# Bespoke finite difference schemes that preserve multiple conservation laws

Timothy J. Grant

## Abstract

Conservation laws provide important constraints on the solutions of partial differential equations (PDEs), therefore it is important to preserve them when discretizing such equations. In this paper, a new systematic method for discretizing a PDE, so as to preserve the local form of multiple conservation laws, is presented. The technique, which uses symbolic computation, is applied to the Korteweg–de Vries (KdV) equation to find novel explicit and implicit schemes that have finite difference analogues of its first and second conservation laws and its first and third conservation laws. The resulting schemes are numerically compared with a multisymplectic scheme.

## 1. Introduction

In recent years there has been much interest in geometric numerical integration. The philosophy of geometric integration seeks numerical schemes that preserve geometric features of differential equations rather than focusing on the control of local errors of generic methods [8]. Geometric structures provide constraints for the behaviour of the system, hence it is desirable that numerical schemes possess analogues of the same constraints. Such schemes will then replicate the desired qualitative behaviour, and may have improved stability and accuracy compared with generic methods applied to the same differential equations.

The main focus for geometric integration has been on ordinary differential equations (ODEs), as has been documented in the monograph [16]; the lecture notes of McLachlan and Quispel [24] provide a good introduction. Numerical schemes have been developed to preserve various geometric structures. Perhaps the most celebrated example is the use of symplectic integrators for Hamiltonian ODEs (see, for example, [21]). Some well known methods have been shown to be symplectic integrators, for example, the implicit midpoint rule, the Stömer–Verlet method and the Gauss collocation methods. Alternatively, one might like to preserve a first integral of a problem such as the energy. This is done by rewriting the ODE as a skew gradient system which is discretized using a discrete gradient [24]. However, Ge and Marsden [31] proved, for non-integrable equations, that a numerical method with fixed time steps cannot be symplectic and exactly preserve energy at the same time.

The focus of this paper is the preservation of conservation laws (CLaws) of partial differential equations (PDEs). De Frutos and Sanz-Serna [12] demonstrated the benefits of preserving CLaws. They showed that, if a numerical method conserves the momentum for the Korteweg–de Vries (KdV) equation, then that numerical method can perform better than a non-conservative scheme with more accurate local truncation errors (LTEs). This is because, for the one-soliton solution, the numerical error for a conservative scheme occurs in the phase, rather than the amplitude, of the soliton. This highlights the desirability of creating finite difference methods that preserve CLaws.

The geometric integration of PDEs has been studied less than for ODEs; however, there exist various methods that relate to preserving CLaws. These methods generally require the PDE to have some additional structure other than the CLaw itself. The most famous of these methods are multisymplectic schemes for Hamiltonian PDEs [5]. In [4] Bridges and Reich prove

> the remarkable result that abstract linear Hamiltonian PDEs in multisymplectic form — discretized with the centered box scheme — conserve energy and momentum exactly; moreover, it is the local energy and momentum conservation that is preserved by the discretization.

However, they state [5] that, with a uniform discretization, it is not possible in general to preserve energy and momentum exactly along with the symplectic structure. Furihata [13] constructs finite difference schemes for equations of the form

$$\frac{\partial u}{\partial t} = \left( \frac{\partial}{\partial x} \right)^{2n+1} \frac{\delta G}{\delta u}, \quad n \in \mathbb{N},$$

that inherit the energy conservation property, that is,

$$\frac{d}{dt} \int G(u, u_x) \, dx = 0.$$

He does this by discretizing the energy function, $G$, and then applying a discrete variational derivative to construct the difference scheme. Using this method (the discrete variational derivative method, DVDM) he provides a scheme for the KdV equation that preserves the energy (1.5). In [19] Koide and Furihata generate schemes for the regularized long wave equation (Benjamin–Bona–Mahony (BBM) equation) that preserve the mass and momentum and the mass and energy using the DVDM. For their nonlinear momentum preserving scheme they show that, if the step sizes satisfy a certain condition, then solutions exist and are unique. The BBM equation, unlike the KdV equation which has an infinite number of CLaws, has only three, physically relevant, CLaws [11, 26]; therefore the DVDM is applicable to equations that are not integrable. In addition to the above methods, which preserve the divergence expression, McLachlan [23] constructs spatial discretizations of PDEs so that the resulting ODE systems have, as first integrals, the conserved quantities of the PDEs. The ODE system can then be integrated using a discrete gradient method to preserve the conserved quantities.

The goal of this paper is to introduce a method for discretizing a PDE so that the resulting discretization has difference analogues of multiple local conservation laws, without referring to any special structures the PDE may possess apart from the conservation laws themselves. In so doing, it is hoped that the method will be widely applicable. The only restriction on the applicability of the method is that the PDE should only have polynomial nonlinearities. For clarity, only scalar PDEs with two independent variables $x$ and $t$ are considered in this paper, though there is no reason why the methodology cannot be applied to systems of equations with more independent variables.

Finally, only discretizations of the KdV equation are studied in this paper. The technique finds new and known schemes, which are discussed in §4. The advantage of using the KdV equation is that the exact solution of the initial value problem is known for soliton solutions, and there are lots of discretizations to compare the resulting schemes with. The most famous finite difference scheme for solving the KdV equation is Zabusky and Kruskal's scheme (Z–K), from their paper [30] in which they coined the term 'soliton'. Their scheme is a two-step explicit method that has finite difference analogues of the mass and momentum CLaws. Sanz-Serna [28] showed that the Z–K scheme is subject to nonlinear instability; he then provided a scheme with an adaptive time step that preserves the mass and momentum exactly with periodic boundary conditions; however, the divergence form of the momentum CLaw is lost

(see [**14**]). In [**2**, **3**], Ascher and McLachlan investigate multisymplectic schemes for the KdV equation. They seek to understand the smooth behaviour of the multisymplectic schemes by studying the numerical dispersion of the linearized equations. From this, they suggest that box schemes are better than schemes that have a non-compact spatial discretization because the latter my introduce artificial wiggles into the solution. In §5 the new (conservation law preserving) schemes are compared with the Z–K scheme and the compact schemes used by Ascher and McLachlan.

### 1.1.   *Conservation laws of differential and difference equations*

A *conservation law* of a differential equation $\Delta = \mathbf{0}$ is a divergence expression that vanishes on solutions of the equation,

$$\text{Div}\,\mathbf{F} \equiv D_t(G) + D_x(F) = 0 \quad \text{when } \Delta = 0, \tag{1.1}$$

where

$$D_x \equiv \frac{\partial}{\partial x} + u_x \frac{\partial}{\partial u} + u_{xx} \frac{\partial}{\partial u_x} + u_{xt} \frac{\partial}{\partial u_t} + \dots$$

is the total $x$ derivative and $D_t$ is the total $t$ derivative. The terms $G$ and $F$ are referred to as the density and flux respectively and are functions of the independent variables, the dependent variable and its derivatives. For example, the KdV equation,

$$\Delta \equiv u_t + u u_x + u_{xxx} = 0, \tag{1.2}$$

has an infinite number of CLaws. In particular, it has the CLaws

$$0 = D_t(u) + D_x(\tfrac{1}{2}u^2 + u_{xx}) = \Delta, \tag{1.3}$$

$$0 = D_t(\tfrac{1}{2}u^2) + D_x(\tfrac{1}{3}u^3 + u u_{xx} - \tfrac{1}{2}u_x^2) = u\Delta, \tag{1.4}$$

$$0 = D_t(\tfrac{1}{3}u^3 - u_x^2) + D_x(\tfrac{1}{4}u^4 + u^2 u_{xx} - 2u_x u_{xxx} + u_{xx}^2 - 2u_x^2 u) = (u^2 - 2u_x D_x)\Delta. \tag{1.5}$$

Drazin and Johnson [**10**] state that, when applied to the water wave problem, (1.3) describes the conservation of mass, (1.4) the conservation of momentum and (1.5) the conservation of energy. A CLaw is *trivial of the first kind* if $\mathbf{F}$ vanishes on solutions of the PDE; it is *trivial of the second kind* if $\text{Div}\,\mathbf{F} \equiv 0$ (a null divergence; the divergence expression is zero without needing to be on solutions of the differential equation). A CLaw is *trivial* if it is a linear combination of the two kinds of trivial CLaws. Two CLaws are equivalent if they differ by a trivial CLaw. If the PDE is in Kovalevskaya form (see [**27**]; note that evolution equations are in Kovalevskaya form) then integrating the CLaw by parts (possibly repeatedly) yields an equivalent CLaw,

$$\text{Div}\,\tilde{\mathbf{F}} = Q \cdot \Delta, \tag{1.6}$$

which is said to be in characteristic form. The multiplier $Q$ is called a *characteristic* of the CLaw. For instance, the characteristic form of (1.5) is

$$\text{Div}\,\tilde{\mathbf{F}} = (u^2 + 2u_{xx})\Delta,$$

so $Q = u^2 + 2u_{xx}$ is the characteristic and (1.5) is equivalent to the CLaw

$$D_t(\tfrac{1}{3}u^3 - u_x^2) + D_x(\tfrac{1}{4}u^4 + u^2 u_{xx} + 2u_x u_t + u_{xx}^2) = 0. \tag{1.7}$$

A characteristic is trivial if it vanishes on solutions of the PDE. Two characteristics that differ by a trivial characteristic are said to be equivalent. Given a system of PDEs in Kovalevskaya

form, Alonso showed [**1**, **27**], there is a one-to-one correspondence between equivalence classes of characteristics and equivalence classes of conservation laws. Therefore characteristics can be used to identify when two seemingly different CLaws are equivalent. Characteristics have their most celebrated application in Noether's theorem [**25**, **27**] where they are used to construct CLaws from variational symmetries. However, the most important fact, for our purposes, is that total divergences form the kernel of the Euler operator [**27**],

$$\mathrm{E}(\mathrm{Div}\,\mathbf{F}) \equiv 0 \quad \text{where } \mathrm{E} \equiv \sum_{i,j}(-D_x)^i(-D_t)^j\frac{\partial}{\partial u_{x^i t^j}}.$$

Therefore, if $Q$ is a function such that

$$\mathrm{E}(Q \cdot \Delta) \equiv 0, \tag{1.8}$$

then $Q$ must be the characteristic of a CLaw.

CLaws are very important when modelling physical phenomena. The definition of a CLaw (1.1) as a divergence expression is a local property, and if integrated over the spatial domain (assuming vanishing boundary conditions) results in a quantity that is constant on solutions. Because (1.1) is a local constraint, preserving it, in a numerical method, provides a greater constraint on the behaviour than conserving the quantity that results from the spatial integration. A finite difference scheme preserves a given CLaw if it has a finite difference analogue of the CLaw of the differential equation.

Having restricted our attention to scalar PDEs with just two independent variables, let us now consider scalar difference equations with only two independent variables defined on a lattice $\mathbf{n} = (m, n)$. The dependent variable is evaluated at the grid points, so $u_{i,j} := u(m + i, n + j) \equiv u_{m+i,n+j}$ for $i, j \in \mathbb{Z}$, and $u$ take values in $\mathbb{R}$. The natural operators on the lattice are the shift operators, which are defined by

$$S_m : (m, n) \mapsto (m + 1, n), \quad S_n : (m, n) \mapsto (m, n + 1), \quad I : (m, n) \mapsto (m, n)$$
$$S_m : u_{i,j} \mapsto u_{i+1,j}, \quad S_n : u_{i,j} \mapsto u_{i,j+1} \quad \text{and} \quad I : u_{i,j} \mapsto u_{i,j}.$$

A difference equation is written as

$$\Delta(m, n, [u]) = 0,$$

where $[u]$ denotes a finite number of shifts of the dependent variables.

A *conservation law* of a partial difference equation (PΔE) is a divergence expression that vanishes on solutions of the system:

$$\mathrm{Div}\,\mathbf{F} := (S_m - I)F + (S_n - I)G = 0 \quad \text{when } [\Delta] = \mathbf{0}, \tag{1.9}$$

where $[\Delta]$ denotes any finite shifts of the difference equation. The functions $F$ and $G$ are known as the densities of the CLaw and may have the independent variables and shifts of the dependent variables as arguments. The key result exploited in this paper, due to Kuperschmidt [**18**, **20**], is that, just as for continuous equations (equation (1.8)), divergence expressions form the kernel of the discrete Euler operator,

$$\mathrm{E} \equiv \sum_{i,j} S_m^{-i} S_n^{-j}\frac{\partial}{\partial u_{i,j}}. \tag{1.10}$$

In the same way as for PDEs, a CLaw of a PΔE is trivial if and only if it is a linear combination of the following two kinds of trivial CLaws:

*First kind*  $\mathbf{F}|_{[\Delta]=\mathbf{0}} = \mathbf{0}$, that is, all of the densities vanish on solutions.

*Second kind*  $\operatorname{Div} \mathbf{F} \equiv 0$, without reference to the equation $[\Delta] = \mathbf{0}$ and its shifts. For instance, this occurs if $\mathbf{F}$ is the difference analogue of a total curl (see [**27**]).

For brevity, densities of a trivial CLaw are referred to as *trivial densities*. Just as for continuous equations, a CLaw is said to be in *characteristic form* if

$$\operatorname{Div} \mathbf{F} = Q(m, n, [u]) \cdot \Delta, \tag{1.11}$$

and the *characteristic* $Q$ is trivial if it vanishes when $[\Delta] = \mathbf{0}$. Analogous to PDEs in Kovalevskaya form, for explicit difference equations, there is a one-to-one correspondence between equivalence classes of characteristics and equivalence classes of characteristics [**15**].

 The discretization method, presented here, finds schemes with uniform steps that have finite difference analogues of the PDE's CLaws. To be explicit, the discretization has CLaws

$$\frac{(S_n - I)}{\nu} \widetilde{G_i} + \frac{(S_m - I)}{\mu} \widetilde{F_i} = 0 \quad \text{when } [\widetilde{\Delta}] = \mathbf{0}, \tag{1.12}$$

where $\mu$ is the spatial step and $\nu$ is the time step, and tildes represent discretizations of the corresponding continuous terms.

## 2.  *Highlights*

The discretization method proposed is very simple (and is discussed fully in §3). Form the most general discretizations of the PDE and the characteristic of the desired CLaw, so that there are undetermined coefficients in the discretizations. Apply the discrete Euler operator to the product of these discretizations. To have a discretization that preserves the CLaw, this expression needs to vanish. This results in a large system of polynomial equations in the undetermined coefficients. The solutions of this system will specify the coefficients required in the general discretization to preserve the given CLaw.

 In §4.3 the method is used to search for explicit two-step discretizations of the KdV equation that preserve the first and second CLaws together. The result is a three-parameter family of schemes, in which the famous Z–K scheme resides. Figures 1 and 2 show the results of using the Z–K scheme $(0, \frac{1}{2}, 0)$ and another scheme from the family $(0, \frac{1}{6}, 0)$ for solving the single soliton problem (with speed $c = 4$), and the two-soliton problem on a periodic domain (see §5 for a discussion of the numerics). In both figures, it is clear that the numerical solitons produced by the $(0, \frac{1}{6}, 0)$ scheme match the actual solution profile far better than the Z–K scheme. Moreover, in Tables 1 and 2, we see that the $(0, \frac{1}{6}, 0)$ scheme makes a slightly smaller error in preserving the second conserved quantity than the Z–K scheme, but it is significantly better at preserving the third conserved quantity. Thus it is clear that the discretization method can find well known schemes and new schemes that may be superior to existing methods.

TABLE 1. *Maximum absolute errors in preserving the conserved quantities, for the one-soliton problem, and the times at which they occurred. The problem was numerically solved for* $t \in [0, 2]$, $\mu = \frac{2}{15}$, *and* $\nu = \frac{1}{3}\mu^3$.

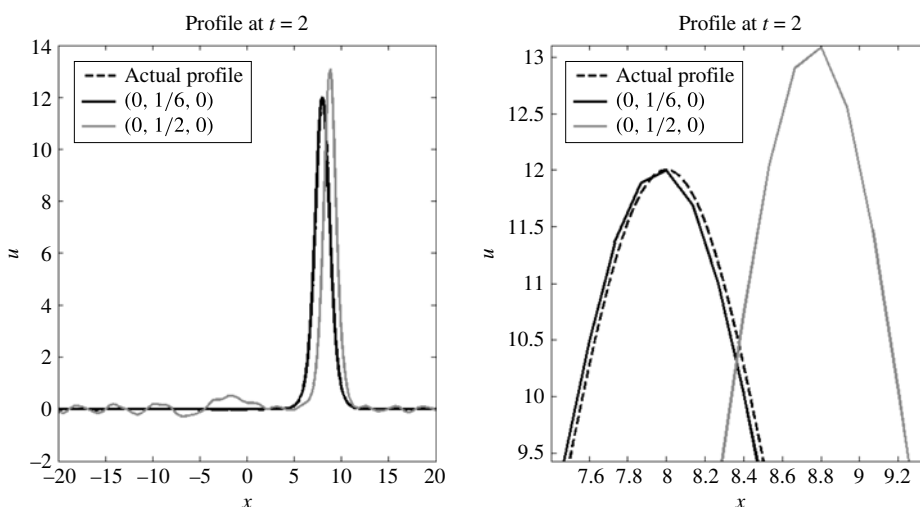| Scheme | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ |
|--------|----------|-----|----------|-----|----------|-----|
| $(0, \frac{1}{2}, 0)$ | 5.6843e−14 | 1.1714 | 0.0012 | 1.3088 | 14.3570 | 1.9179 |
| $(0, \frac{1}{6}, 0)$ | 5.3291e−14 | 0.7630 | 7.6271e−4 | 7.8989e−4 | 0.0077 | 0.2180 |

FIGURE 1. *The Zabusky–Kruskal scheme* $(0, \frac{1}{2}, 0)$ *versus the* $(0, \frac{1}{6}, 0)$ *scheme for the single-soliton problem* ($\mu = \frac{2}{15}$ *and* $\nu = \frac{1}{3}\mu^3$).

