

STRUCTURE-PRESERVING FINITE DIFFERENCE METHODS FOR LINEARLY DAMPED  
DIFFERENTIAL EQUATIONS

by

ASHISH BHATT

M.S. University of Central Florida, 2013

M.Sc. Indian Institute of Technology Dhanbad, 2009

B.Sc. Hemwati Nandan Bahuguna Garhwal University, 2007

A dissertation submitted in partial fulfilment of the requirements  
for the degree of Doctor of Philosophy  
in the Department of Mathematics  
in the College of Sciences  
at the University of Central Florida  
Orlando, Florida

Fall Term  
2016

Major Professor: Brian E. Moore

© 2016 Ashish Bhatt

## ABSTRACT

Differential equations (DEs) model a variety of physical phenomena in science and engineering. Many physical phenomena involve conservative or dissipative forces, which manifest themselves as qualitative properties of DEs that govern these phenomena. Since only a few and simplistic models are known to have exact solutions, approximate solution techniques, such as numerical integration, are used to reveal important insights about solution behavior and properties of these models. Numerical integrators generally result in undesirable quantitative and qualitative errors. Standard numerical integrators aim to reduce quantitative errors, whereas geometric (numerical) integrators aim to reduce or eliminate qualitative errors, as well, in order to improve the accuracy of numerical solutions. It is now widely recognized that geometric (or structure-preserving) integrators are advantageous compared to non-geometric integrators for DEs, especially for long time integration.

Geometric integrators for conservative DEs have been proposed, analyzed, and investigated extensively in the literature. The motif of this thesis is to extend the idea of structure preservation to linearly damped DEs. More specifically, we develop, analyze, and implement geometric integrators for linearly damped ordinary and partial differential equations (ODEs and PDEs) that possess conformal invariants, which are qualitative properties that decay exponentially along any solution of the DE as the system evolves over time. In particular, we derive restrictions on the coefficient functions of exponential Runge-Kutta (ERK) numerical methods for preservation of certain conformal invariants of linearly damped ODEs. An important class of these methods is shown to preserve the damping rate of solutions of damped linear ODEs. Linearly stability and order of accuracy for some specific cases of ERK methods are investigated. Geometric integrators for PDEs are designed using structure-preserving ERK methods in space, time, or both. These integrators for PDEs are also shown to preserve additional structure in certain special cases. Numerical ex-

periments illustrate higher order accuracy and structure preservation properties of various ERK based methods, demonstrating clear advantages over non-structure-preserving methods, as well as usefulness for solving a wide range of DEs.

To my parents

## ACKNOWLEDGMENTS

This work is the culmination of research done by me and my collaborators over the course of my graduate studies. During this period, many have extended their support towards this endeavor. It would not have been as fulfilling and enduring without the support of all these people. I must thank them now.

The most instrumental person in helping me throughout the preparation and completion of this dissertation has been my Ph.D. advisor Dr. Brian E. Moore. He has been a fountainhead of original ideas and has augmented and enriched my own ideas. He has also helped open doors for me to the academic community by introducing me to other experts in the discipline and getting me invited to conference talks. I am grateful for his guidance, supervision, and punctiliousness toward my writing. Parts of this thesis were completed when I visited Norges teknisk-naturvitenskapelige universitet (NTNU), Norway. I am thankful to Dr. Elena Celledoni and Dr. Brynjulf Owren for the invitation to visit NTNU and stimulating discussions.

I thank Dr. S. R. Choudhury and Dr. G. S. Oztepe for a fruitful collaboration. I would also like to thank Dr. S. R. Choudhury, Dr. Basak Gurel, and Dr. Jeffrey Kauffman for serving on my committee and carrying out duties that it entails. I also wish to express my gratitude to all the past mathematicians whose founding principles make the bedrock of this thesis and whose terminology, nomenclature, and notations I have invariably inherited.

I am thankful to Haider, Sulalit, Elliot, Cheng, Mangalagama, Aritra, Arielle, Pawan, and Ted for sitting through my practice talks and proofreading my drafts. I also thank Pawan for making time for so many memorable outdoor adventures that we did together. Thank you to Sumit for our continuing friendship and his unconditional support during all these years.

# TABLE OF CONTENTS

LIST OF FIGURES . . . . .	x
LIST OF TABLES . . . . .	xiii
CHAPTER 1: INTRODUCTION . . . . .	1
1.1 Damped differential equations . . . . .	2
1.2 Finite difference methods . . . . .	8
1.2.1 Finite difference operators and their properties . . . . .	11
1.3 Structure preservation background . . . . .	15
1.4 Outline . . . . .	20
CHAPTER 2: STRUCTURE-PRESERVING EXPONENTIAL RUNGE-KUTTA METH- ODS FOR ODES . . . . .	21
2.1 ERK and partitioned ERK methods . . . . .	22
2.2 Preservation of conformal invariants . . . . .	28
2.3 Preservation of conformal symplecticness . . . . .	34
2.4 Accuracy and stability of ERK and PERK methods . . . . .	40

CHAPTER 3: ODE APPLICATIONS AND EXPERIMENTS . . . . .	48
3.1 Linear oscillators . . . . .	48
3.1.1 Constant damping . . . . .	49
3.1.2 Time-dependent damping . . . . .	51
3.2 Damped pendulum . . . . .	52
3.2.1 Damped pendulum . . . . .	52
3.2.2 Damped driven pendulum . . . . .	54
3.3 N-body ODE . . . . .	57
3.4 Rigid body with periodic perturbation . . . . .	59
CHAPTER 4: STRUCTURE-PRESERVING METHODS FOR PDES . . . . .	61
4.1 Multi-conformal-symplectic PDEs . . . . .	61
4.1.1 Local conservation laws . . . . .	64
4.1.2 Multi-conformal-symplectic numerical methods . . . . .	67
4.2 Non-standard finite difference methods . . . . .	70
CHAPTER 5: PDE APPLICATIONS AND EXPERIMENTS . . . . .	72
5.1 A damped Klein-Gordon equation . . . . .	72
5.1.1 Numerical solutions . . . . .	73



5.1.2	Structure-preservation . . . . .	75
5.2	A Modified Burgers' Equation . . . . .	78
5.2.1	Numerical solutions . . . . .	79
5.2.2	Structure-preservation . . . . .	82
5.3	Damped driven nonlinear Schrödinger equation . . . . .	84
5.3.1	Numerical solutions . . . . .	86
5.3.1.1	Integrating factor method . . . . .	87
5.3.1.2	Exponential time differencing method . . . . .	90
5.3.1.3	Implicit midpoint method . . . . .	91
5.3.2	Numerical results . . . . .	92
5.3.2.1	Linear Schrödinger equation . . . . .	93
5.3.2.2	Damped NLS . . . . .	94
5.3.2.3	Damped driven NLS . . . . .	96
CHAPTER 6: CONCLUSIONS AND FUTURE WORK . . . . .		99
APPENDIX A: DIFFERENTIAL FORMS AND THE WEDGE PRODUCT . . . . .		101
LIST OF REFERENCES . . . . .		106

## LIST OF FIGURES

1.1	Solutions of a PDE without and with damping (Left to right). . . . .	18
1.2	Deformation of the phase space volume (left) and the corresponding flow (right), which corresponds to the flow of a differential equation [41]. Blue and red boxes represent the initial and current phase space volumes, respectively. Blue lines denote solution trajectories of the DE. . . . .	19
3.1	A comparison of the average absolute solution error for the conformal symplectic Euler methods given in (2.18) and (2.19) for solving eq. (3.1) with $\gamma = 0.01$ . Initial condition: $\theta(0) = 0, \omega(0) = 10$ ; final time: $T = 50$ . . . . .	50
3.2	Local absolute error in solution over the step size for IFRK methods of stages 1, 2 and 3 applied to ODE (3.1). Dashed lines represent the slopes with which they are labeled. . . . .	50
3.3	Error $E_n$ (3.2) in conformal symplecticness for the GL-IFRK methods (left) and the standard Gauss-Legendre methods (right) applied to eq. (3.1) with $\gamma(t) = \frac{1}{2}\epsilon \cos(2t)$ . . . . .	51
3.4	The residual (3.4) for three numerical solutions of (3.3). The methods used are IFRK (2.11); ETDRK (2.13); and IFPRK (2.20) denoted here by IFRK, ETDRK, and IFRK, respectively. Left: rapid oscillation with imaginary $\gamma$ ; Right: strong damping with real $\gamma$ . . . . .	53

3.5	Left to right: time series, phase space and Poincare sections of damped driven oscillator, eq. (3.5), with the parameter values mentioned in the title. $T$ is the final time. CIMP and Heun's stand for eqs. (3.6) and (3.8) . . . . .	56
3.6	Left to right: Error in linear momentum, error in angular momentum and corresponding solution trajectories of N-body system. . . . .	58
3.7	Casimir and energy errors (3.10) for simulations of the system (1.7). Left: GL-IFRK methods; Right: standard Gauss-Legendre methods. . . . .	60
5.1	Error in the solution of (4.3) due to methods (5.3), (5.4) and (5.5). Parameter values are given in the figure title. The maximum value of the exact solution at time $T = 50$ is approximately $7 \times 10^{-9}$ . . . . .	74
5.2	Drift in the rate of dissipation for the three methods (5.3), (5.4) and (5.5) with the parameter values mentioned in the figure title. Only every sixth drift vector component is plotted for clarity and CSV1 eclipses CSV2. . . . .	75
5.3	Total conformal momentum $I^i$ and residual $r^i$ due to (5.3), (5.4) and (5.5) with the parameters mentioned in the figure title . . . . .	78
5.4	Snapshots of the numerical solution of eq. (4.5) using (5.10), (5.12) and (4.18) at different times. . . . .	81
5.5	Residual (5.15) due to conformal symplectic methods (5.10) and (5.12) and NSFD (4.18). . . . .	83

5.6	Plane wave solution, momentum, and norm and energy residuals. The second column gives $\mathbf{I}_1$ because $\mathbf{R}_1$ is undefined when the $x$ -derivative of the solution is zero. . . . .	93
5.7	$L^\infty$ error due to the methods of eqs. (4.15) and (5.27). . . . .	95
5.8	Plane wave solution, momentum, and invariant residual. . . . .	95
5.9	Soliton collision and invariant residuals. For IF and ETD methods, residual $\mathbf{R}_1$ is close to machine precision except near the time of collision when solution profile is steep at the spatial location of the collision. . . . .	96
5.10	Periodic and chaotic attractors of damped driven NLS along with imaginary versus real parts of numerical solution at all times and $x = 0$ . . . . .	98

## LIST OF TABLES

1.1	Equations and some of their conformal invariants, under suitable boundary conditions where applicable. Setting $\gamma = 0$ in an equation gives the conservative counterpart of that equation because the corresponding conformal invariants become integral invariants. . . . .	8
3.1	Total linear momentum and total angular momentum for the three methods. Here $q_i^n \approx q_i(t_n)$ etc. . . . .	58

## CHAPTER 1: INTRODUCTION

Some examples of physical phenomena in science and engineering that are governed by DEs include rigid body problem, N-body problem, water and sound wave propagation, non-relativistic quantum mechanics, and superconductivity. From the perspective of structure preservation, DEs can be classified in two broad categories: conservative and damped DEs. While conservative DEs have their own importance, damped DEs are also important in applications because of the presence of resistive or attenuating forces in physical systems governed by the damped DEs. In this thesis, our focus will be on the latter category. Damped DEs are characterized by possession of qualitative properties that decay along any solution. Those qualitative properties that decay exponentially along any solution are referred to as *conformal invariants* and will be defined more precisely later. Both ordinary and partial differential equations (ODEs and PDEs) that possess conformal invariants are considered in this exposition.

Since not all DEs are amenable to exact solutions, approximate solution techniques are indispensable. Numerical methods are one of the approximate solution techniques used to solve DEs. Finite difference, finite element, finite volume, and spectral methods are types of numerical methods used to solve DEs numerically. Finite difference methods are the earliest numerical methods used among all the numerical methods and their structure-preservation properties for conservative DEs are well known. In this thesis, our focus will be on establishing structure-preservation properties of finite difference methods for linearly damped DEs.

In the next section, we define conformal invariants and discuss some motivating examples of DEs and their conformal invariants. In Section 1.2, we discuss fundamentals of finite difference methods and properties of finite difference operators that will be used later in this thesis to design structure-preserving numerical methods. Then we discuss some of the previous work done in the

direction of structure-preservation in Section 1.3. We conclude this chapter with an outline of the rest of the thesis in Section 1.4.

## 1.1 Damped differential equations

DEs, whose qualitative properties, such as energy or momentum, remain constant along any solution, are referred to as conservative DEs. In contrast, DEs whose qualitative properties decay along any solution are referred to as damped differential equations. Other commonly used names for damped DEs are dissipative DEs or non-conservative DEs. This decay in the solution or qualitative properties of a DE is often the result of the presence of resistive forces in the system that is being modeled by the DE.

Consider the Cauchy problem

$$\dot{z}(t) = N(z(t)) - \gamma(t)z(t), \quad z(0) = z_0 \tag{1.1}$$

where  $z \in \mathbb{R}^d$  with  $d \in \mathbb{N}$ ,  $N : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a smooth nonlinear function of  $z$ , and  $\dot{z}$  denotes the derivative of  $z$  with respect to  $t$ . We require that  $\gamma$  be scalar, and we allow it to depend on time, i.e.  $\gamma(t) : \mathbb{R} \rightarrow \mathbb{R}$ . The DE in eq. (1.1) is a generalization of a more prevalent special case

$$\dot{z}(t) = N(z(t)) - \gamma_0 z(t),$$

where  $\gamma_0$  is a real constant. The term involving  $\gamma_0$  is linear and is often responsible for damping in the system. In this thesis, we consider the generalization (1.1) of this linearly damped system. The solution  $z$  of (1.1) can be thought of a map taking the initial condition  $z_0$  to a later point  $z(t)$  after time  $t$  along a solution trajectory. To emphasize this dependence of the solution on the initial condition, one often writes  $z = z(t; z_0)$ .

When  $d$  is even, the solution  $z(t)$  of the system (1.1) can be partitioned into two vector variables of dimension  $d/2 \times 1$ . The system thus obtained in terms of these new variables is referred to as a partitioned system. Now, suppose the partitioned system thus obtained from (1.1) has the form

$$\begin{bmatrix} \dot{q} \\ \dot{p} \end{bmatrix} = \begin{bmatrix} F(q, p) \\ G(q, p) \end{bmatrix} - \begin{bmatrix} \gamma^1(t)q \\ \gamma^2(t)p \end{bmatrix}, \quad \begin{bmatrix} q(0) \\ p(0) \end{bmatrix} = \begin{bmatrix} q_0 \\ p_0 \end{bmatrix} \quad (1.2)$$

where  $q, p \in \mathbb{R}^{d/2}$  with  $d$  even, and the functions  $F, G : \mathbb{R}^d \rightarrow \mathbb{R}^{d/2}$  and  $\gamma^k : \mathbb{R} \rightarrow \mathbb{R}$  for  $k = 1, 2$  are smooth. We have suppressed the dependence of the variables  $q, p$  on  $t$ .

The following definition of a conformal invariant is of fundamental importance to our discussion.

**Definition 1.1.** A non constant function  $\mathcal{I} : \mathbb{R}^d \rightarrow \mathbb{R}$  is a conformal invariant of eq. (1.1) if

$$\frac{d}{dt}\mathcal{I}(z) = -2\gamma(t)\mathcal{I}(z) \quad (1.3)$$

for all  $z = z(t, z_0)$ , all  $z_0 \in \mathbb{R}^d$ , and all  $t \in \mathbb{R}$ . Similarly, a non constant function  $\mathcal{I}(t) : \mathbb{R}^{2d} \rightarrow \mathbb{R}$  is a conformal invariant of (1.2) if

$$\frac{d}{dt}\mathcal{I} = -(\gamma^1(t) + \gamma^2(t))\mathcal{I},$$

for all  $q = q(t, q_0, p_0), p = p(t, q_0, p_0)$ , all  $(q_0, p_0) \in \mathbb{R}^{2d}$ , and all  $t \in \mathbb{R}$

Notice that eq. (1.3) is equivalent to

$$\frac{d}{dt} \left( e^{\int_0^t 2\gamma(s)ds} \mathcal{I}(t) \right) = 0 \quad \iff \quad \mathcal{I}(t) = e^{-\int_0^t 2\gamma(s)ds} \mathcal{I}(0),$$

where  $\mathcal{I}(t) = \mathcal{I}(z(t))$ . This last equation means that conformal invariants decay exponentially along all solutions when  $\gamma$  is a constant. If  $\gamma = 0$ , then the function  $\mathcal{I}$  remains unchanged along all solutions of eq. (1.1) and is referred to as a *first integral, constant of motion, or conserved quantity* of the equation. Similar statements are true for the conformal invariant of the partitioned system. The next chapter has more details on conformal invariants and their preservation. For now, let



us motivate the discussion by giving the following examples of non-conservatively perturbed DEs and their conformal invariants.

**Example 1.2.** Consider the following system

$$\begin{aligned}\dot{\theta} &= \omega, \\ \dot{\omega} &= -\kappa^2\theta - 2\gamma\omega.\end{aligned}\tag{1.4}$$

Notice that this is just the governing equation of a damped oscillator where  $\theta$  is the displacement,  $\kappa$  is the frequency, and  $\gamma$  is the damping parameter. Here, overdot denotes the time derivative. Now defining

$$H_\gamma = \frac{1}{2}(\kappa^2\theta^2 + \omega^2) + \gamma\theta\omega$$

and differentiating with respect to  $t$  gives

$$\begin{aligned}\frac{d}{dt}H_\gamma &= \frac{d}{dt} \left( \frac{1}{2}(\kappa^2\theta^2 + \omega^2) + \gamma\theta\omega \right) \\ &= \kappa^2\theta\dot{\theta} + \omega\dot{\omega} + \gamma(\theta\dot{\omega} + \omega\dot{\theta}) \\ &= (\kappa^2\theta + \gamma\omega)\dot{\theta} + (\omega + \gamma\theta)\dot{\omega} \\ &= (\kappa^2\theta + \gamma\omega)\omega + (\omega + \gamma\theta)(-\kappa^2\theta - \gamma\omega) \\ &= -\gamma\omega^2 - \gamma\kappa^2\theta^2 - 2\gamma^2\theta\omega \\ &= -2\gamma \left( \frac{1}{2}(\kappa^2\theta^2 + \omega^2) + \gamma\theta\omega \right) \\ &= -2\gamma H_\gamma.\end{aligned}$$

Here, we have used system (1.4) to replace time derivatives of  $\theta$  and  $\omega$ . Therefore  $H_\gamma$  is a conformal invariant of the system (1.4). With  $\gamma = 0$ , eq. (1.4) reduces to a conservative harmonic oscillator with energy  $H_0 = H_\gamma|_{\gamma=0}$ . Notice that the energy  $H_0$  of the conservative harmonic oscillator remains unchanged along all solutions, and, hence,  $H_0$  is a first integral.

**Example 1.3.** Consider the following equations of motion for an n-body system

$$\partial_t q_i = \frac{1}{m_i} p_i, \quad (1.5)$$

$$\partial_t p_i = - \sum_{j \neq i} \tau_{ij} (q_i - q_j) - 2\gamma p_i \quad (1.6)$$

for  $i = 1, 2, \dots, N$ ; where  $m_i$  is mass of the  $i^{\text{th}}$  particle,  $\phi_{ij}(\|q_i - q_j\|)$  is the interaction potential (pair-potential) between particles  $i$  and  $j$  at the distance  $\|q_i - q_j\|$ , and

$$\tau_{ij} = \frac{\phi'_{ij}(\|q_i - q_j\|)}{\|q_i - q_j\|}.$$

Here  $\partial_t$  denotes the time derivative. Vectors  $q_i \in \mathbb{R}^3$  and  $p_i \in \mathbb{R}^3$  denote position and linear momentum, respectively, of the  $i^{\text{th}}$  particle. Taking the cross product of (1.5) and (1.6) with  $p_i$  and  $q_i$ , respectively, we have

$$\begin{aligned} \partial_t q_i \times p_i &= 0, \\ \partial_t p_i \times q_i &= - \sum_{j \neq i} \tau_{ij} (-q_j \times q_i) - 2\gamma p_i \times q_i. \end{aligned}$$

Now, summing the second equation over  $i$  gives

$$\begin{aligned} \sum_i \partial_t p_i \times q_i &= - \sum_i \sum_{j \neq i} \tau_{ij} (-q_j \times q_i) - 2\gamma \sum_i p_i \times q_i \\ &= -2\gamma \sum_i p_i \times q_i, \end{aligned}$$

which gives

$$\partial_t \left( \sum_i p_i \times q_i \right) = -2\gamma \sum_i p_i \times q_i.$$

Similarly summing equation (1.6) over  $i$  we arrive at

$$\sum_i \partial_t p_i = -2\gamma \sum_i p_i, \quad \iff \quad \sum_i p_i(t) = e^{-2\gamma t} \sum_i p_i(0).$$

Therefore, total angular momentum  $\sum_i p_i(t) \times q_i(t)$  and total linear momentum  $\sum_i p_i(t)$  are con-

formal invariants for the system of eqs. (1.5) and (1.6). Notice that for  $\gamma = 0$ , this simply means that the total linear and total angular momentum are first integrals of the n-body problem given by eqs. (1.5) and (1.6) with  $\gamma = 0$ .

**Example 1.4.** Consider the following system of equations

$$\begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \\ \dot{z}_3 \end{bmatrix} = \begin{bmatrix} 0 & z_3/I_3 & -z_2/I_2 \\ -z_3/I_3 & 0 & z_1/I_1 \\ z_2/I_2 & -z_1/I_1 & 0 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} - \gamma(t) \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}, \quad (1.7)$$

where  $I_1, I_2, I_3$  are nonzero real constants and  $\gamma(t) = \frac{\epsilon}{2} \cos(2t)$  is a time dependent damping term. When  $\epsilon = 0$ , this system defines the motion of a free rigid body with center of mass at the origin, the solution vector  $z = (z_1, z_2, z_3)^T$  represents the angular momentum, and  $I_1, I_2, I_3$  are principal moments of inertia. It is straightforward to show that the system has two conformal invariants, one for the Casimir

$$\frac{dC}{dt} = -2\gamma(t)C \quad \text{with Casimir} \quad C(z) = z_1^2 + z_2^2 + z_3^2,$$

and one for the energy

$$\frac{dH}{dt} = -2\gamma(t)H \quad \text{with energy} \quad H(z) = \frac{1}{2} \left( \frac{z_1^2}{I_1} + \frac{z_2^2}{I_2} + \frac{z_3^2}{I_3} \right).$$

**Example 1.5.** Consider the nonlinear Schrödinger equation

$$i\psi_t + \psi_{xx} + V'(|\psi|^2)\psi + 2i\gamma\psi = 0$$

where  $\psi = \psi(x, t)$  is a complex valued wave function of space  $x$  and time  $t$ , the nonnegative real number  $\gamma$  is the damping parameter, and subscripts denote the usual partial derivatives. The equation models a variety of physical phenomena including propagation of the envelop of modulated water wave groups. To show that the equation has a conformal invariant, let us multiply the PDE

by  $\bar{\psi}$ , the complex conjugate of  $\psi$ , and integrate to get

$$i \int \psi_t \bar{\psi} dx + \int \psi_{xx} \bar{\psi} dx + \int V'(|\psi|^2) |\psi|^2 dx + 2i\gamma \int |\psi|^2 dx = 0. \quad (1.8)$$

After integration by parts, the second term of this equation becomes

$$\int \psi_{xx} \bar{\psi} dx = [\bar{\psi} \psi_x] - \int |\psi_x|^2 dx$$

where  $[\cdot]$  denotes difference of the enclosed function evaluated at the upper and lower limit of integration. This difference vanishes under appropriate boundary conditions and hence eq. (1.8) becomes

$$i \int \psi_t \bar{\psi} dx - \int |\psi_x|^2 dx + \int V'(|\psi|^2) |\psi|^2 dx + 2i\gamma \int |\psi|^2 dx = 0 \quad (1.9)$$

Taking the imaginary part of this equation we get the linear ODE

$$\partial_t \int |\psi|^2 dx + 4\gamma \int |\psi|^2 dx = 0,$$

which implies

$$\int |\psi(x, t)|^2 dx = e^{-4\gamma t} \int |\psi(x, 0)|^2 dx$$

i.e. the norm  $\int |\psi|^2 dx$  of the solution decays exponentially along solutions of the PDE, or the norm is a conformal invariant.

Some of the above and other examples of damped DEs, along with corresponding conformal invariants, are given in Table 1.1. Conformal Hamiltonian ODE and its conformal invariant in the table are discussed in detail in the next chapter. Notice that for all the examples of the table, setting  $\gamma = 0$  renders the damped DEs conservative and corresponding conformal invariants become constants of motion i.e. they remain unchanged along all solutions. Much research has been done to develop numerical methods that preserve constants of motion of a DE. On the other hand, preservation of conformal invariants is a comparatively less researched area but important nonetheless

because of its physical implications. The main motif of this thesis is to develop numerical methods that preserve conformal invariants, such as those in Table 1.1.

Table 1.1: Equations and some of their conformal invariants, under suitable boundary conditions where applicable. Setting  $\gamma = 0$  in an equation gives the conservative counterpart of that equation because the corresponding conformal invariants become integral invariants.

Equation	Conformal Invariant
Damped harmonic oscillator $\ddot{\theta} + 2\gamma\dot{\theta} + \kappa^2\theta = 0$	$\mathcal{I} = \frac{1}{2}(\kappa^2\theta^2 + \dot{\theta}^2) + \gamma\theta\dot{\theta}$
Lorenz equations $\dot{x} = \sigma(y - x), \dot{y} = rx - y - xz, \dot{z} = xy - bz$	$\mathcal{I} = \int_{\Omega} dV$
Conformal Hamiltonian ODE $\dot{z} = \mathbf{J}^{-1}\nabla_z H(z) - \gamma(t)z$	$\mathcal{I} = \omega = dz \wedge \mathbf{J}dz$
Damped wave equation $u_{tt} - u_{xx} + cu + 2\gamma u_t = 0$	$\mathcal{I} = \int u_t u_x dx$
Damped KdV equation $u_t + uu_x + u_{xxx} + 2\gamma u = 0$	$\mathcal{I} = \int u dx$
Damped nonlinear Schrödinger equation $i\psi_t + \psi_{xx} + V'( \psi ^2)\psi + 2i\gamma\psi = 0$	$\mathcal{I} = \int  \psi ^2 dx$
Damped Camassa Holm equation $u_t - u_{xxt} + 3uu_x + \gamma(u - u_{xx}) = 2u_x u_{xx} + uu_{xxx}$	$\mathcal{I} = \int (u^2 + u_x^2) dx$

## 1.2 Finite difference methods

Methods used to find numerical solutions of differential equations are referred to as *numerical methods*. A numerical method is also referred to as a *scheme*, an *integrator*, or simply a *discretization* and we use all these terms interchangeably throughout this thesis. Finite difference methods make up a class of numerical methods which replaces terms of a continuous equation with finite difference operators.

**Remark.** Discretizing a continuous equation often comes at the cost of *quantitative and qualitative errors*. The quantitative error refers to the error introduced by a numerical method in approximating a solution of a DE, whereas the qualitative error is the error in approximating qualitative properties such as first integrals, conformal invariants, conservation laws, limit cycles, equilibrium points, periodic orbits, chaos, etc. of a DE. We want to develop numerical methods that preserve some qualitative properties. We build upon methods that preserve qualitative properties such as first integrals and conservation laws of conservative DEs. Our methods instead preserve qualitative properties such as conformal invariants and conformal conservation laws of linearly damped DEs.

By choosing appropriate finite difference operators, one is able to reduce or eliminate qualitative and quantitative errors. The next section puts qualitative properties and their preservation in perspective. In this section, we discuss a measurement of quantitative errors, present examples of finite difference methods, and establish properties of the operators used in numerical discretizations. To this end, let us recall that  $z(t; z_0)$  denotes the solution trajectory starting at the initial value  $z_0$  and describe the flow map of an ODE in the following definition.

**Definition 1.6.** *The map  $\psi_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is the flow map of the initial value problem*

$$\dot{z} = f(z), \quad z(0) = z_0 \tag{1.10}$$

*if*

$$\psi_t(z_0) = z(t; z_0), \quad z_0 \in \mathbb{R}^d,$$

*i.e.  $\psi_t$  takes initial data to later points along solution trajectories.*

Similar to the flow map of a continuous process, one can define the flow map of a discrete process. Let  $\Psi_h$  denote the flow map of a numerical method for (1.10) and  $\Psi_h(\bar{z})$  be the approximation of the solution  $z(h; \bar{z})$  through a given point  $\bar{z}$  of the phase space. The numerical approximation

$\Psi_h(\bar{z})$  is often not equal to the solution  $z(h; \bar{z})$  and the order of a method is a measure of the distance between the two, as given in the following definition.

**Definition 1.7.** *The order of a numerical one-step method  $\Psi_h$  is defined to be the largest integer  $p \geq 1$  such that*

$$\|\Psi_h(\bar{z}) - \psi_h(\bar{z})\| \leq Ch^{p+1},$$

for  $\bar{z}$  in the domain of interest, where  $C > 0$  is a constant.

In other words, order is a measure of the quantitative error due to the discretization. Since numerical solutions are only approximate in general, there is always some error (quantitative error) in the solutions. Evidently, numerical methods with less quantitative error (higher order methods) may be more desirable.

Some examples of finite difference methods are perhaps in order. Runge-Kutta (RK) methods for the differential equation

$$\dot{z}(t) = f(z, t), \tag{1.11}$$

where  $z \in \mathbb{R}^d$  with  $d \in \mathbb{N}$ ,  $f : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$  is a smooth function of  $z$ , are given by

$$\begin{aligned} Z_i &= z_n + h \sum_{j=1}^s \alpha_{ij} f(Z_j, t_n + c_j h), \quad i = 1, \dots, s, \\ z_{n+1} &= z_n + h \sum_{i=1}^s \beta_i f(Z_i, t_n + c_i h), \end{aligned} \tag{1.12}$$

where  $s$  is the number of *stages*,  $h$  denotes the step size, and  $t_n = nh$  for  $n = 0, 1, 2, \dots$ ,  $z_n$  is the numerical solution,  $Z_i$ 's are the *stage variables*, and  $\alpha_{ij}, \beta_i$  are referred to as *coefficients* of the methods. RK methods are succinctly represented by the Butcher tableau

$$\begin{array}{c|c} c & A \\ \hline & b^T \end{array}$$

where

$$c = \{c_i\}_{i=1}^s, \quad c_i = \sum_{j=1}^s \alpha_{ij}, \quad b = \{\beta_i\}_{i=1}^s, \quad \text{and} \quad A = \{\alpha_{ij}\}_{i,j=1}^s.$$

Notice that the flow map  $z_{n+1} = \Psi_h(z_n)$  of an RK method is not explicit in general and one has to employ fixed point iterations to compute an implicit flow map to obtain the numerical solution  $z_n$ .

