

## Volume preservation by Runge–Kutta methods



Philipp Bader<sup>a</sup>, David I. McLaren<sup>a</sup>, G.R.W. Quispel<sup>a</sup>, Marcus Webb<sup>b,\*</sup>

<sup>a</sup> Department of Mathematics and Statistics, La Trobe University, 3086 Bundoora VIC, Australia

<sup>b</sup> DAMTP, University of Cambridge, Wilberforce Rd, Cambridge CB3 0WA, UK

### ARTICLE INFO

#### Article history:

Received 6 July 2015

Received in revised form 19 May 2016

Accepted 29 June 2016

Available online 12 July 2016

#### Keywords:

Volume preservation

Runge–Kutta method

Measure preservation

Kahan's method

### ABSTRACT

It is a classical theorem of Liouville that Hamiltonian systems preserve volume in phase space. Any symplectic Runge–Kutta method will respect this property for such systems, but it has been shown by Iserles, Quispel and Tse and independently by Chartier and Murua that no B-Series method can be volume preserving for all volume preserving vector fields. In this paper, we show that despite this result, symplectic Runge–Kutta methods can be volume preserving for a much larger class of vector fields than Hamiltonian systems, and discuss how some Runge–Kutta methods can preserve a modified measure exactly.

© 2016 The Author(s). Published by Elsevier B.V. on behalf of IMACS. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

### 1. Introduction

The construction of numerical schemes for solving ordinary differential equations (ODEs) such that some qualitative geometrical property of the analytical solution is preserved exactly by the numerical solution is an area of great interest and active research today as part of the field of Geometric Integration. The most developed topic in this context is that of integrating Hamiltonian systems while preserving the symplecticity of the flow, and it was found that a class of Runge–Kutta (RK) methods, now called symplectic Runge–Kutta (SRK) methods, provides a convenient way to achieve this [6, §VI.4].

It is a classical theorem due to Liouville that Hamiltonian systems are also volume preserving: for all bounded open sets  $\Omega$  of phase space, the flow map  $\varphi_t$  satisfies  $\text{vol}(\varphi_t(\Omega)) = \text{vol}(\Omega)$  for all  $t$ . Equivalently, the Jacobian determinant,  $\det(\varphi'_t(x))$ , is 1 for all  $x$  and  $t$  [6, VI.9]. Any symplectic mapping of phase space has this property, and therefore SRK methods are volume preserving for Hamiltonian systems. Beyond Hamiltonian systems, an ODE  $\dot{x} = f(x)$  is volume preserving if and only if  $f$  is divergence free (sometimes called source free). General volume preservation like this can be found in applications involving incompressible fluid flows and vorticities, ergodic theory and statistical mechanics, and problems in electromagnetism [5,7,12].

One can ask if any SRK methods are volume preserving for all divergence free vector fields  $f$ , and it has been known for 20 years that the answer is no. Feng and Shang showed that no RK method can be volume preserving even for the class of linear divergence free vector fields [5]. It was later shown by Iserles, Quispel and Tse and independently by Chartier and Murua that no B-Series method can be volume preserving for all divergence free vector fields [3,7]. However, Hairer, Lubich and Wanner have considered separable divergence free vector fields of the form

$$f(x, y) = (u(y), v(x))^{\top}, \quad (\text{HLW})$$

\* Corresponding author.

E-mail addresses: [p.bader@latrobe.edu.au](mailto:p.bader@latrobe.edu.au) (P. Bader), [d.mclaren@latrobe.edu.au](mailto:d.mclaren@latrobe.edu.au) (D.I. McLaren), [r.quispel@latrobe.edu.au](mailto:r.quispel@latrobe.edu.au) (G.R.W. Quispel), [m.d.webb@maths.cam.ac.uk](mailto:m.d.webb@maths.cam.ac.uk) (M. Webb).

<http://dx.doi.org/10.1016/j.apnum.2016.06.010>

0168-9274/© 2016 The Author(s). Published by Elsevier B.V. on behalf of IMACS. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

for functions  $u : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $v : \mathbb{R}^m \rightarrow \mathbb{R}^n$  [6, Thm. 9.4]. There the authors prove that any SRK method with at most two stages (and any composition of such methods) is volume preserving for these systems, giving a hint at the fact that SRK methods can be volume preserving for a much larger class of vector fields than just Hamiltonian systems.

As we will show in the introduction, vector fields  $f$  in that class must satisfy the determinant condition

$$\det \left( I + \frac{h}{2} f'(x) \right) = \det \left( I - \frac{h}{2} f'(x) \right) \text{ for all } h > 0, x \in \mathbb{R}^n, \tag{det}$$

where  $I$  denotes the  $n \times n$  identity matrix. In order to substantiate this claim and in anticipation of some of the results to be discussed later, we consider the following three Runge–Kutta methods  $x \mapsto \phi_h(x)$  which have been shown to preserve certain measures  $\mu(x)dx$  for quadratic Hamiltonian vector fields [2]:

1. The implicit midpoint rule

$$\frac{\phi_h(x) - x}{h} = f \left( \frac{\phi_h(x) + x}{2} \right), \text{ with } \mu(x) = 1,$$

for which the condition (det) was already studied in [11],

2. the trapezoidal rule

$$\frac{\phi_h(x) - x}{h} = \frac{1}{2} \left( f(x) + f(\phi_h(x)) \right), \text{ with } \mu(x) = \det \left( 1 - \frac{h}{2} f'(x) \right), \tag{1.1}$$

3. and Kahan’s method (restricted to quadratic vector fields)

$$\frac{\phi_h(x) - x}{h} = 2f \left( \frac{x + \phi_h(x)}{2} \right) - \frac{1}{2} f(x) - \frac{1}{2} f(\phi_h(x)), \text{ with } \mu(x) = \det \left( 1 - \frac{h}{2} f'(x) \right)^{-1}. \tag{1.2}$$

These quadratic Hamiltonian vector fields satisfy the determinant condition (det) and we will establish in section 5 that this condition is essential for these measure preservation properties. Indeed, using the chain rule, we compute the Jacobian matrix of the midpoint rule as

$$\phi'_h(x) = I + \frac{h}{2} f' \left( \frac{x + \phi_h(x)}{2} \right) \left( I + \phi'_h(x) \right),$$

which in turn gives the condition for volume preservation

$$\det(\phi'_h(x)) = \frac{\det(I + \frac{h}{2} f'((x + \phi_h(x))/2))}{\det(I - \frac{h}{2} f'((x + \phi_h(x))/2))} = 1.$$

Note that in agreement with [5], it is clear that for the implicit midpoint rule we cannot consider a class of vector fields any larger than this and realistically expect volume preservation. Hence we restrict our discussion to vector fields satisfying this determinant condition (det). These functions, as we show later, are divergence free and include Hamiltonian systems and HLW separable systems described above.

The contributions of this paper are to highlight the relevance of the determinant condition (det) for volume preservation by Runge–Kutta methods, and to introduce and prove results regarding volume preservation for some classes of vector fields lying between Hamiltonian vector fields and those satisfying the determinant condition (det). Not only does this further the understanding of Runge–Kutta methods and volume preservation of numerical methods in general, but it gives examples of where in applications one could in principle use Runge–Kutta methods and preserve volume for a non-Hamiltonian system. Furthermore, we discuss how Runge–Kutta methods can also preserve a modified measure exactly. The importance of such methods is that the dynamics of the numerical solution lie in the class of measure preserving systems, giving a qualitative advantage over methods lacking this property [9]. It should be noted that there are general approaches to constructing volume preserving splitting methods for a general divergence free vector field [5,6,12], but Runge–Kutta methods offer practical and theoretical simplicity and familiarity.

## 2. Properties of Runge–Kutta methods

This section is fairly technical, but it provides us with the necessary tools for the discussion in sections 3 and 4. We use the following notation to describe a Runge–Kutta method for the autonomous system  $\dot{x} = f(x)$ . We assume  $f$  is continuously differentiable throughout the paper. For each step-size  $h$ , a  $s$ -stage Runge–Kutta method provides a map  $\phi_h : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , defined by

$$\phi_h(x) = x + h \sum_{i=1}^s b_i f(k_i),$$

where the stages  $k_i$  satisfy

$$k_i = x + h \sum_{j=1}^s a_{ij} f(k_j), \text{ for } i = 1, \dots, s.$$

As usual, we consolidate the  $b_i$ 's and  $a_{ij}$ 's into the Butcher tableau consisting of the vector  $b$  and the matrix  $A$ . We make use of the Kronecker product throughout, which for  $A \in \mathbb{R}^{n \times n}$  and  $B \in \mathbb{R}^{m \times m}$  is defined to be

$$A \otimes B = \begin{pmatrix} a_{11}B & \cdots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{n1}B & \cdots & a_{nn}B \end{pmatrix} \in \mathbb{R}^{nm \times nm}.$$

**Lemma 2.1.** *The Jacobian matrix of a RK method can be written as*

$$\phi'_h(x) = I + h(b^\top \otimes I)F(I_s \otimes I - h(A \otimes I)F)^{-1}(\mathbb{1} \otimes I), \tag{2.1}$$

with determinant

$$\det(\phi'_h(x)) = \frac{\det(I_s \otimes I - h((A - \mathbb{1}b^\top) \otimes I)F)}{\det(I_s \otimes I - h(A \otimes I)F)}, \tag{2.2}$$

where  $F = \text{diag}(f'(k_1), \dots, f'(k_s))$ ,  $\mathbb{1}$  is an  $s \times 1$  vector of 1's and  $I_s$  is the  $s \times s$  identity matrix.

