

Convergent numerical schemes for the compressible hyperelastic rod wave equation

David Cohen · Xavier Raynaud

Received: 5 May 2011 / Revised: 13 December 2011 / Published online: 8 March 2012
© Springer-Verlag 2012

Abstract We propose a fully discretised numerical scheme for the hyperelastic rod wave equation on the line. The convergence of the method is established. Moreover, the scheme can handle the blow-up of the derivative which naturally occurs for this equation. By using a time splitting integrator which preserves the invariants of the problem, we can also show that the scheme preserves the positivity of the energy density.

Mathematics Subject Classification (2000) 65M06 · 65M12 · 35B99 · 35Q53

1 Introduction

We consider the compressible hyperelastic rod wave equation

$$u_t - u_{xxt} + 3uu_x - \gamma(2u_x u_{xx} + uu_{xxx}) = 0, \quad u|_{t=0} = u_0. \quad (1)$$

This equation is obtained by Dai in [13] as a model equation for an infinitely long rod composed of a general compressible hyperelastic material. The author considers a far-field, finite length, finite amplitude approximation for a material where the first order dispersive terms vanish. The function $u = u(t, x)$ represents the radial stretch relative to a prestressed state. The parameter $\gamma \in \mathbb{R}$ is a constant which depends on the material and the prestress of the rod, whose physical values lie between -29.4760

D. Cohen (✉)
Mathematisches Institut, Universität Basel, 4051 Basel, Switzerland
e-mail: David.Cohen@unibas.ch

X. Raynaud
Center of Mathematics for Applications, University of Oslo, 0316 Oslo, Norway
e-mail: xavierra@cma.uio.no

and 3.4174. For materials where first order dispersive terms cannot be neglected, the KdV equation

$$u_t + uu_x + u_{xxx} = 0$$

applies and only smooth solitary waves exist. In contrast, the hyperelastic rod equation (1) admits sharp crested solitary waves.

The Cauchy problems of the hyperelastic rod wave equation on the line and on the circle are studied in [12] and [25], respectively. The stability of a class of solitary waves for the hyperelastic rod wave equation on the line is investigated in [12]. In [21], Lenells provides a classification of all traveling waves. In [12, 25], the authors establish, for a special class of initial data, the global existence in time of strong solutions. However, in the same papers, they also present conditions on the initial data for which the solutions blow up and, in that case, global classical solutions no longer exist. The way the solution blows up is known: In the case $\gamma > 0$, there is a point $x \in \mathbb{R}$ and a blow-up time T for which $\lim_{t \rightarrow T} u_x(t, x) = -\infty$ for (for $\gamma < 0$, we have $\lim_{t \rightarrow T} u_x(t, x) = \infty$).

To handle the blow-up, weak solutions have to be considered but they are no longer unique. For smooth solutions, the energy $\int_{\mathbb{R}} (u^2 + u_x^2) dx$ is preserved and $H^1(\mathbb{R})$ is a natural space for studying the solutions. After blow-up, there exist two consistent ways to prolong the solutions, which lead to *dissipative* and *conservative* solutions. In the first case, the energy which is concentrated at the blow-up point is dissipated while, in the second case, the same energy is restored. The global existence of dissipative solution is established in [6]. In the present article, we consider the conservative solutions, whose global existence is established in [18].

For $\gamma = 1$, the hyperelastic rod wave equation yields the Camassa–Holm equation

$$u_t - u_{xxt} + 3uu_x - 2u_x u_{xx} - uu_{xxx} = 0.$$

There is by now an important literature on the numerical discretisation of the Camassa–Holm equation. Let us review some of these works. In [4] and [5], a particle method was presented and convergence, as the number of particles tends to infinity, is proved for smooth solutions. In [17, 19], numerical schemes relying on a discretisation based on multipeakons are proved to converge for non-smooth solutions. Note that these schemes depend on a special type of solutions, the multipeakons, which do not exist for the hyperelastic rod equation. An adaptive finite volume method was derived in [1] for peakon-type solutions. Still related to the spatial discretisation of the Camassa–Holm equation, are the works [20] and [24], where a collocation method, respectively a local discontinuous Galerkin method, are presented and spatial convergence is proved for smooth solutions. Following a more geometrical approach, the first multi-symplectic schemes for the Camassa–Holm equation were presented in [10], however without proofs of convergence. A convergent finite difference scheme is studied in [16] for a special class of initial data whose properties are preserved by the equation. For non-smooth solutions, beside of [19], the only schemes with proof of convergence are given by [7] and [9], where finite difference schemes are used. But these schemes converge to the dissipative solutions of the Camassa–Holm equation.

In comparison with the Camassa–Holm equation, there are only a few works in the literature which are concerned with numerical methods for the hyperelastic rod wave equation. In [22], the authors consider a Galerkin approximation which preserves a discretisation of the energy. In [11], a Hamiltonian-preserving numerical method and a multisymplectic scheme are derived. In both works, no convergence proofs are provided and the schemes cannot handle the natural blow-up of the solution. The present paper fills these lack.

In this paper, we propose a fully discretised numerical scheme which can compute the solution on any finite time interval. In particular, it approximates solutions which have locally unbounded derivatives (the condition $u_x \in L^2(\mathbb{R})$ allows for an unbounded derivative in $L^\infty(\mathbb{R})$). A standard spatial discretisation of (1) cannot give us global solutions and the proofs of convergence for such schemes become highly nontrivial when the solution approaches blow-up time. The main achievement of this paper is the full convergence proof (both with respect to time and space) of the scheme. To compute the global solutions, we follow the framework given in [18]. With a coordinate transformation into Lagrangian coordinates, we first rewrite the problem as a system of ordinary differential equations in a Banach space (Sects. 2, 3). We establish new decay estimates (Sect. 4) which allow us to consider solutions defined on the whole real line. We discretise the system of equations in space (Sect. 5) and time (Sect. 7) and study the convergence of the numerical solution in Sect. 8. In Sect. 6, we explain how to define a converging sequence of initial data. This construction can be applied to any initial data in $H^1(\mathbb{R})$. Finally, in Sect. 9, numerical experiments demonstrate the validity of our theoretical results. Moreover, the time splitting discretisation enables the scheme to preserve invariants and we can use this property to prove that the scheme preserves the positivity of a discretisation of the energy density $u^2 + u_x^2 dx$, see Theorem 17. A Lagrangian formalism is also used in [5, 17, 19] to derive numerical schemes for the Camassa–Holm equation. However they rely on a particular class of solutions, the multipeakons, which is not available for the hyperelastic rod wave equation, that is when $\gamma \neq 1$.

The main difficulty in the numerical simulation of the hyperelastic rod wave equation is to find an appropriate spatial discretisation which can handle the discontinuities in the first derivative and the loss of regularity. This is reflected in the papers [5, 8, 16, 17, 20, 24] where the focus is clearly set on the spatial discretisation and the integration in time is done by standard numerical schemes. The spatial discretisation we propose is radically different from those proposed previously as it is based on a reformulation of the problem.

We want to emphasize that other schemes (based on more standard spatial discretisations as, for example finite differences) will inevitably experience difficulties when the solution becomes irregular and they will not be able to handle a peakon–antipeakon collision, as described in Sect. 9.3. This is illustrated in the peakon and cuspon traveling solutions presented in Figs. 3 and 5. These solutions are irregular in the sense that they have a discontinuous (in the case of the peakon) and unbounded (in the case of the cuspon) derivative at the peak. A direct spatial discretisation of the partial differential equation (1) naturally induces numerical dissipation. Indeed, when discretising in space, there is an upper bound on the frequencies that a finite discrete space step can represent and the frequencies above this bound are simply ignored.

For our scheme, by using a reformulation of the equation as an ordinary differential equation in a Banach space, the cut-off of high frequencies becomes harmless and we observe indeed very little numerical dissipation. We implemented the upwind scheme presented in [1]—without the mesh refinement—to compare our results. In [24], a local discontinuous Galerkin method is derived and analysed. The spatial discretisation is suited for solutions with discontinuities and the peakon solution are indeed well approximated. However, there is no proof of convergence and the scheme cannot handle collisions, as for example a peakon–antipeakon collision.

The results of this paper are also valid for the generalised hyperelastic rod wave equation

$$u_t - u_{xxt} + \frac{1}{2}g(u)_x - \gamma(2u_x u_{xx} + uu_{xxx}) = 0, \quad u|_{t=0} = u_0. \quad (2)$$

However, for simplicity only the numerical discretisation of Eq. (1) will be analysed. Equation (2) was first introduced in [6]; it defines a whole class of equations, depending on the choice of the (locally uniformly Lipschitz) function g and the value of the parameter γ , which contains several well-known nonlinear dispersive equations. Taking $\gamma = 1$ and $g(u) = 2\kappa u + 3u^2$ (with $\kappa \geq 0$), Eq. (2) reduces to the Camassa–Holm equation [3]; For $g(u) = 3u^2$, Eq. (2) becomes the hyperelastic rod wave equation (1); For $g(u) = 2u + u^2$ and for $\gamma = 0$, Eq. (2) leads to the Benjamin–Bona–Mahony (BBM) equation (or regularised long wave) [2].

2 The semigroup of conservative solutions

The purpose of this section is to recall the main results of [18] about the conservative solutions of the hyperelastic rod wave equation (1). The total energy for the hyperelastic rod wave equation is given by the H^1 norm, which is preserved in time for smooth solutions. An important feature of this equation is that it allows for the concentration of the energy density $(u^2 + u_x^2) dx$ on sets of zero measure. To construct a semigroup of conservative solution, it is necessary to keep track of the energy when it concentrates. This justifies the introduction of the set \mathcal{D} defined as follows.

Definition 1 The set \mathcal{D} is composed of all pairs (u, μ) such that u belongs to $H^1(\mathbb{R})$ and μ is a positive finite Radon measure whose absolute continuous part, μ_{ac} , satisfies

$$\mu_{ac} = (u^2 + u_x^2) dx.$$

The measure μ represents the energy density and the set \mathcal{D} allows μ to have a singular part. The solutions of (1) are constructed via a change of coordinates, from Eulerian to Lagrangian coordinates. An extra variable which account for the energy is necessary. Let us sketch this construction. We apply the inverse Helmholtz operator $(\text{Id} - \partial_{xx})^{-1}$ to (1) and obtain the system of equations

$$u_t + \gamma uu_x + P_x = 0 \quad (3a)$$

$$P - P_{xx} = \frac{3 - \gamma}{2}u^2 + \frac{\gamma}{2}u_x^2. \quad (3b)$$

By using the Green function of the Helmholtz operator, we can write P in an explicit form, i.e.,

$$P(t, x) = \frac{1}{2} \int_{\mathbb{R}} e^{-|x-z|} \left(\frac{3-\gamma}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) (t, z) dz. \tag{4}$$

We also define

$$Q(t, x) := P_x(t, x) = -\frac{1}{2} \int_{\mathbb{R}} \operatorname{sgn}(x-z) e^{-|x-z|} \left(\frac{3-\gamma}{2} u^2 + \frac{\gamma}{2} u_x^2 \right) (t, z) dz. \tag{5}$$

Next, we introduce the characteristics $y(t, \xi)$ defined as the solutions of

$$y_t(t, \xi) = \gamma u(t, y(t, \xi))$$

with $y(0, \xi)$ given. The variable $y(t, \xi)$ corresponds to the trajectory of a particle in the velocity field γu . However, the Lagrangian velocity will be defined as

$$U(t, \xi) = u(t, y(t, \xi))$$

and the cumulative energy $H(t, \xi)$ as

$$H(t, \xi) := \int_{-\infty}^{y(t, \xi)} (u^2 + u_x^2) dx.$$

We next express (4) and (5) in terms of the new variables $X = (y, U, H)$ (see [18] for the details) and we obtain

$$P(t, \xi) = \frac{1}{2} \int_{\mathbb{R}} e^{-\operatorname{sgn}(\xi-\eta)(y(\xi)-y(\eta))} \left(\frac{3-2\gamma}{2} U^2 y_\xi + \frac{\gamma}{2} H_\xi \right) (\eta) d\eta,$$

$$Q(t, \xi) = -\frac{1}{2} \int_{\mathbb{R}} \operatorname{sgn}(\xi-\eta) e^{-\operatorname{sgn}(\xi-\eta)(y(\xi)-y(\eta))} \left(\frac{3-2\gamma}{2} U^2 y_\xi + \frac{\gamma}{2} H_\xi \right) (\eta) d\eta.$$

Finally, we obtain the following system of differential equations

$$y_t = \gamma U \tag{6a}$$

$$U_t = -Q \tag{6b}$$

$$H_t = U^3 - 2PU, \tag{6c}$$

which we rewrite in the compact form

$$X_t = F(X).$$

The mapping F is a mapping from E to E , where E is a Banach space that we now define. We denote by V the space defined as

$$V = \{f \in C_b(\mathbb{R}) \mid f_\xi \in L^2(\mathbb{R})\},$$

where $C_b(\mathbb{R}) = C(\mathbb{R}) \cap L^\infty(\mathbb{R})$. The space V is a Banach space for the norm $\|f\|_V := \|f\|_{L^\infty} + \|f_\xi\|_{L^2}$. The Banach space E is then defined as

$$E = V \times H^1 \times V$$

with norm $\|f\|_E := \|f\|_V + \|f\|_{H^1} + \|f\|_V$. In [18], the existence of short-time solutions of (6) is established by a standard contraction argument in E . The solutions of (6) are not in general global in time but for initial data (ζ_0, U_0, H_0) which belongs to the set \mathcal{F} , which we now define, they are.

Definition 2 The set \mathcal{F} consists of all $(\zeta, U, H) \in E$ such that

$$(\zeta, U, H) \in [W^{1,\infty}(\mathbb{R})]^3 \quad \text{and} \quad \lim_{\xi \rightarrow -\infty} H(\xi) = 0 \quad (7a)$$

$$y_\xi \geq 0, H_\xi \geq 0, y_\xi + H_\xi \geq c \quad \text{almost everywhere, for some constant } c > 0 \quad (7b)$$

$$y_\xi H_\xi = y_\xi^2 U^2 + U_\xi^2 \quad \text{almost everywhere.} \quad (7c)$$

The set \mathcal{F} is preserved by the flow, that is, if $X(0) \in \mathcal{F}$ and $X(t)$ is the solution to (6) corresponding to this initial value, then $X(t) \in \mathcal{F}$ for all time t . The properties of the set \mathcal{F} can then be used to establish a priori estimates on the solutions and show that they exit globally in time, see [18] for more details. We denote by S_t the semigroup of solutions in \mathcal{F} given by the solutions of (6).

Given an initial data $(u, \mu) \in \mathcal{D}$, we have to find the corresponding initial data in \mathcal{F} ; we have to define a mapping between Eulerian and Lagrangian variables. To do so, we set

$$y(\xi) = \sup\{y \mid \mu((-\infty, y)) + y < \xi\}, \quad (8a)$$

$$H(\xi) = \xi - y(\xi), \quad (8b)$$

$$U(\xi) = u \circ y(\xi). \quad (8c)$$

We define $X = L(u, \mu)$ and L maps Eulerian to Lagrangian variables. When $\mu = \mu_{ac}$ (no energy is concentrated), Eq. (8a) simplifies and we get

$$y(\xi) + \int_{-\infty}^{y(\xi)} (u^2 + u_x^2)(x) dx = \xi.$$

Reciprocally, we define the mapping M from Lagrangian to Eulerian variables: Given $X = (y, U, H) \in \mathcal{F}$, we recover $(u, \mu) = M(X) \in \mathcal{D}$ by setting

$$u(x) = U(\xi) \quad \text{for any } \xi \text{ such that } x = y(\xi), \tag{9a}$$

$$\mu = y_{\#}(H_{\xi} d\xi). \tag{9b}$$

Here, $y_{\#}(H_{\xi} d\xi)$ denotes the push-forward of the measure $H_{\xi} d\xi$ by the mapping y .

Finally, we recall the following main result from [18].

Theorem 1 *The mapping $T : \mathcal{D} \times \mathbb{R}_+ \rightarrow \mathcal{D}$, where \mathcal{D} is defined by Definition 1, defines a continuous semigroup of conservative solutions of the hyperelastic rod wave equation (1), that is, given $(\bar{u}, \bar{\mu}) \in \mathcal{D}$, if we denote by $t \mapsto (u(t), \mu(t)) = T_t(\bar{u}, \bar{\mu})$ the corresponding trajectory, then u is a weak solution of the hyperelastic rod wave equation (3).*

The function $y(t, \xi)$ gives the trajectory of a particle which evolves in the velocity field given by $\gamma u(t, x)$. If u is smooth, then it is Lipschitz in the second variable and the mapping $\xi \rightarrow y(t, \xi)$ remains a diffeomorphism. We denote its inverse by $x \rightarrow y^{-1}(t, x)$. In this case, the density $\rho(t, x)$ is given by

$$\rho(t, x) = \frac{1}{y_{\xi}(t, y^{-1}(t, x))}. \tag{10}$$

We can also recover the energy density as

$$(u^2 + u_x^2)(t, x) = \frac{H_{\xi}}{y_{\xi}}(t, y^{-1}(t, x)). \tag{11}$$

In the following sections, we design numerical schemes which preserve the positivity of the particle and energy densities as defined in (10) and (11).

3 Equivalent system of ODEs in a Banach space

In this section, we reformulate the hyperelastic rod wave equation (3) as a system of ordinary differential equations in a Banach space as this was done in [18] but where we decouple the functions y, U and H and their derivatives y_{ξ}, U_{ξ} and H_{ξ} . Thus, after differentiating (6), we obtain

$$\zeta_{\xi t} = \gamma U_{\xi} \quad (\text{or } y_{\xi t} = \gamma U_{\xi}), \tag{12a}$$

$$U_{\xi t} = \frac{\gamma}{2} H_{\xi} + \left(\frac{3 - 2\gamma}{2} U^2 - P \right) y_{\xi}, \tag{12b}$$

$$H_{\xi t} = -2QU y_{\xi} + (3U^2 - 2P)U_{\xi}, \tag{12c}$$

where we set $\zeta(t, \xi) := y(t, \xi) - \xi$. The system (12) is quasilinear with respect to the first derivative $(\zeta_{\xi}, U_{\xi}, H_{\xi})$. This is an essential feature which leads to the stability of the solutions. Convergence proofs of numerical schemes rely generally on stability results and this is also the case here. It explains why the scheme we propose here is constructed upon the first derivatives. Let

$$q = y_\xi, \quad w = U_\xi, \quad h = H_\xi \quad \text{and} \quad v = q - 1 \quad (13)$$

then (6) and (12) rewrite

$$\zeta_t = y_t = \gamma U, \quad (14a)$$

$$U_t = -Q, \quad (14b)$$

$$H_t = U^3 - 2PU, \quad (14c)$$

$$v_t = q_t = \gamma w, \quad (14d)$$

$$w_t = \frac{\gamma}{2}h + \left(\frac{3-2\gamma}{2}U^2 - P \right)q, \quad (14e)$$

$$h_t = -2QUq + (3U^2 - 2P)w, \quad (14f)$$

where P and Q are given by

$$P = \frac{1}{2} \int_{\mathbb{R}} e^{-\operatorname{sgn}(\xi-\eta)(y(\xi)-y(\eta))} \left(\frac{3-2\gamma}{2}U^2q + \frac{\gamma}{2}h \right) (\eta) d\eta \quad (15)$$

and

$$Q = -\frac{1}{2} \int_{\mathbb{R}} \operatorname{sgn}(\xi - \eta) e^{-\operatorname{sgn}(\xi-\eta)(y(\xi)-y(\eta))} \left(\frac{3-2\gamma}{2}U^2q + \frac{\gamma}{2}h \right) (\eta) d\eta. \quad (16)$$

The main variables among the new variables (y, U, H, q, w, h) are the *derivative* variables q, w and h with respect to which the system is linear. The remaining variables, that is, y, U and H as well as the coefficients P and Q can be seen as integrals depending on q, w and h . Thus, the system (14d)–(14f) can be considered as a system of integro partial differential equations (where the integrals are space integrals and the derivatives are time derivative). However, for accuracy, we compute the variables y, U and H through their time evolution, as given by the first three equations in (14), instead of computing them as integrals. In this way we can also prove the convergence of the scheme, which is the main goal of this article.

Since the terms P and Q have similar structure, in the remaining of the paper most of the proofs will be established just for one of them. Now, we do not require (13) to hold any longer and, setting $Y := (\zeta, U, H, v, w, h)$, we obtain the system of differential equations

$$Y_t(t) = G(Y(t)),$$

where G is defined by (14). In the remaining, we will sometimes abuse the notation and write $Y = (y, U, H, q, w, h)$ instead of $Y = (\zeta, U, H, v, w, h)$. Then, we implicitly assume the relations $y(\xi) = \zeta(\xi) + \xi$ and $q = v + 1$. The variables y and q are the physical ones but do not have the proper decay/boundedness properties at infinity and this is why ζ and v have to be introduced. The system (14) is defined in the Banach space F , where F is given by

$$F := L^\infty(\mathbb{R}) \times (L^\infty(\mathbb{R}) \cap L^2(\mathbb{R})) \times L^\infty(\mathbb{R}) \times L^2(\mathbb{R}) \times L^2(\mathbb{R}) \times L^2(\mathbb{R}).$$

For any $Y = (\zeta, U, H, v, w, h) \in F$ we use the following norm on F :

$$\|Y\|_F = \|\zeta\|_{L^\infty} + \|U\|_{L^2} + \|U\|_{L^\infty} + \|H\|_{L^\infty} + \|v\|_{L^2} + \|w\|_{L^2} + \|h\|_{L^2}.$$

The following proposition holds.