For example, the explicit Euler method and the implicit midpoint method are given by

$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array} \quad \text{and} \quad \begin{array}{c|c} 1/2 & 1/2 \\ \hline & 1 \end{array}, \quad (1.13)$$

respectively. The explicit Euler method is an order one and the implicit midpoint method is an order two RK method. Indeed, the latter method is a type of *Gauss-Legendre* RK method. GL-RK methods of stage  $s$  are known to have order of accuracy equal to  $2s$ , highest possible order achievable by a stage  $s$  RK method. For a partitioned system, it may be desirable to apply two different RK methods on each part of the system. Methods thus obtained are referred to as partitioned RK (PRK) methods. Since we are on the topic of finite difference methods, it is efficient to discuss the operators used in these methods here.

### 1.2.1 Finite difference operators and their properties

It is often more convenient to write and work with a numerical method for a differential equation by writing the method more succinctly using discrete analogues of the continuous operators appearing in the equation. Here we introduce some of these discrete analogues and their properties to be used later. We begin by defining the following finite difference operators for  $\phi^k$ , a numerical



approximation of  $\phi(\zeta^k) = \phi(k\Delta\zeta)$ ,  $k \in \mathbb{Z}$ .

$$\begin{aligned} D_\zeta^\alpha \phi^k &= \frac{e^{\alpha\Delta\zeta}\phi^{k+1} - e^{-\alpha\Delta\zeta}\phi^k}{\Delta\zeta}, & A_\zeta^\alpha \phi^k &= \frac{e^{\alpha\Delta\zeta}\phi^{k+1} + e^{-\alpha\Delta\zeta}\phi^k}{2}, \\ T_\zeta \phi^k &= \phi^{k-1}, & \delta_\zeta \phi^k &= \frac{\phi^{k+1} - \phi^{k-1}}{2\Delta\zeta}, & \delta_\zeta^2 \phi^k &= \frac{\phi^{k+1} - 2\phi^k + \phi^{k-1}}{\Delta\zeta^2}. \end{aligned} \quad (1.14)$$

Depending on whether  $\zeta$  denotes space  $x$  or time  $t$ , the operators of eq. (1.14) are spatial or temporal operators, respectively. The operators  $D_\zeta^\alpha$  and  $A_\zeta^\alpha$  are often referred to as *discrete derivative* and *discrete averaging* operators, respectively. Usually, the superscript  $\alpha$  is a function of the damping parameter in the system being discretized and should not be confused with coefficients  $\alpha_{ij}$  of the RK methods (1.12). Thus these operators subsume part of the damping and distribute it evenly over a discrete computational mesh. This absorption and uniform distribution of damping has important ramifications which will be discussed in later chapters. When  $\alpha = 0$ , the derivative and averaging operators simply reduce to standard forward difference and forward averaging operators and are denoted by  $D_\zeta$  and  $A_\zeta$ , respectively. Operator  $T_\zeta$  is a shift operator whereas  $\delta_\zeta$  and  $\delta_\zeta^2$  are second order accurate finite difference approximations of first and second order derivatives. In general, we drop the superscript on  $\phi^k$  when using these operators for the sake of simplicity.

For example, the implicit midpoint method, given by the second tableau of eq. (1.13), for eq. (1.11) can be written in two different ways:

$$\frac{z_{n+1} - z_n}{h} = f\left(\frac{z_{n+1} + z_n}{2}, \frac{t_{n+1} + t_n}{2}\right),$$

or using the operators of eq. (1.14)

$$D_t z_n = f(A_t z_n, A_t t_n).$$

Among these two portrayals of the implicit midpoint method, the later one is more succinct and arguably easier to work with in light of the following lemma stating properties of the operators of eq. (1.14). The succinctness property of the discrete operators becomes even more worthwhile

for PDEs. The following lemma will be used frequently in the following chapters to prove certain properties of numerical methods.

**Lemma 1.8.** *The operators of eq. (1.14) have the following properties, [40, 5]:*

(i) *The derivative and averaging operators commute:*

$$D_\zeta^\alpha A_\eta^\beta \phi = A_\eta^\beta D_\zeta^\alpha \phi, D_\zeta^\alpha D_\eta^\beta \phi = D_\eta^\beta D_\zeta^\alpha \phi, A_\zeta^\alpha A_\eta^\beta \phi = A_\eta^\beta A_\zeta^\alpha \phi.$$

(ii) *They satisfy the following discrete product rule:*

$$D_\zeta^\alpha(\phi * \xi) = D_\zeta^{\alpha/2} \phi * A_\zeta^{\alpha/2} \xi + A_\zeta^{\alpha/2} \phi * D_\zeta^{\alpha/2} \xi.$$

Where  $*$  stands for the standard inner product, the cross product of vectors in  $\mathbb{R}^3$ , or the wedge product of differential one-forms.

(iii) *For two periodic sequences  $\{\phi^k\}$  and  $\{\xi^k\}$  of the same period,*

$$\sum_k \phi^k \delta_\zeta \xi^k = - \sum_k \delta_\zeta \phi^k \xi^k, \sum_k \phi^k \delta_\zeta^2 \xi^k = \sum_k \delta_\zeta^2 \phi^k \xi^k.$$

Where summation index  $k$  ranges over the period of the sequences.

*Proof.* The first item can be obtained by expanding and rearranging [40, 5]. For the second item, using definitions of the discrete operators and properties of the product  $*$ , we get

$$\begin{aligned} D_\zeta^{\alpha/2} \phi * A_\zeta^{\alpha/2} \xi &= \frac{e^{\alpha\Delta\zeta/2} \phi^{k+1} - e^{-\alpha\Delta\zeta/2} \phi^k}{\Delta\zeta} * \frac{e^{\alpha\Delta\zeta/2} \xi^{k+1} + e^{-\alpha\Delta\zeta/2} \xi^k}{2} \\ &= \frac{1}{2\Delta\zeta} \left( e^{\alpha\Delta\zeta} \phi^{k+1} * \xi^{k+1} + \phi^{k+1} * \xi^k - \phi^k * \xi^{k+1} - e^{-\alpha\Delta\zeta} \phi^k * \xi^k \right). \end{aligned} \quad (1.15)$$

And similarly

$$\begin{aligned}
A_\zeta^{\alpha/2} \phi * D_\zeta^{\alpha/2} \xi &= \frac{e^{\alpha\Delta\zeta/2} \phi^{k+1} + e^{-\alpha\Delta\zeta/2} \phi^k}{2} * \frac{e^{\alpha\Delta\zeta/2} \xi^{k+1} - e^{-\alpha\Delta\zeta/2} \xi^k}{\Delta\zeta} \\
&= \frac{1}{2\Delta\zeta} (e^{\alpha\Delta\zeta} \phi^{k+1} * \xi^{k+1} - \phi^{k+1} * \xi^k + \phi^k * \xi^{k+1} - e^{-\alpha\Delta\zeta} \phi^k * \xi^k). \quad (1.16)
\end{aligned}$$

Adding eqs. (1.15) and (1.16) we get

$$\begin{aligned}
D_\zeta^{\alpha/2} \phi * A_\zeta^{\alpha/2} \xi + A_\zeta^{\alpha/2} \phi * D_\zeta^{\alpha/2} \xi &= \frac{1}{\Delta\zeta} (e^{\alpha\Delta\zeta} \phi^{k+1} * \xi^{k+1} - e^{-\alpha\Delta\zeta} \phi^k * \xi^k) \\
&= D_\zeta^\alpha (\phi * \xi)
\end{aligned}$$

as desired. Differential forms and wedge product are discussed in Appendix A.

The last item can be proved by expanding and rearranging terms of the expansion and using periodicity of the sequences. Indeed, assuming  $\phi^{k+M-1} = \phi^k$  and  $\xi^{k+M-1} = \xi^k$  for all  $k$ , we get

$$\begin{aligned}
\sum_{k=1}^{M-1} \phi^k \delta_\zeta \xi^k &= \frac{1}{2\Delta\zeta} (\phi^1(\xi^2 - \xi^{M-1}) + \phi^2(\xi^3 - \xi^1) + \phi^3(\xi^4 - \xi^2) + \dots + \phi^{M-1}(\xi^1 - \xi^{M-2})) \\
&= -\frac{1}{2\Delta\zeta} (\xi^1(\phi^2 - \phi^{M-1}) + \xi^2(\phi^3 - \phi^1) + \xi^3(\phi^4 - \phi^2) + \dots + \xi^{M-1}(\phi^1 - \phi^{M-2})) \\
&= -\sum_{k=1}^{M-1} \delta_\zeta \phi^k \xi^k
\end{aligned}$$

Similarly, using the summation by parts formula, a discrete analog of integration by parts formula,

$$\sum_{k=m}^n f^k D_\zeta g^k = [f^{n+1} g^{n+1} - f^m g^m] - \sum_{k=m}^n g^{k+1} D_\zeta f^k$$

for two sequences  $\{f^k\}$  and  $\{g^k\}$  and periodicity of the sequences  $\{\phi^k\}$  and  $\{\xi^k\}$  we obtain

$$\begin{aligned}
\sum_{k=1}^{M-1} \phi^k \delta_\zeta^2 \xi^k &= \sum_{k=1}^{M-1} \phi^k T_\zeta D_\zeta D_\zeta \xi^k \\
&= - \sum_{k=1}^{M-1} D_\zeta \phi^k T_\zeta D_\zeta \xi^{k+1} \\
&= \sum_{k=1}^{M-1} D_\zeta D_\zeta \phi^k T_\zeta \xi^{k+2} \\
&= \sum_{k=1}^{M-1} T_\zeta D_\zeta D_\zeta \phi^{k+1} \xi^{k+1} \\
&= \sum_{k=1}^{M-1} T_\zeta D_\zeta D_\zeta \phi^k \xi^k \\
&= \sum_{k=1}^{M-1} \delta_\zeta^2 \phi^k \xi^k,
\end{aligned}$$

as desired. □

Numerical methods that use finite difference operators, such as operators of (1.14), are called finite difference methods.

### 1.3 Structure preservation background

First integrals, conformal invariants, and conservation laws of a DE are usually referred to as qualitative properties of the DE. A Numerical method (or integrator), that satisfies a discrete version of a qualitative property of a DE, is referred to as a *geometric integrator* or a *structure-preserving numerical method*. Since they have an extra property of structure-preservation, geometric integrators have been shown to be advantageous when compared to non-geometric integrators especially for long-time simulations of a problem. We are interested in geometric integrators which preserve conformal invariants of DEs.

A qualitative property of a DE that remains constant along any solution of the DE is referred to as a first integral or constant of motion. A lot of work has been done in the direction of developing geometric integrators that preserve first integrals. Indeed, there is a class of numerical methods which preserve first integral of conservative DEs. For example, Runge-Kutta methods given by eq. (1.12) preserve linear, quadratic, and symplectic invariants under certain restrictions on their coefficient functions [15, 26, 42, 45, 17]. More precisely, these methods preserve linear first integrals of the form

$$\mathcal{I} = \sigma^T z, \text{ with } \sigma \in \mathbb{R}^d,$$

i.e. they satisfy

$$\sigma^T z_{n+1} = \sigma^T z_n.$$

They also preserve quadratic first integrals of the form

$$\mathcal{I} = z^T W z,$$

where  $W \in \mathbb{R}^{d \times d}$  is a constant symmetric matrix, and symplectic 2-form

$$\mathcal{I} = dz \wedge \mathbf{J} dz$$

of the Hamiltonian system

$$\dot{z} = \mathbf{J}^{-1} \nabla_z H(z); \tag{1.17}$$

i.e. they satisfy

$$z_{n+1}^T W z_{n+1} = z_n^T W z_n \text{ and } dz_{n+1} \wedge \mathbf{J} dz_{n+1} = dz_n \wedge \mathbf{J} dz_n$$

provided their coefficients satisfy

$$\beta_i \alpha_{ij} + \beta_j \alpha_{ji} - \beta_i \beta_j = 0 \tag{1.18}$$

for all  $i, j$ . The skew-symmetric matrix  $\mathbf{J}^{-1}$  is referred to as the structure matrix of the Hamiltonian system with Hamiltonian  $H$  such that  $H : \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth function. We assume throughout this thesis that the phase space of a Hamiltonian system is even dimensional. Please see Appendix A for a review of differential forms and the wedge product. Similarly, Nyström methods preserve quadratic invariants under certain restrictions on their coefficient functions [19].

The integrators that preserve the symplectic 2-form  $dz \wedge \mathbf{J}dz$  are referred to as the *symplectic* integrators. Such integrators are volume preserving [28, 19]. Indeed, let  $\psi_t$  be the flow of a symplectic map for the Hamiltonian system (1.17) and let  $\Omega$  be a region in the phase space which is transported to another region  $\psi_t(\Omega)$  by the flow  $\psi_t$ . Then the change of variables formula for integrals gives

$$\begin{aligned} \int_{\Omega} dz &= \int_{\psi_t(\Omega)} \det(\psi'_t) dz \\ &= \int_{\psi_t(\Omega)} dz. \end{aligned}$$

The last equality follows because  $\det(\psi'_t) = 1$  for a symplectic map (see appendix A). Therefore, symplectic methods preserve phase space volume of Hamiltonian systems. In other words, a set of initial conditions occupying a solid region in phase space retain their original volume as the system evolves even though the shape of the region may change.

The references cited in this section so far are mostly concerned with ODEs although structure-preserving techniques therein can sometimes be extended to PDEs. There are other approaches, however, which take a different route to the structure-preserving discretization of PDEs. Some of these approaches include multi-symplectic discretizations, discrete variational methods, and average vector field methods. The first approach discretizes both space and time with symplectic geometric integrators, thus producing a multi-symplectic geometric integrator [8, 9, 28]. Multi-symplectic integrators aim to preserve local conservation law(s) which may result in the preser-

vation of certain first integrals of the PDEs. The second approach discretizes the Lagrangian of a PDE and then uses a discrete Lagrange principle to obtain a numerical integrator which is automatically multi-symplectic [30, 31]. While the multi-symplectic and the discrete variational methods guarantee the preservation of the symplectic structure, average vector field methods focus on preservation of energy of the system instead [20, 34, 14, 13].

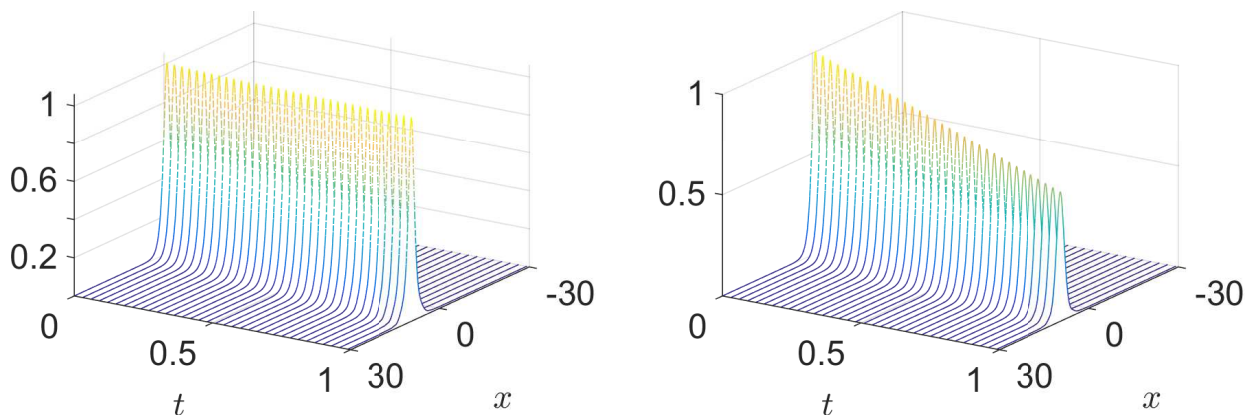


Figure 1.1: Solutions of a PDE without and with damping (Left to right).

There are damped (or dissipative) DEs, however, which have conformal invariants. Many physical systems have damping, dissipative, drag, resistive, or attenuating forces which result in conformal invariants of corresponding differential equation models. Figure 1.1 shows a typical example where a wave solution of a PDE preserves its shape without damping but decays in magnitude in the presence of damping as the time progresses. Figure 1.2 shows an example where a box of initial conditions in the phase space of a differential equation changes its shape as the time progresses. If the volume of the red box is same as the blue box, then the flow of differential equation preserves the phase space volume. The DE flow contracts the phase space volume if the volume of the red box is less than the blue box. Exponential decay in the magnitude of a solution of a DE and exponential phase space volume contraction along the flow are often the result of damping, which results in such conformal invariants. The aim of this thesis is to identify conformal invariants of

DEs and design numerical methods that preserve them.

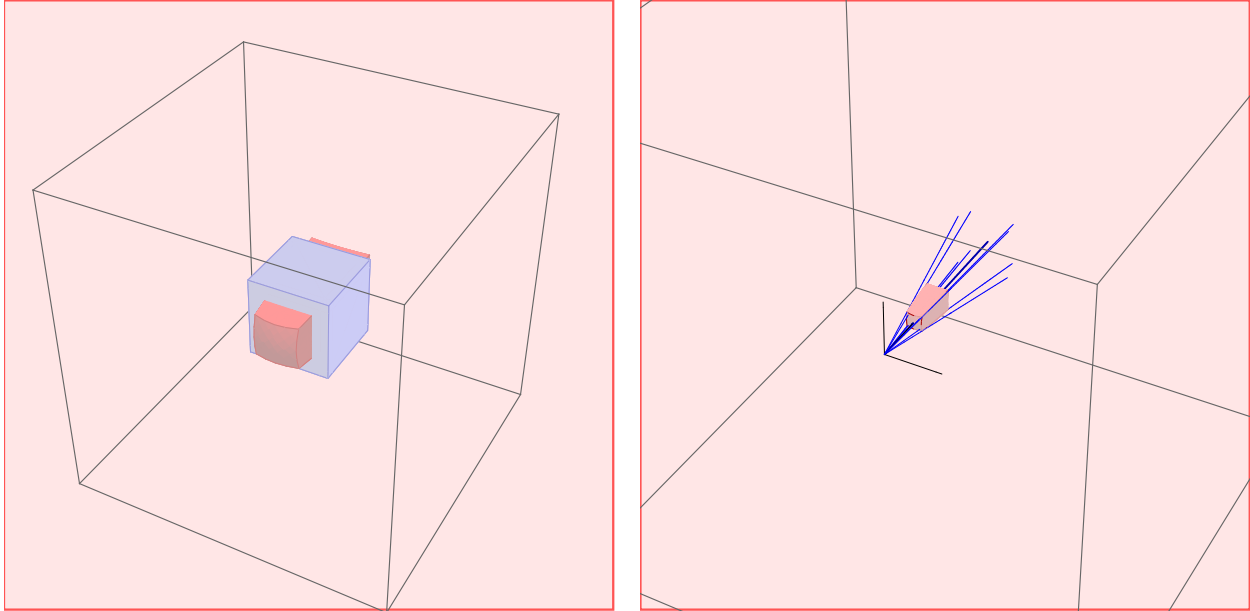


Figure 1.2: Deformation of the phase space volume (left) and the corresponding flow (right), which corresponds to the flow of a differential equation [41]. Blue and red boxes represent the initial and current phase space volumes, respectively. Blue lines denote solution trajectories of the DE.

Here, we mention some of the work that has already been done toward structure-preservation of dissipative differential equations. In [32, 33], authors have used differential geometric framework to define conformal Hamiltonian ODEs and constructed numerical methods which preserve conformal invariants. Dissipative systems were formulated as Birkhoffian systems in [44, 43], where authors used Birkhoffian framework to construct structure-preserving methods for the systems. Authors of [38, 43, 37, 40] generalized the multi-symplectic integration approach to dissipative PDEs which resulted in methods that preserve conformal invariants and local conservation laws. Authors of [31, 13] have suggested structure-preserving discrete gradient and average vector field methods, respectively, for a variety of damped PDEs. We remark at this point that all the reference in this section for geometric integration of conservative and dissipative DEs are systematic and methodical as they follow a strict prescription for obtaining structure-preserving numerical methods.



Moreover, one may be able to discretize a given DE with only some of these numerical integration approaches depending on the type and form of the DE.

Furthermore, we constructed conformal invariant preserving methods for damped DEs in [5]. These methods are based on the famous implicit midpoint and Störmer-Verlet methods. In [6], we derived structure-preserving conditions for ERK methods, specialized methods for linearly damped ODEs. We further constructed structure-preserving numerical methods for a damped driven non-linear Schrödinger equation in [4]. This thesis expounds on the results of [5, 6, 4].

## 1.4 Outline

The main body of this thesis can be divided into two parts. The first part, Chapters 2 and 3, is concerned with ODEs and their structure-preservation. We introduce a framework for numerical methods for damped linear ODEs in Chapter 2. We derive conditions under which the methods satisfy conformal invariants. These conditions are given as restrictions on coefficient functions of the methods. We also do accuracy and stability analysis of some of these methods where we show that some of the methods are unconditionally stable whereas others are only conditionally stable. We conduct some experiments on ODEs in Chapter 3 where we illustrate structure-preserving properties of geometric integrators and their advantages over non-geometric integrators. The second part, Chapters 4 and 5, is concerned with PDEs and their structure-preservation. Structure-preserving methods are provided for damped linear PDEs in Chapter 4, where we also describe conservation laws associated with these PDEs and their preservation by the methods. These methods are primarily obtained by discretizing space, time, or both using structure preserving methods of Chapter 2. We conduct some numerical experiments on PDEs in Chapter 5 to demonstrate advantages of structure-preserving integrators against other integrators. In Chapter 6, we give concluding remarks and future directions.

## CHAPTER 2: STRUCTURE-PRESERVING EXPONENTIAL RUNGE-KUTTA METHODS FOR ODES

Let us recall the initial value problem (1.1)

$$\dot{z}(t) = N(z(t)) - \gamma(t)z(t), \quad z(0) = z_0. \quad (1.1)$$

A variety of DEs can be put in the form of (1.1). Indeed, all the equations of Table 1.1 can be put in the form of eq. (1.1) by discretizing any spatial derivative(s). It is worth noticing that eq. (1.1) is obtained from eq. (1.11) when  $f(z, t) = N(z) - \gamma(t)z$  i.e. when the vector field  $f(z, t)$  can be split in nonlinear and linear components. Linear, quadratic, and symplectic first integrals of this IVP, with  $\gamma = 0$ , can be preserved using RK methods as discussed in Chapter 1. Here, we discuss preservation of corresponding conformal invariants of the equation. This will be achieved using the framework of exponential Runge-Kutta (ERK) methods.

ERK methods are a type of finite difference methods. They are specialized numerical methods for ODEs of the form (1.1) and are based on RK methods [28, 19]. In the following we define two common approaches of constructing ERK methods for eq. (1.1). The first approach uses a transformation to convert the equation into another equation which is then discretized using standard RK methods. Methods for the original equation (1.1) obtained by converting the methods for the transformed equation, using the original transformation, are referred to as integrating factor methods. The second approach uses approximations of the variation of constants formula for the solution of the IVP and the resulting methods are referred to as exponential time differencing methods.

## 2.1 ERK and partitioned ERK methods

Following the approach designed by Lawson [27], define the change of variables (like a Lawson transformation)

$$y(t) = e^{x_0(t)} z(t), \quad \text{with} \quad x_0(t) := \int_0^t \gamma(s) ds. \quad (2.1)$$

Then, eq. (1.1) becomes

$$\dot{y} = e^{x_0(t)} N(e^{-x_0(t)} y). \quad (2.2)$$

Notice, the same system of equations is achieved by multiplying eq. (1.1) through by the integrating factor. In this way, a method for solving eq. (1.1) can be constructed through standard methods that might be applied to eq. (2.2), and the resulting methods are typically called integrating factor methods. More specifically, applying a Runge-Kutta method (1.12) to eq. (2.2) gives

$$\begin{aligned} Y_i &= y_n + h \sum_{j=1}^s \alpha_{ij} e^{x_0(t_n+c_j h)} N(e^{-x_0(t_n+c_j h)} Y_j), \quad i = 1, \dots, s, \\ y_{n+1} &= y_n + h \sum_{i=1}^s \beta_i e^{x_0(t_n+c_i h)} N(e^{-x_0(t_n+c_i h)} Y_i), \end{aligned} \quad (2.3)$$

where  $s$  is the number of stages,  $h$  denotes the step size, and  $t_n = nh$  for  $n = 0, 1, 2, \dots$ ,  $y_n$  is the numerical solution, and  $Y_i$ 's are the stage variables. To write this in terms of the original variables, notice that

$$\int_{t_n+c_i h}^{t_n+h} \gamma(s) ds = \int_{t_n}^{t_n+h} \gamma(s) ds - \int_{t_n}^{t_n+c_i h} \gamma(s) ds = x_n(h) - x_n(c_i h),$$

where we define

$$x_n(t) = \int_0^t \gamma_n(s) ds \quad \text{with} \quad \gamma_n(s) := \gamma(s + t_n). \quad (2.4)$$

Thus, after manipulating the exponentials, the discretization can be rewritten in terms of the original variables to give a class of ERK methods for solving eq. (1.1), which are often called integrating

factor Runge-Kutta (IFRK) methods, given by

$$\begin{aligned} Z_i &= e^{-x_n(c_i h)} z_n + h \sum_{j=1}^s \alpha_{ij} e^{x_n(c_j h) - x_n(c_i h)} N(Z_j), \quad i = 1, \dots, s, \\ z_{n+1} &= e^{-x_n(h)} z_n + h \sum_{i=1}^s \beta_i e^{-x_n(h) + x_n(c_i h)} N(Z_i), \end{aligned} \quad (2.5)$$

where  $z_n \approx z(t_n)$  is the numerical solution.

A common alternative approach for constructing ERK methods is known as exponential time differencing, leading to the so called ETDRK methods. To construct methods of this type, we use the variation of constants formula and write the solution of eq. (2.2) as

$$y(t) = y(0) + \int_0^t e^{x_0(\tau)} N(e^{-x_0(\tau)} y(\tau)) d\tau$$

where  $y(0)$  is the initial value and  $x_0(t)$  is defined in eq. (2.1). Then using the transformation (2.1), the solution of eq. (1.1) becomes

$$z(t) = e^{-x_0(t)} z(0) + e^{-x_0(t)} \int_0^t e^{x_0(\tau)} N(z(\tau)) d\tau, \quad (2.6)$$

Following [23], the integral here can be approximated using a polynomial interpolation of  $N$ , particularly when  $\gamma$  is constant. In cases where  $\gamma$  is truly time-dependent, we may also require an approximation of the integral defined by  $x_n(t)$ . A simple and likely approach, which is rooted in the work of Hipp et al. [22], is to use an approximation, such as  $x_n(h) \approx h\gamma_n(h/2)$ .

In general, an  $s$ -stage ERK method for solving eq. (1.1), which includes both IFRK and ETDRK formulations, can be stated

$$\begin{aligned} Z_i &= \phi_i(h; \gamma_n) z_n + h \sum_{j=1}^s a_{ij}(h; \gamma_n) N(Z_j), \quad i = 1, \dots, s, \\ z_{n+1} &= \phi_0(h; \gamma_n) z_n + h \sum_{i=1}^s b_i(h; \gamma_n) N(Z_i). \end{aligned} \quad (2.7)$$

This formulation should be compared with eqs. (2.5) and (2.6). The *coefficients*,  $\phi_i$ ,  $\phi_0$ ,  $a_{i,j}$ , and  $b_i$ , are scalar functions of the constant step-size  $h$ , which also depend on the damping coefficient  $\gamma$  and the time index  $n$ , and they satisfy

$$\phi_i(h; 0) = \phi_0(h; 0) = 1, \quad a_{ij}(h; 0) = \alpha_{ij}, \quad b_i(h; 0) = \beta_i \quad (2.8)$$

for all  $i, j = 1, 2, \dots, s$ . The coefficients  $\phi_i$  and  $\phi_0$  are either exponential functions or rational approximations of such functions. Here and throughout this thesis we assume, for all  $i$ ,

$$\sum_{j=1}^s \alpha_{ij} = c_i, \quad \text{and} \quad \sum_{i=1}^s \beta_i = 1. \quad (2.9)$$

The RK method with coefficients  $\alpha_{ij}, \beta_i$  is obtained from the ERK method by setting  $\gamma = 0$  and is often referred to as the *underlying RK method*. An ERK method can be succinctly represented by a Butcher-like tableau, given by

$$\begin{array}{c|cc} c & A & \phi \\ \hline & b^T & \phi_0 \end{array}. \quad (2.10)$$

Entries  $c$ ,  $\phi$  and  $b$  of the tableau are column vectors and  $A$  is a square matrix, such that

$$c = \{c_i\}_{i=1}^s, \quad \phi = \{\phi_i\}_{i=1}^s, \quad b = \{b_i\}_{i=1}^s, \quad A = \{a_{ij}\}_{i,j=1}^s.$$

Notice that an ERK method is explicit if and only if the matrix  $A$  is lower triangular and implicit otherwise. The advantage of explicit methods is that they are computationally less expensive compared to implicit methods which require solution of an algebraic system of equations at every time step. The trade-off being that the former are generally conditionally stable whereas the latter are often unconditionally stable.

**Example 2.1.** Examples of some importance in the following exposition are listed here.

- Since our focus is on structure-preservation, natural choices for the underlying RK method are the Gauss-Legendre collocation methods, which are known to have order of accuracy  $2s$ ,

which is maximum for any s-stage RK method. ERK methods obtained from (2.5) using Gauss-Legendre collocation methods as underlying RK methods will be referred to as *GL-IFRK* methods. A one stage GL-IFRK method is given by (cf. [3, 12, 5])

$$\begin{array}{c|c|c} \frac{1}{2} & \frac{1}{2} & e^{-\int_{t_n}^{t_{n+1/2}} \gamma(s) ds} \\ \hline & e^{-\int_{t_{n+1/2}}^{t_{n+1}} \gamma(s) ds} & e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds} \end{array}. \quad (2.11)$$

A two stage GL-IFRK, with constant  $\gamma$ , is

$$\begin{array}{c|cc|c} \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{4} & \left(\frac{1}{4} - \frac{\sqrt{3}}{6}\right) e^{\frac{\sqrt{3}}{3}\gamma h} & e^{-\left(\frac{1}{2} - \frac{\sqrt{3}}{6}\right)\gamma h} \\ \frac{1}{2} + \frac{\sqrt{3}}{6} & \left(\frac{1}{4} + \frac{\sqrt{3}}{6}\right) e^{-\frac{\sqrt{3}}{3}\gamma h} & \frac{1}{4} & e^{-\left(\frac{1}{2} + \frac{\sqrt{3}}{6}\right)\gamma h} \\ \hline & \frac{1}{2} e^{-\left(\frac{1}{2} + \frac{\sqrt{3}}{6}\right)\gamma h} & \frac{1}{2} e^{-\left(\frac{1}{2} - \frac{\sqrt{3}}{6}\right)\gamma h} & e^{-\gamma h} \end{array}. \quad (2.12)$$

Methods of tableaux 2.11 and 2.12 are constructed using 1 and 2-stage Gauss-Legendre collocation methods, respectively, as underlying RK methods. A sixth-order, 3-stage, IFRK method can be constructed by using a 3-stage Gauss-Legendre collocation method as an underlying RK method. Order of accuracy of GL-IFRK methods is  $2s$  [12].