## 3. Method

### 3.1. Discretization approach

The method is based on forming the most general discretization of terms in the PDE and the characteristics of the CLaws with a given set of points. The discretizations are based on Taylor series expansions of the grid function about the point $(x_m, t_n)$,

$$u_{i,j} \approx u(x_m + i\mu, t_n + j\nu),$$

$$= u + i\mu u_x + j\nu u_t + \frac{(i\mu)^2}{2!}u_{xx} + \frac{(j\nu)^2}{2!}u_{tt} + i\mu j\nu u_{xt} + \dots \Big|_{(x_m, t_n)}.$$

It is assumed that $\nu = \lambda\mu^r$ for some fixed $r > 0$ and $\lambda > 0$, and the discretization of the PDE must be consistent, so that as $\mu \to 0$ the continuous terms being discretized are recovered. For example, using the box of points defined by $i = A, \dots, B$ and $j = C, \dots, D$ the discretizations (tildes denote discretizations of continuous terms throughout this paper) of the linear terms in the KdV equation are

$$\widetilde{u_{xxx}} = \frac{1}{\mu^3}\sum_{i=A}^{B}\sum_{j=C}^{D}\alpha_{i,j}u_{i,j}, \quad \widetilde{u_t} = \frac{1}{\nu}\sum_{i=A}^{B}\sum_{j=C}^{D}\beta_{i,j}u_{i,j}, \tag{3.1}$$

TABLE 2. *Maximum absolute errors in preserving the conserved quantities, for the two-soliton problem, and the times at which they occurred. The problem was numerically solved for* $t \in [0, 2]$, $\mu = \frac{2}{25}$, *and* $\nu = \frac{1}{3}\mu^3$.

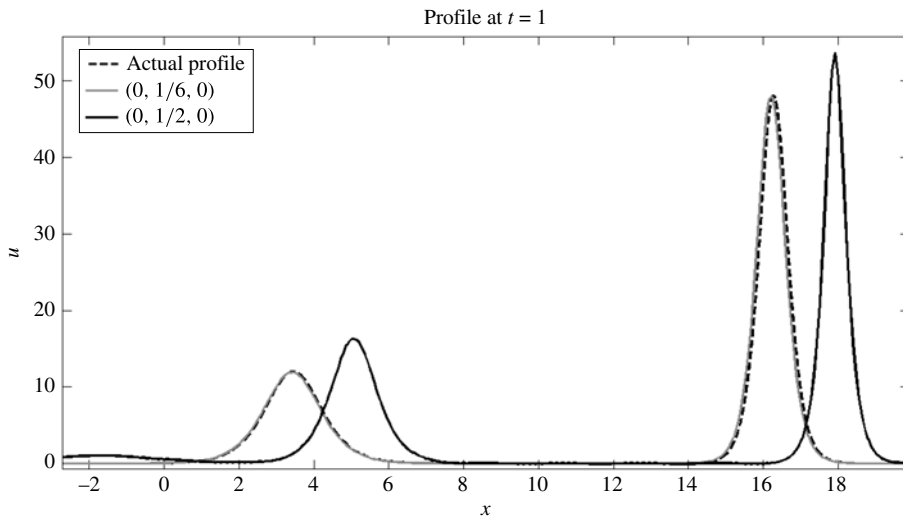| Scheme | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ |
|---|---|---|---|---|---|---|
| $(0, \frac{1}{2}, 0)$ | 2.4158e−13 | 1.4986 | 0.0250 | 1.5747 | 1.8102e+3 | 1.9068 |
| $(0, \frac{1}{6}, 0)$ | 2.8422e−13 | 1.4329 | 0.0152 | 1.1318 | 59.4218 | 1.5203 |

FIGURE 2. *The Zabusky–Kruskal scheme* $(0, \frac{1}{2}, 0)$ *versus the* $(0, \frac{1}{6}, 0)$ *scheme for the two-soliton problem* $(\mu = \frac{2}{25}$ *and* $\nu = \frac{1}{3}\mu^3)$.

where the necessary conditions on the coefficients are

$$0 = \sum_{i=A}^{B}\sum_{j=C}^{D}\alpha_{i,j}, \quad 0 = \sum_{i=A}^{B}\sum_{j=C}^{D}i\alpha_{i,j}, \quad 0 = \sum_{i=A}^{B}\sum_{j=C}^{D}i^2\alpha_{i,j}, \quad 3! = \sum_{i=A}^{B}\sum_{j=C}^{D}i^3\alpha_{i,j}, \quad (3.2)$$

$$0 = \sum_{i=A}^{B}\sum_{j=C}^{D}\beta_{i,j}, \quad 1 = \sum_{i=A}^{B}\sum_{j=C}^{D}j\beta_{i,j}. \quad (3.3)$$

However, for the scheme to be consistent, there are additional conditions that need to be satisfied. For $\widetilde{u_{xxx}} \to u_{xxx}$ as $\mu \to 0$, it is necessary that, for all $k, l \in \mathbb{N}$ (excluding the case $k = 3, l = 0$, otherwise the scheme will not converge to $u_{xxx}$; see (3.2)),

$$\lim_{\mu \to 0}\frac{1}{\mu^3}\sum_{ij}\alpha_{i,j}(i\mu)^k(j\nu)^l\frac{\partial^{k+l}}{\partial x^k \partial t^l}u = \left(\lambda^l\frac{\partial^{k+l}}{\partial x^k \partial t^l}u\right)\lim_{\mu \to 0}\sum_{ij}i^k j^l \alpha_{i,j}\mu^{lr+k-3} = 0.$$

Hence the additional conditions are

$$lr + k - 3 > 0 \quad \text{or} \quad \sum_{ij}i^k j^l \alpha_{i,j} = 0,$$

for all $k, l \in \mathbb{N}$, excluding the case $k = 3, l = 0$. For $l = 0$, the conditions (3.2) ensure that these conditions are satisfied. To avoid extra constraints on the coefficients $\alpha_{i,j}$, in addition to (3.2), $r$ must satisfy $r > (3 - k)/l$ for all $k \geq 0$ and $l \geq 1$. Clearly, as both $l$ and $k$ increase (so, as the order, $n \equiv k + l$, of the terms in the Taylor expansions increases), the restriction on $r$ decreases. Fixing $r > 3$ will ensure no additional constraints beyond (3.2) are required. To relax this restriction on $r$, the additional constraints from the Taylor expansions are imposed (starting from the lowest order terms and working upwards as necessary).

In order to have $\widetilde{u_t} \to u_t$ as $\mu \to 0$ it is necessary that, for all $k, l \in \mathbb{N}$, excluding $l = 1, k = 0$,

$$0 = \lim_{\mu \to 0}\frac{1}{\nu}\sum_{ij}\beta_{i,j}(i\mu)^k(j\nu)^l\frac{\partial^{k+l}}{\partial x^k \partial t^l}u = \left(\lambda^{l-1}\frac{\partial^{k+l}}{\partial x^k \partial t^l}u\right)\lim_{\mu \to 0}\sum_{ij}i^k j^l \beta_{i,j}\mu^{k+(l-1)r};$$

therefore

$$k + (l-1)r > 0 \quad \text{or} \quad \sum_{ij} i^k j^l \beta_{i,j} = 0.$$

The constraint on $r$ is immediately satisfied for $l \geqslant 1$. Thus the only remaining cases have $l = 0$, in which case $r < k$ for all $k \in \mathbb{N}$ is needed to avoid any extra conditions on the coefficients, hence $r < 1$. As before, to relax this constraint the conditions from successively higher order terms need to be applied.

From the above discussion, if (3.2) and (3.3) are the only conditions imposed on the discretizations, then for consistency $r > 3$ and $r < 1$ which is not possible. Therefore extra conditions need to be imposed on the coefficients $\alpha_{i,j}$ and $\beta_{i,j}$ so that there is a set of values of $r$ where both $\widetilde{u_t} \to u_t$ and $\widetilde{u_{xxx}} \to u_{xxx}$. There is not a unique way of doing this.

The KdV equation also contains the quadratic term $uu_x$. To discretize this term, products of Taylor series need to be examined:

$$u_{i,j}u_{k,l} = u^2 + \mu(i+k)uu_x + \nu(j+l)uu_t + \mu^2 ik u_x^2 + \mu\nu(jk+il)u_x u_t + \nu^2 jl u_t^2$$
$$+ \frac{\mu^2(i^2+k^2)}{2!}uu_{xx} + \mu\nu(ij+kl)uu_{xt} + \frac{\nu^2(j^2+l^2)}{2!}uu_{tt} + H.O.T.\bigg|_{(x_m,t_n)}. \quad (3.4)$$

Therefore the $uu_x$ term is discretized as

$$\widetilde{uu_x} = \frac{1}{\mu}\sum_{j=C}^{D}\left(\sum_{i=A}^{B}\sum_{k=i}^{B}\gamma_{i,j,k,j}u_{i,j}u_{k,j} + \sum_{i=A}^{B}\sum_{l=j+1}^{D}\sum_{k=A}^{B}\gamma_{i,j,k,l}u_{i,j}u_{k,l}\right), \quad (3.5)$$

with the necessary conditions

$$0 = \sum \gamma_{i,j,k,l}, \quad 1 = \sum(i+k)\gamma_{i,j,k,l}, \quad (3.6)$$

where $\sum$ is shorthand for the summation used in (3.5). Alternatively $u^2$ could be discretized and then the difference operator in the $m$ direction applied to it, though this will give a slightly less general discretization. Just as for the linear terms, extra conditions may be required on the coefficients to ensure that the method is consistent. The only terms in (3.4) that can remain as $\mu \to 0$, and so can cause the discretization to be inconsistent, are those that are purely time derivatives, and so are multiples of $\nu^n$. The first and second order terms are

$$\lim_{\mu \to 0}\frac{1}{\mu}\sum \gamma_{i,j,k,l}(j+l)\nu = \lim_{\mu \to 0}\lambda\sum\gamma_{i,j,k,l}(j+l)\mu^{r-1} = 0$$
$$\text{so } r > 1 \text{ or } \sum\gamma_{i,j,k,l}(j+l) = 0,$$
$$\lim_{\mu \to 0}\frac{1}{\mu}\sum \gamma_{i,j,k,l}(j^2+2jl+l^2)\nu^2 = \lim_{\mu \to 0}\lambda^2\sum\gamma_{i,j,k,l}(j^2+2jl+l^2)\mu^{2r-1} = 0$$
$$\text{so } r > \frac{1}{2} \text{ or } \sum\gamma_{i,j,k,l}(j^2+2jl+l^2) = 0,$$

and as $n$ increases the restriction on $r$ decreases. If $\frac{3}{2} < r < 3$ it is clear that no extra conditions are required on the $\gamma_{i,j,k,l}$ other than (3.6) to have a consistent discretization of the KdV equation.

### 3.2. Groebner bases

The method, described below, for finding discretizations results in large systems of polynomial equations to solve. These polynomials are in terms of the undetermined coefficients in the

discretizations. These large systems of equations are solved using Groebner bases; the reader is referred to [**6**, **7**, **9**, **22**] for further information. In principle the method should, for a given set of points, find any finite difference schemes that have analogues of the desired CLaws. However, in practice the method is limited by the amount of memory the computer has, because calculating the Groebner basis can require a large amount of memory. Calculating a Groebner basis is an expensive operation: for a set of polynomials, in $n$ variables, with total degree not exceeding $d$, the degree of the polynomials in the Groebner basis is bounded by $2(\frac{1}{2}d^2 + d)^{2^{n-1}}$. It can be shown that for sufficiently large $n$ there exist a constant $c$ and a set of polynomials such that every Groebner basis of the set contains an element that exceeds $2^{2^{cn}}$ [**7**]. However, the efficiency of the algorithm for a given problem is affected by the ordering of the polynomials and the ordering of the variables, so changing these can affect whether a problem is tractable. Due to the expense of finding Groebner bases, the method outlined in this paper is limited by the ability to calculate the Groebner basis. Thus, rather than searching for the most general discretizations, additional assumptions, such as that terms in the discretization are symmetric or antisymmetric about the centre of the discretization, can be imposed to reduce the number of variables and, in so doing, attempt to reduce the memory required to calculate the Groebner basis (see Appendix). Doing this may miss possible solutions; however, by choosing sensible ansätze the problem can be considerably simplified.

### 3.3.  *Recipe for finding schemes*

(1) Choose points for each term in the PDE to depend on. Then discretize each term in the PDE, using Taylor series approximations as described in § 3.1. At this stage the discretization may not be consistent. As there is not a unique way of imposing consistency, it seems best to search for methods that preserve CLaws and then, if any methods are found, impose consistency.

For example, the nonlinear term in the KdV equation can be discretized using two points[†], centred at $(x_m + \frac{1}{2}\mu, t_n)$,

$$\widetilde{uu_x} = \frac{1}{\mu}(\gamma_{0,0}u_0^2 + \gamma_{0,1}u_0u_1 + \gamma_{1,1}u_1^2),$$

with the necessary conditions (from equation (3.6)) that

$$\gamma_{0,0} + \gamma_{0,1} + \gamma_{1,1} = 0 \quad \text{and} \quad -\tfrac{1}{2}\gamma_{0,0} + \tfrac{1}{2}\gamma_{1,1} = 1. \tag{3.7}$$

(2) Choose the points on which each term of the characteristic of the desired CLaw will depend and then form the most general discretization of the characteristic with the chosen points. The discretization of the linear terms in the characteristic and the PDE should be centered in the same place (see Theorem A.1). Continuing the example, by choosing the characteristic of the second CLaw to depend on the same to points, let

$$\widetilde{u^2} = \eta_{0,0}u_0^2 + \eta_{0,1}u_0u_1 + \eta_{1,1}u_1^2,$$

with the necessary condition that

$$\eta_{0,0} + \eta_{0,1} + \eta_{1,1} = 1. \tag{3.8}$$

(3) In practice, one may not wish to seek the most general discretization of a term in the characteristic or the PDE. The complexity of the problem can be reduced (by reducing the

---

[†]As we are considering a discretization that only includes points from one time step, for clarity, the additional subscripts have been dropped.

number of undetermined coefficients in the discretizations) with an assumption about the discretization, for example, that the discretization is symmetric in space or time (see A.2). Imposing, in the example, that $\widetilde{uu_x}$ is antisymmetric about its centre and $\widetilde{u^2}$ is symmetric about its centre is achieved by setting

$$\gamma_{0,0} = -\gamma_{1,1}, \quad \gamma_{0,1} = 0 \tag{3.9}$$

and

$$\eta_{0,0} = \eta_{1,1}. \tag{3.10}$$

(4) Apply the discrete Euler operator (1.10) to the product of the discretized characteristic, $\widetilde{Q_i}$, and the discretized PDE, $\widetilde{\Delta}$. The difference scheme has a conservation law with the given characteristic if

$$0 = \mathrm{E}(\widetilde{Q_i} \cdot \widetilde{\Delta}). \tag{3.11}$$

In the example this yields

$$
\begin{aligned}
0 &= \mathrm{E}(\widetilde{u^2}\,\widetilde{uu_x}) \\
&= \frac{1}{\mu}u_{-1}^3(\eta_{0,0}\gamma_{0,1} + \eta_{0,1}\gamma_{0,0}) + \frac{2}{\mu}u_{-1}^2 u_0(\eta_{0,0}\gamma_{1,1} + \eta_{0,1}\gamma_{0,1} + \eta_{1,1}\gamma_{0,0}) \\
&\quad + \frac{3}{\mu}u_{-1}u_0^2(\eta_{0,1}\gamma_{1,1} + \eta_{1,1}\gamma_{0,1}) + \frac{4}{\mu}u_0^3(\eta_{0,0}\gamma_{0,0} + \eta_{1,1}\gamma_{1,1}) + \frac{3}{\mu}u_0^2 u_1(\eta_{0,0}\gamma_{0,1} + \eta_{0,1}\gamma_{0,0}) \\
&\quad + \frac{2}{\mu}u_0 u_1^2(\eta_{0,0}\gamma_{1,1} + \eta_{0,1}\gamma_{0,1} + \eta_{1,1}\gamma_{0,0}) + \frac{1}{\mu}u_1^3(\eta_{0,1}\gamma_{1,1} + \eta_{1,1}\gamma_{0,1}). \tag{3.12}
\end{aligned}
$$

Splitting (3.11) according to the coefficients of the $u_{i,j}$, $\nu$ and $\mu$ terms yields an overdetermined system of quadratic equations in the coefficients of the discretizations. In the example, the resulting system of equations is

$$
\begin{aligned}
\eta_{0,0}\gamma_{0,1} + \eta_{0,1}\gamma_{0,0} &= 0, \\
\eta_{0,0}\gamma_{1,1} + \eta_{0,1}\gamma_{0,1} + \eta_{1,1}\gamma_{0,0} &= 0, \\
\eta_{0,1}\gamma_{1,1} + \eta_{1,1}\gamma_{0,1} &= 0, \\
\eta_{0,0}\gamma_{0,0} + \eta_{1,1}\gamma_{1,1} &= 0. \tag{3.13}
\end{aligned}
$$

If symmetry conditions had been imposed (equations (3.9) and (3.10)) then this system would reduce to a single equation

$$2\eta_{0,1}\gamma_{0,0} = 0, \tag{3.14}$$

so $\eta_{0,1} = 0$, and hence $\widetilde{u^2} = \frac{1}{2}(u_1^2 + u_0^2)$ and $\widetilde{uu_x} = (1/\mu)(u_1^2 - u_0^2)$.

