v = -cu$ , whose solution is  $u(t) = u(0)$ ,  $v(t) = -cu(0)t$ , which is linear. So the mode  $k=0$  is always stable, no matter which scheme is utilized. Therefore, if the total error of a numerical scheme is  $\mathcal{O}(h^n)$ ,  $n \geq 1$  in norm  $\|\mathbf{P}_2 \cdot\|$ , then the total error will be  $\mathcal{O}(h^{n-1/2})$  in the  $L^2$ -norm.

We have to emphasize that statement (a) and (b) in the preceding theorem are not sharp. This is because that they merely ensure zero stability of a scheme while some of them can be even unconditionally absolutely stable. Despite of this theoretic unsharpness, these two statements are still general and practical because they are conclusive for any Runge-Kutta method. Moreover, they provide necessary guidance for numerical simulations of water wave equations in Section 6.

To illustrate this theorem, we analyze two examples of Runge-Kutta methods: a weighted Euler method which is semi-implicit, and the explicit RK-4 method. Both methods are stable under the sufficient condition  $\tau \leq Ch$  as stated in the first part of the theorem. Nevertheless, this CFL condition is unnecessary for the former method under certain weights. For the second method, the condition  $\tau \leq C\sqrt{h}$  claimed in the second statement of Theorem 3.1 is observed.

(a) Semi-implicit RK2: weighted Euler method. The tableau of a weighted Euler method is:

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1-\delta & \delta \\ \hline & 1-\delta & \delta \end{array},$$

where

$$\mathbf{G} = \begin{pmatrix} 0 & 0 \\ 1-\delta & \delta \end{pmatrix}, \quad \mathbf{w}^T = (1-\delta, \delta).$$

Substitute them into  $\psi(\tau, \nu)$  yields

$$\psi(\tau, \nu) = \sqrt{\frac{1+(1-\delta)^2\tau^2\nu}{1+\delta^2\tau^2\nu}}. \quad (3.16)$$

Suppose that  $\tau \leq Ch$  for some positive constant  $C$ , one can verify that  $\psi(\tau, \nu) \leq 1 + C_1\tau$  holds as  $\tau \rightarrow 0$  for some constant  $C_1$ . Thus by Theorem 3.1, this scheme is stable.

However, if  $\delta \geq \frac{1}{2}$ , then  $\psi(\tau, \nu) \leq 1$  holds without any requirement on  $\tau, h$ . This implies unconditionally strong stability. This example illustrates that condition  $\tau \leq Ch$  is sufficient but not necessary.

(b) Explicit RK4. The tableau of explicit RK4 is:

$$\begin{array}{c|cccc} 0 & & & & \\ \frac{1}{2} & \frac{1}{2} & & & \\ \frac{1}{2} & 0 & \frac{1}{2} & & \\ 1 & 0 & 0 & 1 & \\ \hline & \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & \frac{1}{6} \end{array},$$

where

$$G = \begin{pmatrix} 0 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad w^T = \left( \frac{1}{6}, \frac{1}{3}, \frac{1}{3}, \frac{1}{6} \right).$$

Remember that the absolute stable region of explicit RK-4 contains a part of the imaginary axis, thus by Theorem 3.1, we expect to acquire a CFL condition of form  $\tau \leq C\sqrt{h}$ . Substitute  $G, w^T$  into  $\psi(\tau, \nu)$  and denote  $z = \tau^2\nu$ , we have

$$\psi(\tau, \nu) = \sqrt{\det(1 + z(G - eb^T))^2} = \sqrt{\frac{z^4}{576} - \frac{z^3}{72} + 1}. \tag{3.17}$$

Thus  $\psi(\tau, \nu) \leq 1$  if and only if  $z \leq 8$ . Therefore, by Theorem 3.1, RK-4 is strongly stable if and only if  $z \leq 8$ . Moreover, if

$$\tau \leq \frac{2\sqrt{2}}{\sqrt{\pi\sigma c}} \cdot \sqrt{h}$$

holds, then  $z \leq 8$ . Here  $c, \sigma$  is the constant in the system (3.2). Thus one can take  $C_2 = \frac{2\sqrt{2}}{\sqrt{\pi\sigma c}}$  in Theorem 3.1 to ensure stability.

### 3.2 Energy stable methods

A corollary of Theorem 3.1 is that the discretization of energy  $E$  (see (2.1)) is stable for any explicit Runge-Kutta method whose absolute stable region contains part of the imaginary axis. Remember that in proving Theorem 3.1, we focused on vector  $(\hat{u}', \hat{w})$  in (3.4). Also notice that  $\|\hat{u}'\|_2^2 + \|\hat{w}\|_2^2 = \|\mathbf{P}_2(u, v)^T\|^2$  and

$$\|w\|_2^2 = \frac{\sigma}{c} \sum_{k \in [N]^*} |k| |\hat{v}_k|^2 = \frac{\sigma}{c} \langle \mathcal{L}v, v \rangle,$$

so the stability analysis in Theorem 3.1 is actually established on the  $L^2$ -norm of  $(u', w)$  (or  $\mathbf{P}_2(u, v)^T$ ), which is a combination of  $\dot{H}^{1/2}$ -norm of  $v$  and  $L^2$ -norm of  $u$ . This  $L^2$ -norm is exactly the discretization of energy

$$E = \int_{\mathbb{T}} \left( |v_t|^2 + \frac{1}{2} \mu v \Lambda v \right) dx$$

with a constant factor on the Fourier side by substituting  $v_t = -cu$ . Thus strong stability implies that the energy is non-increasing which is exactly the following corollary:

**Corollary 3.1.** *For any RK method whose absolute stable region contains a part of the imaginary axis  $[-C_1 i, C_1 i]$ , there exists a positive constant  $C_2$  such that when*

$$\tau \leq C_2 \sqrt{h}, \quad (3.18)$$

the discretization of the energy  $E$  in (2.1) is non-increasing, i.e.

$$E_1 = \sum_{k \in [N]^*} \frac{\sigma}{c} |k| |\hat{v}_k|^2 + |\hat{u}_k|^2 \quad (3.19)$$

is non-increasing.

In fact, one has

$$\begin{aligned} \sum_{k \in [N]^*} \|\mathbf{P}_2(u, v)^T\|^2 &= \sum_{k \in [N]^*} \frac{\sigma}{c} |k| |\hat{v}_k|^2 + |\hat{u}_k|^2 \\ &= E_1. \end{aligned}$$

So by the second statement in Theorem 3.1, we know that the scheme is strongly stable and  $E_1$  is non-increasing. This proves the corollary.

### 3.3 High frequency regime

To investigate the behavior of high frequency waves in (3.1), we consider the rescaling  $(x, t) \rightarrow (\epsilon x, \epsilon t)$  which yields the following system:

$$\begin{cases} u_t = \sigma \Lambda v, & (x, t) \in \mathbb{T}_1 \times [0, \infty), \\ \epsilon v_t = -cu, & (x, t) \in \mathbb{T}_1 \times [0, \infty), \end{cases} \quad (3.20)$$

where  $\mathbb{T}_1 = \mathbb{R}/\mathbb{Z}$ . Here  $\epsilon = \ell/K$  is a small number,  $\ell$  is the considered length scale and  $K \gg 1$  is the typical wave number. Eq. (3.20) is equivalent to the following equation of second order in time:

$$\epsilon \partial_{tt} u = -\mu H \partial_x u. \quad (3.21)$$

Again, this rescaled system is Hamiltonian: the energy  $\int_{\mathbb{T}_1} (\epsilon|u_t|^2 + \frac{1}{2}\mu u \Lambda u) d\theta$  is conserved. To see this, one can multiply by  $u_t$  on both sides of (3.21) and integrate over  $\theta$  to derive

$$\frac{d}{dt} \int_{\mathbb{T}_1} \left( \epsilon|u_t|^2 + \frac{1}{2}\mu u \Lambda u \right) d\theta = 0.$$

On the Fourier side, the first order system reads as

$$\partial_t \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} = \begin{pmatrix} 0 & \sigma|k| \\ -\frac{c}{\epsilon} & 0 \end{pmatrix} \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix} := \mathbf{A}_1 \begin{pmatrix} \hat{u} \\ \hat{v} \end{pmatrix}. \tag{3.22}$$

Analysis of RK methods for this system can be conducted in the same way as in Section 3.1 if we replace  $c$  by  $c/\epsilon$ . This concludes the following theorem of stability:

**Theorem 3.2.** *There are 2 cases in this theorem:*

1. *For any Runge-Kutta scheme, if there exists a positive real number  $C$  such that  $\tau \leq \epsilon Ch$  holds, the scheme for system (3.22) is stable. More specifically,*

- (a) *if  $\text{tr}(\mathbf{G}^2) > \text{tr}((\mathbf{G} - \mathbf{e}\boldsymbol{\omega}^T)^2)$ , then the scheme is strongly stable;*
- (b) *if  $\text{tr}(\mathbf{G}^2) \leq \text{tr}((\mathbf{G} - \mathbf{e}\boldsymbol{\omega}^T)^2)$ , then the scheme is weakly stable.*

*Here  $\mathbf{e}$  is the vector of ones and  $\psi(\tau, \nu)$  is defined in (3.3).*

2. *For any RK methods whose absolute stable region contains a part of the imaginary axis  $[-C_1 i, C_1 i]$ , there exists a positive constant  $C_2$  such that when*

$$\tau \leq C_2 \sqrt{\epsilon h}, \tag{3.23}$$

*the scheme is strongly stable.*

Energy is also non-increasing if the scheme is strongly stable, i.e.,

**Corollary 3.2.** *For any RK method whose absolute stable region contains a part of the imaginary axis  $[-C_1 i, C_1 i]$ , there exists a positive constant  $C_2$  such that when*

$$\tau \leq C_2 \sqrt{\epsilon h}, \tag{3.24}$$

*the discretization of the energy  $E$  in (2.1) is non-increasing, i.e.*

$$E_1 = \sum_{k \in [N]^*} \frac{\sigma}{c\epsilon} |k| |\hat{v}_k|^2 + |\hat{u}_k|^2$$