- A second order ETDRK method (with constant  $\gamma$ ), based on the implicit midpoint rule, is

$$\begin{array}{c|c|c} \frac{1}{2} & \frac{1}{2\gamma h}(e^{\gamma h/2} - e^{-\gamma h/2}) & e^{-\gamma h/2} \\ \hline & \frac{1}{\gamma h}(1 - e^{-\gamma h}) & e^{-\gamma h} \end{array}. \quad (2.13)$$

To solve the partitioned system (1.2) it may be desirable to employ one ERK method for the first equation and a different ERK method for the second. This approach yields a partitioned

exponential Runge-Kutta (PERK) method of the form

$$\begin{aligned}
Q_i &= \widehat{\phi}_i(h; \gamma_n^1)q_n + h \sum_{j=1}^s \widehat{a}_{ij}(h; \gamma_n^1)F(Q_j, P_j), \quad i = 1, \dots, s, \\
P_i &= \widetilde{\phi}_i(h; \gamma_n^2)p_n + h \sum_{j=1}^s \widetilde{a}_{ij}(h; \gamma_n^2)G(Q_j, P_j), \quad i = 1, \dots, s, \\
q_{n+1} &= \widehat{\phi}_0(h; \gamma_n^1)q_n + h \sum_{i=1}^s \widehat{b}_i(h; \gamma_n^1)F(Q_i, P_i), \\
p_{n+1} &= \widetilde{\phi}_0(h; \gamma_n^2)p_n + h \sum_{i=1}^s \widetilde{b}_i(h; \gamma_n^2)G(Q_i, P_i),
\end{aligned} \tag{2.14}$$

where *coefficients*,  $\widehat{\phi}_0, \widehat{\phi}_i, \widehat{a}_{ij}, \widehat{b}_j, \widetilde{\phi}_0, \widetilde{\phi}_i, \widetilde{a}_{ij}$ , and  $\widetilde{b}_j$  are scalar functions and they must satisfy the conditions required of an ERK method, namely eqs. (2.8) and (2.9). Here,  $Q_i, P_i$  are stage variables,  $[q_n, p_n] \approx [q(t_n), p(t_n)]$  is the numerical solution for  $n = 0, 1, 2, \dots$  and  $t_n = nh$ , and  $\gamma_n^k$  for  $k = 1, 2$  are defined according to eq. (2.4). In this case, the *underlying method* is a partitioned Runge-Kutta (PRK) method, obtained by setting  $\gamma^k = 0$  for  $k = 1, 2$ . (For our purposes, the Lobatto IIIA-IIIIB methods are natural choices for the underlying PRK methods.) A PERK method can be succinctly represented by a pair of Butcher-like tableaux, given by

$$\begin{array}{c|c|c} \widehat{c} & \widehat{A} & \widehat{\phi} \\ \hline \widehat{\phantom{c}} & \widehat{b}^T & \widehat{\phi}_0 \end{array} \quad \begin{array}{c|c|c} \widetilde{c} & \widetilde{A} & \widetilde{\phi} \\ \hline \widetilde{\phantom{c}} & \widetilde{b}^T & \widetilde{\phi}_0 \end{array} \tag{2.15}$$

one for each ERK method used. Notice that setting

$$\begin{aligned}
\widehat{\phi}_0 &= \widetilde{\phi}_0 = \phi_0, \quad \widehat{\phi}_i = \widetilde{\phi}_i = \phi_i, \quad \widehat{a}_{ij} = \widetilde{a}_{ij} = a_{ij}, \quad \widehat{b}_j = \widetilde{b}_j = b_j, \\
Z_i &= \begin{bmatrix} Q_i \\ P_i \end{bmatrix}, \quad z_n = \begin{bmatrix} q_n \\ p_n \end{bmatrix}, \quad N = \begin{bmatrix} F \\ G \end{bmatrix} \quad \text{for all } i, j
\end{aligned} \tag{2.16}$$

in a PERK method gives an ERK method.

As an example, consider an important special case of the ODE (1.2), given by

$$\dot{q} = \nabla_p T(p), \quad \dot{p} = -\nabla_q V(q) - \gamma p, \quad (2.17)$$

with  $0 < \gamma \in \mathbb{R}$ , which is known as a conformal Hamiltonian system [32]. If one or both equations of this system are discretized with an IFRK or ETD RK method, we refer to such a PERK method as *IFPRK* or *ETDPRK* method, respectively.

**Example 2.2.** Some examples of PERK methods for solving the system are:

- A first-order IFPRK method based on the symplectic Euler method

$$q_{n+1} = q_n + h \nabla_p T(p_{n+1}), \quad p_{n+1} = e^{-\gamma h} p_n - h \nabla_q V(q_n). \quad (2.18)$$

- A first-order ETDPRK method based on the symplectic Euler method

$$q_{n+1} = q_n + h \nabla_p T(p_{n+1}), \quad p_{n+1} = e^{-\gamma h} p_n + \frac{1}{\gamma} (e^{-\gamma h} - 1) \nabla_q V(q_n). \quad (2.19)$$

- A second-order IFPRK method based on the 2-stage Lobatto IIIA-III B (Störmer-Verlet) method

$$\begin{array}{c|ccc|c} 0 & 0 & 0 & 1 & \frac{1}{2} & \frac{1}{2} & 0 & e^{-\gamma h/2} \\ 1 & \frac{1}{2} & \frac{1}{2} & 1 & \frac{1}{2} & \frac{1}{2} & 0 & e^{-\gamma h/2} \\ \hline \wedge & \frac{1}{2} & \frac{1}{2} & 1 & \sim & \frac{1}{2} e^{-\gamma h/2} & \frac{1}{2} e^{-\gamma h/2} & e^{-\gamma h} \end{array} . \quad (2.20)$$

- A second-order IFPRK method based on the 2-stage Lobatto IIIA-III B method

$$\begin{array}{c|ccc|c} 0 & 0 & 0 & 1 & \frac{2}{4+\gamma h} & \frac{2}{4+\gamma h} & 0 & \frac{4}{4+\gamma h} e^{-\gamma h/4} \\ \frac{4}{4-\gamma h} & \frac{2}{4-\gamma h} & \frac{2}{4-\gamma h} & \frac{4+\gamma h}{4-\gamma h} e^{-\gamma h/2} & \frac{2}{4+\gamma h} & \frac{2}{4+\gamma h} & 0 & \frac{4}{4+\gamma h} e^{-\gamma h/4} \\ \hline \wedge & \frac{2}{4-\gamma h} & \frac{2}{4-\gamma h} & \frac{4+\gamma h}{4-\gamma h} e^{-\gamma h/2} & \sim & \frac{4-\gamma h}{4+\gamma h} \frac{e^{-\gamma h/4}}{2} & \frac{e^{-\gamma h/4}}{2} & \frac{4-\gamma h}{4+\gamma h} e^{-\gamma h/2} \end{array} . \quad (2.21)$$



- A second-order ETDPRK method based on the 2-stage Lobatto IIIA-IIIB method

$$\begin{array}{c|cc|c}
0 & 0 & 0 & 1 \\
1 & \frac{1}{2} & \frac{1}{2} & 1 \\
\hline
\hat{\phantom{0}} & \frac{1}{2} & \frac{1}{2} & 1
\end{array}
\quad
\begin{array}{c|cc|c}
\frac{1}{2} & \frac{1}{\gamma h}(1 - e^{-\gamma h/2}) & 0 & e^{-\gamma h/2} \\
\frac{1}{2} & \frac{1}{\gamma h}(1 - e^{-\gamma h/2}) & 0 & e^{-\gamma h/2} \\
\hline
\sim & \frac{1}{2}e^{-\gamma h/2} & \frac{1}{2}e^{-\gamma h/2} & e^{-\gamma h}
\end{array}
. \quad (2.22)$$

Order of accuracy of some of these methods will be proved in Section 2.4. Methods of eqs. (2.20) and (2.21) have been analyzed in some detail in [5, 37], where one can find their applications to ODEs and PDEs, in addition to their linear stability analysis and structure preservation properties for conformal Hamiltonian systems.

Though this discussion has been somewhat limited to integrating factor methods and exponential time differencing methods, other exponential integrators may be included in the general ERK and PERK formulations given in eqs. (2.7) and (2.14). It is important to keep this in mind, as the proofs concerning structure-preservation for ERK and PERK methods in the following sections give restrictions on the coefficient functions, which include, but are not necessarily limited to, integrating factor and exponential time differencing methods.

## 2.2 Preservation of conformal invariants

Conformal invariants were defined in Definition 1.1 and their examples were provided in Examples 1.2 to 1.5 and Table 1.1. It is natural to expect that numerical methods which preserve conformal invariants have certain advantages. This section is devoted to deriving sufficient conditions for preservation of conformal invariants by ERK and PERK methods. Some of the methods of the previous section are shown to satisfy these conditions. It is worth mentioning that setting  $\gamma = 0$  in eq. (1.1) and Definition 1.1 makes conformal invariants constants of motion. So structure-preservation of conservative systems becomes a special case of this exposition.

Similar to the preservation of first integrals, we define preservation of conformal invariants in the following definition. This definition should be compared with Definition 1.1 of the conformal invariants.

**Definition 2.3.** A numerical method  $z_{n+1} = \Psi_h(z_n)$  preserves a conformal invariant  $\mathcal{I}$  of eq. (1.1) if it satisfies

$$\mathcal{I}_{n+1} = e^{-\int_{t_n}^{t_{n+1}} 2\gamma(s)ds} \mathcal{I}_n$$

where  $\mathcal{I}_n = \mathcal{I}(t_n)$  (cf. [5]). Similarly, a numerical method  $z_{n+1} = \Psi_h(z_n)$  preserves a conformal invariant  $\mathcal{I}$  of eq. (1.2) if it satisfies

$$\mathcal{I}_{n+1} = e^{-\int_{t_n}^{t_{n+1}} (\gamma^1(s) + \gamma^2(s))ds} \mathcal{I}_n$$

where  $\mathcal{I}_n = \mathcal{I}(t_n)$ .

Invariants of the form  $\sigma z$ , for constant vector  $\sigma \in \mathbb{R}^d$ , are referred to as *linear invariants*, whereas invariants of the form  $z^T W z$ , for constant matrix  $W \in \mathbb{R}^{d \times d}$ , are referred to as *quadratic invariants* of eq. (1.1). For example, invariant  $H_\gamma$  of Example 1.2 is a quadratic invariant because

$$\begin{aligned} H_\gamma &= \frac{1}{2}(\kappa^2 \theta^2 + \omega^2) + \gamma \theta \omega \\ &= \begin{bmatrix} \theta & \omega \end{bmatrix} \begin{bmatrix} \frac{1}{2}\kappa^2 & \frac{1}{2}\gamma \\ \frac{1}{2}\gamma & \frac{1}{2} \end{bmatrix} \begin{bmatrix} \theta \\ \omega \end{bmatrix}. \end{aligned}$$

Total linear momentum  $\sum_{i=1}^N p_i$  is a linear invariant of eq. (1.5) because

$$\sum_{i=1}^N p_i^k = \mathbf{1}^T p^k, \text{ for } k = 1, 2, 3,$$

where  $\mathbf{1} \in \mathbb{R}^N$  is a column vector of ones and  $p^k = [p_1^k \ p_2^k \ \dots \ p_N^k]^T$ .

In the following, we derive sufficient conditions for quadratic and linear conformal invariant preservation by PERK methods in line with similar conditions for PRK methods.

**Theorem 2.4.** Suppose that the system (1.2) has a conformal invariant  $\mathcal{I} = q^T W p$ , with  $W \in \mathbb{R}^{d/2 \times d/2}$ . Then a PERK method applied to such a system satisfies  $\mathcal{I}_{n+1} = \widehat{\phi}_0 \widetilde{\phi}_0 \mathcal{I}_n$ , provided its coefficients satisfy

$$\widehat{b}_i \frac{\widetilde{\phi}_0}{\widetilde{\phi}_i} = \widetilde{b}_i \frac{\widehat{\phi}_0}{\widehat{\phi}_i}, \quad \widehat{b}_i \widetilde{a}_{ij} \frac{\widetilde{\phi}_0}{\widetilde{\phi}_i} + \widetilde{b}_j \widehat{a}_{ji} \frac{\widehat{\phi}_0}{\widehat{\phi}_j} - \widehat{b}_i \widetilde{b}_j = 0 \quad (2.23)$$

for all  $i, j$ . Moreover, the method preserves  $\mathcal{I}$  provided its coefficients also satisfy

$$\widehat{\phi}_0 \widetilde{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} (\gamma^1(s) + \gamma^2(s)) ds}.$$

*Proof.* Using the Kronecker product  $\otimes$ , one can write the system (2.14) as

$$\begin{aligned} Q &= \widehat{\phi} \otimes q_n + h(\widehat{A} \otimes I)F, \\ P &= \widetilde{\phi} \otimes p_n + h(\widetilde{A} \otimes I)G, \\ q_{n+1} &= \widehat{\phi}_0 q_n + h(\widehat{b}^T \otimes I)F, \\ p_{n+1} &= \widetilde{\phi}_0 p_n + h(\widetilde{b}^T \otimes I)G, \end{aligned} \quad (2.24)$$

where  $I \in \mathbb{R}^{d/2 \times d/2}$  is the identity matrix, and we define the vectors  $Q = \{Q_i\}_{i=1}^s$ ,  $P = \{P_i\}_{i=1}^s$ ,  $F = \{F_i\}_{i=1}^s$ , and  $G = \{G_i\}_{i=1}^s$ , with  $F_i = F(Q_i, P_i)$  and  $G_i = G(Q_i, P_i)$ . This implies that

$$\begin{aligned} q_{n+1}^T W p_{n+1} &= \widehat{\phi}_0 \widetilde{\phi}_0 q_n^T W p_n + h \widehat{\phi}_0 q_n^T W (\widetilde{b}^T \otimes I)G + h \widetilde{\phi}_0 ((\widehat{b}^T \otimes I)F)^T W p_n \\ &\quad + h^2 ((\widehat{b}^T \otimes I)F)^T W (\widetilde{b}^T \otimes I)G, \end{aligned}$$

which is equivalent to

$$\begin{aligned} q_{n+1}^T W p_{n+1} &= \widehat{\phi}_0 \widetilde{\phi}_0 q_n^T W p_n + h \widehat{\phi}_0 q_n^T (\widetilde{b}^T \otimes W)G + h \widetilde{\phi}_0 F^T (\widehat{b} \otimes W)p_n \\ &\quad + h^2 F^T (\widehat{b} \widetilde{b}^T \otimes W)G. \end{aligned} \quad (2.25)$$

Let  $\tilde{B}, \hat{B} \in \mathbb{R}^{s \times s}$  be diagonal matrices such that

$$\tilde{B}\hat{\phi} = \tilde{b} \quad \text{and} \quad \hat{B}\tilde{\phi} = \hat{b}, \quad (2.26)$$

respectively. Then using eq. (2.24) once again

$$\begin{aligned} Q^T(\tilde{B} \otimes W)G &= (\hat{\phi}^T \otimes q_n^T)(\tilde{B} \otimes W)G + hF^T(\hat{A}^T \otimes I)(\tilde{B} \otimes W)G \\ &= q_n^T(\tilde{b}^T \otimes W)G + hF^T(\hat{A}^T \tilde{B} \otimes W)G, \end{aligned} \quad (2.27)$$

and

$$F^T(\hat{B} \otimes W)P = F^T(\hat{b}^T \otimes W)p_n + hF^T(\hat{B}\tilde{A} \otimes W)G. \quad (2.28)$$

Using (2.27)-(2.28) in eq. (2.25), one gets

$$\begin{aligned} q_{n+1}^T W p_{n+1} &= \hat{\phi}_0 \tilde{\phi}_0 q_n^T W p_n + h\hat{\phi}_0 Q^T(\tilde{B} \otimes W)G + h\tilde{\phi}_0 F^T(\hat{B} \otimes W)P \\ &\quad + h^2 F^T((\tilde{b}\tilde{b}^T - \hat{\phi}_0 \hat{A}^T \tilde{B} - \tilde{\phi}_0 \hat{B} \tilde{A}) \otimes W)G. \end{aligned} \quad (2.29)$$

On the other hand, since  $\mathcal{I}$  is a conformal invariant, it follows that

$$0 = q^T W G(q, p) + F(q, p)^T W p,$$

for all  $q, p$ . Thus, provided  $\hat{b}_i \frac{\tilde{\phi}_0}{\tilde{\phi}_i} = \tilde{b}_i \frac{\hat{\phi}_0}{\hat{\phi}_i}$  for all  $i$ ,

$$0 = Q_i^T W G_i + F_i^T W P_i = Q_i^T \tilde{b}_i \frac{\hat{\phi}_0}{\hat{\phi}_i} W G_i + F_i^T \hat{b}_i \frac{\tilde{\phi}_0}{\tilde{\phi}_i} W P_i$$

which implies

$$0 = \sum_{i=1}^s Q_i^T \tilde{b}_i \frac{\hat{\phi}_0}{\hat{\phi}_i} W G_i + F_i^T \hat{b}_i \frac{\tilde{\phi}_0}{\tilde{\phi}_i} W P_i = \hat{\phi}_0 Q^T(\tilde{B} \otimes W)G + \tilde{\phi}_0 F^T(\hat{B} \otimes W)P.$$

Using this and eq. (2.23) in eq. (2.29), we get

$$q_{n+1}^T W p_{n+1} = \widehat{\phi}_0 \widetilde{\phi}_0 q_n^T W p_n.$$

Provided that  $\widehat{\phi}_0 \widetilde{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} (\gamma^1(s) + \gamma^2(s)) ds}$ , this implies

$$q_{n+1}^T W p_{n+1} = e^{-\int_{t_n}^{t_{n+1}} (\gamma^1(s) + \gamma^2(s)) ds} q_n^T W p_n,$$

i.e. the method preserves  $\mathcal{I}$ . □

Among the PERK methods of Example 2.2, only methods (2.18) and (2.20) satisfy the hypotheses of Theorem 2.4, and hence for these methods,

$$q_{n+1}^T W p_{n+1} = e^{-\gamma h} (q_n^T W p_n)$$

i.e. they preserve the conformal quadratic invariant  $q^T W p$ .

**Theorem 2.5.** *Let the function  $\mathcal{I} = \sigma_1^T q + \sigma_2^T p$ , with  $\sigma_1, \sigma_2 \in \mathbb{R}^{d/2}$ , be a conformal invariant of the system (1.2), and assume one of the following three conditions is satisfied: (i)  $\gamma^1(t) = \gamma^2(t) = \gamma(t)$ , (ii)  $\sigma_1 = 0$ , or (iii)  $\sigma_2 = 0$ . Then, a PERK method for such a system satisfies*

$$\sigma_1^T q_{n+1} + \sigma_2^T p_{n+1} = \widehat{\phi}_0 \sigma_1^T q_n + \widetilde{\phi}_0 \sigma_2^T p_n,$$

*provided its coefficients satisfy  $\widetilde{b}_i = \widehat{b}_i$ . Moreover, the method preserves  $\mathcal{I}$  provided its coefficients also satisfy*

$$\widehat{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} \gamma^1(s) ds}, \quad \widetilde{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} \gamma^2(s) ds}.$$

*Proof.* Formulation (2.24) of the PERK method implies

$$\sigma_1^T q_{n+1} + \sigma_2^T p_{n+1} = \widehat{\phi}_0 \sigma_1^T q_n + \widetilde{\phi}_0 \sigma_2^T p_n + \sigma_1^T (\widehat{b}^T \otimes I) F + \sigma_2^T (\widetilde{b}^T \otimes I) G.$$

Thus, to obtain the desired result, we must show that

$$\sigma_1^T(\widehat{b}^T \otimes I)F + \sigma_2^T(\widetilde{b}^T \otimes I)G = 0.$$

But, this follows from the fact that

$$0 = \sigma_1^T F_i + \sigma_2^T G_i = \sigma_1^T \widehat{b}_i F_i + \sigma_2^T \widetilde{b}_i G_i = \sum_{i=1}^s \sigma_1^T \widehat{b}_i F_i + \sigma_2^T \widetilde{b}_i G_i,$$

because  $\mathcal{I}$  is a conformal invariant of the system (1.2) with  $\gamma^1(t) = \gamma^2(t) = \gamma(t)$ , meaning  $\sigma_1^T F(q, p) + \sigma_2^T G(q, p) = 0$  for all  $q, p$ . This implies that

$$\sigma_1^T q_{n+1} + \sigma_2^T p_{n+1} = e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds} (\sigma_1^T q_n + \sigma_2^T p_n),$$

provided  $\widehat{\phi}_0 = \widetilde{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds}$ . The result for cases (ii) and (iii) follows automatically.  $\square$

Among the PERK methods of Example 2.2, only method (2.18) satisfies the hypotheses of this theorem and hence preserves the conformal linear invariants  $\sigma_1^T q + \sigma_2^T p$ .

The following result about the structure preserving properties of the ERK method can be derived in a manner analogous to those of the PERK method.

**Theorem 2.6.** *Suppose the system (1.1) has a conformal invariant  $\mathcal{I} = z^T W z$  where  $W \in \mathbb{R}^{d \times d}$  is a symmetric matrix. Then an ERK method applied to such a system satisfies  $\mathcal{I}_{n+1} = \phi_0^2 \mathcal{I}_n$  provided its coefficients satisfy*

$$b_i a_{ij} \frac{\phi_0}{\phi_i} + b_j a_{ji} \frac{\phi_0}{\phi_j} - b_i b_j = 0 \tag{2.30}$$

for all  $i, j$ . Moreover, the method preserves  $\mathcal{I}$  provided its coefficients also satisfy

$$\phi_0 = e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds}. \tag{2.31}$$

Indeed, one can informally use eqs. (2.7) and (2.16) in the proof of Theorem 2.4 and get this result.

It is worth noticing that the condition (2.30) reduces to eq. (1.18) when  $\gamma = 0$ . All the methods of Example 2.1 satisfy the hypotheses of this corollary, and hence they have the property

$$z_{n+1}^T W z_{n+1} = e^{-2\gamma h} z_n^T W z_n$$

i.e. these methods preserve conformal invariants of the form  $z^T W z$ . The following theorem follows directly from the definition of ERK methods.

**Theorem 2.7.** *Suppose that the system (1.1) has a conformal invariant  $\mathcal{I} = \sigma^T z$ , with  $\sigma \in \mathbb{R}^d$ . Then an ERK method applied to this system satisfies  $\mathcal{I}_{n+1} = \phi_0 \mathcal{I}_n$ . Moreover, the method preserves  $\mathcal{I}$  provided its coefficients satisfy eq. (2.31).*

Since all the methods of Example 2.1 satisfy the hypotheses of this theorem, they preserve conformal linear invariants.

### 2.3 Preservation of conformal symplecticness

In this section, we define conformal symplecticness and derive conformal symplecticness conditions for the ERK and the PERK methods. A special case of eq. (1.1) that we are particularly interested in occurs when  $d$  is even and the vector field  $N$  is of the form  $\mathbf{J}^{-1} \nabla_z H(z)$ . Substituting  $N(z) = \mathbf{J}^{-1} \nabla_z H(z)$  in eq. (1.1), we get the conformal Hamiltonian system [32]

$$\dot{z} = \mathbf{J}^{-1} \nabla_z H(z) - \gamma(t) z \tag{2.32}$$

where

$$z = \begin{bmatrix} q \\ p \end{bmatrix}, \quad \mathbf{J}^{-1} = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}$$

is a constant skew-symmetric matrix and  $H(z) : \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth function. Here  $I \in \mathbb{R}^{d/2 \times d/2}$  is the identity matrix. When  $\gamma(t) = 0$ , this system reduces to a *Hamiltonian system* with Hamil-

tonian  $H$ , which is a first integral, whereas function  $H$  in conformal Hamiltonian system (2.32) is not a first integral or conformal invariant of the system. Nonetheless, we refer to the function  $H$  in eq. (2.32) as Hamiltonian of the system.

From the variational equation associated with eq. (2.32)

$$d\dot{z} = \mathbf{J}^{-1}H_{zz}(z)dz - \gamma(t)dz, \quad (2.33)$$

where  $H_{zz}(z)$  is the Hessian matrix, one can easily obtain  $\dot{\omega} = -2\gamma(t)\omega$ , where  $\omega = dz \wedge \mathbf{J}dz$ , assuming  $H_{zz}(z)$  is symmetric. Indeed, taking the wedge product of eq. (2.33) with  $\mathbf{J}dz$  and using the properties of the wedge product from Appendix A we get

$$\begin{aligned} d\dot{z} \wedge \mathbf{J}dz &= \mathbf{J}^{-1}H_{zz}(z)dz \wedge \mathbf{J}dz - \gamma(t)dz \wedge \mathbf{J}dz, \\ \frac{1}{2} \frac{d}{dt}(dz \wedge \mathbf{J}dz) &= \mathbf{J}^{-1}\mathbf{J}^T H_{zz}(z)dz \wedge dz - \gamma(t)dz \wedge \mathbf{J}dz, \\ \frac{1}{2} \frac{d}{dt}(dz \wedge \mathbf{J}dz) &= -\mathbf{J}^{-1}\mathbf{J}H_{zz}(z)dz \wedge dz - \gamma(t)dz \wedge \mathbf{J}dz, \\ \frac{d}{dt}(dz \wedge \mathbf{J}dz) &= -2\gamma(t)dz \wedge \mathbf{J}dz, \end{aligned}$$

because  $H_{zz}(z)dz \wedge dz = 0$  as  $H_{zz}$  is a symmetric matrix.

Substituting  $F = \nabla_p \mathcal{H}(q, p)$ ,  $G = -\nabla_q \mathcal{H}(q, p)$ , along with  $\gamma^1 = 0$  and  $\gamma^2 = \gamma$  in eq. (1.2) we get

$$\begin{bmatrix} \dot{q} \\ \dot{p} \end{bmatrix} = \begin{bmatrix} \nabla_p \mathcal{H}(q, p) \\ -\nabla_q \mathcal{H}(q, p) \end{bmatrix} - \begin{bmatrix} 0 \\ \gamma(t)p \end{bmatrix}, \quad \begin{bmatrix} q(0) \\ p(0) \end{bmatrix} = \begin{bmatrix} q_0 \\ p_0 \end{bmatrix}. \quad (2.34)$$

From this equation, one can obtain

$$\dot{\omega} = -\gamma(t)\omega$$

for  $\omega = dq \wedge dp$ . Thus, we arrive at the following definition.

**Definition 2.8.** *Differential equations (1.1) and (1.2) are called conformal symplectic if the differ-*



ential two-forms

$$\omega = dz \wedge \mathbf{J}dz \text{ and } \omega = dq \wedge dp,$$

respectively, are conformal invariants of the corresponding equations.

Conformal symplecticness has been defined for time-independent perturbations in previous works [38, 40]. Definition 2.8 generalizes conformal symplecticness to eq. (1.1) which has time-dependent non-conservative perturbation. We want to develop numerical methods which preserve conformal symplecticness of conformal Hamiltonian systems.

**Definition 2.9.** *A numerical method  $z_{n+1} = \Psi_h(z_n)$  for solving differential equations eqs. (1.1) and (1.2) is said to preserve conformal symplecticity, and we call such a method conformal symplectic, if it preserves the corresponding conformal invariants of Definition 2.8.*

Please refer to Definitions 1.1 and 2.3 for the definitions of conformal invariants of a differential equation and their preservation by a numerical method. In particular, differential equation (1.1), with constant  $\gamma$ , is conformal symplectic if the symplectic 2-form  $\omega$  decays exponentially with time along a solution of the differential equation. A conformal symplectic method for a conformal symplectic ODE becomes a symplectic method for a symplectic ODE when  $\gamma = 0$ .

The following theorems concerning conformal symplecticness of PERK and ERK methods adopt the proof strategy presented in [28] for Hamiltonian systems.

**Theorem 2.10.** *A PERK method for eq. (2.34) satisfies*

$$dq_{n+1} \wedge dp_{n+1} = \widehat{\phi}_0 \widetilde{\phi}_0 (dq_n \wedge dp_n)$$

*provided its coefficients satisfy eq. (2.23). Moreover, the method is conformal symplectic if its coefficients also satisfy*

$$\widehat{\phi}_0 \widetilde{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds}.$$

*Proof.* In Kronecker product notation, the variational equation associated with the PERK method (2.24) applied to system (2.34) is

$$\begin{aligned}
dQ &= \widehat{\phi} \otimes dq_n + h(\widehat{A} \otimes I)dF, \\
dP &= \widetilde{\phi} \otimes dp_n + h(\widetilde{A} \otimes I)dG, \\
dq_{n+1} &= \widehat{\phi}_0 dq_n + h(\widehat{b}^T \otimes I)dF, \\
dp_{n+1} &= \widetilde{\phi}_0 dp_n + h(\widetilde{b}^T \otimes I)dG.
\end{aligned} \tag{2.35}$$

where  $F_i = \nabla_p \mathcal{H}(Q_i, P_i)$ ,  $G_i = -\nabla_q \mathcal{H}(Q_i, P_i)$ ,  $dF = F_Q dQ + F_P dP$  and  $dG = G_Q dQ + G_P dP$ .

This implies that

$$\partial_{Q_i} F_i + \partial_{P_i} G_i = \nabla_{pq} \mathcal{H}(Q_i, P_i) - \nabla_{qp} \mathcal{H}(Q_i, P_i) = 0, \quad \text{for all } i,$$

and hence  $F_Q^T + G_P = 0$ . Also  $F_P = F_P^T$  and  $G_Q = G_Q^T$ . Now, system (2.35) implies

$$\begin{aligned}
&dq_{n+1} \wedge dp_{n+1} - \widehat{\phi}_0 \widetilde{\phi}_0 dq_n \wedge dp_n \\
&= h\widehat{\phi}_0 dq_n \wedge (\widetilde{b}^T \otimes I)dG - h\widetilde{\phi}_0 dp_n \wedge (\widehat{b}^T \otimes I)dF + h^2 dF \wedge (\widetilde{b}\widehat{b}^T \otimes I)dG.
\end{aligned} \tag{2.36}$$

Using the first and the second equations of the system (2.35) and letting  $\widetilde{B}$  and  $\widehat{B}$  be diagonal matrices such that they satisfy eq. (2.26), we get

$$\begin{aligned}
dQ \wedge (\widetilde{B} \otimes I)dG &= \widehat{\phi} \otimes dq_n \wedge (\widetilde{B} \otimes I)dG + h(\widehat{A} \otimes I)dF \wedge (\widetilde{B} \otimes I)dG \\
&= dq_n \wedge (\widetilde{b}^T \otimes I)dG + h dF \wedge (\widehat{A}^T \widetilde{B} \otimes I)dG
\end{aligned} \tag{2.37}$$

and

$$dP \wedge (\widehat{B} \otimes I)dF = dp_n \wedge (\widehat{b}^T \otimes I)dF + h dG \wedge (\widetilde{A}^T \widehat{B} \otimes I)dF. \tag{2.38}$$

Using eq. (2.37)-eq. (2.38) in eq. (2.36) we get

$$\begin{aligned} dq_{n+1} \wedge dp_{n+1} - \widehat{\phi}_0 \widetilde{\phi}_0 dq_n \wedge dp_n &= hdQ \wedge (\widetilde{B} \widehat{\phi}_0 \otimes I) dG - hdP \wedge (\widehat{B} \widetilde{\phi}_0 \otimes I) dF \\ &\quad - h^2 dF \wedge ((\widehat{\phi}_0 \widehat{A}^T \widetilde{B} + \widetilde{\phi}_0 \widehat{B}^T \widetilde{A} - \widetilde{b} \widetilde{b}^T) \otimes I) dG. \end{aligned} \quad (2.39)$$

Since  $G_P + F_Q^T = 0$  and eq. (2.23) implies  $\widetilde{B} \widehat{\phi}_0 = \widehat{B} \widetilde{\phi}_0$ , we have

$$dQ \wedge (\widetilde{B} \widehat{\phi}_0 \otimes I) dG - dP \wedge (\widehat{B} \widetilde{\phi}_0 \otimes I) dF = dQ \wedge (\widetilde{B} \widehat{\phi}_0 \otimes I) (G_P + F_Q^T) dP = 0.$$

Using this and eq. (2.23) in eq. (2.39) yields

$$dq_{n+1} \wedge dp_{n+1} - \widehat{\phi}_0 \widetilde{\phi}_0 dq_n \wedge dp_n = 0.$$

This implies that

$$dq_{n+1} \wedge dp_{n+1} = e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds} dq_n \wedge dp_n$$

provided  $\widehat{\phi}_0 \widetilde{\phi}_0 = e^{-\int_{t_n}^{t_{n+1}} \gamma(s) ds}$ . □

One can get the following result for ERK methods by using eqs. (2.7) and (2.16) in the proof of Theorem 2.10.

