

Local Mass Conservation and Velocity Splitting in PV-Based Balanced Models. Part I: The Hyperbalance Equations

ALI R. MOHEBALHOJEH

School of Mathematics and Statistics, University of St Andrews, St Andrews, United Kingdom, and Institute of Geophysics, University of Tehran, Tehran, Iran

MICHAEL E. MCINTYRE

Centre for Atmospheric Science, Department of Applied Mathematics and Theoretical Physics, Cambridge University, Cambridge, United Kingdom*

(Manuscript received 11 January 2006, in final form 27 September 2006)

ABSTRACT

This paper considers stratified and shallow water non-Hamiltonian potential-vorticity-based balanced models (PBMs). These are constructed using the exact (Rossby or Rossby–Ertel) potential vorticity (PV). The most accurate known PBMs are those studied by McIntyre and Norton and by Mohebalhojeh and Dritschel. It is proved that, despite their astonishing accuracy, these PBMs all fail to conserve mass locally. Specifically, they exhibit velocity splitting in the sense of having two velocity fields, \mathbf{v} and \mathbf{v}_m , the first to advect PV and the second to advect mass. The difference $\mathbf{v} - \mathbf{v}_m$ is nonzero in general, even if tiny. Unlike the different velocity splitting found in all Hamiltonian balanced models, the present splitting can be healed. The result is a previously unknown class of balanced models, here called “hyperbalance equations,” whose formal orders of accuracy can be made as high as those of any other PBM. The hyperbalance equations use a single velocity field \mathbf{v} to advect mass as well as to advect and evaluate the exact PV.

1. Introduction

A longstanding question in the theory of non-Hamiltonian balanced models based on potential vorticity (PV) inversion has been whether they can be expected to conserve mass locally, that is, pointwise, at the highest possible accuracies, where accuracy is judged by comparison with primitive equation evolution. McIntyre and Norton (2000, hereafter MN00) conjectured that the answer might be no, for reasons connected with the subtleties of spontaneous imbalance and adjustment. Here we establish that the answer is definitely no for all of the most accurate known models,

but definitely yes for a new class of accurate models, the hyperbalance equations, to be introduced here, when accuracy means formal order of accuracy. The question of whether the answer is yes for the actual numerical accuracy is addressed in Mohebalhojeh and McIntyre (2007, hereafter Part II). Formal order of accuracy is related to the number of diagnostic estimates of time derivatives retained in the system of equations defining the model’s balance relation. It is also related to, but not the same as, the asymptotic order of accuracy when the Rossby and Froude numbers go to zero, or Richardson number to infinity; see the discussion in Mohebalhojeh and Dritschel (2001), hereafter MD01.

The term balanced model is used here in its standard generic sense, that is, to mean a model from which inertia–gravity waves have been completely eliminated, whether by asymptotic procedures or by other means. The most accurate known models use other means (high-order formal truncation, sections 4 and 5 below). In all cases, the initial conditions for such a model require specification of just one scalar field or “master variable,” for instance the PV, to which all the other fields are “slaved” by the model’s balance relation and

* The Centre for Atmospheric Science is a joint initiative of the Department of Chemistry and the Department of Applied Mathematics and Theoretical Physics (more information available online at www.atm.damtp.cam.ac.uk/).

Corresponding author address: A. R. Mohebalhojeh, School of Mathematics and Statistics, University of St Andrews, North Haugh, St Andrews KY16 9SS, United Kingdom.
E-mail: arm@mcs.st-and.ac.uk

from which they can therefore be deduced at any instant by the diagnostic process called inversion.¹ This definition of balanced model excludes, for instance, the model denoted BEM by Allen et al. (1990) and Allen (1991), the initial conditions for which require specification of two scalar fields, and for which PV inversion is therefore not possible. It similarly excludes the so-called modified balance equation derived by Charney (Charney 1962; Moura 1976), which Gent and McWilliams (1983a) call the global balance equation.

A variety of balanced models have been described in the literature to date, some more accurate than others. We may distinguish two particular approaches to their formulation. The two approaches lead to models that will be labeled “locally mass conserving” and “PV-based,” respectively. The formulation of locally mass-conserving balanced models (MBMs) begins, quite naturally, with the idea that exact local mass conservation is so fundamental that it is to be incorporated from the start (e.g., Allen et al. 1990). The formulation of PV-based balanced models (PBMs) begins by contrast with the idea that, because primitive equation evolution has an exact material invariant, the PV, for frictionless, adiabatic flow—accurately invertible in a surprising range of circumstances—it might be useful to make the balanced model represent PV evolution as accurately as possible. This leads to the idea of trying to incorporate exact material PV conservation from the start, by formulating the model in such a way that its velocity field \mathbf{v} , obtained by PV inversion, both advects and evaluates the PV (e.g., Bleck 1974; Norton 1988; Lynch 1989; Warn et al. 1995; Mundt et al. 1997; MN00).

Despite certain difficulties with the PBM concept [Mohebalhojeh 2002; Eq. (5.6) below], it is PBMs rather than MBMs that have so far proved the most accurate. This was demonstrated by the work of Norton

(1988) and MN00, in which astonishingly accurate PV inversion and balanced-model evolution was discovered in the case of a complicated, chaotic shallow-water flow with Froude numbers exceeding 0.7 and Rossby numbers ranging up to ∞ , in a hemispheric model. This was a major surprise, contradicting the expectation that balance should be an asymptotic property dependent on the smallness of such parameters. For further discussion and examples, see MN00 and MD01, also McIntyre (2001).

This paper systematically explores the extent to which the two approaches can be reconciled and made exactly consistent without losing accuracy. The plan of the paper is as follows. Sections 2 and 3 characterize the complete subset of all PBMs for which the reconciliation is possible. The defining property is Eq. (3.2). A known example is the Bolin–Charney model or so-called balance equations, in their shallow water or isentropic-coordinate versions. Another is the linear balance model discussed for instance by Whitaker (1993). Sections 4 and 5 analyze the generic structure of the most accurate known PBMs, those of MN00 and MD01. Section 4 proves that this structure—which can produce the astonishing accuracy already referred to—precludes the reconciliation. In other words, each such PBM, however accurate it may be, exhibits velocity splitting in the sense that it has two distinct velocity fields, one to advect PV and one to advect mass. That splitting is the source of the difficulties noted in Mohebalhojeh (2002) and in Eq. (5.6) below.

Section 6 shows how to heal the splitting and overcome the difficulties by constructing locally mass conserving PBMs that have arbitrarily high formal orders of accuracy. The governing equations of such models can be regarded as generalizations of the linear-balance and Bolin–Charney balance equations; we therefore refer to them as hyperbalance equations. Section 7 presents some brief concluding remarks, discussing MN00’s conjecture and preparing for the numerical accuracy tests in Part II.

2. Υ -models: The simplest example

We need to analyze the general question of whether there exist accurate PBMs that are also locally mass conserving. We continue to restrict attention to frictionless, adiabatic flow. What would such a balanced model be like? It would have to be equivalent to a modified primitive equation system in which the spontaneous-adjustment emission of inertia–gravity waves (e.g., Lighthill 1952; Errico 1982; Warn and Ménéard 1986; Ford 1994a,b,c; Warn 1997; Ford et al. 2000, 2002; Saujani and Shepherd 2002; Vanneste and Yavneh

¹ To say more carefully what “one scalar field” means, one must distinguish between shallow water models and continuously stratified models. For a shallow water model (with a given total mass or mean depth), a balanced model in the standard sense can take as initial conditions the PV field alone, a single function of horizontal position, as compared with the three such functions required as initial conditions by the shallow water primitive equations. For a continuously stratified model (with a given mass under each isentropic surface), a balanced model in the standard sense can take as initial conditions a single infinity of functions of horizontal position, namely, all the isentropic distributions of PV in the interior together with the potential temperature θ or “PV delta function” at the lower and/or upper boundary (Bretherton 1966; Hoskins et al. 1985; Schneider et al. 2003). This can be compared with, for instance, the three interior fields, e.g., vorticity, divergence, and temperature, required together with the surface pressure as initial conditions for the primitive equations, in the sigma-coordinate form used in numerical weather prediction.