*is non-increasing.*

An implication of Theorem 3.2 is that formula (3.23) theoretically provides guidance on selecting step size when simulating the behavior system (3.22). A notorious difficulty in numerical simulation of high-frequency waves is the spatial aliasing error [24] and its consequential strict requirement of stability. To avoid aliasing error, one needs to ensure  $h = \mathcal{O}(1/K)$  [24]. Recall that  $\epsilon = \ell/K$  where  $\ell$  is the considered length scale, thus  $h = \mathcal{O}(1/K) = \mathcal{O}(\epsilon)$ . As in Theorem 3.2, although  $\tau \leq C\epsilon h$  can ensure stability, if this CFL condition is rigorously obeyed, the numerical simulation would suffer from a heavy load of computation since the time step size  $\tau$  satisfies

$$\tau = \mathcal{O}(\epsilon h) = \mathcal{O}(\epsilon^2) = \mathcal{O}(1/K^2).$$

This small time step size indicates that  $\mathcal{O}(K^2)$  times of computation will be conducted which could be unacceptable. However, still in Theorem 3.2, we have proved that  $\tau \leq C\epsilon h$  can be improved to  $\tau \leq C_2\sqrt{\epsilon h}$  for typical Runge-Kutta schemes whose stable regions cover a part of the imaginary axis. In this case, the requirement on  $\tau$  is reduced from  $\tau = \mathcal{O}(1/K^2)$  to

$$\tau = \mathcal{O}(\sqrt{\epsilon h}) = \mathcal{O}(\epsilon) = \mathcal{O}(1/K).$$

Hence, both time steps and spatial grid sized only need to resolve the wave oscillation to ensure stability, i.e.:

$$\tau = \mathcal{O}(1/K), \quad h = \mathcal{O}(1/K). \quad (3.25)$$

This result is sharp in the view that one cannot capture the accurate wave function without resolving its oscillations. Therefore, Theorem 3.2 and Eq. (3.23) direct design of numerical schemes of (3.20) by validating an optimal stability condition. Moreover, we expect that (3.23) can be employed in simulations of water waves, which is in fact realized in Section 6.

## 4 Fully variable-coefficient system

In this section, we consider system (1.6) with variable coefficients. We will prove the convergence of the semi-discretization and provide some evidence for the validity of the CFL condition, namely  $\tau \leq C\sqrt{h}$  in the variable-coefficient case. However, we have to admit that rigorous proof of convergence for fully discrete approximation is not given, this would be considered in our further work. Nevertheless, careful numerical experiments are conducted to verify the convergence and stability of  $\tau \leq C\sqrt{h}$  for both system (1.6) and the water wave equation. Thus, our analysis is still meaningful in the sense that it enables us to shed light on the convergence and stability conditions of the numerical simulation of the water wave equation.

We assume that all the coefficients are smooth on  $\mathbb{T}$ ,  $\sigma \geq \sigma_0 > 0$  and  $c(\theta, t) \geq c_0 > 0$ . Consider the energy functional

$$E = \frac{1}{2} \left( \langle \Lambda v, v \rangle + \langle v, v \rangle + \left\langle \frac{c}{\sigma} u, u \right\rangle \right). \tag{4.1}$$

Taking the derivative and using (1.7), we find

$$\begin{aligned} \dot{E} &= \langle \Lambda v, v_t \rangle + \langle v, v_t \rangle + \frac{1}{2} \left\langle \frac{d}{dt} \left( \frac{c}{\sigma} \right) u, u \right\rangle + \left\langle \frac{c}{\sigma} u, u_t \right\rangle \\ &= \langle \Lambda v, f_2 \rangle + \langle v, -cu + f_2 \rangle + \frac{1}{2} \left\langle \frac{d}{dt} \left( \frac{c}{\sigma} \right) u, u \right\rangle + \left\langle \frac{c}{\sigma} u, \lambda_1 u + \lambda_2 v + f_1 \right\rangle. \end{aligned}$$

Compared with the energy in constant-coefficient case (2.1), the additional  $\langle v, v \rangle$  term assists on estimating the energy: it is needed to control the linear terms like  $\lambda_2 v$ . With  $L^2$ -norms of  $v$ , one can use  $E$  to control the  $H^{1/2}$ -norm instead of  $\dot{H}^{1/2}$ -norm of  $v$ .

Now we use  $E$  and  $\sqrt{E}$  to control the R.H.S.. First, by the Cauchy-Schwartz inequality, we have

$$\langle \Lambda v, f_2 \rangle = \langle v, \Lambda f_2 \rangle \leq \|\Lambda f_2\|_2 \|v\|_2, \quad \langle v, f_2 \rangle \leq \|v\| \|f_2\|, \quad \left\langle \frac{c}{\sigma} u, f_1 \right\rangle \leq \|u\| \left\| \frac{c}{\sigma} f_1 \right\|$$

hold, thus these terms can be controlled by  $\sqrt{E}$ . Meanwhile,  $E$  can be used to control  $\langle v, -cu \rangle$ ,  $\frac{1}{2} \left\langle \frac{d}{dt} \left( \frac{c}{\sigma} \right) u, u \right\rangle$  and  $\left\langle \frac{c}{\sigma} u, \lambda_1 u + \lambda_2 v \right\rangle$ . Thus, one has:

$$\dot{E} \leq C_1 E + C_2 \sqrt{E}. \tag{4.2}$$

By Young's inequality, we have

$$\dot{E} \leq C_1 E + C_2 \sqrt{E} \leq \left( C_1 + \frac{C_2}{2} \right) E + \frac{C_2}{2}.$$

Let  $C = C_1 + C_2/2$ . The Grönwall inequality yields

$$E(t) \leq e^{Ct} E(0) + \frac{e^{Ct} C_2}{2C}.$$

This implies that  $E$  is bounded in finite time. This is consistent with the hyperbolicity.

### 4.1 Convergence of the semi-discretization

In this subsection, we discretize spatial variables first and then show that the semi-discretization is convergent. As mentioned in Section 2.2, we use the Fourier pseudo-spectral method for spatial discretization. However, for the variable-coefficient case, we desire to damp high-frequency modes to gain sufficient smoothness so that some desired

properties hold (see Lemma 4.4 and Section 4.3 below). Therefore, we introduce the filter function  $\rho$  on the Fourier side [6].

We use  $N$ -vectors  $U_h$  and  $V_h$  to approximate  $u(\tau, t)$  and  $v(\tau, t)$  respectively. Let  $\sigma, g_1, g_2, c$  be restricted to the grid points. Given a filter function  $\rho: (-\pi, \pi] \rightarrow \mathbb{R}$ , we denote the operator with symbol  $\rho_h(k) = \rho(hk)$  by  $\check{\rho}_h$ , i.e.,

$$g = \check{\rho}_h f \quad \text{means} \quad \hat{g}_k = \rho(hk) \hat{f}_k, \quad k \in [N]^*.$$

Then, we have the filtered version of action of operators on  $N$ -vectors:

$$\mathcal{D}_\rho = \check{\rho}_h \mathcal{D} f, \quad \mathcal{L}_\rho = \check{\rho}_h \mathcal{L} f. \tag{4.3}$$

The hugest advantage of this filter function is its generalization and flexibility. First, it includes some typical finite difference schemes. For instance, the centered difference on torus  $(D_c u)_j = \frac{1}{2h}(u_{j+1} - u_{j-1})$  can be regarded as a filtered Fourier differentiation with filter  $\rho(\xi) = \frac{\sin(\xi)}{\xi}$ . Second, if a filter is unnecessary, one can simply set  $\rho = 1$ .

We will assume the following conditions for the filter function:

- Condition 4.1.**
- $\rho \geq 0$ , even and  $\rho \in C^2(-\pi, \pi]$  (Note that  $\rho$  may not be  $C^2$  on torus).
  - There exists  $r \in \mathbb{N}_+$  such that

$$\sup_{\xi \in (0, \pi)} |\xi|^{-r} |\rho(\xi) - 1| < \infty. \tag{4.4}$$

Because  $\rho$  is non-negative, we can then define the natural discrete Sobolev norms associated with  $\rho$  to be

$$\|f\|_{H_h^1}^2 := \|f\|_2^2 + \|\mathcal{D}_\rho f\|_2^2, \quad \|f\|_{H_h^{1/2}}^2 := \sum_{k \in [N]^*} (1 + |k| \rho(kh)) |\hat{f}_k|^2. \tag{4.5}$$

From Lemma 2.1, we have following properties:

**Lemma 4.1.** *Suppose  $f, g$  are two  $N$ -vectors. Then integration by parts formulas hold:*

$$\langle f, \mathcal{D}_\rho g \rangle = -\langle \mathcal{D}_\rho f, g \rangle, \quad \langle f, \mathcal{H} g \rangle = -\langle \mathcal{H} f, g \rangle, \quad \langle f, \mathcal{L}_\rho g \rangle = \langle \mathcal{L}_\rho f, g \rangle.$$

These formulas are guaranteed by Parserval’s equality. In fact:

$$\langle f, \mathcal{D}_\rho g \rangle = 2\pi \sum_{k \in [N]^*} \hat{f} \overline{\rho(kh) ik \hat{g}_k} = -2\pi \sum_{k \in [N]^*} (ik \rho(kh) \hat{f}_k) \overline{\hat{g}_k} = -\langle \mathcal{D}_\rho f, g \rangle.$$

Other equalities can be similarly checked and we omit the details.

With the filter function, we discretize the system in space with the filtered pseudo Fourier spectral method, while keeping time continuous:

$$\begin{cases} \frac{dU_h}{dt} = \sigma \mathcal{L}_\rho V_h + \lambda_1 U_h + \lambda_2 V_h + f_1, \\ \frac{dV_h}{dt} = -c U_h + f_2. \end{cases} \tag{4.6}$$

Here  $\sigma, c, \lambda_1, \lambda_2, f_1, f_2$  are also discretized in space, but continuous in time, i.e., they are evaluated at the grid points.