Proposition 1 *The mappings $P : F \rightarrow H^1(\mathbb{R})$ and $Q : F \rightarrow H^1(\mathbb{R})$ belongs to $C^1(F, H^1(\mathbb{R}))$ and $G : F \rightarrow F$ belongs to $C^1(F, F)$. Moreover, given $M > 0$, let*

$$B_M = \{X \in F \mid \|X\|_F \leq M\}.$$

There exists a constant $C(M)$ which only depends on M such that

$$\|P(Y)\|_{H^1} + \|Q(Y)\|_{H^1} + \left\| \frac{\partial P}{\partial Y}(Y) \right\|_{L(F, H^1)} + \left\| \frac{\partial Q}{\partial Y}(Y) \right\|_{L(F, H^1)} \leq C(M) \tag{17}$$

and

$$\|G(Y)\|_F + \left\| \frac{\partial G}{\partial Y}(Y) \right\|_{L(F, F)} \leq C(M) \tag{18}$$

for all $Y \in B_M$.

Here, abusing slightly the notations, we denote by the same letter P the function $P(t, \xi)$ and the mapping $Y \mapsto P$. The same holds for Q . The norms $L(F, H^1(\mathbb{R}))$ and $L(F, F)$ are the operator norms.

Proof First we prove that the mappings $Y \mapsto P$ and $Y \mapsto Q$ as given by (15) and (16) belong to $C^1(F, L^\infty(\mathbb{R}) \cap L^2(\mathbb{R}))$. We rewrite Q as

$$\begin{aligned} Q(X)(\xi) &= -\frac{e^{-\zeta(\xi)}}{2} \int_{\mathbb{R}} \chi_{\{\eta < \xi\}}(\eta) e^{-(\xi-\eta)} e^{\zeta(\eta)} \\ &\quad \times \left(\frac{3-2\gamma}{2} U^2 q + \frac{\gamma}{2} h \right) (\eta) d\eta \\ &\quad + \frac{e^{\zeta(\xi)}}{2} \int_{\mathbb{R}} \chi_{\{\eta > \xi\}}(\eta) e^{(\xi-\eta)} e^{-\zeta(\eta)} \\ &\quad \times \left(\frac{3-2\gamma}{2} U^2 q + \frac{\gamma}{2} h \right) (\eta) d\eta, \end{aligned} \tag{19}$$

where χ_B denotes the indicator function of a given set B . We decompose Q into the sum $Q_1 + Q_2$, where Q_1 and Q_2 are the operators corresponding to the two terms in

the sum on the right-hand side of (19). Let $h(\xi) = \chi_{\{\xi>0\}}(\xi)e^{-\xi}$ and A be the map defined by $A: v \mapsto h \star v$. Then, Q_1 can be rewritten as

$$Q_1 = -\frac{e^{-\zeta(\xi)}}{2} A \circ R(Y)(\xi), \tag{20}$$

where R is the operator from F to $L^2(\mathbb{R})$ given by

$$R(Y)(\xi) = e^{\zeta(\xi)} \left(\frac{3-2\gamma}{2} U^2(1+v) + \frac{\gamma}{2} h \right) (\xi).$$

The mapping A is a continuous linear mapping from $L^2(\mathbb{R})$ into $L^2(\mathbb{R}) \cap L^\infty(\mathbb{R})$ as, from Young inequalities, we have

$$\|h \star v\|_{L^2} \leq \|h\|_{L^1} \|v\|_{L^2} \quad \text{and} \quad \|h \star v\|_{L^\infty} \leq \|h\|_{L^2} \|v\|_{L^2}. \tag{21}$$

For any $Y \in B_M$, we have

$$\|Q_1\|_{L^2 \cap L^\infty} \leq C(M) \|A \circ R\|_{L^2 \cap L^\infty} \leq C(M) \|R\|_{L^2} \leq C(M)$$

for some constant $C(M)$ which depends only on M . From now on, we denote generically by $C(M)$ such constant even if its value may change from line to line. The same result holds for Q and P . Since R is composed of sums and products of C^1 maps, the fact that $R: F \rightarrow L^2$ is C^1 follows directly from the following short lemma whose proof is essentially the same as the proof of the product rule for derivatives in \mathbb{R} .

Lemma 1 *Let $1 \leq p \leq \infty$. If $K_1 \in C^1(F, L^\infty(\mathbb{R}))$ and $K_2 \in C^1(F, L^p(\mathbb{R}))$, then the product $K_1 K_2$ belongs to $C^1(F, L^p(\mathbb{R}))$ and*

$$\frac{\partial(K_1 K_2)}{\partial Y}(Y)[\bar{Y}] = K_1(Y) \frac{\partial K_2}{\partial Y}(Y)[\bar{Y}] + K_2(Y) \frac{\partial K_1}{\partial Y}(Y)[\bar{Y}].$$

With this lemma in hands, we thus obtain that

$$\frac{\partial R}{\partial Y}(Y)[\bar{Y}] = e^\zeta \left(\frac{3-2\gamma}{2} (\bar{\zeta} U^2(1+v) + 2U\bar{U}(1+v) + U^2\bar{v}) + \frac{\bar{\zeta}\gamma}{2} h + \frac{\gamma}{2} \bar{h} \right)$$

and

$$\left\| \frac{\partial R}{\partial Y}(Y) \right\|_{L(F, L^2)} \leq C(M).$$

Then, Q_1 is in $C^1(F, L^2(\mathbb{R}) \cap L^\infty(\mathbb{R}))$,

$$\frac{\partial Q_1}{\partial Y}(Y)[\bar{Y}] = \frac{e^{-\zeta}}{2} \left(\bar{\zeta} A(R(Y)) - A \left(\frac{\partial R}{\partial Y}(Y)[\bar{Y}] \right) \right)$$

and

$$\left\| \frac{\partial Q_1}{\partial Y}(Y) \right\|_{L(F, L^2 \cap L^\infty)} \leq C(M).$$

We obtain the same result for Q_2 , Q and P . We differentiate Q and get

$$Q_\xi = \frac{\gamma}{2}h + \frac{3 - 2\gamma}{2}U^2q - Pq. \tag{22}$$

Hence, the mapping $Y \mapsto Q_\xi$ is differentiable,

$$\frac{\partial Q_\xi}{\partial Y}(Y)[\bar{Y}] = \frac{\gamma}{2}\bar{h} + \frac{3 - 2\gamma}{2}(2U\bar{U} + U^2\bar{q}) - \frac{\partial P}{\partial Y}(Y)[\bar{Y}]q - P\bar{q},$$

and

$$\left\| \frac{\partial Q_\xi}{\partial Y}(Y) \right\|_{L(F, L^2)} \leq C(M).$$

It follows that Q belongs to $C^1(F, H^1(\mathbb{R}))$ and $\left\| \frac{\partial Q}{\partial Y} \right\|_{L(F, H^1)} \leq C(M)$. The same result holds for P and (17) is proved. By using Lemma 1, we get that $G \in C^1(F, F)$ and this proves (18). \square

By using Proposition 1 and the standard contraction argument, we prove the existence of short-time solutions to (14):

Theorem 2 *For any initial values $Y_0 = (\zeta_0, U_0, H_0, v_0, w_0, h_0) \in F$, there exists a time T , only depending on the norm of the initial values, such that the system of differential equations (14) admits a unique solution in $C^1([0, T], F)$. Moreover, for any two solutions Y_1 and Y_2 such that $\sup_{t \in [0, T]} \|Y_1(t)\|_F \leq M$ and $\sup_{t \in [0, T]} \|Y_2(t)\|_F \leq M$, then*

$$\sup_{t \in [0, T]} \|Y_1(t) - Y_2(t)\|_F \leq C(M) \|Y_1(0) - Y_2(0)\|_F, \tag{23}$$

where the constant $C(M)$ depends only on M .

Proof The stability result (23) is a direct application of Proposition 1 and Gronwall’s lemma. \square

The system of differential equations (14) in the Banach space F has an interesting geometric property: it possesses an invariant. In fact, the following quantity

$$I(Y) := U^2q^2 + w^2 - qh$$

is conserved along the exact solution of the problem as we show now. For any $Y(t)$ solution of (14), we have

$$\begin{aligned} \frac{d}{dt} I(Y(t)) &= 2UU_t q^2 + 2U^2 q q_t + 2w w_t - q_t h - q h_t \\ &= -2U Q q^2 + 2U^2 q \gamma w + 2w \left(\frac{\gamma}{2} h + \left(\frac{3-2\gamma}{2} U^2 - P \right) q \right) \\ &\quad - \gamma w h - q(-2QUq + (3U^2 - 2P)) = 0. \end{aligned} \quad (24)$$

Additionally, we have

Lemma 2 *The following properties are preserved (independently one of each other) by the governing equations (14)*

- (i) q, w, h belongs to $L^\infty(\mathbb{R})$.
- (ii) $qh = U^2 q^2 + w^2$ (or $I(Y) = 0$).
- (iii) $qh = U^2 q^2 + w^2$ (or $I(Y) = 0$) and $q \geq 0, h \geq 0, q+h \geq c$ almost everywhere for some constant $c > 0$.
- (iv) The functions y, U and H are differentiable and $y_\xi = q, U_\xi = w$ and $H_\xi = h$.

The proof of Lemma 2 follows the lines of [18, Lemma 2.7]. Having a closer look at Lemma 2, we now define the following set.

Definition 3 The set \mathcal{G} consists of the elements $(y, U, H, q, w, h) \in F$ which satisfy the conditions (i), (iii) and (iv).

As a consequence of Lemma 2, the set \mathcal{G} is preserved by the system. For any initial data in \mathcal{G} , the solution of (14) coincide with the solutions that are obtained in [18]. In particular, we prove in the same way as in [18] that

Theorem 3 *For initial data in \mathcal{G} , the solutions to (14) are global in time.*

We denote by S_t the semigroup of solutions to (14) in \mathcal{G} . Note that global existence can only be established for initial data in \mathcal{G} and do not hold in general for initial data in F .

4 Decay at infinity

The terms P and Q , as given by (15) and (16) which appear in the governing equations (14), are global in the sense that they are not compactly supported even if Y is. Consequently the set of compactly supported functions is not preserved by the system. However, we identify in this section decay properties which are preserved by the system. These are new results which allow us to compute solutions on the full real line. In comparison, most numerical schemes for the Camassa–Holm equation [1, 10, 11, 16, 20, 22, 24] consider periodic solutions. In [7, 9], the case of the real line is considered but a grid of infinite length is used. In the present article, by using the decay estimate of this section, we prove the convergence of the scheme for a grid of finite length.

We denote by F^{exp} , the subspace of F of functions with exponential decay defined as

$$F^{\text{exp}} = \{Y \in F \mid q, w, h \in L^\infty(\mathbb{R}), \quad e^{|\xi|}U, e^{|\xi|}w \in L^2(\mathbb{R}), e^{|\xi|}h \in L^1(\mathbb{R})\}.$$

We define the following norm on F^{exp}

$$\begin{aligned} \|Y\|_{F^{\text{exp}}} &= \|Y\|_F + \|q\|_{L^\infty} + \|w\|_{L^\infty} + \|h\|_{L^\infty} \\ &\quad + \|e^{\frac{|\xi|}{2}}U\|_{L^2} + \|e^{\frac{|\xi|}{2}}w\|_{L^2} + \|e^{|\xi|}h\|_{L^1}. \end{aligned}$$

Given $\alpha > 1$, we denote by F^α , the subspace of F of functions with polynomial decay defined as

$$F^\alpha = \{Y \in F \mid q, w, h \in L^\infty(\mathbb{R}), \quad (1 + |\xi|)^{\frac{\alpha}{2}}U, (1 + |\xi|)^{\frac{\alpha}{2}}w \in L^2(\mathbb{R}), \\ (1 + |\xi|)^\alpha h \in L^1(\mathbb{R})\}.$$

We define the following norm on F^α

$$\begin{aligned} \|Y\|_{F^\alpha} &= \|Y\|_F + \|q\|_{L^\infty} + \|w\|_{L^\infty} + \|h\|_{L^\infty} \\ &\quad + \|(1 + |\xi|)^{\frac{\alpha}{2}}U\|_{L^2} + \|(1 + |\xi|)^{\frac{\alpha}{2}}w\|_{L^2} + \|(1 + |\xi|)^\alpha h\|_{L^1}. \end{aligned}$$

Theorem 4 *The spaces F^{exp} and F^α are preserved by the flow of (14). Considering the short-time solutions given by Theorem 2, we have that*

- (i) *If $Y_0 \in F^{\text{exp}}$, then $\sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^{\text{exp}}} \leq C$,*
- (ii) *If $Y_0 \in F^\alpha$, then $\sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^\alpha} \leq C$,*

for a constant C which only depends on T and $\|Y_0\|_{F^{\text{exp}}}$ [case (i)] or T and $\|Y_0\|_{F^\alpha}$ [case (ii)].

Proof Let us prove the case (i). First, we establish L^1 bounds on the solutions. By applying the Cauchy–Schwartz inequality, we get

$$\int_{\mathbb{R}} |U_0(\xi)| \, d\xi = \int_{\mathbb{R}} e^{-\left|\frac{\xi}{2}\right|} \left| e^{\left|\frac{\xi}{2}\right|} |U_0(\xi)| \right| \, d\xi \leq \sqrt{2} \left\| e^{|\xi|} U_0^2(\xi) \right\|_{L^1}^{\frac{1}{2}},$$

which implies that $U_0 \in L^1(\mathbb{R})$ and $\|U_0\|_{L^1} \leq C$ for some constant C which depends only on $\|e^{|\xi|}U_0^2(\xi)\|_{L^1}$. Similarly we get that $w_0 \in L^1(\mathbb{R})$ and $\|w_0\|_{L^1} \leq C$ for some constant C which depends only on $\|e^{|\xi|}w_0^2(\xi)\|_{L^1}$. We denote generically by C such a constant, which depends only on T and $\|Y_0\|_{F^{\text{exp}}}$. From Theorem 2 and Lemma 2, we get that

$$\|q(t, \cdot)\|_{L^\infty} + \|w(t, \cdot)\|_{L^\infty} + \|h(t, \cdot)\|_{L^\infty} \leq C.$$

By following the same argument as in the proof of Proposition 1, from (20) to (21), but, instead, using the Young inequality $\|\kappa \star r\|_{L^1} \leq \|\kappa\|_{L^1} \|r\|_{L^1}$, we obtain that

$$\|Q(t, \cdot)\|_{L^1} \leq C(\|h(t, \cdot)\|_{L^1} + 1) \tag{25}$$

for a constant C which depends only on $\|Y(t)\|_{F^{\text{exp}}}$ and, therefore, only on $\|Y_0\|_{F^{\text{exp}}}$ and T . The same estimate holds for P , that is,

$$\|P(t, \cdot)\|_{L^1} \leq C(\|h(t, \cdot)\|_{L^1} + 1). \tag{26}$$

Let us denote

$$J(t) := \|U(t, \cdot)\|_{L^1} + \|w(t, \cdot)\|_{L^1} + \|h(t, \cdot)\|_{L^1}.$$

From the governing equations (14), after using (25) and (26), we get

$$J(t) \leq J(0) + C + C \int_0^t J(\tau) d\tau.$$

Hence, by applying Gronwall’s lemma, we get that, for $t \in [0, T]$,

$$J(t) = \|U(t, \cdot)\|_{L^1} + \|w(t, \cdot)\|_{L^1} + \|h(t, \cdot)\|_{L^1} \leq C \tag{27}$$

for another constant C . Let $L(t)$ denotes

$$L(t) = \left\| e^{|\xi|} U^2(t, \cdot) \right\|_{L^1} + \left\| e^{|\xi|} w^2(t, \cdot) \right\|_{L^1} + \left\| e^{|\xi|} h(t, \cdot) \right\|_{L^1}. \tag{28}$$

From the definition of Q , we get that

$$Q(t, \xi) \leq C \int_{\mathbb{R}} e^{-|\xi-\eta|} (U^2 + h)(t, \eta) d\eta \tag{29}$$

so that

$$\begin{aligned} e^{|\xi|} Q(t, \xi) &\leq C \int_{\mathbb{R}} e^{|\xi|} e^{-|\xi-\eta|} e^{-|\eta|} e^{|\eta|} (U^2 + h)(t, \eta) d\eta \\ &\leq CL(t) \end{aligned}$$

because $|\xi| - |\eta| \leq |\xi - \eta|$ and therefore

$$\left\| e^{|\xi|} Q(t, \cdot) \right\|_{L^\infty} \leq CL(t). \tag{30}$$

Similarly, we get that

$$\left\| e^{|\xi|} P(t, \cdot) \right\|_{L^\infty} \leq CL(t). \tag{31}$$

From the governing equations (14), we get that

$$\begin{aligned} \left\| e^{|\xi|} U^2(t, \xi) \right\|_{L^1} &\leq \left\| e^{|\xi|} U_0^2 \right\|_{L^1} + \int_0^t \left\| 2e^{|\xi|} QU(\tau, \cdot) \right\|_{L^1} d\tau \\ &\leq \left\| e^{|\xi|} U_0^2 \right\|_{L^1} + 2 \int_0^t \left\| e^{|\xi|} Q(\tau, \cdot) \right\|_{L^\infty} \|U(\tau, \cdot)\|_{L^1} d\tau \\ &\leq \left\| e^{|\xi|} U_0^2 \right\|_{L^1} + C \int_0^t I(\tau) d\tau, \end{aligned} \tag{32}$$

by using the L^1 a priori estimates (27) and (30). From (14), we also obtain that

$$\begin{aligned} \left\| e^{|\xi|} h(t, \xi) \right\|_{L^1} &\leq \left\| e^{|\xi|} h_0 \right\|_{L^1} + \int_0^t \left(2 \left\| e^{|\xi|} Q(\tau, \cdot) \right\|_{L^\infty} \|U(\tau, \cdot)\|_{L^1} \right) d\tau \\ &\quad + \int_0^t \left(C \left\| e^{|\xi|} U^2(\tau, \cdot) \right\|_{L^1} + \left\| e^{|\xi|} P(\tau, \cdot) \right\|_{L^\infty} \|w(\tau, \cdot)\|_{L^1} \right) d\tau \end{aligned}$$

which, after using the L^1 estimates (27), (30) and (31), yields

$$\left\| e^{|\xi|} h(t, \xi) \right\|_{L^1} \leq \left\| e^{|\xi|} h_0 \right\|_{L^1} + C + C \int_0^t I(\tau) d\tau. \tag{33}$$

Similarly we get that

$$\begin{aligned} \left\| e^{|\xi|} w^2(t, \xi) \right\|_{L^1} &\leq \left\| e^{|\xi|} w_0^2 \right\|_{L^1} + \int_0^t \frac{\gamma}{2} \left\| e^{|\xi|} h(\tau, \xi) \right\|_{L^1} d\tau \\ &\quad + \int_0^t \left(C \left\| e^{|\xi|} U^2(\tau, \cdot) \right\|_{L^1} + C \left\| e^{|\xi|} P(\tau, \cdot) \right\|_{L^\infty} \|w(\tau, \cdot)\|_{L^1} \right) d\tau. \\ &\leq \left\| e^{|\xi|} w_0^2 \right\|_{L^1} + C + C \int_0^t I(\tau) d\tau. \end{aligned} \tag{34}$$

After summing (32), (33) and (34), we get $L(t) \leq L(0) + C + C \int_0^t L(\tau) d\tau$ and the result follows by applying Gronwall's inequality. We now turn to case (ii). We introduce the quantity

$$K(t) = \|(1 + |\xi|)^\alpha U^2(t, \cdot)\|_{L^1} + \|(1 + |\xi|)^\alpha w^2(t, \cdot)\|_{L^1} + \|(1 + |\xi|)^\alpha h(t, \cdot)\|_{L^1}.$$

From (29), we get

$$(1 + |\xi|)^\alpha Q \leq C \int_{\mathbb{R}} (1 + |\xi|)^\alpha e^{-|\xi - \eta|} (1 + |\eta|)^{-\alpha} (1 + |\eta|)^\alpha (U^2 q + h) d\eta. \quad (35)$$

Since $|\xi| \leq |\xi - \eta| + |\eta| \leq (1 + |\xi - \eta|)(1 + |\eta|)$, we have $(1 + |\xi|) \leq 2(1 + |\xi - \eta|)(1 + |\eta|)$ and

$$(1 + |\xi|)^\alpha \leq 2^\alpha (1 + |\xi - \eta|)^\alpha (1 + |\eta|)^\alpha. \quad (36)$$

Then, it follows from (35) that

$$\begin{aligned} (1 + |\xi|)^\alpha Q &\leq C \int_{\mathbb{R}} e^{-|\xi - \eta|} (1 + |\xi - \eta|)^\alpha (1 + |\eta|)^\alpha (U^2 q + h) d\eta \\ &\leq C \|e^{-z} (1 + |z|)^\alpha\|_{L^\infty} K(t) \leq CK(t) \end{aligned} \quad (37)$$

so that $\|(1 + |\xi|)^\alpha Q\|_{L^\infty} \leq CK(t)$. We have to estimate $\|(1 + |\xi|)^\alpha Q\|_{L^1}$. We have

$$\begin{aligned} \|(1 + |\xi|)^\alpha Q\|_{L^1} &\leq \int_{\mathbb{R}^2} (1 + |\xi|)^\alpha e^{-|\xi - \eta|} (1 + |\eta|)^{-\alpha} (1 + |\eta|)^\alpha (U^2 q + h) d\eta d\xi \\ &= \int_{\mathbb{R}^2} (1 + |\eta + z|)^\alpha e^{-|z|} (1 + |\eta|)^{-\alpha} (1 + |\eta|)^\alpha (U^2 q + h) d\eta dz \\ &\leq 2^\alpha \int_{\mathbb{R}^2} (1 + |z|)^\alpha e^{-|z|} (1 + |\eta|)^\alpha (U^2 q + h) d\eta dz \quad [\text{by (36)}] \\ &\leq C \int_{\mathbb{R}} (1 + |z|)^\alpha e^{-|z|} dz \int_{\mathbb{R}} (1 + |\eta|)^\alpha (U^2 + h) d\eta \\ &\leq CK(t). \end{aligned} \quad (38)$$

Hence,

$$\|(1 + |\xi|)^\alpha Q\|_{L^1 \cap L^\infty} \leq CK(t) \quad (39)$$

and the same bound holds for P . From the governing equations, we obtain

$$\begin{aligned} \left\| (1 + |\xi|)^\alpha U^2(t, \xi) \right\|_{L^1} &\leq \left\| (1 + |\xi|)^\alpha U_0^2 \right\|_{L^1} + \int_0^t \left\| 2(1 + |\xi|)^\alpha QU(\tau, \cdot) \right\|_{L^1} d\tau \\ &\leq \left\| (1 + |\xi|)^\alpha U_0^2 \right\|_{L^1} + 2 \int_0^t \left\| (1 + |\xi|)^\alpha Q^2(\tau, \cdot) \right\|_{L^1} d\tau \\ &\quad + 2 \int_0^t \left\| (1 + |\xi|)^\alpha U^2(\tau, \cdot) \right\|_{L^1} d\tau \\ &\leq \left\| (1 + |\xi|)^\alpha U_0^2 \right\|_{L^1} + C \int_0^t K(\tau) d\tau, \end{aligned}$$

by (39), as $\|Q\|_{L^\infty} \leq C$, see (17). In a similar way, one proves that

$$\left\| (1 + |\xi|)^\alpha w^2(t, \xi) \right\|_{L^1} \leq \left\| (1 + |\xi|)^\alpha w_0^2 \right\|_{L^1} + C \int_0^t K(\tau) d\tau$$

and

$$\left\| (1 + |\xi|)^\alpha h(t, \xi) \right\|_{L^1} \leq \left\| (1 + |\xi|)^\alpha h_0 \right\|_{L^1} + C \int_0^t K(\tau) d\tau$$

so that

$$K(t) \leq K(0) + C \int_0^t K(\tau) d\tau$$

and the result follows from Gronwall’s lemma. □

For later use, we note that, in this proof, we have established that

$$\left\| e^{|\xi|} Q \right\|_{L^\infty} + \left\| e^{|\xi|} P \right\|_{L^\infty} \leq C(\|Y\|_{F^{\text{exp}}}) \tag{40}$$

and

$$\left\| (1 + |\xi|)^\alpha Q \right\|_{L^\infty \cap L^1} + \left\| (1 + |\xi|)^\alpha P \right\|_{L^\infty \cap L^1} \leq C(\|Y\|_{F^\alpha}) \tag{41}$$

for some given increasing function C , see (30), (31) and (39).