2004) is totally suppressed by applying some artificial, layerwise-irrotational horizontal force per unit mass, ∇V say, where $V(\mathbf{x}, t)$ is a scalar field to be specified, \mathbf{x} is position in the fluid domain, and t is time. Layerwise refers either to the isentropic surfaces $\theta = \text{const}$ in a stratified system, where θ is potential temperature, or to the single horizontal layer of a shallow water system. The ∇ operator contains horizontal derivatives only, taken at constant θ in the stratified case. The model must take this form because a layerwise-irrotational, but otherwise arbitrary, force field is the most general way of forcing the primitive equations that conserves the PV of a fluid element and involves no mass or heat sources.

The simplest case is that of the f -plane shallow-water equations, which we take in their standard vorticity-divergence form. Then only the horizontal Laplacian $\Upsilon = \nabla^2 V$ of the field $V(\mathbf{x}, t)$ enters the analysis, via the divergence equation. The extra term Υ can, though need not, be regarded as a given forcing. The equations are

$$\frac{\partial \delta}{\partial t} = f\zeta - \nabla^2 \Phi - \nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v}) + \Upsilon, \quad (2.1a)$$

$$\frac{\partial \zeta}{\partial t} = -f\delta - \nabla \cdot (\mathbf{v}\zeta), \quad (2.1b)$$

$$\frac{\partial \Phi}{\partial t} = -gH\delta - \nabla \cdot (\mathbf{v}\Phi), \quad (2.1c)$$

where \mathbf{v} is the velocity, δ the horizontal divergence, ζ the vertical vorticity, Φ the shallow water surface geopotential anomaly, ∇ the horizontal nabla operator, f the constant Coriolis parameter, g the acceleration due to gravity, and H the area-mean layer depth; Φ is defined for convenience to have zero mean. Thus $g^{-1}\Phi$ is the free surface deviation from the mean (domain-averaged) height or layer depth H , and the PV is given by the exact formula of Rossby (1936),

$$Q = \frac{f + \zeta}{H + g^{-1}\Phi}. \quad (2.2)$$

This is a material invariant of (2.1); that is, with $D/Dt = \partial/\partial t + \mathbf{v} \cdot \nabla$ it satisfies

$$DQ/Dt = 0 \quad (2.3)$$

for any $\Upsilon(\mathbf{x}, t)$.

We use the term Upsilon model or Υ model to denote (2.1) and its stratified counterparts. There is an infinity of such models because of the infinity of ways of choosing $\Upsilon(\mathbf{x}, t)$. Most such models will not be balanced. That is, solutions will generically occupy a phase space as big as the phase space of the primitive equations and will

generically exhibit spontaneous-adjustment emission. The question of whether there are exceptional choices of Υ that can totally suppress spontaneous-adjustment emission is nontrivial, at least from the standpoint of physical intuition. Such suppression might be expected to require nondiagnostic information, as illustrated by the history integral appearing in the analysis of Ford et al. (2000) expressing the arrow of time inherent in the spontaneous wave-emission process. Nevertheless, and perhaps counterintuitively, the suppression is possible for a significantly wide range of choices of Υ . Indeed, one such choice is that leading to the well-known Bolin–Charney balance equations,

$$\Upsilon(\mathbf{x}, t) = \frac{\partial \delta}{\partial t} + \nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v}) - \nabla \cdot (\mathbf{v}_\psi \cdot \nabla \mathbf{v}_\psi), \quad (2.4)$$

reducing the divergence equation (2.1a) to the familiar approximation (Bolin 1955; Charney 1955, 1962),

$$\nabla^2 \Phi - f\zeta = -\nabla \cdot (\mathbf{v}_\psi \cdot \nabla \mathbf{v}_\psi) \quad (2.5a)$$

$$= 2J(u_\psi, v_\psi), \quad (2.5b)$$

where J is the horizontal Jacobian and $\mathbf{v}_\psi = (u_\psi, v_\psi)$ with $\mathbf{v}_\psi = \mathbf{z} \times \nabla \psi$, the nondivergent part of the horizontal velocity field \mathbf{v} in the standard Helmholtz decomposition

$$\mathbf{v} = \mathbf{v}_\psi + \mathbf{v}_\chi \quad (2.6a)$$

where $\mathbf{v}_\psi = \mathbf{z} \times \nabla \nabla^{-2} \zeta$ and $\mathbf{v}_\chi = \nabla \nabla^{-2} \delta$, (2.6b)

and where \mathbf{z} is a unit vertical vector. The notation ∇^{-2} signifies as usual that the streamfunction ψ and velocity potential χ satisfy $\nabla^2 \psi = \zeta$ and $\nabla^2 \chi = \delta$ with suitable boundary conditions. The purpose in approaching a standard set of equations via the nonstandard route (2.4) is to suggest that other exceptional choices of $\Upsilon(\mathbf{x}, t)$ are possible, pointing toward our generalization to the new hyperbalance equations.

To prepare for the generalization we need to recall, in addition, why replacing (2.1a) by (2.5) gives a balanced model in the standard generic sense, despite the two time derivatives remaining in (2.1b) and (2.1c). The reason is that (2.1b) and (2.1c) can be replaced by (2.3) together with an elliptic partial differential equation of the form

$$\mathcal{L}\delta = \mathcal{G}, \quad (2.7)$$

where \mathcal{L} denotes the modified Helmholtz operator

$$\mathcal{L} = gH\nabla^2 - f^2, \quad (2.8)$$

and where \mathcal{G} is a known function of Φ , ζ , and δ . Crucially, \mathcal{G} contains no time derivatives. That is, the Bo-

lin–Charney model is a balanced model because it can be described by the single prognostic equation (2.3) together with the definition (2.2) of the PV, the Helmholtz decomposition (2.6), and the two diagnostic equations (2.5) and (2.7).

The diagnostic equation (2.7), a refinement of the standard omega equation for vertical motion, can be derived through the following procedure (e.g., Charney 1962; Whitaker 1993). Again we state it in a way that will generalize. Indeed the following applies word-for-word to the shallow water and continuously stratified cases alike:

- Take $\partial/\partial t$ of the divergence equation after substituting for Υ . (2.9a)
- Eliminate $\partial\zeta/\partial t$ and $\partial\mathbf{v}_\psi/\partial t$ using (2.6) and the exact vorticity equation. (2.9b)
- Eliminate $\partial\Phi/\partial t$ using the remaining primitive equations, that is, using those other than the vorticity and divergence equations. (2.9c)

For the shallow water case, “the remaining primitive equations” means only (2.1c), the exact equation of local mass conservation. If we use (2.5a) in the first step (2.9a), then we find for the rhs of (2.7)

$$\begin{aligned} \mathcal{G} = & f\nabla \cdot (\mathbf{v}\zeta) - \nabla^2 \nabla \cdot (\mathbf{v}\Phi) \\ & + \nabla \cdot [\mathbf{v}_\psi \cdot \nabla \{ \mathbf{z} \times \nabla \nabla^{-2} [-f\delta - \nabla \cdot (\mathbf{v}\zeta)] \}] \\ & + \{ \mathbf{z} \times \nabla \nabla^{-2} [-f\delta - \nabla \cdot (\mathbf{v}\zeta)] \} \cdot \nabla \mathbf{v}_\psi. \end{aligned} \quad (2.10)$$

In virtue of (2.6), this expression (2.10) is a known function of Φ , ζ , and δ , as required. There are no time derivatives. Observe that the functional dependence is nonlocal because of the inverse Laplacians.

The functional form of (2.10) means that (2.5)–(2.8) define a balance condition, or balance relation, in the general sense of the term. That is, (2.5)–(2.8) together define a diagnostic functional relation between the mass field and the full velocity field; that is, between Φ and \mathbf{v} . Symbolically,

$$\mathbf{v}(\mathbf{x}, t) = \mathbf{v}_B\{\mathbf{x}; \Phi(\cdot)\}, \quad (2.11)$$

where the notation $\Phi(\cdot)$ denotes nonlocal functional dependence, in \mathbf{x} though not in t ; that is, (\cdot) should be read as shorthand for (\cdot, t) . As before, $\mathbf{v} = \mathbf{v}_\psi + \mathbf{v}_\chi$ from (2.6).