(5) Repeat steps (2)–(4) for any additional CLaws that one wishes to preserve, to create a system of constraints on the coefficients. In the example, if symmetry assumptions are not imposed, the resulting system of equations consists of (3.7), (3.8) and (3.13).

(6) Calculate the Groebner basis of the large system to see if there are any solutions.

(7) Solve the system. This may lead to several disjoint families of discretizations.

(8) Use the direct construction method [17] to construct the densities (this is effectively by inspection as the characteristic is known). Note that these densities may not be direct analogues of the continuous densities as there may be terms that vanish in the limit. Also, it may be necessary to add trivial densities to ensure that the discrete flux and densities tend to

the continuous flux and density as the step sizes tend to zero. For our simple example (with the symmetry conditions),

$$\widetilde{u^2}\widetilde{uu_x} = \frac{1}{2}(u_1{}^2 + u_0{}^2)\frac{1}{\mu}(u_1{}^2 - u_0{}^2) = \frac{(S_m - I)}{\mu}\left(\frac{1}{2}u_0{}^4\right).$$

(9) Check that the discretization of the PDE is consistent. If it is not consistent then impose additional constraints, as described in §3.1.

## 4.   Finding schemes

Having outlined the method for finding finite difference schemes that preserve CLaws, the results of applying the method to the KdV equation are shown below. Many discretizations have been found by the method and are discussed in [14]; these include the norm preserving scheme used in [3] and Furihata's scheme [13], that preserves the first and third CLaws. However, in the interests of space, only some of the novel implicit schemes found, that are as compact as possible (to eliminate the occurrence of spurious waves [3]), are presented here, as well some explicit schemes[†].

### 4.1.   Three compact implicit schemes that preserve the first and second CLaws

4.1.1.   *First scheme.*   The first scheme presented here (equation (4.1), referred to as the 12scheme) is actually a one-parameter family of schemes that preserves the first and second CLaws,

$$
\begin{aligned}
\widetilde{\Delta} &= \frac{1}{\nu}(S_n - I)\widetilde{G_1} + \frac{1}{\mu}(S_m - I)\widetilde{F_1} \\
&= u_t + uu_x + u_{xxx}|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2), \\
\widetilde{Q_2} &= \tfrac{1}{2}(S_n + I)(\tfrac{1}{2}u_{-1,0} + \tfrac{1}{2}u_{0,0}) = u|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2), \\
\widetilde{G_1} &= \tfrac{1}{2}u_{-1,0} + \tfrac{1}{2}u_{0,0} = u|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu^2), \\
\widetilde{F_1} &= \frac{1}{2}\epsilon(u_{0,0}(u_{0,0} + 2u_{0,1} + u_{-2,0} + u_{-2,1}) + u_{0,1}(u_{0,1} + u_{-2,0} + u_{-2,1}) \\
&\quad + u_{-2,0}(u_{-2,0} + 2u_{-2,1}) + u_{-2,1}{}^2) \\
&\quad + \left(\frac{1}{24} - \frac{1}{2}\epsilon\right)(u_{-1,0} + u_{-1,1})(u_{0,0} + u_{0,1} + u_{-2,0} + u_{-2,1} + u_{-1,0} + u_{-1,1}) \\
&\quad + \frac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0}) \\
&= \frac{1}{2}u^2 + u_{xx}\bigg|_{(x_m-\mu,t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2).
\end{aligned}
\tag{4.1}
$$

This discretization contains a parameter, $\epsilon$, in the $\widetilde{uu_x}$ term that we are free to chose. From Figure 3 it is clear that setting $\epsilon = 0$ (removing the dashed lines in the figure) will give the most compact discretization, and setting $\epsilon = \frac{1}{12}$ (removing the solid lines in the figure) will give the least compact discretization, for $\frac{1}{2}u^2$ in $\widetilde{F_1}$. Thus $\epsilon = 0$ will give the most compact discretizations of $uu_x$ in the KdV equation.

To find the densities of the second CLaw, given the characteristic, the direct construction method is used. To simplify the construction, note that the expression $\widetilde{Q_2}\widetilde{\Delta}$ can be split

---

[†]The simplifying assumptions used to find all these schemes are described in [14].
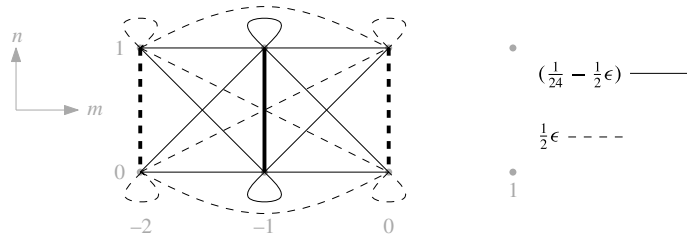
FIGURE 3. *The discretization of $\frac{1}{2}u^2$ in $\widetilde{F_1}$ for the first scheme. A line joining two points indicates that the product of the variable at the two endpoints is included in the discretization. The line style designates the coefficient of the product in the discretization; a bold line indicates double the value of the coefficient.*

according to the coefficients $\mu$ and $\nu$. As neither the density nor the flux of the second CLaw contain any $t$ derivatives, it is desirable that any $\nu$ terms are the result of applying the difference operator in the $n$ direction. Thus, the density is found by searching for a function $\widetilde{G_2}(u_{-2,0}, u_{-1,0}, u_{0,0}, u_{1,0})$ that satisfies

$$(S_n - I)\widetilde{G_2} = \text{coeff}\left(\widetilde{Q_2}\widetilde{\Delta}, \frac{1}{\nu}\right),$$

where $\text{coeff}(\widetilde{Q_2}\widetilde{\Delta}, 1/\nu)$ denotes the coefficient of $1/\nu$ in the expression $\widetilde{Q_2}\widetilde{\Delta}$. The discrete flux, $\widetilde{F_2}$, is then found by solving

$$\frac{(S_m - I)}{\mu}(\widetilde{F_2}(u_{-2,0}, u_{-1,0}, u_{0,0}, u_{-2,1}, u_{-1,1}, u_{0,1})) = \widetilde{Q_2}\widetilde{\Delta} - \frac{(S_n - I)}{\nu}\widetilde{G_2}.$$

The resulting densities are

$$\widetilde{G_2} = \tfrac{1}{8}(u_{0,0} + u_{-1,0})^2 = \tfrac{1}{2}u^2|_{(x_m - (\mu/2), t_n)} + \mathcal{O}(\mu^2),$$

$$\begin{aligned}
\widetilde{F_2} &= \frac{1}{96}(u_{-1,1} + u_{-2,0} + u_{-2,1} + u_{-1,0})(u_{-1,0} + u_{0,0} + u_{-1,1} + u_{0,1}) \\
&\quad \times (u_{-1,1} + u_{-1,0} + 12\epsilon(u_{0,0} + u_{-2,0} - 2u_{-1,1} + u_{0,1} + u_{-2,1} - 2u_{-1,0})) \\
&\quad + \frac{1}{8\mu^2}[(u_{-1,0} + u_{0,0} + u_{-1,1} + u_{0,1})(u_{-2,0} + u_{-2,1}) \\
&\quad + (u_{-1,0} + u_{-1,1})(u_{0,0} + u_{0,1}) - 3(u_{-1,0} + u_{-1,1})(u_{-1,0} + u_{-1,1})] \\
&= \tfrac{1}{3}u^3 + uu_{xx} - \tfrac{1}{2}u_x^2|_{(x_m - \mu, t_n + (\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2).
\end{aligned}$$

4.1.2. *Second scheme.* The second scheme discretizes $u_{xxx}$ in the same way as the first scheme (as does the third). The clearest difference between the first and second schemes is that the second scheme uses a wider computational stencil to discretize $u_t$ and $Q_2$. Another noteworthy feature is that the discretization for $uu_x$ factors into a discretization for $u$ and $u_x$, unlike for the first scheme. The second scheme is:

$$\begin{aligned}
\widetilde{\Delta} &= \frac{1}{4\nu}(S_n - I)(u_{-2,0} + u_{-1,0} + u_{0,0} + u_{1,0}) \\
&\quad + \frac{1}{\mu}(S_m - I)\left(\frac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0})\right) \\
&\quad - \frac{1}{24\mu}(u_{-2,0} - u_{1,1} - u_{1,0} + u_{-2,1})(u_{0,1} + u_{-1,1} + u_{0,0} + u_{-1,0}) \\
&= u_t + uu_x + u_{xxx}|_{(x_m - (\mu/2), t_n + (\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2),
\end{aligned} \tag{4.2}$$

$$\widetilde{Q_2} = \tfrac{1}{8}(S_n + I)(u_{-2,0} + u_{-1,0} + u_{0,0} + u_{1,0}) = u|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2),$$

$$\widetilde{G_1} = \tfrac{1}{4}(u_{-2,0} + u_{-1,0} + u_{0,0} + u_{1,0}) = u|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu^2),$$

$$\begin{aligned}
\widetilde{F_1} &= \frac{1}{24}u_{-1,1}u_{0,1} + \frac{1}{24}u_{-1,1}u_{0,0} + \frac{1}{24}u_{-1,0}u_{0,1} + \frac{1}{24}u_{-1,0}u_{0,0} + \frac{1}{24}u_{-2,1}u_{0,1} \\
&\quad + \frac{1}{24}u_{-2,1}u_{0,0} + \frac{1}{24}u_{-2,1}u_{-1,1} + \frac{1}{24}u_{-2,1}u_{-1,0} + \frac{1}{24}u_{-2,0}u_{0,1} + \frac{1}{24}u_{-2,0}u_{0,0} \\
&\quad + \frac{1}{24}u_{-2,0}u_{-1,1} + \frac{1}{24}u_{-2,0}u_{-1,0} + \frac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0}) \\
&= \frac{1}{2}u^2 + u_{xx}|_{(x_m-\mu,t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2),
\end{aligned}$$

$$\widetilde{G_2} = \tfrac{1}{32}(u_{-2,0} + u_{-1,0} + u_{0,0} + u_{1,0})^2 = \tfrac{1}{2}u^2|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu^2),$$

$$\begin{aligned}
\widetilde{F_2} &= \frac{1}{192}(u_{-2,0} + u_{-2,1} + u_{-1,0} + u_{-1,1})(u_{0,1} + u_{-1,1} + u_{0,0} + u_{-1,0}) \\
&\quad \times (u_{0,0} + u_{-2,0} + u_{-2,1} + u_{0,1}) \\
&\quad + \frac{1}{16\mu^2}(-2u_{-2,0}u_{-1,0} + 4u_{-2,0}u_{0,0} - 2u_{-2,0}u_{-1,1} + 4u_{-2,0}u_{0,1} + u_{-2,0}{}^2 \\
&\quad + 2u_{0,0}u_{0,1} + u_{0,1}{}^2 + 2u_{-2,1}u_{-2,0} - 2u_{-2,1}u_{-1,0} + 4u_{-2,1}u_{0,0} - 2u_{-2,1}u_{-1,1} \\
&\quad + 4u_{-2,1}u_{0,1} - 2u_{-1,1}{}^2 + u_{0,0}{}^2 + u_{-2,1}{}^2 - 2u_{-1,0}u_{0,1} - 2u_{-1,0}u_{0,0} - 4u_{-1,0}u_{-1,1} \\
&\quad - 2u_{-1,0}{}^2 - 2u_{-1,1}u_{0,1} - 2u_{-1,1}u_{0,0}) \\
&= \frac{1}{3}u^3 + uu_{xx} - \frac{1}{2}u_x^2|_{(x_m-\mu,t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2).
\end{aligned}$$

4.1.3. *Third scheme.* The third scheme is like the second scheme in that it uses a wide computational stencil to discretize $u_t$ and $Q_2$; however, unlike the second scheme, this scheme places a greater weight on grid points closer to the centre of the discretization. Finally, like the second scheme, its discretization of $uu_x$ neatly factorizes. The discretizations for the third scheme are:

$$\begin{aligned}
\widetilde{\Delta} &= \frac{1}{6\nu}(S_n - I)(u_{-2,0} + 2u_{-1,0} + 2u_{0,0} + u_{1,0}) \\
&\quad + \frac{1}{\mu}(S_m - I)\left(\frac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0})\right) \\
&\quad + \frac{1}{32\mu}(u_{0,0} + u_{-1,1} + u_{-1,0} + u_{0,1}) \\
&\quad \times (u_{0,0} - u_{-2,0} - u_{-2,1} + u_{0,1} + u_{1,1} + u_{1,0} - u_{-1,1} - u_{-1,0}), \quad\quad (4.3) \\
&= u_t + uu_x + u_{xxx}|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2),
\end{aligned}$$

$$\begin{aligned}
\widetilde{Q_2} &= \tfrac{1}{6}(S_n + I)(\tfrac{1}{2}u_{-2,0} + u_{-1,0} + u_{0,0} + \tfrac{1}{2}u_{1,0}) \\
&= u|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2),
\end{aligned}$$

$$\widetilde{G_1} = \tfrac{1}{6}(u_{-2,0} + 2u_{-1,0} + 2u_{0,0} + u_{1,0}) = u|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu^2),$$

$$\begin{aligned}
\widetilde{F_1} &= \frac{1}{32}(u_{-1,0} + u_{-2,1} + u_{-2,0} + u_{-1,1})(u_{0,0} + u_{-1,1} + u_{-1,0} + u_{0,1}) \\
&\quad + \frac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0}) \\
&= \frac{1}{2}u^2 + u_{xx}|_{(x_m-\mu,t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2),
\end{aligned}$$

$$\widetilde{G_2} = \tfrac{1}{72}(u_{-2,0} + 2u_{-1,0} + 2u_{0,0} + u_{1,0})^2 = \tfrac{1}{2}u^2|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu^2),$$

$$\widetilde{F_2} = \frac{1}{384}(u_{-1,1} + u_{-1,0} + u_{-2,1} + u_{-2,0})(u_{0,0} + u_{0,1} + u_{-1,0} + u_{-1,1})$$

$$\times \ (u_{0,0} + u_{-2,1} + 2u_{-1,0} + u_{-2,0} + 2u_{-1,1} + u_{0,1})$$

$$+ \frac{1}{\mu^2}\left(\frac{5}{24}u_{-2,0}u_{0,0} - \frac{1}{24}u_{-2,0}u_{-1,0} - \frac{1}{24}u_{-2,0}u_{-1,1} + \frac{5}{24}u_{-2,0}u_{0,1} + \frac{1}{24}u_{0,0}{}^2\right.$$

$$+ \frac{1}{24}u_{0,1}{}^2 - \frac{5}{24}u_{-1,1}{}^2 + \frac{1}{24}u_{-2,0}{}^2 + \frac{1}{12}u_{-2,1}u_{-2,0} - \frac{1}{24}u_{-2,1}u_{-1,0} + \frac{5}{24}u_{-2,1}u_{0,0}$$

$$- \frac{1}{24}u_{-2,1}u_{-1,1} - \frac{1}{24}u_{-1,1}u_{0,1} + \frac{1}{12}u_{0,0}u_{0,1} + \frac{5}{24}u_{-2,1}u_{0,1} + \frac{1}{24}u_{-2,1}{}^2 - \frac{1}{24}u_{-1,0}u_{0,0}$$

$$- \frac{5}{24}u_{-1,0}{}^2 - \frac{5}{12}u_{-1,0}u_{-1,1} - \frac{1}{24}u_{-1,0}u_{0,1} - \frac{1}{24}u_{-1,1}u_{0,0}\right)$$

$$= \frac{1}{3}u^3 + uu_{xx} - \frac{1}{2}u_x^2|_{(x_m - \mu, t_n + (\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2).$$

4.1.4. *A nonlinear stability consideration.* It is noteworthy that all three schemes (equations (4.1)–(4.3)) preserve the density of the second CLaw as $\widetilde{G_2} = \frac{1}{2}\widetilde{G_1}{}^2$. Therefore, the results of the discretization above can be averaged, $\bar{u}_{m-1/2,n} = \widetilde{G_1}$, to give an alternative discretization of the KdV equation that preserves the first and second conserved quantities exactly. It is tempting to think that this will give unconditional stability in the $\ell_2$-norm with periodic boundary conditions. However, the original schemes, from which the averages are calculated, may not themselves be unconditionally stable.