**Proof.** Computing directly, we find

$$\phi'_h(x) = I + h \sum_{i=1}^s b_i f'(k_i) k'_i(x) = I + h(b^\top \otimes I)F(k'_1(x), \dots, k'_s(x))^\top. \tag{2.3}$$

By definition of the stages  $k_i$ , the derivatives  $k'_i(x)$  satisfy

$$\begin{pmatrix} I - ha_{11}f'(k_1) & -ha_{12}f'(k_2) & \cdots & -ha_{1s}f'(k_s) \\ -ha_{21}f'(k_1) & I - ha_{22}f'(k_2) & \cdots & -ha_{2s}f'(k_s) \\ \vdots & \vdots & \ddots & \vdots \\ -ha_{s1}f'(k_1) & -ha_{s2}f'(k_2) & \cdots & I - ha_{ss}f'(k_s) \end{pmatrix} \begin{pmatrix} k'_1(x) \\ k'_2(x) \\ \vdots \\ k'_s(x) \end{pmatrix} = \begin{pmatrix} I \\ I \\ \vdots \\ I \end{pmatrix}.$$

Written more compactly using Kronecker products, this is

$$(I_s \otimes I - h(A \otimes I)F)(k'_1(x), \dots, k'_s(x))^\top = \mathbb{1} \otimes I. \tag{2.4}$$

The form of the Jacobian matrix can now be found by substituting (2.4) into (2.3).

For the determinant, use the block determinant identity

$$\det(U) \det(X - WU^{-1}V) = \det \begin{pmatrix} U & V \\ W & X \end{pmatrix} = \det(X) \det(U - VX^{-1}W) \tag{2.5}$$

on the expression (2.1) with  $U = I_s \otimes I - h(A \otimes I)F$ ,  $V = (\mathbb{1} \otimes I)$ ,  $W = -h(b^\top \otimes I)F$  and  $X = I$ .  $\square$

We wish to understand for which vector fields  $f$  and which Runge–Kutta methods defined by  $A$  and  $b$ , the determinant (2.2) is unity. As one might expect, this turns out to be simpler for symplectic Runge–Kutta methods. Now, for the purpose of exposition, we restrict to methods described in the following definition and instruct the reader in how certain results can be proven for general SRK methods at the end of the section.

**Definition 2.2.** A SRK method is said to be a *special symplectic Runge–Kutta method (SSRK)* if  $b_j \neq 0$  for all  $j$ , so that the Butcher tableau may be written  $A = \frac{1}{2}(\Omega + \mathbb{1}\mathbb{1}^\top)B$ , where  $B = \text{diag}(b)$  and  $\Omega$  is a skew-symmetric matrix.

This definition is reasonable because if  $b_j \neq 0$  for all  $j$ , then the matrix  $M = BA + A^\top B - bb^\top$  is zero (which implies the method is symplectic) if and only if  $\Omega$  is skew-symmetric. The expression  $\frac{1}{2}(\Omega + \mathbb{1}\mathbb{1}^\top)B$  therefore constitutes a normal form for most SRK methods of interest [6].

**Lemma 2.3.** *An  $s$ -stage SSRK method is volume preserving for  $\dot{x} = f(x)$  if and only if*

$$\det(I_s \otimes I - h(A \otimes I)F) = \det(I_s \otimes I + h(A^\top \otimes I)F), \tag{2.6}$$

where  $F = \text{diag}(f'(k_1), \dots, f'(k_s))$ .

**Proof.** The equation  $M = 0$  can be written  $-A^\top = B(A - \mathbb{1}b^\top)B^{-1}$ . Hence

$$\begin{aligned} \det(I_s \otimes I + h(A^\top \otimes I)F) &= \det(I_s \otimes I - h(B(A - \mathbb{1}b^\top)B^{-1} \otimes I)F) \\ &= \det(I_s \otimes I - h(B \otimes I)(A - \mathbb{1}b^\top) \otimes I)F(B \otimes I)^{-1}) \\ &= \det(I_s \otimes I - h((A - \mathbb{1}b^\top) \otimes I)F) \end{aligned}$$

The result now follows from Lemma 2.1.  $\square$

When  $s = 1$ , the only SSRK method is the implicit midpoint rule. In this case, Lemma 2.3 gives the determinant condition (det) from the introduction.

When  $s = 2$ , we have a three-parameter family of SSRK methods, which reduces to two-parameter if we impose the consistency condition  $b_1 + b_2 = 1$ . Now Lemma 2.3 gives the condition

$$\det \begin{pmatrix} I - ha_{11}f'(k_1) & -ha_{12}f'(k_2) \\ -ha_{21}f'(k_1) & I - ha_{22}f'(k_2) \end{pmatrix} = \det \begin{pmatrix} I + ha_{11}f'(k_1) & ha_{21}f'(k_2) \\ ha_{12}f'(k_1) & I + ha_{22}f'(k_2) \end{pmatrix}. \tag{2.7}$$

Applying the block determinant identity (2.5), this boils down to

$$\begin{aligned} \det(I - ha_{11}f'(k_1) - ha_{22}f'(k_2) + h^2 \det(A)f'(k_1)f'(k_2)) \\ = \det(I + ha_{11}f'(k_1) + ha_{22}f'(k_2) + h^2 \det(A)f'(k_1)f'(k_2)). \end{aligned} \tag{2.8}$$

We were able here to simplify the identity (2.5) because the top-left block  $(I - ha_{11}f'(k_1))$  and the bottom-left block  $(-ha_{21}f'(k_1))$  commute. This cannot be done for  $s \geq 3$ .

These next three lemmata give some basic operations that can be performed on the vector field which send volume preserving ODEs to volume preserving ODEs, and effect a simple change in the Jacobian determinant of some RK methods for general vector fields.

**Lemma 2.4.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and define a linear change of variables  $\tilde{f}(x) = Pf(P^{-1}x)$  for some invertible matrix  $P$ . Then the RK map  $\tilde{\phi}_h$  for solving  $\dot{x} = \tilde{f}(x)$  satisfies

$$\tilde{\phi}_h(x) = P\phi_h(P^{-1}x), \quad \tilde{\phi}'_h(x) = P\phi'_h(P^{-1}x)P^{-1}. \tag{2.9}$$

**Lemma 2.5.** Let  $u : \mathbb{R}^m \rightarrow \mathbb{R}^m, v : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^n$  and define  $f : \mathbb{R}^{m+n} \rightarrow \mathbb{R}^{m+n}$  by

$$f(x, y) = \begin{pmatrix} u(x) \\ v(x, y) \end{pmatrix}. \tag{2.10}$$

Now let  $\phi_h : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^{n+m}$  be a one-stage RK map for solving  $(\dot{x}, \dot{y})^\top = f(x, y)$ ,  $\psi_h : \mathbb{R}^m \rightarrow \mathbb{R}^m$  that for solving  $\dot{x} = u(x)$ , and  $\chi_h : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^n$  that for solving  $\dot{y} = v(x, y)$  where  $x$  is treated as a parameter. Then

$$\phi_h(x, y) = \begin{pmatrix} \psi_h(x) \\ \chi_h(k_1(x), y) \end{pmatrix}, \tag{2.11}$$

and consequently

$$\det(\phi'_h(x, y)) = \det(\psi'_h(x)) \det(\partial_y \chi_h(k_1(x), y)),$$

where  $k_1(x) = x + ha_{11}u(k_1(x))$  is the internal stage of the RK method  $\psi_h(x)$  and  $\partial_y$  denotes the derivative with respect to the  $y$  coordinate.

**Proof.** The full method is

$$\phi_h(x, y) = \begin{pmatrix} x \\ y \end{pmatrix} + hb_1 \begin{pmatrix} u(k_1(x)) \\ v(k_1(x), l_1(k_1(x), y)) \end{pmatrix}, \tag{2.12}$$

with the internal stages

$$\begin{pmatrix} k_1(x) \\ l_1(k_1(x), y) \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix} + ha_{11} \begin{pmatrix} u(k_1(x)) \\ v(k_1(x), l_1(k_1(x), y)) \end{pmatrix}.$$

The methods applied to each component of the  $\mathbb{R}^{n+m}$  dimensional system are given by

$$\begin{pmatrix} \psi_h(x) \\ \chi_h(x, y) \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix} + hb_1 \begin{pmatrix} u(k_1(x)) \\ v(x, l_1(x, y)) \end{pmatrix}. \tag{2.13}$$

Comparing (2.12) with (2.13) yields the result (2.11). To prove the last part, note that the Jacobian matrix has block structure

$$\phi'_h(x, y) = \begin{pmatrix} \psi'_h(x) & 0 \\ \partial_x(\chi_h(k_1(x), y)) & \partial_y(\chi_h(k_1(x), y)) \end{pmatrix}$$

and so the determinant  $\det(\phi'_h(x, y))$  is the product of the determinants of the diagonal blocks.  $\square$

For some simple vector fields, this can be generalised to certain  $s$ -stage methods. Note that the notation for  $\chi_h$  is different to that for Lemma 2.5.

