

On Exactly Conservative Integrators

John C. Bowman

Max-Planck-Institut für Plasmaphysik

EURATOM Association, D 85748 Garching, Germany

and

B. A. Shadwick

Department of Physics, University of California at Berkeley

366 Le Conte Hall, Berkeley, CA 94720-7300

and

P. J. Morrison

Institute for Fusion Studies, The University of Texas at Austin

Austin, Texas 78712

Traditional explicit numerical discretizations of conservative systems generically predict artificial secular drifts of nonlinear invariants. These algorithms are based on polynomial functions of the time step. We discuss a general approach for developing explicit algorithms that conserve such invariants exactly. We illustrate the method by applying it to the truncated two-dimensional Euler equations.

Introduction

It is often desirable on physical grounds that numerical discretizations of initial value problems respect certain nonlinear conservation properties. This is the case in inviscid fluid simulations, where relaxation to a statistical mechanical equilibrium parameterized by quadratically nonlinear invariants is expected. In this paper we illustrate a general method (Shadwick *et al.* 1997) for deriving *explicit* conservative integration algorithms for truncations \mathcal{K} of the Fourier-transformed Euler equations for a two-dimensional fluid,

$$\frac{du_{\mathbf{k}}}{dt} = S_{\mathbf{k}}(u), \quad (1)$$

where $u = \{u_{\mathbf{k}} : \mathbf{k} \in \mathcal{K}\}$. For simplicity we will restrict our discussion to the case where the $u_{\mathbf{k}}$'s are real; the general case follows immediately upon splitting complex amplitudes into real and imaginary parts. The source functions $S_{\mathbf{k}}$ satisfy the properties

$$\sum_{\mathbf{k}} u_{\mathbf{k}} S_{\mathbf{k}} = 0 \quad \text{and} \quad \sum_{\mathbf{k}} k^2 u_{\mathbf{k}} S_{\mathbf{k}} = 0, \quad (2)$$

which lead to the the conservation of two nonlinear invariants, the total energy and enstrophy,

$$E = \frac{1}{2} \sum_{\mathbf{k}} u_{\mathbf{k}}^2, \quad \text{and} \quad Z = \frac{1}{2} \sum_{\mathbf{k}} k^2 u_{\mathbf{k}}^2, \quad (3)$$

respectively.

Unfortunately, when (1) is integrated numerically using standard explicit methods [or even with symplectic integrators; see Ge & Marsden (1988)], neither E nor Z are exactly conserved. This behaviour is elucidated upon applying Euler's method with a time step τ :

$$u_{\mathbf{k}}(t + \tau) = u_{\mathbf{k}}(t) + \tau S_{\mathbf{k}}. \quad (4)$$

The energy at the new time is seen to be

$$\begin{aligned} E(t + \tau) &= \frac{1}{2} \sum_{\mathbf{k}} [u_{\mathbf{k}}(t) + \tau S_{\mathbf{k}}]^2 = \frac{1}{2} \sum_{\mathbf{k}} [u_{\mathbf{k}}^2 + 2\tau S_{\mathbf{k}} u_{\mathbf{k}} + \tau^2 S_{\mathbf{k}}^2] \\ &= E(t) + \frac{1}{2} \tau^2 \sum_{\mathbf{k}} S_{\mathbf{k}}^2, \end{aligned} \quad (5)$$

upon using (2). The total energy is thus always increasing. A similar calculation for the enstrophy gives

$$Z(t + \tau) = Z(t) + \frac{1}{2} \tau^2 \sum_{\mathbf{k}} k^2 S_{\mathbf{k}}^2, \quad (6)$$

which likewise is always increasing. For extremely long runs these results imply that a very small time step will be required to bound the accumulated error by a given value.

One technique for enforcing the preservation of a constant of motion is to use the invariant to reduce the number of equations that must be solved. If the constants are in involution, then an entire degree of freedom (one coordinate and one momenta) can be removed from the dynamics for each such constant. This is seldom practical since the relationship between the constants of motion and a given dynamical variable may well be noninvertible [see the discussion in Gear (1986)]. The net result is that the reduced equations tend to be more complicated than the original system (hence the “force” terms are more expensive to compute); in a system with a large number of degrees of freedom, little advantage is gained. Furthermore, if the constants of motion are not in involution, the system obtained by eliminating these invariants will be noncanonical (Morrison 1993; Marsden 1992), resulting in even greater complexity.