5 Semi-discretisation in space

The first step towards a discretisation of (14) is to consider step-functions. We consider an equally-spaced grid on the real line defined by the points

$$\xi_i = i \Delta \xi,$$

where $\Delta \xi$ is the grid step and $i = 0, \pm 1, \pm 2, \dots$. We introduce the space

$$F_{\Delta \xi} = \{Y \in F : \text{each component of } Y \text{ consists of} \\ \text{piecewise constant functions in each intervals } [\xi_i, \xi_{i+1})\}.$$

The system (14) does not preserve the set $F_{\Delta \xi}$ of piecewise constant function. Thus, we define

$$P_{\Delta \xi}(Y)(\xi) = \sum_{i=-\infty}^{\infty} P(Y)(\xi_i) \chi_{[\xi_i, \xi_{i+1})}(\xi), \quad (42)$$

$$Q_{\Delta \xi}(Y)(\xi) = \sum_{i=-\infty}^{\infty} Q(Y)(\xi_i) \chi_{[\xi_i, \xi_{i+1})}(\xi) \quad (43)$$

and consider a second system of differential equations

$$\begin{aligned} \zeta_t &= \gamma U \\ U_t &= -Q_{\Delta \xi} \\ H_t &= U^3 - 2P_{\Delta \xi} U \\ q_t &= \gamma w \\ w_t &= \frac{\gamma}{2} h + \left(\frac{3-2\gamma}{2} U^2 - P_{\Delta \xi} \right) q \\ h_t &= -2Q_{\Delta \xi} U q + (3U^2 - 2P_{\Delta \xi}) w, \end{aligned} \quad (44)$$

or, shortly,

$$Y_t(t) = G_{\Delta \xi}(Y(t)).$$

Like in the preceding section, we show that this system of differential equations possesses a short-time solution, an invariant and that its solution converges to the solution of (14) as $\Delta \xi \rightarrow 0$. In the next theorem we prove, by a contraction argument, the short-time existence of solutions to (44).

Theorem 5 *For any initial value $Y_0 = (y_0, U_0, H_0, q_0, w_0, h_0) \in F$, there exists a time T , only depending on $\|Y_0\|_F$, such that the system of differential equations (44) admits a unique solution in $C^1([0, T], F)$.*

This theorem is a consequence of point (i) in the following lemma.

Lemma 3 *The following statements hold*

(i) *The mapping $G_{\Delta\xi} : F \rightarrow F$ belongs to $C^1(F, F)$ and*

$$\|G_{\Delta\xi}(Y)\|_F + \left\| \frac{\partial G_{\Delta\xi}}{\partial Y}(Y) \right\|_{L(F,F)} \leq C(M), \tag{45}$$

for any $Y \in B_M$.

(ii) *For any $Y \in F$, we have*

$$\|G(Y) - G_{\Delta\xi}(Y)\|_F \leq C\sqrt{\Delta\xi} \tag{46}$$

for some constant C which only depends on $\|Y\|_F$.

Proof For any function $f \in H^1(\mathbb{R})$, let $\mathbf{P}(f)$ be the function defined as $\mathbf{P}(f)(\xi) = \sum_{i=-\infty}^{\infty} f(\xi_i)\chi_{[\xi_i, \xi_{i+1})}(\xi)$. Thus, we can rewrite $Q_{\Delta\xi}(Y)$ and $P_{\Delta\xi}(Y)$ as

$$Q_{\Delta\xi}(Y) = \mathbf{P}[Q(Y)] \quad \text{and} \quad P_{\Delta\xi}(Y) = \mathbf{P}[P(Y)].$$

Let us prove that \mathbf{P} is a continuous mapping from $H^1(\mathbb{R})$ to $L^\infty(\mathbb{R}) \cap L^2(\mathbb{R})$. By using the Sobolev embedding theorem of $H^1(\mathbb{R})$ into $L^\infty(\mathbb{R})$, we get

$$\|\mathbf{P}(f)\|_{L^\infty} \leq \|f\|_{L^\infty} \leq C \|f\|_{H^1}$$

for some constant C , so that \mathbf{P} is continuous from $H^1(\mathbb{R})$ into $L^\infty(\mathbb{R})$. The L^2 norm of $\mathbf{P}(f)$ is given by

$$\|\mathbf{P}(f)\|_{L^2}^2 = \sum_{i=-\infty}^{\infty} \Delta\xi f(\xi_i)^2.$$

We have, for all $\xi \in [\xi_i, \xi_{i+1})$, that

$$\begin{aligned} f(\xi_i)^2 &= f(\xi)^2 - 2 \int_{\xi_i}^{\xi} f(\eta) f_\xi(\eta) d\eta \\ &\leq f(\xi)^2 + \int_{\xi_i}^{\xi_{i+1}} f^2(\eta) d\eta + \int_{\xi_i}^{\xi_{i+1}} f_\xi^2(\eta) d\eta \end{aligned}$$

which, after integration over $[\xi_i, \xi_{i+1})$, yields

$$\Delta\xi f(\xi_i)^2 \leq \int_{\xi_i}^{\xi_{i+1}} f(\eta)^2 d\eta + \Delta\xi \left(\int_{\xi_i}^{\xi_{i+1}} f^2(\eta) d\eta + \int_{\xi_i}^{\xi_{i+1}} f_\xi^2(\eta) d\eta \right).$$

Hence,

$$\|\mathbf{P}(f)\|_{L^2}^2 \leq (1 + \Delta\xi) \|f\|_{L^2}^2 + \Delta\xi \|f_\xi\|_{L^2}^2$$

and the mapping \mathbf{P} is continuous from $H^1(\mathbb{R})$ to $L^2(\mathbb{R})$. Since $Q_{\Delta\xi}$ and $P_{\Delta\xi}$ are compositions of a continuous linear map \mathbf{P} and a C^1 map, they are also C^1 and

$$\frac{\partial P_{\Delta\xi}}{\partial Y}(\bar{Y}) = \mathbf{P} \left(\frac{\partial P}{\partial Y}(Y)[\bar{Y}] \right)$$

for all $\bar{Y} \in F$. The same holds for Q so that (45) follows from Lemma 1. Let us prove point (ii). First we note that (46) follows directly from the definitions of $G, G_{\Delta\xi}$ and the estimate

$$\|Q(Y) - Q_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} + \|P(Y) - P_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} \leq C\sqrt{\Delta\xi}. \tag{47}$$

Let us prove (47). We estimate $\|\text{Id} - \mathbf{P}\|_{L(H^1, L^\infty \cap L^2)}$, where the norm here is the operator norm from $H^1(\mathbb{R})$ to $L^\infty(\mathbb{R}) \cap L^2(\mathbb{R})$. Let us consider $f \in H^1(\mathbb{R})$, we have

$$\|f - \mathbf{P}(f)\|_{L^\infty} \leq \sup_i \|f(\xi) - f(\xi_i)\|_{L^\infty([\xi_i, \xi_{i+1}])}.$$

For any $\xi \in [\xi_i, \xi_{i+1}]$, we have $|f(\xi) - f(\xi_i)| \leq \sqrt{\Delta\xi} \|f_\xi\|_{L^2}$, by the Cauchy–Schwartz inequality. Hence,

$$\|f - \mathbf{P}(f)\|_{L^\infty} \leq \sqrt{\Delta\xi} \|f_\xi\|_{L^2} \leq \sqrt{\Delta\xi} \|f\|_{H^1}.$$

We have

$$\begin{aligned} \int_{\xi_i}^{\xi_{i+1}} |f(\xi) - \mathbf{P}(f)(\xi)|^2 d\xi &= \int_{\xi_i}^{\xi_{i+1}} \left| \int_{\xi_i}^{\xi} f_\xi(\eta) d\eta \right|^2 d\xi \\ &\leq \int_{\xi_i}^{\xi_{i+1}} ((\xi - \xi_i) \int_{\xi_i}^{\xi} f_\xi^2(\eta) d\eta) d\xi \\ &\leq \int_{\xi_i}^{\xi_{i+1}} f_\xi^2(\eta) d\eta \int_{\xi_i}^{\xi_{i+1}} (\xi - \xi_i) d\xi \\ &= \frac{(\Delta\xi)^2}{2} \int_{\xi_i}^{\xi_{i+1}} f_\xi^2(\eta) d\eta. \end{aligned}$$

Hence,

$$\|f - \mathbf{P}(f)\|_{L^2} \leq \frac{\Delta\xi}{\sqrt{2}} \|f\|_{H^1} \tag{48}$$

and we have proved that $\|\text{Id} - \mathbf{P}\|_{L^2 \cap L^\infty} \leq C\sqrt{\Delta\xi}$ for some constant C . Then, we have

$$\|Q(Y) - Q_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} \leq C\sqrt{\Delta\xi} \|Q(Y)\|_{H^1} \leq C'\sqrt{\Delta\xi}$$

for another constant C' which depends only on $\|Y\|_F$. One proves in the same way the same estimate for P and thus we obtain (47). \square

Concerning our new system of equations (44), it is not difficult to show in the same way as in (24) that

$$I_{\Delta\xi}(Y) := U^2 q^2 + w^2 - qh$$

is also a conserved quantity along the exact solution of our problem. The system (44) is introduced because it allows for a spatial discretisation of the original system (14). Indeed, the set of piecewise constant functions is preserved:

Lemma 4 *The set $F_{\Delta\xi}$ is preserved, that is, if $Y_0 \in F_{\Delta\xi}$ and $Y(t)$ is the solution of (44) with initial data Y_0 , then $Y(t) \in F_{\Delta\xi}$ for all $t \in [0, T]$.*

The proof of this lemma is straightforward. We can now compare solutions of (44) and of the original system (14).

Theorem 6 *Given $M > 0$ and $Y_0, Y_{0,\Delta\xi} \in F$. Let $Y(t)$ be the short-time solution of (14) with initial data Y_0 and $Y_{\Delta\xi}(t)$ be the short-time solution of (44) with initial data $Y_{0,\Delta\xi}$ in the interval $[0, T]$. If we have*

$$\|Y(t)\|_F \leq M \quad \text{and} \quad \|Y_{\Delta\xi}(t)\|_F \leq M \quad \text{for all } t \in [0, T],$$

then we also have

$$\|Y(t) - Y_{\Delta\xi}(t)\|_F \leq \left(\|Y_0 - Y_{0,\Delta\xi}\| + CT\sqrt{\Delta\xi} \right) e^{CT} \quad \text{for all } t \in [0, T] \tag{49}$$

with some constant C which depends only on M .

Proof The proof of this theorem is a consequence of Lemma 3 and of Gronwall’s lemma. We have

$$\begin{aligned}
 Y(t) - Y_{\Delta\xi}(t) &= Y_0 - Y_{0,\Delta\xi} + \int_0^t \left(G(Y(\tau)) - G_{\Delta\xi}(Y_{\Delta\xi}(\tau)) \right) d\tau \\
 &= Y_0 - Y_{0,\Delta\xi} + \int_0^t \left(G(Y(\tau)) - G(Y_{\Delta\xi}(\tau)) + G(Y_{\Delta\xi}(\tau)) \right. \\
 &\quad \left. - G_{\Delta\xi}(Y_{\Delta\xi}(\tau)) \right) d\tau
 \end{aligned}$$

which yields, after using Proposition 1 and Lemma 3,

$$\|Y(t) - Y_{\Delta\xi}(t)\|_F \leq \|Y_0 - Y_{0,\Delta\xi}\|_F + C \int_0^t \|Y(\tau) - Y_{\Delta\xi}(\tau)\|_F d\tau + CT\sqrt{\Delta\xi},$$

for some constant C which depends only on M . Then, (49) follows from Gronwall’s lemma. □

Lemma 2 and Theorem 4 show that there exist properties of the initial data that are preserved by the system (14). The same results—with the exception of property (iv) in Lemma 2—hold for the system (44). This is the content of the following theorem.

Theorem 7 *We consider an initial datum $Y_0 \in F$ and the corresponding short time solution $Y(t)$ of (44) given by Theorem 5.*

(i) *If q_0, w_0, h_0 belongs to $L^\infty(\mathbb{R})$ then*

$$\sup_{t \in [0, T]} \left(\|q(t, \cdot)\|_{L^\infty} + \|w(t, \cdot)\|_{L^\infty} + \|h(t, \cdot)\|_{L^\infty} \right) \leq C$$

for some constant C which depends only on T and $\|Y_0\|_{F^{\text{exp}}}$.

- (ii) *If we have $qh = U^2q^2 + w^2$ for $t = 0$ (or $I(Y_0) = 0$) then this holds for all $t \in [0, T]$.*
- (iii) *If we have $qh = U^2q^2 + w^2$ (or $I(Y) = 0$) and $q \geq 0, h \geq 0, q + h \geq c$ almost everywhere for some constant $c > 0$, then the same relations holds for all $t \in [0, T]$.*
- (iv) *If $Y_0 \in F^{\text{exp}}$, then*

$$\sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^{\text{exp}}} \leq C, \tag{50}$$

if $Y_0 \in F^\alpha$, then

$$\sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^\alpha} \leq C, \tag{51}$$

where the constant C depends only on T and $\|Y_0\|_{F^{\text{exp}}}$, and T and $\|Y_0\|_{F^\alpha}$, respectively.

Proof The system (44) is obtained from (14) by simply replacing P and Q by $P_{\Delta\xi}$ and $Q_{\Delta\xi}$ as defined in (42) and (43). Therefore, the proofs of points (i), (ii) and (iii) in Lemma 2, which do not require any special properties of P and Q , apply directly to (44). After introspection of the proof of Theorem 4, we can see that in order to prove (50), we need to prove that the estimates (25), (26), (30), (31), which hold for P and Q , also hold for $P_{\Delta\xi}$ and $Q_{\Delta\xi}$, namely,

$$\|Q_{\Delta\xi}(t, \cdot)\|_{L^1} \leq C(\|h(t, \cdot)\|_{L^1} + 1), \quad \|P_{\Delta\xi}(t, \cdot)\|_{L^1} \leq C(\|h(t, \cdot)\|_{L^1} + 1) \tag{52}$$

and

$$\|e^{|\xi|} Q_{\Delta\xi}(t, \cdot)\|_{L^\infty} \leq CL(t), \quad \|e^{|\xi|} P_{\Delta\xi}(t, \cdot)\|_{L^\infty} \leq CL(t), \tag{53}$$

where $L(t)$ is defined in (28) and C is a constant which depends only on T and $\|Y_0\|_{F^{\text{exp}}}$. We denote generically by C such constant. In the same way that we obtained (48), we now get that, for any $f \in W^{1,1}(\mathbb{R})$,

$$\begin{aligned} \int_{\xi_i}^{\xi_{i+1}} |f(\xi) - \mathbf{P}(f)(\xi)| \, d\xi &= \int_{\xi_i}^{\xi_{i+1}} \left| \int_{\xi_i}^{\xi} f_\xi(\eta) \, d\eta \right| \, d\xi \\ &\leq \Delta\xi \int_{\xi_i}^{\xi_{i+1}} |f_\xi(\eta)| \, d\eta \end{aligned}$$

and therefore

$$\|f - \mathbf{P}(f)\|_{L^1} \leq \Delta\xi \|f_\xi\|_{L^1}. \tag{54}$$

We obtain, after using successively (54), (25), (22) and (26), that

$$\begin{aligned} \|Q_{\Delta\xi}\|_{L^1} &\leq \|Q_{\Delta\xi} - Q\|_{L^1} + \|Q\|_{L^1} \\ &\leq \Delta\xi \|Q_\xi\|_{L^1} + C(\|h\|_{L^1} + 1) \\ &= \Delta\xi \left\| \frac{\gamma}{2}h + \frac{3-2\gamma}{2}U^2q - Pq \right\|_{L^1} + C(\|h\|_{L^1} + 1) \\ &\leq C \|P\|_{L^1} + C(\|h\|_{L^1} + 1) \\ &\leq C(\|h\|_{L^1} + 1). \end{aligned}$$

We handle in the same way $\|P_{\Delta\xi}\|_{L^1}$ and this concludes the proof of (52). For any $\xi \in \mathbb{R}$, we have $\xi \in [\xi_i, \xi_{i+1})$ for some i . Then,

$$e^{|\xi|} Q_{\Delta\xi}(t, \xi) = e^{|\xi| - |\xi_i|} e^{|\xi_i|} Q(t, \xi_i) \leq e^{\Delta\xi} \|e^{|\xi|} Q(t, \xi)\|_{L^\infty} \leq CL(t)$$

by (30) and, therefore, $\|e^{|\xi|} Q_{\Delta\xi}(t, \cdot)\|_{L^\infty} \leq CL(t)$. Similarly, we obtain the corresponding result for $P_{\Delta\xi}$ so that (53) is proved. Again, after introspection of the proof of Theorem 4, we can check that, in order to prove (51), we need to prove that

$$\|(1 + |\xi|)^\alpha Q_{\Delta\xi}(t, \cdot)\|_{L^\infty \cap L^1} + \|(1 + |\xi|)^\alpha P_{\Delta\xi}(t, \cdot)\|_{L^\infty \cap L^1} \leq CK(t). \tag{55}$$

We have

$$\|(1 + |\xi|)^\alpha Q_{\Delta\xi}(t, \cdot)\|_{L^\infty} \leq \|(1 + |\xi|)^\alpha Q(t, \cdot)\|_{L^\infty} \leq CK(t)$$

by (37). Since $e^{\xi-\eta} \leq e^{\Delta\xi} e^{\xi_i-\eta}$ for any $(\xi, \eta) \in [\xi_i, \xi_{i+1}]^2$, we get

$$\begin{aligned} \|(1 + |\xi|)^\alpha Q_{\Delta\xi}(t, \cdot)\|_{L^1} &\leq \sum_{i=-\infty}^{\infty} \int_{\xi_i}^{\xi_{i+1}} \int_{\mathbb{R}} (1 + |\xi|)^\alpha e^{-|\xi_i-\eta|} (U^2q + h) \, d\eta \, d\xi \\ &\leq e^{\Delta\xi} \sum_{i=-\infty}^{\infty} \int_{\xi_i}^{\xi_{i+1}} \int_{\mathbb{R}} (1 + |\xi|)^\alpha e^{-|\xi-\eta|} (U^2q + h) \, d\eta \, d\xi \\ &= e^{\Delta\xi} \int_{\mathbb{R}} \int_{\mathbb{R}} (1 + |\xi|)^\alpha e^{-|\xi-\eta|} (U^2q + h) \, d\eta \, d\xi \leq CK(t), \end{aligned}$$

by (38). The corresponding results for P are established in the same way and this concludes the proof of (55). □

In order to complete the discretisation in space, we have to consider a finite subspace of $F_{\Delta\xi}$. Given any integer N , we denote $R = N\Delta\xi$ and we introduce the subset F_R of F defined as

$$\begin{aligned} F_R = \{Y \in F : \\ U(\xi) = q(\xi) = w(\xi) = h(\xi) = 0, & \text{ for all } \xi \in (-\infty, -R) \cup [R, \infty), \\ \zeta(\xi) = \zeta_\infty, H(\xi) = H_\infty, & \text{ for all } \xi \in [R, \infty), \\ \zeta(\xi) = \zeta_{-\infty}, H(\xi) = 0 & \text{ for all } \xi \in (-\infty, -R), \\ \text{where } \zeta_{\pm\infty} \text{ and } H_\infty \text{ are constants}\}. \end{aligned}$$

The set F_R basically corresponds to functions with compact support (U, q, w and h vanish outside a compact set). We do not require that the functions ζ and H have compact support (ζ and H belongs to L^∞ with no extra decay condition) but we impose that they are constant outside the compact interval $[-R, R]$. We denote $F_{\{\Delta\xi, R\}} = F_R \cap F_{\Delta\xi}$. The set $F_{\{\Delta\xi, R\}}$ is not preserved by the flow of (44) because, as mentioned earlier, P and Q do not preserve compactly supported functions. That is why we introduce the cut-off versions of P and Q given by

$$\begin{aligned}
 P_{\{\Delta\xi, R\}}(Y)(\xi) &= \sum_{i=-N}^{N-1} P(Y)(\xi_i)\chi_{[\xi_i, \xi_{i+1})}(\xi), \\
 Q_{\{\Delta\xi, R\}}(Y)(\xi) &= \sum_{i=-N}^{N-1} Q(Y)(\xi_i)\chi_{[\xi_i, \xi_{i+1})}(\xi)
 \end{aligned}$$

and define a third system of differential equations

$$\begin{aligned}
 \zeta_t &= \gamma U, \\
 U_t &= -Q_{\{\Delta\xi, R\}}, \\
 H_t &= U^3 - 2P_{\{\Delta\xi, R\}}U, \\
 q_t &= \gamma w, \\
 w_t &= \frac{\gamma}{2}h + \left(\frac{3 - 2\gamma}{2}U^2 - P_{\{\Delta\xi, R\}}\right)q, \\
 h_t &= -2Q_{\{\Delta\xi, R\}}Uq + (3U^2 - 2P_{\{\Delta\xi, R\}})w,
 \end{aligned} \tag{56}$$

or, shortly,

$$Y_t = G_{\{\Delta\xi, R\}}(Y).$$

It is clear from the definition that the system (56) preserves $F_{\{\Delta\xi, R\}}$ and therefore, since $F_{\{\Delta\xi, R\}}$ is of finite dimension, the system (56) is a spatial discretisation of (14) which allows for numerical computations. To emphasize that we are now working in finite dimension, we denote

$$Y_i(t) = Y_{\{\Delta\xi, R\}}(t, \xi_i),$$

$\zeta_i = \zeta_{\{\Delta\xi, R\}}(t, \xi_i)$, $U_i = U_{\{\Delta\xi, R\}}(t, \xi_i)$ and so on for H_i , q_i , w_i , h_i , P_i and Q_i for $i = \{-N, \dots, N - 1\}$. We have

$$Y_{\{\Delta\xi, R\}}(t, \xi) = \sum_{i=-N}^{N-1} Y_i(t)\chi_{[\xi_i, \xi_{i+1})}(\xi).$$

Again, we can show that

$$I_{\{\Delta\xi, R\}}^i(Y) := U_i^2 q_i^2 + w_i^2 - q_i h_i \tag{57}$$

are conserved quantities along the exact solution of problem (56). Finally, note that $F_{\{\Delta\xi, R\}}$ is contained in F^{exp} and F^α . Concerning the exact solution of (56), we have the following theorem.