**Lemma 2.6.** Let  $u : \mathbb{R}^m \rightarrow \mathbb{R}^m, v : \mathbb{R}^n \rightarrow \mathbb{R}^n, w : \mathbb{R}^m \rightarrow \mathbb{R}^n$  and define  $f : \mathbb{R}^{m+n} \rightarrow \mathbb{R}^{m+n}$  by

$$f(x, y) = \begin{pmatrix} u(x) \\ w(x) + v(y) \end{pmatrix}. \tag{2.14}$$

Now let  $\phi_h(x, y)$  be the RK map for solving  $(\dot{x}, \dot{y})^\top = f(x, y), \psi_h(x)$  that for solving  $\dot{x} = u(x)$ , and  $\chi_h(c, y)$  that for solving  $\dot{y} = c + v(y)$ . Define  $c_i = \sum_j a_{ij}$ . If the Butcher tableau is such that

$$\delta_j(i, k) = \frac{a_{ij} - a_{kj}}{c_i - c_k} \quad 1 \leq i, j, k \leq s, \tag{2.15}$$

is finite and independent of distinct  $i$  and  $k$  for each  $j$  then there exist functions  $d_h, e_h, c_h : \mathbb{R}^m \rightarrow \mathbb{R}^n$  such that

$$\phi_h(x, y) = \begin{pmatrix} \psi_h(x) \\ \chi_h(d_h(x), y + he_h(x)) + hc_h(x) \end{pmatrix} \text{ for all } y. \tag{2.16}$$

Consequently,

$$\det(\phi'_h(x, y)) = \det(\psi'_h(x)) \det(\partial_y \chi_h(d_h(x), y + he_h(x))). \tag{2.17}$$

**Proof.** Write  $\phi_h(x, y) = (\psi_h(x), \sigma_h(x, y))$ . Note that  $\sigma_h(x, y) \neq \chi_h(w(x), y)$ , but they are related as follows.

$$\sigma_h(x, y) = y + h \sum_{i=1}^s b_i w(k_i) + h \sum_{i=1}^s b_i v(l_i(w(k_1), \dots, w(k_s), y)), \tag{2.18}$$

$$\chi_h(c, y) = y + h \sum_{i=1}^s b_i c + h \sum_{i=1}^s b_i v(l_i(c, \dots, c, y)),$$

with stage values

$$l_i(\zeta_1, \dots, \zeta_s, y) = y + h \sum_{j=1}^s a_{ij} \zeta_j + h \sum_{j=1}^s a_{ij} v(l_j(\zeta_1, \dots, \zeta_s, y)).$$

Now let  $d$  be an arbitrary number. Then we have for each  $i$ ,

$$l_i(w(k_1), \dots, w(k_s), y) = y + he_i + h \sum_{j=1}^s a_{ij} d + h \sum_{j=1}^s a_{ij} v(l_j(w(k_1), \dots, w(k_s), y)),$$

where  $e_i = \sum_{j=1}^s a_{ij}(w(k_j) - d)$ . Hence

$$l_i(w(k_1), \dots, w(k_s), y) = l_i(d, \dots, d, y + he_i). \tag{2.19}$$

We want to choose  $d$  such that  $e_i = e_k \forall i, k$ . Equivalently,

$$\sum_{j=1}^s a_{ij}(w(k_j) - d) = \sum_{j=1}^s a_{kj}(w(k_j) - d) \text{ for all } i \neq k.$$

Solving for  $d$  we find

$$d = \sum_{j=1}^s w(k_j) \left( \frac{a_{ij} - a_{kj}}{c_i - c_k} \right) \text{ for all } i \neq k.$$

This will only give us a unique finite value of  $d$  no matter what values  $w(k_i)$  take if the value of  $\delta_j(i, k)$  is finite and independent of distinct  $i$  and  $k$  for every  $j$ , which is given by assumption. Hence we can set  $d_h(x) = d$ ,  $e_h(x) = e_1$  and by (2.19), we write (2.18) as

$$\begin{aligned} \sigma_h(x, y) &= y + h \sum_{i=1}^s b_i w(k_i) + h \sum_{i=1}^s b_i v(l_i(d_h, \dots, d_h, y + he_h)) \\ &= (y + he_h) + h \sum_{i=1}^s b_i d_h + h \sum_{i=1}^s b_i v(l_i(d_h, \dots, d_h, y + he_h)) + h \left( \sum_{i=1}^s b_i (w(k_i) - d_h) - e_h \right) \\ &= \chi_h(d_h(x), y + he_h(x)) + hc_h(x), \end{aligned}$$

where  $c_h(x) = (\sum_{i=1}^s b_i (w(k_i) - d_h) - e_h)$ . The factorisation of the determinant is evident from the block structure of the Jacobian matrix

$$\phi'_h(x, y) = \begin{pmatrix} \psi'_h(x) & 0 \\ \star & \partial_y(\chi_h(d_h(x), y + he_h(x)) + hc_h(x)) \end{pmatrix}. \quad \square$$

**Remark 2.7.** Let us shed some light on the meaning of (2.15) being finite and independent of distinct  $i$  and  $k$  for each  $j$ . The finiteness implies that the method has  $c_i \neq c_k$  for all  $i \neq k$ , which is known as *nonconfluency* [6]. For two-stage SSRK methods, condition (2.15) is satisfied if the method is consistent and  $\Omega \neq 0$ . There is a one-parameter family of self-adjoint three-stage SSRK methods of order four that satisfy the condition. The three-stage Gauss–Legendre method, however, does not belong to this class.

**Definition 2.8.** A vector field  $f : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^{n+m}$  possesses a *linear foliation* if there exists a linear change of variables as in Lemma 2.4 such that  $f$  is as in (2.10) from Lemma 2.5 for some functions  $u$  and  $v$ . Such vector fields are called *linearly foliate*. See [8] for general Lie group foliations in the context of Geometric Integration.

**Remark 2.9.** For general SRK methods, the condition in Lemma 2.3 along with the condition with  $A$  replaced by  $A - \mathbb{1}b^\top$  is sufficient for volume preservation. This result can be obtained along the lines of [6, Thm. 9.4] regarding separable systems (HLW), as follows. Consider the foliation  $\dot{x} = f(x)$ ,  $\dot{y} = -f'(x)^\top y$ , which is Hamiltonian with respect to  $H(x, y) = y^\top f(x)$ . Then, using the notation of Lemma 2.6, the Jacobian matrix of the Runge–Kutta map has block structure  $\begin{pmatrix} \phi'_h(x) & 0 \\ \star & \partial_y \sigma_h(x, y) \end{pmatrix}$ . As in [6, Thm. 9.4], since the vector field is Hamiltonian, a SRK method will produce a symplectic map, which implies  $\det(\phi'_h(x)) \det(\partial_y \sigma_h(x, y)) = 1$ . Hence to show that  $\det(\phi'_h(x)) = 1$  it suffices to show that  $\det(\phi'_h(x)) = \det(\partial_y \sigma_h(x, y))$ . Computing these two sides as in Lemma 2.1, using the block determinant relation and equating numerators and denominators, gives the 2 conditions mentioned above.

### 3. Classification of volume preserving vector fields

**Definition 3.1.** We define the following classes of vector fields on Euclidean space using vector fields  $f(x, y) = (u(x), v(x, y))^\top$  possessing linear foliations as in Definition 2.8. The classes  $\mathcal{F}^{(\infty)}$  and  $\mathcal{F}^{(2)}$  are defined recursively – for complete rigour, an inductive construction, beginning from the trivial base class containing only zero vector fields, can be easily performed.

$$\begin{aligned} \mathcal{H} &= \left\{ f \text{ such that there exists } P \text{ such that for all } x, Pf'(x)P^{-1} = -f'(x)^\top \right\}, \\ \mathcal{S} &= \left\{ f \text{ such that there exists } P \text{ such that for all } x, Pf'(x)P^{-1} = -f'(x) \right\}, \\ \mathcal{F}^{(\infty)} &= \left\{ f(x, y) = (u(x), v(x, y))^\top \text{ where } u \in \mathcal{H} \cup \mathcal{F}^{(\infty)} \text{ and there exists } P \text{ such that for all } x, y \right. \\ &\quad \left. P \partial_y v(x, y) P^{-1} = -\partial_y v(x, y)^\top \right\}, \\ \mathcal{F}^{(2)} &= \left\{ f(x, y) = (u(x), v(x, y))^\top \text{ where } u \in \mathcal{S} \cup \mathcal{H} \cup \mathcal{F}^{(2)} \text{ and there exists } P \text{ such that for all } x, y \right. \\ &\quad \left. \text{either } P \partial_y v(x, y) P^{-1} = -\partial_y v(x, y)^\top \text{ or } P \partial_y v(x, y) P^{-1} = -\partial_y v(x, y) \right\}, \\ \mathcal{D} &= \left\{ \text{vector fields satisfying } \det\left(I + \frac{h}{2} f'(x)\right) = \det\left(I - \frac{h}{2} f'(x)\right) \text{ for all } h > 0 \text{ and all } x \right\}. \end{aligned}$$

**Lemma 3.2.** The set  $\mathcal{H}$  contains all vector fields of the form  $f(x) = J^{-1} \nabla H(x)$  where  $J$  is constant and skew-symmetric. All SRK methods are volume preserving for vector fields in  $\mathcal{H}$ .