**Theorem 2.11.** *An ERK method for eq. (2.32) satisfies*

$$dz_{n+1} \wedge \mathbf{J} dz_{n+1} = \phi_0^2 (dz_n \wedge \mathbf{J} dz_n)$$

*provided its coefficients satisfy eq. (2.30). Moreover, the method is conformal symplectic if its coefficients also satisfy eq. (2.31).*

Let us illustrate these theorems with some examples. PERK methods (2.18) and (2.20) satisfies the hypotheses of Theorem 2.10 and ERK methods (2.11)–(2.12) satisfy the hypotheses of Theorem 2.11; hence, these methods are conformal symplectic. Although methods (2.19) and (2.21) do

not satisfy the condition (2.23) it is a quick calculation to show that they are also conformal symplectic. Moreover method (2.21) has a conformal quadratic invariance [5]. This shows that a PERK method does not necessarily need to satisfy the condition (2.23) in order to preserve these geometric properties. Method (2.22) neither satisfies the condition (2.23) nor is conformal symplectic, so it is not enough to use any ERK method which has a conservative underlying RK method. It is interesting to note that the conditions (2.23) and (2.30) are part of the sufficient conditions for both conformal quadratic invariance and conformal symplecticness of the respective methods (cf. [7]). Also note that combining the hypotheses of these theorems with the conditions (2.8) implies that the underlying RK and PRK methods are always symplectic.

It is possible to define conformal symplecticness entirely in terms of flow maps rather than the wedge product. To this end, let  $\psi_t(z_0)$  be the flow map of the system (2.32) at time  $t$  with initial condition  $z_0$ , then an alternative (equivalent) formulation of the conformal symplecticness given in Definition 2.8 is

$$(\psi'_t(z_0))^T \mathbf{J}^{-1} \psi'_t(z_0) = e^{-2 \int_0^t \gamma(s) ds} \mathbf{J}^{-1}, \quad (2.40)$$

where  $\psi'_t(z_0)$  denotes the Jacobian and is a solution of the variational eq. (2.33). Indeed, let  $\psi_t(z_0) : U \rightarrow \mathbb{R}^{2d}$  be the flow map of the system (2.32) at time  $t$  with initial condition  $z_0$ . Since the solution  $z(t)$  of the system satisfies  $z(t, z_0) = \psi_t(z_0)$ , then

$$dz = \psi'_t(z_0) dz_0$$

and hence Definition 2.8 is equivalent to

$$\begin{aligned}
& \psi'_t(z_0)dz_0 \wedge \mathbf{J}^{-1}\psi'_t(z_0)dz_0 = e^{-2\int_0^t \gamma(s)ds} dz_0 \wedge \mathbf{J}^{-1}dz_0, \\
& \iff (\psi'_t(z_0))^T \mathbf{J}^{-1}\psi'_t(z_0)dz_0 \wedge dz_0 = e^{-2\int_0^t \gamma(s)ds} \mathbf{J}^{-1}dz_0 \wedge dz_0, \\
& \iff \left( (\psi'_t(z_0))^T \mathbf{J}^{-1}\psi'_t(z_0) - e^{-2\int_0^t \gamma(s)ds} \mathbf{J}^{-1} \right) dz_0 \wedge dz_0 = 0, \\
& \iff (\psi'_t(z_0))^T \mathbf{J}^{-1}\psi'_t(z_0) = e^{-2\int_0^t \gamma(s)ds} \mathbf{J}^{-1}. \tag{2.41}
\end{aligned}$$

Where we have used the the skew-symmetry of the matrix  $(\psi'_t(z_0))^T \mathbf{J}^{-1}\psi'_t(z_0) - e^{-2\int_0^t \gamma(s)ds} \mathbf{J}^{-1}$  and Lemma A.1 to get eq. (2.41). One can also obtain formula (2.40) using variational eq. (2.33) of eq. (2.32) by taking the time derivative of  $(\psi'_t(z_0))^T \mathbf{J}^{-1}\psi'_t(z_0)$ .

Thus, a numerical method with flow map  $\Psi_h$  is conformal symplectic if

$$(\Psi'_h(z_n))^T \mathbf{J}^{-1}(\Psi'_h(z_n)) = e^{-2\int_0^{t_n+1} \gamma(s)ds} \mathbf{J}^{-1}. \tag{2.42}$$

## 2.4 Accuracy and stability of ERK and PERK methods

In the last two sections, we derived conditions under which ERK and PERK methods preserve conformal invariants and conformal symplecticness of a differential equation. We showed that some methods of Examples 2.1 and 2.2 are structure-preserving. In this section, we take a closer look at some methods of these examples, their order of accuracy, and their linear stability.

Conditions under which numerical methods have certain order of accuracy are called order conditions. Order conditions for RK methods are well known but the theory is not developed as much for ERK methods partly because of increased complexity due to additional coefficient functions. Some of the work that has been in the direction of order of accuracy analysis of ERK methods includes [2, 10]. The aim of this section is not to derive general order conditions for ERK meth-

ods but to derive these conditions for specific cases of ERK and PERK methods. In the following theorem, we derive order of accuracy of some of the methods of Examples 2.1 and 2.2.

**Theorem 2.12.** *Method of eqs. (2.11), (2.20) and (2.21) are second order for constant  $\gamma$ .*

*Proof.* In the notation of discrete operators of Lemma 1.8, eq. (2.11) becomes

$$\mathbf{J}D_t^{\gamma/2}z = N(A_t^{\gamma/2}z). \quad (2.43)$$

A Taylor series expansion reveals

$$e^{\gamma h}z(t+h) = z(t) + h(\partial_t + \gamma)z(t) + \frac{h^2}{2}(\partial_t + \gamma)^2z(t) + \frac{h^3}{6}(\partial_t + \gamma)^3z(t) + \dots,$$

which implies, assuming  $z_n = z(t_n)$  for all  $n$ ,

$$\begin{aligned} \mathbf{J}D_t^{\gamma/2}z &= \mathbf{J} \left( \frac{1}{h}(e^{\gamma h/2}z_{n+1} - e^{-\gamma h/2}z_n) \right) \\ &= \mathbf{J} \left( \frac{1}{h}(e^{\gamma h/2}z(t_n+h) - e^{-\gamma h/2}z(t_n)) \right) \\ &= \mathbf{J} \left( (\partial_t + \frac{\gamma}{2})z(t_n) + \frac{h}{2}(\partial_t + \frac{\gamma}{2})^2z(t_n) + \frac{\gamma}{2}z(t_n) - \frac{h\gamma^2}{2 \cdot 4}z(t_n) \right) + \mathcal{O}(h^2) \\ &= \mathbf{J} \left( (\partial_t + \gamma)z(t_n) + \frac{h}{2}(\partial_t^2 + \gamma\partial_t)z(t_n) \right) + \mathcal{O}(h^2) \end{aligned} \quad (2.44)$$

and

$$\begin{aligned} N(A_t^{\gamma/2}z_n) &= N \left( \frac{1}{2}(e^{\gamma h/2}z_{n+1} + e^{-\gamma h/2}z_n) \right) \\ &= N \left( \frac{1}{2}(e^{\gamma h/2}z(t_n+h) + e^{-\gamma h/2}z(t_n)) \right) \\ &= N(z(t_n)) + \frac{h}{2}\partial_z N(z(t_n))\dot{z}(t_n) + \mathcal{O}(h^2). \end{aligned} \quad (2.45)$$

Substituting eqs. (2.44) and (2.45) into eq. (2.43), we get

$$\begin{aligned}\mathbf{J}\dot{z}(t_n) &= N(z(t_n)) - \gamma\mathbf{J}z(t_n) - \frac{h}{2} (\mathbf{J}(\ddot{z}(t_n) + \gamma\dot{z}(t_n)) - \partial_z N(z(t_n))\dot{z}(t_n)) + \mathcal{O}(h^2) \\ &= N(z(t_n)) - \gamma\mathbf{J}z(t_n) + \mathcal{O}(h^2).\end{aligned}$$

Comparing this equation to ODE  $\mathbf{J}\dot{z}(t) = N(z(t)) - \gamma\mathbf{J}z(t)$ , we see that the method is second order.

Notice that the method

$$\begin{aligned}(1 + \gamma_1 h/2)p_{n+1/2} &= e^{-\gamma_3 h/2} p_n - \frac{h}{2} \nabla_q V(q_n), \\ (1 - \gamma_1 h/2)q_{n+1} &= e^{-\gamma_1 h/2} [(1 + \gamma_1 h/2)e^{-\gamma_1 h/2} q_n + h \nabla_p T(p_{n+1/2})] \\ e^{\gamma_2 h/2} p_{n+1} &= e^{-\gamma_1 h/2} \left[ (1 - \gamma_1 h/2)p_{n+1/2} - \frac{h}{2} \nabla_q V(q_{n+1}) \right]\end{aligned}\tag{2.46}$$

gives method (2.20) on setting  $\gamma_1 = 0, \gamma_2 = 2\gamma, \gamma_3 = 2\gamma$ , and method (2.21) on setting  $\gamma_1 = \gamma, \gamma_2 = 0, \gamma_3 = \gamma$ . Assuming  $q_n = q(t_n), p_n = p(t_n)$  and expanding first equation of method (2.46) in its Taylor series about  $h = 0$ , we obtain

$$\dot{p}(t_n) = -\nabla_q V(q(t_n)) - (\gamma_1 + \gamma_3)p(t_n) + \mathcal{O}(h).\tag{2.47}$$

Now rewriting the second equation of method (2.46), we obtain

$$\frac{1}{h} [(1 - \gamma_1 h/2)e^{\gamma_1 h/2} q_{n+1} - (1 + \gamma_1 h/2)e^{-\gamma_1 h/2} q_n] = \nabla_p T(p_{n+1/2}).$$

Replacing  $p_{n+1/2}$  by its definition in this equation and doing Taylor expansions about  $h = 0$  we get

$$\dot{q}(t_n) = \nabla_p T(p(t_n)) - \frac{h}{2} (T_{pp}(p(t_n)) \cdot ((\gamma_1 + \gamma_3)p(t_n) + \nabla_q V(q(t_n)))) + \ddot{q}(t_n) + \mathcal{O}(h^2).$$

In this equation, using eq. (2.47) and  $\ddot{q} = T_{pp}(p) \cdot \dot{p} + \mathcal{O}(h)$ , we get

$$\dot{q}(t_n) = \nabla_p T(p(t_n)) + \mathcal{O}(h^2).\tag{2.48}$$

Now rewriting third equation of method (2.46)

$$\frac{2}{h} [e^{(\gamma_1+\gamma_2)h/2} p_{n+1} - (1 - \gamma_1 h/2) p_{n+1/2}] = -\nabla_q V(q_{n+1}),$$

replacing  $p_{n+1/2}$  and  $q_{n+1}$  by their definitions in this equation and expanding about  $h = 0$  we get

$$\begin{aligned} \dot{p}(t_n) &= -\nabla_q V(q(t_n)) - \frac{1}{2}(3\gamma_1 + \gamma_2 + \gamma_3)p(t_n) \\ &\quad - \frac{h}{2} (\ddot{p}(t_n) + (\gamma_1 + \gamma_2)\dot{p}(t_n) + V_{qq}(q(t_n)) \cdot \nabla_p T(p(t_n))) \\ &\quad - \frac{1}{4} (\gamma_1 - \gamma_2 + \gamma_3) (3\gamma_1 + \gamma_2 + \gamma_3) p(t_n) - \gamma_1 \nabla_q V(q(t_n)) + \mathcal{O}(h^2). \end{aligned}$$

In the this equation, substituting

$$\dot{q} = \nabla_p T(p) + \mathcal{O}(h^2),$$

$$\dot{p} = -\nabla_q V(q) - \frac{1}{2}(3\gamma_1 + \gamma_2 + \gamma_3)p + \mathcal{O}(h),$$

and

$$\ddot{p} = -V_{qq}(q) \cdot \dot{q} - \frac{1}{2}(3\gamma_1 + \gamma_2 + \gamma_3)\dot{p} + \mathcal{O}(h)$$

and simplifying we get

$$\begin{aligned} \dot{p}(t_n) &= -\nabla_q V(q(t_n)) - \frac{1}{2}(3\gamma_1 + \gamma_2 + \gamma_3)p(t_n) + \frac{h}{4} (\gamma_1 + \gamma_2 - \gamma_3) \nabla_q V(q(t_n)) + \mathcal{O}(h^2) \\ &= -\nabla_q V(q(t_n)) - 2\gamma p(t_n) + \mathcal{O}(h^2) \end{aligned} \tag{2.49}$$

because  $\gamma_1 + \gamma_2 - \gamma_3 = 0$  for both methods (2.20) and (2.21). Comparing eqs. (2.48) and (2.49) to the ODE being discretized

$$\dot{q}(t) = \nabla_p T(p(t)),$$

$$\dot{p}(t) = -\nabla_q V(q(t)) - 2\gamma p(t),$$

we see that both methods (2.20) and (2.21) are second order accurate □



**Remark.** Other methods that are closely related to method (2.43) have been considered before [25, 40, 44]. The difference between these methods and the method (2.43) is how they spread out the damping over a stencil. The discrete operators for some of these methods read

$$\mathcal{D}_t z_n = \frac{1}{h}(z_{n+1} - e^{-\gamma h} z_n), \quad \mathcal{A}_t z_n = \frac{1}{2}(z_{n+1} + e^{-\gamma h} z_n).$$

The methods that use these discrete operators are only first order accurate. Method (2.43), however, spreads out the damping more evenly over the stencil. The uniform distribution of the damping results in improved accuracy of the method (2.43) over the methods that use non-uniform distribution.

Other methods can be similarly shown to have a certain order of accuracy. However, complexity in doing order analysis using Taylor series multiplies quickly as the order increases. It was shown in [12] that GL-IFRK methods, e.g. eqs. (2.11) and (2.12) with constant  $\gamma$ , have order of accuracy  $2s$ .

Next, we derive the stability function of the ERK method (2.7) and then derive stability condition for certain special cases. Substituting  $N(z) = \lambda z$ ,  $\lambda \in \mathbb{C}$ , in method (2.7), we get the method in the vector form

$$\begin{aligned} Z &= \phi z_n + h\lambda AZ, \\ z_{n+1} &= \phi_0 z_n + h\lambda b^T Z \end{aligned}$$

where

$$Z = \begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_s \end{bmatrix} \quad \text{and} \quad \phi = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_s \end{bmatrix}.$$

Assuming the matrix  $(I - h\lambda A)$  is invertible, this gives

$$Z = (I - h\lambda A)^{-1} \phi z_n,$$

which in turn implies

$$z_{n+1} = R(h\lambda, h\gamma) z_n. \quad (2.50)$$

where the stability function  $R(w, x)$  is given by

$$R(w, x) = \phi_0(x) + wb(x)^T (I - wA(x))^{-1} \phi(x). \quad (2.51)$$

If  $|R(w, x)| < 1$  then the numerical solution  $z_n$  remains bounded and the method is called A-stable.

The stability function  $R$  is reminiscent of the stability function

$$R_0(w) = R(w, 0) = 1 + wb(0)^T (I - wA(0))^{-1} \mathbf{1} \quad (2.52)$$

of the underlying RK method where  $\mathbf{1} \in \mathbb{R}^s$  is a column vector of ones. One, therefore, wonders if there is a simple relationship between the two, at least for some special cases of ERK. This brings us to the next theorem.

The following theorem establishes a simple relationship between stability functions of IFRK methods and their underlying RK methods, thus giving complete characterization of the stability function of the former methods.

**Theorem 2.13.** *The stability function  $R$  of IFRK methods, eq. (2.5), satisfies*

$$R(w, x) = e^{-x} R_0(w)$$

where  $R_0(w)$  is the stability function of the underlying RK method.

*Proof.* Since  $R_0$  is the stability function of the underlying method, it follows from (2.3)

$$y_{n+1} = R_0(\lambda h)y_n.$$

Using the transformation (2.1), it follows that

$$z_{n+1} = R_0(\lambda h)e^{-x_0(h)}z_n = e^{-\gamma h}R_0(\lambda h)z_n.$$

Therefore the result follows. □

A-stability of a numerical method ensures that the numerical solution of a linear differential equation does not grow to infinity. A-stability of GL-IFRK methods can be quickly established from Theorem 2.13. Indeed, since the underlying Gauss-Legendre method is A-stable i.e. its stability function satisfies

$$|R_0(w)| < 1 \text{ for all } w \text{ such that } \Re(w) < 0,$$

the GL-IFRK methods are also A-stable because

$$|R(w, x)| = e^{-x}|R_0(w)| < 1 \text{ for all } x > 0 \text{ and all } w \text{ such that } \Re(w) < 0$$

by Theorem 2.13. Since  $|R(w, x)| < 1$ , the numerical solution  $z_n$  remains bounded.

Let us now turn to the stability of PERK methods. Consider the exact solution of the damped harmonic oscillator (1.4)

$$\begin{bmatrix} \theta(t) \\ \omega(t) \end{bmatrix} = e^{-\gamma t} \begin{bmatrix} \cos(\beta t) + \frac{\gamma}{\beta} \sin(\beta t) & \frac{1}{\beta} \sin(\beta t) \\ -\frac{\kappa^2}{\beta} \sin(\beta t) & \cos(\beta t) - \frac{\gamma}{\beta} \sin(\beta t) \end{bmatrix} \begin{bmatrix} \theta_0 \\ \omega_0 \end{bmatrix}, \quad (2.53)$$

where  $\beta = \sqrt{\kappa^2 - \gamma^2}$ . The transition matrix, which is the square matrix on the right hand side of eq. (2.53), has eigenvalues

$$\lambda_{\pm} = e^{-\gamma t} \left( \mu \pm \sqrt{\mu^2 - 1} \right), \quad (2.54)$$

where  $\mu = \cos(\beta t)$ . Notice, the dissipative part of the solution is completely described by the exponential  $e^{-\gamma t}$ , and the conservative part is completely described by the complex conjugate pairs  $\mu \pm \sqrt{\mu^2 - 1}$ , which lie on the unit circle because  $|\mu| \leq 1$ .

It is desirable that our numerical methods reproduce this behavior. Eigenvalues of the transition matrix of a numerical method generally depend on step-size  $h$  also. We consider the method stable if  $|\mu(h)| \leq 1$ . Now, apply the explicit PERK method (2.20) to eq. (1.4) to obtain

$$\begin{bmatrix} \theta_{n+1} \\ \omega_{n+1} \end{bmatrix} = \begin{bmatrix} 1 - \frac{\kappa^2 h^2}{2} & h e^{-\gamma h} \\ h \kappa^2 e^{-\gamma h} \left( \frac{\kappa^2 h^2}{4} - 1 \right) & e^{-2\gamma h} \left( 1 - \frac{\kappa^2 h^2}{2} \right) \end{bmatrix} \begin{bmatrix} \theta_n \\ \omega_n \end{bmatrix}.$$

Hence, the eigenvalues of the transition matrix are (2.54) with

$$\mu = \left( 1 - \frac{\kappa^2 h^2}{2} \right) \cosh(\gamma h).$$

Thus, requiring  $|\mu| \leq 1$  to ensure stability implies the stability condition

$$1 - \operatorname{sech}(\gamma h) \leq \frac{\kappa^2 h^2}{2} \leq 1 + \operatorname{sech}(\gamma h),$$

which gives a restriction on the step size  $h$  for explicit method (2.20) to be stable. One can similarly show that the explicit PERK method (2.21) is also conditionally stable [5]. In contrast, implicit methods of eqs. (2.11) and (2.12) are unconditionally stable. A step size restriction is often the price one pays for more straightforward implementation and less computational complexity of explicit methods compared to implicit methods.

## CHAPTER 3: ODE APPLICATIONS AND EXPERIMENTS

In this chapter, ERK and PERK methods of Examples 2.1 and 2.2 are applied to various ODEs with constant damping and time-dependent non-conservative perturbation terms to demonstrate their properties of structure preservation [5, 37]. Studies have also performed numerical simulations using various first order ERK methods on very similar problems [25, 38, 40, 44]. Our purpose in this chapter is to demonstrate the effectiveness of ERK and PERK methods from a few points of view that are different from previous studies. First, we demonstrate preservation of conformal symplecticness. Second, we consider problems with time-dependent damping. Third, we conduct experiments using methods of higher orders (four and six). Fourth, we illustrate the advantages of such methods for a damped Poisson (non-canonical) system. Fifth, we implement structure-preserving exponential time differencing methods and compare the results to more commonly used integrating factor methods.

### 3.1 Linear oscillators

In this section, we discretize linear oscillators with constant damping and time-dependent non-conservative perturbation terms with ERK and PERK methods of Examples 2.1 and 2.2. To this end, consider the following generalization of eq. (1.4)

$$\ddot{\theta} + 2\gamma(t)\dot{\theta} + \kappa^2\theta = 0 \tag{3.1}$$

where  $\kappa \in \mathbb{R}$  is constant frequency. This equation can be put in the form of conformal Hamiltonian system (2.32) by setting

$$H = \frac{1}{2}(\kappa^2\theta^2 + \omega^2) + \gamma(t)\theta\omega.$$

Thus the methods of Example 2.1 are applicable to this equation. In the following, we analyze numerical simulations in two cases:  $\gamma = \text{const.}$  and  $\gamma(t) = \frac{1}{2}\epsilon \cos(2t)$  with  $\epsilon \in \mathbb{R}$ .

### 3.1.1 Constant damping

When the damping parameter  $\gamma$  is constant, ODE (3.1) becomes a constant coefficient linear DE of order 2 and can be rewritten as the conformal Hamiltonian system (2.17) by setting

$$T(\theta) = \frac{1}{2}\kappa^2\theta^2, \quad V(\omega) = \frac{1}{2}\omega^2, \quad \gamma(t) = 2\gamma,$$

where  $\gamma$  on the right hand side of the last equation is constant. In the form of eq. (2.17), numerical methods of Example 2.2 are applicable to the oscillator. In this case, we can compare our numerical solutions against the exact solution. To begin, we compare the integrating factor and exponential time differencing methods given in eqs. (2.18) and (2.19), respectively. Both methods are first order and conformal symplectic. Figure 3.1 shows the average absolute error for each method, as  $\gamma$  is fixed, while the frequency and the step size are varied. Notice, the exponential time differencing method exhibits clear advantages over the integrating factor method as the frequency increases, and these advantages are more pronounced as the step size decreases, even for problems with high frequencies.

Next, we present an example illustrating order of accuracy and structure preservation by GL-IFRK methods, IFRK methods (2.5) having the Gauss-Legendre schemes as the underlying RK methods. To illustrate higher order convergence, we apply stage 1, 2, and 3 GL-IFRK methods to eq. (3.1) with  $\gamma$  constant. Figure 3.2 shows the ratio of local absolute error in solution and the step size as a function of the step size, which agrees with the theoretical local order of the methods.

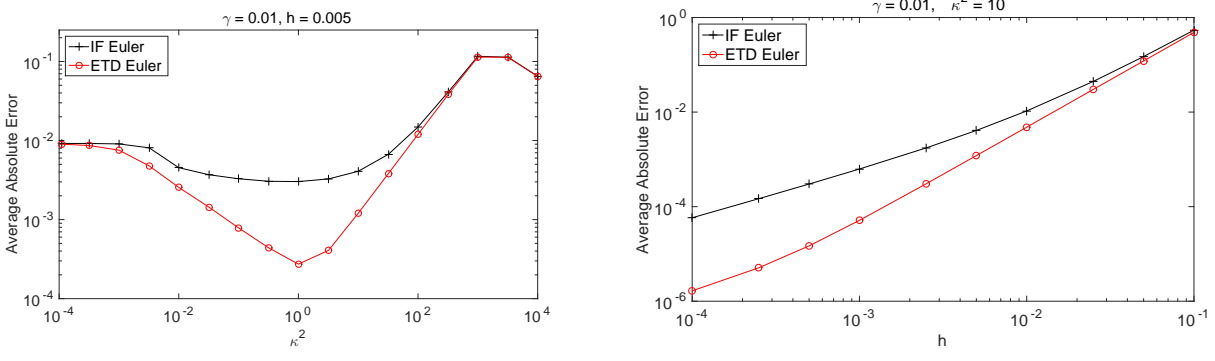


Figure 3.1: A comparison of the average absolute solution error for the conformal symplectic Euler methods given in (2.18) and (2.19) for solving eq. (3.1) with  $\gamma = 0.01$ . Initial condition:  $\theta(0) = 0$ ,  $\omega(0) = 10$ ; final time:  $T = 50$ .

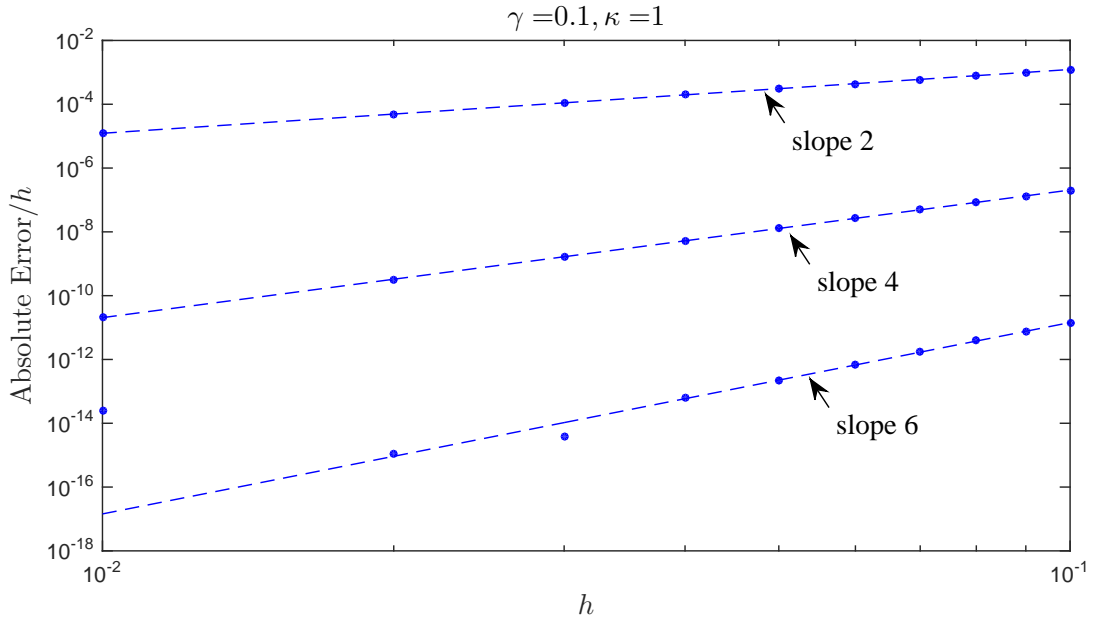


Figure 3.2: Local absolute error in solution over the step size for IFRK methods of stages 1, 2 and 3 applied to ODE (3.1). Dashed lines represent the slopes with which they are labeled.

### 3.1.2 Time-dependent damping

Setting  $\gamma(t) = \frac{1}{2}\epsilon \cos(2t)$  with  $\epsilon \in \mathbb{R}$  in eq. (3.1) yields a special case of Hill's equation, which is used to model rain-wind induced vibrations in an oscillator. Depending on parameter values, the solutions ( $q : \mathbb{R} \rightarrow \mathbb{R}$ ) in this case may be periodic, bounded, or unbounded [21], providing richer solution behavior than the linear oscillator with constant damping.

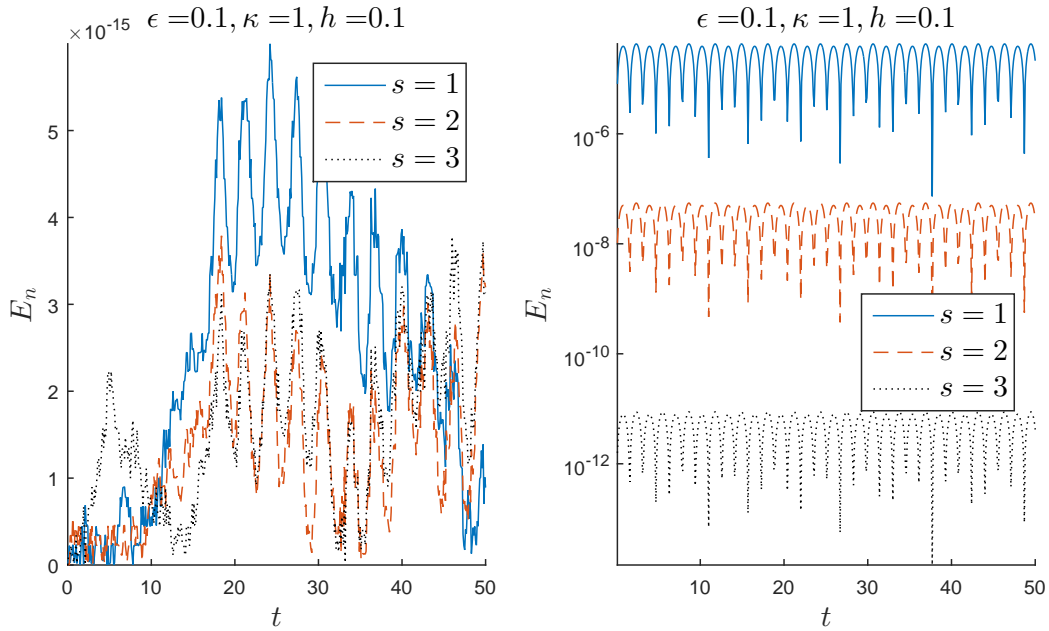


Figure 3.3: Error  $E_n$  (3.2) in conformal symplecticness for the GL-IFRK methods (left) and the standard Gauss-Legendre methods (right) applied to eq. (3.1) with  $\gamma(t) = \frac{1}{2}\epsilon \cos(2t)$ .