To prove convergence, we first check the consistency of this discretization. By Fourier analysis and the aliasing formula, we have the following lemma:

**Lemma 4.2.** *Let  $\varphi \in C^\infty(\mathbb{T})$  and  $N \in \mathbb{N}$ . Then, the restriction  $f = (f_j) = (\varphi(\theta_j))$  of  $\varphi$  to the grid points satisfies*

$$(\mathcal{D}_\rho f)_j - \varphi'(\theta_j) = R_1(\theta_j, h, r)h^r, \quad (\mathcal{L}_\rho f)_j - (\Lambda \varphi)(\theta_j) = R_2(\theta_j, h, r)h^r, \quad j \in [N], \quad (4.7)$$

where  $R_i: \mathbb{T} \rightarrow \mathbb{R}$  ( $i = 1, 2$ ) are functions with  $|\partial_\theta^\alpha R_i(\theta, h, r)|$  bounded uniformly in  $\theta$  and  $h$ , for any  $\alpha \in \mathbb{N}$ .

Proof of this lemma can be found in [26]. As a corollary of Lemma 4.2, we have the following consistency result which is direct:

**Lemma 4.3.** *Assume that the exact solution to (1.6) is  $(u, v) \in C^\infty(\mathbb{T} \times [0, T])$  and the filter satisfies (4.4). Let  $U_e = (u(\theta_j, t))$ ,  $V_e = (v(\theta_j, t))$ , i.e. the restriction of exact solutions on grids. Then we have*

$$\begin{cases} \frac{dU_e}{dt} = \sigma \mathcal{L}_\rho V_e + \lambda_1 U_e + \lambda_2 V_e + f_1 + R_3(\theta_j, t; h)h^r, \\ \frac{dV_e}{dt} = -cU_e + f_2 + R_4(\theta_j, t; h)h^r, \end{cases} \quad (4.8)$$

where  $R_i(\cdot, \cdot; h)h^r$  ( $i = 3, 4$ ) are the local truncations errors and  $R_i(\cdot, \cdot; h)$  ( $i = 3, 4$ ) are two smooth functions on  $\mathbb{T} \times [0, T]$  with  $W^{\alpha, \infty}$  norms uniformly bounded in  $h$  for any  $\alpha \in \mathbb{N}$ .

Now we show the convergence of the semi-discretized equations (4.6).

**Proposition 4.1.** Consider (1.6) with  $\sigma \geq \sigma_0 > 0$  and  $c \geq c_0 > 0$ . Assume that all the coefficients are smooth. Let the exact solution to (1.6) be  $(u, v) \in C^\infty(\mathbb{T} \times [0, T])$ . Let  $r$  be the constant in (4.4). Let  $(U_e, V_e)$  be the restriction of the exact solution to grid and  $(U_h, V_h)$  be the numerical solution given by the pseudo-spectral method (4.6) with the same initial values. Then there exists a constant  $M(T) > 0$ , such that  $\forall t \in [0, T]$ :

$$\begin{aligned} \|U_h(t) - U_e(t)\|_2 &\leq M(T)h^r, \\ \|V_h(t) - V_e(t)\|_{H_h^{1/2}} &\leq M(T)h^r. \end{aligned} \quad (4.9)$$

*Proof.* Define the error vectors

$$e_u = U_h - U_e, \quad e_v = V_h - V_e. \quad (4.10)$$

Taking the difference of Eqs. (4.6) and (4.8), we find the error functions satisfy the following equations

$$\begin{cases} \frac{de_u}{dt} = \sigma \mathcal{L}_\rho e_v + \lambda_1 e_u + \lambda_2 e_v + R_3 h^r, \\ \frac{de_v}{dt} = -c e_u + R_4 h^r. \end{cases} \quad (4.11)$$

Consider the energy functional for this ODE system (analogy to (4.1))

$$E = \frac{1}{2} \left( \langle \mathcal{L}_\rho e_v, e_v \rangle + \|e_v\|_2^2 + \langle \frac{c}{\sigma} e_u, e_u \rangle \right). \quad (4.12)$$

Note that  $\langle \mathcal{L}_\rho e_v, e_v \rangle + \langle e_v, e_v \rangle = \|e_v\|_{H_h^{1/2}}^2$  and  $\langle \frac{c}{\sigma} e_u, e_u \rangle$  is equivalent to  $\|e_u\|_2^2$  (i.e. there exist  $C_1 > 0, C_2 > 0$  such that  $C_1 \|e_u\|_2^2 \leq \langle \frac{c}{\sigma} e_u, e_u \rangle \leq C_2 \|e_u\|_2^2$ ).

Therefore,

$$\frac{dE}{dt} = \left\langle \mathcal{L}_\rho e_v, \frac{de_v}{dt} \right\rangle + \left\langle \frac{c}{\sigma} e_u, \frac{de_u}{dt} \right\rangle + \frac{1}{2} \left\langle \frac{d}{dt} \left( \frac{c}{\sigma} \right) e_u, e_u \right\rangle + \left\langle e_v, \frac{d}{dt} e_v \right\rangle.$$

According to Eq. (4.11),

$$\begin{aligned} \left\langle \mathcal{L}_\rho e_v, \frac{de_v}{dt} \right\rangle + \left\langle \frac{c}{\sigma} e_u, \frac{de_u}{dt} \right\rangle &= \langle R_4 h^r, \mathcal{L}_\rho e_v \rangle + \left\langle \frac{c}{\sigma} e_u, \lambda_1 e_u + \lambda_2 e_v + R_3 h^r \right\rangle \\ &= \langle (\mathcal{L}_\rho R_4) h^r, e_v \rangle + \left\langle \frac{c}{\sigma} e_u, \lambda_1 e_u + \lambda_2 e_v + R_3 h^r \right\rangle \\ &\leq M_1 \sqrt{E} h^r + M_2 E. \end{aligned} \quad (4.13)$$

In the first estimate, we used the fact that  $\mathcal{L}_\rho R_4$  is uniformly bounded by the smoothness of the error (by Lemma 4.3). The last term  $\frac{d\langle e_v, e_v \rangle}{dt}$  is straightforward:

$$\begin{aligned} \frac{d\langle e_v, e_v \rangle}{dt} &= 2 \langle e_v, -c e_u + R h^r \rangle \\ &\leq M_2 (\|e_u\|_2^2 + \|e_v\|_2^2 + \|e_v\|_2 h^r). \end{aligned}$$

We have

$$\frac{dE}{dt} \leq M(E + \sqrt{E} h^r) \Rightarrow \frac{d}{dt} \sqrt{E} \leq \frac{M}{2} (\sqrt{E} + h^r).$$

By Grönwall's inequality, we finally obtain

$$\sqrt{E} \leq M(T) h^r, \quad \forall 0 \leq t \leq T,$$

which leads to our estimate for the error directly.  $\square$

## 4.2 Time discretization

In this section, we aim to study the spatial operators on the right-hand side of (4.6). In particular, we are interested in the eigenvalues of the operators, which will shed light on the time discretization of the ODE system (4.6). The strategy is similar as that in Section 3, i.e. we introduce a similar transformation given by  $P(t)$  so that the stiff part of the operator scaling as  $\frac{1}{\sqrt{h}}$  becomes anti-symmetric. For the convenience of further discussion, we introduce the notion of smoothing operators which is an analogy to the big- $O$  notation introduced in [6]:

**Definition 4.1.** Let  $\mathcal{A} = \{A_N\}$  be a family of operators indexed by  $N$ . We define its action on  $N$ -vector  $f$  as  $\mathcal{A}(f) := A_N(f)$ . We say that  $\mathcal{A}$  is  $m$ -th order smoothing, if there exists  $C > 0$  independent of  $N$  such that for any vector  $f$  we have

$$\|\mathcal{A}(\mathcal{D}^p f)\|_2 \leq C\|f\|_2, \quad \|\mathcal{D}_\rho^p(\mathcal{A}(f))\|_2 \leq C\|f\|_2, \quad \forall 0 \leq p \leq m.$$

If  $\mathcal{A}$  is  $m$ -th order smoothing, we denote it as  $\mathcal{A}_{-m}$ .

We note that  $h\mathcal{D}_\rho = \mathcal{A}_0$  since  $|kh| \leq \pi$ . Recall a lemma from [6]

**Lemma 4.4.** For  $\varphi \in C^\infty$ , let  $[\varphi, \mathcal{H}] \cdot = \varphi\mathcal{H} \cdot - \mathcal{H}(\varphi \cdot)$  be the commutator between  $\varphi$  and  $\mathcal{H}$  ( $\mathcal{H}$  is the discrete Hilbert Transform defined in (2.8)). Assume that Condition 4.1 holds for  $\rho$ . Let  $\mathcal{E}_N$  represent the set of  $N$ -vectors. If  $\rho(\pi) = 0$ , then

$$[\varphi, \mathcal{H}](\check{\rho}_h \omega) = \mathcal{A}_{-1}(\omega), \quad \forall \omega \in \mathcal{E}_N; \tag{4.14}$$

if  $\rho(\pi) = 0$  and  $\rho'(\pi) = 0$  hold, then

$$[\varphi, \mathcal{H}](\check{\rho}_h \omega) = \mathcal{A}_{-2}(\omega), \quad \forall \omega \in \mathcal{E}_N. \tag{4.15}$$

We denote the operator on the right hand side of (4.6) as  $A(t) : \mathcal{E}_N^2 \rightarrow \mathcal{E}_N^2$ :

$$A(t) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \sigma \mathcal{L}_\rho v \\ -cu \end{pmatrix}, \tag{4.16}$$

so that (4.6) can be rewritten as

$$\frac{d}{dt} \begin{pmatrix} U_h \\ V_h \end{pmatrix} = A(t) \begin{pmatrix} U_h \\ V_h \end{pmatrix} + \begin{pmatrix} \lambda_1 U_h + \lambda_2 V_h \\ 0 \end{pmatrix} + \begin{pmatrix} f_1 \\ f_2 \end{pmatrix}.$$