The balance relation (2.11) defined by (2.5)–(2.8) forms a closed system of equations when combined with just one scalar prognostic equation. This is a balanced model in the standard sense. The prognostic equation can, for instance, be the exact mass-conservation equation, or the exact vorticity equation, or the exact PV equation together with the definition of

PV. With any of these choices of prognostic equation the entire system of equations is closed, and can be time-marched starting with just one scalar field as initial condition. If we choose the prognostic equation to be the PV equation, including boundary delta functions where necessary, then by construction the system is a PBM because the exact PV is advected by the full velocity field $\mathbf{v} = \mathbf{v}_\psi + \mathbf{v}_\chi$.

3. Υ models in general

What other choices of Υ give balanced models in the standard sense? The answer can be given in a completely general way, applicable to the continuously stratified and variable- f cases. Note first that in all cases the modified divergence equation takes the form

$$\frac{\partial\delta}{\partial t} = f\zeta - \nabla^2\Phi - \nabla \cdot (\mathbf{v} \cdot \nabla\mathbf{v}) - (\mathbf{z} \times \mathbf{v}) \cdot \nabla f + \Upsilon. \quad (3.1)$$

In a stratified system Φ is the Montgomery potential, and ζ is Rossby’s isentropic vorticity $\mathbf{z} \cdot \nabla \times \mathbf{v}$, with derivatives taken at constant θ ; the PV is given by Rossby’s formula $Q = (f + \zeta)(-\partial\theta/\partial p)$ where p is pressure altitude. In all cases the function Φ completely specifies the mass field when use is made of the hydrostatic equation, the equation of state, and other thermodynamic relations as necessary (e.g., Hoskins et al. 1985, p. 900). By considering the procedure (2.9) we may straightforwardly show that, in order for the Υ model to be a balanced model, it is necessary and sufficient that Υ be chosen to have the functional form

$$\begin{aligned} \Upsilon(\mathbf{x}, t) = & \frac{\partial\delta}{\partial t} + \nabla \cdot (\mathbf{v} \cdot \nabla\mathbf{v}) \\ & + (\mathbf{z} \times \mathbf{v}) \cdot \nabla f + \mathcal{R}\{\mathbf{x}, t; \Phi(\cdot), \zeta(\cdot)\}, \end{aligned} \quad (3.2)$$

or a form in which \mathcal{R} has equivalent or reduced information content regarding the instantaneous state of the system. Here \mathcal{R} is an arbitrary, but prescribed, function of its arguments. Equivalent information content means for instance that $\zeta(\cdot)$ can be replaced by $\mathbf{v}_\psi(\cdot)$ since the first of (2.6b) allows us to regard the fields \mathbf{v}_ψ and ζ as containing the same information about the instantaneous state of the system. Any choice of the form (3.2) leads to a balanced model because, apart from the prognostic equation, the remaining equations are all diagnostic and are collectively equivalent to a balance relation (2.11). This is straightforward to verify by following the procedure (2.9). For continuously stratified systems the operator \mathcal{L} in (2.7) becomes a 3-dimensional elliptic operator.

Conversely, a form of \mathcal{R} with greater information content prevents (3.2) from defining a balanced model. Greater information content means that \mathcal{R} depends on δ or \mathbf{v}_χ as well as on some or all of Φ , ζ , and \mathbf{v}_ψ . The procedure (2.9) then produces a version of (2.7) containing time derivatives of δ , or equivalently of \mathbf{v}_χ , again straightforward to verify from (2.9). That is, (2.7) is no longer diagnostic.

Apart from the restriction on information content, we may choose $\mathcal{R}\{\mathbf{x}, t; \Phi(\cdot), \zeta(\cdot)\}$ in any way consistent with reasonable boundary conditions on the force potential $V = \nabla^{-2}\Upsilon$, for example, with zero domain area average for periodic boundary conditions. Of course many of the choices of $\mathcal{R}\{\mathbf{x}, t; \Phi(\cdot), \zeta(\cdot)\}$ are absurd in the sense that they cannot correspond to a usefully accurate balanced model, or absurd in the sense that they imply an explicitly time-dependent balance relation. The latter absurdity can be eliminated by further restricting $\mathcal{R}\{\mathbf{x}, t; \Phi(\cdot), \zeta(\cdot)\}$ to have no explicit time dependence; that is, $\mathcal{R} = \mathcal{R}\{\mathbf{x}; \Phi(\cdot), \zeta(\cdot)\}$, or forms with equivalent information content. With Υ chosen as in (3.2) and with the further restriction just mentioned, the divergence equation (3.1) takes the generic form

$$\nabla^2\Phi - f\zeta = \mathcal{R}\{\mathbf{x}; \Phi(\cdot), \zeta(\cdot)\}, \tag{3.3}$$

and (2.7) the generic form

$$\mathcal{L}\delta = \mathcal{G}\{\mathbf{x}; \Phi(\cdot), \zeta(\cdot), \delta(\cdot)\}. \tag{3.4}$$

We note for later convenience that in (3.3) another form of \mathcal{R} with equivalent or reduced information content is²

$$\mathcal{R} = \mathcal{R}\{\mathbf{x}; Q(\cdot), \zeta(\cdot)\}. \tag{3.5}$$

Now, if $\partial/\partial t$ of this \mathcal{R} is taken, we obtain terms involving $\partial Q/\partial t$. So the adaptation of the procedure (2.9) then required is to replace its third step, (2.9c), by the elimination of $\partial Q/\partial t$ using the exact equation (2.3) along with the elimination of $\partial\Phi/\partial t$ using the remaining primitive equations. Then \mathcal{G} in (3.4) acquires an additional argu-

² More precisely, the information content of (3.5) is equivalent to that on the right of (3.3) in all hydrostatic models for which the pressure p is zero at the top of the model. Otherwise, the information content is reduced. This follows from Rossby's PV formulae together with the hydrostatic equation, the equation of state, and other thermodynamic relations as necessary. For the shallow-water equations, Rossby's formula (2.2) with given Q and ζ immediately gives Φ . For continuous stratification, Rossby's formula $Q = (f + \zeta)(-\partial\theta/\partial p)$ gives only the $\partial p/\partial\theta$ field, which contains less information than the Φ field unless p is zero at the top of the model. In the latter case, it is a simple exercise to deduce p and Φ along with the other thermodynamic fields by vertical integration at fixed horizontal position (again see Hoskins et al. 1985, p. 900).

ment $Q(\cdot)$. This will be an important stepping stone toward the new hyperbalance equations.

We note in passing why the so-called modified balance equations of Charney (1962) and Moura (1976) cannot define a balanced model, in the standard sense used here. Their equations correspond to taking

$$\mathcal{R} = -\nabla \cdot (\mathbf{v}_\psi \cdot \nabla \mathbf{v}_\psi) - (\mathbf{z} \times \mathbf{v}) \cdot \nabla f. \tag{3.6}$$

The ∇f term introduces into \mathcal{R} a dependence on \mathbf{v}_χ and thus on δ , which represents greater information content than in (3.3) or (3.5) and causes time derivatives to appear in (3.4). As with the BEM of Allen et al. (1990) and Allen (1991), initialization requires more than one scalar field, and the model exhibits spurious high-frequency solutions (Moura 1976; Gent and McWilliams 1983b).