Borrowing from the ideas of backward error analysis (Sanz-Serna & Calvo 1994), we instead construct a new system of equations that, under the conventional (nonconservative) integrator, yields a conservative numerical approximation to the original equations. Consider the alternative problem described by equations of the form

$$\frac{du_{\mathbf{k}}}{dt} = S_{\mathbf{k}}(u) + f_{\mathbf{k}}. \quad (7)$$

Our objective is to find an $f_{\mathbf{k}}$ that guarantees exact energy and enstrophy conservation and that vanishes in the small time-step limit. The form of $f_{\mathbf{k}}$ will depend on the integration algorithm. For pedagogical reasons, we begin by illustrating the method by deriving $f_{\mathbf{k}}$ for Euler’s method. We then proceed to construct a more practical second-order conservative predictor–corrector scheme.

Conservative Euler algorithm

Application of Euler's method to the modified system (7) yields

$$u_{\mathbf{k}}(t + \tau) = u_{\mathbf{k}}(t) + \tau(S_{\mathbf{k}} + f_{\mathbf{k}}). \quad (8)$$

The energy at the new time,

$$\begin{aligned} E(t + \tau) &= \frac{1}{2} \sum_{\mathbf{k}} [u_{\mathbf{k}}(t) + \tau(S_{\mathbf{k}} + f_{\mathbf{k}})]^2 \\ &= E(t) + \frac{1}{2} \sum_{\mathbf{k}} [2\tau f_{\mathbf{k}} u_{\mathbf{k}} + \tau^2(S_{\mathbf{k}} + f_{\mathbf{k}})^2], \end{aligned} \quad (9)$$

will be conserved provided

$$\sum_{\mathbf{k}} [2f_{\mathbf{k}} u_{\mathbf{k}} + \tau(S_{\mathbf{k}} + f_{\mathbf{k}})^2] = 0. \quad (10)$$

There is considerable freedom in satisfying (10). To ensure that our discrete solution approaches the exact solution of the original differential equation in the limit $\tau \rightarrow 0$, it is necessary that $f_{\mathbf{k}}$ vanish in this limit. That is, in the limit of an infinitesimal time step, we must recover the original integration algorithm (to first order in τ). Moreover, one would prefer that $f_{\mathbf{k}}$ not introduce additional couplings into the differential equations. In light of this observation, let us try to satisfy (10) with the more restrictive condition that each term in the sum must independently vanish:

$$2f_{\mathbf{k}} u_{\mathbf{k}} + \tau(S_{\mathbf{k}} + f_{\mathbf{k}})^2 = 0. \quad (11)$$

There is an additional motivation for this *ansatz*, namely that for $f_{\mathbf{k}}$ satisfying (11), the enstrophy will also be conserved. These equations are easily solved, yielding

$$\tau f_{\mathbf{k}} = -(u_{\mathbf{k}} + \tau S_{\mathbf{k}}) + \sigma_{\mathbf{k}} \sqrt{u_{\mathbf{k}}^2 + 2\tau S_{\mathbf{k}} u_{\mathbf{k}}}, \quad (12)$$

where $\sigma_{\mathbf{k}} = \sigma_{\mathbf{k}}(t, \tau)$ is so far an unknown sign. Evaluation of (12) at $\tau = 0$ implies that $\sigma_{\mathbf{k}}(t, 0) = \text{sgn}(u_{\mathbf{k}}(t))$. Upon substituting (12) into the Euler integrator, (8), we obtain

$$u_{\mathbf{k}}(t + \tau) = \sigma_{\mathbf{k}} \sqrt{u_{\mathbf{k}}^2 + 2\tau S_{\mathbf{k}} u_{\mathbf{k}}}. \quad (13)$$

It is now clear that $\sigma_{\mathbf{k}}(t, \tau)$ must in fact be the sign of $u_{\mathbf{k}}(t + \tau)$.

If $u_{\mathbf{k}}(t) \neq 0$, then for sufficiently small τ the sign can be expressed explicitly as $\sigma_{\mathbf{k}} = \text{sgn}(u_{\mathbf{k}}(t))$. In the $\tau \rightarrow 0$ limit, $f_{\mathbf{k}}$ then vanishes, or equivalently, (13) reduces to Euler's method:

$$u_{\mathbf{k}}(t + \tau) = \text{sgn}(u_{\mathbf{k}}(t)) \sqrt{u_{\mathbf{k}}^2 + 2\tau S_{\mathbf{k}} u_{\mathbf{k}}} \approx u_{\mathbf{k}} + \tau S_{\mathbf{k}}. \quad (14)$$

In this case the new algorithm predicts values of $u_{\mathbf{k}}(t + \tau)$ that are quite close to those given by Euler's method—this is exactly what one would expect. The energy and enstrophy errors arising from (4) are the result of small (but nontrivial) errors in $u_{\mathbf{k}}(t + \tau)$ that can be corrected by making only a slight modification to the algorithm.