We also define the operators  $P(t) : \mathcal{E}_N^2 \rightarrow \mathcal{E}_N \times \mathcal{Q} \subset \mathcal{E}_N^2$  and  $P^{-1}(t) : \mathcal{E}_N \times \mathcal{Q} \rightarrow \mathcal{E}_N^2$  as

$$\begin{aligned} P(t) \begin{pmatrix} u \\ v \end{pmatrix} &:= \begin{pmatrix} u \\ \Lambda_\rho^{1/2}(\sqrt{\frac{\sigma}{c}}v) \end{pmatrix}, \\ P^{-1}(t) \begin{pmatrix} u \\ v \end{pmatrix} &:= \begin{pmatrix} u \\ \sqrt{\frac{c}{\sigma}}\Lambda_\rho^{-1/2}v \end{pmatrix}. \end{aligned} \tag{4.17}$$

Here, the set  $\mathcal{Q}$  is the following subspace of  $\mathcal{E}_N$ :

$$\mathcal{Q} = \left\{ v \in \mathcal{E}_N : \|v\|_{\mathcal{Q}} := \sum_{k \in [N]^*} \frac{1}{\rho(kh)|k|} |\hat{v}_k|^2 < \infty \right\}.$$

We have the following theorem:

**Theorem 4.1.** *Suppose that the filter satisfies Condition 4.1 and  $\rho(\pi) = 0$ . Then, we can decompose the operator  $A(t)$  (Eq. (4.16)) as*

$$A(t) = \frac{1}{\sqrt{h}} P(t)^{-1} A_1(t) P(t) + P(t)^{-1} A_2(t) P(t), \tag{4.18}$$

(recall  $h = 2\pi/N$ ), where the linear operators  $A_1(t) : \mathcal{E}_N^2 \rightarrow \mathcal{E}_N^2$  and  $A_2(t) : \mathcal{E}_N \times \mathcal{Q} \rightarrow \mathcal{E}_N^2$  satisfy

- (i) *The ranges of  $A_i(t)$  are contained in  $\mathcal{E}_N \times \mathcal{Q}$ .  $A_1(t)$  is anti-symmetric and there exist constants  $N_0 > 0, C > 0$  independent of  $h$  such that*

$$\left\| A_1(t) \begin{pmatrix} u \\ v \end{pmatrix} \right\|_2 \leq C(\|u\|_2 + \|v\|_2), \quad \left\| A_2(t) \begin{pmatrix} u \\ v \end{pmatrix} \right\|_2 \leq C\|v\|_{\mathcal{Q}}, \quad \forall N \geq N_0.$$

- (ii) *The eigenvalues of  $P(t)^{-1} A_1(t) P(t)$  are purely imaginary whose norms are bounded by a constant  $C > 0$  independent of  $N$ . The eigenvalues of  $P(t)^{-1} A_2(t) P(t)$  are bounded by a constant  $C > 0$ s independent of  $N$ .*

*Proof.* We consider the operator  $B(t)$  whose domain is  $\mathcal{Q}$ , defined by

$$B(t) := P(t)A(t)P(t)^{-1}.$$

Then, it is given by:

$$B(t) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \sigma \Lambda_\rho \left( \sqrt{\frac{c}{\sigma}} \Lambda_\rho^{-1/2} v \right) \\ -\Lambda_\rho^{1/2} (\sqrt{\sigma c} u) \end{pmatrix}.$$

We then define  $A_1(t)$  as

$$\begin{pmatrix} u \\ v \end{pmatrix} \mapsto A_1(t) \begin{pmatrix} u \\ v \end{pmatrix} := \sqrt{h} \begin{pmatrix} \sqrt{\sigma c} \Lambda_\rho^{1/2} v \\ -\Lambda_\rho^{1/2} (\sqrt{\sigma c} u) \end{pmatrix},$$

and  $A_2(t) := B(t) - \frac{1}{\sqrt{h}} A_1(t)$  is given by

$$A_2(t) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \sigma \left[ \Lambda_{\rho'} \sqrt{\frac{c}{\sigma}} \right] (\Lambda_\rho^{-1/2} v) \\ 0 \end{pmatrix}.$$

We can directly verify that the ranges of  $A_i$  are in  $\mathcal{E}_N \times \mathcal{Q}$ . Moreover,  $A_1$  is bounded and anti-symmetric.

Now we focus on  $A_2$ . Note that

$$\left[ \Lambda_{\rho'} \sqrt{\frac{c}{\sigma}} \right] (\Lambda_\rho^{-1/2} v) = \mathcal{D}_\rho \left[ \mathcal{H}, \sqrt{\frac{c}{\sigma}} \right] (\Lambda_\rho^{-1/2} v) + \left[ \mathcal{D}_{\rho'} \sqrt{\frac{c}{\sigma}} \right] \mathcal{H} \Lambda_\rho^{-1/2} v.$$

Denote  $w = \Lambda_\rho^{-1/2} v$  and it is clear that

$$\|w\|_2 \leq C\|v\|_{\mathcal{Q}}.$$

By Lemma 4.4, the first term is

$$\mathcal{D}_\rho \mathcal{A}_{-1}(w) = \mathcal{A}_0(w).$$

The second term, by the discrete product rule in [6] is also  $\mathcal{A}_0(w)$ . This then verifies (i).

For (ii), we see that the action of  $P(t)^{-1}A_i(t)P(t)$  ( $i = 1, 2$ ) are well-defined for all  $(u, v) \in \mathcal{E}_N^2$ . Hence, they can be understood as matrices. For  $P(t)^{-1}A_1P(t)$ , because  $A_1$  is antisymmetric and bounded, so (ii) holds. We now focus on  $P(t)^{-1}A_2P(t)$ . Suppose that  $(u, v)$  is a complex eigenvector in  $\mathcal{E}_N^2$ , so that

$$P(t)^{-1}A_2P(t) \begin{pmatrix} u \\ v \end{pmatrix} = \lambda \begin{pmatrix} u \\ v \end{pmatrix}.$$

Denote  $\begin{pmatrix} u_1 \\ v_1 \end{pmatrix} = P(t) \begin{pmatrix} u \\ v \end{pmatrix}$

$$|\lambda|(\|u\|_2 + \|v\|_2) = \left\| P(t)^{-1}A_2P(t) \begin{pmatrix} u \\ v \end{pmatrix} \right\|_2 = \left\| A_2P(t) \begin{pmatrix} u \\ v \end{pmatrix} \right\|_2 \leq C\|v_1\|_Q \leq C\|v\|_2.$$

This then shows (ii). □

By Theorem 4.1, we prove that the leading order structure of  $A(t)$  is an anti-symmetric operator, whose eigenvalues are purely imaginary, and scales as  $1/\sqrt{h}$ . Therefore, if we use ODE solvers whose stability region contains some part of the imaginary axis, such as the explicit RK- $p$  methods with  $p \geq 3$ , then we expect the stability condition is still

$$\frac{\tau}{\sqrt{h}} \leq C$$

for variable coefficient cases.

### 4.3 Comments on linear systems with transport terms

Recall the following linear systems:

$$\begin{aligned} u_t &= \sigma(\theta, t)\Lambda v + b(\theta, t)\partial_\theta u + g_1, \\ v_t &= -c(\theta, t)u + b(\theta, t)\partial_\theta v + g_2, \end{aligned} \tag{4.19}$$

where  $\sigma, c, b$  are given coefficient functions, and  $g_1, g_2$  include lower order terms. The transport terms affect the discretization in two aspects:

- (i) First of all, a filtered version of  $\mathcal{L}_\rho$  and  $\mathcal{D}_\rho$  which satisfies  $\rho(\pi) = 0, \rho'(\pi) = 0$  is needed to dampen high frequency modes so that the discretized energy is still stable.

To see this, let us consider the continuous version of the equations, and consider the same energy functional (4.1). One can estimate  $\dot{E}$  similarly as before except for term  $\langle \Lambda v, b \partial_\theta v \rangle$ . To estimate this term, we find

$$\int_{\mathbb{T}} (\Lambda v) b \partial_\theta v d\theta = -\frac{1}{2} \int_{\mathbb{T}} \partial_\theta v [H, b] (\partial_\theta v) d\theta = \frac{1}{2} \int_{\mathbb{T}} v \partial_\theta ([H, b] \partial_\theta v) d\theta,$$

where  $[H, b] = Hb - bH$  is the commutator. Note that  $b$  is smooth. Thus the commutator  $[H, b] \partial_\theta v$  gives a convolution between a smooth function and  $\partial_\theta v$ . It follows that

$$\frac{1}{2} \int_{\mathbb{T}} v \partial_\theta ([H, b] \partial_\theta v) d\theta \leq C \int_{\mathbb{T}} v^2 d\theta.$$

Therefore,  $\langle \Lambda v, b \partial_\theta v \rangle$  can also be bounded by the  $L_2$ -norm of  $v$  which ensures that  $\dot{E} \leq C_1 E + C_2 \sqrt{E}$  still holds.

Unfortunately, the same estimate does not work again for the discrete case. A filter function is necessary for a stable energy. In fact, for the discretized Hilbert transform,  $[\mathcal{H}, b]$  is not smooth in general. In [6], the authors found that  $[\mathcal{H}, b]$  may not even be  $\mathcal{A}_{-1}$ . By Lemma 4.4, one needs a filter  $\rho$  so that the commutator  $[\mathcal{H}, b] \check{\rho} = \mathcal{A}_{-2}$ , which has the smoothing effect to ensure that

$$\dot{E} \leq C_1 E + C_2 \sqrt{E}$$

still holds.

- (ii) On the other side, conventional numerical treatments on the transport terms require the CFL condition of the form

$$\frac{\tau}{h} \leq C.$$

This condition is much more restrictive compared with (3.9), which could be resolved using semi-Lagrangian method [16, 21].

## 5 Regarding water wave simulation

In Section 1, we explained that the nonlocal hyperbolic systems (1.6) and (4.19) are closely related to water wave equations with infinite depth. In this section, we provide sufficient insight into how our results of the nonlocal system imply the stability conditions for the simulation of water wave equations.