Finally, in the shallow-water f -plane case we note the explicit expressions for \mathcal{G} arising from (3.3) and its variant (3.5). With (3.3), the procedure (2.9) delivers (3.4) with

$$\mathcal{G} = f\nabla \cdot (\mathbf{v}\zeta) - \nabla^2\nabla \cdot (\mathbf{v}\Phi) - \partial\mathcal{R}/\partial t \tag{3.7a}$$

$$= f\nabla \cdot (\mathbf{v}\zeta) - \nabla^2\nabla \cdot (\mathbf{v}\Phi) - \frac{\mathcal{D}\mathcal{R}}{\mathcal{D}\Phi} \odot \frac{\partial\Phi}{\partial t} - \frac{\mathcal{D}\mathcal{R}}{\mathcal{D}\zeta} \odot \frac{\partial\zeta}{\partial t} \tag{3.7b}$$

$$= f\nabla \cdot (\mathbf{v}\zeta) - \nabla^2\nabla \cdot (\mathbf{v}\Phi) + \frac{\mathcal{D}\mathcal{R}}{\mathcal{D}\Phi} \odot [gH\delta + \nabla \cdot (\mathbf{v}\Phi)] + \frac{\mathcal{D}\mathcal{R}}{\mathcal{D}\zeta} \odot [f\delta + \nabla \cdot (\mathbf{v}\zeta)], \tag{3.7c}$$

where \mathcal{D} signifies variational differentiation and \odot the corresponding inner product, involving integration over the physical domain \mathbb{D} . For instance,

$$\frac{\mathcal{D}\mathcal{R}}{\mathcal{D}\Phi} \odot \frac{\partial\Phi}{\partial t} = \int_{\mathbb{D}} \frac{\mathcal{D}\mathcal{R}\{\mathbf{x}; \Phi(\cdot), \zeta(\cdot)\}}{\mathcal{D}\Phi(\mathbf{x}')} \frac{\partial\Phi(\mathbf{x}')}{\partial t} d^2\mathbf{x}'. \tag{3.8}$$

Observe from (3.7c) that \mathcal{G} does, indeed, have the required form, free of time derivatives.

If we start instead from (3.5) on the right of (3.3), then the $\mathcal{D}\mathcal{R}/\mathcal{D}\Phi$ term in (3.7c) is replaced by $\mathcal{D}\mathcal{R}/\mathcal{D}Q \odot (\mathbf{v} \cdot \nabla Q)$ so that (3.7c) is replaced by

$$\begin{aligned} \mathcal{G} &= \mathcal{G}\{\mathbf{x}; Q(\cdot), \Phi(\cdot), \zeta(\cdot), \delta(\cdot)\} \\ &= f\nabla \cdot (\mathbf{v}\zeta) - \nabla^2\nabla \cdot (\mathbf{v}\Phi) + \frac{\mathcal{D}\mathcal{R}}{\mathcal{D}Q} \odot (\mathbf{v} \cdot \nabla Q) \\ &\quad + \frac{\mathcal{D}\mathcal{R}}{\mathcal{D}\zeta} \odot [f\delta + \nabla \cdot (\mathbf{v}\zeta)]. \end{aligned} \tag{3.9}$$

Once again the functional form is free of time derivatives, and so once again we have a balanced model.

4. None of the known accurate PBMs are Υ models

We now consider the most accurate known PBMs, those delivering the astonishing accuracies already referred to even when Froude and Rossby numbers are not numerically small. They are all constructed by what we will call “high-order formal truncation,” or high-order truncation for brevity, extending the pattern of eliminations beyond those in the procedure (2.9).

In the extended procedure one takes the divergence equation (3.1) with $\Upsilon = 0$, together with a finite number of time derivatives $\partial/\partial t, \partial^2/\partial t^2, \dots, \partial^M/\partial t^M$ of that equation. For normal-mode formulations like that of MN00 one takes the fast-normal-mode equations and time derivatives thereof. One then truncates the system either by neglecting, or by otherwise getting rid of, the highest two time derivatives of δ , or the highest fast-normal-mode time derivatives. For $M \geq 2$ or its normal-mode equivalent (which is to take at least one time derivative of the fast normal-mode equations), it will prove convenient to call the procedure a high-order truncation.

A count of equations and variables shows that all remaining time derivatives can then, in principle, be eliminated without further truncation to give a closed diagnostic system defining a balance relation, in the sense of (2.11). That is, the system is closed in the sense that, given a balanceable mass field Φ , one can solve for everything else, including \mathbf{v} . Symbolically, $\mathbf{v} = \mathbf{v}_B\{\mathbf{x}; \Phi(\cdot)\}$ as before. This, in turn, gives us a balanced model whose single prognostic equation can again be taken as (2.3), making it into a PBM. As remarked in MN00 and MD01, the technique of high-order truncation was first proposed by K. H. Hinkelmann in the context of forecast initialization; the normal-mode counterpart was later introduced by B. Machenhauer, F. Baer, J. Tribbia, and others.

In practice, one does not carry out the eliminations, but leaves the truncated equations as a closed system of simultaneous partial differential equations in which, by implication, the remaining time derivatives have been replaced by auxiliary variables internal to the system. MN00 called these auxiliary variables “diagnostic estimates” of the time derivatives that they replace. To keep clear the distinction between, for instance, a true time derivative $\partial\delta/\partial t$ and a diagnostic estimate of it we use the special notation $\delta_A^{(1)}$ for the diagnostic estimate. The suffix A reminds us that $\delta_A^{(1)}$ is an auxiliary variable within the system. Similarly, $\delta_A^{(n)}$ will denote a diagnostic estimate of $\partial^n\delta/\partial t^n$.

Since the notional eliminations must make use of the mass and vorticity equations and their time derivatives, the system also contains auxiliary variables in the form of diagnostic estimates $\Phi_A^{(n)}$ and $\zeta_A^{(n)}$ of $\partial^n\Phi/\partial t^n$ and

$\partial^n\zeta/\partial t^n$. In the case of the f -plane shallow-water equations, for instance, the complete system before truncation is

$$\delta_A^{(n)} = f\zeta_A^{(n-1)} - \nabla^2\Phi_A^{(n-1)} - \nabla \cdot (\mathbf{v} \cdot \nabla\mathbf{v})_A^{(n-1)} \quad (n = 1, \dots, M + 1) \quad (4.1a)$$

$$\zeta_A^{(n)} = -f\delta_A^{(n-1)} - \nabla \cdot (\mathbf{v}\zeta)_A^{(n-1)} \quad (n = 1, \dots, M) \quad (4.1b)$$

$$\Phi_A^{(n)} = -gH\delta_A^{(n-1)} - \nabla \cdot (\mathbf{v}\Phi)_A^{(n-1)} \quad (n = 1, \dots, M) \quad (4.1c)$$

$$\mathbf{v}_A^{(n)} = \mathbf{v}_{A\psi}^{(n)} + \mathbf{v}_{A\chi}^{(n)} = \mathbf{z} \times \nabla\nabla^{-2}\zeta_A^{(n)} + \nabla\nabla^{-2}\delta_A^{(n)} \quad (n = 0, \dots, M), \quad (4.1d)$$

where for convenience we have renoted \mathbf{v} as $\mathbf{v}_A^{(0)}$ and similarly $\zeta, \delta, \mathbf{v}_\psi, \mathbf{v}_\chi$, and Φ as $\zeta_A^{(0)}, \delta_A^{(0)}, \mathbf{v}_{A\psi}^{(0)}, \mathbf{v}_{A\chi}^{(0)}$, and $\Phi_A^{(0)}$. The notation $(\cdot)_A^{(n)}$ outside a product stands for the Leibniz formula for the n th derivative of the product, with diagnostic estimates substituted. Thus, $(\mathbf{v} \cdot \nabla\mathbf{v})_A^{(1)} = \mathbf{v}_A^{(1)} \cdot \nabla\mathbf{v}_A + \mathbf{v}_A \cdot \nabla\mathbf{v}_A^{(1)}$ and $(\mathbf{v} \cdot \nabla\mathbf{v})_A^{(2)} = \mathbf{v}_A^{(2)} \cdot \nabla\mathbf{v}_A + 2\mathbf{v}_A^{(1)} \cdot \nabla\mathbf{v}_A^{(1)} + \mathbf{v}_A \cdot \nabla\mathbf{v}_A^{(2)}$, and so on; $\mathbf{v}_A^{(n)}$ is defined consistently with (2.6) and (4.1d) via

$$\mathbf{v}_A^{(n)} = \mathbf{v}_{A\psi}^{(n)} + \mathbf{v}_{A\chi}^{(n)} \quad (4.2a)$$

with

$$\mathbf{v}_{A\psi}^{(n)} = \mathbf{z} \times \nabla\nabla^{-2}\zeta_A^{(n)} \quad \text{and} \quad \mathbf{v}_{A\chi}^{(n)} = \nabla\nabla^{-2}\delta_A^{(n)}. \quad (4.2b)$$

For variable f we must add, respectively, $-(\mathbf{z} \times \mathbf{v}_A^{(n-1)}) \cdot \nabla f$ and $-\mathbf{v}_A^{(n-1)} \cdot \nabla f$ to the rhs of (4.1a), (4.1b). The simplest choice of truncation is that already mentioned,

$$\begin{cases} \delta_A^{(M)} = 0, \\ \delta_A^{(M+1)} = 0, \end{cases} \quad (4.3a) \quad (4.3b)$$

equivalent to neglecting the highest two time derivatives altogether. The complete system of equations, (4.1)–(4.3), forms a closed set when the mass field Φ is given. Other truncations are possible in place of (4.3). Two of them were thoroughly studied in MD01 and will be specified in the next section.