Theorem 8 *For an initial values $Y_0 = (y_0, U_0, H_0, q_0, w_0, h_0) \in F$, there exists a time T , only depending on the norm of the initial values, such that the system of differential equations (56) admits a unique solution in $C^1([0, T], F)$.*

This theorem is a consequence of point (i) in the following lemma.

Lemma 5 *The following statements holds*

(i) *The mapping $G_{\{\Delta\xi, R\}} : F \rightarrow F$ belongs to $C^1(F, F)$ and*

$$\|G_{\{\Delta\xi, R\}}(Y)\|_F + \left\| \frac{\partial G_{\{\Delta\xi, R\}}}{\partial Y}(Y) \right\|_{L(F, F)} \leq C(M), \tag{58}$$

for any $Y \in B_M$.

(ii) *For any $Y \in F^{\text{exp}}$, we have*

$$\|G_{\{\Delta\xi, R\}}(Y) - G_{\Delta\xi}(Y)\|_F \leq Ce^{-R}, \tag{59}$$

for some constant C which only depends on $\|Y\|_{F^{\text{exp}}}$.

(iii) *For any $Y \in F^\alpha$, we have*

$$\|G_{\{\Delta\xi, R\}}(Y) - G_{\Delta\xi}(Y)\|_F \leq C \left(\sqrt{\Delta\xi} + \frac{1}{R^{\alpha/2}} \right), \tag{60}$$

for some constant C which only depends on $\|Y\|_{F^\alpha}$.

Note that for $Y(t)$ solution of (56), we have

$$\sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^{\text{exp}}} \leq C \quad \text{and} \quad \sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^\alpha} \leq C,$$

where C depends on $\|Y_0\|_{F^{\text{exp}}}$ and $\|Y_0\|_{F^\alpha}$, respectively. This follows from (50), (51), (59) and (60).

Proof of Lemma 5 For any function $f \in L^\infty(\mathbb{R}) \cap L^2(\mathbb{R})$, let $\mathbf{P}_R(f)$ be the function defined as $\mathbf{P}_R(f)(\xi) = f(\xi)\chi_{[-R, R]}$. Thus, we can rewrite $Q_{\{\Delta\xi, R\}}(Y)$ and $P_{\{\Delta\xi, R\}}(Y)$ as

$$Q_{\{\Delta\xi, R\}}(Y) = \mathbf{P}_R[Q_{\Delta\xi}(Y)] \quad \text{and} \quad P_{\{\Delta\xi, R\}}(Y) = \mathbf{P}_R[P_{\Delta\xi}(Y)].$$

The operator \mathbf{P}_R is a projection from $L^\infty(\mathbb{R}) \cap L^2(\mathbb{R})$ into itself and therefore its norm is smaller than one. Hence, (58) follows from (45). Let us prove (ii). We consider $Y \in F^{\text{exp}}$. We have to prove

$$\|Q_{\{\Delta\xi, R\}}(Y) - Q_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} + \|P_{\{\Delta\xi, R\}}(Y) - P_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} \leq Ce^{-R}. \tag{61}$$

By (40), we have $\|e^{|\xi|}Q\|_{L^\infty} + \|e^{|\xi|}P\|_{L^\infty} \leq C$. Hence,

$$\|Q_{\{\Delta\xi, R\}} - Q_{\Delta\xi}\|_{L^\infty} = \sup_{|\xi_i| \geq R} |Q(\xi_i)| \leq C \sup_{|\xi_i| \geq R} e^{-|\xi_i|} = Ce^{-R}.$$

We have

$$\|Q_{\{\Delta\xi, R\}} - Q_{\Delta\xi}\|_{L^2}^2 = \Delta\xi \sum_{|\xi_i| \geq R} Q(\xi_i)^2 \leq C \Delta\xi \sum_{|i \Delta\xi| \geq R} e^{-2|i \Delta\xi|} \leq C \frac{2\Delta\xi}{1 - e^{-2\Delta\xi}} e^{-2R}$$

and therefore $\|Q_{\{\Delta\xi, R\}} - Q_{\Delta\xi}\|_{L^2} \leq Ce^{-R}$. We prove in the same way the corresponding result for P and it concludes the proof of (61). The estimate (59) follows from (61). Let us prove (iii). We consider $Y \in F^\alpha$. We have to prove that

$$\begin{aligned} & \|Q_{\{\Delta\xi, R\}}(Y) - Q_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} \\ & + \|P_{\{\Delta\xi, R\}}(Y) - P_{\Delta\xi}(Y)\|_{L^2 \cap L^\infty} \leq C \left(\sqrt{\Delta\xi} + \frac{1}{R^{\alpha/2}} \right). \end{aligned} \tag{62}$$

By (41), we have $\|(1 + |\xi|)^\alpha Q\|_{L^\infty} + \|(1 + |\xi|)^\alpha P\|_{L^\infty} \leq C$. Hence,

$$\|Q_{\{\Delta\xi, R\}} - Q_{\Delta\xi}\|_{L^\infty} = \sup_{|\xi_i| \geq R} |Q(\xi_i)| \leq C \sup_{|\xi_i| \geq R} (1 + |\xi_i|)^{-\alpha} = C(1 + R)^{-\alpha}. \tag{63}$$

We have

$$\begin{aligned} \|Q_{\{\Delta\xi, R\}} - Q_{\Delta\xi}\|_{L^2(\mathbb{R})} &= \|Q_{\Delta\xi}\|_{L^2(\mathbb{R} \setminus [-R, R])} \\ &\leq \|Q_{\Delta\xi} - Q\|_{L^2(\mathbb{R} \setminus [-R, R])} + \|Q\|_{L^2(\mathbb{R} \setminus [-R, R])} \\ &\leq C \left(\sqrt{\Delta\xi} + \|Q\|_{L^2(\mathbb{R} \setminus [-R, R])} \right), \end{aligned} \tag{64}$$

from (47). Since

$$\begin{aligned} \|Q\|_{L^2(\mathbb{R} \setminus [-R, R])}^2 &\leq (1 + |R|)^{-\alpha} \int_{\mathbb{R} \setminus [-R, R]} (1 + |\xi|)^\alpha Q^2 d\xi \\ &\leq C(1 + R)^{-\alpha}, \quad \text{by (41),} \end{aligned}$$

the estimate (62) follows from (63) and (64). □

Again, the system (56) preserves properties of the initial data:

Theorem 9 *We consider an initial datum $Y_0 \in F$ and the corresponding short time solution $Y(t)$ of (56) given by Theorem 8. Then, $Y(t)$ satisfy points (i)–(iv) as given in Theorem 7.*

Finally, for any initial datum in $Y_0 \in F^{\text{exp}}$, resp. $Y_0 \in F^\alpha$, we obtain the following error estimate for bounded solutions.

Theorem 10 *Given Y_0 and $Y_{0, \Delta\xi, R}$ in F^{exp} , let $Y(t)$ and $Y_{\{\Delta\xi, R\}}(t)$ be the short-time solutions of (14) and (56), respectively, with initial datum Y_0 and $Y_{0, \Delta\xi, R}$, respectively. If we have*

$$\|Y(t)\|_{F^{\text{exp}}} \leq M \quad \text{and} \quad \|Y_{\{\Delta\xi, R\}}(t)\|_F \leq M \quad \text{for all } t \in [0, T],$$

then we have

$$\sup_{t \in [0, T]} \|Y(t, \cdot) - Y_{\{\Delta\xi, R\}}(t, \cdot)\|_F \leq C \left(\|Y_0 - Y_{0, \Delta\xi, R}\|_F + \sqrt{\Delta\xi} + e^{-R} \right), \tag{65}$$

where the constant C depends only on M . For Y_0 and $Y_{0,\Delta\xi,R}$ in F^α , we have that if

$$\|Y(t)\|_{F^\alpha} \leq M \quad \text{and} \quad \|Y_{\{\Delta\xi,R\}}(t)\|_F \leq M \quad \text{for all } t \in [0, T],$$

then

$$\sup_{t \in [0,T]} \|Y(t, \cdot) - Y_{\{\Delta\xi,R\}}(t, \cdot)\|_F \leq C \left(\|Y_0 - Y_{0,\Delta\xi,R}\|_F + \sqrt{\Delta\xi} + \frac{1}{R^{\alpha/2}} \right). \quad (66)$$

Proof We have

$$\begin{aligned} \|Y(t, \cdot) - Y_{\{\Delta\xi,R\}}(t, \cdot)\|_F &\leq \|Y_0 - Y_{0,\Delta\xi,R}\|_F \\ &+ \int_0^t \|G(Y(\tau, \cdot)) - G_{\{\Delta\xi,R\}}(Y_{\{\Delta\xi,R\}}(\tau, \cdot))\|_F \, d\tau. \end{aligned} \quad (67)$$

By Proposition 1 and Lemmas 3 and 5, we get

$$\begin{aligned} &\|G(Y(\tau, \cdot)) - G_{\{\Delta\xi,R\}}(Y_{\{\Delta\xi,R\}}(\tau, \cdot))\|_F \\ &\leq \|G(Y(\tau, \cdot)) - G_{\Delta\xi}(Y(\tau, \cdot))\|_F \\ &\quad + \|G_{\Delta\xi}(Y(\tau, \cdot)) - G_{\{\Delta\xi,R\}}(Y(\tau, \cdot))\|_F \\ &\quad + \|G_{\{\Delta\xi,R\}}(Y(\tau, \cdot)) - G_{\{\Delta\xi,R\}}(Y_{\{\Delta\xi,R\}}(\tau, \cdot))\|_F \\ &\leq C \left((\Delta\xi)^{\frac{1}{2}} + e^{-R} + \|Y(\tau, \cdot) - Y_{\{\Delta\xi,R\}}(\tau, \cdot)\|_F \right) \end{aligned}$$

for a constant C which depends only on M . Hence, (65) follows from (67) after applying Gronwall’s lemma. The proof of (66) is similar. \square

6 Approximation of the initial data and convergence of the semi-discrete solutions

6.1 Approximation of the initial data

The construction of the initial data $Y_{0,\Delta\xi,R}$ for (56) is done in two steps. First, we change variable from Eulerian to Lagrangian, that is, we compute $Y_0 \in \mathcal{G}$ such that $X = (y_0, U_0, H_0) \in \mathcal{F}$ satisfies

$$U_0 = u_0 \circ y_0. \quad (68)$$

In the new set of variables, we can solve (14) or, rather, its discretisation (56). Note that, given $u_0 \in H^1(\mathbb{R})$, there exists several $Y_0 \in \mathcal{G}$ such that (68) holds (this is a consequence of relabeling invariance, see [18] and this fact will be used in the numerical examples of Sect. 9). Here, we present a framework valid for general initial data in $H^1(\mathbb{R})$. In Sect. 2, we defined the mapping L from \mathcal{D} to \mathcal{F} . For $(u_0, \mu_0) \in \mathcal{D}$, i.e., for

$u_0 \in H^1(\mathbb{R})$ and $\mu_0 = (u_0^2 + u_{0,x}^2) dx$ absolutely continuous, this mapping simplifies and reads

$$y_0(\xi) + \int_{-\infty}^{y_0(\xi)} (u_0^2 + u_{0,x}^2) dx = \xi, \tag{69a}$$

$$U_0 = u_0 \circ y \quad \text{and} \quad H_0 = \text{Id} - y_0. \tag{69b}$$

Then, we set

$$q_0 = y_{0,\xi}, \quad w = U_{0,\xi}, \quad h = H_{0,\xi}. \tag{69c}$$

As earlier, we denote $v_0 = 1 - q_0$ and $\zeta_0 = \text{Id} - y_0$. We have

$$h_0 q_0 = q_0^2 U_0^2 + w_0^2, \quad q_0 + h_0 = 1, \quad q_0 > 0, \quad h_0 \geq 0 \quad \text{for almost every } \xi \in \mathbb{R}. \tag{70}$$

The element $Y_0 = (y_0, U_0, H_0, q_0, w_0, h_0)$ belongs to \mathcal{G} . The second step consists of computing an approximation of Y_0 in $F_{\{\Delta\xi, R\}}$. In the following theorem, we show how the change of variable given by (69) deal with the decay conditions. For simplicity, we drop the subscript zero in the notation. Let us introduce the Banach spaces $H^{1,\text{exp}}$ and $H^{1,\alpha}$ as the subspaces of H^1 with respective norms

$$\|u\|_{H^{1,\text{exp}}}^2 = \left\| e^{|\cdot|^{\frac{\alpha}{2}}} u \right\|_{L^2}^2 + \left\| e^{|\cdot|^{\frac{\alpha}{2}}} u_x \right\|_{L^2}^2$$

and

$$\|u\|_{H^{1,\alpha}}^2 = \left\| (1 + |\xi|)^{\frac{\alpha}{2}} u \right\|_{L^2}^2 + \left\| (1 + |\xi|)^{\frac{\alpha}{2}} u_x \right\|_{L^2}^2.$$

Theorem 11 *Given u and Y as given by (69), we have*

- (i) $u \in H^{1,\text{exp}}$ if and only if $Y \in F^{\text{exp}}$,
- (ii) $u \in H^{1,\alpha}$ if and only if $Y \in F^\alpha$.

Proof Let us assume that $u \in H^{1,\text{exp}}$. By definition, we have $h = (u^2 + u_x^2) \circ y y_\xi$. Hence,

$$\begin{aligned} \int_{\mathbb{R}} e^{|\xi|} h(\xi) d\xi &= \int_{\mathbb{R}} e^{|\xi|} (u^2 + u_x^2) \circ y(\xi) y_\xi(\xi) d\xi \\ &= \int_{\mathbb{R}} e^{|y^{-1}(x)|} (u^2 + u_x^2)(x) dx \\ &= \int_{\mathbb{R}} e^{|y^{-1}(x)-x|} e^{|x|} (u^2 + u_x^2)(x) dx \\ &\leq e^{\|y(\xi)-\xi\|_{L^\infty}} \int_{\mathbb{R}} e^{|x|} (u^2 + u_x^2)(x) dx < \infty. \end{aligned}$$

Using (70), we get

$$\int_{\mathbb{R}} e^{|\xi|} w^2(\xi) d\xi \leq \|q\|_{L^\infty} \int_{\mathbb{R}} e^{|\xi|} h(\xi) d\xi < \infty.$$

In order to prove that $\int_{\mathbb{R}} e^{|\xi|} U^2 d\xi$ is finite, we decompose the integral as follows:

$$\int_{\mathbb{R}} e^{|\xi|} U^2 d\xi = \int_{\{\xi \in \mathbb{R} | q < \frac{1}{2}\}} e^{|\xi|} U^2 d\xi + \int_{\{\xi \in \mathbb{R} | q > \frac{1}{2}\}} e^{|\xi|} U^2 d\xi.$$

We have

$$\begin{aligned} \int_{\{\xi \in \mathbb{R} | q < \frac{1}{2}\}} e^{|\xi|} U^2 d\xi &\leq \|U\|_{L^\infty}^2 \int_{\{\xi \in \mathbb{R} | q < \frac{1}{2}\}} e^{|\xi|} d\xi \\ &\leq \|U\|_{L^\infty}^2 \int_{\{\xi \in \mathbb{R} | h > \frac{1}{2}\}} e^{|\xi|} d\xi, \quad \text{as } q + h = 1, \\ &\leq 2 \|U\|_{L^\infty}^2 \int_{\{\xi \in \mathbb{R} | h > \frac{1}{2}\}} h e^{|\xi|} d\xi \leq C \int_{\mathbb{R}} e^{|\xi|} h d\xi < \infty \end{aligned}$$

and

$$\begin{aligned} \int_{\{\xi \in \mathbb{R} | q > \frac{1}{2}\}} e^{|\xi|} U^2 d\xi &\leq 2 \int_{\{\xi \in \mathbb{R} | q > \frac{1}{2}\}} e^{|\xi|} \frac{U^2}{q} d\xi \\ &\leq 2 \int_{\{\xi \in \mathbb{R} | q > \frac{1}{2}\}} e^{|\xi|} q h d\xi, \quad \text{as } U^2 \leq qh \text{ by (70),} \\ &< \infty. \end{aligned}$$

Hence, $\int_{\mathbb{R}} e^{|\xi|} U^2 d\xi < \infty$. Let us now assume that $Y \in F^{\text{exp}}$. Then,

$$\begin{aligned} \int_{\mathbb{R}} e^{|\xi|} (u^2 + u_x^2)(x) dx &= \int_{\mathbb{R}} e^{|\gamma(\xi)|} (u^2 + u_x^2)(\gamma(\xi)) \gamma_\xi(\xi) d\xi \\ &= \int_{\mathbb{R}} e^{|\gamma(\xi)|} h(\xi) d\xi \end{aligned}$$

$$\begin{aligned} &\leq \int_{\mathbb{R}} e^{|\gamma(\xi)-\xi|} e^{|\xi|} h(\xi) d\xi \\ &\leq e^{(\|\gamma(\xi)-\xi\|_{L^\infty})} \int_{\mathbb{R}} e^{|\xi|} h(\xi) d\xi < \infty \end{aligned}$$

and $u_0 \in H^{1,\text{exp}}$. The case (ii) is proved in the same way. □

As a consequence of this theorem and Theorem 4, we obtain

Theorem 12 *The spaces $H^{1,\text{exp}}$ and $H^{1,\alpha}$ are preserved by the hyperelastic rod equation: If $u_0 \in H^{1,\text{exp}}$, then $u(t, \cdot) \in H^{1,\text{exp}}$ for all positive time and, similarly, if $u_0 \in H^{1,\alpha}$, then $u(t, \cdot) \in H^{1,\alpha}$ for all positive time.*

To the best of our knowledge, these decay results are new, even for the Camassa–Holm equation (case $\gamma = 1$). They have to be compared with [15] where it is established that the only solution which has compact support for all positive time is the zero solution, i.e., the compactness of the support (which is a kind of decay condition) is *not* preserved by the equation.

Let us now construct the approximating sequence for the initial data. From (70), we get that

$$0 \leq q \leq 1, \quad 0 \leq h \leq 1$$

and

$$U_\xi = w \leq \sqrt{hq} \leq \frac{1}{2}(h + q) = \frac{1}{2}. \tag{71}$$

Given an integer n , we consider $\Delta\xi$ and R such that $\frac{1}{n} = \frac{1}{R} + \Delta\xi = \frac{1}{R} + \frac{R}{N}$ so that $n \rightarrow \infty$ if and only if $\Delta\xi \rightarrow 0$ and $R \rightarrow \infty$. We introduce the mapping $\mathbf{I}_{\Delta\xi} : L^2 \rightarrow L^2$ which approximates L^2 functions by piecewise constant functions, that is, given $f \in L^2$, let

$$\bar{f}_i = \frac{1}{\Delta\xi} \int_{\xi_i}^{\xi_{i+1}} f(\xi) d\xi$$

and set

$$\mathbf{I}_{\Delta\xi}(f)(\xi) = \sum_{i=-\infty}^{\infty} \bar{f}_i \cdot \chi_{[\xi_i, \xi_{i+1})}(\xi).$$

We define $Y_n = (y_n, U_n, H_n, q_n, w_n, h_n)$ as follows. Let

$$v_n(\xi) = \mathbf{P}_R \mathbf{I}_{\Delta\xi}(v), \quad w_n(\xi) = \mathbf{P}_R \mathbf{I}_{\Delta\xi}(w), \quad h_n(\xi) = \mathbf{P}_R \mathbf{I}_{\Delta\xi}(h).$$

As usual, we denote $q = 1 + v$ and $q_n = 1 + v_n$. Moreover, let us define the weighted integrals

$$U_{i,n} = \frac{\int_{\xi_i}^{\xi_{i+1}} q_n^2 U_n d\xi}{\int_{\xi_i}^{\xi_{i+1}} q_n^2 d\xi}.$$

We set

$$U_n(\xi) = \sum_{i=-N}^{N-1} U_{i,n} \cdot \chi_{[\xi_i, \xi_{i+1})}(\xi), \quad \text{for } i = -N, \dots, N - 1.$$

We define

$$H_n(\xi) = \mathbf{P} \left(\int_{-\infty}^{\xi} h_n d\eta \right) \quad \text{if } \xi \in [-R, R]$$

and $H_n(\xi) = \int_{-\infty}^{-R} h_n d\eta$ if $\xi \in (-\infty, -R)$, $H_n(\xi) = \int_R^{\infty} h_n d\eta$ if $\xi \in (R, \infty)$. For y_n , we set

$$y_n(\xi) = \xi - H_n(\xi) \quad \text{if } \xi \in [-R, R]$$

and $y_n(\xi) = \xi - H_n(-R)$ if $\xi \in (-\infty, -R)$, $y_n(\xi) = \xi - H_n(R)$ if $\xi \in (R, \infty)$. The definition of \mathbf{P} is given in the proof of Lemma 3. The following theorem states that Y_n approximates Y in $F_{\{\Delta\xi, R\}}$ and satisfies additional properties which will be useful in Theorem 17, where we prove that the positivity of the energy is preserved by the numerical scheme.