Recall the linearization of (1.2), which leads to (1.3). Indeed, this linearization happens for the numerical schemes as well. In [6], the authors proposed a filtered pseudo-spectral differentiation method to discretize the spatial variables. On the basis of condition 1, they also assumed that the filter satisfies (i)  $\rho(\pi) = 0$  and  $\rho'(\pi) = 0$ ; (ii) there exists

$r \geq 4$  such that  $|\rho(\xi) - 1| \leq C|\xi|^r$  for  $\xi \in (-\pi, \pi]$ . Let  $j \in [N]$  and  $(z_j, \phi_j, \gamma_j)$  be the numerical solutions at the grid points. Then, the discretization to (1.2) is given by (Eqs. (7)-(9) in [6])

$$\begin{cases} \frac{d}{dt} \bar{z}_j = \frac{1}{4\pi i} \sum_{p=-N/2+1, p-j \text{ odd}}^{N/2} \gamma_p \cot\left(\frac{\check{\rho}z_j - \check{\rho}z_p}{2}\right) 2h + \frac{\gamma_j}{2(1 + \mathcal{D}_\rho(z_j - \alpha_j))} =: w_j, \\ \frac{d}{dt} \phi_j = \frac{1}{2} |w_j|^2 - g y_j, \\ \mathcal{D}_\rho \phi_j = \frac{\gamma_j}{2} + \operatorname{Re} \left[ \frac{1 + \mathcal{D}_\rho(z_j - \alpha_j)}{4\pi i} \sum_{p=-N/2+1, p-j \text{ odd}}^{N/2} \gamma_p \cot\left(\frac{\check{\rho}z_j - \check{\rho}z_p}{2}\right) 2h \right]. \end{cases} \quad (5.1)$$

To control the numerical error, the authors in [6] introduced the following functions:

$$\begin{aligned} \eta_j(t) &= \operatorname{Im} \left[ (z_j(t) - z(\alpha_j, t)) \frac{\bar{z}_\alpha(\alpha_j, t)}{|z_\alpha(\alpha_j, t)|} \right], \\ \delta_j(t) &= \operatorname{Re} \left[ (z_j(t) - z(\alpha_j, t)) \frac{\bar{z}_\alpha(\alpha_j, t)}{|z_\alpha(\alpha_j, t)|} \right] + (\mathcal{H}\eta)_j, \\ \zeta_j &= (\phi_j(t) - \phi(\alpha_j, t)) - \operatorname{Re}(w_j(z_j(t) - z(\alpha_j, t))). \end{aligned} \quad (5.2)$$

Moreover, they assumed the strong Taylor sign condition (Eq. (88) in [6])

$$c(\alpha, t) := -\partial_n p \geq c_0 > 0$$

and that the true solutions are smooth with

$$\sigma(\alpha, t) := \frac{1}{|z_\alpha|} \geq \sigma_0 > 0.$$

Based on these assumptions, the authors in [6] found that these variables for errors satisfy the following semi-linear non-local hyperbolic system (Eqs. (89)-(91)):

$$\begin{aligned} \partial_t \eta_j &= \sigma(\alpha_j, t) (\Lambda \zeta)_j + \mathcal{A}_0(\eta, \delta) + \mathcal{A}_0(\zeta) + R_5(h) h^r, \\ \partial_t \zeta_j &= -c(\alpha_j, t) \eta_j + \frac{1}{2} |w_j(t) - w(\alpha_j, t)|^2, \\ \partial_t \delta_j &= \mathcal{A}_0(\eta, \delta, \zeta) + R_6(h) h^r. \end{aligned} \quad (5.3)$$

See Definition 4.1 for  $\mathcal{A}_0$ . The leading order behavior of the semi-linear system (5.3) is exactly (1.3), the nonlocal hyperbolic system. Making use of this hyperbolic structure, the authors showed that the semi-discrete system (5.1) converges to the original water wave problem (1.2). However, there was no discussion about time discretization.

Therefore, given the same linearization effect on both continuous equations and numerical schemes, we expect that similar stability conditions hold for both (1.6) and the water wave equation. Because the errors satisfy the leading order equation of the non-local hyperbolic system (1.3) or (1.6), previous discussion in (4.2) convinced us that the

stability conditions for time discretization would be similar to those for (1.6). If one uses the scheme (5.1) and RK4 for time discretization, a relaxed constraint

$$\frac{\tau}{\sqrt{h}} \leq C,$$

for stability should be observed. We will examine this numerically in Section 6.

**Remark 5.1.** According to (1.5), if we discretize the water wave equation based on the conformal mapping method in [28] instead of the discretization used in the Lagrangian formulation as in [6], we will have a transport term for the numerical error.

We discuss briefly what happens if the water wave problem is of finite depth  $H_0 > 0$ , i.e. the fluid is bounded by a rigid plane  $y = -H_0$ . One of the main interests of water waves with finite depth is the shallow water wave. Recent research focuses of numerical simulations and design of numerical schemes of shallow water waves [2, 9, 22]. In this section, we also comment on the generalization of our results to finite-depth water waves, not only for shallow water waves.

In this finite-depth case, the potential  $\Phi = \Phi(x, y, t)$  satisfies the boundary condition

$$\frac{\partial \Phi}{\partial y} = 0 \text{ on } y = -H_0;$$

the half Laplacian  $\Lambda = (-\Delta)^{1/2}$  does not have Fourier symbol  $|k|$  any more. Now its Fourier symbol should be  $|k| \tanh(H_0|k|)$ . Therefore, another nonlocal hyperbolic system similar to equation (1.7) can be derived:

$$\begin{aligned} u_t &= \sigma(\theta, t) \Lambda_1 v + g_1, \\ v_t &= -c(\theta, t) u + g_2. \end{aligned} \tag{5.4}$$

Here  $\Lambda_1$  is the non-local operator with Fourier symbol  $|k| \tanh(H_0|k|)$ . The derivation of symbol  $|k| \tanh(H_0|k|)$  can be found in [1].

When  $c, \sigma$  are positive constants and  $\lambda_1 = \lambda_2 = g_1 = g_2 = 0$  in (5.4), this system reduces to a non-local linear equation of second order in time:

$$u_{tt} = -\mu \Lambda_1 u. \tag{5.5}$$

The dispersion relation of (5.5) is then given by  $\omega^2 = \mu |k| \tanh(H_0|k|)$ .

Notice that  $\tanh(H_0|k|)$  is bounded, thus analogue of Theorem 3.1 and Corollary 3.1 can be proved for discretizations of this new non-local system. If we again adopt Runge-Kutta methods whose absolute stable region contain some part of the imaginary axis, then by analogue of Theorem 3.1, the corresponding CFL condition is again  $\tau \leq C\sqrt{h}$ ; by analogue of Corollary 3.1, the discretization is energy stable as well.

For the high-frequency regime, one can rescale in time and apply the result for  $\epsilon = 1$ . The conclusions are again similar.

In the variable-coefficient case, because  $\Lambda_1$  can be regarded as the composition of the Hilbert transform and the pseudo-derivative with symbol  $k \tanh(H_0|k|)$ , the commutator estimates in Theorem 4.1 will not change either. This new commutator is smooth at  $k=0$  but it does not help us to damp the high frequency for controlling aliasing errors. Hence, all the conclusions will be similar to the infinite depth case. We expect that when using the energy stable Runge-Kutta methods to solve them, the CFL condition is again

$$\tau \leq C\sqrt{h}.$$

This will be left for future numerical study.

## 6 Numerical examples

In this section, we present some simulations to verify our conclusion and carry out careful numerical experiments. We will introduce our results in the following way: in Section 6.1 and 6.4, we utilize the explicit RK-4 method for the temporal discretization to discretize both the simplified nonlocal hyperbolic system (1.6) and the water wave equation. We verify that the stability condition agrees with (3.9) for both systems which are compatible with our previous analysis. Convergence of the discretization of the nonlocal hyperbolic system is demonstrated in Section 6.2 and exploration of the nonlocal system in the high-frequency regime is performed in Section 6.3. Moreover, a turn-over wave example in [6] is recovered to verify the correctness of our simulation program run in Section 6.4.

For all simulations in this section, periodic boundary conditions are selected. The (filtered) Fourier spectral method is conducted for spatial discretization.

### 6.1 Stability conditions of RK methods for the nonlocal hyperbolic system

In this example, we test stability conditions for the nonlocal hyperbolic system (1.6) with  $g_1=0, g_2=0$ . We consider both the constant-coefficient case with

$$c=3, \quad \sigma=1$$

and the variable-coefficient case with

$$c(\theta, t) = \exp(\cos(\theta+t)), \quad \sigma(\theta, t) = 2 + \sin(\theta+t).$$

We perform simulations for various step sizes using Fourier spectral method in space and forward Euler (FE) and Runge-Kutta 4 (RK4) for temporal discretization. The numerical solutions are calculated up to  $T=10$ . Results are presented in Fig. 1 where the blue part represents the unstable region while the yellow part represents the stable region. The stability in both Section 6.1 and Section 6.4 was determined by the amplitude at a certain breaking point  $T_0$ , before the terminal time  $T$ . For a certain group of temporal step size and grid size, If the  $L_\infty$  norm of the solution at  $T_0$  is larger than a threshold (or even diverges, say 'NaN'), then it is determined as instable. Otherwise, it is stable.