However, if $u_{\mathbf{k}}(t) = 0$, it is seen from (12) that $f_{\mathbf{k}} = -S_{\mathbf{k}}$. Consequently, (13) has a spurious fixed point at $u_{\mathbf{k}}(t) = 0$. Moreover, given a fixed time step τ , (14) will break down when $|u_{\mathbf{k}}| < 2\tau |S_{\mathbf{k}}|$. A related problem with (13) is that the argument of the radical can become negative. The condition $u_{\mathbf{k}}(u_{\mathbf{k}} + 2\tau S_{\mathbf{k}}) < 0$ implies that Euler's method predicts a sign change of $u_{\mathbf{k}}$ between t and $t + 2\tau$; hence $u_{\mathbf{k}}$ is in this case also in the vicinity of zero. While a modification to (13) that circumvents these problems is given in Shadwick *et al.* (1997), the second-order algorithm discussed next does not share these difficulties and is in any case of greater practical value.

Conservative predictor–corrector algorithm

For most applications, it is preferable to use a scheme that is of higher order and has better stability properties than Euler's method. Let us apply a simple second-order predictor–corrector (PC) scheme to (1):

$$\tilde{u}_{\mathbf{k}} = u_{\mathbf{k}} + \tau S_{\mathbf{k}}, \quad (15a)$$

$$u_{\mathbf{k}}(t + \tau) = u_{\mathbf{k}} + \frac{\tau}{2} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}}), \quad (15b)$$

where $\tilde{S}_{\mathbf{k}} = S_{\mathbf{k}}(\tilde{u})$ and $\tilde{u} = \{\tilde{u}_{\mathbf{k}} : \mathbf{k} \in \mathcal{K}\}$. Using a second-order method overcomes the fixed-point problem that we encountered with Euler's method.

The energy will evolve according to

$$E(t + \tau) = \frac{1}{2} \sum_{\mathbf{k}} \left[u_{\mathbf{k}}^2 + \tau u_{\mathbf{k}} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}}) + \frac{\tau^2}{4} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}})^2 \right]$$

$$= E(t) + \frac{\tau^2}{8} \sum_{\mathbf{k}} (S_{\mathbf{k}} - \tilde{S}_{\mathbf{k}})^2. \quad (16)$$

A similar calculation gives

$$Z(t + \tau) = Z(t) + \frac{\tau^2}{8} \sum_{\mathbf{k}} k^2 (S_{\mathbf{k}} - \tilde{S}_{\mathbf{k}})^2. \quad (17)$$

Again, the numerical method yields an ever increasing energy and enstrophy.

To obtain a conservative version of this algorithm, we apply the predictor–corrector method to the modified equation of motion, (7), yielding

$$\tilde{u}_{\mathbf{k}} = u_{\mathbf{k}} + \tau (S_{\mathbf{k}} + f_{\mathbf{k}}), \quad (18a)$$

$$u_{\mathbf{k}}(t + \tau) = u_{\mathbf{k}} + \frac{\tau}{2} (S_{\mathbf{k}} + f_{\mathbf{k}} + \tilde{S}_{\mathbf{k}} + \tilde{f}_{\mathbf{k}}). \quad (18b)$$

Only a small correction to (15) is required to enforce the desired conservation properties. It turns out that these properties can be achieved by modifying only the corrector part of the integrator; since the predictor is merely an intermediate approximation, there is no need for it to be conservative. We can thus replace (18) with the simpler prescription

$$\tilde{u}_{\mathbf{k}} = u_{\mathbf{k}} + \tau S_{\mathbf{k}}, \quad (19a)$$

$$u_{\mathbf{k}}(t + \tau) = u_{\mathbf{k}} + \frac{\tau}{2} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}} + g_{\mathbf{k}}). \quad (19b)$$