**Proof.** For the first part, note that if  $f(x) = J^{-1}\nabla H(x)$ , then  $Jf'(x)J^{-1} = \nabla^2 H(x)J^{-1} = -f'(x)^\top$ . For the second part, let  $A \in \mathbb{R}^{s \times s}$  and  $P$  be such that for all  $x$ ,  $Pf'(x)P^{-1} = -f'(x)^\top$ . Then using the notation of Lemma 2.3,

$$\det(I_s \otimes I - h(A \otimes I)F) = \det(I_s \otimes I - h(I_s \otimes P)(A \otimes I)(I_s \otimes P^{-1})(I_s \otimes P)F(I_s \otimes P^{-1})) \tag{3.1}$$

$$= \det(I_s \otimes I + h(A \otimes I)F^\top) \tag{3.2}$$

$$= \det(I_s \otimes I + hF(A^\top \otimes I) \text{ (transpose)}) \tag{3.3}$$

$$= \det(I_s \otimes I + h(A^\top \otimes I)F) \text{ (Sylvester's law)}. \tag{3.4}$$

By Lemma 2.3 and Remark 2.9, all SRK methods are volume preserving.  $\square$

**Lemma 3.3.** *The set  $\mathcal{S}$  contains all separable HLW systems i.e.  $f(x, y) = (u(y), v(x))^\top$ . All SRK methods with at most 2 stages, and compositions thereof, are volume preserving for vector fields in  $\mathcal{S}$ .*

**Proof.** For the first part, note that if  $f(x, y) = (u(y), v(x))^\top$ , then  $Df'(x, y)D^{-1} = -f'(x, y)$  where  $D = \text{diag}(I_m, -I_n)$ . For the second part, let  $A \in \mathbb{R}^{2 \times 2}$  and  $P$  be such that for all  $x$   $Pf'(x)P^{-1} = -f'(x)$ . Then for the two stages  $k_1, k_2$  of the SRK method,

$$\begin{aligned} \det(I - ha_{11}f'(k_1) - ha_{22}f'(k_2) + h^2 \det(A)f'(k_1)f'(k_2)) \\ = \det(I - ha_{11}Pf'(k_1)P^{-1} - ha_{22}Pf'(k_2)P^{-1} + h^2 \det(A)Pf'(k_1)P^{-1}Pf'(k_2)P^{-1}) \\ = \det(I + ha_{11}f'(k_1) + ha_{22}f'(k_2) + h^2 \det(A)f'(k_1)f'(k_2)). \end{aligned}$$

By (2.8) and Remark 2.9, all 2-stage SRK methods are volume preserving. To complete the proof, note that a 1-stage SRK method is equivalent to a 2-stage SRK method with two equal stages, and compositions of volume preserving maps are also volume preserving.  $\square$

**Lemma 3.4.** *The inclusions  $\mathcal{H} \subset \mathcal{F}^{(\infty)} \subset \mathcal{F}^{(2)} \subset \mathcal{D}$  and  $\mathcal{S} \subset \mathcal{F}^{(2)} \subset \mathcal{D}$  hold.*

**Proof.**  $\mathcal{H} \subset \mathcal{F}^{(\infty)} \subset \mathcal{F}^{(2)}$  and  $\mathcal{S} \subset \mathcal{F}^{(2)}$  are clear by considering trivial foliations in which  $n + m = m$ . We will show that  $\mathcal{S} \subset \mathcal{D}$ ,  $\mathcal{H} \subset \mathcal{D}$  and that  $\mathcal{D}$  is closed under the employed recursive process leading to linearly foliate systems.

For  $f \in \mathcal{S}$ ,  $\det(I + \frac{h}{2}f'(x)) = \det(I + \frac{h}{2}Pf'(x)P^{-1}) = \det(I - \frac{h}{2}f'(x))$ .

For  $f \in \mathcal{H}$ ,  $\det(I + \frac{h}{2}f'(x)) = \det(I + \frac{h}{2}Pf'(x)P^{-1}) = \det(I - \frac{h}{2}f'(x)^\top) = \det(I - \frac{h}{2}f'(x))$ .

Let  $f \in \mathcal{D}$  and define  $\tilde{f}(x) = Pf(P^{-1}x)$  for an invertible matrix  $P$ . Then  $\det(I + \frac{h}{2}\tilde{f}(x)) = \det(I + \frac{h}{2}Pf'(P^{-1}x)P^{-1}) = \det(I + \frac{h}{2}f'(P^{-1}x))$ . Doing the same with a  $-$  instead of a  $+$  shows that  $\tilde{f} \in \mathcal{D}$ .

Let  $f(x, y) = (u(x), v(x, y))^\top$  where  $u \in \mathcal{D}$  and  $y \mapsto v(x, y) \in \mathcal{D}$  for all  $x$ . Then

$$\det(I + \frac{h}{2}f'(x, y)) = \det \begin{pmatrix} I + \frac{h}{2}u'(x) & 0 \\ \frac{h}{2}\partial_x v(x, y) & I + \frac{h}{2}\partial_y v(x, y) \end{pmatrix} \tag{3.5}$$

$$= \det(I + \frac{h}{2}u'(x)) \det(I + \frac{h}{2}\partial_y v(x, y)). \tag{3.6}$$

Doing the same with a  $-$  instead of a  $+$  shows that  $f \in \mathcal{D}$ .  $\square$

**Theorem 3.5.** *The following are equivalent.*

- (i)  $f \in \mathcal{D}$
- (ii)  $\det(I + zf'(x)) = \det(I - zf'(x))$  for all  $z \in \mathbb{C}$  and all  $x$
- (iii) The non-zero eigenvalues of  $f'(x)$ , counting multiplicities, come in positive-negative pairs
- (iv)  $\text{tr}(f'(x)^{2k+1}) = 0$  for all  $x$  and  $k = 0, 1, 2, \dots$

**Proof.** (i)  $\iff$  (ii): Assuming (i), for every  $x$ ,  $p(z) = \det(I + zf'(x)) - \det(I - zf'(x))$  is a polynomial in  $z$  that is zero for infinitely many values of  $z = h/2 \in \mathbb{R}_+$ . By the Fundamental Theorem of Algebra,  $p(z) = 0$  for all  $z \in \mathbb{C}$ . The converse follows from setting  $h = 2z \in \mathbb{R}_+ \subset \mathbb{C}$ .

(ii)  $\implies$  (iii): By triangularisation we can see that for every  $x$ , the polynomial  $q(z) = \det(I - zf'(x))$  is equal to  $(1 - z\lambda_1) \cdots (1 - z\lambda_r)$  where  $\lambda_1, \dots, \lambda_r$  are the non-zero eigenvalues of  $f'(x)$ . If (ii) holds, then  $q(z) = q(-z)$ , and the roots  $1/\lambda_i$  of  $q$  come in positive-negative pairs. Hence the eigenvalues  $\lambda_i$  do too.

(iii)  $\implies$  (iv): For all  $x$ ,  $\text{tr}(f'(x)^{2k+1}) = \lambda_1^{2k+1} + \dots + \lambda_r^{2k+1}$  where  $\lambda_1, \dots, \lambda_r$  are the non-zero eigenvalues of  $f'(x)$ . Hence if the non-zero eigenvalues come in positive-negative pairs then  $\text{tr}(f'(x)^{2k+1}) = 0$  for  $k = 0, 1, 2, \dots$

(iv)  $\implies$  (ii): Newton’s identity gives

$$e_{2k+1}(\lambda_1, \dots, \lambda_n) = \frac{1}{2k+1} \sum_{i=1}^{2k+1} (-1)^{i-1} e_{2k+1-i}(\lambda_1, \dots, \lambda_n) \text{tr}(f'(x)^i), \tag{3.7}$$

where  $e_j(\lambda_1, \dots, \lambda_n)$  is the elementary symmetric polynomial in  $\lambda_1, \dots, \lambda_n$  and, incidentally, the coefficient of  $z^j$  in  $q(z) = \det(I - zf'(x))$ . Since for any  $k$  and  $i$ , either  $2k + 1 - i$  is odd or  $i$  is odd, we can use an induction argument to show that all the coefficients of  $z^{2k+1}$  in  $q(z)$  are zero. Hence  $q(z) = q(-z)$ .  $\square$

**Corollary 3.6.** All elements of  $\mathcal{D}$  are divergence free. Restricted to 2 dimensional vector fields,  $\mathcal{D}$ ,  $\mathcal{H}$  and  $\mathcal{S}$  are all equal to divergence free vector fields.