Theorem 13 *Given $Y \in \mathcal{G}$, there exist a sequence $Y_n \in F_{\{\Delta\xi, R\}}$ such that*

$$\lim_{n \rightarrow \infty} \|Y_n - Y\|_F = 0, \tag{72a}$$

and

$$q_n h_n \geq U_n^2 q_n^2 + w_n^2, \quad q_n + h_n = 1, \quad \text{for all } n \geq 0 \text{ and for all } \xi. \tag{72b}$$

Moreover, we have

$$\|Y_n\|_{F^{\text{exp}}} \leq C \|Y\|_{F^{\text{exp}}} \quad \text{and} \quad \|Y_n\|_{F^\alpha} \leq C \|Y\|_{F^\alpha} \tag{72c}$$

for $Y \in F^{\text{exp}}$, resp. $Y \in F^\alpha$, and where the constant C which does not depend on Y and n .

Proof Let us first prove (72b). Since $q + h = 1$ [see (70)], we obtain $q_n + h_n = 1$ from the definitions of v_n (recall that $q_n = 1 + v_n$) and h_n . We consider a fix given interval $I = [\xi_i, \xi_{i+1}]$ and, for convenience, denote by an integral without boundary the weighted integral $\int f(\xi) d\xi = \frac{1}{\Delta\xi} \int_{\xi_i}^{\xi_{i+1}} f(\xi) d\xi$ so that, for $\xi \in I$, $q_n = \int q d\xi$, $w_n = \int w d\xi$ and $h_n = \int h d\xi$. Using Jensen’s inequality, we get that

$$\begin{aligned} q_n^2 + U_n^2 q_n^2 + w_n^2 &= \left(\int q d\xi \right)^2 + U_n^2 \left(\int q d\xi \right)^2 + \left(\int w d\xi \right)^2 \\ &\leq \int q^2 d\xi + U_n^2 \int q^2 d\xi + \int w^2 d\xi \\ &= \int q^2 d\xi + U_n^2 \int q^2 d\xi + \int (q(1 - q) - q^2 U^2) d\xi \\ &= q_n + U_n^2 \int q^2 d\xi - \int (q^2 U^2) d\xi. \end{aligned} \tag{73}$$

Using the Cauchy–Schwarz inequality and the definition of U_n , we obtain

$$U_n^2 \int q^2 d\xi = \frac{(\int q^2 U)^2 d\xi}{\int q^2 d\xi} \leq \frac{\int q^2 d\xi \int q^2 U^2 d\xi}{\int q^2 d\xi} = \int q^2 U^2 d\xi.$$

Hence, (73) yields

$$q_n^2 + U_n^2 q_n^2 + w_n^2 \leq q_n$$

which, as $q_n + h_n = 1$, is equivalent to $q_n h_n \geq U_n^2 q_n^2 + w_n^2$. Let us now prove (72a). A direct computation shows that

$$\| \mathbf{P}_R \mathbf{I}_{\Delta\xi}(f) \|_{L^2} \leq \| f \|_{L^2}, \tag{74}$$

for any $f \in L^2(\mathbb{R})$ and any n . Since $\lim_{n \rightarrow \infty} \| \mathbf{P}_R \mathbf{I}_{\Delta\xi}(f) - f \|_{L^2} = 0$ for any smooth function f with compact support, we obtain, by density and (74), that the same result holds for any $f \in L^2(\mathbb{R})$. Hence,

$$\lim_{n \rightarrow \infty} \| q_n - q \|_{L^2} = 0, \quad \lim_{n \rightarrow \infty} \| w_n - w \|_{L^2} = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \| h_n - h \|_{L^2} = 0.$$

On the interval $I = [\xi_i, \xi_{i+1}]$, we have

$$|U_n(\xi) - U(\xi)| = \left| \frac{\int q^2(\eta)(U(\eta) - U(\xi)) d\eta}{\int q^2 d\eta} \right| \leq \frac{\Delta\xi}{2}$$

as $|U_\xi| \leq \frac{1}{2}$, see (71). Hence, $\|U_n - U\|_{L^\infty(-R, R)} \leq \frac{\Delta\xi}{2}$ and

$$\begin{aligned} \|U_n - U\|_{L^\infty} &\leq \|U_n - U\|_{L^\infty(-R, R)} + \|U\|_{L^\infty((-\infty, -R) \cup (R, \infty))} \\ &\leq \frac{\Delta\xi}{2} + \|U\|_{L^\infty((-\infty, -R) \cup (R, \infty))}. \end{aligned} \tag{75}$$

Since $U \in H^1(\mathbb{R})$, $\lim_{\xi \rightarrow \pm\infty} U_n = 0$ and (75) yields $\lim_{n \rightarrow \infty} \|U_n - U\|_{L^\infty} = 0$. We have

$$\begin{aligned} \|U_n - \mathbf{P}_R(U)\|_{L^2}^2 &= \sum_{i=-N}^{N-1} \int_{\xi_i}^{\xi_{i+1}} \left(\frac{\int_{\xi_i}^{\xi_{i+1}} q^2 U \, d\eta}{\int_{\xi_i}^{\xi_{i+1}} q^2 \, d\eta} - U(\xi) \right)^2 d\xi \\ &\leq \sum_{i=-N}^{N-1} \frac{1}{\int_{\xi_i}^{\xi_{i+1}} q^2 \, d\eta} \int_{\xi_i}^{\xi_{i+1}} \int_{\xi_i}^{\xi_{i+1}} q^2(\eta) (U(\xi) - U(\eta))^2 d\xi \, d\eta, \end{aligned} \tag{76}$$

after applying Cauchy–Schwarz. For $\xi, \eta \in I$, we have

$$(U(\xi) - U(\eta))^2 = \left(\int_{\eta}^{\xi} U_{\xi}(\bar{\eta}) \, d\bar{\eta} \right)^2 \leq \Delta\xi \int_{\eta}^{\xi} U_{\xi}(\bar{\eta})^2 \, d\bar{\eta} \leq \Delta\xi \int_{\xi_i}^{\xi_{i+1}} U_{\xi}^2 \, d\bar{\eta}.$$

Hence, (76) yields

$$\|U_n - \mathbf{P}_R(U)\|_{L^2}^2 \leq \sum_{i=-N}^{N-1} \frac{(\Delta\xi)^2}{\int_{\xi_i}^{\xi_{i+1}} q^2 \, d\eta} \int_{\xi_i}^{\xi_{i+1}} q^2 \, d\eta \int_{\xi_i}^{\xi_{i+1}} U_{\xi}^2 \, d\bar{\eta} \leq (\Delta\xi)^2 \|U_{\xi}\|_{L^2}^2.$$

It follows that

$$\begin{aligned} \|U_n - U\|_{L^2} &\leq \|U_n - \mathbf{P}_R(U)\|_{L^2} + \|U - \mathbf{P}_R(U)\|_{L^2} \\ &\leq \Delta\xi \|U_{\xi}\|_{L^2} + \|U\|_{L^2((-\infty, -R) \cup (R, \infty))} \end{aligned}$$

and therefore $\lim_{n \rightarrow \infty} \|U_n - U\|_{L^2} = 0$. The function h belongs to $L^1(\mathbb{R})$ because $h = h^2 + U^2 q^2 + w^2$, by (70). A direct computation shows that

$$\|\mathbf{P}_R \mathbf{I}_{\Delta\xi}(f)\|_{L^1} \leq \|f\|_{L^1}, \tag{77}$$

for any $f \in L^1(\mathbb{R})$ and any n . Since $\lim_{n \rightarrow \infty} \|\mathbf{P}_R \mathbf{I}_{\Delta\xi}(f) - f\|_{L^1} = 0$ for any smooth function f with compact support, we obtain, by density and (77), that the same result holds for any $f \in L^1(\mathbb{R})$. Hence, $\lim_{n \rightarrow \infty} \|h_n - h\|_{L^1} = 0$ and therefore

$$\lim_{n \rightarrow \infty} \|H_n - H\|_{L^\infty} = 0.$$

Since $y_n = \xi - H_n$ and $y = \xi - H$, we get also that $\lim_{n \rightarrow \infty} \|y_n - y\|_{L^\infty} = 0$. Let us look at the bounds on the decay of Y . We assume $Y \in F^e$. We have

$$\begin{aligned} \int_{\mathbb{R}} e^{|\xi|} |h_n| \, d\xi &= \frac{1}{\Delta\xi} \sum_{i=-N}^{N+1} \int_{\xi_i}^{\xi_{i+1}} \int_{\xi_i}^{\xi_{i+1}} e^{|\xi|} |h(\eta)| \, d\eta \, d\xi \\ &= \frac{1}{\Delta\xi} \sum_{i=-N}^{N+1} \int_{\xi_i}^{\xi_{i+1}} \int_{\xi_i}^{\xi_{i+1}} e^{|\xi|} e^{-|\eta|} e^{|\eta|} |h(\eta)| \, d\eta \, d\xi \end{aligned}$$

$$\begin{aligned} &\leq \frac{1}{\Delta\xi} \sum_{i=-N}^{N+1} \int_{\xi_i}^{\xi_{i+1}} \int_{\xi_i}^{\xi_{i+1}} e^{|\xi-\eta|} e^{|\eta|} |h(\eta)| \, d\eta \, d\xi \\ &\leq e^{\Delta\xi} \sum_{i=-N}^{N+1} \int_{\xi_i}^{\xi_{i+1}} e^{|\eta|} |h(\eta)| \, d\eta \leq 3 \|Y\|_{F^{\text{exp}}} , \end{aligned}$$

after assuming, without loss of generality, that $\Delta\xi \leq 1$. Similarly one proves that $\int_{\mathbb{R}} e^{|\xi|} w^2 \, d\xi \leq C \|Y\|_{F^{\text{exp}}}$. It remains to estimate $\int_{\mathbb{R}} U_n^2 e^{|\xi|} \, d\xi$. For any $\eta, \xi \in [\xi_i, \xi_{i+1}]$, we have

$$\begin{aligned} U^2(\eta) &= U^2(\xi) + 2 \int_{\xi}^{\eta} U U_{\xi}(\bar{\xi}) \, d\bar{\xi} \\ &\leq U^2(\xi) + \int_{\xi_i}^{\xi_{i+1}} (U^2 + (U_{\xi})^2)(\bar{\xi}) \, d\bar{\xi} = U^2(\xi) + \int_{\xi_i}^{\xi_{i+1}} (U^2 + w^2)(\bar{\xi}) \, d\bar{\xi} . \end{aligned}$$

Hence,

$$\begin{aligned} U_{i,n}^2 &= \left(\frac{\int_{\xi_i}^{\xi_{i+1}} q^2(\eta) U(\eta) \, d\eta}{\int_{\xi_i}^{\xi_{i+1}} q^2(\eta) \, d\eta} \right)^2 \\ &\leq \frac{\int_{\xi_i}^{\xi_{i+1}} q^2(\eta) U^2(\eta) \, d\eta}{\int_{\xi_i}^{\xi_{i+1}} q^2(\eta) \, d\eta} \quad (\text{by Cauchy-Schwarz}) \\ &\leq U^2(\xi) + \int_{\xi_i}^{\xi_{i+1}} (U^2 + w^2)(\bar{\xi}) \, d\bar{\xi} \end{aligned}$$

for any $\xi \in [\xi_i, \xi_{i+1}]$. Then,

$$\begin{aligned} \int_{\mathbb{R}} e^{|\xi|} U_n^2 \, d\xi &= \sum_{i=-N}^{N-1} \int_{\xi_i}^{\xi_{i+1}} e^{|\xi|} U_{i,n}^2 \, d\xi \\ &\leq \sum_{i=-N}^{N-1} \int_{\xi_i}^{\xi_{i+1}} \left(e^{|\xi|} \left(U^2(\xi) + \int_{\xi_i}^{\xi_{i+1}} (U^2 + w^2)(\bar{\xi}) \, d\bar{\xi} \right) \right) \, d\xi \\ &\leq \int_{\mathbb{R}} e^{|\xi|} U^2(\xi) \, d\xi + \sum_{i=-N}^{N-1} \int_{\xi_i}^{\xi_{i+1}} \int_{\xi_i}^{\xi_{i+1}} (U^2 + w^2)(\bar{\xi}) e^{|\xi|} \, d\xi \, d\bar{\xi} \end{aligned}$$

$$\begin{aligned} &\leq \|Y\|_{F^{\text{exp}}} + \sum_{i=-N}^{N-1} e^{\Delta\xi} \int_{\xi_i}^{\xi_{i+1}} \int_{\xi_i}^{\xi_{i+1}} (U^2 + w^2)(\bar{\xi}) e^{|\bar{\xi}|} d\xi d\bar{\xi} \\ &\leq (1 + 2\Delta\xi e^{\Delta\xi}) \|Y\|_{F^{\text{exp}}} \leq (1 + 2e) \|Y\|_{F^{\text{exp}}} . \end{aligned}$$

Thus we have proved that $\|Y_n\|_{F^{\text{exp}}} \leq C \|Y\|_{F^\alpha}$ for a constant C which does not depend on Y and n . One proves in the same way that $\|Y_n\|_{F^\alpha} \leq C \|Y\|_{F^\alpha}$. \square

6.2 Convergence of the semi-discrete solutions

Let $Y(t)$ and $Y_{\{\Delta\xi, R\}}(t)$ be respectively the solution of (14) with initial data Y_0 and the solution of (56) with initial data $Y_{0, \Delta\xi, R}$. We assume $Y_0 \in F^{\text{exp}}$. Given $T > 0$, we consider the fixed time interval $[0, T]$. Since $Y_0 \in \mathcal{G}$, the solution $Y(t)$ exists globally and

$$\sup_{t \in [0, T]} \|Y(t, \cdot)\|_{F^{\text{exp}}} \leq M$$

for a constant M which depends only on T and $\|Y_0\|_{F^{\text{exp}}}$, see Theorems 3 and 4. The solution $Y_{\{\Delta\xi, R\}}$ does not necessarily exist globally in time. However, we claim that there exists $n > 0$ such that for any $\Delta\xi$ and R such that $\Delta\xi + \frac{1}{R} \leq \frac{1}{n}$, we have

$$\sup_{t \in [0, T]} \|Y_{\{\Delta\xi, R\}}(t, \cdot)\| < 2M. \tag{78}$$

It implies in particular that the solution $Y_{\{\Delta\xi, R\}}$ is defined on $[0, T]$. Let us assume the opposite. Then, there exists a sequence $\Delta\xi_k, R_k$ and $t_k < T$ such that $\lim_{k \rightarrow \infty} \Delta\xi_k = 0, \lim_{k \rightarrow \infty} R_k = \infty$,

$$\sup_{t \in [0, t_k]} \|Y_{\{\Delta\xi, R\}}(t, \cdot)\| = 2M.$$

From (65), we get

$$\sup_{t \in [0, t_k]} \|Y(t, \cdot) - Y_{\{\Delta\xi_k, R_k\}}(t, \cdot)\|_F \leq C(M) \left(\|Y_0 - Y_{0, \Delta\xi_k, R_k}\|_F + \sqrt{\Delta\xi_k} + e^{-R_k} \right). \tag{79}$$

The constant $C(M)$ depends on M but not on $\Delta\xi_k$ and R_k . Thus, we have

$$\begin{aligned} 2M &= \sup_{t \in [0, t_k]} \|Y_{\{\Delta\xi_k, R_k\}}(t, \cdot)\| \leq \|Y(t_k, \cdot)\| + \|Y(t_k, \cdot) - Y_{\{\Delta\xi_k, R_k\}}(t_k, \cdot)\| \\ &\leq M + C \left(\|Y_0 - Y_{0, \Delta\xi, R}\|_F + \sqrt{\Delta\xi_k} + e^{-R_k} \right) \end{aligned}$$

which leads to a contradiction as the right-hand side in the last inequality above tends to M when k tends to infinity. Once (78) is established, Theorem 14 follows from (65).

The same estimates can be obtained for $Y_0 \in F^\alpha$. Without loss of generality, we assume that the approximating sequence satisfies $\|Y_0 - Y_{0,\Delta\xi,R}\|_F \leq \frac{C(M)}{2M}$ where $C(M)$ is given in (79), so that $Y_{\{\Delta\xi,R\}}$ exists on $[0, T]$. Then, we have the following theorem.

Theorem 14 *Given $Y_0 \in F^{\text{exp}}$, for any $T > 0$, there exists a constant $n > 0$ such that, for all $\Delta\xi$ and R such that $\Delta\xi + \frac{1}{R} \leq \frac{1}{n}$, we have*

$$\sup_{t \in [0, T]} \|Y(t, \cdot) - Y_{\{\Delta\xi,R\}}(t, \cdot)\|_F \leq C \left(\|Y_0 - Y_{0,\Delta\xi,R}\|_F + \sqrt{\Delta\xi} + e^{-R} \right).$$

The constant C depends only on $\|Y_0\|_{F^{\text{exp}}}$ and T . Correspondingly, given $Y_0 \in F^\alpha$, we have

$$\sup_{t \in [0, T]} \|Y(t, \cdot) - Y_{\{\Delta\xi,R\}}(t, \cdot)\|_F \leq C \left(\|Y_0 - Y_{0,\Delta\xi,R}\|_F + \sqrt{\Delta\xi} + \frac{1}{R^{\alpha/2}} \right)$$

and C depends only on $\|Y_0\|_{F^\alpha}$ and T .

7 Discretisation in time

In this section, we deal with the numerical integration in time of the system of differential equations (56) which corresponds to the semi-discretisation in space of system (14). The flow of this system of differential equations has some geometric properties and it is of interest to derive numerical schemes that preserve these properties. Such integrators are called geometric numerical schemes, see for example the monograph [14]. Thus we will look for numerical schemes preserving the invariants (57) of our system of differential equations. Moreover, this last property will enable us to show that the numerical schemes preserve the positivity of the energy density. These invariants are quartic functions of Y and we are not aware of schemes preserving quartic polynomials, this is why we first split the system of equations (56) into two pieces. Each sub-system will then have quadratic invariants and we can use a numerical scheme preserving these invariants. The following sub-systems read

$$\begin{aligned} \zeta_{i,t} &= 0 \\ U_{i,t} &= 0 \\ H_{i,t} &= 0 \\ q_{i,t} &= \gamma w_i \quad i = -N, \dots, N - 1 \\ w_{i,t} &= \frac{\gamma}{2} h_i + \left(\frac{3 - 2\gamma}{2} U_i^2 - P_i \right) q_i \\ h_{i,t} &= (3U_i^2 - 2P_i) w_i, \end{aligned} \tag{80}$$

or shortly

$$\bar{Y}_t = \bar{G}_1(\bar{Y}),$$

where $\bar{Y}(t) = (Y_{\{\Delta\xi, R\}}(t, \xi_i))_{i=-N}^{N-1}$ and similarly for \bar{G}_1 . We also define the system of differential equations

$$\begin{aligned}
 \zeta_{i,t} &= \gamma U_i \\
 U_{i,t} &= -Q_i \\
 H_{i,t} &= U_i^3 - 2P_i U_i \\
 q_{i,t} &= 0 \quad i = -N, \dots, N - 1 \\
 w_{i,t} &= 0 \\
 h_{i,t} &= -2Q_i U_i q_i,
 \end{aligned}
 \tag{81}$$

or shortly

$$\bar{Y}_t = \bar{G}_2(\bar{Y}).$$

The space $F_{\{\Delta\xi, R\}}$ is finite dimensional. We denote $\bar{F} = \mathbb{R}^{2N \times 6}$. The mapping from \bar{F} to $F_{\{\Delta\xi, R\}}$

$$\{\bar{Y}_i = (\bar{\zeta}_i, \bar{U}_i, \bar{H}_i, \bar{q}_i, \bar{w}_i, \bar{h}_i)\}_{i=-N}^{N-1} \mapsto Y = (\zeta, U, H, q, w, h)$$

is a bijection, where we define

$$\zeta(\xi) = \sum_{i=-N}^{N-1} (\bar{\zeta}_i \chi_{[\xi_i, \xi_{i+1})}(\xi)) + \bar{\zeta}_{-N} \chi_{(-\infty, -R]}(\xi) + \bar{\zeta}_N \chi_{[R, \infty)}(\xi)$$

and similar definitions for the other components of Y . This mapping is in addition an isometry if we consider the norm

$$\begin{aligned}
 \|\bar{Y}\|_{\bar{F}} &= \|\bar{\zeta}\|_{l^\infty(\mathbb{R}^{2N})} + \|\bar{U}\|_{l^2(\mathbb{R}^{2N})} + \|\bar{H}\|_{l^\infty(\mathbb{R}^{2N})} + \|\bar{H}\|_{l^\infty(\mathbb{R}^{2N})} \\
 &\quad + \|\bar{v}\|_{l^2(\mathbb{R}^{2N})} + \|\bar{w}\|_{l^2(\mathbb{R}^{2N})} + \|\bar{h}\|_{l^2(\mathbb{R}^{2N})},
 \end{aligned}
 \tag{82}$$

where

$$\|\bar{z}\|_{l^2(\mathbb{R}^{2N})} = \left(\Delta\xi \sum_{i=-N}^{N-1} \bar{z}_i^2 \right)^{\frac{1}{2}}$$

for any $\bar{z} \in \mathbb{R}^{2N}$. In the remaining, we will always consider the norm given by (82) for \bar{F} so that the bounds found in the previous sections directly apply. In particular, we have the following lemma, which is a consequence of Proposition 1 and the same arguments that lead to Lemmas 3 and 5.