As before, we determine $g_{\mathbf{k}}$ by demanding conservation of energy and enstrophy. The energy and enstrophy at $t + \tau$ will be simultaneously conserved if $g_{\mathbf{k}} u_{\mathbf{k}} - \tau S_{\mathbf{k}} \tilde{S}_{\mathbf{k}} + \frac{\tau}{4} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}} + g_{\mathbf{k}})^2 = 0$. Some straightforward algebra gives

$$\frac{\tau}{2} g_{\mathbf{k}} = - \left[u_{\mathbf{k}} + \frac{\tau}{2} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}}) \right] + \sigma_{\mathbf{k}} \sqrt{u_{\mathbf{k}}^2 + \tau (u_{\mathbf{k}} S_{\mathbf{k}} + \tilde{u}_{\mathbf{k}} \tilde{S}_{\mathbf{k}})}, \quad (20)$$

where we choose $\sigma_{\mathbf{k}} = \pm 1$ such that as $\tau \rightarrow 0$, $g_{\mathbf{k}}$ vanishes. We consider the limit of small τ in two cases. If $u_{\mathbf{k}}$ is nonzero, then for small enough τ , both $u_{\mathbf{k}}$ and $\tilde{u}_{\mathbf{k}}$ have the same sign and we can expand the radical to give

$$\frac{\tau}{2} g_{\mathbf{k}} = -u_{\mathbf{k}} - \frac{\tau}{2} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}}) + \sigma_{\mathbf{k}} \operatorname{sgn}(u_{\mathbf{k}}) \left[u_{\mathbf{k}} + \frac{\tau}{2} (S_{\mathbf{k}} + \tilde{S}_{\mathbf{k}}) \right] + O(\tau^2), \quad (21)$$

leading us to choose $\sigma_{\mathbf{k}} = \text{sgn}(u_{\mathbf{k}})$. Otherwise, if $u_{\mathbf{k}} = 0$, then $\tilde{u}_{\mathbf{k}} = \tau S_{\mathbf{k}}$ and $\tilde{S}_{\mathbf{k}} = S_{\mathbf{k}} + O(\tau)$, so that

$$\frac{\tau}{2} g_{\mathbf{k}} = -\tau S_{\mathbf{k}} + \sigma_{\mathbf{k}} \sqrt{\tau^2 S_{\mathbf{k}}^2} + O(\tau^2) = -\tau S_{\mathbf{k}} + \tau \sigma_{\mathbf{k}} \text{sgn}(S_{\mathbf{k}}) S_{\mathbf{k}} + O(\tau^2). \quad (22)$$

In this case we take $\sigma_{\mathbf{k}} = \text{sgn}(S_{\mathbf{k}}) = \text{sgn}(\tilde{u}_{\mathbf{k}})$. In the previous case, we noted, for small τ , that $u_{\mathbf{k}}$ and $\tilde{u}_{\mathbf{k}}$ have the same sign. Therefore, the choice $\sigma_{\mathbf{k}} = \text{sgn}(\tilde{u}_{\mathbf{k}})$ will always provide the correct limiting behaviour.

Using the expression (22) for $g_{\mathbf{k}}$ in our modified predictor–corrector algorithm, (19), we obtain the following conservative integrator:

$$\tilde{u}_{\mathbf{k}} = u_{\mathbf{k}} + \tau S_{\mathbf{k}}, \quad (23a)$$

$$u_{\mathbf{k}}(t + \tau) = \tilde{\sigma}_{\mathbf{k}} \sqrt{u_{\mathbf{k}}^2 + \tau (u_{\mathbf{k}} S_{\mathbf{k}} + \tilde{u}_{\mathbf{k}} \tilde{S}_{\mathbf{k}})}, \quad (23b)$$

where $\tilde{\sigma}_{\mathbf{k}} = \text{sgn}(\tilde{u}_{\mathbf{k}})$. Unlike (13), this algorithm, which we call “conservative predictor–corrector,” (C–PC), does not suffer from fixed points: as $\tau \rightarrow 0$, it reduces to the conventional predictor–corrector, (15), even in the case $u_{\mathbf{k}} = 0$. Thus, C–PC may be seen as finite-time step generalization of PC. Both methods agree with the exact solution to second order in the time step. It is still possible that the argument of the radical can become negative; however, this merely indicates that the step size is too large. We now show that a finite number of time-step reductions can be used to integrate the system (1) through a region where the argument of the radical in (23b) is negative, provided that $S_{\mathbf{k}}$ has continuous (and hence bounded) first derivatives on the closed interval of integration.