An IFRK method for solving eq. (3.1) is given by eq. (2.5). Since  $\phi_0 = e^{-x_n(h)}$ , we know that such methods satisfy the hypotheses of Theorem 2.11. A numerical method with flow map  $\Psi_h(z_n)$ , which solves the system (2.32), is conformal symplectic if

$$E_n := \|(\Psi'_h(z_n))^T \mathbf{J}^{-1} \Psi'_h(z_n) - e^{-2 \int_0^{t_{n+1}} \gamma(s) ds} \mathbf{J}^{-1}\| \quad (3.2)$$

vanishes, see eq. (2.42). Since  $z_{n+1} = \Psi_h(z_n)$  implies  $dz_{n+1} = \Psi'_h(z_n) dz_n$ , the Jacobian  $\Psi'_h(z_n)$



can be computed by numerically solving the system for  $dz_{n+1}$ . For instance, the first two equations of the system (2.35) can be numerically solved for  $dQ, dP$  using exact methods or fixed point iterations and the resulting solutions can be substituted in the last two equations of the system to find  $dq_{n+1}, dp_{n+1}$ . In Figure 3.3, we plot the error  $E_n$  for both the GL-IFRK methods and the standard Gauss-Legendre methods of stages 1, 2 and 3, illustrating preservation of conformal symplecticness by the GL-IFRK methods, but not by the standard methods.

## 3.2 Damped pendulum

Here, we implement geometric integrators on a damped pendulum. Damping, nonlinearity, and external driving force result in chaotic solutions in certain parameter regimes of a damped driven pendulum. Chaotic solutions are those which have sensitive dependence on the initial condition: changing the initial condition slightly results in a comparatively large change in solution trajectory of the system.

### 3.2.1 Damped pendulum

To begin with, consider the pendulum equation with constant linear damping

$$\ddot{\theta} + 2\gamma\dot{\theta} + \sin(\theta) = 0. \quad (3.3)$$

We consider the problem in two cases. First, when  $\gamma$  is purely imaginary the differential equation has rapid oscillations resulting from the first order linear term when  $|\gamma| > 1$ . Second, using a real constant satisfying  $\gamma > 1$  the pendulum is strongly damped. It can be shown that the differential equation satisfies

$$\frac{d}{dt}(\frac{1}{2}\omega^2 - \cos(\theta)) = -2\gamma\omega^2$$

which could be called the energy balance. Though this is not what we have called a conformal invariant for the system, it does provide a way to measure the accuracy of a method.

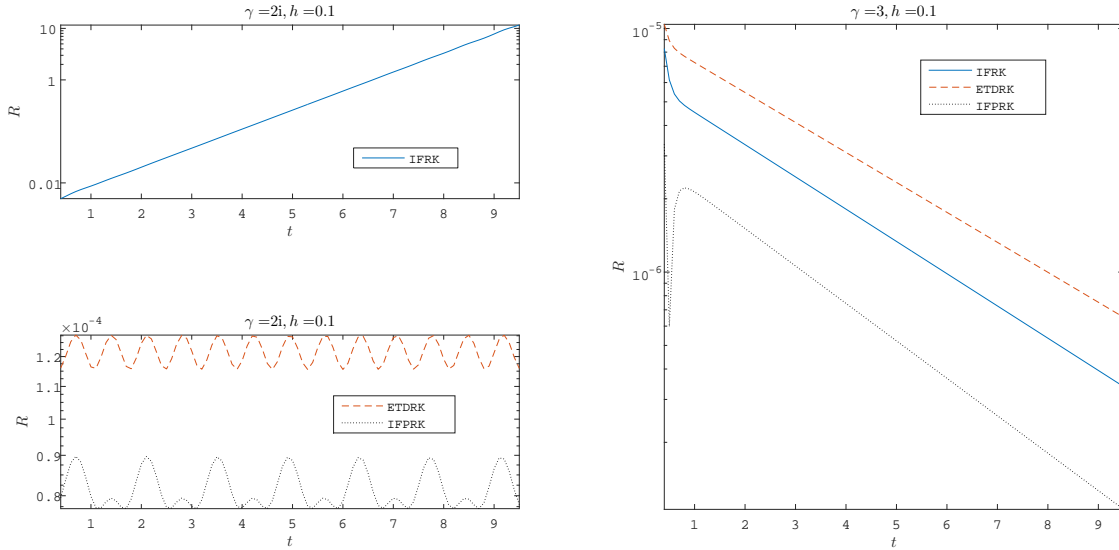


Figure 3.4: The residual (3.4) for three numerical solutions of (3.3). The methods used are IFRK (2.11); ETDRK (2.13); and IFPRK (2.20) denoted here by IFRK, ETDRK, and IFPRK, respectively. Left: rapid oscillation with imaginary  $\gamma$ ; Right: strong damping with real  $\gamma$ .

To this model problem, we apply the IFRK methods, eqs. (2.11) and (2.20), and the ETDRK method (2.13). All three methods are second order accurate and conformal symplectic. The expression  $(e^{\gamma h/2} - e^{-\gamma h/2})$  in tableau 2.13 is evaluated by computing  $2 \sinh(\gamma h/2)$  instead. Denoting the eighth order central finite difference operator by  $\mathcal{D}_t$ , we plot the residual

$$R(t_n) = |\mathcal{D}_t(\frac{1}{2}\omega_n^2 - \cos(\theta_n)) + 2\gamma\omega_n^2| \quad (3.4)$$

for each of the three methods in Figure 3.4. For imaginary  $\gamma$  method (2.11) produces a serious drift in the energy balance, while the other two methods more accurately maintain the energy balance. For  $\gamma \in \mathbb{R}$ , all the methods show rapid decay. In each case, the PERK method produces smaller

residuals, and the fact that it is explicit gives it a strong advantage.

It is important to note here that these differences between integrating factor and exponential time differencing methods are a result of choosing relatively large values of  $\gamma$ . Other comparisons of interest with small damping coefficients do not often reveal such obvious differences between the methods. As a result, the integrating factor methods may often be preferable, because they are generally easier to construct and analyze.

### 3.2.2 Damped driven pendulum

A damped driven pendulum is governed by the following ODE

$$\ddot{\theta} + 2\Gamma\dot{\theta} + \frac{g}{l}\sin(\theta) = F\sin(\Omega t),$$

where  $\theta$  is the angle the pendulum makes with the vertical,  $\Gamma$  is the damping parameter,  $g$  is the acceleration due to gravity,  $l$  is the length of the pendulum;  $F$  and  $\Omega$  are the amplitude and the angular frequency of the driving force. Choosing  $\Omega^{-1}$  to be the new units of time, the last equation becomes

$$\ddot{\theta} + 2\gamma\dot{\theta} + \lambda^2\sin(\theta) = f\sin(t), \tag{3.5}$$

with

$$\gamma = \Omega^{-1}\Gamma, \quad \lambda^2 = \Omega^{-2}\frac{g}{l}, \quad f = \Omega^{-2}F.$$

The damping, driving force, and the sinusoidal terms in this equation are the ones responsible for chaos in the system. In the numerical experiments that follow in this section, we will discretize eq. (3.5) with the numerical methods developed in earlier sections and a second order RK (Heun's) method and compare the numerical results.

To this end, let us write equation eq. (3.5) as a first order system

$$\begin{aligned}\dot{\theta} &= \omega, \\ \dot{\omega} &= -\lambda^2 \sin(\theta) + f \sin(t) - 2\gamma\omega.\end{aligned}$$

When  $f = 0$ , this is a conformal Hamiltonian system (2.32) with

$$H = \frac{1}{2}(\omega^2 - \lambda^2 \cos(\theta)) - \gamma\theta\omega.$$

Even though eq. (3.5) is not conformal Hamiltonian, it is conformal symplectic, nonetheless.

Therefore, we extend the application of conformal symplectic methods of Chapter 2 to this system.

An implicit method based on (2.11) for (3.5) is

$$\begin{aligned}D_t^{\gamma/2}\theta &= A_t^{\gamma/2}\omega + \gamma A_t^{\gamma/2}\theta, \\ D_t^{\gamma/2}\omega &= -\lambda^2 \sin(A_t^{\gamma/2}\theta) + f \sin(A_t t) - \gamma A_t^{\gamma/2}\omega.\end{aligned}\tag{3.6}$$

An explicit PERK method for (3.5) reads

$$\begin{aligned}(1 + \gamma_1 \Delta t/2)\omega_{n+1/2} &= e^{-\gamma_3 \Delta t/2}\omega_n + \frac{\Delta t}{2}(-\lambda^2 \sin(\theta_n) + f \sin(t_n)), \\ (1 - \gamma_1 \Delta t/2)\theta_{n+1} &= e^{-\gamma_1 \Delta t/2}[(1 + \gamma_1 \Delta t/2)e^{-\gamma_1 \Delta t/2}\theta_n + \Delta t\omega_{n+1/2}], \\ e^{\gamma_2 \Delta t/2}\omega_{n+1} &= e^{-\gamma_1 \Delta t/2}\left[(1 - \gamma_1 \Delta t/2)\omega^{i+1/2} + \frac{\Delta t}{2}(-\lambda^2 \sin(\theta_{n+1}) + f \sin(t_{n+1}))\right]\end{aligned}\tag{3.7}$$

which gives a method based on tableau (2.20) on setting  $\gamma_1 = 0, \gamma_2 = 2\gamma, \gamma_3 = 2\gamma$ , which we call

CSV1, and a method based on tableau (2.21) on setting  $\gamma_1 = \gamma, \gamma_2 = 0, \gamma_3 = \gamma$ , which we call

CSV2. Finally, Heun's method for this system is

$$\begin{aligned}\omega_{n+1/2} &= \omega_n + \Delta t(-\lambda^2 \sin(\theta_n) + f \sin(t_n) - 2\gamma\omega_n), \\ \theta_{n+1} &= \theta_n + \frac{\Delta t}{2}(\omega_n + \omega_{n+1/2}), \\ \omega_{n+1} &= \frac{1}{2}(\omega_n + \omega_{n+1/2}) + \frac{\Delta t}{2}(-\lambda^2 \sin(\theta_n + \Delta t\omega_n) + f \sin(t_{n+1}) - 2\gamma\omega_{n+1/2}).\end{aligned}\tag{3.8}$$

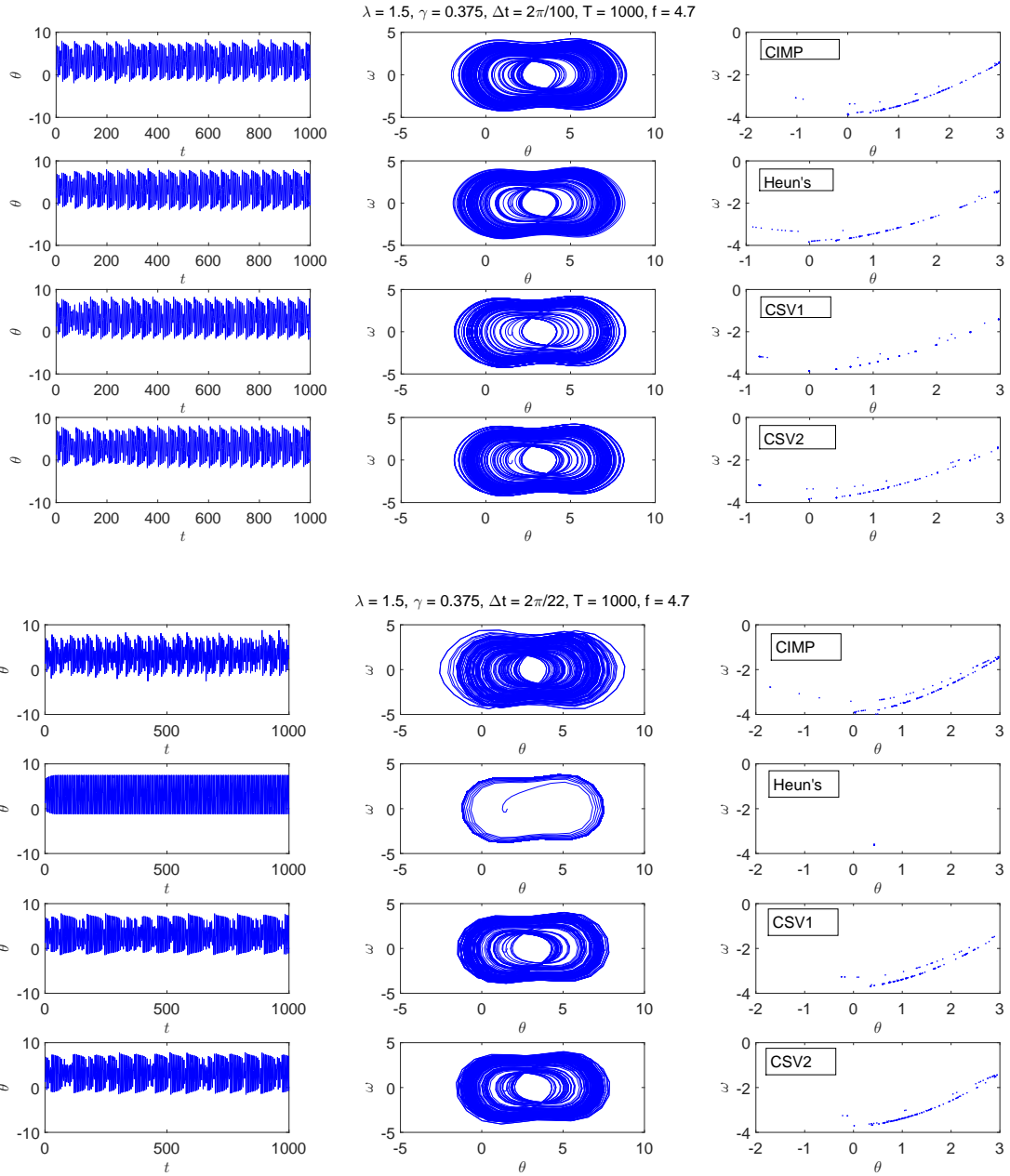


Figure 3.5: Left to right: time series, phase space and Poincaré sections of damped driven oscillator, eq. (3.5), with the parameter values mentioned in the title.  $T$  is the final time. CIMP and Heun's stand for eqs. (3.6) and (3.8)

Figure 3.5 gives an example of a bounded orbit which is neither periodic nor convergent. The figure shows the failure of Heun's method to reproduce the chaotic orbit with a larger time-step size whereas other methods, which are conformal symplectic, successfully do so.

### 3.3 N-body ODE

Let us define the following Hamiltonian function

$$\mathcal{H}(q, p) = \sum_{i=1}^N \frac{\|p_i\|^2}{2m_i} + \sum_{i=1}^{N-1} \sum_{j=i+1}^N \phi_{ij}(\|q_i - q_j\|) \quad (3.9)$$

where  $q = [q_1, q_2, \dots, q_N]$  and  $p = [p_1, p_2, \dots, p_N]$ . With this Hamiltonian, one can rewrite the N-body system (1.5)-(1.6) as conformal Hamiltonian system (2.17) given by

$$\begin{aligned} \partial_t q_i &= \nabla_{p_i} \mathcal{H}(q, p), \\ \partial_t p_i &= -\nabla_{q_i} \mathcal{H}(q, p) - 2\gamma p_i. \end{aligned}$$

This conformal Hamiltonian system is now amenable for numerical treatment by methods of tableaux (2.20) and (2.21). Defining

$$H(q, p) = \mathcal{H}(q, p) + \gamma \sum_{i=1}^N q_i \cdot p_i \text{ and } z_i = \begin{bmatrix} q_i \\ p_i \end{bmatrix},$$

we can also rewrite the N-body system as conformal Hamiltonian system (2.32) using

$$\partial_t z_i = \mathbf{J}^{-1} \nabla_{z_i} H(q, p) - \gamma z_i,$$

which lends itself to the method of Example 2.1 including method (2.11).

In this experiment, we use methods of tableaux (2.20), (2.21), and (2.11), referred to as CSV1, CSV2, and CIMP in Table 3.1 and Figure 3.6, for the N-body system (1.5)-(1.6) of Example 1.3

with interaction potential

$$\phi_{ij}(r) = -\frac{Gm_i m_j}{r}$$

to simulate the system and graphically illustrate momentum preserving properties of these methods that were proved in Section 2.2. Table 3.1 encapsulates the said properties of the three methods. Notice that CSV2 and CIMP do not preserve linear momentum.

Table 3.1: Total linear momentum and total angular momentum for the three methods. Here  $q_i^n \approx q_i(t_n)$  etc.

Method	$\sum_i p_i^{n+1} =$	$\sum_i q_i^{n+1} \times p_i^{n+1} =$
CSV1	$e^{-2\gamma\Delta t} \sum_i p_i^n$	$e^{-2\gamma\Delta t} \sum_i q_i^n \times p_i^n$
CSV2	$\frac{2-\gamma\Delta t}{2+\gamma\Delta t} e^{-\gamma\Delta t} \sum_i p_i^n$	$e^{-2\gamma\Delta t} \sum_i q_i^n \times p_i^n$
CIMP	$\frac{2-\gamma\Delta t}{2+\gamma\Delta t} e^{-\gamma\Delta t} \sum_i p_i^n$	$e^{-2\gamma\Delta t} \sum_i q_i^n \times p_i^n$

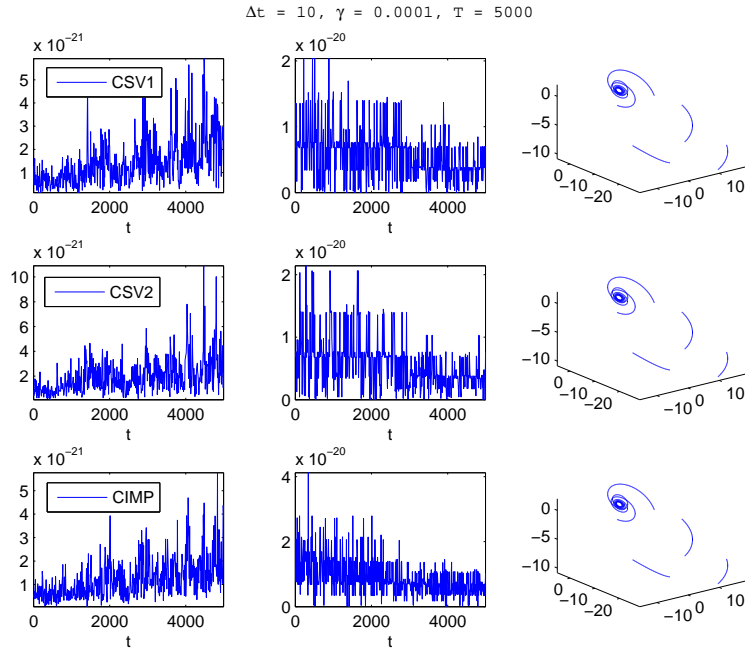


Figure 3.6: Left to right: Error in linear momentum, error in angular momentum and corresponding solution trajectories of N-body system.

Figure 3.6 shows the norm of the difference between headings and corresponding entries of Table 3.1. This difference is less than the machine precision as expected. The innermost objects spiral inward and the rest of the objects are in orbits which are fast converging toward the fixed point.

### 3.4 Rigid body with periodic perturbation

Here, we present an example illustrating conformal quadratic invariant preservation by the ERK methods for eq. (1.7), which is not a conformal Hamiltonian system. Figure 3.7 shows plots of both

$$E_n = |C(z_n) - C(z_0)e^{-\frac{\epsilon}{2} \sin(2t_n)}|, \quad \text{and} \quad \mathcal{E}_n = |H(z_n) - H(z_0)e^{-\frac{\epsilon}{2} \sin(2t_n)}|, \quad (3.10)$$

which are the residuals in Casimir and energy, respectively, of (1.7) and  $z_n$  is the numerical solution. The residuals due to the standard Gauss-Legendre methods are proportional to the order of the methods, i.e. the residuals decrease as the order of the methods increases. The figure verifies that the conformal quadratic invariants are preserved by the GL-IFRK methods.

We have applied various ERK and PERK methods to a variety of ODEs. An ETD method is seen to be advantageous compared to an IF method in Figure 3.1. PERK methods exhibit better performance (with respect to accuracy and efficiency) compared to ERK methods in Figures 3.4 to 3.6. Figure 3.4 shows that structure-preserving methods can be beneficial even in chaotic regimes. Figures 3.3, 3.5 and 3.7 illustrate structure-preservation properties of ERK methods and their advantages over non-structure preserving methods.



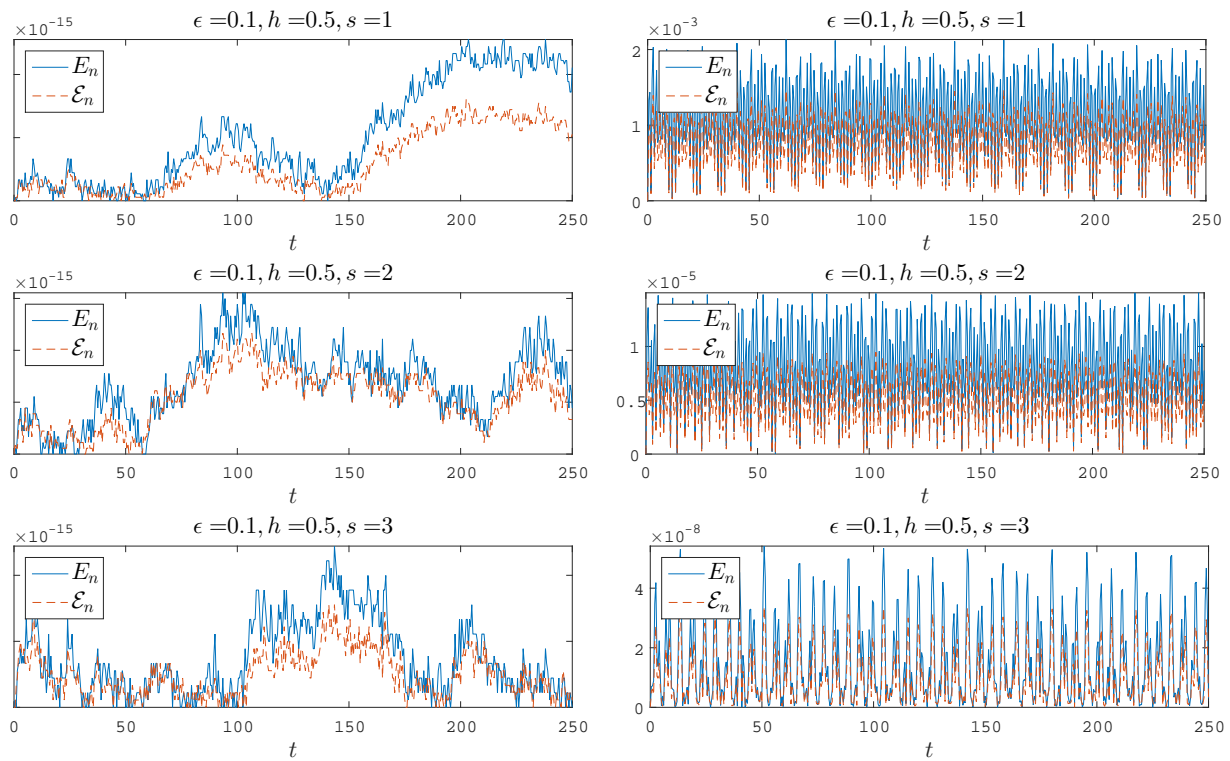


Figure 3.7: Casimir and energy errors (3.10) for simulations of the system (1.7). Left: GL-IFRK methods; Right: standard Gauss-Legendre methods.

## CHAPTER 4: STRUCTURE-PRESERVING METHODS FOR PDES

Some physical phenomena have not only time dependence but spatial dependence also. Such phenomena are modeled by partial differential equations. Much like ODEs, PDEs also have qualitative properties such as integral invariants and conservation laws. Therefore, structure-preserving methods or geometric integrators for PDEs are desirable. In this chapter, we discretize a given PDE with a structure-preserving method in space, time, or both to get a structure-preserving method for the PDE. The geometric integrators to discretize space or time may be chosen from the ERK methods presented in Chapter 2.

### 4.1 Multi-conformal-symplectic PDEs

Multi-conformal-symplectic (MCS) methods are structure-preserving numerical methods for partial differential equations. They can be seen as a generalization of conformal symplectic idea for ODEs to PDEs. The central concept behind MCS integrators is that two symplectic integrators, one in space and time each, work in tandem to preserve MCS structure of a PDE.

It was first noted in [8, 9] that some PDEs can be put in the following form

$$\mathbf{K}z_t + \mathbf{L}z_x = \nabla S(z) \tag{4.1}$$

where  $\mathbf{K}$  and  $\mathbf{L}$  are constant skew-symmetric matrices,  $S(z)$  is a smooth scalar function,  $z$  is a vector of field variables, and subscripts denote usual partial derivatives. Equation (4.1) can be seen as a PDE equivalent of the Hamiltonian system (1.17). In the form of eq. (4.1), a PDE automatically satisfies certain local conservation laws. Many conservative PDEs can be put in the form of (4.1). See [28] and references therein for a discussion on this PDE, associated conservation laws, and

some examples.

In this thesis, we consider the following generalization of PDE (4.1)

$$\mathbf{K}z_t + \mathbf{L}z_x = \nabla S(z) - \frac{a}{2}\mathbf{K}z + \mathbf{F}(t), \quad (4.2)$$

where  $a$  is a non-negative real number and  $\mathbf{F}(t)$  is a time dependent vector function. Additional terms on the right hand side of this PDE, as compared to PDE (4.1), usually represent dissipation and forcing terms in the PDE. This equation was first introduced in [38, 40] where its properties were also discussed. Some examples of eq. (4.2) follow.

**Example 4.1.** Consider a damped Klein-Gordon equation

$$u_{tt} = u_{xx} - cu - 2\gamma u_t. \quad (4.3)$$

Here  $u = u(x, t)$  is the solution of the equation,  $c$  is a real constant, and  $\gamma$  is the damping parameter. Subscripts denote the usual partial derivatives. This equation arises in relativistic mechanics and has been discussed extensively in the literature, including [16, 40]. The case  $\gamma = 0$  corresponds to the conservative counterpart of the equation.

We can write eq. (4.3) in the form of (4.2) as [40]

$$\begin{aligned} -v_t - w_x &= cu + 2\gamma v \\ u_t - p_x &= v \\ -p_t + u_x &= -w + 2\gamma p \\ w_t + v_x &= -cp \end{aligned}$$

which can be written in short as

$$\mathbf{K}z_t + \mathbf{L}z_x = \nabla S_\gamma(z) - \gamma\mathbf{K}z \quad (4.4)$$

with

$$\mathbf{K} = \begin{bmatrix} 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, \quad z = \begin{bmatrix} u \\ v \\ w \\ p \end{bmatrix},$$

and  $S_\gamma(z) = \frac{1}{2}(2\gamma(uv + wp) + v^2 - w^2 + cu^2 - cp^2)$ .

**Example 4.2.** Consider a modified Burgers' equation

$$u_t + uu_x = -2\gamma u. \quad (4.5)$$

This equation can also be put in the form of (4.2)

$$\mathbf{K}z_t + \mathbf{L}z_x = \nabla S(z) - 2\gamma \mathbf{K}z \quad (4.6)$$

with

$$\mathbf{K} = \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix}, \quad z = \begin{bmatrix} u \\ v \\ w \end{bmatrix},$$

and  $S(z) = -uv + \frac{u^3}{3}$ .

**Example 4.3.** Consider the following generalization of the NLS equation of Table 1.1

$$i\psi_t + \psi_{xx} + V'(|\psi|^2)\psi + 2i\gamma\psi = F(t), \quad (4.7)$$

where  $\psi = \psi(x, t)$  is a complex valued wave function of space  $x$  and time  $t$ , the nonnegative real number  $\gamma$  is a damping parameter, and subscripts denote the usual partial derivatives. The time dependent term on the right hand side  $F(t)$  is an external driving force. We can put eq. (4.7) in the form of (4.2) as

$$\mathbf{K}z_t + \mathbf{L}z_x = \nabla S(z) - 2\gamma \mathbf{K}z + \mathbf{F}(t), \quad (4.8)$$

with

$$\mathbf{K} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{L} = \begin{bmatrix} 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, z = \begin{bmatrix} v \\ w \\ p \\ q \end{bmatrix}, \mathbf{F}(t) = \begin{bmatrix} -\Re F(t) \\ -\Im F(t) \\ 0 \\ 0 \end{bmatrix},$$

$S = \frac{1}{2}(p^2 + q^2 + V(v^2 + w^2))$ , and  $\psi = v + iw$ . Here  $\Re$  and  $\Im$  denote real and imaginary parts, respectively, of a complex number.

Notice that the damping part of the equations in the last three examples is absorbed by either the function  $S$  or the term involving parameter  $a$  or both. We also remark that setting  $\gamma = F = 0$  in these three examples gives a multi-symplectic formulation (4.1) of the corresponding PDEs instead. One advantage of writing PDEs in the form of eq. (4.2) is that the equation automatically satisfies certain local conservation laws. Setting  $a = b = 0$  and  $\mathbf{F} = 0$  in these conservation laws gives corresponding conservation laws for eq. (4.1).

#### 4.1.1 Local conservation laws

A conservation law of a PDE assumes the following form

$$\partial_t P + \partial_x Q = 0$$

where  $P$  and  $Q$  depend on  $z$ . For example, a conservation law associated with the modified Burgers' equation (4.5) is

$$\partial_t(e^{2\gamma t}u) + \partial_x(e^{2\gamma t}\frac{1}{2}u^2) = 0.$$

Indeed, the left hand side of this equation gives

$$\begin{aligned}
\partial_t(e^{2\gamma t}u) + \partial_x(e^{2\gamma t}\frac{1}{2}u^2) &= e^{2\gamma t}u_t + 2\gamma e^{2\gamma t}u + e^{2\gamma t}uu_x \\
&= e^{2\gamma t}(u_t + 2\gamma u + uu_x) \\
&= 0,
\end{aligned}$$

because  $u$  solves eq. (4.5).