**Theorem 3.7.** The set  $\mathcal{F}^{(\infty)}$  contains all

- (i) Affine vector fields  $f(x) = Lx + d$  such that  $\det(I + \frac{h}{2}L) = \det(I - \frac{h}{2}L)$  for all  $h > 0$
- (ii) Vector fields such that  $f'(x) = JS(x)$  where  $J$  is skew-symmetric and  $S(x)$  is symmetric

**Proof.** (i) Let  $L$  satisfy the determinant condition (det). By the Jordan normal form, and the fact that the eigenvalues must come in positive-negative pairs by Theorem 3.5, we can find an invertible matrix  $P$  such that

$$PLP^{-1} = \text{diag}(\lambda_1 I + N_1, -\lambda_1 I + N_{-1}, \lambda_2 I + N_2, -\lambda_2 I + N_{-2}, \dots, \lambda_r I + N_r, -\lambda_r I + N_{-r}, N_0),$$

where the  $N_k$  are matrices that are zero everywhere except for possible 1’s on the first subdiagonal  $(N_k)_{i+1,i}$ . Hence  $f$  is a tower of linear foliations of affine functions with Jacobian matrices either  $N_0$  or  $\text{diag}(\lambda I + N_1, -\lambda I + N_{-1})$ . If  $f(x) = N_0 x + d$  then this is clearly a tower of foliations of zero systems i.e.  $u(x) = 0, v(x, y) = x$ . Now consider the case  $f(x) = \text{diag}(\lambda I + N_1, -\lambda I + N_{-1})x + d$ . There is a simple permutation of variables so that the Jacobian matrix becomes

$$\begin{pmatrix} \lambda & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -\lambda & 0 & 0 & 0 & 0 & 0 \\ \star & 0 & \lambda & 0 & 0 & 0 & 0 \\ 0 & \star & 0 & -\lambda & 0 & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \star & 0 & \lambda & 0 \\ 0 & 0 & \dots & 0 & \star & 0 & -\lambda \end{pmatrix},$$

where the  $\star$ ’s are possible 1’s (0 otherwise). Hence  $f$  is a tower of linear foliations of harmonic oscillators,  $u(x_1, x_2) = (\lambda x_1, -\lambda x_2), v(x_1, x_2, y_1, y_2) = (\star x_1 + \lambda y_1, \star x_2 - \lambda y_2)$ .

(ii) By a linear orthogonal change of variables, we can assume  $J = \text{diag}(0, K^{-1})$ , where  $K$  is skew-symmetric. In this case there is symmetric  $T(x)$  and  $V(x)$  such that

$$f'(x) = \begin{pmatrix} 0 & 0 \\ 0 & K^{-1} \end{pmatrix} \begin{pmatrix} T(x) & U(x) \\ U(x)^\top & V(x) \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ K^{-1}U(x)^\top & K^{-1}V(x) \end{pmatrix}. \tag{3.8}$$

This shows that  $f$  possesses a linear foliation with a zero system  $u \in \mathcal{H}$  and a system  $v$  with  $\partial_y v(x, y) = K^{-1}V(x)$  so that  $y \mapsto v(x, y) \in \mathcal{H}$  with the same  $P = K$  for all  $x, y$ .  $\square$

**Theorem 3.8.** Consider an  $s$ -stage SRK method that is volume preserving for the vector field  $u : \mathbb{R}^m \rightarrow \mathbb{R}^m$ , and let  $v : \mathbb{R}^{m+n} \rightarrow \mathbb{R}^{m+n}$  be such that there exists an invertible matrix  $P$  such that for all  $x, y$ ,

$$P \partial_y v(x, y) P^{-1} = -\partial_y v(x, y)^\top.$$

Then the SRK method is volume preserving for the vector field

$$f(x, y) = (u(x), v(x, y))^\top. \tag{3.9}$$

**Proof.** Let  $A \in \mathbb{R}^{s \times s}$  and take  $P$  from the assumption. By Lemma 2.3 and Remark 2.9, a SRK method is volume preserving if

$$\det(I_s \otimes I - h(A \otimes I)F) = \det(I_s \otimes I + h(A^\top \otimes I)F), \tag{2.6}$$

where  $F = \text{diag}(f'(k_1), \dots, f'(k_s))$ . For (3.9), the Jacobian matrix becomes

$$f'(x, y) = \begin{pmatrix} u'(x) & 0 \\ \star & \partial_y v(x, y) \end{pmatrix}$$



and using a similarity transformation, we can bring  $\det(I \otimes I - h(A \otimes I)F)$  to the form

$$\det \begin{pmatrix} I - a_{11}u'_1 & \cdots & -a_{1s}u'_s & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ -a_{s1}u'_1 & \cdots & I - a_{ss}u'_s & 0 & \cdots & 0 \\ \star & \cdots & \star & I - a_{11}v'_1 & \cdots & a_{1s}v'_s \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \star & \cdots & \star & I - a_{s1}v'_1 & \cdots & a_{ss}v'_s \end{pmatrix}, \tag{3.10}$$

where  $u'_i, v'_i$  are shorthand for  $\partial_x u(k_i)$  and  $\partial_y v(k_i)$ , respectively. Thus, the condition (2.6) factorises to

$$\det(I_S \otimes I - h(A \otimes I)U) \det(I_S \otimes I - h(A \otimes I)V) = \det(I_S \otimes I + h(A^\top \otimes I)U) \det(I_S \otimes I + h(A^\top \otimes I)V),$$

with  $U = \text{diag}(u'_1, \dots, u'_s)$ ,  $V = \text{diag}(v'_1, \dots, v'_s)$ . We compute

$$\begin{aligned} \det(I_S \otimes I - h(A \otimes I)V) &= \det(I_S \otimes I - h(I_S \otimes P)(A \otimes I)V(I_S \otimes P)^{-1}) \\ &= \det(I_S \otimes I - h(A \otimes I)(I_S \otimes P)V(I_S \otimes P^{-1})) \\ &= \det(I_S \otimes I + h(A \otimes I)V^\top) \\ &= \det(I_S \otimes I + hV(A^\top \otimes I)) \\ &= \det(I_S \otimes I + h(A^\top \otimes I)V). \end{aligned}$$

The last line comes from Sylvester's determinant identity. The proof is completed noticing that  $\det(I_S \otimes I - h(A \otimes I)U) = \det(I_S \otimes I + h(A^\top \otimes I)U)$  is satisfied by the assumption that the method is volume preserving for  $u$ .  $\square$

**Corollary 3.9.** All SRK methods are volume preserving for vector fields in  $\mathcal{F}^{(\infty)}$ .

**Proof.** SRK methods are volume preserving for vector fields in  $\mathcal{H}$  by Lemma 3.2 and volume preservation for the recursive constructions of  $\mathcal{F}^{(\infty)}$  is assured by Theorem 3.8.  $\square$

**Theorem 3.10.** Consider a SRK method with at most two stages (or a composition of such methods) that is volume preserving for the vector field  $u : \mathbb{R}^m \rightarrow \mathbb{R}^m$ , and let  $v : \mathbb{R}^{m+n} \rightarrow \mathbb{R}^{m+n}$  be such that there exists an invertible matrix  $P$  such that for all  $x, y$ ,

$$P \partial_y v(x, y) P^{-1} = -\partial_y v(x, y).$$

Then the SRK method is volume preserving for the vector field

$$f(x, y) = (u(x), v(x, y))^\top.$$

**Proof.** Let  $A \in \mathbb{R}^{2 \times 2}$  and take  $P$  from assumption. As in Theorem 3.8, the Jacobian matrix is block triangular

$$f'(x, y) = \begin{pmatrix} u'(x) & 0 \\ \partial_x v(x, y) & \partial_y v(x, y) \end{pmatrix}.$$

For 2-stage methods, the condition for volume preservation from equation (2.8) is

$$\begin{aligned} \det(I - ha_{11}f'(k_1) - ha_{22}f'(k_2) + h^2 \det(A)f'(k_1)f'(k_2)) \\ = \det(I + ha_{11}f'(k_1) + ha_{22}f'(k_2) + h^2 \det(A)f'(k_1)f'(k_2)). \end{aligned}$$

Now, because of the block-triangular structure of  $f'(k_i)$  and

$$f'(k_1)f'(k_2) = \begin{pmatrix} u'(k_1)u'(k_2) & 0 \\ \star & \partial_y v(k_1)\partial_y v(k_2) \end{pmatrix}, \quad f'(k_1) + f'(k_2) = \begin{pmatrix} u'(k_1) + u'(k_2) & 0 \\ \star & \partial_y v(k_1) + \partial_y v(k_2) \end{pmatrix},$$

where we have used the convention that  $u(k_i)$  has used the  $x$  component of  $k_i$ . The condition (2.8) then factorises into

$$\begin{aligned} \det(I - h(a_{11}f'(k_1) + a_{22}f'(k_2)) + h^2 \det(A)f'(k_1)f'(k_2)) \\ = \det(I - h(a_{11}u'(k_1) + a_{22}u'(k_2)) + h^2 \det(A)u'(k_1)u'(k_2)) \\ \cdot \det(I - h(a_{11}\partial_y v(k_1) + a_{22}\partial_y v(k_2)) + h^2 \det(A)\partial_y v(k_1)\partial_y v(k_2)). \end{aligned}$$

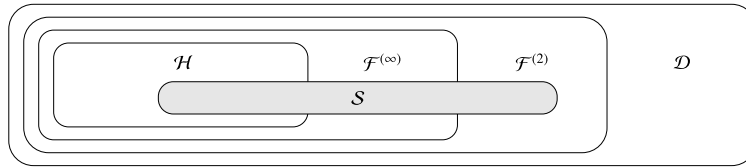


Fig. 1. Venn diagram illustrating the relationships established by Lemma 3.4.