Suppose that at some time t there exists a mode \mathbf{k} (not necessarily unique) such that

$$u_{\mathbf{k}}^2 + \tau (u_{\mathbf{k}} S_{\mathbf{k}} + \tilde{u}_{\mathbf{k}} \tilde{S}_{\mathbf{k}}) < 0. \quad (24)$$

Since the left-hand side of this expression is a continuous function of τ , for each such mode \mathbf{k} there exists a $\tau_1 \geq 0$ such that $u_{\mathbf{k}}^2 + \tau_1 (u_{\mathbf{k}} S_{\mathbf{k}} + \tilde{u}_{\mathbf{k}} \tilde{S}_{\mathbf{k}}) = 0$. If more than one mode satisfies this condition, we choose the one with the smallest τ_1 . One can replace the original time step with τ_1 so that, at time $t_1 = t + \tau_1$, (23b) becomes $u_{\mathbf{k}}(t_1) = 0$. From here, one may resume the integration with the original time step:

$$\tilde{u}_{\mathbf{k}}(t_1 + \tau) = \tau S_{\mathbf{k}}(t_1), \quad (25a)$$

$$\begin{aligned}
u_{\mathbf{k}}(t_1 + \tau) &= \text{sgn}(\tilde{u}_{\mathbf{k}}) \left(\tau^2 S_{\mathbf{k}}(t_1) \tilde{S}_{\mathbf{k}}(t_1) \right)^{1/2} \\
&= \text{sgn}(\tilde{u}_{\mathbf{k}}) \tau \left(S_{\mathbf{k}}^2(t_1) + \tau S_{\mathbf{k}}(t_1) \sum_j S_j(t_1) \left. \frac{\partial S_{\mathbf{k}}}{\partial u_j} \right|_{t_\xi} \right)^{1/2}, \quad (25b)
\end{aligned}$$

for some $t_\xi \in (t_1, t_1 + \tau)$, by the Mean Value Theorem. If $S_{\mathbf{k}}$ and its derivatives are bounded on the interval of integration, for sufficiently small τ the argument of the radical will be non-negative and a single reduction of the time step will suffice. In practice, it is not necessary to reduce τ exactly to τ_1 . Instead, successive reductions by a constant factor will eventually make the argument of the radical non-negative so that the system can be further integrated. Under the stated conditions, it is consequently never necessary to reduce the time step all the way to zero (a circumstance that has never been encountered in our implementations of C-PC).

In Fig. (1) we contrast the evolution of a system of three real modes under PC and C-PC, choosing $k^2 = 3$, $p^2 = 9$, $q^2 = 6$, $S_{\mathbf{k}} = u_{\mathbf{p}}u_{\mathbf{q}}$, $S_{\mathbf{p}} = u_{\mathbf{q}}u_{\mathbf{k}}$, and $S_{\mathbf{q}} = -2u_{\mathbf{k}}u_{\mathbf{p}}$ (Kraichnan 1963). This problem is integrable; the exact solution is a simple closed curve. The solid line is the orbit computed with the conservative integrator, while the dots represent the solution obtained from the conventional predictor-corrector. We see that the conservative integrator correctly reproduces the topological structure of the trajectory; in contrast, the conventional method exhibits a drift corresponding to a 4% gain in the total energy.

Discussion

A simple interpretation of (13) and (23) sheds light both on their form and on the existence of the two branches, labeled by $\sigma_{\mathbf{k}}$. Most traditional numerical methods conserve the linear invariants of a system. Consequently, one might be led to consider the possibility of transforming $u_{\mathbf{k}}$ to new variables, in terms of which the invariants are linear. For the Euler equations, this can be accomplished by making the transformation $\phi_{\mathbf{k}} = u_{\mathbf{k}}^2$. Upon applying the Euler method in the $\phi_{\mathbf{k}}$ space and transforming back by taking the square root, one immediately obtains (13). This indicates that our restriction of the general constraint (10) to the condition (11) merely ensures that the modal energies evolve in a manner consistent with the Euler discretization of the energy equations. The C-PC algorithm can be viewed in the same light,