Let us consider the stability of the 12scheme (4.1) on a periodic domain so that $u_{1,j} = u_{M,j}$ for all $j$. More generally, we can consider schemes that exactly preserve the structure of the second CLaw so that it can be written in the form

$$(S_n - I)\left(\frac{1}{2}\left(\sum_{i=1}^{p} \alpha_i u_{i,j}\right)^2\right) + (S_m - I)\widetilde{F_2} = 0, \tag{4.4}$$

where $0 < \alpha_i \in \mathbb{R}$ and $p < M$ (note that for (4.1) $p = 2$, and that the other two schemes, (4.2) and (4.3), both have $p = 4$). Summing over the domain then yields the conserved quantity

$$\sum_{m=0}^{M-1}\left(\sum_{i=1}^{p} \alpha_i u_{m+i,n+j}\right)^2 = A, \tag{4.5}$$

where $A$ is a constant. We now seek to show that, as a consequence of (4.5) and the periodic boundary conditions, if $p$ divides $M$ then $|u_{i,j}|$ must be bounded.

The case $p = 1$ is trivial and occurs for the ten-point norm preserving scheme in [3]. So let us begin with $p = 2$. We wish to show that there exists some number $B_i \in \mathbb{R}$, such that, if $u_{1,j} \geqslant B_1 > 0$, then there exist numbers $B_i > 0$ such that $u_{i,j} \leqslant -B_i$ for $i$ even, and $u_{i,j} \geqslant B_i$ for $i$ odd. Therefore if $M$ is even, $u_{M,j} < 0$. But the periodic boundary conditions imply that $u_{M,j} = u_{1,j} > 0$; hence, the assumption that $u_{1,j} \geqslant B_1$ must be false.

For $p = 2$ we use $\alpha$ and $\beta$ in preference to $\alpha_1$ and $\alpha_2$ respectively and, as we are at a fixed time level, the $j$ have been dropped, so (4.5) implies that

$$-\sqrt{A} - \alpha u_m \leqslant \beta u_{m+1} \leqslant \sqrt{A} - \alpha u_m, \quad \forall m. \tag{4.6}$$

Let us assume that $m$ is odd and $u_m \geqslant B_m > 0$. Then (4.6) implies that

$$u_{m+1} \leqslant \frac{1}{\beta}(\sqrt{A} - \alpha u_m) \leqslant \frac{1}{\beta}(\sqrt{A} - \alpha B_m).$$

So, to ensure that $u_{m+1} \leqslant -B_{m+1} \leqslant 0$, $B_m$ must satisfy

$$\frac{1}{\beta}(\sqrt{A} - \alpha B_m) \leqslant -B_{m+1} \quad \text{so } B_m \geqslant \frac{1}{\alpha}(\sqrt{A} + \beta B_{m+1}).$$

Similarly, suppose that $m$ is even and that $-u_m \geqslant B_m > 0$. Then (4.6) yields

$$\frac{1}{\beta}(-\sqrt{A} + \alpha B_m) \leqslant \frac{1}{\beta}(-\sqrt{A} - \alpha u_m) \leqslant u_{m+1},$$

so for $u_{m+1} \geqslant B_{m+1}$, $B_m$ must satisfy

$$B_m \geqslant \frac{1}{\alpha}(\sqrt{A} + \beta B_{m+1}).$$

As the condition is the same for both cases, we have a well defined sequence to satisfy to ensure that $u_{i,j}$ oscillates about zero. Given $B_m > 0$, we can then find a suitable $B_{m-1} > 0$, etc. until we find $B_1 > 0$. This sequence yields

$$B_1 \geqslant \frac{1}{\alpha}(\sqrt{A} + \beta B_2) \geqslant \frac{\sqrt{A}}{\alpha} + \frac{\beta}{\alpha}\left(\frac{\sqrt{A}}{\alpha} + \frac{\beta}{\alpha}B_3\right) \geqslant \ldots \geqslant \frac{\sqrt{A}}{\alpha}\left(\sum_{i=0}^{M-2}\left(\frac{\beta}{\alpha}\right)^i\right) + \left(\frac{\beta}{\alpha}\right)^{M-1}B_M.$$

Hence if $B_1 \geqslant (\sqrt{A}/\alpha)(\sum_{i=0}^{M-2}(\beta/\alpha)^i)$ then $B_M \geqslant 0$, and so if $M$ is even then $u_{M,j} \leqslant 0$, which is a contradiction. Thus we have shown the solution must be bounded by

$$|u_{i,j}| < \begin{cases} \dfrac{\sqrt{A}}{\alpha}\dfrac{(1 - (\beta/\alpha)^{M-1})}{1 - (\beta/\alpha)}, & \alpha \neq \beta; \\[2ex] \dfrac{\sqrt{A}}{\alpha}(M - 1), & \alpha = \beta. \end{cases}$$

For $p > 2$, the proof should work as follows. Equation (4.5) implies that

$$\left(\sum_{i=1}^{p} \alpha_i u_{m+i,n+j}\right)^2 \leqslant A, \quad \forall m. \tag{4.7}$$

If $u_{1,j} \to \infty$ then

$$\sum_{i=2}^{p} \alpha_i u_{i,j} \to -\infty, \tag{4.8}$$

in order to satisfy (4.7). Now we also have the condition that $(\sum_{i=2}^{p+1} \alpha_{i-1} u_{i,j})^2 \leqslant A$. Because all the $\alpha_i$ are positive, (4.7) implies that $\sum_{i=2}^{p} \alpha_{i-1} u_{i,j} \to -\infty$, and so $u_{p+1,j} \to \infty$. In the same manner, we can move along the domain and use (4.7) to show that $u_{p+i,j}$ must have the same asymptotic behaviour as $u_{i,j}$, either tending to plus infinity, minus infinity, or remaining bounded. We can now relabel our periodic domain so that $u_{1,j} \to \infty$ and $u_{p,j} \to -\infty$ or is bounded. Thus if $p$ divides $M$ we have a contradiction, because $u_{i,j} = u_{M,j}$. Therefore, the solution must remain bounded. However, this gives no information on the size of this bound, which could be very large.

Therefore, the presence of a CLaw in the form (4.4), combined with periodic boundary conditions with the correct number of points, prevents the modes that cause the scheme to blow up from growing without limit.

### 4.2. An implicit scheme that preserves the first and third claws

Not only is it possible to find compact schemes that preserve the first two CLaws of the KdV equation, but also the first and third CLaws, However, they are less common and some have unintuitive discretizations (see [14]). Perhaps the best scheme found is the 13scheme,

$$
\begin{aligned}
0 = {} & \frac{1}{2\nu}(S_n - I)(u_{-1,0} + u_{0,0}) \\
& + \frac{1}{\mu}(S_m - I)\left( u_{-2,0}\left( \frac{1}{96}u_{-2,0} + \frac{1}{24}u_{-1,0} + \frac{1}{48}u_{0,0} + \frac{1}{96}u_{-2,1} + \frac{1}{48}u_{-1,1} + \frac{1}{96}u_{0,1} \right) \right. \\
& + u_{-1,0}\left( \frac{1}{24}u_{-1,0} + \frac{1}{24}u_{0,0} + \frac{1}{48}u_{-2,1} + \frac{1}{24}u_{-1,1} + \frac{1}{48}u_{0,1} \right) \\
& + u_{0,0}\left( \frac{1}{96}u_{0,0} + \frac{1}{96}u_{-2,1} + \frac{1}{48}u_{-1,1} + \frac{1}{96}u_{0,1} \right) \\
& + u_{-2,1}\left( \frac{1}{96}u_{-2,1} + \frac{1}{24}u_{-1,1} + \frac{1}{48}u_{0,1} \right) + u_{-1,1}\left( \frac{1}{24}u_{-1,1} + \frac{1}{24}u_{0,1} \right) + \left. \frac{1}{96}{u_{0,1}}^2 \right) \\
& + \frac{1}{2\mu^3}(S_n + I)(-u_{-2,0} + 3u_{-1,0} - 3u_{0,0} + u_{1,0}) \\
= {} & u_t + uu_x + u_{xxx}|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\mu^2) + \mathcal{O}(\nu^2), \hspace{2cm} (4.9)
\end{aligned}
$$

with characteristic and densities shown in Table 3. Two things should be noted about the density and flux of the third CLaw: firstly, that they are only first order approximations about their central point; and secondly, that the first term in the flux vanishes as the step sizes tend to zero, and so does not correspond to an expression in the continuous flux.

### 4.3. The three-parameter family

Inspired by the Zabusky–Kruskal scheme, the following assumptions were made to find explicit schemes that preserve the first and second CLaws. Their method is a two-step method with $\widetilde{u_t} = (1/\nu)(u_{0,1} - u_{0,-1})$, so for $\widetilde{u_{xxx}}$ to have the same centre point, and be symmetric, it must be $(1/\mu^3)(-u_{-2,0} + 2u_{-1,0} - 2u_{1,0} + u_{2,0})$. The discrete characteristic chosen was $\widetilde{Q_2} = u_{0,0}$ and finally $\widetilde{uu_x}$ was chosen to be the most general discretization of $uu_x$ with the five horizontal points centred at $(0,0)$. The result of these choices is a three-parameter $(\alpha, \beta, \gamma)$ family of finite difference methods for the KdV equation (shown in Table 4) that all preserve the first two CLaws. Setting $\alpha = 0$, $\beta = \frac{1}{2}$ and $\gamma = 0$ yields the Z–K scheme, the most compact scheme in the family.

A possible approach to choosing parameter values is to attempt to minimize the LTE of

$$
\begin{aligned}
\widetilde{uu_x} = {} & uu_x|_{(x_m,t_n)} + \left(\tfrac{5}{6} - \tfrac{2}{3}\alpha - \tfrac{4}{3}\beta - \tfrac{2}{3}\gamma\right)(u_{xxx}u + 2u_{xx}u_x)\mu^2 \\
& + \left(\tfrac{1}{2} - 2\alpha - \beta - \gamma\right)(u_{xxx}u_x + u_{xx}^2)\mu^3 + \mathcal{O}(\mu^4).
\end{aligned}
$$

To make this term a fourth-order approximation in space,

$$
\alpha = \alpha, \quad \beta = \tfrac{3}{4} + \alpha, \quad \gamma = -\tfrac{1}{4} - 3\alpha,
$$

and one such choice is $\alpha = \gamma = -\frac{1}{16}$ and $\beta = \frac{11}{16}$. Despite the free parameter, it is not possible to make the discretization for this term a fifth order approximation. Moreover, the $\widetilde{u_{xxx}}$ term is a second order approximation so there may not be any advantage in increasing the accuracy of the nonlinear term beyond this.

TABLE 3. *Densities, fluxes and characteristic for the 13scheme.*

$$\widetilde{Q_3} = \tfrac{5}{96}u_{0,0}{}^2 + \tfrac{1}{96}u_{-2,1}{}^2 + \tfrac{1}{96}u_{1,0}{}^2 + \tfrac{1}{96}u_{-2,0}{}^2 + \tfrac{5}{96}u_{-1,0}{}^2 + \tfrac{5}{96}u_{-1,1}{}^2 + \tfrac{5}{96}u_{0,1}{}^2 + \tfrac{1}{96}u_{1,1}{}^2 + \tfrac{1}{12}u_{-1,0}u_{0,0} + \tfrac{1}{48}u_{-1,0}u_{1,0} + \tfrac{1}{24}u_{0,0}u_{1,0} + \tfrac{1}{48}u_{-2,0}u_{-1,1}$$
$$+ \tfrac{1}{96}u_{-2,0}u_{0,1} + \tfrac{1}{48}u_{-1,0}u_{-2,1} + \tfrac{5}{96}u_{-1,0}u_{-1,1} + \tfrac{1}{24}u_{-1,0}u_{0,1} + \tfrac{1}{96}u_{-1,0}u_{1,1} + \tfrac{1}{96}u_{0,0}u_{-2,1} + \tfrac{1}{24}u_{0,0}u_{-1,1} + \tfrac{5}{96}u_{0,0}u_{0,1} + \tfrac{1}{48}u_{0,0}u_{1,1} + \tfrac{1}{96}u_{1,0}u_{-1,1}$$
$$+ \tfrac{1}{48}u_{1,0}u_{0,1} + \tfrac{1}{24}u_{-2,1}u_{-1,1} + \tfrac{1}{48}u_{-2,1}u_{0,1} + \tfrac{1}{12}u_{-1,1}u_{0,1} + \tfrac{1}{48}u_{-1,1}u_{1,1} + \tfrac{1}{24}u_{0,1}u_{1,1} + \tfrac{1}{48}u_{-2,0}u_{0,0} + \tfrac{1}{96}u_{-2,0}u_{-2,1} + \tfrac{1}{96}u_{1,0}u_{1,1} + \tfrac{1}{24}u_{-2,0}u_{-1,0}$$
$$+ \tfrac{1}{2\mu^2}(S_n + I)(u_{-2,0} - u_{-1,0} - u_{0,0} + u_{1,0})$$
$$= u^2 + 2u_{xx}|_{(x_m-(\mu/2),t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2)$$

$$\widetilde{G_1} = \tfrac{1}{2}u_{-1,0} + \tfrac{1}{2}u_{0,0} = u|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu^2)$$
$$\widetilde{F_1} = (u_{-2,0}(\tfrac{1}{96}u_{-2,0} + \tfrac{1}{24}u_{-1,0} + \tfrac{1}{48}u_{0,0} + \tfrac{1}{96}u_{-2,1} + \tfrac{1}{48}u_{-1,1} + \tfrac{1}{96}u_{0,1}) + u_{-1,0}(\tfrac{1}{24}u_{-1,0} + \tfrac{1}{24}u_{0,0} + \tfrac{1}{48}u_{-2,1} + \tfrac{1}{24}u_{-1,1} + \tfrac{1}{48}u_{0,1})$$
$$+ u_{0,0}(\tfrac{1}{96}u_{0,0} + \tfrac{1}{96}u_{-2,1} + \tfrac{1}{48}u_{-1,1} + \tfrac{1}{96}u_{0,1}) + u_{-2,1}(\tfrac{1}{96}u_{-2,1} + \tfrac{1}{24}u_{-1,1} + \tfrac{1}{48}u_{0,1}) + u_{-1,1}(\tfrac{1}{24}u_{-1,1} + \tfrac{1}{24}u_{0,1}) + \tfrac{1}{96}u_{0,1}{}^2)$$
$$+ \tfrac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0})$$
$$= \tfrac{1}{2}u^2 + u_{xx}|_{(x_m-\mu,t_n+(\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2)$$