A similarity transformation with  $P$  leads to

$$\begin{aligned} & \det(I - h(a_{11}\partial_y v(k_1) + a_{22}\partial_y v(k_2)) + h^2 \det(A)\partial_y v(k_1)\partial_y v(k_2)) \\ &= \det(I - h(a_{11}P\partial_y v(k_1)P^{-1} + a_{22}P\partial_y v(k_2)P^{-1}) + h^2 \det(A)P\partial_y v(k_1)P^{-1}P\partial_y v(k_2)P^{-1}) \\ &= \det(I + h(a_{11}\partial_y v(k_1) + a_{22}\partial_y v(k_2)) + h^2 \det(A)\partial_y v(k_1)\partial_y v(k_2)). \end{aligned} \tag{3.11}$$

Condition (2.8) for  $f$  is now satisfied by considering (2.8) for the vector field  $u$  (which holds because we assume the SRK method is volume preserving) and (3.11). This proves the result for 2-stage SRK methods. To complete the proof, note that a 1-stage SRK method is equivalent to a 2-stage SRK method with two equal stages, and compositions of volume preserving maps are also volume preserving.  $\square$

**Corollary 3.11.** All SRK methods with at most two stages (and compositions thereof) are volume preserving for vector fields in  $\mathcal{F}^{(2)}$ .

**Proof.** All SRK methods with at most two stages (and compositions thereof) are volume preserving for vector fields in  $\mathcal{H}$  by Lemma 3.2 and  $\mathcal{S}$  by Lemma 3.3. Volume preservation for the recursive construction of  $\mathcal{F}^{(2)}$  is assured by Theorems 3.8 and 3.10.  $\square$

We already saw in the introduction that the implicit midpoint rule (which is the only 1-stage SRK method) is volume preserving for all  $f \in \mathcal{D}$ , and that all such vector fields must lie in  $\mathcal{D}$ . However, does the set  $\mathcal{F}^{(2)}$  contain all vector fields such that all 2-stage SRK methods are volume preserving? And does the set  $\mathcal{F}^{(infty)}$  contain all vector fields such that all SRK methods are volume preserving? We do not yet know the answers to these questions, but will illustrate applications that lie within these sets of vector fields and relevant counterexamples in the following sections.

#### 4. Examples and counterexamples

In this section, we present vector fields from different intersections of the Venn diagram in Fig. 1. The first counterexample shows that Corollary 3.11 is not true for three-stage methods. In the second example, we show that  $\mathcal{D} \setminus \mathcal{F}^{(2)} \neq \emptyset$  and that only the midpoint rule can be volume preserving for all methods in  $\mathcal{D}$ . This counterexample is of the lowest possible dimension (3) but one might argue that volume preservation is hindered in this example because the vector field is not completely smooth at  $x = 0$ . The third example clarifies the matter: we give a way to construct a class of (smooth) vector fields in  $\mathcal{D}$  for which two-stage methods cannot be expected to preserve volume.

**Counterexample 4.1.** Hairer, Lubich and Wanner [6, VI.9] used the vector field

$$\dot{x} = \sin z, \quad \dot{y} = \cos z, \quad \dot{z} = \sin y + \cos x,$$

to show that the three-stage Gauss–Legendre method is not volume preserving, despite the vector field lying in  $\mathcal{S}$ . What could be interesting is to find some class of functions  $\mathcal{F}^{(3)}$  – if it exists – such that all three-stage SRK methods are volume preserving, but not all four-stage SRK methods.

**Counterexample 4.2.** Consider the continuously differentiable vector field

$$f(x, y, z) = \begin{cases} (\frac{1}{3}x^3 - c, -x^2y, 0) & \text{if } x \geq 0 \\ (\frac{1}{3}x^3 - c, 0, -x^2z) & \text{if } x < 0 \end{cases},$$

$$f'(x, y, z) = \begin{pmatrix} x^2 & 0 & 0 \\ -2xy & -x^2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \text{ if } x \geq 0, \quad \begin{pmatrix} x^2 & 0 & 0 \\ 0 & 0 & 0 \\ -2xz & 0 & -x^2 \end{pmatrix} \text{ if } x < 0.$$

Then  $f \in \mathcal{D}$ , but not all 2-stage SRK methods are volume preserving. The principle here is that if  $k_1$  and  $k_2$  have  $x$ -components with different signs, then  $f'(k_1)$  and  $f'(k_2)$  will violate the condition in (2.8). For instance, the two-stage Gauss–Legendre method with initial value  $(1/2, 0, 0)$ , drift  $c = 1$  and step size  $h = 1/2$  has stage values with different signs in the  $x$ -coordinate and hence, does not preserve volume.

**Counterexample 4.3.** The following example illustrates that SRK methods cannot preserve simple systems in  $\mathcal{D}$  that do not belong to the classes  $\mathcal{F}^{(\infty)}$  or  $\mathcal{F}^{(2)}$ . Let  $g \in \mathcal{D}$  and let  $A(x)$  be skew-symmetric (and invertible) and  $S(y)$  be symmetric matrices. Then, any vector field with Jacobian matrix

$$f'(x, y) = \begin{pmatrix} g'(x) & 0 \\ \star & A(x)S(y) \end{pmatrix}$$

will satisfy the determinant condition, however, the similarity transform  $P$  to yield  $P\partial_y f(x, y)P^{-1} = -\partial_y f(x, y)^\top$  is now  $P = A(x)^{-1}$  and this dependence on  $x$  hinders a crucial step in the above proof. For volume preservation of SRK methods, it is thus essential to have a constant transform  $P$  for all values of  $x, y$  or at least in a region of interest for the numerical integration. We give the following concrete example,

$$A(x) = \begin{pmatrix} 0 & x_1 & x_1 \\ -x_1 & 0 & x_1x_2 \\ -x_1 & -x_1x_2 & 0 \end{pmatrix}, \quad S(y) = \begin{pmatrix} y_1^2 & 0 & 0 \\ 0 & y_2^2 & 0 \\ 0 & 0 & y_3 \end{pmatrix},$$

which is combined with the harmonic oscillator  $g(x_1, x_2) = (x_2, -x_1)^\top$  and could originate from  $f(x, y) = (g(x), A(x)(\frac{1}{3}y_1^3, \frac{1}{3}y_2^3, \frac{1}{2}y_1^2)^\top)^\top$ . Integrating with step size  $h = 1/2$  from  $(x_0, y_0) = (1, 1/2, 1/3, 1/4, 1/5)$  leads to a change of volume for the two-stage Gauss–Legendre method. The implicit midpoint rule preserves volume as expected.

The following examples come from applications and show the richness of the sets  $\mathcal{F}^{(2)}$  and  $\mathcal{F}^{(\infty)}$  in comparison with the previously known Hamiltonian or separable case. Furthermore, foliations give rise to new splitting integrators that preserve volume.

**Example 4.4** (Volume preserving splitting using Runge–Kutta methods). Consider the ODE in  $\mathbb{R}^{m+n}$  with

$$\dot{x} = F(x) + G(y), \quad \dot{y} = H(x, y), \tag{4.1}$$

with Jacobian

$$f'(x, y) = \begin{pmatrix} F'(x) & G'(y) \\ \partial_x H(x, y) & \partial_y H(x, y) \end{pmatrix}.$$

Then, a splitting into the vector fields

$$(A) : \begin{cases} \dot{x} = F(x), \\ \dot{y} = H(x, y), \end{cases} \quad \text{and} \quad (B) : \begin{cases} \dot{x} = G(y), \\ \dot{y} = 0, \end{cases}$$

with corresponding flows  $\varphi_h^{(A)}, \varphi_h^{(B)}$ , yields a block-tridiagonal Jacobian for system (A). If  $f$  and  $h$  lie in  $\mathcal{F}^{(2)}$  (or  $\mathcal{F}^{(\infty)}$ ), so will system (A). Part (B) corresponds to a trivial shift and substituting  $\varphi_h^{(A)}$  with a symplectic two-stage method  $\phi^{(A)} + h$  (any stage for  $\mathcal{F}^{(\infty)}$ ) gives a volume preserving splitting integrator  $\prod_i \phi_{a_i h}^{(A)} \circ \varphi_{b_i h}^{(B)}$ . An example of such a system (4.1) is the generalised Arnold–Beltrami–Childress (ABC) flow,

$$\dot{x} = A(y, z), \quad \dot{y} = B(x, z), \quad \dot{z} = C(x, y),$$

in the case where  $A$  and  $B$  do not depend on  $z$ .

**Example 4.5** (ABC flow). The following realisation of the ABC flow [4]

$$\dot{x} = A \sin(z) + C \cos(y), \quad \dot{y} = B \sin(x) + A \cos(z), \quad \dot{z} = C \sin(y) + B \cos(x)$$

with  $C = 0$  lies in  $\mathcal{F}^{(2)}$  and is of type (4.1).