Closedness implies that all the diagnostic estimates are implicitly defined, by the system of equations, as nonlocal functions of the mass field Φ . We may therefore write

$$\delta_A^{(n)} = \delta_A^{(n)}\{\mathbf{x}; \Phi(\cdot)\} \quad (4.4)$$

and similarly for the other diagnostic estimates including $\mathbf{v} = \mathbf{v}_A^{(0)}$, which we can now identify, as before, with the function $\mathbf{v}_B\{\mathbf{x}; \Phi(\cdot)\}$ of (2.11).

The question now arises: do high-order truncations exist, in the foregoing sense, that can be represented by a choice of Υ of the form (3.2) or equivalent? That is, could any of the accurate PBMs in question be Υ models? We now prove that the answer is no, for all such high-order truncations including (4.3). That in turn implies that every PBM based on such a truncation exhibits velocity splitting, in the non-Hamiltonian sense already explained.

The proof is as follows. In every such PBM—recall that high-order means $M \geq 2$ —the truncation causes the divergence equation to be replaced by

$$\delta_A^{(1)} = f\zeta - \nabla^2\Phi - \nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v}) - (\mathbf{z} \times \mathbf{v}) \cdot \nabla f, \tag{4.5}$$

where \mathbf{v} denotes the full velocity field $\mathbf{v} = \mathbf{v}_\psi + \mathbf{v}_\chi$ as before. In place of the time derivative $\partial\delta/\partial t$ we have the diagnostic estimate $\delta_A^{(1)}$. In an Υ model, by contrast, the corresponding equation is (3.1); that is,

$$\frac{\partial\delta}{\partial t} = f\zeta - \nabla^2\Phi - \nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v}) - (\mathbf{z} \times \mathbf{v}) \cdot \nabla f + \Upsilon. \tag{4.6}$$

Suppose that a high-order truncation, in the foregoing sense, can be found such that the resulting PBM is also an Υ model. This leads to a contradiction. For, if the supposition were correct, then (4.5) for that PBM would have to agree with (4.6). In other words, we would have to choose

$$\Upsilon = \partial\delta/\partial t - \delta_A^{(1)}. \tag{4.7}$$

But then (3.2) and (4.7) imply that

$$\begin{aligned} \mathcal{R} &= -\delta_A^{(1)} - \nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v}) - (\mathbf{z} \times \mathbf{v}) \cdot \nabla f \\ &= \mathcal{R}\{\mathbf{x}; \Phi(\cdot), \zeta(\cdot), \delta(\cdot)\}, \end{aligned} \tag{4.8}$$

since $\delta_A^{(1)} = \delta_A^{(1)}\{\mathbf{x}; \Phi(\cdot)\}$ as already noted, and $\mathbf{v} = \mathbf{v}_\psi + \mathbf{v}_\chi$ as always. The dependence on δ violates the restriction on the information content of \mathcal{R} . Therefore, as with (3.6), the Υ model defined by (4.7) cannot be a balanced model. This contradicts the original supposition.

It follows that all PBMs defined by high-order truncations in the foregoing sense, including (4.3) and (5.1)–(5.2) below, do indeed exhibit non-Hamiltonian velocity splitting.

5. The $\delta\delta$, $\gamma\gamma$, and $\delta\gamma$ truncations

Besides (4.3), two other types of truncation will be of special interest. They are

$$\begin{cases} \delta_A^{(M)} = -\nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v})_A^{(M-1)} - (\mathbf{z} \times \mathbf{v}_A^{(M-1)}) \cdot \nabla f, & (5.1a) \\ \delta_A^{(M+1)} = -\nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v})_A^{(M)} - (\mathbf{z} \times \mathbf{v}_A^{(M)}) \cdot \nabla f & (5.1b) \end{cases}$$

and

$$\begin{cases} \delta_A^{(M)} = 0, & (5.2a) \\ \delta_A^{(M+1)} = -\nabla \cdot (\mathbf{v} \cdot \nabla \mathbf{v})_A^{(M)} - (\mathbf{z} \times \mathbf{v}_A^{(M)}) \cdot \nabla f. & (5.2b) \end{cases}$$

The truncations (5.2) were shown in MD01 to be precisely equivalent, in the f -plane case, to normal-mode truncations of various orders. The expressions on the right of (5.1) and (5.2b) were motivated by the idea of neglecting high time derivatives not of the divergence but of the ageostrophic vorticity $f^{-1}\gamma$ where $\gamma = f\zeta - \nabla^2\Phi$, the first two terms on the right of (3.1) (with $\Upsilon = 0$). We therefore refer to (5.1) as $\gamma\gamma$ truncations of various orders and, similarly, (5.2) as $\delta\gamma$ truncations and (4.3) as $\delta\delta$ truncations. Of course, “order” has no absolute meaning since we are dealing with formal rather than asymptotic truncations. For consistency with MN00 and MD01, we define the orders respectively as $M - 1$ for the $\gamma\gamma$, M for the $\delta\gamma$, and $M + 1$ for the $\delta\delta$ truncations. Notice that exactly the right number of diagnostic estimates is left free, in each case, to give closedness, as has been checked by the well-behaved numerical computations reported in MN00 and MD01.

Special cases of note include the case $M = 1$ of (5.1) (zeroth-order $\gamma\gamma$), which is the linear balance model discussed, for example, by Whitaker (1993). This is the same as the Bolin–Charney model except that $\mathcal{R} = 0$ in (3.3) and (3.7c). The case $M = 0$ of (5.2) (zeroth-order $\delta\gamma$) corresponds to the model introduced by Bleck (1974) and sometimes called the geostrophic PV (GPV) model (e.g., Mundt et al. 1997). The case $M = 1$ of (5.2) (first-order $\delta\gamma$) corresponds to the “slow equations” of Lynch (1989).

To prepare for the hyperbalance equations, let us slightly extend the notational conventions below (4.1) so that

$$\mathbf{v} = \mathbf{v}_A = \mathbf{v}_A^{(0)}, \tag{5.3}$$

and similarly for $\mathbf{v}_{A\psi}$ and $\mathbf{v}_{A\chi}$ and the other variables. As already pointed out for the $\delta\delta$ truncations, solving systems of the kind in question for a given mass field Φ determines all other variables, diagnostically, as nonlocal functions of Φ . In particular, we have $\mathbf{v} = \mathbf{v}_A\{\mathbf{x}; \Phi(\cdot)\} = \mathbf{v}_B\{\mathbf{x}; \Phi(\cdot)\}$ —defining a balance relation—as well as

$$\begin{aligned} \mathbf{v}_\psi &= \mathbf{v}_{A\psi}\{\mathbf{x}; \Phi(\cdot)\}, \quad \mathbf{v}_\chi = \mathbf{v}_{A\chi}\{\mathbf{x}; \Phi(\cdot)\}, \quad \text{and} \\ \delta_A^{(1)} &= \delta_A^{(1)}\{\mathbf{x}; \Phi(\cdot)\}, \end{aligned} \tag{5.4}$$

which will be referred to in the next section. For want of a better name, the balance relations $\mathbf{v} = \mathbf{v}_A = \mathbf{v}_B$ resulting from the $\delta\delta$, $\gamma\gamma$, and $\delta\gamma$ truncations will be

called “plain” $\delta\delta$, $\gamma\gamma$, and $\delta\gamma$ balance relations of various orders.