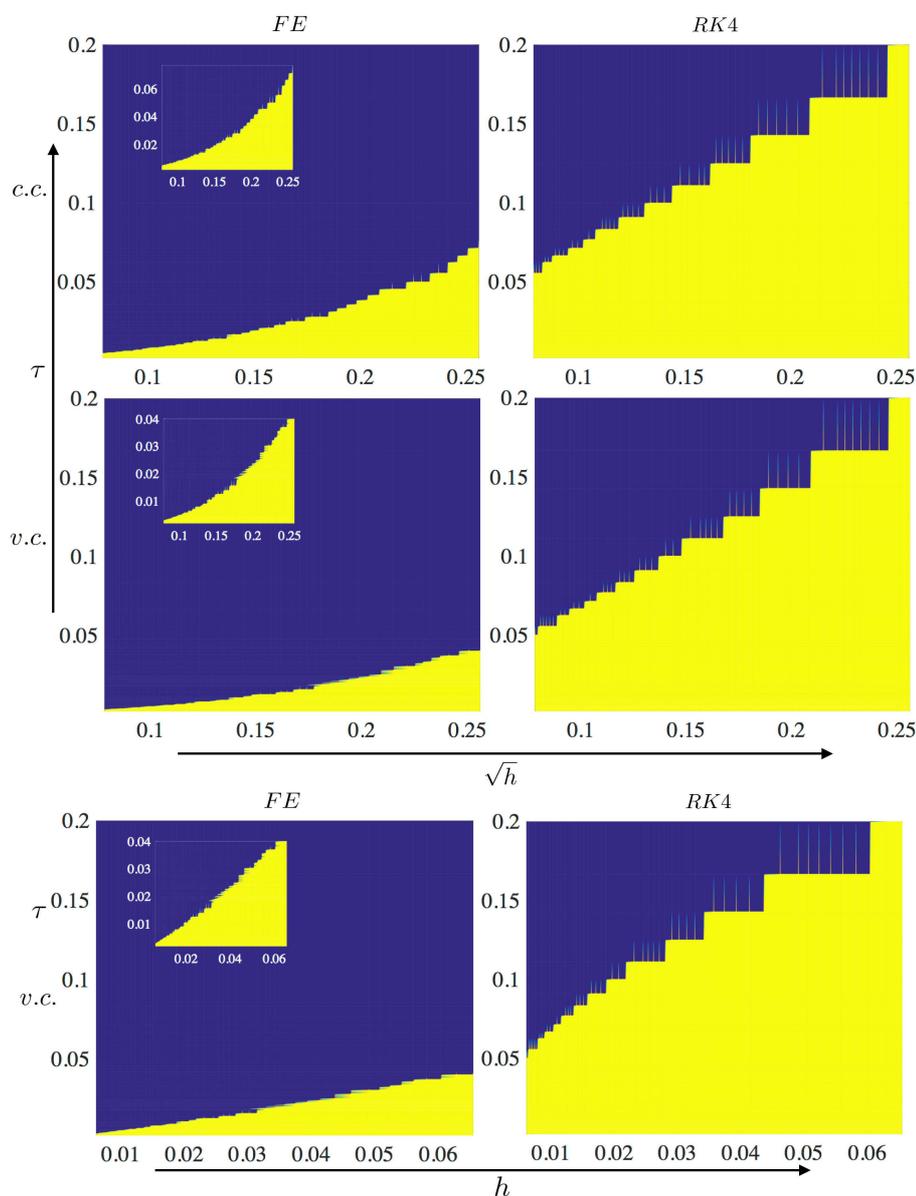


Figure 1: Stability conditions for nonlocal hyperbolic system (1.6). The first row shows the results for the constant coefficient (c.c.) case, while the second row shows the results for the variable coefficient case (v.c.). The two columns are for forward Euler and RK4 respectively. Bottom: Re-plot of second row into  $h$ - $\tau$  plane.

In the top half of Fig. 1, we plot stability in the  $\sqrt{h}$ - $\tau$  plane. Notice that the borders for RK4 are more flat and similar to lines. This plot indicates that the stability condition for RK4 is really (3.9). Meanwhile, the borders for FE are some convex curves similar to a parabola. Thus the stability condition for FE should be the linear CFL condition  $\tau \leq Ch$ .

In order to illustrate this point clearly, we zoom in on the figures and re-plot the stable region in  $h$ - $\tau$  plane. To check the condition for FE, we re-plot the variable-coefficient case as shown at the bottom of Fig. 1. The new figure shows that the stability condition for FE is  $\tau \leq Ch$ . In all cases, the stable region of RK4 is much larger than that of FE, which indicates that RK4 is more stable than FE.

These results verify our analysis in Sections 3 and 4. Recall that the stability region for FE only intersects the imaginary axis at  $z = 0$ , while RK4 contains some part of the imaginary axis. Therefore, by Theorem 3.1, the CFL condition for RK4 should be (3.9); the linear CFL condition sufficiently ensures stability of FE.

## 6.2 Convergence analysis

In this subsection, we verify the convergence numerically for the nonlocal hyperbolic system (1.6). For transport terms, we take  $g_1 = 0$ ,  $g_2 = 0$ . The constant coefficients are given by

$$c = 3, \quad \sigma = 1.$$

The initial conditions are given by

$$u_0(\theta) = e^{\sin(\theta)} + \cos(\theta), \quad v_0(\theta) = \cos^2(\theta).$$

For the temporal discretization, forward Euler (FE), backward Euler (BE), Crank-Nicolson (CN) and Runge-Kutta 4 (RK4) are used.

All simulations are computed up to  $T = 2$ . The reference solution (or ‘accurate solution’) is computed using Runge-Kutta 4 with  $h = 2\pi/2^7$  and  $\tau = 10^{-5}$ . The error plots are shown in Fig. 2. In Fig. 2(a), spectral accuracy is clearly observed in spatial discretization: when  $h \approx 0.2$ , the errors have already been dominated by the temporal error. Fig. 2(b) indicates that the temporal errors are of the order as expected. Therefore, our discretization schemes indeed converge.

## 6.3 The system in the high-frequency regime

In this section, we investigate the system in the high-frequency regime. In particular, we aim at verifying whether there is a caustic phenomenon that is similar to [4].

In research of high-frequency waves, a WKB kind initial value is typically considered. Meanwhile, we rescale the system as in Section 3.3. Therefore, we consider (3.21) with a selected initial value:

$$\begin{cases} \epsilon \partial_{tt} u = -\mu H \partial_x u, & (x, t) \in \mathbb{T}_1 \times \mathbb{R}^+, \\ u(x, 0) = e^{-100(x-0.5)^2} e^{i \log(20 \cosh(5x-2.5)) / \epsilon}, \\ u_t(x, 0) = e^{-100(x-0.5)^2} e^{i \log(20 \cosh(5x-2.5)) / \epsilon}. \end{cases} \quad (6.1)$$

The initial value consists of a Gaussian function and a high-frequency term. The former is used to control the support and the latter is a WKB type function.

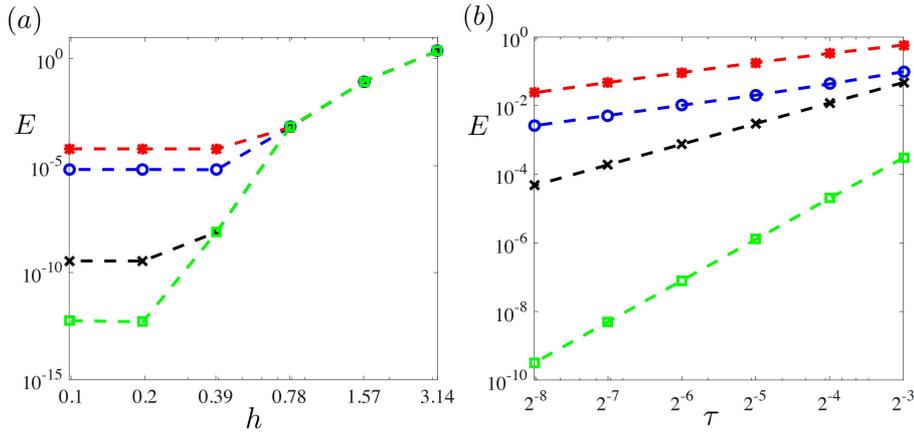


Figure 2: Convergence study for forward Euler (blue circles), backward Euler (red stars), Crank-Nicolson (black crosses) and Runge-Kutta 4 (green squares). (a). Spectral convergence in spatial with  $\tau=10^{-5}$  (b). Time convergence, with  $h=2\pi/2^7$ .

In numerical experiments, we select  $\mu = 1$  in (6.1) and different rescaling factors:  $\epsilon = 2^{-i}$ ,  $i = 4, 5, \dots, 12$ . We plot the snapshot of the amplitude at  $t = 0.0625$  for  $\epsilon = 2^{-4}, 2^{-6}, 2^{-8}, 2^{-10}$  in Fig. 3. The figure indicates that the amplitude increases when  $\epsilon$  decreases. Remember that the initial value  $u(x, 0)$  is of an amplitude no greater than 1 for all different  $\epsilon$ 's, thus this growing trend of the amplitude provides evidence for caustic phenomenon. To obtain a more careful observation, we check the maximum amplitude before time  $t = 0.0625$  in the whole domain  $[0, 1]$ , and plot the following log-log figure in Fig. 4. From Fig. 4, we observe that when  $\epsilon$  is sufficiently small, the curve is almost a line whose slope is 1, which indicates that the maximum amplitude is approximately proportional to  $1/\epsilon$ . This observation also supports the existence of the caustic phenomenon.

### 6.4 Stability conditions for water wave simulation

In this section, we simulate the water wave equation (Eq. (1.2) with  $\alpha \in \mathbb{T}$ ). As in previous simulations, the spatial discretization is implemented using the filtered Fourier spectral method. The filter function we use in this section is the same as in [6, Section 6]:

$$\rho(\xi) = \exp(-10(|\xi|/\pi)^{25}), \quad \xi \in (-\pi, \pi]. \tag{6.2}$$

This filter function numerically satisfies the condition  $\rho(\pi) = 0$  and  $\rho'(\pi) = 0$ :  $|\rho(\pi)|$  and  $|\rho'(\pi)|$  are sufficiently small.