$$\widetilde{G_3} = \tfrac{1}{192}u_{1,0}(18u_{0,0}u_{1,0} + 18u_{0,0}{}^2 + 10u_{1,0}{}^2 + 3u_{-1,0}u_{1,0} + 3u_{-1,0}{}^2 + 12u_{-1,0}u_{0,0}) - \tfrac{1}{2\mu^2}(u_{0,0} - u_{-1,0})(u_{0,0} - u_{-2,0})$$
$$= \tfrac{1}{3}u^3 - u_x^2|_{(x_m-(\mu/2),t_n)} + \mathcal{O}(\mu)$$
$$\widetilde{F_3} = -\tfrac{\mu}{192\nu}(2u_{-2,0}u_{-1,0}u_{-1,1} - 6u_{-2,0}u_{-1,0}u_{0,0} + 3u_{-2,0}u_{-1,0}u_{0,1} + u_{-2,0}u_{0,0}u_{0,1} + 3u_{-2,0}u_{-1,1}u_{0,1} - 3u_{-2,1}u_{-1,0}u_{0,0} - 3u_{-2,1}u_{0,0}u_{-1,1} - u_{-2,1}u_{0,0}u_{0,1}$$
$$- 2u_{-2,1}u_{-1,0}u_{-1,1} + 6u_{-2,1}u_{-1,1}u_{0,1} + u_{-1,0}u_{0,0}u_{0,1} + 2u_{-1,0}u_{-1,1}u_{0,1} - u_{-1,1}u_{0,0}u_{0,1} + u_{-1,1}u_{-2,0}{}^2 - u_{0,0}u_{-2,0}{}^2 - 2u_{-1,0}u_{0,0}u_{-1,1} + u_{-2,0}{}^2u_{0,1}$$
$$- 14u_{-1,0}{}^2u_{0,0} - u_{-2,1}{}^2u_{-1,0} + 2u_{-1,0}{}^2u_{0,1} + 14u_{-1,1}{}^2u_{0,1} + u_{-2,1}{}^2u_{-1,1} - u_{-2,1}{}^2u_{0,0} - 2u_{-1,1}{}^2u_{0,0} + u_{-2,1}{}^2u_{0,1} + 5u_{-1,1}{}^3 - 10u_{0,0}{}^3 - 5u_{-1,0}{}^3$$
$$+ 10u_{0,1}{}^3 + u_{-2,0}u_{-1,1}u_{-2,1} - u_{-2,0}{}^2u_{-1,0} - u_{-2,0}u_{0,0}u_{-2,1} - 4u_{-2,0}u_{-1,0}{}^2 - 2u_{-2,0}u_{0,0}{}^2 + 2u_{-2,0}u_{-1,1}{}^2 - u_{-2,0}u_{-1,0}u_{-2,1} + u_{-2,0}u_{0,1}{}^2$$
$$- u_{-2,1}u_{0,0}{}^2 - 2u_{-2,1}u_{-1,0}{}^2 + 4u_{-2,1}u_{-1,1}{}^2 + 2u_{-2,1}u_{0,1}{}^2 + u_{-1,0}u_{0,1}{}^2 - 17u_{-1,0}u_{0,0}{}^2 + 17u_{-1,1}u_{0,1}{}^2 - u_{-1,1}u_{0,0}{}^2 + u_{-2,0}u_{0,1}u_{-2,1})$$
$$+ \tfrac{1}{\mu\nu}[\tfrac{1}{4}(u_{0,1}(u_{-2,1} - u_{-2,0} + u_{-1,1} - u_{-1,0}) + u_{-1,1}(u_{-1,1} - u_{-1,0} - u_{-2,0}) - \tfrac{1}{2}u_{-2,1}u_{-1,1} + 2u_{-1,0}u_{-2,0})$$
$$+ \tfrac{1}{4}(u_{0,0}(u_{-2,1} - u_{-2,0} + u_{-1,1} - u_{-1,0}) + u_{-1,0}(u_{-2,1} + u_{-1,1} - u_{-1,0} + u_{-2,0}) - u_{-2,1}u_{-1,1})$$
$$+ \tfrac{1}{48}(u_{0,0}(u_{0,0} + u_{0,1} + u_{-2,1} + 2u_{-2,0} + 2u_{-1,1} + 4u_{-1,0}) + u_{0,1}{}^2 + 4u_{-1,0}{}^2 + 4u_{-1,1}{}^2 + 2u_{-2,1}u_{0,1} + u_{-2,0}u_{0,1} + 4u_{-1,1}u_{0,1} + 2u_{-1,0}u_{0,1} + u_{-2,0}{}^2 + u_{-2,1}{}^2$$
$$+ 2u_{-2,0}u_{-1,1} + 2u_{-1,0}u_{-2,1} + 4u_{-2,1}u_{-1,1} + 4u_{-1,0}u_{-1,1} + 4u_{-2,0}u_{-1,0} + u_{-2,0}u_{-2,1})\tfrac{1}{2\mu^2}(S_n + I)(u_{-2,0} - 2u_{-1,0} + u_{0,0})$$
$$+ \tfrac{1}{4\mu^4}(-(u_{-2,0} + u_{-2,1})(-(u_{-2,0} + u_{-2,1}) + 4(u_{-1,0} + u_{-1,1}) - 2(u_{0,0} + u_{0,1})) + (u_{0,0} + u_{0,1})((u_{0,0} + u_{0,1}) - 4(u_{-1,0} + u_{-1,1}))$$
$$+ 4(u_{-1,0} + u_{-1,1})^2) + \tfrac{1}{9216}(u_{0,0}{}^2 + u_{0,0}u_{0,1} + u_{0,0}u_{-2,1} + 2u_{-2,0}u_{0,0} + 2u_{0,0}u_{-1,1} + 4u_{-1,0}u_{0,0} + u_{0,1}{}^2 + 4u_{-1,0}{}^2 + 4u_{-1,1}{}^2$$
$$+ 2u_{-2,1}u_{0,1} + u_{-2,0}u_{0,1} + 4u_{-1,1}u_{0,1} + 2u_{-1,0}u_{0,1} + u_{-2,0}{}^2 + u_{-2,1}{}^2 + 2u_{-2,0}u_{-1,1} + 2u_{-1,0}u_{-2,1} + 4u_{-2,1}u_{-1,1} + 4u_{-1,0}u_{-1,1}$$
$$+ 4u_{-2,0}u_{-1,0} + u_{-2,0}u_{-2,1})^2$$
$$= u^2u_{xx} + 2u_tu_x + \tfrac{1}{4}u^4 + u_{xx}^2 - \tfrac{1}{4}uu_t\tfrac{\nu}{\mu^4}|_{(x_m-\mu,t_n+(\nu/2))} + \mathcal{O}(\mu) + \mathcal{O}(\nu)$$

The implicit schemes, when linearized about the constant solution, are all unconditionally stable [**14**]. However, this is not true for the explicit schemes found. Therefore, to provide a step size restriction for implementing the three-parameter family, a linear stability analysis (following [**29**]) is performed about the constant solution $u = \rho$. The resulting linear scheme depends on only a single parameter $\theta = \alpha + 2\beta + \gamma - 1$. If we then assume that $|\rho| \leqslant u_{\max}$ where $|u(x,t)| \leqslant u_{\max}$ then for linear stability it is necessary that

$$\nu \left( \mu^2 u_{\max} \left( 1 + \sqrt{3}|\theta| \right) + \frac{3\sqrt{3}}{2} \right) \leqslant \mu^3. \tag{4.10}$$

This suggests that increasing $|\theta|$ requires a slightly more severe step size restriction. However, its contribution is part of the $\mu^2$ term so there is very little difference between the linear stability of the schemes provided $|\theta|$ is not large.

### 4.4. The three-step explicit scheme

Since it is possible to preserve the first and second CLaws with an explicit scheme, can an explicit scheme be found to preserve the first and third CLaws? The answer is yes. Instead of two time steps, such a scheme was found by searching for a scheme with three time steps. The resulting difference scheme is

$$\begin{aligned} 0 &= \frac{1}{3\nu}(u_{0,1} - u_{0,-2}) + \frac{1}{4\mu}(u_{1,-1}u_{1,0} - u_{-1,-1}u_{-1,0}) \\ &\quad + \frac{1}{4\mu^3}(S_n + I)(-u_{-2,-1} + 2u_{-1,-1} - 2u_{1,-1} + u_{2,-1}) \\ &= u_t + uu_x + u_{xxx}|_{(x_m, t_n - (\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2). \end{aligned} \tag{4.11}$$

TABLE 4. *The three-parameter family and its densities.*

$u_{0,1} = u_{0,-1} + \frac{\nu}{\mu^3}(u_{-2,0} - 2u_{-1,0} + 2u_{1,0} - u_{2,0})$
$\qquad + \frac{\nu}{3\mu}(u_{-2,0}u_{0,0} + u_{1,0}{}^2 - u_{1,0}u_{-1,0} + u_{-2,0}u_{-1,0} - u_{-1,0}u_{0,0} - u_{2,0}{}^2)$
$\qquad + \frac{2\alpha\nu}{3\mu}(2u_{-1,0}u_{0,0} - u_{0,0}u_{2,0} - 2u_{-2,0}u_{-1,0} - 2u_{1,0}{}^2 + u_{2,0}{}^2 + u_{-2,0}{}^2 - u_{-2,0}u_{0,0} + 2u_{-1,0}u_{1,0})$
$\qquad + \frac{2\beta\nu}{3\mu}(2u_{-1,0}u_{0,0} - u_{-2,0}u_{-1,0} + u_{-1,0}u_{1,0} - 2u_{1,0}{}^2 - u_{-2,0}u_{0,0} + u_{2,0}{}^2 - u_{0,0}u_{1,0} + u_{-1,0}{}^2)$
$\qquad + \frac{2\gamma\nu}{3\mu}(u_{-1,0}u_{0,0} - u_{1,0}u_{2,0} + u_{-1,0}u_{1,0} - u_{1,0}{}^2 - u_{-2,0}u_{0,0} + u_{2,0}{}^2)$

---

$\widetilde{G_1} = \frac{1}{2}u_{0,0} + \frac{1}{2}u_{0,-1}$
$\qquad = u|_{(x_m, t_n - (\nu/2))} + \mathcal{O}(\nu^2)$
$\widetilde{F_1} = \frac{u_{-2,0} - u_{-1,0} - u_{0,0} + u_{1,0}}{2\mu^2} + \frac{1}{6}(u_{1,0}{}^2 + u_{-2,0}u_{0,0} + u_{-2,0}u_{-1,0})$
$\qquad + \frac{\alpha}{3}(u_{-2,0}{}^2 - 2u_{-2,0}u_{-1,0} - u_{-2,0}u_{0,0} + u_{-1,0}{}^2 + u_{-1,0}u_{1,0} - u_{1,0}{}^2 + u_{0,0}{}^2)$
$\qquad + \frac{\beta}{3}(u_{-1,0}{}^2 - u_{1,0}{}^2 - u_{-2,0}u_{0,0} - u_{-2,0}u_{-1,0} + u_{-1,0}u_{0,0} + u_{0,0}{}^2)$
$\qquad + \frac{\gamma}{3}(u_{-1,0}u_{0,0} - u_{1,0}{}^2 - u_{-2,0}u_{0,0} + u_{0,0}u_{1,0})$
$\qquad = \frac{1}{2}u^2 + u_{xx}|_{(x_m - (\mu/2), t_n)} + \mathcal{O}(\mu^2)$

---

$\widetilde{G_2} = \frac{1}{2}u_{0,0}u_{0,-1}$
$\qquad = \frac{1}{2}u^2|_{(x_m, t_n - (\nu/2))} + \mathcal{O}(\nu^2)$
$\widetilde{F_2} = \frac{u_{-2,0}u_{0,0} - 2u_{-1,0}u_{0,0} + u_{-1,0}u_{1,0}}{2\mu^2} + \frac{1}{6}(u_{-2,0}u_{-1,0}u_{0,0} + u_{-1,0}u_{1,0}{}^2 + u_{-2,0}u_{0,0}{}^2 - u_{-1,0}u_{0,0}{}^2)$
$\qquad + \frac{\alpha}{3}(u_{-1,0}{}^2u_{1,0} - 2u_{-2,0}u_{-1,0}u_{0,0} + u_{-2,0}{}^2u_{0,0} + 2u_{-1,0}u_{0,0}{}^2 - u_{-1,0}u_{1,0}{}^2 - u_{-2,0}u_{0,0}{}^2)$
$\qquad + \frac{\beta}{3}(u_{-1,0}{}^2u_{0,0} - u_{-1,0}u_{1,0}{}^2 - u_{-2,0}u_{0,0}{}^2 + 2u_{-1,0}u_{0,0}{}^2 - u_{-2,0}u_{-1,0}u_{0,0})$
$\qquad + \frac{\gamma}{3}(u_{-1,0}u_{0,0}u_{1,0} - u_{-1,0}u_{1,0}{}^2 - u_{-2,0}u_{0,0}{}^2 + u_{-1,0}u_{0,0}{}^2)$
$\qquad = uu_{xx} - \frac{1}{2}u_x{}^2|_{(x_m - (\mu/2), t_n)} + \mathcal{O}(\mu^2)$

The characteristic is

$$\widetilde{Q_3} = u_{0,-1}u_{0,0} + \frac{1}{\mu^2}(S_n + I)(u_{-1,-1} - 2u_{0,-1} + u_{1,-1})$$

$$= u^2 + 2u_{xx}|_{(x_m, t_n - (\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2).$$

The density and flux for the first CLaw are

$$\widetilde{G_1} = \tfrac{1}{3}u_{0,-2} + \tfrac{1}{3}u_{0,-1} + \tfrac{1}{3}u_{0,0} = u|_{(x_m, t_n - \nu)} + \mathcal{O}(\nu^2),$$

$$\widetilde{F_1} = \frac{1}{4}(u_{0,-1}u_{0,0} + u_{-1,-1}u_{-1,0}) + \frac{1}{4\mu^2}(S_n + I)(u_{-2,-1} - u_{-1,-1} - u_{0,-1} + u_{1,-1})$$

$$= \frac{1}{2}u^2 + u_{xx}|_{(x_m - (\mu/2), t_n - (\nu/2))} + \mathcal{O}(\nu^2) + \mathcal{O}(\mu^2).$$

The density and flux for the third CLaw are

$$\widetilde{G_3} = \frac{1}{3}u_{0,-1}u_{0,0}u_{0,-2} + \frac{1}{3\mu^2}(u_{1-1}u_{0,-2} + u_{0,-1}u_{1-2} - u_{1,-1}u_{1-2} - 2u_{0,0}u_{0-2}$$

$$\quad - 2u_{0,-1}u_{0,0} - u_{0,-1}u_{0,-2} + u_{1-2}u_{0,0} + u_{0,0}u_{-1,-1} + u_{0,0}u_{1,-1} + u_{-1-2}u_{0,0})$$

$$= \frac{1}{3}u^3 + \frac{2}{3}uu_{xx} - \frac{1}{3}u_x^2|_{(x_m, t_n - \nu)} + \mathcal{O}(\mu) + \mathcal{O}(\nu),$$

$$\widetilde{F_3} = \frac{1}{3\mu\nu}(u_{-1,-1}u_{0,-2} + u_{0,-1}u_{0,0} - u_{0,0}u_{-1,-1} - u_{0,-1}u_{0,-2} - u_{-1-2}u_{0,0} + u_{-1,0}u_{0,-2})$$

$$\quad + \frac{1}{4}u_{0,-1}u_{0,0}u_{-1,-1}u_{-1,0} + \frac{1}{2\mu^2}\left( \frac{1}{2}(u_{0,0}u_{0,-1}(u_{0,-1} + u_{-2,0})\right.$$

$$\quad + u_{-1,-1}u_{-1,0}(u_{-1,0} + u_{1,-1})) - u_{-1,-1}u_{-1,0}u_{0,-1} - u_{-1,-1}u_{-1,0}u_{0,0}$$

$$\quad + \frac{1}{2}(u_{0,-1}u_{0,0}(u_{-2,-1} + u_{0,0}) + u_{-1,-1}u_{-1,0}(u_{-1,-1} + u_{1,0}))$$

$$\quad \left. - u_{0,-1}u_{0,0}u_{-1,0} - u_{-1,-1}u_{0,-1}u_{0,0} \right)$$

$$\quad + \frac{1}{4\mu^4}((S_n + I)(u_{-1,-1} - 2u_{0,-1} + u_{1,-1}))((S_n + I)(u_{-2,-1} - 2u_{-1,-1} + u_{0,-1}))$$

$$= \frac{1}{4}u^4 + u_{xx}^2 + u^2u_{xx} + \frac{4}{3}u_xu_t - \frac{2}{3}uu_{xt}|_{(x_m - (\mu/2), t_n - (\nu/2))} + \mathcal{O}(\nu) + \mathcal{O}(\mu).$$

Note that $\widetilde{F_3}$ contains $\nu$ terms and to ensure these terms did not blow up in the limit, as the step sizes tend to zero, trivial CLaws were added during the reconstruction.

As this is an explicit scheme, just as for the three-parameter family, a linear stability analysis is performed about the constant solution $u = \rho$ to obtain any necessary step size restrictions. The resulting necessary but not sufficient condition for stability is that

$$\frac{3\nu}{\mu^3}\left( \left|\frac{2u_{\max}}{4}\right|\mu^2 + \frac{3\sqrt{3}}{4} \right) \leqslant 3, \tag{4.12}$$

where it is assumed that $|\rho| \leqslant u_{\max}$. This is very similar to the stability condition found for the three-parameter family of explicit schemes (4.10), even though this is a three-step rather than a two-step scheme. From numerical experiments it was found that $\nu = \frac{1}{4}\mu^3$ was needed for stability. This is of the same order of magnitude as the necessary condition found in (4.12), and is similar to the condition for the three-parameter family, where using $\nu = \frac{1}{3}\mu^3$ produced stable results.