To convert any plain balance relation into a PBM we need to append to (4.1) not only the prognostic equation (2.3) and a pair of truncation equations but also an approximation to (2.2), the definition of Q , in a sense to be explained shortly. Now the diagnostic part of such a system defines a PV inversion operator. That is, given the Q field, the extended system can be solved for all the other variables, thus defining them as nonlocal functions of Q . We make this explicit by writing in place of (5.4)

$$\begin{aligned} \mathbf{v}_\psi &= \mathbf{v}_{AI\psi}\{\mathbf{x}; Q(\cdot)\}, \quad \mathbf{v}_\chi = \mathbf{v}_{AI\chi}\{\mathbf{x}; Q(\cdot)\}, \quad \text{and} \\ \delta_{AI}^{(1)} &= \delta_{AI}^{(1)}\{\mathbf{x}; Q(\cdot)\}, \end{aligned} \quad (5.5)$$

and similarly for $\delta_{AI}^{(n)}$, $\zeta_{AI}^{(n)}$, and all the other diagnostic estimates, including, now, $\Phi_{AI}^{(n)}\{\mathbf{x}; Q(\cdot)\}$. The suffix AI means that these are auxiliary variables, for the purposes of the next section, all computed by PV inversion. PV inversion operators of this kind will be called plain $\delta\delta$, $\gamma\gamma$, and $\delta\gamma$ inversion operators, and the corresponding PBMs will be called plain- $\delta\delta$, $-\gamma\gamma$, and $-\delta\gamma$ PBMs. The sense in which (2.2) is approximated is to replace it by

$$f + \zeta_{AI} = (Q + q_{AI})(H + g^{-1}\Phi_{AI}), \quad (5.6)$$

where q_{AI} is another auxiliary quantity in the form of a constant to be solved for as part of the inversion procedure. It is numerically small, the more so the higher the accuracy. However, for the plain PV inversion operators and PBMs under consideration, $q_{AI} \neq 0$ generically if we exclude unphysical model behavior such as H changing with time (violating global as well as local mass conservation) or the Kelvin circulation changing in the far field or around the boundary of a finite domain. This follows from the arguments given in Mohebalhojeh (2002)—noting that for given mass and circulation we are free to specify $Q(\mathbf{x})$ only up to an additive constant—together with the fact that the plain PBMs are not locally mass conserving, as was proved at the end of section 4. The small quantity q_{AI} represents a slight inconsistency that section 4 has shown to be unavoidable in the plain PBM dynamics, since under the model evolution q_{AI} is free to vary in time, even though it is constant in space, implying that DQ/Dt and $D(Q + q_{AI})/Dt$ cannot both be exactly zero. These difficulties are, however, overcome by the hyperbalance equations, as will now be shown.

6. The hyperbalance equations

Most of the following applies quite generally; that is, it applies to the shallow-water and stratified and to the constant- f and variable- f cases. Suppose we have aux-

iliary diagnostic estimates in the sense just established for the plain balance relations and plain PV inversion operators. In particular, the diagnostic estimates of $\partial\delta/\partial t$ and of \mathbf{v}_ψ and \mathbf{v}_χ are available as nonlocal functions either of Φ or of Q . The idea leading to the hyperbalance equations is to use these diagnostic estimates, either (5.4) or the plain PV inversions (5.5), to improve the accuracy of the divergence equation without violating the rules about information content established in section 3. The plain PV inversions are preferable for reasons of numerical well-conditionedness. We therefore use (5.5) in

$$\begin{aligned} \mathcal{R} &= -\delta_{AI}^{(1)} - \nabla \cdot [(\mathbf{v}_{AI\psi} + \mathbf{v}_{AI\chi}) \cdot \nabla(\mathbf{v}_{AI\psi} + \mathbf{v}_{AI\chi})] \\ &\quad - [\mathbf{z} \times (\mathbf{v}_{AI\psi} + \mathbf{v}_{AI\chi})] \cdot \nabla f. \end{aligned} \quad (6.1)$$

The accuracy of the divergence equation $\nabla^2\Phi - f\zeta = \mathcal{R}$ will now be comparable to the accuracy of whichever plain PV inversion is used to generate (5.5). The functional form of \mathcal{R} is now

$$\mathcal{R} = \mathcal{R}\{\mathbf{x}; Q(\cdot)\} \quad (6.2)$$

instead of (3.5). Dropping the dependence on $\zeta(\cdot)$ has reduced the information content of \mathcal{R} relative to that in (3.5). Therefore, the form (6.2) is permissible as a special case of (3.5), within the procedure (2.9) after modification as described below (3.5). The reduced information content implies in turn that (3.9) holds with its last term deleted. The deletion represents an immediate gain in computational economy.

Moreover, this is not the only such gain arising from the choice (6.1). The point will become clear after stating the resulting hyperbalance equations, which from here on will be referred to as the hyperbalance equations based on (5.5). They consist of the prognostic equation (2.3), that is, $\partial Q/\partial t = -\mathbf{v} \cdot \nabla Q$, together with the following set of diagnostic equations:

- (i) (6.1) defining $\mathcal{R}\{\mathbf{x}; Q(\cdot)\}$;
- (ii) (4.1) with subscripts A replaced by subscripts AI ; thus, e.g., $\mathbf{v}_{AI\chi}^{(n)} = \nabla\nabla^{-2}\delta_{AI}^{(n)}$;
- (iii) *either*: one of the three pairs of equations defining the truncations that close the system, (4.3) for $\delta\delta$, (5.1) for $\gamma\gamma$, or (5.2) for $\delta\gamma$, with suffixes A replaced by suffixes AI , *or*: a pair of equations defining some other truncation;
- (iv) $\mathcal{L}\delta = \mathcal{G}$ where \mathcal{G} is defined by (3.9), or its stratified counterpart, with the term in $\mathcal{D}\mathcal{R}/\mathcal{D}\zeta$ deleted and with $\zeta = \zeta_{AI}$ and $\Phi = \Phi_{AI}$;
- (v) $\mathbf{v} = \mathbf{v}_{AI\psi} + \nabla\nabla^{-2}\delta$ where δ satisfies (iv); this defines the velocity field \mathbf{v} that enters into (3.9) and advects the PV and mass fields;
- (vi) the exact definition of PV in the relevant form, for example, $Q = (f + \zeta_{AI})/(H + g^{-1}\Phi_{AI})$ for the shallow water equations.

This completes the definition of the hyperbalance equations or, rather, the simplest possible class of such equations. They will be solved numerically in Part II. The reader might wonder what has happened to the divergence equation $\nabla^2\Phi - f\zeta = \mathcal{R}$; the answer is that it is already hidden within (ii) above, as the first ($n = 1$) equation of (4.1a) with suffix A replaced by suffix AI and with the ∇f term restored if necessary, as noted below (4.2). The fact that we do not have to take $\mathbf{v}_\psi \neq \mathbf{v}_{AI\psi}$ and append a separate divergence equation represents a further gain in computational economy.

For brevity, when using one of the three truncations (4.3), (5.1), (5.2) we will refer to the corresponding locally mass conserving PBMs defined by the hyperbalance equations as the hyper- $\delta\delta$, hyper- $\gamma\gamma$, or hyper- $\delta\gamma$ PBMs of various orders, as the case may be, as distinct from the corresponding plain PBMs defined in sections 4 and 5.

The definition (v) of the advecting velocity is crucial to the local-mass-conserving property of the hyperbalance equations. It may be contrasted with the advecting velocity of a plain PBM:

$$\mathbf{v} = \mathbf{v}_{AI\psi} + \nabla\nabla^{-2}\delta_{AI} \quad (\text{plain}), \quad (6.3)$$

as distinct from

$$\mathbf{v} = \mathbf{v}_{AI\psi} + \nabla\nabla^{-2}\delta \quad (\text{hyperbalance}). \quad (6.4)$$

As already implied, the plain PBM is the dynamical system defined by (6.3) together with $\partial Q/\partial t = -\mathbf{v} \cdot \nabla Q$, (ii) and (iii) above, and (5.6) in place of (vi) above. In the hyperbalance case the velocity field (6.4) advects the mass field as well as the PV. That is, the local mass-conservation equation, (2.1c) in the shallow water case, is satisfied exactly when \mathbf{v} is given by (6.4). To verify this, one reverses the procedure (2.9) and uses exact PV conservation $\partial Q/\partial t = -\mathbf{v} \cdot \nabla Q$ to derive from (i)–(vi) an equation equivalent to operating on (2.1c) with an elliptic operator $\tilde{\mathcal{L}}$, which, in the shallow water case, is given by $\tilde{\mathcal{L}} = gH\nabla^2 - fHQ$. With reasonable boundary conditions this implies that (2.1c) is itself satisfied. From (ii), only (4.1a) with $n = 1$ is needed.³

In the plain case the mass is advected by a slightly different velocity field \mathbf{v}_m , which is defined only implicitly but which can be calculated at any time by inverting a Poisson equation. This is the non-Hamiltonian velocity splitting phenomenon, which is also the reason why

(vi) above must be replaced by (5.6) for a plain PBM. The fact that $\mathbf{v} - \mathbf{v}_m$ is nonzero, even if tiny, like q_{AI} in (5.6)—recall the astonishing accuracies attainable—follows from the result proved at the end of section 4 for all plain PBMs.