To verify that our simulation programs, we first recover the same example in [6, Section 6] with turn-over phenomenon. The initial data are given by:

$$\begin{aligned} x(\alpha, 0) &= \alpha, \\ y(\alpha, 0) &= 0.6\cos(\alpha), \\ \gamma(\alpha, 0) &= 1 + 0.6\sin(\alpha). \end{aligned} \tag{6.3}$$

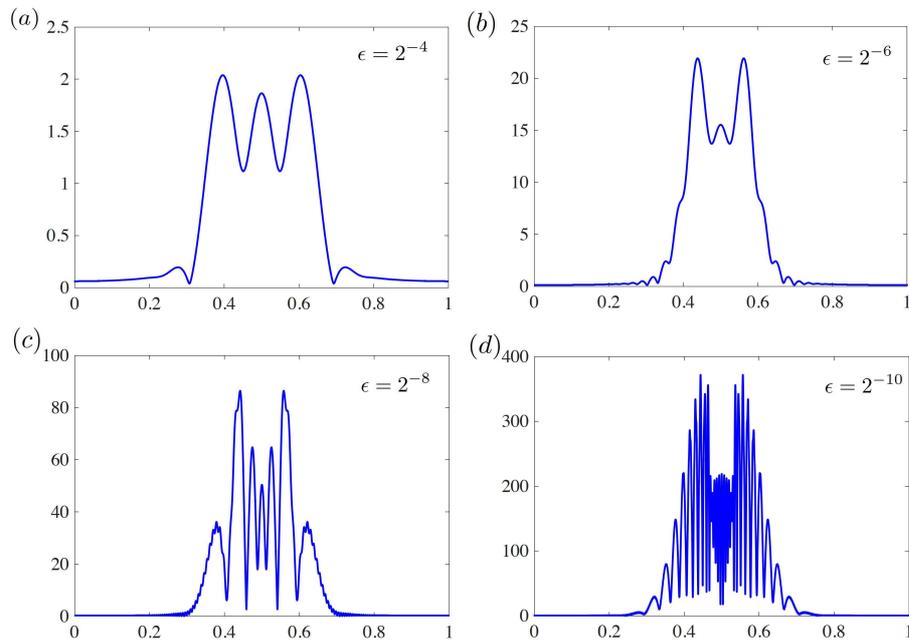


Figure 3: Snapshot for amplitude  $|u|$  versus  $x$  at  $t=0.0625$  for different  $\epsilon$ s. (a)~(d) for  $\epsilon=2^{-4}, 2^{-6}, 2^{-8}, 2^{-10}$  respectively.

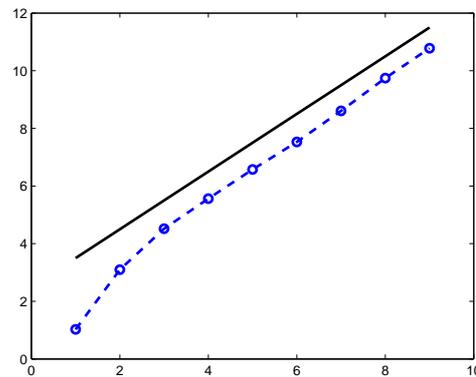


Figure 4: Blue circles represent the plot of  $y = \max_{t \leq 0.0625, x \in [0,1]} \log_2 |u(x,t)|$  versus  $z = -\log_2 \epsilon - 3$ , where  $\epsilon = 2^{-i}$ ,  $i = 4, 5, \dots, 12$ . Meanwhile, a reference line with a slope of 1 is also drawn to show the quantitative relation between  $y$  (or amplitude) and  $z$  (or  $\epsilon$ ).

In the simulation, the spatial grid size is  $h = 1/512$  and the time step size is  $\tau = 1/4000$ . we use RK4 for time discretization here. The snapshots of the waves at different times are shown in Fig. 5. As one can see in the figure below, the same numerical results in [6] are recovered.

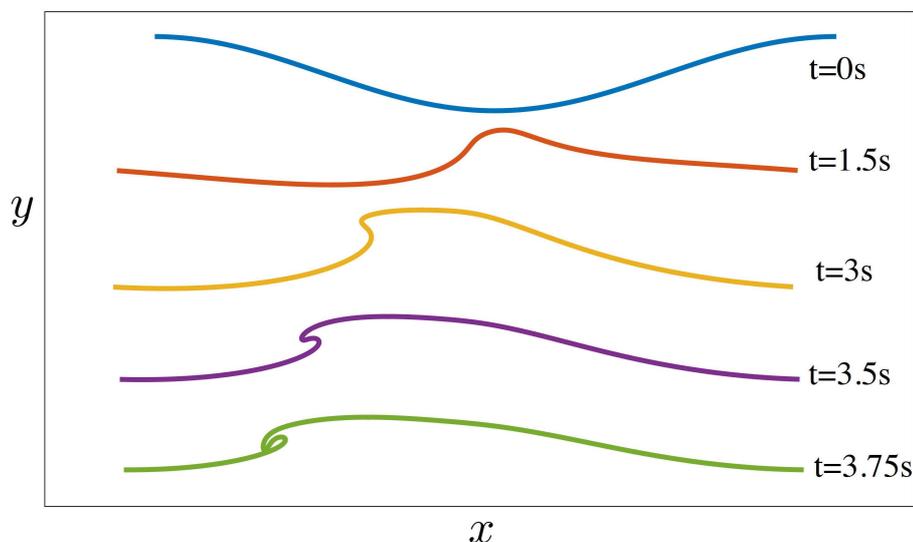


Figure 5: Turn-over of water waves. As time evolves, the water wave turns over gradually. When the time is close to 3.75, the wave tends to break.

Now, we numerically investigate the stability condition for the discretization for the water wave equation. We use the following initial data:

$$\begin{aligned} x(\alpha, 0) &= \alpha, \\ y(\alpha, 0) &= 0.3\cos(\alpha), \\ \gamma(\alpha, 0) &= 1 + 0.3\sin(\alpha). \end{aligned} \quad (6.4)$$

The spatial discretization is performed using the filtered Fourier spectral method with the same filter in Eq. (6.2). For temporal discretization, we select the FE method and the RK-4 method. The simulations are performed up to time  $T = 4$ .

The results are presented in Fig. 6. Same as in Section 6.1, the blue part represents the unstable region while the yellow part represents the stable region. Also, partially zoomed-in versions of the figure are presented in the case of FE to clearly check the stability condition. From this figure, one can tell that the stability condition for RK4 is close to  $\tau \lesssim \sqrt{h}$ , while the stability condition for FE is more similar to  $\tau \lesssim h$ . This observation is in accordance with our conclusions in Section 5 and previous analysis.

## Acknowledgements

L. Li was supported by National Natural Science Foundation of China (Grant No. 31571071) and Shanghai Sailing Program 19YF1421300. J.-G.Liu was supported in part by DMS-2106988. Z. Zhou was supported by the National Key R&D Program of China, Project Number 2021YFA001200, and the NSFC, grant number 12171013.

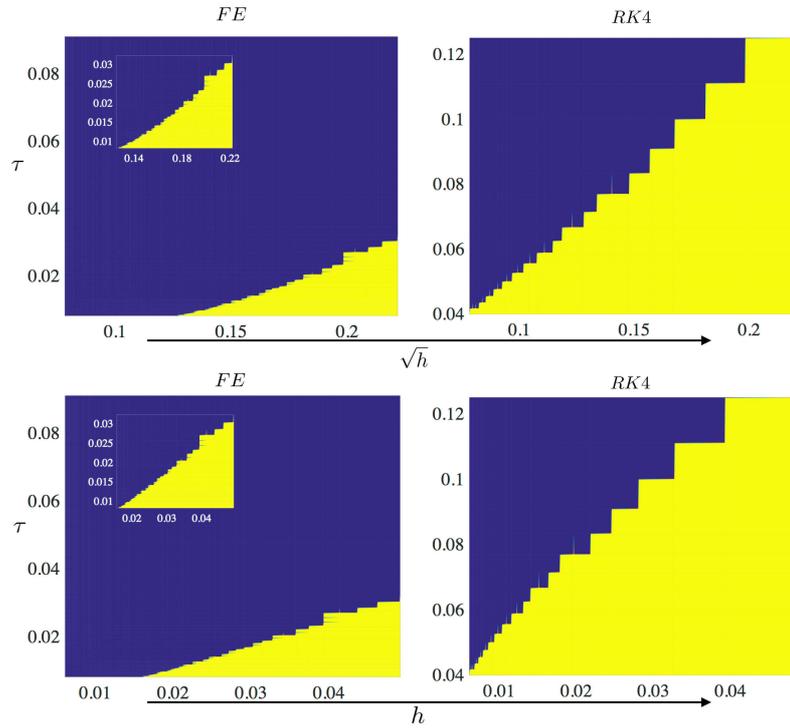


Figure 6: Stability for water wave equations. The two columns are for forward Euler and RK4 respectively. Top: Stable regions plotted in  $\sqrt{h}-\tau$  plane. Bottom: Re-plot of the top into  $h-\tau$  plane.

### A Derivation of formulas in section 2.1

On the Fourier side:

$$\begin{cases} \hat{u}_{tt} + \mu|\xi|\hat{u} = 0, & (\xi, t) \in \mathbb{R} \times \mathbb{R}^+, \\ \hat{u}(\xi, 0) = 0, \\ \hat{u}_t(\xi, 0) = 1, \end{cases} \quad \begin{cases} \hat{u}_{tt} + \mu|\xi|\hat{u} = 0, & (\xi, t) \in \mathbb{R} \times \mathbb{R}^+, \\ \hat{u}(\xi, 0) = 1, \\ \hat{u}_t(\xi, 0) = 0. \end{cases} \quad (\text{A.1})$$

The two equations in (A.1) are second order ODE initial value problems, thus the unique solutions are respectively

$$\hat{G}(\xi, t) = \frac{\sin(\sqrt{\mu|\xi|}t)}{\sqrt{\mu|\xi|}}, \quad \hat{F}(\xi, t) = \cos(\sqrt{\mu|\xi|}t).$$

Performing inverse Fourier transform on  $\hat{G}, \hat{F}$  derives the solutions to equations in (2.6), namely the Green function  $G(x, t)$  and its time derivative  $F(x, t)$ :

$$G(x, t) = \int_{\mathbb{R}} \frac{\sin(\sqrt{\mu|\xi|}t)}{\sqrt{\mu|\xi|}} e^{2\pi i \xi x} d\xi, \quad F(x, t) = \int_{\mathbb{R}} \cos(\sqrt{\mu|\xi|}t) e^{2\pi i \xi x} d\xi. \quad (\text{A.2})$$

With some computation,  $G(x,t), F(x,t)$  can be represented by Fresnel integral  $C(x)$  and  $S(x)$ , which are defined in the following way:

$$C(x) = \int_0^x \cos(z^2) dz, \quad S(x) = \int_0^x \sin(z^2) dz. \quad (\text{A.3})$$