Local conservation laws of multi-conformal-symplecticness and conformal momentum will be derived in this section. Let us begin by defining  $\mathbf{L}_+$  and  $\mathbf{L}_-$  such that

$$\mathbf{L} = \mathbf{L}_+ + \mathbf{L}_- \text{ and } \mathbf{L}_+^T = -\mathbf{L}_-.$$

Then eq. (4.2) becomes

$$\mathbf{K}z_t + \mathbf{L}_+z_x + \mathbf{L}_-z_x = \nabla S(z) - \frac{a}{2}\mathbf{K}z + \mathbf{F}(t). \quad (4.9)$$

Local conservation laws for this PDE can be derived in a manner analogous to the conservation laws of eqs. (4.1) and (4.2). Indeed, the variational equation associated with this PDE is given by

$$\mathbf{K}dz_t + \mathbf{L}_+dz_x + \mathbf{L}_-dz_x = S_{zz}(z)dz - \frac{a}{2}\mathbf{K}dz.$$

Taking the wedge product of this equation with  $dz$  and using

$$\begin{aligned}
dz \wedge \mathbf{K}dz_t &= \partial_t(\frac{1}{2}(dz \wedge \mathbf{K}dz)), \\
dz \wedge \mathbf{L}_-dz_x + dz \wedge \mathbf{L}_+dz_x &= \partial_x(dz \wedge \mathbf{L}_+dz), \\
dz \wedge S_{zz}(z)dz &= 0,
\end{aligned}$$

we get the following local conservation law

$$\partial_t\omega + \partial_x\kappa = -a\omega, \quad (4.10)$$

which can be written more compactly [39] as

$$\partial_t(e^{at}\omega) + \partial_x(e^{at}\kappa) = 0,$$

where  $\omega = \frac{1}{2}dz \wedge \mathbf{K}dz$  and  $\kappa = dz \wedge \mathbf{L}_+dz$ . That is, differential two-forms  $\omega$  and  $\kappa$  associated with PDE (4.9) satisfy a linear partial differential equation. Another way to interpret the above equation is that changes in space and time mutually annihilate each other. This local conservation law is called a *multi-conformal-symplectic conservation law*, after which the PDE is referred to as a *multi-conformal-symplectic PDE*. When  $a = 0$  and  $\mathbf{F}(t) = \mathbf{0}$ , eq. (4.9) reduces to eq. (4.1) which satisfies a *multi-symplectic-conservation law* given by eq. (4.10) with  $a = 0$ . It can be shown that eq. (4.10) holds for PDE (4.2) with the following differential 2-forms also:

$$\omega = dz \wedge \mathbf{K}dz \text{ and } \kappa = dz \wedge \mathbf{L}dz.$$

Therefore eqs. (4.3), (4.5) and (4.7) are MCS PDEs.

Special form of eqs. (4.2) and (4.9) guarantee another local conservation law. This conservation law, analogous to momentum conservation law for eq. (4.1), can be obtained for eq. (4.9) when  $\mathbf{F} = 0$ . Indeed, taking the inner product of  $z_x$  with eq. (4.9), we get

$$\langle \mathbf{K}z_t, z_x \rangle = \partial_x S(z) - \frac{a}{2} \langle \mathbf{K}z, z_x \rangle,$$

if and only if  $\mathbf{F} = 0$ . In this equation, using

$$\partial_t \langle \mathbf{K}z, z_x \rangle - \partial_t \langle z_t, \mathbf{K}z \rangle = 2 \langle \mathbf{K}z_t, z_x \rangle,$$

we get

$$\partial_t I + \partial_x G = -aI, \tag{4.11}$$

where  $G = -S(z) - \frac{1}{2} \langle z_t, \mathbf{K}z \rangle$  and  $I = \frac{1}{2} \langle z, \mathbf{K}z_x \rangle$ . Equation (4.11) is referred to as the *conformal*

*momentum conservation law*. From this local conservation law, one can obtain the following global conservation law by integrating in space and assuming vanishing or periodic boundary conditions, so that

$$\partial_x \int G dx = 0$$

and

$$\partial_t \int I dx = -a \int I dx. \quad (4.12)$$

This implies that the property  $\mathcal{I} = \int I dx$  is a conformal invariant (Definition 1.1) of the MCS PDE. It is desirable for a numerical integrator of a MCS PDE to preserve as many of these conservation laws as possible.

#### 4.1.2 Multi-conformal-symplectic numerical methods

Numerical methods which satisfy a discrete version of the multi-conformal-symplectic conservation law, eq. (4.10), are called *multi-conformal-symplectic numerical methods*. We present two examples of MCS methods for MCS PDE (4.2). The first one of these examples is for a specific MCS PDE, eq. (4.7) and the second example is for the general PDE, eq. (4.2). Both these equations will be shown to satisfy a discrete version of the MCS conservation law (4.10).

**Example 4.4.** Discretizing eq. (4.8) in space using notation of eq. (1.14), we obtain

$$\mathbf{K}\partial_t z^n + \mathbf{L}_+ D_x z^n + \mathbf{L}_- D_x T_x z^n = \nabla S(z^n) - 2\gamma \mathbf{K} z^n + \mathbf{F}(t), \quad (4.13)$$



where

$$\mathbf{L}_+ = \begin{bmatrix} 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \text{ and } \mathbf{L}_- = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

Discretizing eq. (4.13) in time with the method of eq. (2.11) we obtain

$$\mathbf{K}D_t^\gamma z + \mathbf{L}_+ D_x A_t^\gamma z + \mathbf{L}_- D_x T_x A_t^\gamma z = \nabla S(A_t^\gamma z) + \mathbf{F}(A_t t) \quad (4.14)$$

where we have suppressed superscripts in the numerical solution  $z^{n,i}$  for brevity. Notice that both spatial and time indices are denoted by superscripts for PDE numerical solutions as compared to subscripts for ODE numerical solutions. In terms of the original variable, this equation becomes

$$\mathbf{i}D_t^\gamma \psi + A_t^\gamma \delta_x^2 \psi + V'(|A_t^\gamma \psi|^2) A_t^\gamma \psi = F(A_t t), \quad (4.15)$$

which should be compared with its continuous counterpart, eq. (4.7).

That (4.14) is a multi-conformal-symplectic method can be seen from the following computation.

The variational equation associated with this method is

$$\mathbf{K}D_t^\gamma dz + \mathbf{L}_+ D_x A_t^\gamma dz + \mathbf{L}_- D_x T_x A_t^\gamma dz = S_{zz}(A_t^\gamma z) A_t^\gamma dz$$

where  $S_{zz}$  is the Hessian of  $S$ . Taking the wedge product of  $A_t^\gamma dz$  with this equation we get

$$A_t^\gamma dz \wedge \mathbf{K}D_t^\gamma dz + A_t^\gamma dz \wedge \mathbf{L}_+ D_x A_t^\gamma dz + A_t^\gamma dz \wedge \mathbf{L}_- D_x T_x A_t^\gamma dz = A_t^\gamma dz \wedge S_{zz}(A_t^\gamma z) A_t^\gamma dz. \quad (4.16)$$

Using Lemma 1.8 and the symmetry of  $S_{zz}$  we get

$$\begin{aligned} A_t^\gamma dz \wedge \mathbf{K} D_t^\gamma dz &= D_t^{2\gamma} \left( \frac{1}{2} (dz \wedge \mathbf{K} dz) \right), \\ A_t^\gamma dz \wedge \mathbf{L}_+ D_x A_t^\gamma dz + A_t^\gamma dz \wedge \mathbf{L}_- D_x T_x A_t^\gamma dz &= D_x (A_t^\gamma T_x dz \wedge \mathbf{L}_+ A_t^\gamma dz), \\ A_t^\gamma dz \wedge S_{zz} (A_t^\gamma z) A_t^\gamma dz &= 0. \end{aligned}$$

These equations along with eq. (4.16) give

$$D_t^{2\gamma} \left( \frac{1}{2} (dz \wedge \mathbf{K} dz) \right) + D_x (A_t^\gamma T_x dz \wedge \mathbf{L}_+ A_t^\gamma dz) = 0$$

which is a discrete version of the multi-conformal-symplectic conservation law (4.10).

**Example 4.5.** For our next example, consider the following method for eq. (4.2) with  $\mathbf{F} = 0$ .

$$\mathbf{K}(D_t^{a/4} A_x z) + \mathbf{L}(D_x A_t^{a/4} z) = \nabla S(A_x A_t^{a/4} z). \quad (4.17)$$

This method is obtained by using conformal symplectic method (2.11) in both space and time.

Variational equation for this method is

$$\mathbf{K}(D_t^{a/4} A_x dz) + \mathbf{L}(D_x A_t^{a/4} dz) = S_{zz} A_x A_t^{a/4} dz.$$

Similar to the previous example, one can get the following discrete version of eq. (4.10)

$$D_t^{a/2} (A_x dz \wedge \mathbf{K} A_x dz) + D_x \left( \mathbf{L} A_t^{a/4} dz \wedge A_t^{a/4} dz \right) = 0.$$

Therefore eq. (4.17) is an MCS integrator.

Of course, conformal symplectic numerical methods other than eq. (2.11) can also be used to discretize space and time to obtain MCS integrators. We explore some of these possibilities and use the two methods discussed in this section to solve PDEs in the next chapter.

## 4.2 Non-standard finite difference methods

Non-standard finite difference methods provide an alternative approach to discretizing differential equations. This approach uses a set of rules to discretize a given DE. The following rules are suggested to design a non-standard finite difference method:

- Denominator function of a discrete derivative is a more complicated function of time/space step-size than denominator function of a standard discrete derivative.
- Non-linear terms are modeled non-locally.
- Order of the discrete derivatives is exactly equal to the corresponding order of the derivatives in the differential equation.

In addition to these rules, a relationship between time and spatial step-size often exists to ensure stability of the method. These rules often result in operators that are *non-standard* i.e. they resemble the standard discrete derivative and averaging operators,  $D_\zeta$  and  $A_\zeta$  of eq. (1.14), but are not exactly the same. Although relatively new compared to standard finite difference methods, non-standard methods have been successfully used to discretize a multitude of ODEs and PDEs [35, 36]. Instead of structure preservation, this approach focuses on providing “best” solutions to a differential equation. For this reason, we will use a non-standard method to discretize a PDE for the purposes of comparison in this thesis.

For example, consider the modified Burgers’ equation (4.5)

$$u_t + uu_x = -2\gamma u.$$

Solving only the linear part

$$u_t = -2\gamma u,$$

one gets the solution

$$u(x, t) = e^{-2\gamma t} u(x, 0).$$

An exact finite difference method, for the linear part, which produces this solution is given by

$$\begin{aligned} u^{n,i+1} &= e^{-2\gamma\Delta t} u^{n,i}, \\ \iff \frac{u^{n,i+1} - u^{n,i}}{\left(\frac{1-e^{-2\gamma\Delta t}}{2\gamma}\right)} &= -2\gamma u^{n,i}. \end{aligned}$$

Notation  $\Delta x$  and  $\Delta t$  is used to denote spatial and time step-sizes in numerical methods for PDEs.

We use superscripts  $\{n, i\}$  for spatial and time indices, respectively. Modeling the nonlinear term  $uu_x$  nonlocally gives the following non-standard finite difference (NSFD) scheme for (4.5)

$$\frac{u^{n,i+1} - u^{n,i}}{\left(\frac{1-e^{-2\gamma\Delta t}}{2\gamma}\right)} + u^{n,i} \left( \frac{u^{n,i+1} - u^{n-1,i+1}}{\Delta x} \right) = -2\gamma u^{n,i}. \quad (4.18)$$

Numerical methods discussed in this chapter will be used to simulate numerical solutions of PDEs in the next chapter, where the efficacy of these methods will also be compared. Consistent with numerical experiments for ODEs, our focus will remain on structure preservation.

## CHAPTER 5: PDE APPLICATIONS AND EXPERIMENTS

We apply and compare the performance of numerical methods for PDEs discussed in the previous chapter and other methods that will be developed in this chapter. We begin by discretizing a linear PDE with the methods of Chapters 2 and 4 and prove their structure-preservation properties. Then we use some of these methods to discretize a nonlinear PDE and compare them against a non-standard finite difference method. In our final PDE experiment we compare some of the methods of Chapters 2 and 4 against a non-structure preserving method to highlight advantages of structure-preservation. Among these methods, MCS methods automatically satisfy certain conservation laws and other structure-preserving methods are shown to preserve structure by direct computations. In the following, superscripts  $n$  and  $i$  denote the spatial and temporal indices, respectively. Notations  $\Delta x$  and  $\Delta t$  denote spatial and time step sizes, respectively.

### 5.1 A damped Klein-Gordon equation

Consider the damped Klein-Gordon eq. (4.3) on the interval  $[-\pi, \pi]$  with periodic boundary conditions. Its exact solution is taken to be

$$u(x, t) = e^{-\gamma t} \cos(Kx - Wt), \quad W = \sqrt{K^2 + c - \gamma^2} \quad (5.1)$$

where  $K$  is the *wavenumber* and  $W$  is called the *frequency* of the wave. Constant  $\gamma$  is the damping parameter. Equation (4.3) can be written as a multi-conformal-symplectic PDE as in eq. (4.4). In the following, we describe different approaches to discretize this PDE.

### 5.1.1 Numerical solutions

We take two different approaches to design numerical solutions for eq. (5.1). Our first approach discretizes a spatially semi-discretized PDE with two different conformal symplectic PERK methods in time. This gives two explicit numerical schemes for the Klein-Gordon equation. The other approach uses an implicit MCS integrator from Chapter 4 to discretize the equation.

One can rewrite eq. (4.3) as a system of equations

$$u_t = v, \quad v_t = -(cu - u_{xx}) - 2\gamma v. \quad (5.2)$$

Using the central finite difference operator of eq. (1.14), we discretize this system using PERK methods (2.20) and (2.21) in time to get numerical methods

$$\begin{aligned} v^{n,i+1/2} &= e^{-\gamma\Delta t} v^{n,i} + \frac{\Delta t}{2} (\delta_x^2 - c) u^{n,i}, \\ u^{n,i+1} &= u^{n,i} + \Delta t v^{n,i+1/2}, \\ v^{n,i+1} &= e^{-\gamma\Delta t} \left[ v^{n,i+1/2} + \frac{\Delta t}{2} (\delta_x^2 - c) u^{n,i+1} \right], \end{aligned} \quad (5.3)$$

and

$$\begin{aligned} (1 + \gamma\Delta t/2)v^{n,i+1/2} &= e^{-\gamma\Delta t/2} v^{n,i} + \frac{\Delta t}{2} (\delta_x^2 - c) u^{n,i}, \\ (1 - \gamma\Delta t/2)u^{n,i+1} &= e^{-\gamma\Delta t/2} \left[ (1 + \gamma\Delta t/2)e^{-\gamma\Delta t/2} u^{n,i} + \Delta t v^{n,i+1/2} \right], \\ v^{n,i+1} &= e^{-\gamma\Delta t/2} \left[ (1 - \gamma\Delta t/2)v^{n,i+1/2} + \frac{\Delta t}{2} (\delta_x^2 - c) u^{n,i+1} \right], \end{aligned} \quad (5.4)$$

respectively. On the other hand, discretizing PDE (4.4) with the multi-conformal-symplectic method (4.17) we get

$$\mathbf{K}(D_t^{\gamma/2} A_x z) + \mathbf{L}(A_t^{\gamma/2} D_x z) = \nabla S_\gamma(A_x A_t^{\gamma/2} z). \quad (5.5)$$

In the following we simulate solutions and compare these three methods.

Figure 5.1 shows the absolute error in numerical solutions produced by methods of eqs. (5.3) to (5.5), denoted by CSV1, CSV2, and CIMP, respectively, where the exact solution is propagated with numerical frequency. Notice that the three numerical solutions are slightly out of phase with one another due to their different numerical frequencies. Notice also that the error decreases for all three methods as time increases, primarily because the solution is dissipating to zero, but the relative error remains close to  $10^{-2}$ .

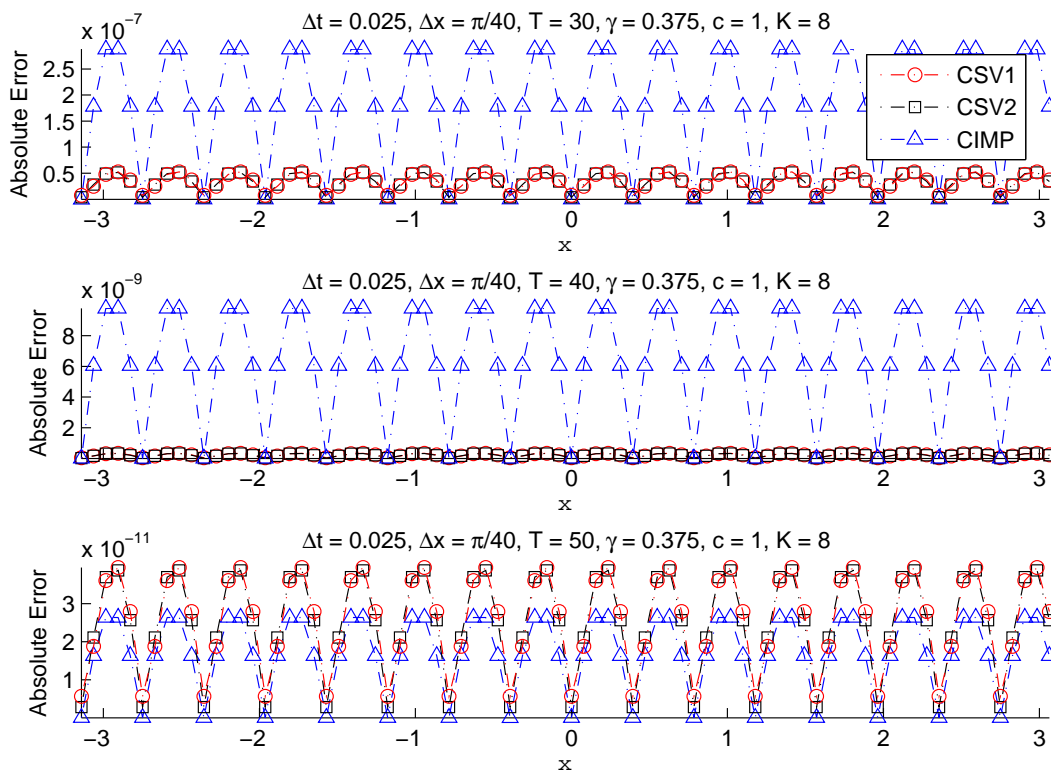


Figure 5.1: Error in the solution of (4.3) due to methods (5.3), (5.4) and (5.5). Parameter values are given in the figure title. The maximum value of the exact solution at time  $T = 50$  is approximately  $7 \times 10^{-9}$ .

Following the approach of [40] consider the function

$$d(t) = \ln \left( \max_{x \in [-\pi, \pi]} u(x, t) \right) + \gamma t \quad (5.6)$$

for measuring the drift in the rate of dissipation for the Klein-Gordon equation. Figure 5.2 shows that there is no drift in the rate of dissipation (5.6) for the three methods.

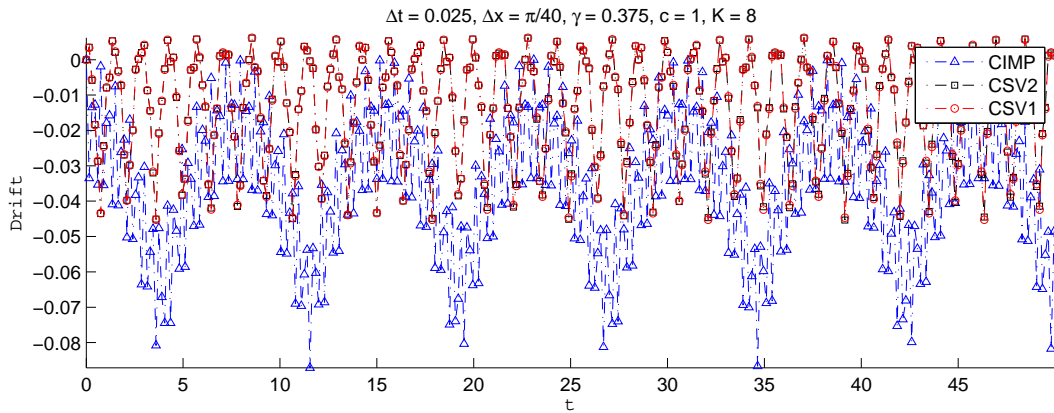


Figure 5.2: Drift in the rate of dissipation for the three methods (5.3), (5.4) and (5.5) with the parameter values mentioned in the figure title. Only every sixth drift vector component is plotted for clarity and CSV1 eclipses CSV2.

### 5.1.2 Structure-preservation

Consistent with the theme of this thesis, these methods preserve more than conformal symplecticity, which helps explain the practical advantages of the methods. To be specific, define the *momentum* for the Klein-Gordon eq. (5.2) by

$$I(t) = \int u_t u_x dx. \quad (5.7)$$



Under appropriate boundary conditions, it can be shown that the PDE has the conformal invariant

$$I(t) = e^{-2\gamma t} I(0). \quad (5.8)$$

In fact, all three methods preserve this property.

**Theorem 5.1.** *The numerical methods (5.3), (5.4), and (5.5) each preserve eq. (5.8).*

*Proof.* Notice that

$$\begin{aligned} (1 + \gamma_1 \Delta t/2)v^{n,i+1/2} &= e^{-\gamma_3 \Delta t/2} v^{n,i} - \frac{\Delta t}{2}(cu^{n,i} - \delta_x^2 u^{n,i}), \\ (1 - \gamma_1 \Delta t/2)u^{n,i+1} &= e^{-\gamma_1 \Delta t/2} \left[ (1 + \gamma_1 \Delta t/2)e^{-\gamma_1 \Delta t/2} u^{n,i} + \Delta t v^{n,i+1/2} \right], \\ e^{\gamma_2 \Delta t/2} v^{n,i+1} &= e^{-\gamma_1 \Delta t/2} \left[ (1 - \gamma_1 \Delta t/2)v^{n,i+1/2} - \frac{\Delta t}{2}(cu^{n,i+1} - \delta_x^2 u^{n,i+1}) \right], \end{aligned}$$

gives method (5.3) on setting  $\gamma_1 = 0, \gamma_2 = 2\gamma, \gamma_3 = 2\gamma$  and method (5.4) on setting  $\gamma_1 = \gamma, \gamma_2 = 0, \gamma_3 = \gamma$ . Using this combined method and Lemma 1.8 we see that

$$\begin{aligned} \sum v^{n,i+1} \delta_x u^{n,i+1} &= e^{-(\gamma_1 + \gamma_2) \Delta t/2} (1 - \gamma_1 \Delta t/2) \sum v^{n,i+1/2} \delta_x u^{n,i+1} \\ &= e^{-(3\gamma_1 + \gamma_2) \Delta t/2} (1 + \gamma_1 \Delta t/2) \sum v^{n,i+1/2} \delta_x u^{n,i} \\ &= e^{-(3\gamma_1 + \gamma_2 + \gamma_3) \Delta t/2} \sum v^{n,i} \delta_x u^{n,i} \\ &= e^{-2\gamma \Delta t} \sum v^{n,i} \delta_x u^{n,i} \end{aligned}$$

for both methods. Here  $\sum$  denotes summation with respect to the spatial index over all the spatial grid points.

Also notice that eq. (5.5) can be re-written in terms of the original variable  $u$  as

$$(D_t^{\gamma/2})^2 A_x^2 u - D_x^2 (A_t^{\gamma/2})^2 u = (\gamma^2 - c) A_x^2 (A_t^{\gamma/2})^2 u$$

where  $A_x^2 u = A_x A_x u$  etc. Now observe that, using Lemma 1.8,

$$\begin{aligned}
0 &= D_t^\gamma \sum D_t^{\gamma/2} A_x^2 u \cdot \delta_x A_t^{\gamma/2} A_x^2 u \\
&= \sum (D_t^{\gamma/2})^2 A_x^2 u \cdot \delta_x (A_t^{\gamma/2})^2 A_x^2 u + \sum D_t^{\gamma/2} A_t^{\gamma/2} A_x^2 u \cdot \delta_x D_t^{\gamma/2} A_t^{\gamma/2} A_x^2 u \\
&= \sum \left( D_x^2 (A_t^{\gamma/2})^2 u + (\gamma^2 - c) A_x^2 (A_t^{\gamma/2})^2 u \right) \delta_x (A_t^{\gamma/2})^2 A_x^2 u \\
&= \sum \left( D_x^2 (A_t^{\gamma/2})^2 u \right) \delta_x (A_t^{\gamma/2})^2 A_x^2 u
\end{aligned}$$

because for a periodic sequence  $U$ ,

$$\begin{aligned}
\sum D_x^2 U \cdot \delta_x A_x^2 U &= \frac{1}{4\Delta x^2} \sum (U^{n+2} - 2U^{n+1} + U^n) \cdot \delta_x (U^{n+2} + 2U^{n+1} + U^n) \\
&= \frac{1}{2\Delta x^2} \sum U^{n+2} \delta_x U^{n+1} + \frac{1}{4\Delta x^2} \sum U^{n+2} \delta_x U^n \\
&\quad - \frac{1}{2\Delta x^2} \sum U^{n+1} \delta_x U^{n+2} - \frac{1}{2\Delta x^2} \sum U^{n+1} \delta_x U^n \\
&\quad + \frac{1}{4\Delta x^2} \sum U^n \delta_x U^{n+2} + \frac{1}{2\Delta x^2} \sum U^n \delta_x U^{n+1} \\
&= \frac{1}{\Delta x^2} \sum U^{n+2} \delta_x U^{n+1} - \frac{1}{\Delta x^2} \sum U^{n+1} \delta_x U^n = 0.
\end{aligned}$$

Therefore,  $P^{i+1} = e^{-2\gamma\Delta t} P^i$  for

$$P^i = \sum D_t^{\gamma/2} A_x^2 u^{n,i} \cdot \delta_x A_t^{\gamma/2} A_x^2 u^{n,i}.$$

□

We define the residual  $r^i$  in preserving conformal momentum as

$$r^i = \ln \left( \frac{I^{i+1}}{I^i} \right) + 2\gamma\Delta t,$$

where

$$I^i = \sum_n v^{n,i} \delta_x u^{n,i} \quad \text{and} \quad I^i = \sum D_t^{\gamma/2} A_x^2 u^{n,i} \cdot \delta_x A_t^{\gamma/2} A_x^2 u^{n,i},$$

for Störmer-Verlet type methods (5.3), (5.4), and implicit midpoint type method (5.5) respectively.

Upon plotting  $I^i$  alongside respective residuals  $r^i$  for the three methods (5.3), (5.4) and (5.5), one obtains Figure 5.3, verifying momentum preservation property of the methods. (Only some of the data points are plotted to prevent overcrowding.)

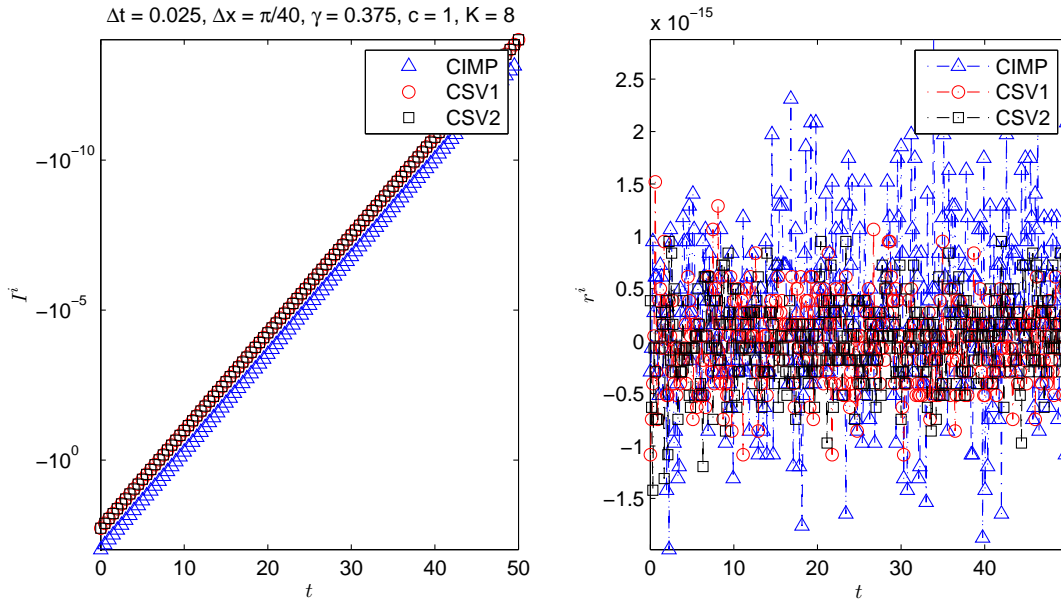


Figure 5.3: Total conformal momentum  $I^i$  and residual  $r^i$  due to (5.3), (5.4) and (5.5) with the parameters mentioned in the figure title

## 5.2 A Modified Burgers' Equation

In this experiment, we consider application of ERK and PERK methods to a nonlinear equation, which we refer to as a modified Burgers' equation, given by eq. (4.5):

$$u_t + uu_x = -2\gamma u,$$

on the interval  $[-\pi, \pi]$  with periodic boundary conditions. This is a fundamental equation that models physical phenomenon in fluid mechanics, acoustics, traffic flow etc. It is also one of the

simplest PDEs whose solutions develop shock: spatial derivative of the solution becomes infinite in finite time. Formation of shocks makes this PDE a challenging equation to solve numerically. We derive three numerical methods for this nonlinear equation. Our first method for this PDE is based on method (2.21). The second method is a multi-conformal-symplectic method obtained from discretizing eq. (4.6) with conformal symplectic method (2.11) in both space and time. The third method is obtained from discretizing the PDE with a non-standard finite difference method obtained by using the rules of Section 4.2.

### 5.2.1 Numerical solutions

To test our methods on this nonlinear problem, we discretize the space  $x \in [-L, L]$ ,  $L = \pi$ , by introducing a uniform spatial grid  $[x_1, x_2, \dots, x_M]$  with gridsize  $\Delta x$  such that  $x_1 = -L$ ,  $x_M = L$ , and  $M$  is even. Then we approximate  $u$  by  $u^n = u(x_n)$ ,  $n = 1, 2, \dots, M$ , with periodic boundary conditions  $u^{n+M} = u^n$ , and define the following vectors.

$$v = [u^1, u^3, u^5, \dots, u^{M-1}]^T \quad \text{and} \quad w = [u^2, u^4, u^6, \dots, u^M]^T.$$

Given this spatial decomposition of the solution vector  $u^n$ , one can semi-discretize (4.5) to get the following system of ODEs

$$\begin{aligned} \frac{dw^n}{dt} &= -\partial_x^+ ((v^n)^2/2) - 2\gamma w^n, \\ \frac{dv^n}{dt} &= -\partial_x^- ((w^n)^2/2) - 2\gamma v^n, \end{aligned} \tag{5.9}$$

where the superscript  $n$  on the vectors  $v$  and  $w$  is the index of these vectors, e.g.  $v^2 = u^3$  etc., and  $\partial_x^+$  and  $\partial_x^-$  are one-half of standard forward and backward difference operators, respectively. Even-odd splitting of the dependent variable  $u$  was suggested in [1], and this is the approach used

here. Notice that the system is a conformal Hamiltonian system of the form

$$z_t^n = \mathbf{D} \nabla_{z^n} H(z^n) - 2\gamma z^n \text{ with } \mathbf{D} = \begin{bmatrix} \mathbf{0} & \partial_x^- \\ \partial_x^+ & \mathbf{0} \end{bmatrix},$$

which is constant and skew-symmetric,  $z^n = [v^n, w^n]^T$  and  $H(z^n) = -\frac{1}{6}((w^n)^3 + (v^n)^3)$ . Although structure matrix  $\mathbf{D}$  of the system above is not same as the matrix  $\mathbf{J}^{-1}$  of (2.32), one can still show that this system is conformal symplectic and hence it is desirable to apply conformal symplectic methods on the system. Using method (2.21) to discretize time in system (5.9) we get the following method

$$\begin{aligned} w^{n,i+\frac{1}{2}} &= e^{-\gamma\Delta t} w^{n,i} - \frac{\Delta t}{2} \partial_x^+ \frac{(v^{n,i})^2}{2}, \\ v^{n,i+1} &= e^{-\gamma\Delta t} \left( e^{-\gamma\Delta t} v^{n,i} - \Delta t \partial_x^- \frac{\left(w^{n,i+\frac{1}{2}}\right)^2}{2} \right), \\ w^{n,i+1} &= e^{-\gamma\Delta t} \left( w^{n,i+\frac{1}{2}} - \frac{\Delta t}{2} \partial_x^+ \frac{(v^{n,i+1})^2}{2} \right). \end{aligned} \quad (5.10)$$

Let us now discretize multi-conformal-symplectic formulation (4.6) of the modified Burgers' equation with method (4.17). Doing so, gives the following multi-conformal-symplectic method

$$\mathbf{K}(D_t^\gamma A_x z) + \mathbf{L}(A_t^\gamma D_x z) = \nabla S(A_x A_t^\gamma z), \quad (5.11)$$

which gives the following method in terms of the original variable.

$$A_t^\gamma D_t^\gamma A_x^2 u + \frac{1}{2} A_t^\gamma D_x (A_x A_t^\gamma u)^2 = 0.$$

An equivalent one-step method is

$$D_t^\gamma A_x u + \frac{1}{2} D_x (A_t^\gamma u)^2 = 0 \quad (5.12)$$

which should be compared with (4.5).