**Example 4.6** (Lotka–Volterra). Another well-known example is the Lotka–Volterra systems in biology and economics [13] in the general form

$$\dot{x}_i = x_i \left( \lambda_i + \sum_{j=1}^n a_{ij} x_j \right),$$

which can be written in new coordinates  $u_i = \log x_i$  for positive values  $x_i > 0$  and in matrix form as

$$\dot{u} = \lambda + Ae^u.$$

If  $\lambda \in \text{range}(A)$ , then McLachlan et al. [10] showed that we can rewrite the system in gradient form

$$\dot{u} = A \nabla V(u), \quad V(u) = \sum_i e^{u_i} + \lambda_i c_i,$$

with  $\sum_j a_{ij} c_j = \lambda_i$ . If  $A$  is skew-symmetric, this corresponds to the Hamiltonian case  $\mathcal{H}$ , however, if  $A$  is of block-triangular form,

$$A = \begin{pmatrix} J_1 & 0 \\ \star & J_2 \end{pmatrix}$$

for any skew-symmetric invertible matrices  $J_1, J_2$  and an arbitrary block  $\star$ , the system belongs to  $\mathcal{F}^{(\infty)}$ .

**Example 4.7 (Skew-persymmetric flow).** Consider the following class of divergence free ODEs,

$$\begin{aligned} \dot{x} &= F(x-z) - Ay \\ \dot{y} &= Az - Bx \\ \dot{z} &= F(x-z) + By \end{aligned} \quad f'(x, y, z) = \begin{pmatrix} F'(x-z) & -A & -F'(x-z) \\ -B & 0 & A \\ F'(x-z) & B & -F'(x-z) \end{pmatrix}.$$

The Jacobian satisfies  $P f'(x, y, z) P^{-1} = -f'(x, y, z)^\top$  for the invertible matrix

$$P = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}. \tag{4.2}$$

Hence any symplectic Runge–Kutta method will be volume preserving for such ODEs, which may not be obvious at first sight.

**Example 4.8 (Skew-centrosymmetric flow).** Consider the following class of divergence free ODEs,

$$\begin{aligned} \dot{x} &= F(x-z) + G(y) \\ \dot{y} &= H(z-x) \\ \dot{z} &= F(x-z) - G(y) \end{aligned} \quad f'(x, y, z) = \begin{pmatrix} F'(x-z) & G'(y) & -F'(x-z) \\ -H'(z-x) & 0 & H'(z-x) \\ F'(x-z) & -G'(y) & -F'(x-z) \end{pmatrix}$$

The Jacobian satisfies  $P f'(x, y, z) P^{-1} = -f'(x, y, z)$  for the invertible matrix

$$P = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}. \tag{4.3}$$

Hence any symplectic Runge–Kutta method with 2 or fewer stages (and compositions thereof) will be volume preserving for such ODEs, which may not be obvious at first sight.

### 5. Measure-preservation of Runge–Kutta methods

In the introduction, we have pointed out that the trapezoidal method is not necessarily volume preserving but instead preserves the measure  $\det\left(I - \frac{h}{2} f'(x)\right) dx$  for quadratic Hamiltonian vector fields [2]. Recall that a map  $\phi$  preserves a measure  $\mu(x) dx$  if

$$\det(\phi'(x)) = \frac{\mu(x)}{\mu(\phi(x))}$$

This result is generalised in the following lemma.

**Lemma 5.1.** *The trapezoidal rule (1.1) preserves the measure  $\mu(x) dx$  with*

$$\mu(x) = \det\left(I \pm \frac{h}{2} f'(x)\right)$$

for vector fields  $f$  that satisfy the determinant condition (det).

**Proof.** We compute the Jacobian matrix  $\phi'_h$  of the trapezoidal rule,

$$\phi'_h(x) = I + \frac{h}{2} f'(x) + \frac{h}{2} f'(\phi_h(x)) \phi'_h(x),$$

and see that

$$\det(\phi'_h(x)) = \frac{\det\left(I + \frac{h}{2} f'(x)\right)}{\det\left(I - \frac{h}{2} f'(\phi_h(x))\right)} = \frac{\mu(x)}{\mu(\phi_h(x))}. \quad \square$$

This means that volume is conserved to order  $\mathcal{O}(h^2)$  globally (by Theorem 3.5), but more importantly that the dynamics of the numerical solution lie in the class of measure preserving systems, giving a qualitative advantage over methods lacking this property [9]. The trapezoidal method is conjugate to the implicit midpoint rule [6, VI.8], which goes some way towards explaining this behaviour. However, the next method we consider has similar measure preserving properties, but doesn't appear likewise to be “conjugate to volume preserving”.

There has been recent interest in the properties of the Kahan method [1,2]. For a quadratic vector field  $f(x) = Q(x) + L(x) + d$  where  $Q$  is quadratically homogeneous,  $L$  is linear and  $d$  is constant, the symmetric bilinear form  $q(x, y)$  is formed by polarisation,

$$q(x, y) = \frac{1}{2} \left( Q(x + y) - Q(x) - Q(y) \right), \tag{5.1}$$

and Kahan's unconventional numerical method is then given by

$$\frac{\phi_h(x) - x}{h} = q(x, \phi_h(x)) + \frac{1}{2} L(x + \phi_h(x)) + d. \tag{5.2}$$

In [2], it was shown that Kahan's method is equivalent to a three-stage Runge–Kutta method restricted to quadratic vector fields. We give the following generalisation.

**Lemma 5.2.** *Restricted to quadratic vector fields, Kahan's method is equivalent to the s-stage Runge–Kutta method*

$$\phi_h(x) = x + h \sum_{i=1}^s b_i f(x + c_i(\phi_h(x) - x)), \tag{5.3}$$

for any  $b$  and  $c$  satisfying  $\sum_{i=1}^s b_i = 1$ ,  $\sum_{i=1}^s b_i c_i = \frac{1}{2}$ ,  $\sum_{i=1}^s b_i c_i^2 = 0$ . This implies that the Butcher tableau satisfies  $A = cb^T$ , but the converse is not true.

**Proof.** Let  $x' = \phi_h(x)$  and write the vector field as  $f(x) = q(x, x) + Lx + d$  with the symmetric bilinear form  $q$ , then, expanding out and setting equal to Kahan's method (5.2)

$$\begin{aligned} \frac{x' - x}{h} &= \sum_{i=1} b_i q \left( x + c_i(x' - x), x + c_i(x' - x) \right) + L(x + c_i(x' - x)) + d \\ &= q(x', x) + \frac{1}{2} L(x + x') + d \end{aligned}$$

yields the above conditions.  $\square$

In [2, Prop. 5], it was shown that for quadratic Hamiltonian vector fields, Kahan's method preserves the measure with density  $\mu(x) = \det(I - \frac{h}{2} f'(x))^{-1}$ . The proof is easily extended to all quadratic vector fields satisfying the determinant condition (det).

**Lemma 5.3.** *Kahan's method preserves the measure  $\mu(x)dx$  with*

$$\mu(x) = \det \left( I \pm \frac{h}{2} f'(x) \right)^{-1} \tag{5.4}$$

for quadratic vector fields  $f$  that satisfy the determinant condition (det).

**Proof.** We compute the Jacobian matrix  $\phi'_h$  of Kahan's method in the form (1.2),

$$\phi'_h(x) = \frac{I - \frac{h}{2} f'(x) + \frac{h}{2} f' \left( \frac{x + \phi_h(x)}{2} \right)}{I + \frac{h}{2} f'(\phi_h(x)) - \frac{h}{2} f' \left( \frac{x + \phi_h(x)}{2} \right)}.$$

Since  $f$  is quadratic,  $f'$  is affine and thus

$$\det(\phi'_h(x)) = \frac{\det(I + \frac{h}{2} f'(\phi_h(x)))}{\det(I - \frac{h}{2} f'(x))} = \frac{\mu(x)}{\mu(\phi_h(x))}. \quad \square$$

Due to the similarity of this measure to that preserved by the trapezoidal method, one might at first glance suggest that the Kahan method is conjugate to some volume preserving method too, but this does not appear to be the case. At least, Kahan's method is not conjugate by B-series to any symplectic method [2]. It may be interesting to investigate how these measure preserving properties of the trapezoidal rule and Kahan's method can be generalised.

From Lemmata 2.4, 2.5 and 2.6 on linear foliations follow similar measure preservation properties generalising the volume preservation properties discussed in the previous section.

**Theorem 5.4.** Suppose that a given Runge–Kutta method preserves the measure  $\mu$  on  $\mathbb{R}^n$  when solving the ODE  $\dot{x} = f(x)$ . Then when solving the ODE  $\dot{x} = \tilde{f}(x)$ , where  $\tilde{f}(x) = Pf(P^{-1}x)$  for some invertible matrix  $P$ , the method preserves the measure with density  $\tilde{\mu}(x) = \mu(P^{-1}x)$ .