Lemma 6 *The mappings $\bar{G}_1 : \bar{F} \rightarrow \bar{F}$ and $\bar{G}_2 : \bar{F} \rightarrow \bar{F}$ belong to $C^1(\bar{F}, \bar{F})$ and*

$$\|\bar{G}_1(\bar{Y})\|_{\bar{F}} + \left\| \frac{\partial \bar{G}_1}{\partial \bar{Y}}(\bar{Y}) \right\|_{L(\bar{F}, \bar{F})} \leq C(M),$$

and

$$\|\bar{G}_2(\bar{Y})\|_{\bar{F}} + \left\| \frac{\partial \bar{G}_2}{\partial \bar{Y}}(\bar{Y}) \right\|_{L(\bar{F}, \bar{F})} \leq C(M),$$

for any $\bar{Y} \in \bar{B}_M$, where

$$\bar{B}_M = \{\bar{Y} \in \bar{F} \mid \|\bar{Y}\|_{\bar{F}} \leq M\}.$$

As this was done in the last sections, one can show that both systems possess $\bar{I}_i(Y) = U_i^2 q_i^2 + w_i^2 - q_i h_i$, see (57), as first integrals. That is $\bar{I}'_i(Y)\bar{G}_k(Y) = 0$ for all Y , for $k = 1, 2$ and for $i = -N, \dots, N - 1$. In particular, this implies that every solutions of (80) or (81) satisfy $\bar{I}_i(\bar{Y}(t)) = \bar{I}_i(\bar{Y}(0))$ for $i = -N, \dots, N - 1$ and $t \geq 0$. Having a closer look at the differential equations (80) and (81), one sees that the invariants are now quadratic functions [\bar{U} is constant for (80) and \bar{q} is constant for (81)] and we therefore use a numerical scheme that preserves quadratic invariants.

Proposition 2 *Let us apply a Runge–Kutta scheme with coefficients satisfying*

$$b_i a_{ij} + b_j a_{ji} = b_i b_j \quad \text{for all } i, j = 1, \dots, s \tag{83}$$

to the system (80), then it conserves exactly the invariants $\bar{I}_i(Y) = U_i^2 q_i^2 + w_i^2 - q_i h_i$ for $i = -N, \dots, N - 1$. The same holds if we apply the scheme to (81).

Proof The proof of this proposition is a simple adaptation of the proof of Theorem 2.2 from [14, Chapter IV]. Let us start with system (80). Dropping the indexes and the bars for ease of notations, we first write the invariant $I(Y)$ as

$$I(Y) = Y^T D(Y)Y + d(Y)^T Y$$

with $Y = (\zeta, U, H, q, w, h)$, $D(Y) = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & U^2 & 0 & -1/2 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1/2 & 0 & 0 \end{pmatrix}$ and $d(Y) = 0^T$.

For the Runge–Kutta method, we write $Y_1 = Y_0 + h \sum_{j=1}^s b_j K_j$ with $K_i = G_1(Y_0 + h \sum_{j=1}^s a_{ij} K_j)$. From the definition of the method, of the matrix $D(Y)$ and of the vector $d(Y)$, it follows that

$$\begin{aligned}
 I(Y_1) &= Y_1^T D(Y_1)Y_1 + d(Y_1)^T Y_1 = \left(Y_0 + h \sum_{i=1}^s b_i K_i \right)^T D(Y_0) \left(Y_0 + h \sum_{j=1}^s b_j K_j \right) \\
 &= Y_0^T D(Y_0)Y_0 + h \sum_{i=1}^s b_i K_i^T D(Y_0)Y_0 + h \sum_{j=1}^s b_j Y_0^T D(Y_0)K_j \\
 &\quad + h^2 \sum_{i,j=1}^s b_i b_j K_i^T D(Y_0)K_j.
 \end{aligned}$$

Writing $K_i = G_1(\tilde{Y}_i)$ with $\tilde{Y}_i = Y_0 + h \sum_{j=1}^s a_{ij} K_j$, we obtain that

$$\begin{aligned}
 I(Y_1) &= Y_0^T D(Y_0)Y_0 + 2h \sum_{i=1}^s b_i \tilde{Y}_i^T D(Y_0)G_1(\tilde{Y}_i) \\
 &\quad + h^2 \sum_{i,j=1}^s (b_i b_j - b_i a_{ij} - b_j a_{ji}) K_i^T D(Y_0)K_j.
 \end{aligned}$$

The last term in the above equation vanishes due to condition (83). By definition of the problem and of the matrix $D(Y)$, we have $D(Y_0) = D(\tilde{Y}_i)$ because U is preserved and since $I(Y)$ is a first integral for (80), we get $\tilde{Y}_i^T D(\tilde{Y}_i)G_1(\tilde{Y}_i) = 0$. It thus follows

$$I(Y_1) = Y_0^T D(Y_0)Y_0 + 0 = I(Y_0)$$

and the Runge–Kutta scheme applied to (80) conserves the invariant $I(Y)$.

The proof for system (81) is similar, take $D(Y) = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & q^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$ and $d(Y) =$

$(0, 0, 0, 0, 0, -q)^T$. □

Let us consider the following differential equation $y_t(t) = f(y(t))$. The implicit midpoint rule

$$y_1 = y_0 + \Delta t f \left(\frac{y_1 + y_0}{2} \right)$$

satisfies the condition (83) and thus preserves quadratic invariants. The implicit midpoint rule will be the building block for the construction of the schemes we will use for the numerical experiments in Sect. 9. For other schemes preserving quadratic invariants, we refer to [14] for example.

As a direct consequence of Proposition 2, we have the following result.

Theorem 15 *Let us apply a Runge–Kutta scheme $\Phi_{\Delta t}^1$, resp. $\Phi_{\Delta t}^2$, with coefficients satisfying (83) to the system (80), resp. (81), with time step size Δt . Then the Lie–Trotter splitting*

$$\Phi_{\Delta t} := \Phi_{\Delta t}^2 \circ \Phi_{\Delta t}^1$$

has order of convergence one and preserves all the invariants \bar{I}_i for $i = -N, \dots, N - 1$. The Strang splitting

$$\Phi_{\Delta t} := \Phi_{\Delta t/2}^1 \circ \Phi_{\Delta t}^2 \circ \Phi_{\Delta t/2}^1$$

is symmetric, has thus order of convergence two and preserves all the invariants \bar{I}_i for $i = -N, \dots, N - 1$.

If we take for $\Phi_{\Delta t}^i, i = 1, 2$, the implicit midpoint rule, we obtain a first order splitting scheme for (56) that preserves exactly the invariants (a second order scheme is obtained using the Strang splitting). This will be the schemes that we will consider in the numerical experiments of Sect. 9.

8 Full discretisation

Our concern is now to combine the results from the last two sections and to show that our numerical schemes converge to the exact solution of the system of Eqs. (14). We integrate $\bar{Y}(t)$ on the time interval $[0, T]$ and obtain \bar{Y}_j for the time steps $j \Delta t, j = 0, \dots, N_T$ where $\Delta t = \frac{T}{N_T}$. We have the following convergence result.

Theorem 16 *Given initial values Y_0 in F^{exp} and $\bar{Y}_0 \in F_R$, for the Lie–Trotter splitting we have*

$$\max_{j \in \{0, \dots, N_T\}} \|S_{j \Delta t}(Y_0) - \Phi_{j \Delta t}(\bar{Y}_0)\|_F \leq C \left(\|Y_0 - \bar{Y}_0\|_F + \sqrt{\Delta \xi} + e^{-R} + \Delta t \right), \tag{84}$$

where we recall that S_t stands for the semigroup of solutions to (14) and, where the constant C depends only on $\|Y_0\|_{F^{\text{exp}}}, \|\bar{Y}_0\|_{F^{\text{exp}}}$ and T . Correspondingly, given initial values Y_0 in F^α and $\bar{Y}_0 \in F_R$, we have

$$\max_{j \in \{0, \dots, N_T\}} \|S_{j \Delta t}(Y_0) - \Phi_{j \Delta t}(\bar{Y}_0)\|_F \leq C \left(\|Y_0 - \bar{Y}_0\|_F + \sqrt{\Delta \xi} + \frac{1}{R^{\alpha/2}} + \Delta t \right), \tag{85}$$

where the constant C depends only on $\|Y_0\|_{F^\alpha}, \|\bar{Y}_0\|_{F^\alpha}$ and T . The same results hold for the Strang splitting with second order accuracy in time, that is, when we replace Δt with Δt^2 in (84).

Let us denote $Y(t) = S_t(Y_0)$ and

$$\Phi_t(\bar{Y}_0) = \frac{((j + 1)\Delta t - t)\Phi_{j\Delta t}(\bar{Y}_0) + (t - j\Delta t)\Phi_{(j+1)\Delta t}(\bar{Y}_0)}{\Delta t}$$

for $t \in [j\Delta t, (j + 1)\Delta t]$. We can rewrite (84) as

$$\max_{t \in [0, T]} \|S_t(Y_0) - \Phi_t(\bar{Y}_0)\|_F \leq C \left(\|Y_0 - \bar{Y}_0\|_F + \sqrt{\Delta\xi} + e^{-R} + \Delta t \right).$$

Proof of Theorem 16 To estimate the total error

$$\|S_{j\Delta t}(Y_0) - \Phi_{j\Delta t}(\bar{Y}_0)\|_F$$

we split it in time and in space. Let us start with the error in time. The proof follows basically the steps of the standard proof of the convergence of numerical scheme for ordinary differential equations. The crucial point is that we guarantee here that the convergence rate in time is independent of the discretisation step in space. Let us first prove the following claim: Given $M > 0$, for any $\bar{Y} \in \bar{B}_M$ and $\bar{Z} \in \bar{B}_M$, we have

$$\begin{aligned} &\Phi_{\Delta t}(\bar{Y}) - \varphi_{\Delta t}(\bar{Z}) \\ &= \bar{Y} - \bar{Z} + \Delta t \left(\bar{G}_1(\bar{Y}) - \bar{G}_1(\bar{Z}) + \bar{G}_2(\bar{Y}) - \bar{G}_2(\bar{Z}) \right) + \mathcal{O}(\Delta t^2), \end{aligned} \tag{86}$$

where $\varphi_{\Delta t}(\bar{Z})$ stands for the exact flow of (56) at time Δt with starting values \bar{Z} . Here, and in the following, the \mathcal{O} -notation stands for an element in \bar{F} satisfying

$$\|\mathcal{O}(\varepsilon)\|_{\bar{F}} \leq C(M)\varepsilon$$

for all $\varepsilon > 0$, where the constant $C(M)$ depends on M but is independent on R and on the space grid size $\Delta\xi$. We first show that the midpoint rule

$$\Phi_{\Delta t}^j(\bar{Y}) = \bar{Y} + \Delta t \bar{G}_j \left(\frac{\Phi_{\Delta t}^j(\bar{Y}) + \bar{Y}}{2} \right),$$

applied to Eq. (80), resp. (81), is at least first order accurate. To do this, let us introduce the mapping $K : \bar{F} \times \bar{F} \rightarrow \bar{F}$ given by

$$K(\bar{Z}, \bar{Y}) = \bar{Z} - \bar{Y} - \Delta t \bar{G}_1 \left(\frac{\bar{Z} + \bar{Y}}{2} \right).$$

We have $K(\Phi_{\Delta t}^1(\bar{Y}), \bar{Y}) = 0$. Since

$$\frac{\partial K}{\partial \bar{Z}}(\bar{Y}) = \text{Id} - \frac{\Delta t}{2} \frac{\partial \bar{G}_1}{\partial \bar{Y}} \left(\frac{\bar{Z} + \bar{Y}}{2} \right)$$

and $\left\| \frac{\partial \bar{G}_1}{\partial \bar{Y}}(\bar{Y}) \right\|_{\bar{F}} \leq C(M)$ (by Lemma 6), there exist $C(M)$ such that, for $\Delta t \leq \frac{1}{C(M)}$, we have that $\frac{\partial K}{\partial Z}(\bar{Y})$ is invertible. By the implicit function Theorem, we get that $\Phi_{\Delta t}^1(\bar{Y})$ is well-defined. Moreover, also following from the implicit function Theorem, we get that

$$\left\| \Phi_{\Delta t}^1(\bar{Y}) \right\|_{\bar{F}} \leq C(M).$$

Then,

$$\begin{aligned} \Phi_{\Delta t}^1(\bar{Y}) &= \bar{Y} + \Delta t \bar{G}_1 \left(\bar{Y} + \frac{\Delta t}{2} \bar{G}_1 \left(\frac{\Phi_{\Delta t}^1(\bar{Y}) + \bar{Y}}{2} \right) \right) \\ &= \bar{Y} + \Delta t \bar{G}_1(\bar{Y}) + \mathcal{O}(\Delta t^2) \end{aligned}$$

by Lemma 6. Using Lemma 6 again, we obtain for the exact flow of (80) that

$$\varphi_{\Delta t}^1(\bar{Z}) = \bar{Z} + \Delta t \bar{G}_1(\bar{Z}) + \mathcal{O}(\Delta t^2).$$

Following the same arguments, we obtain that

$$\Phi_{\Delta t}^2(\Phi_{\Delta t}^1(\bar{Y})) = \Phi_{\Delta t}^1(\bar{Y}) + \Delta t \bar{G}_2(\Phi_{\Delta t}^1(\bar{Y})) + \mathcal{O}(\Delta t^2)$$

and for the composition of the exact flows

$$\varphi_{\Delta t}^2(\varphi_{\Delta t}^1(\bar{Z})) = \varphi_{\Delta t}^1(\bar{Z}) + \Delta t \bar{G}_2(\varphi_{\Delta t}^1(\bar{Z})) + \mathcal{O}(\Delta t^2).$$

Hence,

$$\begin{aligned} &\Phi_{\Delta t}^2(\Phi_{\Delta t}^1(\bar{Y})) - \varphi_{\Delta t}^2(\varphi_{\Delta t}^1(\bar{Z})) \\ &= \Phi_{\Delta t}^1(\bar{Y}) + \Delta t \bar{G}_2(\Phi_{\Delta t}^1(\bar{Y})) - \varphi_{\Delta t}^1(\bar{Z}) - \Delta t \bar{G}_2(\varphi_{\Delta t}^1(\bar{Z})) + \mathcal{O}(\Delta t^2) \\ &= \bar{Y} - \bar{Z} + \Delta t (\bar{G}_1(\bar{Y}) - \bar{G}_1(\bar{Z})) \\ &\quad + \Delta t (\bar{G}_2(\bar{Y} + \Delta t \bar{G}_1(\bar{Y}) + \mathcal{O}(\Delta t^2)) - \bar{G}_2(\bar{Z} + \Delta t \bar{G}_1(\bar{Z}) + \mathcal{O}(\Delta t^2))) \\ &\quad + \mathcal{O}(\Delta t^2) \\ &= \bar{Y} - \bar{Z} + \Delta t (\bar{G}_1(\bar{Y}) - \bar{G}_1(\bar{Z}) + \bar{G}_2(\bar{Y}) - \bar{G}_2(\bar{Z})) + \mathcal{O}(\Delta t^2). \end{aligned} \tag{87}$$

We consider now the splitting error. We have

$$\varphi_{\Delta t}(\bar{Z}) - \bar{Z} = \Delta t \bar{G}(\bar{Z}) + \mathcal{O}(\Delta t^2)$$

and

$$\varphi_{\Delta t}^1(\bar{Z}) - \bar{Z} = \Delta t \bar{G}_1(\bar{Z}) + \mathcal{O}(\Delta t^2)$$

and thus

$$\varphi_{\Delta t}^2(\varphi_{\Delta t}^1(\bar{Z})) = \varphi_{\Delta t}^1(\bar{Z}) + \Delta t \bar{G}_2(\varphi_{\Delta t}^1(\bar{Z})) + \mathcal{O}(\Delta t^2).$$

Hence,

$$\begin{aligned} \varphi_{\Delta t}^2(\varphi_{\Delta t}^1(\bar{Z})) - \varphi_{\Delta t}(\bar{Z}) &= \Delta t \bar{G}(\bar{Z}) - \Delta t \bar{G}_1(\bar{Z}) \\ &\quad - \Delta t \bar{G}_2(\bar{Z} + \Delta t \bar{G}_1(\bar{Z}) + \mathcal{O}(\Delta t^2)) + \mathcal{O}(\Delta t^2) \\ &= \Delta t (\bar{G}(\bar{Z}) - \bar{G}_1(\bar{Z}) - \bar{G}_2(\bar{Z})) + \mathcal{O}(\Delta t^2) \\ &= \mathcal{O}(\Delta t^2), \end{aligned} \quad (88)$$

as $\bar{G} = \bar{G}_1 + \bar{G}_2$. Combining (88) and (87), we obtain (86) and the claim is proved. Let us now set

$$M = \sup_{t \in [0, T]} \|\varphi_t(\bar{Y}_0)\|_F.$$

For a given Δt , we define

$$j_{\Delta t} = \max\{j \in \{0, \dots, N_T - 1\} \mid \|\Phi_{\bar{j}_{\Delta t}}(\bar{Y}_0)\|_{\bar{F}} \leq 2M \text{ for all } \bar{j} \leq j\}. \quad (89)$$

For $j \leq j_{\Delta t}$, we get from (86) that

$$\|\Phi_{(j+1)\Delta t}(\bar{Y}_0) - \varphi_{(j+1)\Delta t}(\bar{Y}_0)\|_F \leq (1 + C(M)\Delta t) \|\Phi_{(j)\Delta t}(\bar{Y}_0) - \varphi_{(j)\Delta t}(\bar{Y}_0)\|_F + \mathcal{O}(\Delta t^2).$$

By induction, it follows that

$$\begin{aligned} \|\Phi_{(j+1)\Delta t}(\bar{Y}_0) - \varphi_{(j+1)\Delta t}(\bar{Y}_0)\|_F &\leq \|\mathcal{O}(\Delta t^2)\| \sum_{k=0}^j (1 + C(M)\Delta t)^k \\ &\leq \|\mathcal{O}(\Delta t^2)\| \frac{1}{C(M)\Delta t} \end{aligned}$$

and therefore

$$\Phi_{(j+1)\Delta t}(\bar{Y}_0) = \varphi_{(j+1)\Delta t}(\bar{Y}_0) + \mathcal{O}(\Delta t). \quad (90)$$

We claim that there exists a constant $C(M)$ such that for all $\Delta t \leq \frac{1}{C(M)}$, we have $j_{\Delta t} = N_T - 1$ and therefore (90) holds for all $j \leq N_T - 1$. Let us assume the opposite. Then, there exists Δt_k such that $\lim_{k \rightarrow \infty} \Delta t_k = 0$ and $j_{\Delta t_k} < N_T - 1$. By definition (89), we have $\|\Phi_{(j_{\Delta t_k}+1)\Delta t_k}(\bar{Y}_0)\|_{\bar{F}} > 2M$. Then, (90) implies

$$\begin{aligned} 2M &\leq \|\Phi_{(j_{\Delta t_k}+1)\Delta t_k}(\bar{Y}_0) - \varphi_{(j_{\Delta t_k}+1)\Delta t_k}(\bar{Y}_0)\|_F + \|\varphi_{(j_{\Delta t_k}+1)\Delta t_k}(\bar{Y}_0)\|_F \\ &\leq \mathcal{O}(\Delta t_k) + M \end{aligned}$$

which leads to a contradiction when k tends to ∞ . Finally, for the total error in space and time, we have:

$$\|S_{j\Delta t}(Y_0) - \Phi_{j\Delta t}(\bar{Y}_0)\|_F \leq \|S_{j\Delta t}(Y_0) - \varphi_{j\Delta t}(\bar{Y}_0)\|_F + \|\varphi_{j\Delta t}(\bar{Y}_0) - \Phi_{j\Delta t}(\bar{Y}_0)\|_F,$$

where all the functions are evaluated at time $j\Delta t$ for $j \leq N_T$. The first term can be estimated using Theorem 14 and we thus obtain

$$\max_{j \in \{0, \dots, N_T\}} \|S_{j\Delta t}(Y_0) - \varphi_{j\Delta t}(\bar{Y}_0)\|_F \leq C \left(\|Y_0 - \bar{Y}_0\|_F + \sqrt{\Delta\xi} + e^{-R} \right).$$

For the second one we use (90) and this concludes the proof of the theorem for the Lie–Trotter splitting. If we had taken the Strang splitting instead, we would have obtained an error in time of order two since this scheme is symmetric. The proof for initial data in F^α is the same. \square

Our next task will be to show that our schemes preserve the positivity of the particle density and of the energy density as does the exact solution of (14) with initial data given by Theorem 13. In order to prove this result, we introduce F^∞ defined as

$$F^\infty = \{Y = (y, U, H, q, w, h) \in F \mid \|q\|_{L^\infty} + \|w\|_{L^\infty} + \|h\|_{L^\infty} < \infty\}$$

with the norm

$$\|Y\|_{F^\infty} = \|Y\|_F + \|q\|_{L^\infty} + \|w\|_{L^\infty} + \|h\|_{L^\infty}.$$

We know that the space F^∞ is preserved by the governing equations (14), see Lemma 2. Using the semilinear structure of (14d)–(14f) with respect to q, w, h , one can show in the same way that (18) was shown, that, for a given $M > 0$,

$$\|G(Y)\|_{F^\infty} + \left\| \frac{\partial G}{\partial Y}(Y) \right\|_{L(F^\infty, F^\infty)} \leq C(M) \tag{91}$$

for any $Y \in B_M^\infty = \{Y \in F^\infty \mid \|Y\|_{F^\infty} \leq M\}$. The same result holds for the mappings $G_{\Delta\xi}, G_{\Delta\xi, R}, \bar{G}_1$ and \bar{G}_2 . In particular we can prove, as in Theorem 16 for the proof of (86), that

$$\begin{aligned} &\Phi_{\Delta t}(\bar{Y}) - \varphi_{\Delta t}(\bar{Z}) \\ &= \bar{Y} - \bar{Z} + \Delta t (\bar{G}_1(\bar{Y}) - \bar{G}_1(\bar{Z}) + \bar{G}_2(\bar{Y}) - \bar{G}_2(\bar{Z})) + \mathcal{O}(\Delta t^2), \end{aligned}$$

where the definition of $\mathcal{O}(\cdot)$ is replaced by

$$\|\mathcal{O}(\varepsilon)\|_{\bar{F}^\infty} \leq C(M)\varepsilon.$$

Here, $\bar{F}^\infty = \bar{F} = \mathbb{R}^{2N \times 6}$ but equipped with the norm derived from $\|\cdot\|_{F^\infty}$, see (82).