## 5. Basic numerics

Having found the above discretizations, they were compared by using them to solve the one-soliton initial value problem for the KdV equation, $0 = u_t + uu_x + u_{xxx}$. The exact solution to this problem on an infinite domain is

$$u(x,t) = 3c \operatorname{sech}^2\left(\left(\frac{\sqrt{c}}{2}\right)(x - ct)\right).$$

For the numerics a periodic domain was used with $x \in [-20, 20]$, and $t = [0, 2]$. The resulting solution profiles were compared with that of the exact, infinite domain, solution. This is fine provided that $ct$ is small enough that the soliton does not get close to the boundary. Fairly coarse discretizations were used to emphasize the qualitative differences between the different schemes. In order to solve the implicit schemes, *fsolve* in MATLAB was allowed to run until the error reached the default tolerance, so that the difference between the schemes should be due to their different discretizations rather than due to the nonlinear solver. The error in preserving the different conservation laws was calculated by approximating the conserved quantities using

$$\mu \sum_m u_{m,n}, \quad \mu \sum_m u_{m,n}^2, \quad \mu \sum_m \frac{1}{3} u_{m,n}^3 + \frac{1}{\mu^2} u_{m,n}(u_{m-1,n} - 2u_{m,n} + u_{m+1,n})$$

at each time step.

The schemes were also compared with the eight-point multisymplectic method of Ascher and McLachlan [2, 3],

$$
\begin{aligned}
0 = {} & \frac{1}{8\nu}(S_n - I)(u_{-2,0} + 3u_{-1,0} + 3u_{0,0} + u_{1,0}) \\
& + \frac{1}{2\mu^3}(S_n + I)(-u_{-2,0} + 3u_{-1,0} - 3u_{0,0} + u_{1,0}) \\
& + \frac{1}{64\mu}((u_{0,0} + u_{1,0} + u_{0,1} + u_{1,1})^2 - (u_{-2,0} + u_{-1,0} + u_{-2,1} + u_{-1,1})^2),
\end{aligned}
$$

and their narrow box scheme

$$
\begin{aligned}
0 = {} & \frac{1}{2\nu}(S_n - I)(u_{-1,0} + u_{0,0}) + \frac{1}{8\mu}(S_m - I)(u_{-1,1} + u_{-1,0})^2 \\
& + \frac{1}{2\mu^3}(S_n + I)(-u_{-2,0} + 3u_{-1,0} - 3u_{0,0} + u_{1,0}),
\end{aligned}
$$

which was found using a finite volume discretization. In particular, Ascher and McLachlan conducted numerics for the alternative form of the KdV equation (note that for this form of the KdV equation the density of the 3rd CLaw is $\frac{1}{3}u^3 + \delta uu_{xx}$),

$$0 = u_t + uu_x + \delta u_{xxx},$$

with $\delta = 0.022^2$ and initial condition $u(x, 0) = \cos(\pi x)$ on a periodic domain with $x \in [0, 2]$. This is the original numerical problem of Zabusky and Kruskal [30] when they discovered solitons, and is referred to here as the Z–K problem. This problem is an interesting numerical test for the various schemes because, for this choice of $\delta$, the nonlinear part of the KdV equation is more significant than in the soliton problem ($\delta = 1$); the problem becomes close to the inviscid Burger's equation which develops shocks; hence this should be a good problem to see how the different schemes handle the nonlinear part of the KdV equation. Ascher and McLachlan showed, when studying this problem, that the multisymplectic scheme stayed remarkably smooth despite a very coarse discretization, whereas the other schemes could not

cope with the discretization, and the norm preserving scheme, even though it did not blow up, resembled noise. They concluded that this was due to the compactness of the discretization of the multisymplectic scheme and that its linearization is the most accurate unconditionally stable box semi-discretization at the steady state [**2**]. Thus it is interesting to replicate this experiment for the schemes found in this paper to see if any of the schemes outperform the multisymplectic scheme.

The results of various numerical tests are shown in Tables 5–9. The tables show the maximum absolute errors in preserving the first three conserved quantities and the times these maximum errors occurred. For the soliton problem, the position of the solution at $t = 2$ is calculated by locating the maximum point at the final time step. $M$ denotes the number of spatial steps and $N$ denotes the number of time steps (so $\mu = 40/M$ and $\nu = 2/N$). For the implicit schemes the errors of the averages given by $u(x_m - 1/2\mu, t_n) \approx \frac{1}{2}(u_{m,n} + u_{m-1,n})$ are also included.

The results in Table 5 show that as $\theta$ (the single parameter in the linearization) decreases, the amplitude and speed of the soliton decrease. For $\theta = -\frac{2}{3}$ the position of the numerical soliton appears to be remarkably close to the actual solution and the third CLaw is approximated better. All the solutions developed a wave-train following the soliton; however, the $(0, 0, 0)$ scheme's wave-train has a very small amplitude and the $(0, \frac{1}{6}, 0)$ scheme's wave-train was smaller still (see, for example, Figure 1). Reducing the number of spatial steps shows how the different schemes cope with coarse discretizations. At $M = 250$ the $(0, \frac{2}{3}, 0)$ scheme developed sawtooth waves, at $M = 200$ the Z–K scheme has clear sawtooth waves and the $(-\frac{1}{16}, \frac{11}{16}, -\frac{1}{16})$ scheme resembles noise, for $M = 100$ the $(0, 0, 0)$ scheme has become a lower slower soliton with a larger wave-train following; however, when $M = 50$ the $(0, \frac{1}{6}, 0)$ scheme still looks like a soliton, coping with the very coarse discretization incredibly well, far better than the Z–K scheme even though its discretization is not compact or symmetric.

Table 6 displays the errors from different schemes (within the three-parameter family), all with $\theta = -\frac{2}{3}$. All these schemes show similar results: their profiles all are very close to the actual soliton and have a very small wave-train following the numerical soliton, the smallest of which belongs to the $(0, \frac{1}{6}, 0)$ scheme. They all preserve the third CLaw better than the Z–K scheme, though there are clear differences between the schemes. However, overall the evidence suggests that schemes with the same linearization behave very similarly. The remarkable behaviour of the $\theta = -\frac{2}{3}$ schemes persists with different speed solitons and with different choices of step sizes. These schemes are not generally symmetric; however, the one-parameter family of schemes, given by $(\frac{1}{6}\beta, \beta, \frac{1}{6} - 3\beta)$, yields discretizations with $\widetilde{uu_x}$ antisymmetric and $\theta = -\frac{2}{3}$. One such scheme is $(\frac{5}{12}, \frac{1}{4}, -\frac{7}{12})$, but its results show that the added symmetry has not led to notably better errors. The good preservation of the third CLaw could be because it is locally preserved; however, using the method of this paper no discrete characteristic for a third CLaw has been found for any of the schemes in the three-parameter family. The good preservation of the third CLaw may be a consequence of preserving the phase speed better than the other schemes. However, in Figure 4 one can see that the profile of the $(-1, 1, -\frac{2}{3})$ scheme is closest to the actual position of the soliton but the $(0, \frac{1}{6}, 0)$ scheme preserves the third CLaw significantly better than it. This new two-parameter family ($\theta = -\frac{2}{3}$ schemes) outperforms the Zabusky–Kruskal scheme (which has the most compact discretization for the nonlinear term) in the numerical tests conducted; it remains unknown why this is the case.

The results of numerics for the other schemes are shown in Tables 7–9. These clearly demonstrate that the averaged 12scheme (see § 4.1) preserves the second conserved quantity exactly. The 12scheme has one parameter in the nonlinear term. It is not clear how to choose this; however, there seem to be two sensible *a priori* choices. The choice $\epsilon = 0$ makes the nonlinear term as compact as possible. This scheme can cope well with very coarse meshes (see Figure 5) for the Z–K problem; however, for the soliton problem, the numerical soliton travels slower than the actual solution and the third CLaw is not preserved very well. The alternative is $\epsilon = -1/24$, which causes the linearization to be as compact as possible. This
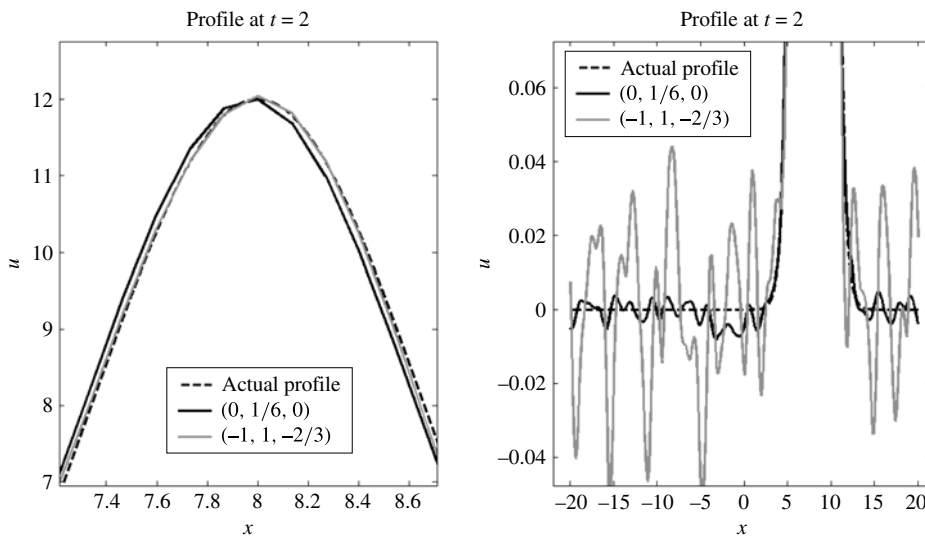
FIGURE 4. Different $(\alpha, \beta, \gamma)$ schemes, $M = 300$ and $N = 5064$, with the same linearization: $\theta = -\frac{2}{3}$.
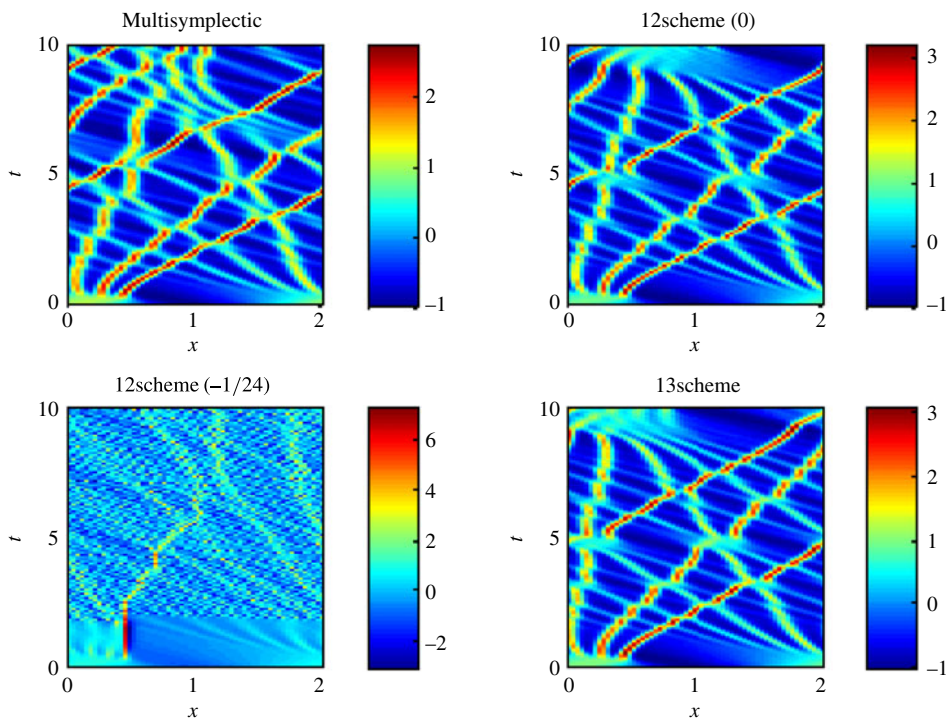


FIGURE 5. The results of simulating the Z–K problem using $\nu = \frac{1}{5}\mu$, $M = 60$, $N = 1500$ with different schemes; note that the 12 scheme with $\epsilon = -1/24$ has failed to cope with the coarse discretization.

scheme's soliton travels faster, closer to the actual position, and the third CLaw is better preserved. The cost of this is that the scheme cannot cope with the very coarse discretization used for the Z–K problem (see Figure 5). An alternative suggestion for choosing the free parameter is to tune it to best preserve the third CLaw, or some other property, for a given problem and mesh size of interest.

TABLE 5. *Absolute errors in preserving the conserved quantities for the three-parameter family, and the location and height of the soliton at the final time step, with different values of $\theta$, when solving the single-soliton problem with $c = 4$. The problem was numerically solved for $t \in [0, 2]$, $M = 300$, $N = 2532$, $\mu = \frac{2}{15}$ and $\nu = \frac{1}{3}\mu^3$.*

| Scheme | $\theta$ | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ | Soliton | $x$ |
|---|---|---|---|---|---|---|---|---|---|
| $(0, \frac{2}{3}, 0)$ | $1/3$ | 4.9738e−14 | 0.2986 | 0.0017 | 1.6603 | 45.4066 | 1.8555 | 13.8973 | 9.3333 |
| $(-\frac{1}{16}, \frac{11}{16}, -\frac{1}{16})$ | $1/4$ | 4.9738e−14 | 1.906 | 0.0016 | 1.7899 | 34.1630 | 1.9005 | 13.7255 | 9.2000 |
| $(0, \frac{1}{2}, 0)$ | $0$ | 5.6843e−14 | 1.1714 | 0.0012 | 1.3088 | 14.3570 | 1.9179 | 13.0876 | 8.8000 |
| $(0, \frac{1}{6}, 0)$ | $-2/3$ | 5.3291e−14 | 0.7630 | 7.6271e−04 | 7.8989e−04 | 0.0077 | 0.2180 | 11.9992 | 8 |
| $(0, 0, 0)$ | $-1$ | 6.0396e−14 | 0.6888 | 6.5568e−04 | 7.8989e−04 | 2.1504 | 1.9589 | 11.5458 | 7.7333 |

TABLE 6. *Absolute errors in preserving the conserved quantities, and the location and height of the soliton at the final time step, for different schemes from the three-parameter family, all with $\theta = -2/3$, for the single-soliton problem with $c = 4$. The problem was numerically solved for $t \in [0, 2]$, $M = 300$, $\mu = \frac{2}{15}$, $N = 5064$ and $\nu = \frac{1}{6}\mu^3$.*

| Scheme | $\theta$ | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ | Soliton | $x$ |
|---|---|---|---|---|---|---|---|---|---|
| $(\frac{1}{3}, -\frac{1}{6}, \frac{1}{3})$ | $-2/3$ | 6.3949e−14 | 1.7528 | 1.9097e−04 | 3.9494e−04 | 0.0595 | 1.9617 | 11.9864 | 8 |
| $(\frac{1}{3}, 0, 0)$ | $-2/3$ | 6.0396e−14 | 1.2251 | 1.9099e−04 | 3.9494e−04 | 0.0345 | 1.9633 | 11.9931 | 8 |
| $(0, \frac{1}{6}, 0)$ | $-2/3$ | 6.7502e−14 | 0.6509 | 1.9068e−04 | 3.9494e−04 | 0.0034 | 0.2062 | 11.9993 | 8 |
| $(0, 0, \frac{1}{3})$ | $-2/3$ | 5.3291e−14 | 0.1896 | 1.9054e−04 | 3.9494e−04 | 0.0062 | 1.7686 | 11.9924 | 8 |
| $(-1, 1, -\frac{2}{3})$ | $-2/3$ | 5.3291e−14 | 0.2192 | 2.0237e−04 | 2 | 0.4731 | 1.9538 | 12.0398 | 8 |
| $(\frac{5}{12}, \frac{1}{4}, -\frac{7}{12})$ | $-2/3$ | 6.3949e−14 | 1.7946 | 7.6485e−04 | 7.8988e−4 | 0.0167 | 1.9645 | 12.0033 | 8 |

TABLE 7. *Absolute errors in preserving the conserved quantities for the Z–K problem for $t \in [0, 10]$, $\nu = \frac{1}{5}\mu$, $M = 60$, $N = 1500$.*

| Scheme | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ |
|---|---|---|---|---|---|---|
| 13scheme | 8.6307e−13 | 9.2667 | 0.1434 | 0.8133 | 823.4132 | 0.8000 |
| Averaged | 8.6284e−13 | 9.26667 | 0.0304 | 0.8400 | 585.8678 | 0.8600 |
| 12scheme (0) | 2.2238e−13 | 6.8867 | 0.1224 | 0.8400 | 0.1861 | 9.0667 |
| Averaged | 2.2225e−13 | 6.8867 | 2.0505e−10 | 9.9467 | 0.0772 | 0.7200 |
| 12scheme (−1/24) | 4.4190e−13 | 9.9467 | 1.2388 | 1.8667 | 2.7482 | 9.8333 |
| Averaged | 4.4169e−13 | 9.9467 | 1.3100e−10 | 5.8067 | 0.8393 | 7.1333 |
| Narrow box | 77.3553 | 8.4333 | 1.8715e+06 | 0.5467 | 3.3493e+06 | 0.5267 |
| Averaged | 77.3553 | 8.3800 | 1.1263e+06 | 0.5467 | 1.3372e+09 | 0.5467 |
| Multisymplectic | 4.2492e−13 | 2.3400 | 0.1271 | 0.7733 | 0.1517 | 0.6733 |
| Averaged | 4.2472e−13 | 2.3400 | 0.0279 | 0.7933 | 0.0889 | 0.6667 |