**Proof.** By assumption,  $\det(\phi'_h(y))\mu(\phi_h(y)) = \mu(y)$  for all  $y \in \mathbb{R}^n$ . Using the notation and results of Lemma 2.4,  $\det(\tilde{\phi}'_h(x))\tilde{\mu}(\tilde{\phi}_h(x)) = \det(\phi'_h(P^{-1}x))\mu(\phi_h(P^{-1}x)) = \mu(P^{-1}x) = \tilde{\mu}(x)$ .  $\square$

**Theorem 5.5.** Suppose that a given 1-stage Runge–Kutta method preserves the measure  $\rho dx$  on  $\mathbb{R}^m$  when solving the ODE  $\dot{x} = u(x)$ , and it preserves the measure  $v(y)dy$  on  $\mathbb{R}^n$  when solving the ODE  $\dot{y} = v(x, y)$  for all  $x \in \mathbb{R}^m$ . Then when solving the ODE  $(\dot{x}, \dot{y}) = (u(x), v(x, y))$ , the method preserves the product measure  $\mu(x, y)dxdy = \rho(x)v(y)dxdy$  on  $\mathbb{R}^{n+m}$ .

**Proof.** By Lemma 2.5,  $\phi_h(x, y) = (\psi_h(x), \chi_h(k_1(x), y))^\top$ , where  $k_1(x)$  is the internal stage of the 1-stage method. Hence by the definition of  $\mu$ ,

$$\mu(\phi_h(x, y)) = \rho(\psi_h(x))v(\chi_h(k_1(x), y)).$$

By assumption, we have for all  $x$  and  $y$ ,

$$\det(\psi'_h(x))\rho(\psi_h(x)) = \rho(x), \quad \det(\partial_y \chi_h(x, y))v(\chi_h(x, y)) = v(y). \tag{5.5}$$

Finally, Lemma 2.5 gives us  $\det(\phi'_h(k_1(x), y)) = \det(\psi'_h(x))\det(\partial_y \chi_h(x, y))$ . Combining all of these results,

$$\begin{aligned} \det(\phi'_h(x, y))\mu(\phi_h(x, y)) &= \det(\psi'_h(x))\rho(\phi_h(x))\det(\partial_y \chi_h(k_1(x), y))v(\chi_h(k_1(x), y)) \\ &= \rho(x)v(y) \\ &= \mu(x, y). \end{aligned}$$

Hence the measure with density  $\mu$  on  $\mathbb{R}^{n+m}$  is conserved.  $\square$

From the results of Lemma 2.6 and using its notation, we deduce that a generalisation for measure preserving RK methods with more stages even for sums  $f(x, y) = (u(x), w(x) + v(y))^\top$  is not trivial since then,

$$\det(\phi'_h(x, y)) = \det(\psi'_h(x))\det(\partial_y \chi_h(d(x), y + he(x))),$$

and a product measure  $\mu(x, y)dxdy = \rho(x)v(y)dxdy$  transforms according to

$$\begin{aligned} \det(\phi'_h(x, y))\mu(\phi_h(x, y)) &= \det(\phi'_h(x, y))\rho(\psi_h(x))v(\chi_h(d(x), y + he(x)) + hc(x)) \\ &= \det(\psi'_h(x))\rho(\psi_h(x)) \cdot \det(\partial_y \chi_h(d(x), y + he(x)))v(\chi_h(d(x), y + he(x)) + hc(x)). \end{aligned}$$

Assume that  $\psi_h, \chi_h$  preserve the measures with densities  $\rho(x)$  and  $v(y)$ , respectively, then, if  $c_h = e_h = 0$ , the product measure is preserved. This additional condition holds, e.g., for the trapezoidal rule for which we get that  $d_h(x) = (w(k_1) + w(k_2))/2$ . Further methods satisfying  $c_h = e_h = 0$  can be constructed easily<sup>1</sup> but they might preserve measures for trivial vector fields only. Kahan’s method derived from Lemma 5.2 does not simplify in this way, however, we can give the following result:

**Theorem 5.6.** Generalised Kahan’s methods from Lemma 5.2 preserve the measure  $\mu(x, y)dxdy$  with  $\mu(x, y) = \det(I + \frac{h}{2}f'(x, y))^{-1}$  for linearly foliate vector fields of the form  $f(x, y) = (u(x), v(y) + w(x))^\top$  where  $w$  is arbitrary, and  $u, v \in \mathcal{D}$  are quadratic.

**Proof.** Let  $z = (x, y)$ , and write  $\phi_h(z) = (\psi_h(x), \sigma_h(x, y))^\top$ . We compute the Jacobian determinant of (5.3)

$$\det(\phi'_h(z)) = \frac{\det(I + h \sum_{i=1}^N b_i(1 - c_i)f'(z + c_i(\phi_h(z) - z)))}{\det(I - h \sum_{i=1}^N b_i c_i f'(z + c_i(\phi_h(z) - z)))}$$

using that  $f'$  is block-diagonal, we arrive at

$$= \frac{\det(I + h \sum_{i=1}^N b_i(1 - c_i)u'(x + c_i(\psi_h(x) - x))) \det(I + h \sum_{i=1}^N b_i(1 - c_i)v'(y + c_i(\sigma_h(x, y) - y))}{\det(I - h \sum_{i=1}^N b_i c_i u'(x + c_i(\psi_h(x) - x))) \det(I - h \sum_{i=1}^N b_i c_i v'(y + c_i(\sigma_h(x, y) - y))}$$

and since  $u', v'$  are affine, we can simplify using the assumptions on the coefficients from Lemma 5.2 to

$$= \frac{\det(I + \frac{h}{2}u'(\psi_h(x))) \det(I + \frac{h}{2}v'(\sigma_h(x, y)))}{\det(I - \frac{h}{2}u'(x)) \det(I - \frac{h}{2}v'(y))} = \frac{\det(I + \frac{h}{2}f'(\phi_h(x, y)))}{\det(I - \frac{h}{2}f'(x, y))} = \frac{\mu(z)}{\mu(\phi_h(z))}. \quad \square$$

<sup>1</sup> Let, e.g.,  $a_{1j} = 0$  and  $a_{ij} = b_j$  for some  $i$  and all  $j$ .

**Remark 5.7.** The theorem is not true for more general foliate vector fields within the class  $\mathcal{F}^{(\infty)}$ , e.g.,  $f(x, y) = (u(x), J^{-1}\nabla_y H(x, y))^T$  where  $u(x)$  is a simple harmonic oscillator and with the Hamiltonian  $H(x, y) = (p_x q_x) p_y q_y$  using the usual notation for the momentum and position coordinates  $x = (q_x, p_x)$ ,  $y = (q_y, p_y)$ . Note that the Hamiltonian is still quadratic in  $y$ !

## Acknowledgements

The authors would like to thank Robert McLachlan for stimulating discussions and suggestions, and the anonymous referees for their useful comments to improve this paper. This research was supported by the Marie Curie International Research Staff Exchange Scheme, grant number DP140100640, within the 7th European Community Framework Programme; by the Australian Research Council grant number 269281; and by the UK EPSRC grant EP/H023348/1 for the Cambridge Centre for Analysis.

## References

- [1] E. Celledoni, R.I. McLachlan, D.I. McLaren, B. Owren, G.R.W. Quispel, Integrability properties of Kahan's method, *J. Phys. A* 47 (36) (2014) 365202.
- [2] E. Celledoni, R.I. McLachlan, B. Owren, G.R.W. Quispel, Geometric properties of Kahan's method, *J. Phys. A* 46 (2) (2013) 025201.
- [3] P. Chartier, A. Murua, Preserving first integrals and volume forms of additively split systems, *IMA J. Numer. Anal.* 27 (2) (2007) 381–405.
- [4] Thierry Dombre, Uriel Frisch, John M. Greene, Michel Hénon, A. Mehr, Andrew M. Soward, Chaotic streamlines in the ABC flows, *J. Fluid Mech.* 167 (1986) 353–391.
- [5] K. Feng, Z.J. Shang, Volume-preserving algorithms for source-free dynamical systems, *Numer. Math.* 71 (4) (1995) 451–463.
- [6] E. Hairer, C. Lubich, G. Wanner, *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*, vol. 31, Springer Science & Business Media, 2006.
- [7] A. Iserles, G.R.W. Quispel, P.S.P. Tse, B-series methods cannot be volume-preserving, *BIT Numer. Math.* 47 (2) (2007) 351–378.
- [8] R.I. McLachlan, M. Perlmutter, G.R.W. Quispel, Lie group foliations: dynamical systems and integrators, *Future Gener. Comput. Syst.* 19 (7) (2003) 1207–1219.
- [9] R.I. McLachlan, G.R.W. Quispel, What kinds of dynamics are there? Lie pseudogroups, dynamical systems and geometric integration, *Nonlinearity* 14 (6) (2001) 1689.
- [10] R.I. McLachlan, G.R.W. Quispel, N. Robidoux, Geometric integration using discrete gradients, *Philos. Trans. R. Soc. Lond. A* 357 (1754) (1999) 1021–1045.
- [11] M.-Z. Qin, W.-J. Zhu, Volume-preserving schemes and numerical experiments, *Comput. Math. Appl.* 26 (4) (1993) 33–42.
- [12] G.R.W. Quispel, Volume-preserving integrators, *Phys. Lett. A* 206 (1) (1995) 26–30.
- [13] V. Volterra, *Leçons sur la théorie mathématique de la lutte pour la vie*, reprinted by Éditions Jacques Gabay, Sceaux (1990), 1931.