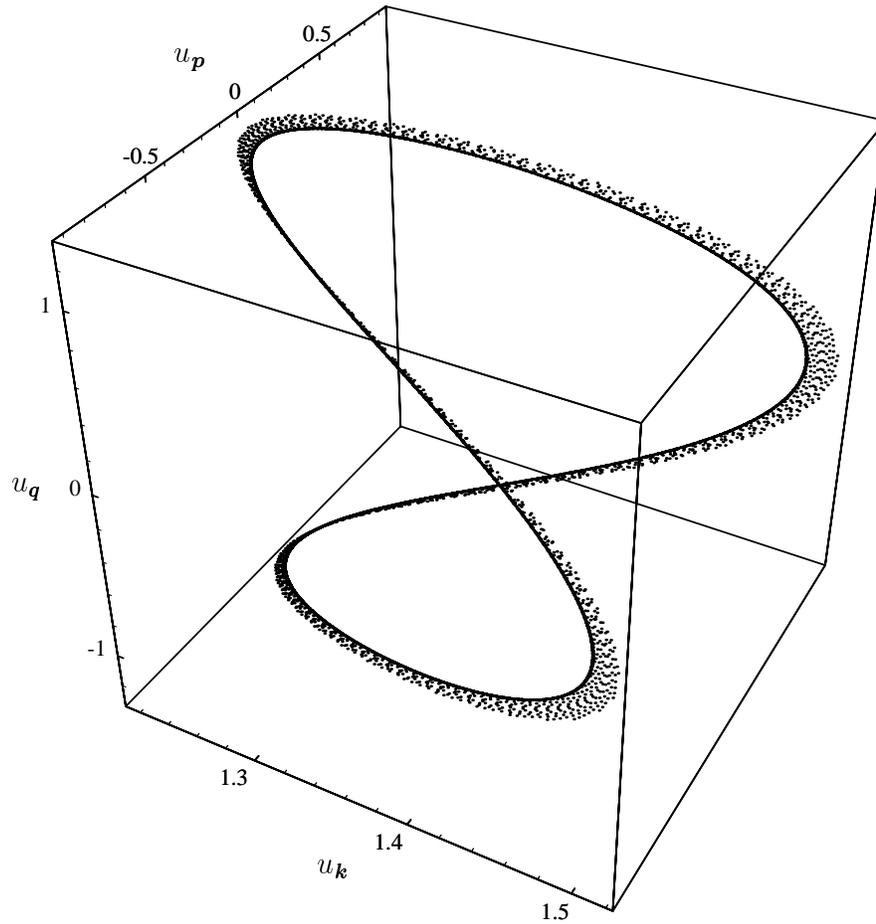


Figure 1: Integration of a three-mode truncation of the Euler equations using a conventional second-order predictor-corrector (dotted line) and the conservative predictor-corrector (solid line). Both methods took approximately 4000 time steps of size 0.05. Initially $u_k = \sqrt{1.5}$, $u_p = 0$, and $u_q = \sqrt{1.5}$.

except that the predictor is taken to have the simpler, nonconservative form. Further applications of this idea to the Lotka–Volterra predator–prey model and to the Kepler problem are given in Shadwick *et al.* (1997) and to high-resolution turbulence computations in Bowman *et al.* (1996).

The authors acknowledge support from United States DoE contract No. DE–FG03–96ER–54346.

REFERENCES

- BOWMAN, J. C., SHADWICK, B. A. & MORRISON, P. J. 1996 Spectral reduction for two-dimensional turbulence. In *Transport, Chaos, and Plasma Physics 2*, edited by Benkadda, S., Doveil, F. & Elskens, Y. pp. 58–73. Institute Méditerranéen de Technologie (Marseille, 1995) World Scientific.
- GE ZHONG & MARSDEN, J. E. 1988 Lie–Poisson Hamilton–Jacobi theory and Lie–Poisson integrators. *Phys. Lett. A* **133**, 134–139.
- GEAR, C. W. 1986 Maintaining solution invariants in the numerical solution of ODEs. *SIAM J. Sci. Stat. Comput.* **7**, 734–743.
- KRAICHNAN, R. H. 1963 Direct-interaction approximation for a system of several interacting simple shear waves. *Phys. Fluids* **6**, 1603–1609.
- MARSDEN, J. E. 1992 *Lectures on Mechanics*. Cambridge University Press, Cambridge.
- MORRISON, P. J. 1993 Hamiltonian description of the ideal fluid. In *Proceeding of the 1993 Geophysical Fluid Dynamics Summer School*, edited by Salmon, R. & Ewing-Deremer, B. Number 94-12 in WHOI Report. Woods Hole, MA Woods Hole Oceanographic Institution WHOI. Also available as Institute for Fusion Studies report number IFSR #640r, University of Texas at Austin, to appear in *Rev. Mod. Phys.*, 1997.
- SANZ-SERNA, J. M. & CALVO, M. P. 1994 *Numerical Hamiltonian Problems*. Number 7 in Applied Mathematics and Mathematical Computation Chapman and Hall, London.
- SHADWICK, B. A., BOWMAN, J. C. & MORRISON, P. J. 1997 Exactly conservative integrators. to appear in *SIAM J. Appl. Math.*