Theorem 17 *We consider an initial datum which satisfy*

$$q_i^0 h_i^0 \geq (U_i^0 q_i^0)^2 + (w_i^0)^2, \quad q_i^0 \geq 0, \quad h_i^0 \geq 0 \quad \text{and} \quad q_i^0 + h_i^0 \geq c$$

for all $i = -N, \dots, N - 1$, for some constant $c > 0$. Then, given $T > 0$, there exists $n > 0$, which depends only on $c, \|\bar{Y}^0\|_{F^\infty}$ and T , such that, if $\Delta\xi + \frac{1}{R} + \Delta t < \frac{1}{n}$, the positivity of the particle density $1/q$ and of the energy density h are preserved by our numerical discretisation, that is,

$$q_i^j \geq 0 \quad \text{and} \quad h_i^j \geq 0,$$

for $i = -N, \dots, N - 1$ and $j = 1, \dots, N_T$.

Proof The main idea of the proof is to control the growth of $1/(q_i^k + h_i^k)$. To do so we adapt the proof of Lemma 2 to this discrete situation. Let $M = 2 \sup_{t \in [0, T]} \|\varphi_t(\bar{Y}_0)\|_{F^\infty}$. As in the proof of Theorem 16, we can prove that for Δt small enough (the bound depending only on M), we have

$$\|\Phi_{k\Delta t}(\bar{Y}_0)\|_{F^\infty} \leq 2M$$

for all $k = 0, \dots, N_T$. For $k < N_T$, we have, by definition of our scheme, that

$$\begin{aligned} \frac{1}{q_i^{k+1} + h_i^{k+1}} - \frac{1}{q_i^k + h_i^k} &= -\frac{q_i^{k+1} - q_i^k + h_i^{k+1} - h_i^k}{(q_i^{k+1} + h_i^{k+1})(q_i^k + h_i^k)} \\ &= -\frac{\Delta t(\gamma w_i^k - 2Q(Y^k)U_i^k q_i^k + (3(U_i^k)^2 - 2P(Y^k))w_i^k) + \mathcal{O}(\Delta t^2)}{(q_i^{k+1} + h_i^{k+1})(q_i^k + h_i^k)}. \end{aligned}$$

Hence, using the bounds (91), we get

$$\left| \frac{1}{q_i^{k+1} + h_i^{k+1}} - \frac{1}{q_i^k + h_i^k} \right| \leq \frac{\Delta t C(M)}{|q_i^{k+1} + h_i^{k+1}|} \left(\frac{|w_i^k| + |q_i^k| + \Delta t}{|q_i^k + h_i^k|} \right). \tag{92}$$

Let us prove by induction that, for Δt small enough (depending only M),

$$\frac{1}{q_i^k + h_i^k} \leq \frac{1}{c} e^{2C(M)T} + 1, \quad q_i^k \geq 0 \quad \text{and} \quad h_i^k \geq 0 \tag{93}$$

for $i = -N, \dots, N - 1$, all $k = 0, \dots, N_T$ and where $C(M)$ is the constant given in (92). By definition of our initial data, these assumptions hold for $k = 0$. We assume now that (93) holds for $k = 0, \dots, j$ and we want to prove that it also holds for $j + 1$. We set $\bar{M} = \frac{1}{c} e^{2C(M)T} + 1$. Since the numerical schemes preserve the invariant $q_i^k h_i^k = (U_i^k q_i^k)^2 + (w_i^k)^2$, we obtain in particular that

$$q_i^k h_i^k \geq (U_i^k q_i^k)^2 + (w_i^k)^2 \tag{94}$$

for all $k = 0, \dots, N_T$. From this, it follows that $|w_i^k| \leq \frac{1}{\sqrt{2}}(q_i^k + h_i^k)$ as $q_i^k \geq 0$ and $h_i^k \geq 0$. For $k \leq j$, we get from (92) and our induction hypothesis that

$$\left| \frac{1}{q_i^{k+1} + h_i^{k+1}} - \frac{1}{q_i^k + h_i^k} \right| \leq \frac{\Delta t C(M)}{|q_i^{k+1} + h_i^{k+1}|} \left(1 + \frac{1}{\sqrt{2}} + \bar{M} \Delta t \right). \tag{95}$$

From the above equation, we get

$$\left| \frac{1}{q_i^{k+1} + h_i^{k+1}} \right| \leq \frac{1}{1 - 2C(M)\Delta t - \bar{M}C(M)\Delta t^2} \left| \frac{1}{q_i^k + h_i^k} \right|$$

and therefore

$$\begin{aligned} \left| \frac{1}{q_i^{j+1} + h_i^{j+1}} \right| &\leq \frac{1}{(1 - 2C(M)\Delta t - \bar{M}C(M)\Delta t^2)^j} \left| \frac{1}{q_i^0 + h_i^0} \right| \\ &\leq \frac{1}{c(1 - 2C(M)\Delta t - \bar{M}C(M)\Delta t^2)^{\frac{T}{\Delta t}}}. \end{aligned}$$

We have

$$\lim_{\Delta t \rightarrow 0} \frac{1}{c(1 - 2C(M)\Delta t - \bar{M}C(M)\Delta t^2)^{\frac{T}{\Delta t}}} = \frac{1}{c} e^{2C(M)T} < \bar{M}.$$

Therefore, by taking Δt small enough, depending only on the value of M and not on the number of induction steps j , we get

$$\left| \frac{1}{q_i^{j+1} + h_i^{j+1}} \right| \leq \bar{M}.$$

Using the above inequality and (95), we obtain

$$-\frac{1}{q_i^{j+1} + h_i^{j+1}} + \frac{1}{q_i^j + h_i^j} \leq \bar{M} \Delta t C(M) \left(1 + \frac{1}{\sqrt{2}} + \bar{M} \Delta t \right)$$

so that $\frac{1}{q_i^{j+1} + h_i^{j+1}} \geq 0$ for a sufficiently small Δt . By (94), we have that $q_i^{j+1} h_i^{j+1} \geq 0$ and therefore

$$q_i^{j+1} \geq 0 \quad \text{and} \quad h_i^{j+1} \geq 0,$$

which concludes our proof by induction. □

Now we go back to the original set of coordinates. Given an initial datum $u_0 \in H^{1,\text{exp}}(\mathbb{R})$ or $H^{1,\alpha}(\mathbb{R})$, we construct the initial datum Y_0 as given by (69). Then the function $u(t, x)$ defined as

$$u(t, x) = U(t, \xi) \quad \text{for } y(t, \xi) = x \tag{96}$$

is well-defined, is a weak solution to (3) which corresponds to the global conservative solution. The definition (96) of $u(t, x)$ means that for any given time t the set of points

$$(y(t, \xi), U(t, \xi)) \in \mathbb{R}^2 \quad \text{for } \xi \in \mathbb{R}$$

is the graph of $u(t, x)$. Let $\frac{1}{n} = \Delta\xi + \frac{1}{R} + \Delta t$ so that n tends to infinity if and only if $\Delta\xi, \Delta t$ tend to zero and R tends to infinity. We consider an approximating sequence $Y_{0,n}$ which satisfies the conditions (72a) and (72c) of the sequence of initial values which is constructed in Sect. 6. Let $Y_n(t) = \Phi(Y_{0,n})$. From Theorem 16, we obtain the following convergence theorem.

Theorem 18 *The full discretised scheme provide us with points which converge to the graph of the exact conservative solution $u(t, x)$. Indeed, if $u_0 \in H^{1,\text{exp}}(\mathbb{R})$, we have*

$$\begin{aligned} & \max_{\substack{i=-N,\dots,N-1 \\ j=0,\dots,N_T}} |(y_n(t_j, \xi_i), U_n(t_j, \xi_i)) - (y(t_j, \xi_i), U(t_j, \xi_i))| \\ & \leq C \left(\|Y_0 - \bar{Y}_0\|_F + \sqrt{\Delta\xi} + e^{-R} + \Delta t \right), \end{aligned}$$

where the constant C depends only on $\|u_0\|_{H^{1,\text{exp}}}$ and, if $u_0 \in H^{1,\alpha}(\mathbb{R})$,

$$\begin{aligned} & \max_{\substack{i=-N,\dots,N-1 \\ j=0,\dots,N_T}} |(y_n(t_j, \xi_i), U_n(t_j, \xi_i)) - (y(t_j, \xi_i), U(t_j, \xi_i))| \\ & \leq C \left(\|Y_0 - \bar{Y}_0\|_F + \sqrt{\Delta\xi} + \frac{1}{R^{\alpha/2}} + \Delta t \right), \end{aligned} \tag{97}$$

where the constant C depends only on $\|u_0\|_{H^{1,\alpha}}$.

Since

$$|y(t, \xi_{i+1}) - y(t, \xi_i)| = \left| \int_{\xi_i}^{\xi_{i+1}} q(t, \xi) d\xi \right| \leq C \Delta\xi,$$

where C depends only on $\|Y_0\|_{F^\infty}$, we have an a priori upper bound on the density of points of the graph of u we can approximate by our scheme.

In the case where u_0 does not belong to $H^{1,\alpha}(\mathbb{R})$, we can approximate u_0 by functions $u_{0,k} \in H^{1,\alpha}(\mathbb{R})$, which converge to u_0 in $H^1(\mathbb{R})$. From [18], we know that the change of variable (69) produces sequences $Y_{0,k}$ and Y_0 such that $\lim_{k \rightarrow 0} \|Y_{0,k} - Y_0\|_F = 0$. In this way, by using the results done for functions in F^α , we can approximate the exact solution $Y(t)$ and prove convergence. However, since $\|Y_{0,k}\|_{F^\alpha}$ is not uniformly bounded with respect to k , we lose the control on the error rate (the term $\frac{1}{R^{\alpha/2}}$) which is given by (97).

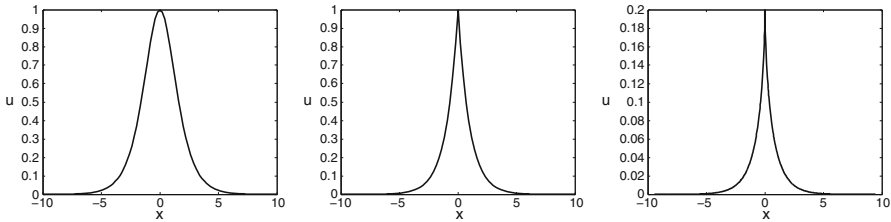


Fig. 1 Traveling waves with decay with speed $c = 1$: smooth ($\gamma = 0.2$), peakon ($\gamma = 1$), cuspon ($\gamma = 5$)

9 Numerical experiments

In this section, we present some numerical experiments for the hyperelastic rod wave equation (1). In order to demonstrate the efficiency of our schemes, we will numerically compute three types of traveling waves with decay, see Fig. 1. The derivation of the cusped ($\gamma > 1$), resp. smooth ($\gamma < 1$), solutions follows the lines of [21]. We refer for example to [23] for a thorough discussion on the peakon case (i.e., $\gamma = 1$).

Let us first start by giving some details related to the implementation of our numerical schemes.

9.1 Algorithm flowchart

- Let us consider a space interval $[-R, R]$ together with an equidistant grid of mesh size $\Delta\xi$. Let Δt denote the time step of our numerical integrator.
- Given $u_0 \in H^1(\mathbb{R})$ an initial value for the hyperelastic rod wave equation (1), we use (69) to compute the initial values

$$Y_{0,\Delta\xi,R} = (y_{0,\Delta\xi,R}, U_{0,\Delta\xi,R}, H_{0,\Delta\xi,R}, q_{0,\Delta\xi,R}, w_{0,\Delta\xi,R}, h_{0,\Delta\xi,R})$$

for the discretised system (56).

- We solve (80), and (81) by using an implicit midpoint rule defined as follows

$$\bar{Y}_{t+\Delta t} := \Phi_{\Delta t}^i(\bar{Y}_t) = \bar{Y}_t + \Delta t \bar{G}_i \left(\frac{\bar{Y}_{t+\Delta t} + \bar{Y}_t}{2} \right) \tag{98}$$

for $i = 1, 2$. We use fixed point iterations to solve the nonlinear system of equations given by (98).

- We finally obtain a symmetric and second-order accurate Strang splitting,

$$\Phi_{\Delta t/2}^1 \circ \Phi_{\Delta t}^2 \circ \Phi_{\Delta t/2}^1,$$

for (56). This numerical integrator preserves all the invariants (57).

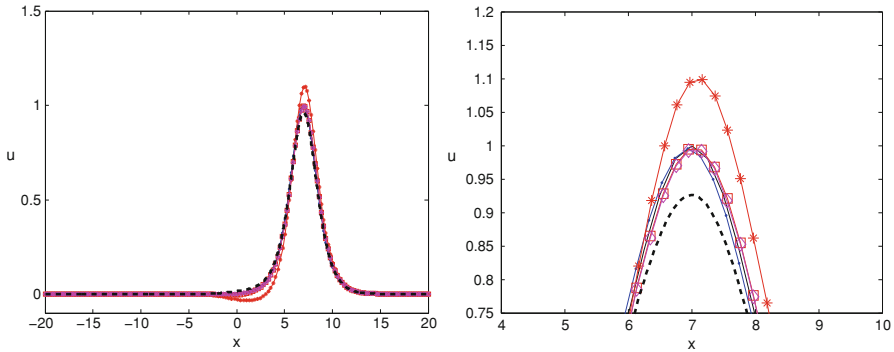


Fig. 2 Exact and numerical solutions of a smooth traveling wave with decay. *Solid line exact, dashed line upwind scheme, dashdotted line ODE45, stars explicit Euler, square Lie–Trotter, diamond Strang*

9.2 Smooth traveling waves with decay ($\gamma < 1$)

According to the classification presented in [21], for a fixed $\gamma \neq 0$, traveling waves $u(x - ct)$ are parametrised by three parameters, M , m and the speed c . Moreover, they are solutions of the following differential equation

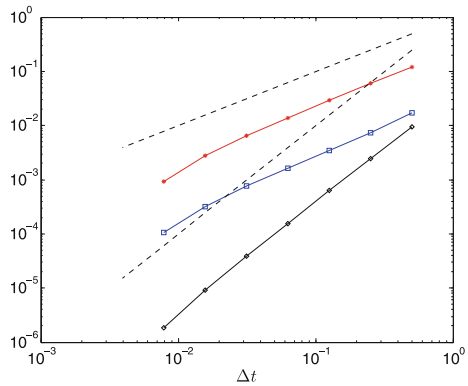
$$u_x^2 = F(u) = \frac{(M - u)(u - m)(u - z)}{c - \gamma u}. \tag{99}$$

For positive values of γ , a smooth traveling wave with decay with $m = \inf_{x \in \mathbb{R}} u(x)$ and $M = \max_{x \in \mathbb{R}} u(x)$ is obtained if $z = m < M < c/\gamma$, where $z := c - M - m$. For our purpose, we have to set $m = 0$ so that the solution decays at infinity. This gives us the conditions $c = M$ and $\gamma < 1$. We thereby obtain the initial values for our system of differential equations (56) by solving (99) numerically. To do this, some care has to be taken as $u \mapsto \sqrt{F(u)}$ is not Lipschitz. We instead solve $u_{xx} = F'(u)/2$. Once this is appropriately done we get the initial values $U_0 = u$, $w_0 = u_x$. We then set $y_0 = \xi$, $q_0 = 1$, $h_0 = U_0^2 + w_0^2$ and $H_0 = \int_{-\infty}^{y_0} h_0$. These initial values do not correspond to the ones defined by (69) but they are equivalent via relabeling and therefore can be used for computation, see [18] for details on the relabeling. We have implemented an upwind scheme based on the original formulation of the Eq. (1), as in [1] but without adaptivity. Figure 2 displays the exact solution together with the numerical solutions given by the upwind scheme, the ODE45 solver from Matlab, the explicit Euler scheme, the Lie–Trotter and the Strang splitting schemes at time $T = 7$. We plot the points

$$(y(t, \xi_i), U(t, \xi_i)), \quad \text{for } i = -N, \dots, N - 1,$$

which approximate the graph of the exact solution $u(t, x)$ for $t = T$. The initial value is a smooth traveling wave with parameters $\gamma = 0.2, m = 0, M = c = 1$, see Fig. 1. We took relatively large discretisation parameters $\Delta\xi = 0.25$ and $\Delta t = 0.1$. For the upwind scheme, we compute the solution $u(t, x)$ in the original space coordinate x . In

Fig. 3 Error in the infinity norm of the explicit Euler scheme (stars), the Lie–Trotter scheme (square) and the Strang scheme (diamond) at time $T = 1$ for the smooth solution. The dashed lines have slopes one, resp. two



this experiment and the others that follow, we consider for this scheme a space discretisation step Δx which is ten times smaller than $\Delta \xi$ and we set $\Delta t = \Delta x / (2 \max(u_0))$. We observe that the explicit Euler scheme gives a less accurate solution than the other schemes and that dissipation occurs for the scheme using the formulation (1). We also observe that, even for these large discretisation parameters, the splitting schemes have the same high as the exact solution, thus following it at the same speed. We do not observe any dissipation. Since both splitting schemes give relative similar results, in what follows, we will only display the results given by the Strang splitting scheme. We finally note that all schemes preserve the positivity of the particle density but only the splitting schemes conserve exactly the invariants from Sect. 7 (these results are not displayed). Let us conclude this subsection with a loglog plot of the temporal order of convergence of the numerical schemes. One can see from Fig. 3 that the order of convergence for the explicit Euler scheme and for the Lie–Trotter splitting scheme is one and the one for the Strang splitting scheme is two, as predicted by Theorem 18. The parameters for this simulation are the same as above, except that $T = 1$ and $\Delta \xi = 0.04$.

We finally want to mention that for negative values of γ , smooth traveling waves with decay also exist. They are obtained if $c/\gamma < m = M < z$.

9.3 Peakon ($\gamma = 1$)

The Camassa–Holm equation, i.e., Eq. (1) with $\gamma = 1$, possesses solutions with a particular shape: the peakons. A single peakon is a traveling wave which is given by

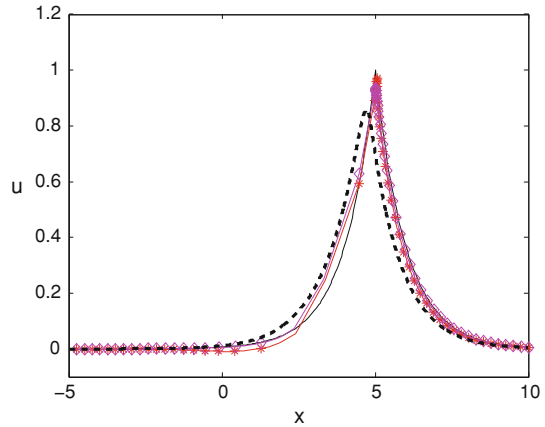
$$u(t, x) = c e^{-|x-ct|}.$$

We note, that at the peak, the derivative of this particular solution is discontinuous. We set the initial values as

$$y_0(\xi) = \xi, \quad U_0(\xi) = u(0, \xi), \quad w_0(\xi) = u_x(0, \xi),$$

$$q_0 = 1, \quad h_0 = U_0^2 + w_0^2, \quad H_0 = \int_{-\infty}^{\xi} h_0(\eta) d\eta.$$

Fig. 4 Exact and numerical solutions at time $T = 5$ for a peakon. *Solid line* exact, *dashed line* upwind scheme, *stars* explicit Euler, *diamond* Strang



In Fig. 4, we display the numerical solutions given by the scheme from [1], the explicit Euler scheme and the Strang splitting for a single peakon traveling from left to right with speed $c = 1$, see Fig. 1. For readability reason, we do not display the solution given by the ODE45 solver, but we note that this numerical solution is very similar to the one given by the splitting scheme. Due to the discontinuity of the derivative, we have to take smaller (in space) discretisation parameters: $\Delta\xi = 0.05$ and $\Delta t = 0.2$. We note more grid-points before the peak and very few just after it, but the speed of the wave is still relatively close to the exact one. This is not the case for schemes based on the Eulerian formulation (1), as illustrated by the numerical solution given by the scheme from [1]. As in the preceding case, only the splitting schemes preserve exactly the invariants of our problem.

The benefit of computing the solutions via an equivalent system in Lagrangian variables becomes clear when comparing the upwind scheme, applied to the original Eq. (1), and an explicit Euler scheme, applied to the system in Lagrangian variables. We compare these two methods as they have the same order of convergence. Then, we observe that the explicit Euler scheme—even with a time and space discretisation step which is ten times larger—gives much better results than the upwind scheme. This has to be balanced with the fact that the system in Lagrangian variables consist of six variables instead of one for the original equation. However this disadvantage becomes marginal as the solution becomes more irregular, as we can see for the cusped traveling wave below.