Within the diagnostic part of the hyperbalance equations, (i)–(vi) above, the fact that $\nabla\nabla^{-2}\delta \neq \nabla\nabla^{-2}\delta_{AI}$ can be regarded as a kind of internal velocity splitting. It is a necessary price to pay for the three desirables: exact PV advection, exact local mass conservation, and high accuracy.

In Part II, intercomparisons are made across all six types of PBM explicitly defined above. These are the hyper- $\delta\delta$, hyper- $\gamma\gamma$, hyper- $\delta\gamma$, plain- $\delta\delta$, plain- $\gamma\gamma$, and plain- $\delta\gamma$ PBMs of various orders. The numerical codes for the plain PBMs, already well checked from the work leading to MD01, are used as building blocks within the new codes for the hyperbalance PBMs. Because of this, we once again need to use (5.6) in place of (vi) above then let $q_{AI} \rightarrow 0$, to within rounding and truncation error, as iterations converge. The new codes have in turn been subjected to various careful checks. One of those checks came from the fact, easily verified from the above, that the first-order hyper- $\delta\delta$ PBM is the same as the Bolin–Charney model,⁴ for which a separate code was available and could be run. Another check, perhaps the most stringent of all, came from the direct testing of local mass conservation itself. Through careful error control and choice of numerical algorithm, we have been able to distinguish numerically between the satisfaction or violation of the local mass-conservation equation, that is, to make evident numerically the distinction between the plain and hyperbalance PBMs as regards local mass conservation. That would hardly have been possible without codes that are not only correct but also unusually accurate.

7. Concluding remarks

Historically, the question “can PV-conserving balanced models be locally mass conserving?” seems to have been first raised by Bleck (1974), who noted that inverting the Rossby–Ertel PV using geostrophic bal-

³ Specifically, one operates on this last equation with $\partial/\partial t$ then adds the result to (iv), recognizing that the terms in \mathcal{R} cancel because $(\mathcal{D}\mathcal{R}/\mathcal{D}Q) \odot (\mathbf{v} \cdot \nabla Q) = -(\mathcal{D}\mathcal{R}/\mathcal{D}Q) \odot (\partial Q/\partial t)$ in virtue of (2.3). Finally one uses (vi) to eliminate ξ_{AI} in favor of Φ_{AI} , remembering that $\Phi_{AI} = \Phi$ by (iv).

⁴ First-order hyper- $\delta\delta$, by definition, means based on first-order plain- $\delta\delta$ PV inversion. With our conventions this in turn means $M = 0$. The first truncation equation (4.3a) is to be read in accordance with the conventions (5.3)ff. and with suffixes A replaced by suffixes AI . Thus (4.3a) means that $\delta_{AI}^{(0)} = \delta_{AI} = 0$, and (4.3b) that $\delta_{AI}^{(1)} = 0$. Since $\{\partial/\partial t, \partial^2/\partial t^2, \dots, \partial^M/\partial t^M\}$ is an empty set when $M = 0$, (4.1a)–(4.1c) collapse to the divergence equation alone. It follows that the first-order hyper- $\delta\delta$ PBM is the same as the Bolin–Charney model.

ance in isentropic coordinates, then advecting PV by the resulting geostrophic velocity, leads to violation of local mass conservation. Bleck’s model is the same as the geostrophic PV (GPV) model of Mundt et al. (1997). Another precedent can be found in the work of Allen et al. (1990), especially their investigation of local and global conservation in the slow equations of Lynch (1989). Allen et al. showed that the velocity field that advects PV in that model also violates local mass conservation.

MN00 presumed that the plain- $\delta\delta$ and plain-normal-mode PBMs they studied also violate local mass conservation, albeit by a tiny margin in the most accurate cases, when they wrote:

It turns out to be simple to modify the second-order direct [i.e., plain] inversion operator to conserve mass locally [to give the Bolin–Charney model, see Corrigendum to MN00] but impossible, as far as we can see, to achieve any such modification at higher order.

But MN00 saw this as a virtue:

Whether local mass conservation is *desirable* as a property of balanced models is another question again. It is strongly arguable that enforcement of exact local mass conservation, and indeed energy and momentum conservation, would be likely to degrade the accuracy of a PV-conserving balanced model. In primitive-equation evolution, the spontaneous-adjustment emission of inertia–gravity waves—involving the spontaneous mutual adjustment of the mass and velocity fields within an unsteady, freely-evolving vortical flow (Ford et al. 2000, and references therein)—must modify the local mass, energy, and momentum budgets in ways that cannot be perfectly captured by a balanced model. Mass adjustments or rearrangements on the timescales of fast gravity-wave motion, in primitive-equation evolution, might be partially mimicked in an accurate balanced model as instantaneous mass rearrangements. By definition, such rearrangements require infinite velocities, and so cannot be exactly compatible with local mass conservation described by a velocity field that remains finite. Something has to give way.

On the strength of that argument, one of us went on to claim (McIntyre 2003), in an encyclopedia article discussing Hamiltonian and non-Hamiltonian balanced models [see also McIntyre (2001); Ford et al. (2002)], that “what gives way is the concept of a unique velocity field.” But it is now necessary to retract that claim in the non-Hamiltonian case. The discovery of the hyperbalance equations has shown the claim to be wrong, as regards formal order of accuracy at least. The ability to convert a plain PBM into the corresponding hyperbal-

ance PBM provides a way to enforce local mass conservation at any formal order of accuracy, while retaining a unique velocity field.

However, MN00’s argument still suggests that the plain-to-hyper conversion might degrade numerical accuracy. The question of numerical accuracy is addressed in Part II. It demands exceedingly careful and delicate numerical tests. An extensive set of such tests is described in Part II. They were carried out for complicated, chaotic vortex flows like those studied in Dritschel et al. (1999). The results strongly indicate that the claim, properly speaking conjecture, arising from MN00’s argument is wrong numerically also, for the shallow water equations at least. It appears that the highest accuracy attainable by plain PBMs is not significantly or systematically greater than that attainable by hyperbalance PBMs. Why this is so remains mysterious. Of course numerical tests can never cover all possible flows. Cases might yet be discovered in which the conjecture is supported. However, the results of Part II would appear to make this unlikely, for the shallow-water equations, since the flows studied evolve through a great variety of vortex configurations. It is possible, on the other hand, that a different result will be found when we go from shallow water to more than one layer. That has yet to be tested and remains an open question.

In any case, the hyperbalance equations have advantages from a purely theoretical standpoint. As well as removing the difficulties pointed out by Mohebalhojeh (2002), related to q_{AI} in (5.6) above, the ability to enforce local mass conservation without significant loss of accuracy means that we have in our possession what had long seemed an unattainable prize—a class of PBMs competitive with all known PBMs in terms of accuracy, yet having an exactly conserved potential enstrophy, as well as all the other non-Hamiltonian Casimir invariants. The Casimir invariants are defined in the same way as in the Hamiltonian case since their conservation depends only on exact local mass conservation and exact material PV conservation. For the stratified and shallow-water equations the general Casimirs take, respectively, the form of domain integrals

$$\int_{\mathbb{D}} F(Q, \theta) dm \quad \text{and} \quad \int_{\mathbb{D}} F(Q) dm \quad (7.1)$$

where $F(\cdot)$ is an arbitrary function, θ is the potential temperature as before, and dm is the mass element. For the shallow-water equations we have $dm = (H + g^{-1}\Phi) d^2\mathbf{x}$ and for the continuously stratified equations $dm = g^{-1}(-\partial p/\partial\theta)d^2\mathbf{x} d\theta$. The potential enstrophy is given by $F = Q^2$. These integrals are exact constants of the mo-

tion defined by the hyperbalance equations as well as of the motion defined by the primitive equations, as is easily verified.