TABLE 8. *The single-soliton problem, $c = 8$, for $t \in [0, 2]$ and $\nu = \mu$. $M = 1600$, $\mu = 0.0250$, $N = 80$.*

| Scheme | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ | Soliton | $x$ |
|---|---|---|---|---|---|---|---|---|
| 13scheme | 1.1369e−13 | 1.0500 | 0.0127 | 2 | 0.0506 | 1.8000 | 23.9234 | 15.7500 |
| Averaged | 1.1369e−13 | 0.3000 | 0.0139 | 2 | 0.0356 | 1.8000 | 23.9198 | 15.7375 |
| 12scheme (0) | 4.7322e−12 | 1.9500 | 4.8120e−04 | 1.1000 | 0.2009 | 1.1250 | 24.0150 | 15.7000 |
| Averaged | 4.6825e−12 | 1.9500 | 1.1579e−10 | 1.9500 | 0.1942 | 1.1250 | 24.0067 | 15.6875 |
| 12scheme (−1/24) | 2.3590e−12 | 0.4000 | 3.3325e−04 | 1.1000 | 0.2045 | 1.1250 | 24.0332 | 15.7000 |
| Averaged | 2.2666e−12 | 0.4000 | 3.5470e−11 | 0.4000 | 0.1994 | 1.1250 | 24.0226 | 15.7125 |
| Multisymplectic | 6.5796e−12 | 1.5000 | 6.1595e−04 | 1.1000 | 0.1932 | 1.1250 | 24.0161 | 15.7000 |
| Averaged | 6.5583e−12 | 1.4750 | 1.5395e−04 | 1.1000 | 0.1868 | 1.1250 | 24.0043 | 15.6875 |
| Narrow box | 5.9828e−12 | 0.3250 | 1.1763e−04 | 1.1000 | 0.2003 | 1.1250 | 24.0299 | 15.7000 |
| Averaged | 5.9899e−12 | 0.3250 | 2.3526e−04 | 1.1000 | 0.1951 | 1.1250 | 24.0213 | 15.7125 |

The eight-point scheme that preserves the first and third CLaws does not perform very well. This scheme does not preserve the energy (third conserved quantity) particularly well, considering that it preserves the third CLaw in characteristic form. However, it is none the less interesting because it shows that it is possible to preserve the third CLaw using a compact stencil and so obtain better behaviour on a coarse grid. In fact, it has the same linearization as the multisymplectic scheme against which it was compared, hence, like the multisymplectic scheme, it copes with the coarse discretization well, remaining smooth (see Figure 5). As the time step increases compared to the spatial step, it preserves the third CLaw better compared to the other schemes (compare Tables 8 and 9). The multisymplectic scheme's performance also improves with this ratio of step sizes, especially when it is compared to the narrow box scheme. This demonstrates that the choice of step sizes affects the relative performance of the schemes.

The explicit schemes have a major problem: if a fine spatial mesh is required then the number of time steps required is prohibitive, as $\nu = \mathcal{O}(\mu^3)$, and so the implicit schemes are far more efficient. This is the case for the Z–K problem. However, Table 9 shows the results of using the explicit schemes to solve the single-soliton problem where the step sizes have been chosen so that the explicit schemes take approximately the same amount of time to solve as the implicit schemes. With this choice of step sizes both the explicit schemes preserve the third CLaw very well, better than any of the implicit schemes.

TABLE 9. *Absolute errors in preserving the conserved quantities, and the location and height of the soliton at the final time step, for the single-soliton problem, $c = 8$, for $t \in [0, 2]$. The implicit schemes used $M = 400$, $N = 100$ ($\nu = \frac{1}{5}\mu$) and the explicit schemes used $M = 400$, $N = 8000$.*

| Scheme | 1st CLaw | $t$ | 2nd CLaw | $t$ | 3rd CLaw | $t$ | Soliton | $x$ |
|---|---|---|---|---|---|---|---|---|
| (0,1/6,0) | 9.9476e−14 | 0.3135 | 7.0056e−04 | 0.1237 | 0.0277 | 1.4775 | 23.9688 | 15.9000 |
| Three-step | 1.0658e−13 | 0.4600 | 0.0041 | 1.7610 | 2.0372e−04 | 0.9270 | 23.9161 | 16.0000 |
| 13scheme | 4.2633e−14 | 0.4600 | 0.0419 | 1.8000 | 1.8919 | 1.8000 | 23.7463 | 15.6000 |
| Averaged | 7.8160e−14 | 0.5400 | 0.0030 | 0.0600 | 1.3874 | 1.8000 | 23.6408 | 15.6500 |
| 12scheme (0) | 1.7764e−13 | 1.6400 | 0.0160 | 1.8000 | 0.5220 | 1.8000 | 23.9307 | 15.7000 |
| Averaged | 1.5632e−13 | 1.5600 | 1.8190e−12 | 1.5600 | 0.3160 | 1.8000 | 23.7588 | 15.6500 |
| 12scheme (−1/24) | 5.6843e−13 | 2 | 0.0300 | 1.8000 | 0.1531 | 0.9800 | 24.1134 | 15.8000 |
| Averaged | 6.0396e−13 | 2 | 9.2655e−12 | 2 | 0.1812 | 1.0200 | 24.0948 | 15.8500 |
| Multisymplectic | 3.3396e−13 | 1.9800 | 0.0178 | 1.7800 | 0.4768 | 1.7800 | 23.8721 | 15.7000 |
| Averaged | 3.3396e−13 | 1.9800 | 0.0044 | 1.7800 | 0.3163 | 1.8200 | 23.7976 | 15.7500 |
| Narrow box | 1.1369e−13 | 0.8600 | 0.0080 | 1.8000 | 0.1226 | 0.0200 | 24.1479 | 15.9000 |
| Averaged | 1.2079e−13 | 0.8400 | 0.0155 | 1.8000 | 0.2996 | 1.0200 | 24.0252 | 15.8500 |

## 6. *Summary and discussion*

The main result of this paper has been to develop a method to symbolically find discretizations of a partial difference equation that locally preserves as many CLaws as possible. The method developed is a brute force approach that is applicable to PDEs, with polynomial nonlinearities, on a fixed mesh. The method requires solving a large overdetermined system of quadratic equations; this is the method's main limiting factor. Currently, the Groebner basis is calculated to solve the system, which is very expensive in time and memory. Thus, an important area for future research is to find a more efficient method for solving this highly structured system of equations.

Despite the computational cost in general, the method is practical for searching for one-step methods in time. The method has been used to find new eight-point implicit discretizations and new explicit discretizations for the KdV equation that preserve the first and second CLaws together, and the first and third CLaws together. The method also finds the existing norm preserving scheme and Furihata's scheme, demonstrating that the method has the potential to find useful discretizations. It is an open problem as to whether the second and third CLaws of the KdV equation can be preserved together. The author suspects that if it is possible, this will require at least two time steps; this has the potential to allow parasitic waves to form. However, the extra constraint on the behaviour, which comes from preserving the additional CLaw, may prevent them from occurring. Since currently it has not been possible to preserve all the CLaws desired, this raises the question: what is the most important CLaw to preserve?

When conducting basic numerics, some of the new methods seem to perform very well, such as the two explicit methods and the 12scheme. In particular, the numerical schemes, discussed here, behave comparably to a multisymplectic scheme, and can outperform it for conserving a given CLaw. Hence if no multisymplectic structure is known for an equation that has CLaws, this brute force approach may yield a good method. However, apart from very coarse discretizations, a compact scheme found by a volume preserving technique performed the best. Nevertheless, some of the new eight-point schemes found can be considered to preserve the second CLaw exactly whilst having a compact stencil, which is a very desirable property.

## Appendix. *Simplifying assumptions*

THEOREM A.1. *The linear terms in the characteristic of a CLaw of a difference equation, and the linear terms of the difference equation, must share the same central point.*

*Proof.* The theorem is proved here for ODEs; the proof for two-dimensional PDEs is found in [**14**]. Let

$$P = \sum_{j=-A}^{A} \alpha_j u_j \quad \text{such that } \alpha_{-A}, \alpha_A \neq 0 \quad \text{and} \quad Q = \sum_{i=B}^{C} \beta_i u_i,$$

so that $P$ has an odd number of terms and is centred at the point $j = 0$ (see Figure A.1). If $P$ has an even number of terms then the following reasoning will still apply by considering $i + \frac{1}{2}, j + \frac{1}{2}, A + \frac{1}{2}, B + \frac{1}{2}, C + \frac{1}{2} \in \mathbb{Z}$ (see Figure A.2). The case $B < -A < 0 < C < A$ is proved here; the other cases follow by symmetry or similar reasoning.

We require that the product of $P$ and $Q$ is in the kernel of the Euler operator, so

$$0 = E(PQ) = \sum_{k=-A}^{A} \alpha_k S^{-k} Q + \sum_{k=B}^{C} \beta_k S^{-k} P \tag{A.1}$$
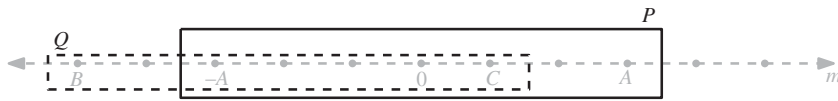
FIGURE A.1. *The linear terms in the characteristic and ordinary difference equation.*
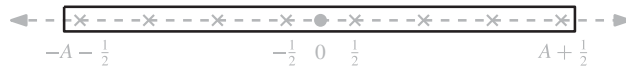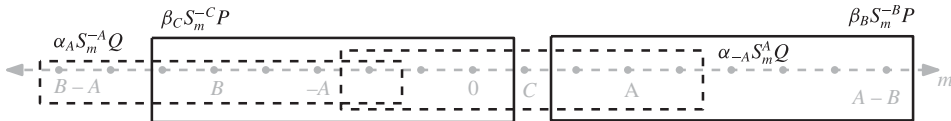


FIGURE A.2. *A term with an even number of points.*



FIGURE A.3. *The Euler operator acting on a product of linear terms. The points in the extreme shifts of $P$ and $Q$ in (A.1) are enclosed.*

$$
= \sum_{k=-A}^{A} \left( \alpha_k \sum_{i=B}^{C} \beta_i u_{i-k} \right) + \sum_{k=B}^{C} \left( \beta_k \sum_{j=-A}^{A} \alpha_j u_{j-m} \right)
$$

$$
= \sum_{k=-A}^{A} \sum_{i=B}^{C} \alpha_k \beta_j (u_{j-k} + u_{k-j})
$$

$$
= \sum_{k=-A}^{A} \left( \sum_{j=B}^{-C-1} \alpha_k \beta_j u_{j-k} + \sum_{j=-C}^{C} (\alpha_k \beta_j + \alpha_{-k} \beta_{-j}) u_{j-k} + \sum_{j=C+1}^{-B} \alpha_{-k} \beta_{-j} u_{j-k} \right). \qquad \text{(A.2)}
$$

The extreme shifts of $P$ and $Q$ occurring in (A.1) are depicted graphically in Figure A.3.

Because $C < |B|$, it is clear from (A.1) and Figure A.3 that $u_{A-B}$ is the rightmost term in the summation ($A - B$ is the maximum grid point that occurs). This term only occurs in the expression $\beta_B S_m^{-B} P$; hence its coefficient is $\beta_B \alpha_A$. Similarly, $u_{B-A}$ is the most extreme shift of the dependent variable in the negative direction, and it only occurs in $\alpha_A S_m^{-A} Q$, so its coefficient is also $\beta_B \alpha_A$. For the summation to be zero the coefficients of these terms must vanish, therefore (as by assumption $\alpha_A \neq 0$) $\beta_B = 0$. Thus $Q = \sum_{i=B+1}^{C} \beta_i u_i$. If $C < |B-1|$ the above reasoning is repeated to show that $\beta_{B+1} = 0$. This process is continued until $Q = \sum_{i=-C}^{C} \beta_i u_i$. $\qquad \square$

It is now apparent that only the middle term from (A.2) remains. This equation is satisfied if $\alpha_k = \alpha_{-k}$ and $\beta_{-j} = -\beta_j$. Therefore if one of the linear terms is symmetric about the centre point and the other is antisymmetric then their product is in the kernel of the Euler operator. This observation generalizes to the following theorem to include the two-dimensional case (the proof is found in [14]).

THEOREM A.2. *The discrete Euler operator applied to the product of a linear sum of terms with 180° rotational symmetry and a linear sum of terms with 180° antisymmetry is zero.*

It must be noted that there exist products of linear terms, not in this form, that are also in the kernel of the Euler operator. For example, when there is only one independent variable, if

$$
P = \beta_{-1} u_{-1} + (\beta_{-1} + \beta_1) u_0 + \beta_1 u_1,
$$

$$
Q = \frac{-\alpha \beta_{-1}}{\beta_1} u_{-1} + \frac{\alpha(\beta_{-1} - \beta_1)}{\beta_1} u_0 + \alpha u_1,
$$

then $E(PQ) = 0$. Theorem A.2 is analogous to the fact, from the continuous theory, that

$$E(u_{,m_1 x, n_1 t} u_{,m_2 x, n_2 t}) = (-1)^{m_1+n_1} D_x^{m_1} D_t^{n_1} u_{,m_2 x, n_2 t} + (-1)^{m_2+n_2} D_x^{m_2} D_t^{n_2} u_{,m_1 x, n_1 t} = 0,$$

for $m_1 + n_1$ odd and $m_2 + n_2$ even, where $u_{,m_i x, n_i t} \equiv D_x^{m_i} D_t^{n_i} u$.

For the KdV equation the first three CLaws in characteristic form satisfy the following conditions. From the first CLaw:

$$\text{(a) } E(u_t) = 0, \quad \text{(b) } E(uu_x), \quad \text{(c) } E(u_{xxx}) = 0. \tag{A.3}$$

From the second CLaw:

$$\text{(a) } E(u(u_t)) = 0, \quad \text{(b) } E(u(uu_x)) = 0, \quad \text{(c) } E(u(u_{xxx})) = 0. \tag{A.4}$$

From the third CLaw:

$$\text{(a) } E(u^2 u_t) = 0, \quad \text{(b) } E(2u_{xx} u_t) = 0, \quad \text{(c) } E(2u_{xx} u_{xxx}) = 0,$$
$$\text{(d) } E(u^2(uu_x)) = 0, \quad \text{(e) } E(u^2 u_{xxx} + 2u_{xx}(uu_x)) = 0. \tag{A.5}$$

Thus, if the discretizations for the linear terms are antisymmetric for odd derivatives and symmetric for even derivatives, equations (A.3a), (A.3c), (A.4a), (A.4c), (A.5b) and (A.5c) are immediately satisfied by Theorem A.2.

Having seen the benefits of symmetry assumptions for the linear terms, one might wish to make assumptions about the nonlinear terms in the discretizations. By assuming rotational symmetry and antisymmetry about the centre of the discretizations some simplification is obtained but no equations from (A.4) and (A.5) are automatically satisfied. A $k$th order polynomial term is discretized by a sum of terms of the form $\alpha_{i_1,j_1,i_2,j_2,\dots,i_k,j_k} u_{i_1,j_1} u_{i_2,j_2} \dots u_{i_k,j_k}$. The term $u_{-i_1,-j_1} u_{-i_2,-j_2} \dots u_{-i_k,-j_k}$ is $180°$ opposite the previously mentioned term and we label its coefficient as $\alpha_{-i_k,-j_k,\dots,-i_1,-j_1}$ and we can assume without loss of generality that the discretization is formed by a summation of pairs of terms that are $180°$ opposite one another. The discretization is symmetric if for each pair

$$\alpha_{i_1,j_1,i_2,j_2,\dots,i_k,j_k} = \alpha_{-i_k,-j_k,\dots,-i_1,-j_1}$$

and antisymmetric if

$$\alpha_{i_1,j_1,i_2,j_2,\dots,i_k,j_k} = -\alpha_{-i_k,-j_k,\dots,-i_1,-j_1}.$$

This is illustrated in Figure A.4 for a quadratic term discretized on a line. (Note that, just as for linear terms, suitable adjustments must be made for schemes not centred at $(0,0)$, that is, those schemes with an even number of points in one direction.)