Another discrete model that can approximate solutions of modified Burgers' equation is the following non-standard finite difference method (NSFD) (4.18) (cf. [35])

$$\frac{u^{n,i+1} - u^{n,i}}{\left(\frac{1-e^{-2\gamma\Delta t}}{2\gamma}\right)} + u^{n,i} \left(\frac{u^{n,i+1} - u^{n-1,i+1}}{\Delta x}\right) = -2\gamma u^{n,i},$$

which is able to exactly reproduce certain solutions of the equation. Since this method is, in some way, also structure-preserving, it provides an interesting comparison to the conformal symplectic methods. The three methods given by eqs. (4.18), (5.10) and (5.12) will be referred to as NSFD, CSV2, and CIMP, respectively, in the following experiments.

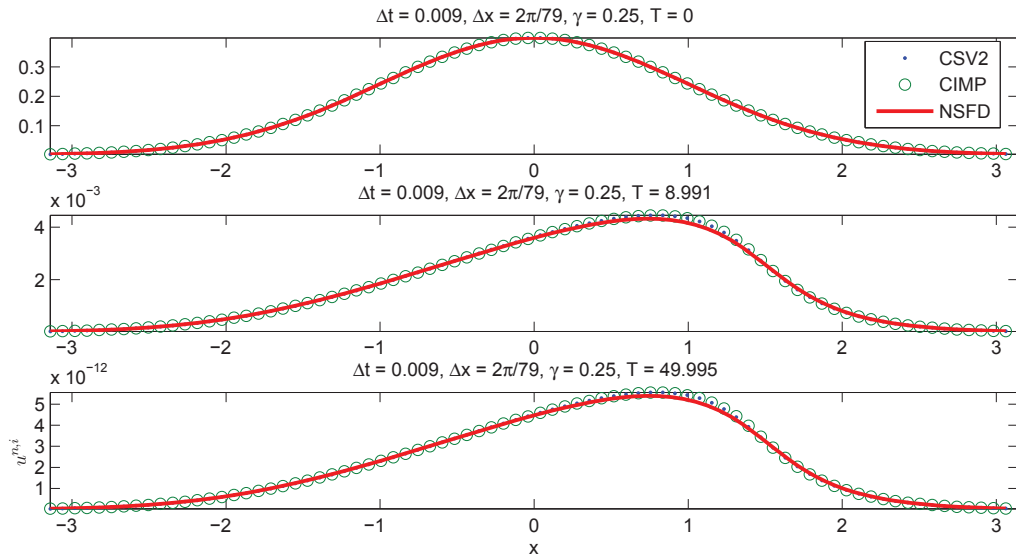


Figure 5.4: Snapshots of the numerical solution of eq. (4.5) using (5.10), (5.12) and (4.18) at different times.

Numerical solutions of eq. (4.5) at different times are given in Figure 5.4. The initial condition is taken to be a normal probability distribution function with mean 0 and standard deviation 1. As time progresses, the waveform becomes steeper, and the NSFD solution is damping at a different rate than the other two methods. In fact, the methods (5.10) and (5.12) are more accurate in this

case, and the results of the following subsection explain why.

### 5.2.2 Structure-preservation

Consider the Casimir

$$C[u] = \int u \, dx \quad (5.13)$$

which corresponds to mass of the wave at time  $t$  and is an integral of motion for the eq. (4.5) with  $\gamma = 0$ . Differentiating eq. (5.13) with respect to  $t$ , we have

$$\frac{dC}{dt} = \int u_t \, dx = - \int uu_x \, dx - 2\gamma \int u \, dx.$$

Now using integration by parts and periodic boundary conditions we get

$$\int_{-\pi}^{\pi} uu_x \, dx = uu \Big|_{-\pi}^{\pi} - \int_{-\pi}^{\pi} u_x u \, dx, \quad \implies \quad \int_{-\pi}^{\pi} uu_x \, dx = 0.$$

Therefore

$$\frac{dC}{dt} = -2\gamma C \quad \iff \quad C(t) = C(t_0)e^{-2\gamma(t-t_0)}, \quad (5.14)$$

which is a conformal property of eq. (4.5).

**Theorem 5.2.** *Methods (5.10) and (5.12) preserve eq. (5.14), but (4.18) does not.*

*Proof.* Define the discrete Casimir function (5.13) by  $C^i = \sum_n u^{n,i}$ . For method (5.10), we get

$$\begin{aligned} C^{i+1} &= \sum_{n=1}^M u^{n,i+1} = \sum_{n=1}^{M/2} (v^{n,i+1} + w^{n,i+1}) \\ &= \sum_{n=1}^{M/2} e^{-2\gamma\Delta t} (v^{n,i} + w^{n,i}) = e^{-2\gamma\Delta t} C^i, \end{aligned}$$

i.e. method (5.10) preserves eq. (5.14). Summing eq. (5.12) over the spatial index  $n$ , we get

$$\begin{aligned} \sum_n D_t^\gamma A_x u = 0, & \implies \sum_n A_x u^{n,i+1} = e^{-2\gamma\Delta t} \sum_n A_x u^{n,i}, \\ & \implies \sum_n u^{n,i+1} = e^{-2\gamma\Delta t} \sum_n u^{n,i}, \end{aligned}$$

i.e. method (5.12) also preserves eq. (5.14). Similarly, summing eq. (4.18) over the spatial index  $n$  we get

$$\begin{aligned} \frac{2\gamma}{1 - e^{-2\gamma\Delta t}} \sum_n u^{n,i+1} &= \frac{2\gamma}{1 - e^{-2\gamma\Delta t}} \sum_n u^{n,i} - 2\lambda \sum_n u^{n,i} \\ &\quad - \frac{1}{\Delta x} \sum_n (u^{n,i} u^{n,i+1} - u^{n,i} u^{n-1,i+1}), \end{aligned}$$

which implies

$$\sum_n u^{n,i+1} = e^{-2\gamma\Delta t} \sum_n u^{n,i} - \frac{1 - e^{-2\gamma\Delta t}}{2\gamma\Delta x} \sum_n (u^{n,i} u^{n,i+1} - u^{n,i} u^{n-1,i+1}).$$

Therefore, method (4.18) does not preserve eq. (5.14). □

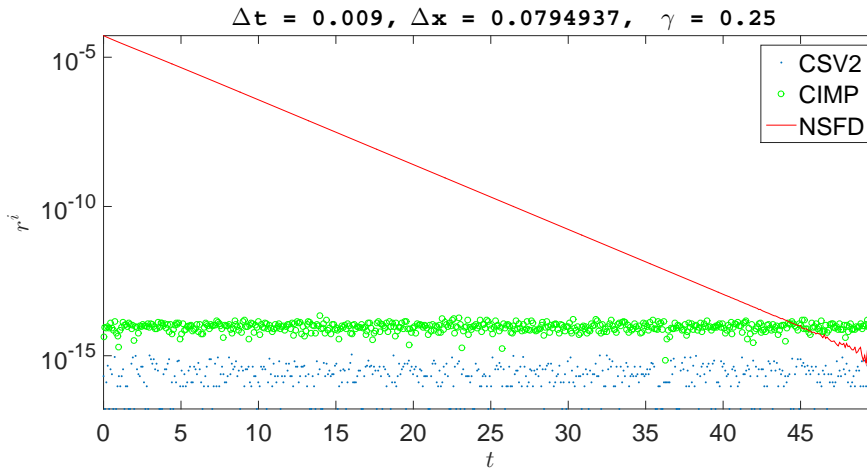


Figure 5.5: Residual (5.15) due to conformal symplectic methods (5.10) and (5.12) and NSFD (4.18).



To measure the error in preservation of (5.14) we define the residual  $r^i$

$$r^i = \ln \left( \frac{C^{i+1}}{C^i} \right) + 2\gamma\Delta t, \quad (5.15)$$

which is plotted in Figure 5.5, verifying that the conformal symplectic methods preserve dissipation of mass, but NSFD does not.

### 5.3 Damped driven nonlinear Schrödinger equation

The damped driven nonlinear Schrödinger eq. (4.7) is discretized by structure-preserving methods in this section. The time dependent term on the right hand side  $F(t)$  is an external driving force which, along with the nonlinear term  $V'(|\psi|^2)\psi$  and damping term  $2i\gamma\psi$ , induces chaos in certain parameter regimes.

Undamped ( $\gamma = 0$ ) and unforced ( $F(t) = 0$ ) eq. (4.7) is a Hamiltonian system and is referred to as *integrable NLS*. Every symmetry of a Hamiltonian system results in a conservation law according to Noether's theorem. Symmetries of the integrable NLS about space, phase, and time result in conservation of the quantities

$$I_1 = \Im \int \bar{\psi}\psi_x dx, \quad I_2 = \int |\psi|^2 dx, \quad I_3 = \int -|\psi_x|^2 + V(|\psi|^2) dx \quad (5.16)$$

with periodic or vanishing boundary conditions. Quantities  $I_1$ ,  $I_2$ , and  $I_3$  are referred to as momentum, norm, and energy, respectively. Bar over a complex variable denotes complex conjugate of the variable.

Equation (4.7) with  $F = 0$  is referred to as *damped NLS*. Although Noether's theorem is not applicable to non-conservative systems, one can nonetheless show by direct computations that a

solution of damped NLS satisfies

$$\partial_t(e^{4\gamma t} I_1) = 0, \quad \partial_t(e^{4\gamma t} I_2) = 0 \quad (5.17)$$

with periodic or vanishing boundary conditions. Indeed,

$$\begin{aligned} \partial_t \int \bar{\psi} \psi_x dx &= \int \bar{\psi}_t \psi_x dx + \int \bar{\psi} \psi_{xt} dx \\ &= \int \bar{\psi}_t \psi_x dx + [\bar{\psi} \psi_t] - \int \bar{\psi}_x \psi_t dx. \end{aligned}$$

Here, we have used integration by parts and  $[\cdot]$  denotes the boundary terms. Assuming periodic boundary conditions, this equation gives

$$\partial_t \int \bar{\psi} \psi_x dx = \int \bar{\psi}_t \psi_x dx - \int \bar{\psi}_x \psi_t dx. \quad (5.18)$$

In the right hand side of this equation, using eq. (4.7) with  $F = 0$ , we get

$$\begin{aligned} \int \bar{\psi}_t \psi_x dx - \int \bar{\psi}_x \psi_t dx &= \mathbf{i} \left[ \int -V'(|\psi|^2) (\bar{\psi} \psi_x + \psi \bar{\psi}_x) dx + 2\mathbf{i}\gamma \int (\bar{\psi} \psi_x - \psi \bar{\psi}_x) dx \right] \\ &= \mathbf{i} \left[ \int -V'(|\psi|^2) (2\Re\{\bar{\psi} \psi_x\}) dx + 2\mathbf{i}\gamma \int (2\mathbf{i}\Im\{\bar{\psi} \psi_x\}) dx \right] \end{aligned}$$

Using this equation in eq. (5.18) and taking imaginary part of both sides we get

$$\partial_t \Im \int \bar{\psi} \psi_x dx = -4\gamma \Im \int \bar{\psi} \psi_x dx.$$

This implies that

$$\partial_t I_1 = -4\gamma I_1 \iff \partial_t(e^{4\gamma t} I_1) = 0.$$

The second equation of (5.17) was proved in Example 1.5.

Substituting  $V(\eta) = \lambda\eta$  and  $F = 0$  in eq. (4.7) we obtain a *linear Schrödinger equation*:

$$\mathbf{i}\psi_t + \psi_{xx} + \lambda\psi + 2\mathbf{i}\gamma\psi = 0. \quad (5.19)$$

A solution of this equation satisfies

$$\partial_t(e^{4\gamma t} I_k) = 0 \text{ for } k = 1, 2, 3, \quad (5.20)$$

where  $I_k$  are as in eq. (5.16) with  $V(\eta) = \lambda\eta$ . It follows that preservation of dissipative properties of linear and nonlinear Schrödinger equations by numerical methods is desirable.

Damped driven NLS is an example of an infinite-dimensional dynamical system. Long time behavior of eq. (4.7) has been studied extensively both analytically and numerically. Its global attractor attracts all nearby trajectories on a compact bounded set. It was shown in [18] that chaotic attractors exist and they are confined in a finite dimensional space. Authors of [11] used high order RK methods to conduct numerical experiments showing a quasi-periodic route to chaos in the dynamical system.

### 5.3.1 Numerical solutions

In [31], authors take a discrete variational derivatives route to derive a linearly implicit finite difference scheme that inherits an energy conservation or dissipation property of a complex valued PDE such as integrable NLS and closely related Ginzburg-Landau equation. In [24], authors construct a symplectic geometric integrator by generalizing the generating functions approach and comparing it to a multi-symplectic geometric integrator for integrable NLS. The multi-symplectic integration technique was generalized to dissipative PDEs in [40] where authors proposed a norm-preserving multi-conformal-symplectic integrator for a damped NLS. Authors of [29] propose a variational integrator, which is naturally multi-symplectic, by first defining a Lagrangian function for a variable coefficient integrable NLS.

As far as we know, structure-preserving methods for damped driven NLS (DDNLS) have not been previously suggested in the literature. In this section, we use an appropriate spatial discretization

and two conformal symplectic ERK methods for time discretizations to assemble MCS methods for the DDNLS. We compare these two methods against a method which employs a closely related discretization in space and time but the resulting method is not structure-preserving. The aim is to construct structure preserving numerical methods for DDNLS, discuss their structure preserving properties, and implement them to show their effectiveness. In the following, we discretize multi-conformal-symplectic formulation (4.9) of DDNLS eq. (4.7). We discretize this formulation with integrating factor method (2.11), exponential time differencing (2.13), and implicit midpoint methods in time.

### 5.3.1.1 Integrating factor method

Discretizing (4.7) with (2.11) we get method (4.14) or equivalently method (4.15). This method will be referred to as the IF method in the numerical plots. Let us point out some salient features of this method.

- (i) That (4.14) is a MCS integrator was shown in Example 4.4.
- (ii) The method preserves the invariant  $e^{4\gamma t} I_2$  for damped NLS. This can be shown from a conformal norm conservation law [40], similar to the conformal momentum conservation law eq. (4.11). Alternatively, one can show preservation of  $e^{4\gamma t} I_2$  by doing computations analogous to their continuous counterpart in Section 5.3. Indeed, assuming  $F(A_t t) = 0$  and multiplying eq. (4.15) with  $A_t^\gamma \bar{\psi}$ , one gets

$$\mathbf{i} D_t^\gamma \psi A_t^\gamma \bar{\psi} + A_t^\gamma \delta_x^2 \psi A_t^\gamma \bar{\psi} + V'(|A_t^\gamma \psi|^2) A_t^\gamma \psi A_t^\gamma \bar{\psi} = 0$$

which gives

$$\sum_n \left( \mathbf{i} D_t^\gamma \psi A_t^\gamma \bar{\psi} + A_t^\gamma \delta_x^2 \psi A_t^\gamma \bar{\psi} + V'(|A_t^\gamma \psi|^2) A_t^\gamma \psi A_t^\gamma \bar{\psi} \right) = 0. \quad (5.21)$$

The first term of this equation gives, using eq. (1.14),

$$\begin{aligned} \sum_n \mathbf{i} D_t^\gamma \psi^{n,i} A_t^\gamma \bar{\psi}^{n,i} &= \sum_n \frac{\mathbf{i}}{2\Delta t} (e^{\gamma\Delta t} \psi^{n,i+1} - e^{-\gamma\Delta t} \psi^{n,i}) (e^{\gamma\Delta t} \bar{\psi}^{n,i+1} + e^{-\gamma\Delta t} \bar{\psi}^{n,i}) \\ &= \sum_n \frac{1}{2\Delta t} \left( \mathbf{i} (e^{2\gamma\Delta t} |\psi^{n,i+1}|^2 - e^{-2\gamma\Delta t} |\psi^{n,i}|^2) - 2\Im(\psi^{n,i+1} \bar{\psi}^{n,i}) \right), \end{aligned} \quad (5.22)$$

the second term gives, using eq. (1.14) and periodic boundary conditions,

$$\begin{aligned} \sum_n A_t^\gamma \delta_x^2 \psi^{n,i} A_t^\gamma \bar{\psi}^{n,i} &= \sum_n \frac{1}{\Delta x^2} \left( A_t^\gamma \psi^{n+1,i} A_t^\gamma \bar{\psi}^{n,i} - 2|A_t^\gamma \psi^{n,i}|^2 + A_t^\gamma \psi^{n-1,i} A_t^\gamma \bar{\psi}^{n,i} \right) \\ &= \sum_n \frac{1}{\Delta x^2} \left( A_t^\gamma \psi^{n+1,i} A_t^\gamma \bar{\psi}^{n,i} - 2|A_t^\gamma \psi^{n,i}|^2 + A_t^\gamma \psi^{n,i} A_t^\gamma \bar{\psi}^{n+1,i} \right) \\ &= \sum_n \frac{1}{\Delta x^2} \left( 2\Re(A_t^\gamma \psi^{n+1,i} A_t^\gamma \bar{\psi}^{n,i}) - 2|A_t^\gamma \psi^{n,i}|^2 \right) \end{aligned} \quad (5.23)$$

and the third term gives

$$\sum_n V'(|A_t^\gamma \psi^{n,i}|^2) A_t^\gamma \psi^{n,i} A_t^\gamma \bar{\psi}^{n,i} = \sum_n V'(|A_t^\gamma \psi^{n,i}|^2) |A_t^\gamma \psi^{n,i}|^2. \quad (5.24)$$

Substituting eqs. (5.22) to (5.24) in eq. (5.21) we get

$$\begin{aligned} &\sum_n \frac{1}{2\Delta t} \left( \mathbf{i} (e^{2\gamma\Delta t} |\psi^{n,i+1}|^2 - e^{-2\gamma\Delta t} |\psi^{n,i}|^2) - 2\Im(\psi^{n,i+1} \bar{\psi}^{n,i}) \right) \\ &+ \sum_n \frac{1}{\Delta x^2} \left( 2\Re(A_t^\gamma \psi^{n+1,i} A_t^\gamma \bar{\psi}^{n,i}) - 2|A_t^\gamma \psi^{n,i}|^2 \right) + \sum_n V'(|A_t^\gamma \psi^{n,i}|^2) |A_t^\gamma \psi^{n,i}|^2 = 0. \end{aligned}$$

Now taking the imaginary part of this equation we get

$$\sum_n (e^{2\gamma\Delta t} |\psi^{n,i+1}|^2 - e^{-2\gamma\Delta t} |\psi^{n,i}|^2) = 0$$

or equivalently

$$\sum_n |\psi^{n,i+1}|^2 = e^{-4\gamma\Delta t} \sum_n |\psi^{n,i}|^2$$

i.e. the method preserves  $e^{4\gamma t} I_2$ .

(iii) For linear Schrödinger eq. (5.19), the method preserves  $e^{4\gamma t} I_k$  for all  $k$ . For preservation of  $e^{4\gamma t} I_1$ , using Lemma 1.8 we obtain

$$\begin{aligned} \mathbf{i} D_t^{2\gamma} \sum_n \bar{\psi} \delta_x \psi &= \mathbf{i} \sum_n D_t^\gamma \bar{\psi} A_t^\gamma \delta_x \psi + \mathbf{i} \sum_n A_t^\gamma \bar{\psi} D_t \delta_x \psi \\ &= \mathbf{i} \sum_n D_t^\gamma \bar{\psi} A_t^\gamma \delta_x \psi - \mathbf{i} \sum_n D_t \psi A_t^\gamma \delta_x \bar{\psi}. \end{aligned}$$

Using eq. (4.15), with  $V(\eta) = \lambda\eta$ , in this equation we get

$$\mathbf{i} D_t^{2\gamma} \sum_n \bar{\psi} \delta_x \psi = \sum_n (A_t^\gamma \delta_x^2 \bar{\psi} A_t^\gamma \delta_x \psi + A_t^\gamma \delta_x^2 \psi A_t^\gamma \delta_x \bar{\psi}) + \lambda \sum_n (A_t^\gamma \bar{\psi} A_t^\gamma \delta_x \psi + A_t^\gamma \psi A_t^\gamma \delta_x \bar{\psi}).$$

Using Lemma 1.8 again, we get

$$\mathbf{i} D_t^{2\gamma} \sum_n \bar{\psi} \delta_x \psi = \sum_n (A_t^\gamma \delta_x^2 \bar{\psi} A_t^\gamma \delta_x \psi - A_t^\gamma \delta_x \psi A_t^\gamma \delta_x^2 \bar{\psi}) + \lambda \sum_n (A_t^\gamma \bar{\psi} A_t^\gamma \delta_x \psi - A_t^\gamma \delta_x \psi A_t^\gamma \bar{\psi}) = 0.$$

Therefore,

$$\sum_n \overline{\psi^{n,i+1}} \delta_x \psi^{n,i+1} = e^{-4\gamma\Delta t} \sum_n \overline{\psi^{n,i}} \delta_x \psi^{n,i}.$$

This shows that the method preserves  $e^{4\gamma t} I_1$ .

Preservation of  $e^{4\gamma t} I_2$  by the method can be shown in a manner similar to preservation of the

same for damped NLS in the previous item. For preservation of  $e^{4\gamma t} I_3$ , notice that

$$\begin{aligned} -D_t^{2\gamma} \sum_n |D_x \psi|^2 &= D_t^{2\gamma} \sum_n \psi \delta_x^2 \bar{\psi} \\ &= \sum_n (D_t^\gamma \psi A_t^\gamma \delta_x^2 \bar{\psi} + A_t^\gamma \psi D_t^\gamma \delta_x^2 \bar{\psi}) \\ &= \sum_n (D_t^\gamma \psi A_t^\gamma \delta_x^2 \bar{\psi} + A_t^\gamma \delta_x^2 \psi D_t^\gamma \bar{\psi}), \end{aligned}$$

where we have used the summation by parts formula, periodic boundary conditions, and Lemma 1.8. Now using the method of eq. (4.15), with  $V(\eta) = \lambda\eta$ , we obtain

$$-D_t^{2\gamma} \sum_n |D_x \psi|^2 = -\lambda \sum_n (D_t^\gamma \psi A_t^\gamma \bar{\psi} + D_t^\gamma \bar{\psi} A_t^\gamma \psi). \quad (5.25)$$

Using this in the following

$$\begin{aligned} D_t^{2\gamma} \sum_n (-|D_x \psi|^2 + \lambda|\psi|^2) &= -\lambda \sum_n (D_t^\gamma \psi A_t^\gamma \bar{\psi} + D_t^\gamma \bar{\psi} A_t^\gamma \psi) + \lambda \sum_n (D_t^\gamma \psi A_t^\gamma \bar{\psi} + D_t^\gamma \bar{\psi} A_t^\gamma \psi) \\ &= 0. \end{aligned}$$

Therefore

$$\sum_n (-|D_x \psi^{n,i+1}|^2 + \lambda|\psi^{n,i+1}|^2) = e^{-4\gamma\Delta t} \sum_n (-|D_x \psi^{n,i}|^2 + \lambda|\psi^{n,i}|^2),$$

i.e. the method preserves  $e^{4\gamma t} I_3$  for the linear Schrödinger equation.

### 5.3.1.2 Exponential time differencing method

Now, discretizing eq. (4.13) in time with exponential time differencing method (2.13), we obtain

$$\frac{\gamma\Delta t}{\sinh(\gamma\Delta t)} \mathbf{K} D_t^\gamma z + \mathbf{L}_+ D_x A_t^\gamma z + \mathbf{L}_- D_x T_x A_t^\gamma z = \nabla S(A_t^\gamma z) + \mathbf{F}(A_t t) \quad (5.26)$$

In terms of the original variable, the method becomes

$$i \frac{\gamma \Delta t}{\sinh(\gamma \Delta t)} D_t^\gamma \psi + A_t^\gamma \delta_x^2 \psi + V'(|A_t^\gamma \psi|^2) A_t^\gamma \psi = F(A_t t). \quad (5.27)$$

This method will be referred to as ETD in the numerical plots and it should be compared with the continuous equation (4.7) and method (4.14). We summarize structure-preservation properties of method (5.27) below. These properties can be derived by informally replacing the time derivative by  $\frac{\gamma \Delta t}{\sinh(\gamma \Delta t)} D_t^\gamma$  in the derivation of corresponding properties of the method (4.14).

- (i) Method (5.26) is multi-conformal-symplectic because it satisfies following discrete multi-conformal-symplectic conservation law

$$\frac{\gamma \Delta t}{\sinh(\gamma \Delta t)} D_t^{2\gamma} \left( \frac{1}{2} (dz \wedge \mathbf{K} dz) \right) + D_x (A_t^\gamma T_x dz \wedge \mathbf{L}_+ A_t^\gamma dz) = 0.$$

- (ii) The method preserves  $e^{4\gamma t} I_2$  for damped NLS i.e.

$$\sum_n |\psi^{n,i+1}|^2 = e^{-4\gamma \Delta t} \sum_n |\psi^{n,i}|^2.$$

- (iii) The method preserves  $e^{4\gamma t} I_k$ ,  $k = 1, 2, 3$ , for linear Schrödinger eq. (5.19):

- $\sum_n \overline{\psi^{n,i+1}} \delta_x \psi^{n,i+1} = e^{-4\gamma \Delta t} \sum_n \overline{\psi^{n,i}} \delta_x \psi^{n,i}.$
- $\sum_n |\psi^{n,i+1}|^2 = e^{-4\gamma \Delta t} \sum_n |\psi^{n,i}|^2.$
- $\sum_n -|D_x \psi^{n,i+1}|^2 + \lambda |\psi^{n,i+1}|^2 = e^{-4\gamma \Delta t} \sum_n -|D_x \psi^{n,i}|^2 + \lambda |\psi^{n,i}|^2.$

### 5.3.1.3 Implicit midpoint method

For the purpose of comparison with a closely related method which is not structure-preserving, let us discretize time in the semi-discretized system (4.13) with the (symplectic) implicit midpoint



rule:

$$\mathbf{K}D_t z + \mathbf{L}_+ D_x A_t z + \mathbf{L}_- D_x T_x A_t z = \nabla S(A_t z) - 2\gamma \mathbf{K} A_t z + \mathbf{F}(A_t z).$$

Rewriting this system in terms of the original variable, one gets

$$iD_t \psi + A_t \delta_x^2 \psi + V'(|A_t \psi|^2) A_t \psi + 2i\gamma A_t \psi = F(A_t z) \quad (5.28)$$

which should be compared with its continuous counterpart eq. (4.7) and methods given by eqs. (4.15) and (5.27). Method (5.28) will be referred to as IMP in the numerical plots. It can be easily shown that the method is neither multi-conformal-symplectic nor does it preserve dissipative properties of the linear Schrödinger equation and damped NLS.

### 5.3.2 Numerical results

We now turn to numerical implementation of the methods introduced in the last section. We refer to methods of eqs. (4.15), (5.27) and (5.28) as IF, ETD, and IMP, respectively. We start by demonstrating structure preservation for linear Schrödinger eq. (5.19) and damped NLS and then show that the methods successfully capture global attractors for damped driven NLS eq. (4.7). We compute errors in preserving the invariants of eqs. (5.17) and (5.20) by computing the residuals

$$R_k^i = \log \left( \frac{I_k^{i+1}}{I_k^i} \right) + 4\gamma \Delta t, \quad (5.29)$$

where  $I_k^i$  is the numerical approximation of  $I_k(t^i)$  for  $k = 1, 2, 3$ :

$$I_1^i = \Im \sum_n \overline{\psi^{n,i}} \mathcal{D}_x \psi^{n,i}, \quad I_2^i = \sum_n |\psi^{n,i}|^2, \quad I_3^i = \sum_n -|\mathcal{D}_x \psi^{n,i}|^2 + V(|\psi^{n,i}|^2).$$

The spectral differentiation matrix operator  $\mathcal{D}_x$  is implemented using MATLAB's FFT routine. We choose  $\mathcal{D}_x$ , instead of finite difference operators, to reduce error in evaluating  $I_k^i$ 's along numerical solutions. Let us denote the vector  $\{R_k^i\}_i$  by  $\mathbf{R}_k$  and  $\{I_k^i\}_i$  by  $\mathbf{I}_k$  for all  $k$ . We shall assume

$V(\eta) = \lambda\eta$  and  $V(\eta) = \frac{1}{2}\lambda\eta^2$ , with  $\lambda = 2$ , for linear and nonlinear NLS, respectively, in all the experiments that follow.

### 5.3.2.1 Linear Schrödinger equation

IF and ETD preserve dissipation in momentum, norm, and energy of linear Schrödinger eq. (5.19).

A plane wave solution of the linear Schrödinger equation is given by

$$\psi(x, t) = Ae^{i\lambda t - 2\gamma t},$$

where  $A$  is amplitude of the solution. Initializing the three methods with this plane wave we obtain Figure 5.6. The figure verifies numerical preservation of dissipation in the properties of the equation by IF and ETD methods.