We would also like to note, that the order of convergence of the numerical schemes are the same as for the smooth solution, see Fig. 3. The results are however not displayed.

9.4 Cusped traveling waves with decay ($\gamma > 1$)

Let us now turn our attention to cusped traveling waves. For $\gamma > 0$, according to the classification given in [21], cusped solutions with $c/\gamma = \max_{x \in \mathbb{R}} u(x)$ and $m = \inf_{x \in \mathbb{R}} u(x)$ are obtained if $z = m = 0 < c/\gamma < M$. This gives us the condition

$c = M$ and thus $\gamma > 1$. The cuspon $u(x)$ satisfies (99), which yields for the indicated values of the parameters

$$u_x = -\sqrt{F(u)} = -\left(\frac{M - u}{c - \gamma u}\right)^{\frac{1}{2}} u \tag{100}$$

for $x \geq 0$ and with the boundary value at zero given by $u(0) = \frac{c}{\gamma}$. For such boundary value, the differential equation (100) is not well-posed and the slope at the top of the cuspon (that is $x = 0$) is indeed equal to infinity. However, we can find a triplet $X = (y, U, H)$ in \mathcal{F} which corresponds to this curve, that is, such that $(u, u^2 + u_x^2 dx) = M(X)$, see (9) for the definition of the map M . The representation of the curve $(x, u(x))$ is not unique: For any diffeomorphism $(\varphi(\xi), u(\varphi(\xi)))$, we obtain an other parameterization of the same curve. Here, we look for a smooth $\varphi(\xi)$ (and we set $y(\xi) = \varphi(\xi)$) such that $U = u(\varphi(\xi)) = u(y(\xi))$ is smooth, even if u is not. We introduce the function

$$g(u) = -\int_{\frac{c}{\gamma}}^u \frac{dz}{\sqrt{F(z)}}.$$

Since $\frac{dx}{du} = -\frac{1}{\sqrt{F(u)}}$, by (100), if we choose

$$U(\xi) = \frac{c}{\gamma} - \xi, \quad y(\xi) = g(U(\xi))$$

then we get, at least for $\xi \in [0, \frac{c}{\gamma}]$, a triplet for which $U(\xi) = u(y(\xi))$. We set the energy density by using (7c) and get

$$H_\xi = U^2 y_\xi + \frac{U_\xi^2}{y_\xi}.$$

However, in this case,

$$y_\xi = g'(U)U_\xi = \left(\frac{c - \gamma U}{M - U}\right)^{\frac{1}{2}} \frac{1}{U}$$

so that $H_\xi(0) = \infty$ and it is incompatible with the requirement that all the derivatives in Lagrangian coordinates are bounded in $L^\infty(\mathbb{R})$, see (7a). Thus, we take

$$U(\xi) = \frac{c}{\gamma} - \xi^2, \quad y(\xi) = g(U(\xi)), \quad H_\xi = U^2 y_\xi + \frac{U_\xi^2}{y_\xi}.$$

In this case, we have

$$y_\xi(\xi) = g'(U)U_\xi = \frac{2}{U(\xi)} \left(\frac{c - \gamma U(\xi)}{M - U(\xi)}\right)^{\frac{1}{2}} \xi = \frac{2\sqrt{\gamma}}{U(\xi)(M - U(\xi))^{\frac{1}{2}}} \xi^2$$

and

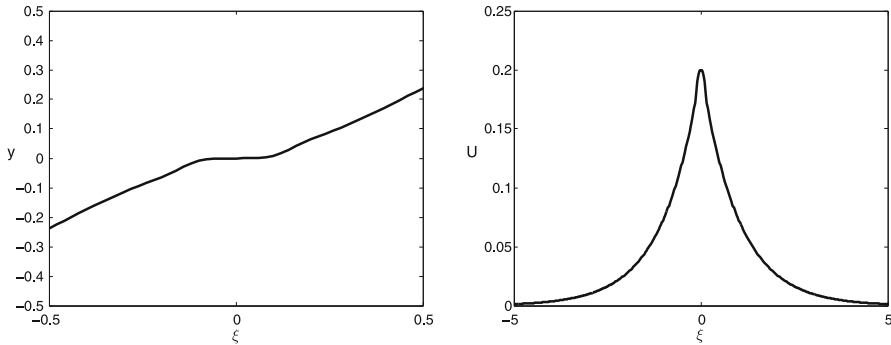


Fig. 5 The function $y(\xi)$ (left picture) and the function $U(\xi)$. Note that these functions are smooth while $u_0(x)$ is not Lipschitz, see Fig. 1

$$H_\xi(0) = \frac{2c}{\gamma^2} (M\gamma - c)^{\frac{1}{2}}$$

is finite. The problem we face now is that the functions are given only on the interval $[0, \frac{c}{\gamma}]$ and $\lim_{\xi \rightarrow \frac{c}{\gamma}} y(\xi) = \infty$. We know that the tail of the cuspon behaves as $u(x) \approx \frac{c}{\gamma} e^{-\sqrt{\frac{M}{c}}x}$ as x tends to ∞ , see [21]. Since we require that $y(\xi) - \xi$ remains bounded, we would like to have $U(\xi) \approx \frac{c}{\gamma} e^{-\sqrt{\frac{M}{c}}\xi}$ for large ξ . Therefore we introduce the following partitions functions χ_1 and χ_2 defined as

$$\chi_1(\xi) = \begin{cases} 1 & \text{if } \xi < a \\ -\frac{1}{b-a}(\xi - b) & \text{for } \xi \in [a, b] \\ 0 & \text{if } x > b \end{cases}$$

and $\chi_2(\xi) = 1 - \chi_1$, where $a < b$ are two parameters. We finally set

$$U(\xi) = \chi_1(\xi) \left(\frac{c}{\gamma} - \xi^2 \right) + \chi_2(\xi) \frac{c}{\gamma} e^{-\sqrt{\frac{M}{c}}\xi}$$

and

$$y(\xi) = g(U(\xi)), \quad H_\xi = U^2 y_\xi + \frac{U_\xi^2}{y_\xi}.$$

By a proper choice of the parameters a and b , we can guarantee that $y_\xi(\xi) \geq 0$ for all $\xi \geq 0$. We extend $X(\xi) = (y(\xi), U(\xi), H(\xi))$ on the whole axis by parity and we obtain an element in \mathcal{F} such that (9) is satisfied. Figure 5 displays $y(\xi)$ and $U(\xi)$. Figure 6 displays the exact solution together with the numerical solutions given by the upwind scheme, the explicit Euler scheme and the Strang splitting scheme at time $T = 6$. As before, we note that the numerical solution given by the ODE45 solver is very similar to the one given by our splitting scheme. The initial value is a cusped traveling wave with parameters $\gamma = 5, m = 0, M = c = 1$, see Fig. 1. For the

Fig. 6 Exact and numerical solutions of a cusped traveling wave with decay. *Solid line* exact, *dashed line* upwind scheme, *stars* explicit Euler, *diamond* Strang

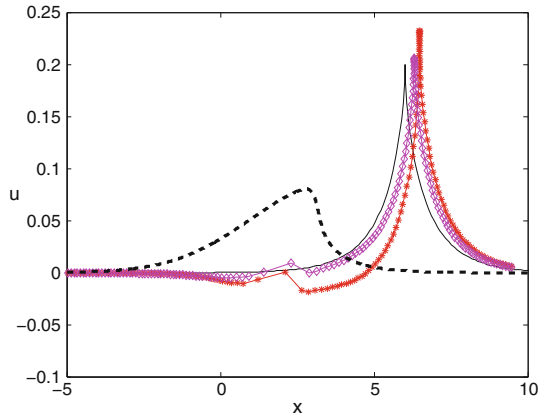
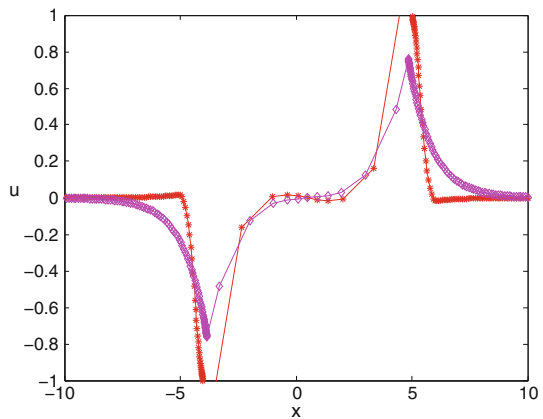


Fig. 7 Peakon–antipeakon collision for $\gamma = 1$. *Stars* explicit Euler, *diamond* Strang



discretisation parameters, we take $\Delta\xi = 0.1$ and $\Delta t = 0.1$. We see that, even for initial data with infinite derivative $u_x(0) = \pm\infty$, the spatial discretisation converges. For the time discretisation, as expected, explicit Euler is less accurate than the other schemes. Note that the oscillation that appears on the left of the peak will disappear as the mesh gets finer. We also remark that only the splitting schemes preserve the positivity of the particle density and conserve the invariants. The upwind scheme performs badly because the solution is not regular. The schemes based on the reformulation in Lagrangian variables do not suffer of that. We also observe that the order of convergence of the numerical schemes are the same as for the smooth solution, see Fig. 3. The results are however not displayed.

We finally note that, for negative values of γ , an anticusped traveling wave with $c/\gamma = \min_{x \in \mathbb{R}} u(x)$ and $m = \sup_{x \in \mathbb{R}} u(x)$ is obtained if $c/\gamma < m = M < z$.

9.5 Peakon–antipeakon collisions

In Fig. 7 we display a collision between a peakon and an antipeakon for $\gamma = 1$. For this problem, the initial value is given by

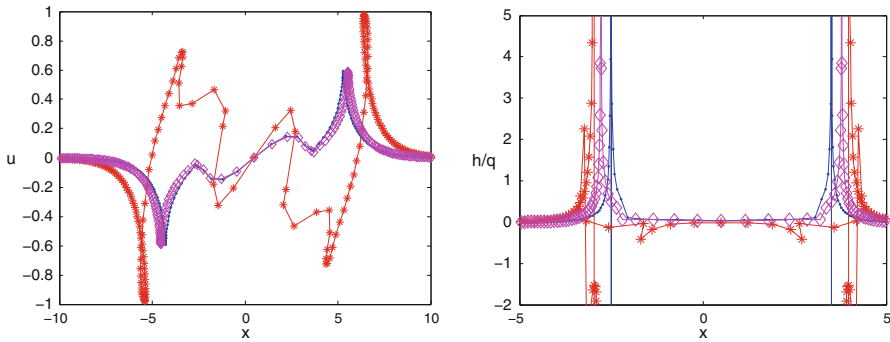


Fig. 8 Peakon–antipeakon collision for $\gamma = 5$ at time $T = 2$ (left) and energy density (right) at the first time, where the numerical solution given by ODE45 is not positive ($q = -1.7394e - 05$). Dashedotted line ODE45, stars explicit Euler, diamond Strang

$$u(0, x) = e^{-|x|} - e^{-|x-1|}.$$

The numerical solutions are computed with grid parameters $\Delta\xi = 0.1$ and $\Delta t = 0.1$ until time $T = 8$. Once again we notice that the spatial discretisation converges. Let us now see what happens for a peakon–antipeakon collision with $\gamma \neq 1$. In Fig. 8 we present a similar experiment as the above one, but where we use $\gamma = 5$ and $T = 2$. Here, we plot the graph given by the points

$$(y(t, \xi_i), \frac{h}{q}(t, \xi_i)), \quad \text{for } i = -N, \dots, N - 1$$

for $t = T$. From the right part of Fig. 8 we see that only the splitting schemes preserve the positivity of the energy density. As always, only the splitting schemes conserve exactly the invariants.

9.6 Collision of smooth traveling waves

We want now to study the behaviour of the numerical schemes when dealing with a collision of smooth traveling waves, as this is an important feature of our numerical scheme to be able to handle such configuration. To do so, we consider the following initial value

$$u(0, x) = -xe^{-x^2/2}.$$

Figure 9 displays the exact solution (i.e., the numerical solution with very small discretisation parameters) for $\gamma = 0.8$. It is remarkable to see that even for such solution, our scheme performs very well. In order to get a better understanding of this problem, we look at the evolution of the waves with time. Figure 10 shows this evolution together with a zoom close to the collision time. We now present the results given by the numerical schemes with grid parameters $\Delta\xi = 0.25$ and $\Delta t = 0.1$ in Fig. 11. We have also

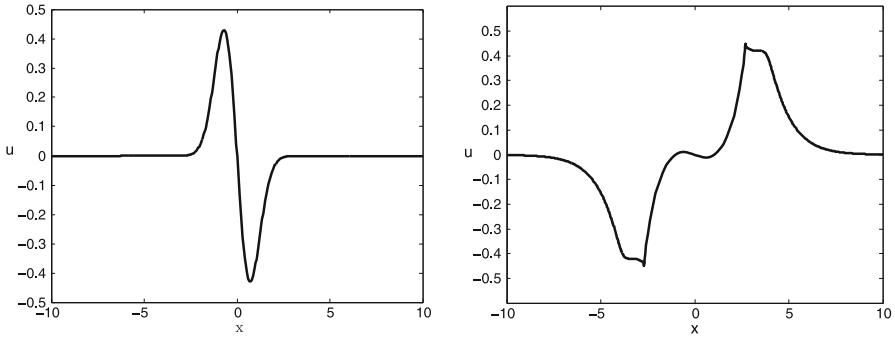


Fig. 9 Collision of smooth traveling waves: Initial datum (*left*) and exact solution at time $T = 11$

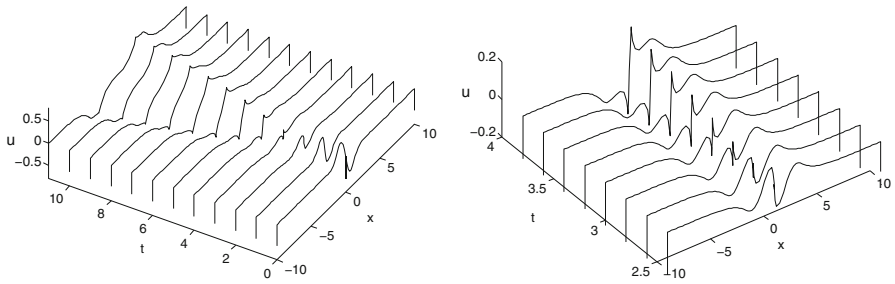
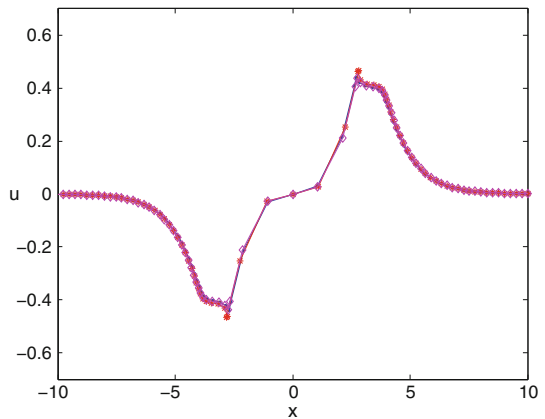


Fig. 10 Collision of smooth traveling waves: evolution in time (*left*) and zoom of the evolution close to the collision

Fig. 11 Collision of smooth traveling waves: numerical solutions at time $T = 11$. *Dashdotted line* ODE45, *stars* explicit Euler, *diamond* Strang



checked that only the splitting schemes preserve the positivity of the particle density and conserve the invariants of our problem. Finally, in Fig. 12 we display, with the same parameter values as above, the evolution in time of the energy density along the numerical solution given by the Strang splitting scheme. We can observe the concentration of the energy and then its separation in two parts, following the waves. With all these numerical observations, we can conclude that the proposed spatial discretisation

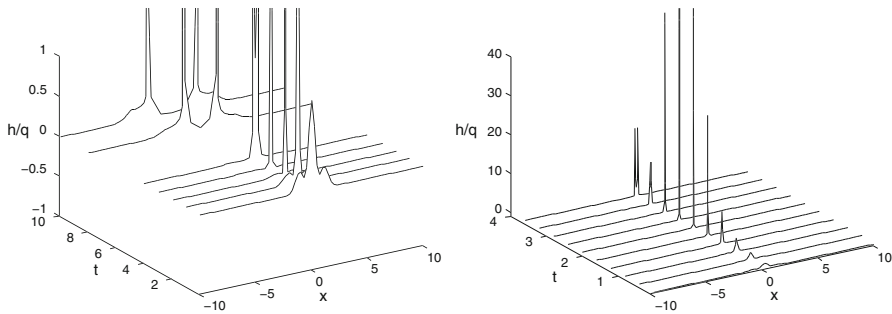


Fig. 12 Evolution of the energy density (*left picture*) along the numerical solution given by the Strang splitting and close up look at the blow up time (*right*)

is robust and qualitatively correct. The time integrators are relatively comparable but only the splitting schemes have the additional properties of maintaining the positivity of the energy density and conserve exactly the invariants of our partial differential equation.

Acknowledgments X. Raynaud wants to thank the institute of mathematics in Basel for its hospitality during a visit in summer 2009 where this work was initiated. A large part of this work was carried out when the authors visited MaGiC. At this place we would also like to thank the CAS in Oslo and the HIM in Bonn. We appreciate the referees' comments on an earlier version.

References

1. Artebrant, R., Schroll, H.J.: Numerical simulation of Camassa–Holm peakons by adaptive upwinding. *Appl. Numer. Math.* **56**(5), 695–711 (2006)
2. Benjamin, T.B., Bona, J.L., Mahony, J.J.: Model equations for long waves in nonlinear dispersive systems. *Philos. Trans. R. Soc. Lond. Ser. A* **272**(1220), 47–78 (1972)
3. Camassa, R., Holm, D.D.: An integrable shallow water equation with peaked solitons. *Phys. Rev. Lett.* **71**(11), 1661–1664 (1993)
4. Camassa, R., Huang, J., Lee, L.: On a completely integrable numerical scheme for a nonlinear shallow-water wave equation. *J. Nonlinear Math. Phys.* **12**(suppl. 1), 146–162 (2005)
5. Camassa, R., Huang, J., Lee, L.: Integral and integrable algorithms for a nonlinear shallow-water wave equation. *J. Comput. Phys.* **216**(2), 547–572 (2006)
6. Coclite, G.M., Holden, H., Karlsen, K.H.: Global weak solutions to a generalized hyperelastic-rod wave equation. *SIAM J. Math. Anal.* **37**(4), 1044–1069 (electronic) (2005)
7. Coclite, G.M., Karlsen, K.H., Risebro, N.H.: A convergent finite difference scheme for the Camassa–Holm equation with general H^1 initial data. *SIAM J. Numer. Anal.* **46**(3), 1554–1579 (electronic) (2008)
8. Coclite, G.M., Karlsen, K.H., Risebro, N.H.: A convergent finite difference scheme for the Camassa–Holm equation with general H^1 initial data. *SIAM J. Numer. Anal.* **46**(3), 1554–1579 (2008)
9. Coclite, G.M., Karlsen, K.H., Risebro, N.H.: An explicit finite difference scheme for the Camassa–Holm equation. *Adv. Differ. Equ.* **13**(7–8), 681–732 (2008)
10. Cohen, D., Owren, B., Raynaud, X.: Multi-symplectic integration of the Camassa–Holm equation. *J. Comput. Phys.* **227**(11), 5492–5512 (2008)
11. Cohen, D., Raynaud, X.: Geometric finite difference schemes for the generalized hyperelastic-rod wave equation. *J. Comput. Appl. Math.* **235**(8), 1925–1940 (2011)
12. Constantin, A., Strauss, W.A.: Stability of a class of solitary waves in compressible elastic rods. *Phys. Lett. A* **270**(3–4), 140–148 (2000)

13. Dai, H.-H.: Exact travelling-wave solutions of an integrable equation arising in hyperelastic rods. *Wave Motion* **28**(4), 367–381 (1998)
14. Hairer, E., Lubich, C., Wanner, G.: *Geometric Numerical Integration, Structure-Preserving Algorithms for Ordinary Differential Equations*. Springer Series in Computational Mathematics, vol. 31. Springer, Berlin (2002)
15. Henry, D.: Compactly supported solutions of the Camassa–Holm equation. *J. Nonlinear Math. Phys.* **12**(3), 342–347 (2005)
16. Holden, H., Raynaud, X.: Convergence of a finite difference scheme for the Camassa–Holm equation. *SIAM J. Numer. Anal.* **44**(4), 1655–1680 (electronic) (2006)
17. Holden, H., Raynaud, X.: A convergent numerical scheme for the Camassa–Holm equation based on multipeakons. *Discret. Contin. Dyn. Syst.* **14**(3), 505–523 (2006)
18. Holden, H., Raynaud, X.: Global conservative solutions of the generalized hyperelastic-rod wave equation. *J. Differ. Equ.* **233**(2), 448–484 (2007)
19. Holden, H., Raynaud, X.: A numerical scheme based on multipeakons for conservative solutions of the Camassa–Holm equation. In: *Hyperbolic problems: theory, numerics, applications*, pp. 873–881. Springer, Berlin (2008)
20. Kalisch, H., Raynaud, X.: Convergence of a spectral projection of the Camassa–Holm equation. *Numer. Methods Partial Differ. Equ.* **22**(5), 1197–1215 (2006)
21. Lenells, J.: Traveling waves in compressible elastic rods. *Discrete Contin. Dyn. Syst. Ser. B* **6**(1), 151–167 (electronic) (2006)
22. Matsuo, T., Yamaguchi, H.: An energy-conserving galerkin scheme for a class of nonlinear dispersive equations. *J. Comput. Phys.* **228**(12), 4346–4358 (2009)
23. Raynaud, X.: On a shallow water wave equation. Ph.D Thesis (2006)
24. Xu, Y., Shu, C.-W.: A local discontinuous Galerkin method for the Camassa–Holm equation. *SIAM J. Numer. Anal.* **46**(4), 1998–2021 (2008)
25. Yin, Z.: On the Cauchy problem for a nonlinearly dispersive wave equation. *J. Nonlinear Math. Phys.* **10**(1), 10–15 (2003)