Before proceeding to study a general polynomial term let us consider a simple example. To avoid subscripts, let a hat denote the coefficient of a term that is $180°$ opposite a term already used in an expression, so, for our example,

$$P = \alpha u_{-1} u_0 + \hat{\alpha} u_1 u_0 + \gamma u_1 u_3 + \hat{\gamma} u_{-1} u_{-3} \quad \text{and} \quad Q = \beta u_1 + \hat{\beta} u_{-1}, \tag{A.6}$$

and

$$P \cdot Q = (\alpha\beta u_{-1} u_0 u_1 + \hat{\alpha}\hat{\beta} u_1 u_0 u_{-1}) + (\alpha\hat{\beta} u_{-1}^2 u_0 + \hat{\alpha}\beta u_1^2 u_0))$$
$$+ (\gamma\beta u_1 u_3 u_1 + \hat{\gamma}\hat{\beta} u_{-1} u_{-3} u_{-1}) + (\hat{\gamma}\beta u_{-1} u_{-3} u_1 + \gamma\hat{\beta} u_1 u_3 u_{-1}). \tag{A.7}$$
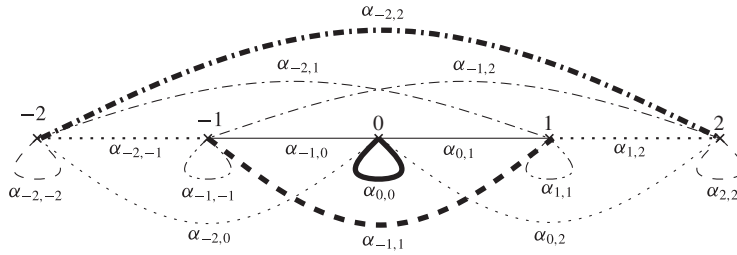
FIGURE A.4. *A graphical representation of the discretization of a quadratic term of an ODE. Each line represents the coefficient of the product of the dependent variable at the two points the line is connecting. Products that are 180° opposite to each other are depicted by matching lines. Coefficients depicted by bold lines will be zero in an antisymmetric discretization.*

The product of the two terms can be split into pairs of terms which are opposite to one another. We require that $E(P \cdot Q) = 0$; expanding this out yields the following conditions on the coefficients:

$$
\begin{aligned}
u_{-1}u_1 &: (\alpha\beta + \hat{\alpha}\hat{\beta}) = 0, \\
u_{-2}u_2 &: (\hat{\gamma}\beta + \gamma\hat{\beta}) = 0, \\
u_1u_2 &: (\alpha\beta + \hat{\alpha}\hat{\beta}) = 0, \quad u_{-1}u_{-2} : (\alpha\beta + \hat{\alpha}\hat{\beta}) = 0, \\
u_{-1}^2 &: \alpha\hat{\beta} = 0, \quad u_1^2 : \hat{\alpha}\beta = 0, \\
u_0u_1 &: 2\alpha\hat{\beta} = 0, \quad u_0u_{-1} : 2\hat{\alpha}\beta = 0, \\
u_{-2}^2 &: \gamma\beta = 0, \quad u_2^2 : \hat{\gamma}\hat{\beta} = 0, \\
u_0u_2 &: 2\gamma\beta = 0, \quad u_{-2}u_0 : 2\hat{\gamma}\hat{\beta} = 0, \\
u_2u_4 &: \hat{\gamma}\beta = 0, \quad u_{-2}u_{-4} : \gamma\hat{\beta} = 0.
\end{aligned}
\tag{A.8}
$$

Adding and subtracting the coefficients of terms that are opposite one another (for example, $u_{-1}^2$ and $u_1^2$) gives an equivalent system of equations:

$$
\begin{aligned}
(\alpha\beta + \hat{\alpha}\hat{\beta}) &= 0, \\
(\hat{\gamma}\beta + \gamma\hat{\beta}) &= 0, \\
(\alpha\beta + \hat{\alpha}\hat{\beta}) &= 0, \\
\alpha\hat{\beta} + \hat{\alpha}\beta &= 0, \quad \alpha\hat{\beta} - \hat{\alpha}\beta = 0, \\
\gamma\beta + \hat{\gamma}\hat{\beta} &= 0, \quad \gamma\beta - \hat{\gamma}\hat{\beta} = 0, \\
\hat{\gamma}\beta + \gamma\hat{\beta} &= 0, \quad \hat{\gamma}\beta - \gamma\hat{\beta} = 0.
\end{aligned}
\tag{A.9}
$$

Now if we assume that $P$ is a symmetric discretization, that is, $\hat{\alpha} = \alpha$ and $\hat{\gamma} = \gamma$, and $Q$ is an antisymmetric discretizaton, that is, $\hat{\beta} = -\beta$, or vice versa ($P$ is antisymmetric and $Q$ is symmetric), then the left hand column of (A.9) is satisfied, so the number of equations that need to be solved has been more than halved. The reason more than half of the equations are satisfied is because the coefficients of self-symmetric terms, such as $u_{-1}u_0u_1$ in $P \cdot Q$, must vanish before the Euler operator is even applied, and the coefficients of self-symmetric terms, such as $u_2u_{-2}$, in $E(P \cdot Q)$ must also vanish. If instead we assume that both $P$ and $Q$ are symmetric or both are antisymmetric then the right hand column of (A.9) is satisfied. This is still a major simplification, but more than half the equations remain because the self-symmetric terms remain.

For the general case, suppose that $P$ has coefficients denoted by $\alpha$ and $Q$ has coefficients denoted by $\beta$ and, for convenience,

$$\alpha = \alpha_{i_1,j_1,i_2,j_2,\dots,i_k,j_k}, \quad \hat{\alpha} = \alpha_{-i_k,-j_k,\dots,-i_1,j_1},$$
$$\beta = \beta_{i_{k+1},j_{k+1},\dots,i_K,j_K}, \quad \hat{\beta} = \beta_{-i_K,-j_K,\dots,-i_{k+1},-j_{k+1}}.$$

As seen in the example (A.7) the product of $P$ and $Q$ can be considered as consisting of pairs of terms that are opposite to one another. Applying the Euler operator to each term in one of these pairs gives

$$E\left(\left(\alpha\prod_{l=1}^{k}u_{i_l,j_l}\right)\left(\beta\prod_{l=k+1}^{K}u_{i_l,j_l}\right)\right) = \alpha\beta\left(\sum_{q=1}^{k}\prod_{\substack{l=1\\l\neq q}}^{K}u_{i_l-i_q,j_l-j_q}\right),$$

$$E\left(\left(\hat{\alpha}\prod_{l=1}^{k}u_{-i_l,-j_l}\right)\left(\hat{\beta}\prod_{l=k+1}^{K}u_{-i_l,-j_l}\right)\right) = \hat{\alpha}\hat{\beta}\left(\sum_{q=1}^{k}\prod_{\substack{l=1\\l\neq q}}^{K}u_{-i_l+i_q,-j_l+j_q}\right).$$

Thus the coefficient of $\prod_{\substack{l=1\\l\neq q}}^{K}u_{i_l-i_q,j_l-j_q}$ must be a sum of multiples of the coefficients of the form $\alpha\beta$ which must vanish if the product is in the kernel of the Euler operator. Similarly, the coefficient of its opposite term, $\prod_{\substack{l=1\\l\neq q}}^{K}u_{-i_l+i_q,-j_l+j_q}$, must be a sum of multiples of coefficients of the form $\hat{\alpha}\hat{\beta}$ which must also vanish. An equivalent system of equations (which must also vanish) is obtained by adding and subtracting these two coefficients together. This gives two sets of equations: one set is formed by a sum of terms of the form $(\alpha\beta + \hat{\alpha}\hat{\beta})$ and the other set by sums of terms of the form $(\alpha\beta - \hat{\alpha}\hat{\beta})$. Thus if $\alpha = \hat{\alpha}$ and $\beta = -\hat{\beta}$, that is, one of the polynomial terms is a symmetric discretization and the other is an antisymmetric discretization, the first set of equations is satisfied. Alternatively, if $\alpha = \hat{\alpha}$ and $\beta = \hat{\beta}$, or $\alpha = -\hat{\alpha}$ and $\beta = -\hat{\beta}$, the second set of equations is satisfied, that is, both discretizations are symmetric or antisymmetric, respectively. Thus by assuming antisymmetry or symmetry on the appropriate polynomial terms the number of equations that need to be satisfied is approximately halved. As in the example, terms that are self-symmetric in $P \cdot Q$ and $E(P \cdot Q)$ will vanish for the symmetric antisymmetric case but not the other cases. Therefore making symmetry assumptions on the polynomial terms in the characteristic and PDE seems a sensible ansatz.

For discretizing the KdV equation, this ansatz is imposed by insisting that $\widetilde{uu_x}$ has $180°$ antisymmetry. This ansatz, combined with the previous ansatz on the linear terms, is equivalent to imposing that the discretization preserves the discrete symmetry of the KdV equation,

$$t \mapsto -\hat{t}, \quad x \mapsto -\hat{x}, \quad \text{KdV} \to -u_{\hat{t}} - uu_{\hat{x}} - u_{\hat{x}\hat{x}\hat{x}} = 0.$$

By adding additional restrictions (more severe ansätze) the number of variables (undetermined coefficients in the discretizations) is reduced, so the computer may be able to find a solution to the Groebner basis. Some possible extra restrictions are as follows:
– If a term in (A.3)–(A.5) is in the kernel of the Euler operator because it a total difference in either $x$ or $t$ then, instead of the usual discrete Euler operator, one can impose that the term is in the kernel of the discrete Euler operator that treats the $u$ at different time or space levels as independent variables, respectively. The appropriate operators, when there are two independent variables, are

$$E_{m+i} := \sum_{j} S_n^{-j}\frac{\partial}{\partial u_{i,j}} \quad \text{and} \quad E_{n+j} := \sum_{i} S_m^{-i}\frac{\partial}{\partial u_{i,j}}.$$

For example, condition (A.5a) is a consequence of the fact that $u^2 u_t = D_t(\frac{1}{3}u^3)$, so, rather than insisting that $E(\widetilde{u^2 \tilde{u}_t}) = 0$, we can require that $E_{m+i}(\widetilde{u^2 \tilde{u}_t}) = 0$ for $i \in \mathbb{Z}$. This amounts to insisting that $\widetilde{u^2 \tilde{u}_t} = (S_n - I)g$ for some function $g([u])$, rather than insisting on the less stringent requirement that $\widetilde{u^2 \tilde{u}_t} = (S_n - I)g + (S_m - I)f$ for some functions $f([u])$ and $g([u])$.

– A more restrictive assumption than $180°$ symmetry or antisymmetry is imposing symmetry in one direction (for example, symmetry in the $x$ direction if discretizing a $t$ derivative) and antisymmetry or symmetry in the other direction, depending on whether an odd or even derivative is being discretized.

– An alternative assumption is to impose that the nonlinear terms, such as $\widetilde{uu_x}$ and $\widetilde{u^2}$, factorize. This will again lead to cubic equations in the coefficients rather than quadratic equations.

## References

**1.** L. M. ALONSO, 'On the Noether map', *Lett. Math. Phys.* 3 (1979) no. 5, 419–424.

**2.** U. M. ASCHER and R. I. MCLACHLAN, 'Multisymplectic box schemes and the Korteweg–de Vries equation', *Appl. Numer. Math.* 48 (2004) no. 3–4, 255–269.

**3.** U. M. ASCHER and R. I. MCLACHLAN, 'On symplectic and multisymplectic schemes for the KdV equation', *J. Sci. Comput.* 25 (2005) no. 1–2, 83–104.

**4.** T. J. BRIDGES and S. REICH, 'Multi-symplectic integrators: numerical schemes for Hamiltonian PDEs that conserve symplecticity', *Phys. Lett.* A 284 (2001) no. 4–5, 184–193.

**5.** T. J. BRIDGES and S. REICH, 'Numerical methods for Hamiltonian PDEs', *J. Phys.* A 39 (2006) no. 19, 5287–5320.

**6.** B. BUCHBERGER and M. KAUERS, 'Groebner basis', *Scholarpedia* 5 (2010) no. 10, 7763.

**7.** B. BUCHBERGER and M. KAUERS, 'Buchberger's algorithm', *Scholarpedia* 6 (2011) no. 10, 7764.

**8.** C. J. BUDD and M. D. PIGGOTT, 'Geometric integration and its applications', *Handbook of numerical analysis*, Handbook of Numerical Analysis XI (North-Holland, Amsterdam, 2003) 35–139.

**9.** D. COX, J. LITTLE and D. O'SHEA, 'An introduction to computational algebraic geometry and commutative algebra', *Ideals, varieties, and algorithms*, 3rd edn, Undergraduate Texts in Mathematics (Springer, New York, 2007).

**10.** P. G. DRAZIN and R. S. JOHNSON, *Solitons: an introduction*, Cambridge Texts in Applied Mathematics (Cambridge University Press, Cambridge, 1989).

**11.** S. V. DUZHIN and T. TSUJISHITA, 'Conservation laws of the BBM equation', *J. Phys.* A 17 (1984) no. 16, 3267–3276.

**12.** J. DE FRUTOS and J. M. SANZ-SERNA, 'Accuracy and conservation properties in numerical integration: the case of the Korteweg–de Vries equation', *Numer. Math.* 75 (1997) no. 4, 421–445.

**13.** D. FURIHATA, 'Finite difference schemes for $\partial u/\partial t = (\partial/\partial x)^\alpha \delta G/\delta u$ that inherit energy conservation or dissipation property', *J. Comput. Phys.* 156 (1999) no. 1, 181–205.

**14.** T. J. GRANT, 'Characteristics of conservation laws for finite difference equations', PhD Thesis, Department of Mathematics, University of Surrey, 2011.

**15.** T. J. GRANT and P. E. HYDON, 'Characteristics of conservation laws for difference equations', *Found. Comput. Math.* 13 (2013) no. 4, 667–692.

**16.** E. HAIRER, C. LUBICH and G. WANNER, 'Structure-preserving algorithms for ordinary differential equations', *Geometric numerical integration*, 2nd edn, Springer Series in Computational Mathematics 31 (Springer, Berlin, 2006).

**17.** P. E. HYDON, 'Conservation laws of partial difference equations with two independent variables', *J. Phys.* A 34 (2001) no. 48, 10347–10355.

**18.** P. E. HYDON and E. L. MANSFIELD, 'A variational complex for difference equations', *Found. Comput. Math.* 4 (2004) no. 2, 187–217.

**19.** S. KOIDE and D. FURIHATA, 'Nonlinear and linear conservative finite difference schemes for regularized long wave equation', *Jpn. J. Ind. Appl. Math.* 26 (2009) no. 1, 15–40.

**20.** B. A. KUPERSCHMIDT, 'Discrete Lax equations and differential-difference calculus', *Astérisque* (1985) no. 123, 212.

**21.** B. LEIMKUHLER and S. REICH, *Simulating Hamiltonian dynamics*, Cambridge Monographs on Applied and Computational Mathematics 14 (Cambridge University Press, Cambridge, 2004).
**22.** E. MANSFIELD, 'Differential Groebner bases', PhD Thesis, 1991.
**23.** R. I. MCLACHLAN, 'Spatial discretization of partial differential equations with integrals', *IMA J. Numer. Anal.* 23 (2003) no. 4, 645–664.
**24.** R. MCLACHLAN and R. QUISPEL, 'Six lectures on the geometric integration of ODEs', *Foundations of computational mathematics (Oxford, 1999)*, London Mathematical Society Lecture Note Series 284 (Cambridge University Press, Cambridge, 2001) 155–210.
**25.** E. NOETHER, 'Invariant variation problems', *Transport Theory Statist. Phys.* 1 (1971) no. 3, 186–207; Translated from the German (*Nachr. Akad. Wiss. Göttingen Math.-Phys. Kl.* II 1918, 235–257).
**26.** P. J. OLVER, 'Euler operators and conservation laws of the BBM equation', *Math. Proc. Cambridge Philos. Soc.* 85 (1979) no. 1, 143–160.
**27.** P. J. OLVER, *Applications of Lie groups to differential equations*, 2nd edn, Graduate Texts in Mathematics 107 (Springer, New York, 1993).
**28.** J. M. SANZ-SERNA, 'An explicit finite-difference scheme with exact conservation properties', *J. Comput. Phys.* 47 (1982) 199–210.
**29.** A. C. VLIEGENTHART, 'On finite-difference methods for the Korteweg–de Vries equation', *J. Engrg. Math.* 5 (1971) 137–155.
**30.** N. J. ZABUSKY and M. D. KRUSKAL, 'Interaction of "solitons" in a collisionless plasma and the recurrence of initial states', *Phys. Rev. Lett.* 15 (1965) no. 6, 240–243.
**31.** G. ZHONG and J. E. MARSDEN, 'Lie–Poisson Hamilton–Jacobi theory and Lie–Poisson integrators', *Phys. Lett.* A 133 (1988) no. 3, 134–139.

*Timothy J. Grant*
*Department of Mathematics*
*University of Surrey*
*Guildford GU2 7XH*
*UK*

timothy_grant@hotmail.co.uk