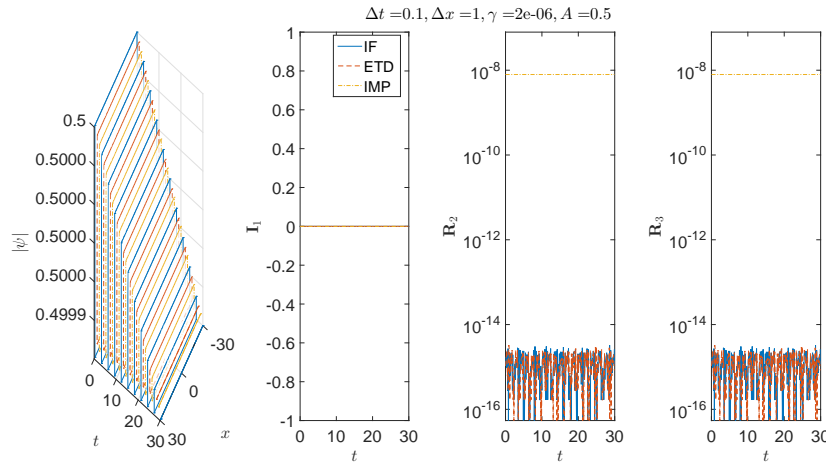


Figure 5.6: Plane wave solution, momentum, and norm and energy residuals. The second column gives  $I_1$  because  $R_1$  is undefined when the  $x$ -derivative of the solution is zero.

### 5.3.2.2 Damped NLS

For this part of the experiments, we set  $\mathbf{F} = 0$  in eq. (4.7), so that there is no driving force. Plane wave solutions of eq. (4.7), with  $\mathbf{F} = 0$ , are given by

$$\psi(x, t) = Ae^{-2\gamma t} \exp\left(\mathrm{i}\lambda|A|^2 \left(\frac{1 - e^{-4\gamma t}}{4\gamma}\right)\right),$$

where  $A$  is the amplitude of the plane wave. We plot  $L^\infty$  errors in Figure 5.7 using the exact plane wave solution and numerical solutions due to the methods. The figure verifies theoretical spatial order of the methods. In Figure 5.8, we plot numerical solutions, along with the residuals defined in eq. (5.29), initialized with the plane wave solution. The figure demonstrates preservation of the invariant  $e^{4\gamma t}I_2$  by IF and ETD methods. For the plane wave solution, which has a spatially flat profile and hence the  $x$ -derivative is zero, all the methods also seem to preserve  $e^{4\gamma t}I_1$ .

It is well known that the integrable NLS with cubic nonlinearity ( $V'(|\psi|^2) = \lambda|\psi|^2$ ) has soliton solutions. Soliton solutions travel and pass through each other maintaining their original shapes. In our next experiment, we demonstrate collision of two waves for the damped NLS equation in Figure 5.9. These waves propagate towards each other, collide, and emerge out of the collision with their original shapes and smaller amplitudes. The initial profile is

$$\psi(x, 0) = e^{5\mathrm{i}x} \operatorname{sech}(x + 1) + 1.5e^{-5\mathrm{i}x} \operatorname{sech}(1.5(x - 5)).$$

The figure also demonstrates preservation of dissipation in the norm,  $e^{4\gamma t}I_2$ , by IF and ETD methods. For IF and ETD methods, residual  $\mathbf{R}_1$  becomes large near the time of the collision of the two peaks but both methods recover after the collision.

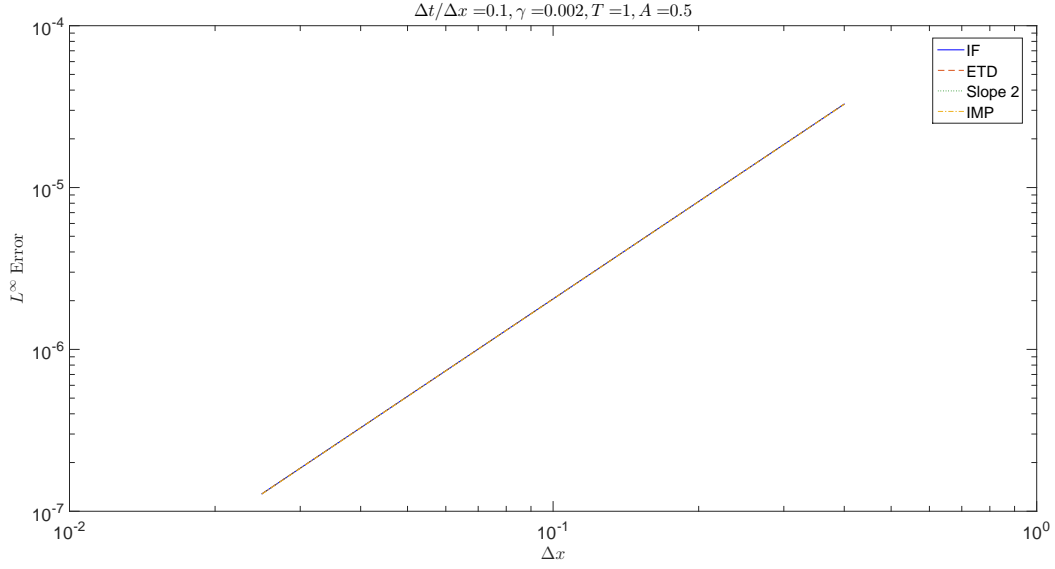


Figure 5.7:  $L^\infty$  error due to the methods of eqs. (4.15) and (5.27).

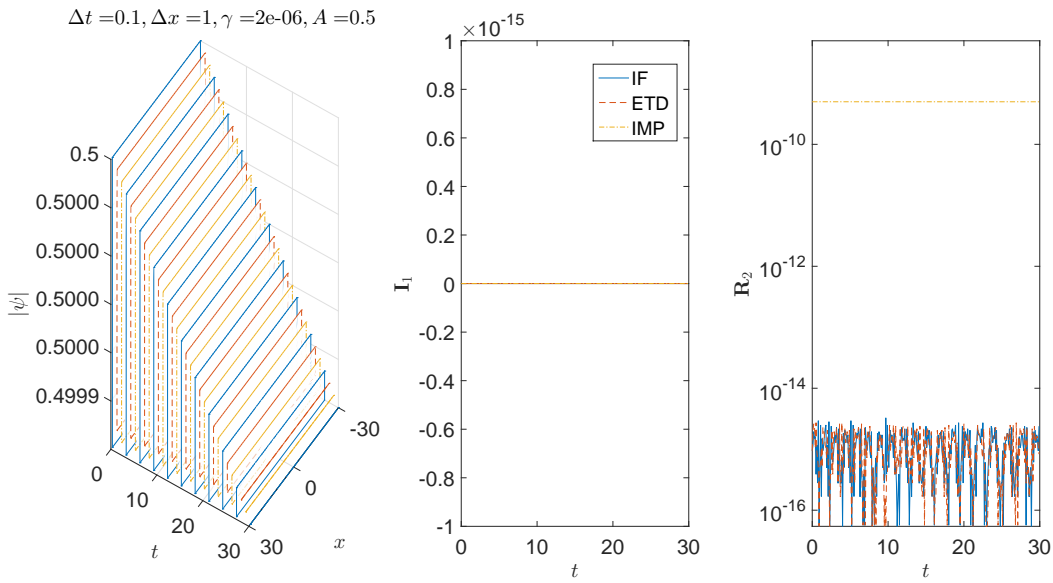


Figure 5.8: Plane wave solution, momentum, and invariant residual.

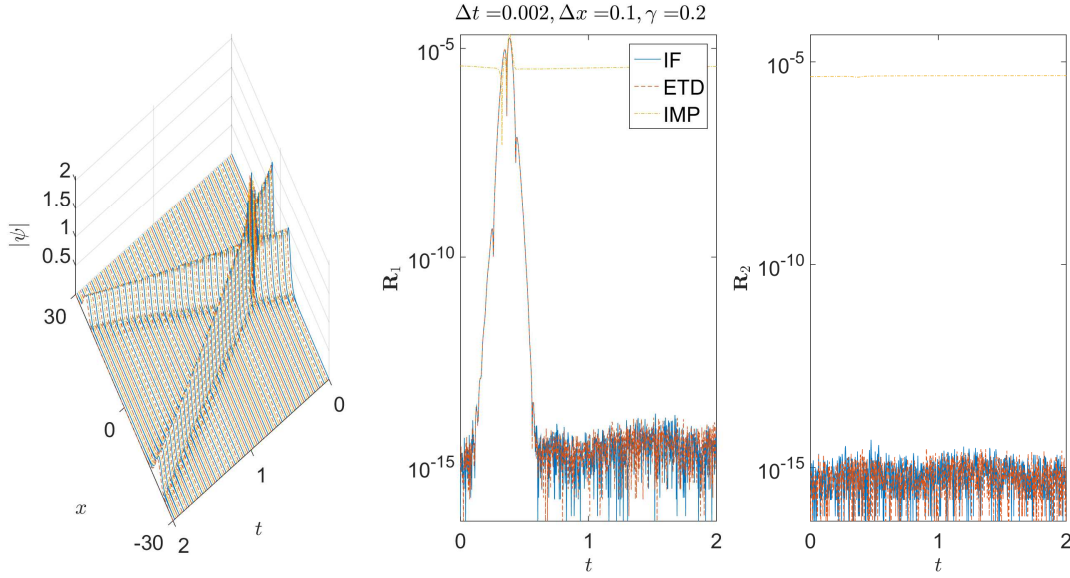


Figure 5.9: Soliton collision and invariant residuals. For IF and ETD methods, residual  $R_1$  is close to machine precision except near the time of collision when solution profile is steep at the spatial location of the collision.

### 5.3.2.3 Damped driven NLS

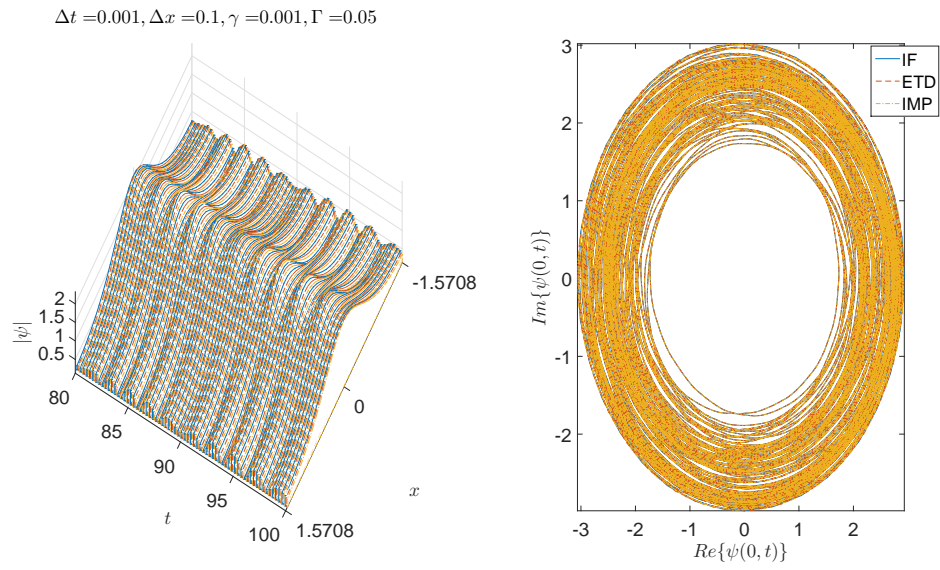
When a damped nonlinear pendulum is driven with external force, it shows chaotic behavior in certain parametric regimes. Similarly, theory predicts chaotic solutions when external driving force  $F(t)$  is included in a damped NLS. For these experiments we assume

$$F(t) = \Gamma e^{i(\omega_0 t + \alpha)}.$$

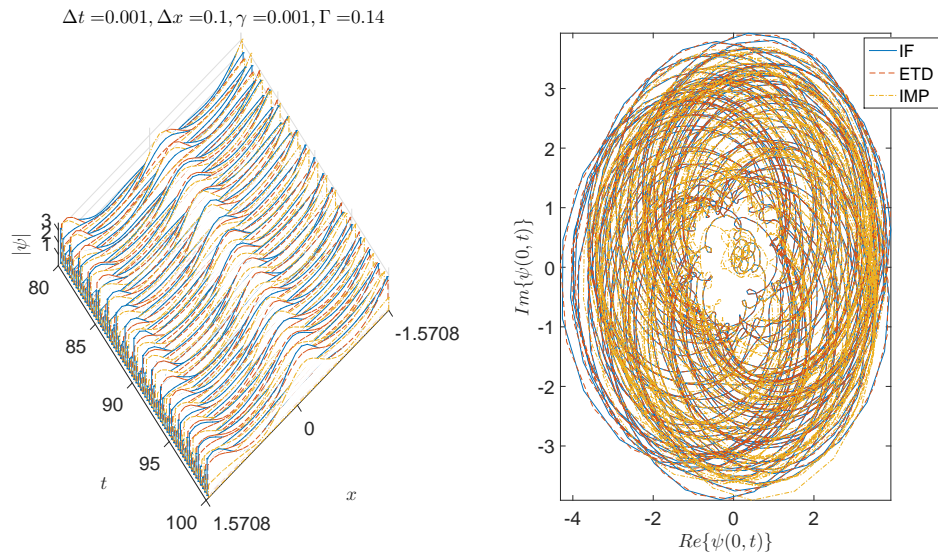
Where  $\Gamma$ ,  $\omega_0$  and  $\alpha$  are amplitude, frequency, and phase, respectively, of the driving force. We set  $\omega_0 = 1$ ,  $\alpha = 0$ , and vary  $\Gamma$ . As  $\Gamma$  varies, we observe periodic and chaotic attractors. The initial condition used is a hyperbolic secant profile

$$\psi(x, 0) = 3 \operatorname{sech}(3x).$$

Figure 5.10 shows the amplitude of the numerical solutions and their imaginary versus real parts at  $x = 0$  for all times. Figure 5.10a shows a periodic attractor and Figure 5.10b shows a temporally chaotic state. For a temporally chaotic state, the peak of the solutions jumps back and forth between two different spatial locations. Numerical solutions with a small amplitude at  $x = 0$  correspond to points near the origin and those with a large amplitude at  $x = 0$  correspond to points far from the origin in imaginary versus real parts subplots of the figure. For the parameter values chosen, all the numerical methods are in agreement as they settle down to the same attractor after the transient phase is over.



(a) Periodic attractor



(b) Chaotic attractor

Figure 5.10: Periodic and chaotic attractors of damped driven NLS along with imaginary versus real parts of numerical solution at all times and  $x = 0$ .

## CHAPTER 6: CONCLUSIONS AND FUTURE WORK

This thesis presents and develops a class of exponential Runge-Kutta and partitioned exponential Runge-Kutta methods. The methods are useful for differential equations with solutions that satisfy properties of the form  $\mathcal{I}(t) = e^{-x_0(t)}\mathcal{I}(0)$ . In many cases of interest  $\mathcal{I}$  is linear, quadratic, or a symplectic two-form. Because the methods produce solutions that satisfy  $\mathcal{I}(t_{n+1}) = e^{-x_0(h)}\mathcal{I}(t_n)$  (under certain restrictions on the coefficient functions), they preserve the properties in a way that is stronger than other methods that simply guarantee  $\mathcal{I}(t_{n+1}) < \mathcal{I}(t_n)$  when  $x_0(h) > 0$ . Our focus is on integrating factor methods and exponential time differencing methods, but the theorems on structure-preservation may also apply to other types of exponential integrators. The strengths of the methods are illustrated for various integrators applied to several model problems through numerical experiments.

We have also developed structure-preserving integrators that preserve conservation laws of the form

$$\partial_t P + \partial_x Q = -aP$$

of a PDE. These methods were applied to PDEs and they were shown to satisfy additional structure in some special cases. When these methods were compared against other non-structure-preserving methods, the strengths and advantages of structure-preservation were demonstrated. In summary, our research on structure-preserving numerical methods extends the existing body of knowledge and provides improvement and deeper understanding of geometric integrators for linearly damped DEs.

In keeping this thesis taut and focused, we had to put off several interesting perspectives and new questions have also emerged out of the study. The methods developed here are interesting for conservative systems that are perturbed with linear, possibly time-dependent, non-conservative terms,



and many aspects of the methods are well-understood, thanks to a wealth of prior research on exponential integrators. Nevertheless, there are several avenues for future research on this topic, including order conditions for PERK methods, backward error analysis, extension and application to partial differential equations, and development of methods that preserve other important properties of mechanical systems that are perturbed by non-conservative terms.

Often times, a DE has several qualitative properties and a geometric integrator is considered better than other geometric integrators if the former preserves more qualitative properties of the DE than the latter. Damped PDEs often have conformal invariants such as momentum, mass, and energy. It is a natural extension of this thesis to develop such integrators for damped PDEs that have not been considered here.

Order of accuracy of Runge-Kutta (RK) methods can be obtained by examining whether their coefficients satisfy certain conditions, referred to as the order conditions. To the contrary, exponential Runge-Kutta methods lack such order conditions except in some specific cases. One way to get around this is to obtain ERK methods in such a way that their order is obvious by design. Generating functions are used to design symplectic RK methods of a specified order. Generating functions approach may also reveal important insights about developing structure-preserving ERK methods of a specified order.

Geometric integrators for some one-dimensional PDEs can be easily generalized to their higher dimensional versions whereas others require fundamentally different approach. Preservation of properties such as volume and measure, of Hamiltonian systems, by numerical methods have been shown to be advantageous in the literature. However, more research needs to be done to develop methods which preserve these properties or their dissipation for non-Hamiltonian systems.

## **APPENDIX A: DIFFERENTIAL FORMS AND THE WEDGE PRODUCT**

Damped DEs that we consider in this thesis can be cast as a conformal Hamiltonian ODE, a multi-conformal-symplectic PDE, or a perturbation of these two DEs. Specific forms of these DEs guarantee differential form(s) that satisfy a damped linear ordinary or partial differential equation. With this in mind, we introduce the concept of differential forms here. A differential or differential 1-form  $df$  of a function  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  is defined to be

$$df = f_{z_1} dz_1 + f_{z_2} dz_2 + \dots + f_{z_m} dz_m$$

where  $z \in \mathbb{R}^m$  and subscripts denote partial derivatives. This is written in a more compact form as

$$df = f_z \cdot dz,$$

where  $dz = [dz_1 \ dz_2 \ \dots \ dz_m]^T$  is a vector of  $m$  differentials and  $f_z$  is a column vector of partial derivatives of  $f$ . A differential operates on a vector  $\zeta \in \mathbb{R}^m$  in the following manner

$$df(\zeta) = f_{z_1} dz_1(\zeta) + f_{z_2} dz_2(\zeta) + \dots + f_{z_m} dz_m(\zeta).$$

Notice that this is simply the directional derivative of  $f$  in the direction of  $\zeta$ . The *wedge product*  $df \wedge dg : \mathbb{R}^{2m} \rightarrow \mathbb{R}$  of two differential forms  $df$  and  $dg$  is defined to be

$$df \wedge dg(\zeta, \eta) = dg(\zeta)df(\eta) - df(\zeta)dg(\eta).$$

The left hand side of the above equation is referred to as a differential 2-form. The wedge product converts a differential  $k$ -form into a  $k + 1$ -form in general. Differential forms are often denoted by Greek letters  $\omega, \kappa, \tau$ , etc.

The wedge product of two vector functions is defined in a similar manner. Let

$$da = [da_1 \ da_2 \ \dots \ da_m]^T \text{ and } db = [db_1 \ db_2 \ \dots \ db_m]^T$$

be two differential 1-forms. Then their wedge product gives a differential 2-form and is defined to

be

$$da \wedge db = \sum_{i=1}^m da_i \wedge db_i.$$

Given three differential 1-forms  $\omega$ ,  $\kappa$ , and  $\tau$ , which are  $m$ -vectors in  $\mathbb{R}^m$ , one can check following properties of the wedge product [28]

- $(\omega + \kappa) \wedge \tau = \omega \wedge \tau + \kappa \wedge \tau$
- $(\omega \wedge \kappa) \wedge \tau = \omega \wedge (\kappa \wedge \tau)$
- $\omega \wedge \kappa = -\kappa \wedge \omega$
- $\omega \wedge (\mathbf{A}\kappa) = (\mathbf{A}^T\omega) \wedge \kappa$

for any  $m \times m$  matrix  $\mathbf{A}$ . It follows from these properties that for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{m \times m}$  and a differential 1-form  $dz \in \mathbb{R}^m$ , we have

$$\mathbf{A}dz \wedge dz = 0.$$

Moreover, the converse is also true if  $\mathbf{A}$  is a constant matrix:

**Lemma A.1.** *If a vector  $z \in \mathbb{R}^m$  and a real matrix  $\mathbf{A} \in \mathbb{R}^{m \times m}$  satisfy*

$$\mathbf{A}dz \wedge dz = 0,$$

*then  $\mathbf{A}$  is symmetric.*

*Proof.* Equation  $\mathbf{A}dz \wedge dz = 0$  implies that

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1m} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2m} \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3m} \\ \vdots & & & & \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mm} \end{bmatrix} \begin{bmatrix} dz_1 \\ dz_2 \\ dz_3 \\ \vdots \\ dz_m \end{bmatrix} \wedge \begin{bmatrix} dz_1 \\ dz_2 \\ dz_3 \\ \vdots \\ dz_m \end{bmatrix} = 0,$$

$$\begin{aligned} & a_{12}dz_2 \wedge dz_1 + a_{13}dz_3 \wedge dz_1 + \cdots + a_{1m}dz_m \wedge dz_1 \\ & + a_{21}dz_1 \wedge dz_2 + a_{23}dz_3 \wedge dz_2 + \cdots + a_{2m}dz_m \wedge dz_2 \\ & + a_{31}dz_1 \wedge dz_3 + a_{32}dz_2 \wedge dz_3 + \cdots + a_{3m}dz_m \wedge dz_3 \\ & \cdots + a_{d1}dz_1 \wedge dz_d + a_{d2}dz_2 \wedge dz_d + \cdots + a_{mm-1}dz_{m-1} \wedge dz_m = 0. \end{aligned}$$

This implies that

$$\sum_{i=2}^m (a_{i1} - a_{1i})(dz_1 \wedge dz_i) + \sum_{i=3}^m (a_{i2} - a_{2i})(dz_2 \wedge dz_i) + \cdots + (a_{mm-1} - a_{m-1m})(dz_{m-1} \wedge dz_m) = 0.$$

But  $\{dz_i \wedge dz_j\}_{1 \leq i < j \leq m}$  forms a basis of vector space  $\wedge^2(\mathbb{R}^m)$  of all differential 2-forms. Therefore

$$a_{ij} = a_{ji} \text{ for all } i, j \quad \text{i.e.} \quad \mathbf{A} = \mathbf{A}^T.$$

□

Let us recall that a numerical method with flow map  $\psi_h$  is conformal symplectic if

$$(\psi'_t(z))^T \mathbf{J}^{-1} \psi'_t(z) = e^{-2 \int_0^t \gamma(s) ds} \mathbf{J}^{-1}$$

by eq. (2.41). Therefore, a numerical method with flow map  $\psi_t$  is symplectic if

$$(\psi'_t(z))^T \mathbf{J}^{-1} (\psi'_t(z)) = \mathbf{J}^{-1}.$$

Taking the determinant of both sides we see that

$$\begin{aligned} \det((\psi'_t(z))^T \mathbf{J}^{-1}(\psi'_t(z))) &= \det(\mathbf{J}^{-1}), \\ \det(\psi'_t(z)) &= 1, \end{aligned}$$

because  $\det(AB) = \det(A)\det(B)$  for any two matrices  $A$  and  $B$ .

## LIST OF REFERENCES

- [1] U. M. Ascher and R. I. McLachlan. Multisymplectic box schemes and the Korteweg–de Vries equation. *Appl. Numer. Math.*, 48(3):255–269, 2004.
- [2] H. Berland, B. Owren, and B. Skaflestad. B-series and order conditions for exponential integrators. *SIAM J. Numer. Anal.*, 43(4):1715–1727, 2005.
- [3] H. Berland, B. Owren, and B. Skaflestad. Solving the nonlinear Schrödinger equation using exponential integrators. *Model. Ident. Control*, 27(4):201–217, 2006.
- [4] A. Bhatt. Multi conformal symplectic integration of damped driven nonlinear Schrödinger equation. *preprint submitted*, 2016.
- [5] A. Bhatt, D. Floyd, and B. E. Moore. Second order conformal symplectic schemes for damped Hamiltonian systems. *J. Sci. Comput.*, 66(3):1234–1259, 2015.
- [6] A. Bhatt and B. E. Moore. Structure preserving exponential Runge-Kutta methods. *preprint submitted*, 2016.
- [7] P. B. Bochev and C. Scovel. On quadratic invariants and symplectic structure. *BIT Numer. Math.*, 34(3):337–345, 1994.
- [8] T. J. Bridges. A geometric formulation of the conservation of wave action and its implications for signature and the classification of instabilities. In *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, volume 453, pages 1365–1395, 1997.
- [9] T. J. Bridges. Multi-symplectic structures and wave propagation. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 121, pages 147–190. Cambridge Univ Press, 1997.

- [10] J. Butcher. Trees, b-series and exponential integrators. *IMA J. Numer. Anal.*, 30(1):131–140, 2010.
- [11] D. Cai and D. W. McLaughlin. Chaotic and turbulent behavior of unstable one-dimensional nonlinear dispersive waves. *J. Math. Phys.*, 41(6):4125–4153, 2000.
- [12] E. Celledoni, D. Cohen, and B. Owren. Symmetric exponential integrators with an application to the cubic Schrödinger equation. *Found. Comput. Math.*, 8(3):303 – 317, 2008.
- [13] E. Celledoni, V. Grimm, R. I. McLachlan, D. McLaren, D. O’Neale, B. Owren, and G. Quispel. Preserving energy resp. dissipation in numerical PDEs using the average vector field method. *J. Comput. Phys.*, 231(20):6770–6789, 2012.
- [14] E. Celledoni, R. I. McLachlan, D. I. McLaren, B. Owren, G. R. W. Quispel, and W. M. Wright. Energy-preserving Runge-Kutta methods. *ESAIM Math. Model. Numer. Anal.*, 43(04):645–649, 2009.
- [15] G. Cooper. Stability of Runge-Kutta methods for trajectory problems. *IMA J. Numer. Anal.*, 7(1):1–13, 1987.
- [16] A. Eden, A. Milani, and B. Nicolaenko. Finite dimensional exponential attractors for semi-linear wave equations with damping. *J. Math. Anal. Appl.*, 169(2):408–419, 1992.
- [17] S. Geng. Symplectic partitioned Runge-Kutta methods. *J. Comput. Math.*, pages 365–372, 1993.
- [18] J.-M. Ghidaglia. *Finite dimensional behavior for weakly damped driven Schrödinger equations*, volume 5. Annales de l’IHP Analyse non linéaire, 1988.
- [19] E. Hairer, C. Lubich, and G. Wanner. *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*. Springer Series in Computational Mathematics. Springer Berlin Heidelberg, 2006.



- [20] A. Harten, P. D. Lax, and B. Van Leer. On upstream differencing and Godunov-type schemes for hyperbolic conservation laws. In *Upwind and High-Resolution Schemes*, pages 53–79. Springer, 1997.
- [21] Hartono and A. van der Burgh. A linear differential equation with a time-periodic damping coefficient: stability diagram and an application. *J. Eng. Math.*, 49(2):99–112.
- [22] D. Hipp, M. Hochbruck, and A. Ostermann. An exponential integrator for non-autonomous parabolic problems. *Electron. T. Numer. Anal.*, 41:497–511, 2014.
- [23] M. Hochbruck and A. Ostermann. Explicit exponential Runge-Kutta methods for semilinear parabolic problems. *SIAM J. Numer. Anal.*, 43(3):1069–1090, 2005.
- [24] A. Islas, D. Karpeev, and C. Schober. Geometric integrators for the nonlinear Schrödinger equation. *J. Comput. Phys.*, 173(1):116–148, 2001.
- [25] X. Kong, H. Wu, and F. Mei. Structure-preserving algorithms for Birkhoffian systems. *J. Geom. Phys.*, 62(5):1157–1166, 2012.
- [26] F. Lasagni. Canonical Runge-Kutta methods. *Z. Angew. Math. Phys.*, 39(6):952–953, 1988.
- [27] J. D. Lawson. Generalized Runge-Kutta processes for stable systems with large Lipschitz constants. *SIAM J. Numer. Anal.*, 4(3):372–380, 1967.
- [28] B. Leimkuhler and S. Reich. *Simulating Hamiltonian Dynamics*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, 2004.
- [29] C.-C. Liao, J.-C. Cui, J.-Z. Liang, and X.-H. Ding. Multi-symplectic variational integrators for nonlinear Schrödinger equations with variable coefficients. *Chinese Phys. B*, 25(1):010205, 2015.

- [30] J. E. Marsden, G. W. Patrick, and S. Shkoller. Multisymplectic geometry, variational integrators, and nonlinear PDEs. *Commun. Math. Phys.*, 199(2):351–395, 1998.
- [31] T. Matsuo and D. Furihata. Dissipative or conservative finite-difference schemes for complex-valued nonlinear partial differential equations. *J. Comput. Phys.*, 171(2):425–447, 2001.
- [32] R. McLachlan and M. Perlmutter. Conformal Hamiltonian systems. *J. Geom. Phys.*, 39(4):276 – 300, 2001.
- [33] R. I. McLachlan and G. Quispel. What kinds of dynamics are there? Lie pseudogroups, dynamical systems and geometric integration. *Nonlinearity*, 14(6):1689, 2001.
- [34] R. I. McLachlan, G. Quispel, and N. Robidoux. Geometric integration using discrete gradients. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 357(1754):1021–1045, 1999.
- [35] R. E. Mickens. *Nonstandard finite difference models of differential equations*, volume 115. World Scientific, 1994.
- [36] R. E. Mickens. *Applications of Nonstandard Finite Difference Schemes*. World Scientific, Singapore, 2000.
- [37] K. Modin and G. Söderlind. Geometric integration of Hamiltonian systems perturbed by Rayleigh damping. *BIT Numer. Math.*, 51(4):977–1007, 2011.
- [38] B. E. Moore. Conformal multi-symplectic integration methods for forced-damped semi-linear wave equations. *Math. Comput. Simulat.*, 80(1):20 – 28, 2009.
- [39] B. E. Moore. Multi-conformal-symplectic PDEs and discretizations. *Preprint submitted*, 2015.

- [40] B. E. Moore, L. Noreña, and C. M. Schober. Conformal conservation laws and geometric integration for damped Hamiltonian PDEs. *J. Comput. Phys.*, 232(1):214 – 233, 2013.
- [41] M. Rogers. Expansion and divergence. Wolfram Demonstrations Project, 2010.
- [42] J. M. Sanz-Serna. Runge-Kutta schemes for Hamiltonian systems. *BIT Numer. Math.*, 28(4):877–883, 1988.
- [43] H. Su, M. Qin, Y. Wang, and R. Scherer. Multi-symplectic Birkhoffian structure for PDEs with dissipation terms. *Phys. Lett. A*, 374(24):2410–2416, 2010.
- [44] Y. Sun and Z. Shang. Structure-preserving algorithms for Birkhoffian systems. *Phys. Lett. A*, 336(4):358–369, 2005.
- [45] Y. B. Suris. The canonicity of mappings generated by Runge-Kutta type methods when integrating the systems  $\ddot{x} = -\partial U / \partial x$ . *USSR Comp. Math. Math+*, 29(1):138–144, 1989.