The expression of  $G(x,t)$  is so complicated that we have to omit it here, but it is similar to the one of  $F(x,t)$ , which is

$$F(x,t) = \frac{\sqrt{\mu t} \left( C \left( \sqrt{\frac{\mu t^2}{2\pi|x|}} \right) + \sin \left( \frac{\mu t^2}{4|x|} \right) S \left( \sqrt{\frac{\mu t^2}{2\pi|x|}} \right) \right)}{|x|^{3/2}}. \quad (\text{A.4})$$

After deriving  $G(x,t)$  and  $F(x,t)$ , we can now move onto the following general case,

Formula (2.5) can be derived in the following way. On the Fourier side, (2.4) turns into

$$\begin{cases} \hat{u}_{tt} + \mu|\xi|\hat{u} = 0, & (\xi, t) \in \mathbb{R} \times \mathbb{R}^+, \\ \hat{u}(\xi, 0) = \hat{f}(\xi), \\ \hat{u}_t(x, 0) = \hat{g}(\xi). \end{cases} \quad (\text{A.5})$$

$\hat{f}(\xi), \hat{g}(\xi)$  are the Fourier Transform of  $f(x), g(x)$  respectively. Recall that  $\hat{G}(\xi, t), \hat{F}(\xi, t)$  in (2.7) are solutions to the two ODEs in (A.1) respectively. Therefore, by principle of superposition, the solution of (A.5) can be written as

$$\hat{u}(\xi, t) = \hat{f}(\xi)\hat{F}(\xi, t) + \hat{g}(\xi)\hat{G}(\xi, t). \quad (\text{A.6})$$

Remember that Fourier Transform turns convolution into multiplication, namely the Fourier Transform of  $f * g$  is  $\hat{f}\hat{g}$  exactly. Hence, if we perform Inverse Fourier Transform, we can get that

$$u(x, t) = f(x) * F(x, t) + g(x) * G(x, t),$$

which is (2.5).

## B Proof of Lemma 3.1

The proof of this lemma employs the following lemma in [14] which provides the explicit formula of  $f(\tau, \nu)$ :

**Lemma B.1.** Denote  $e$  as vector of ones. For the linear test equation  $y' = \lambda y$ , the RK method applied to this equation reduces to  $y_{n+1} = f(\tau\lambda)y_n$ , with  $f(z)$  given by

$$f(z) = 1 + z\mathbf{w}^T (\mathbf{I} - z\mathbf{G})^{-1} \mathbf{e} = \frac{\det(\mathbf{I} - z\mathbf{G} + z\mathbf{e}\mathbf{w}^T)}{\det(\mathbf{I} - z\mathbf{G})}. \quad (\text{B.1})$$

Here  $\mathbf{G}$  is the Runge-Kutta matrix and  $\mathbf{w}$  are the weights.

By Lemma B.1, we have

$$f(\tau, \nu) = \frac{\det(\mathbf{I} - \tau\lambda_1 \mathbf{G} + \tau\lambda_1 \mathbf{e}\mathbf{w}^T)}{\det(\mathbf{I} - \tau\lambda_1 \mathbf{G})}. \quad (\text{B.2})$$

Remember that  $\lambda_1 = i\sqrt{\nu}$ , thus a direct computation yields

$$\begin{aligned} |\det(\mathbf{I} - \tau\lambda_1 \mathbf{G} + \tau\lambda_1 \mathbf{e}\mathbf{w}^T)| &= |\det(\mathbf{I} - i\tau\sqrt{\nu}(\mathbf{G} - \mathbf{e}\mathbf{w}^T))| \\ &= \sqrt{|\det(\mathbf{I} - i\tau\sqrt{\nu}(\mathbf{G} - \mathbf{e}\mathbf{w}^T))| |\det(\mathbf{I} + i\tau\sqrt{\nu}(\mathbf{G} - \mathbf{e}\mathbf{w}^T))|} \\ &= \sqrt{|\det((\mathbf{I} - i\tau\sqrt{\nu}(\mathbf{G} - \mathbf{e}\mathbf{w}^T))(\mathbf{I} + i\tau\sqrt{\nu}(\mathbf{G} - \mathbf{e}\mathbf{w}^T)))|} \\ &= \sqrt{|\det(\mathbf{I} + \tau^2\nu(\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)|}, \end{aligned}$$

and  $|\det(\mathbf{I} - \tau\lambda_1 \mathbf{G})| = \sqrt{|\det(\mathbf{I} + \tau^2\nu\mathbf{G}^2)|}$  holds for the same reason. Therefore,

$$|f(\tau, \nu)| = \frac{|\det(\mathbf{I} - \tau\lambda_1 \mathbf{G} + \tau\lambda_1 \mathbf{e}\mathbf{w}^T)|}{|\det(\mathbf{I} - \tau\lambda_1 \mathbf{G})|} = \sqrt{\frac{|\det(\mathbf{I} + \tau^2\nu(\mathbf{G} - \mathbf{e}\mathbf{w}^T)^2)|}{|\det(\mathbf{I} + \tau^2\nu\mathbf{G}^2)|}}.$$

## References

- [1] D. J. Acheson. *Elementary Fluid Dynamics*, 1991.
- [2] A. Al-Ghosoun, A. S. Osman, and M. Seaid. A computational model for simulation of shallow water waves by elastic deformations in the topography. *Communications in Computational Physics*, 2021.
- [3] G. R. Baker, D. I. Meiron, and S. A. Orszag. Generalized vortex methods for free-surface flow problems. *J. Fluid Mech.*, 123:477–501, 1982.
- [4] W. Bao, S. Jin, and P. A. Markowich. On time-splitting spectral approximations for the Schrödinger equation in the semiclassical regime. *Journal of Computational Physics*, 175(2):487–524, 2002.
- [5] J. T. Beale, T. Y. Hou, and J. S. Lowengrub. Growth rates for the linearized motion of fluid interfaces away from equilibrium. *Communications on Pure and Applied Mathematics*, 46(9):1269–1301, 1993.
- [6] J. T. Beale, T. Y. Hou, and J. S. Lowengrub. Convergence of a boundary integral method for water waves. *SIAM J. Numer. Anal.*, 33(5):1797–1843, 1996.
- [7] L. A. Caffarelli and A. Vasseur. Drift diffusion equations with fractional diffusion and the quasi-geostrophic equation. *Annals of Mathematics*, pages 1903–1930, 2010.
- [8] D. Chalikov and D. Sheinin. Modeling extreme waves based on equations of potential flow with a free surface. *J. Comput. Phys.*, 210(1):247–273, 2005.
- [9] Y. Cheng, H. Dong, M. Li, and W. Xian. A high order central DG method of the two-layer shallow water equations. *Communications in Computational Physics*, 28(4):1437–1463, 2020.
- [10] F. Dias and T. J. Bridges. The numerical computation of freely propagating time-dependent irrotational water waves. *Fluid Dynamics Research*, 38(12):803–830, 2006.
- [11] W.-P. Düll. On the mathematical description of water waves. *arXiv preprint arXiv:1612.06242*, 2016.

- [12] A. I. Dyachenko, V. E. Zakharov, and E. A. Kuznetsov. Nonlinear dynamics of the free surface of an ideal fluid. *Plasma Phys. Rep.*, 22(10):829–840, 1996.
- [13] A. I. Dyachenko, E. A. Kuznetsov, M. D. Spector, and V. E. Zakharov. Analytical description of the free surface dynamics of an ideal fluid (canonical formalism and conformal mapping). *Phys. Lett. A*, 221(1):73–79, 1996.
- [14] E. Hairer and G. Wanner. Solving ordinary differential equations II. Stiff and differential-algebraic problems, 1996.
- [15] S. Jin, C. Sparber, and Z. Zhou. On the classical limit of a time-dependent self-consistent field system: Analysis and computation. *arXiv preprint arXiv:1406.3810*, 2014.
- [16] S. Jin and Z. Zhou. A semi-Lagrangian time splitting method for the Schrödinger equation with vector potentials. *Communications in Information and Systems*, 13(3):247–289, 2013.
- [17] A. Kiselev, F. Nazarov, and A. Volberg. Global well-posedness for the critical 2D dissipative quasi-geostrophic equation. *Inventiones Mathematicae*, 167(3):445–453, 2007.
- [18] R. J. LeVeque. *Finite Difference Methods for Ordinary and Partial Differential Equations: Steady-State and Time-Dependent Problems*. SIAM, 2007.
- [19] J.-G. L. and R. L. Pego. On local singularities in ideal potential flows with free surface. *Chinese Annals of Mathematics, Series B*, 40(6):925–948, 2019.
- [20] J.-G. L. and R. L. Pego. In search of local singularities in ideal potential flows with free surface. *arXiv preprint arXiv:2108.00445*, 2021.
- [21] Z. Ma, Y. Zhang, and Z. Zhou. An improved semi-Lagrangian time splitting spectral method for the semi-classical Schrödinger equation with vector potentials using NUFFT. *Applied Numerical Mathematics*, 111:144–159, 2017.
- [22] X. Meng, T.-T.-P. Hoang, Z. Wang, and L. Ju. Localized exponential time differencing method for shallow water equations: Algorithms and numerical study. *Communications in Computational Physics*, 29(1), 2020.
- [23] L. M. Milne-Thomson. *Theoretical Hydrodynamics*. 4th ed. The Macmillan Co., New York, 1960.
- [24] C. E. Shannon. Communication in the presence of noise. *Proceedings of the IEEE*, 72(9):1192–1201, 1984.
- [25] G. I. Taylor. The instability of liquid surfaces when accelerated in a direction perpendicular to their planes. I. In *P. Roy. Soc. Lond. A Mat.*, volume 201, pages 192–196. The Royal Society, 1950.
- [26] L. N. Trefethen. *Spectral Methods in MATLAB*. SIAM, 2000.
- [27] M. R. Turner and T. J. Bridges. Time-dependent conformal mapping of doubly-connected regions. *Adv. Comput. Math.*, 42(4):947–972, 2016.
- [28] S. Wu. Well-posedness in Sobolev spaces of the full water wave problem in 2-d. *Invent. Math.*, 130(1):39–72, 1997.