We end with a reminder that Hamiltonian velocity splitting is a very different thing. For a thorough discussion the reader may consult McIntyre and Roulstone (2002, and references therein). In particular, such velocity splitting cannot be healed but is, inescapably, part of what “has to give way” in consequence of the relatively severe constraints imposed by the Hamiltonian structure, including the property of conserving energy and absolute momentum as well as mass.

A reviewer has raised the question of what might be called accuracy in the climate sense—of long-time behavior and statistics with and without local mass conservation. That must remain a question for future work, indeed a very difficult question. By comparison, the numerical experiments reported in Part II were less ambitious, having been primarily designed to test MN00’s conjecture on the effect of constraining local mass adjustment upon the accuracy of plain and hyperbalance PBMs. To this end, we paid attention only to measuring cumulative accuracy by pointwise comparison of the solutions of the PBMs and primitive equations over less extensive time scales, albeit of the order of many vortex rotations. A further question for the future is that of the existence and uniqueness of solutions of the hyperbalance equations. Here, nothing is known beyond the numerical evidence to be presented in Part II, which, however, does suggest that well-behaved solutions exist and are unique for the domain and the initial conditions considered.

Acknowledgments. We warmly thank E. Neven, W. A. Norton, S.-Z. Ren, and I. Roulstone, and our late and sadly missed colleague Rupert Ford, for helpful and wide-ranging discussions in the earliest stages of this work. We are grateful to the Isaac Newton Institute for Mathematical Sciences for support and further intellectual stimulation during the 1996 Programme on the Mathematics of Atmosphere and Ocean Dynamics. The original stimulus was, however, provided yet earlier by J. S. Allen, in a personal conversation in 1992 in which he raised the issue of local mass conservation in connection with highly accurate PV inversion. We are grateful to G. K. Vallis and two other reviewers, for suggestions that we hope have helped to make this difficult material more widely intelligible. A. R. M. thanks the Iran Ministry of Science, Research and Technology for support in the form of a research scholarship, the Universities of Tehran and St Andrews, and the U.K. Natural Environment Research Council for Research Fellowships and other support. M. E. M. thanks the

Engineering and Physical Sciences Research Council for crucial support in the form of a Senior Research Fellowship.

REFERENCES

- Allen, J. S., 1991: Balance equations based on momentum equations with global invariants of potential enstrophy and energy. *J. Phys. Oceanogr.*, **21**, 265–276.
- , J. A. Barth, and P. A. Newberger, 1990: On intermediate models for barotropic continental shelf and slope flow fields. Part I: Formulation and comparison of exact solutions. *J. Phys. Oceanogr.*, **20**, 1017–1042.
- Bleck, R., 1974: Short range prediction in isentropic coordinates with filtered and unfiltered numerical models. *Mon. Wea. Rev.*, **102**, 813–829.
- Bolin, B., 1955: Numerical forecasting with the barotropic model. *Tellus*, **7**, 27–49.
- Bretherton, F. P., 1966: Baroclinic instability and the short wavelength cut-off in terms of potential vorticity. *Quart. J. Roy. Meteor. Soc.*, **92**, 335–345.
- Charney, J. G., 1955: The use of the primitive equations of motion in numerical prediction. *Tellus*, **7**, 22–26.
- , 1962: Integration of the primitive and balance equations. *Proc. Int. Symp. on Numerical Weather Prediction*, Tokyo, Japan, Meteor. Soc. Japan, 131–152.
- Dritschel, D. G., L. M. Polvani, and A. R. Mohebalhojeh, 1999: The contour-advection semi-Lagrangian algorithm for the shallow water equations. *Mon. Wea. Rev.*, **127**, 1551–1565.
- Errico, R. M., 1982: Normal mode initialization and the generation of gravity waves by quasi-geostrophic forcing. *J. Atmos. Sci.*, **39**, 573–586.
- Ford, R., 1994a: Gravity wave radiation from vortex trains in rotating shallow water. *J. Fluid Mech.*, **281**, 81–118.
- , 1994b: The instability of an axisymmetric vortex with monotonic potential vorticity in rotating shallow water. *J. Fluid Mech.*, **280**, 303–334.
- , 1994c: The response of a rotating ellipse of uniform potential vorticity to gravity wave radiation. *Phys. Fluids*, **6**, 3694–3704.
- , M. E. McIntyre, and W. A. Norton, 2000: Balance and the slow quasimanifold: Some explicit results. *J. Atmos. Sci.*, **57**, 1236–1254.
- , —, and —, 2002: Reply. *J. Atmos. Sci.*, **59**, 2878–2882.
- Gent, P. R., and J. C. McWilliams, 1983a: Consistent balanced models in bounded and periodic domains. *Dyn. Atmos. Oceans*, **7**, 67–93.
- , and —, 1983b: The equatorial waves of balanced models. *J. Phys. Oceanogr.*, **13**, 1179–1192.
- Hoskins, B. J., M. E. McIntyre, and A. W. Robertson, 1985: On the use and significance of isentropic potential-vorticity maps. *Quart. J. Roy. Meteor. Soc.*, **111**, 877–946; *Corrigendum*, **113**, 402–404.
- Lighthill, M. J., 1952: On sound generated aerodynamically. I. General theory. *Proc. Roy. Soc. London*, **A211**, 564–587.
- Lynch, P., 1989: The slow equations. *Quart. J. Roy. Meteor. Soc.*, **115**, 201–219.
- McIntyre, M. E., 2001: Balance, potential-vorticity inversion, Lighthill radiation, and the slow quasimanifold. *Proceedings of the IUTAM/IUGG/Royal Irish Academy Symposium on Advances in Mathematical Modelling of Atmosphere and Ocean Dynamics*, P. F. Hodnett, Ed., Kluwer Academic, 45–68.

- , 2003: Balanced flow. *Encyclopedia of Atmospheric Sciences*, Vol. 2, J. R. Holton, J. A. Pyle, and J. A. Curry, Eds., Academic Press, 680–685.
- , and W. A. Norton, 2000: Potential vorticity inversion on a hemisphere. *J. Atmos. Sci.*, **57**, 1214–1235; Corrigendum, **58**, 949.
- , and I. Roulstone, 2002: Are there higher-accuracy analogues of semi-geostrophic theory? *Large-scale Atmosphere-Ocean Dynamics: Vol. II. Geometric Methods and Models (Proceedings of the Newton Institute Programme on Mathematics of Atmosphere and Ocean Dynamics)*, J. Norbury and I. Roulstone, Eds., Cambridge University Press, 301–364.
- Mohebalhojeh, A. R., 2002: On shallow-water potential-vorticity inversion by Rossby number expansions. *Quart. J. Roy. Meteor. Soc.*, **128**, 679–694.
- , and D. G. Dritschel, 2001: Hierarchies of balance conditions for the f -plane shallow water equations. *J. Atmos. Sci.*, **58**, 2411–2426.
- , and M. E. McIntyre, 2007: Local mass conservation and velocity splitting in PV-based balanced models. Part II: Numerical results. *J. Atmos. Sci.*, **64**, 1794–1810.
- Moura, A. D., 1976: The eigensolutions of the linearized balance equations over a sphere. *J. Atmos. Sci.*, **33**, 877–907.
- Mundt, M. D., G. K. Vallis, and J. Wang, 1997: Balanced models and dynamics for the large- and mesoscale circulation. *J. Phys. Oceanogr.*, **27**, 1133–1152.
- Norton, W. A., 1988: Balance and potential vorticity inversion in atmospheric dynamics. Ph.D. thesis, University of Cambridge, 167 pp.
- Rossby, C. G., 1936: Dynamics of steady ocean currents in the light of experimental fluid mechanics. *Pap. Phys. Oceanogr. Meteor.*, **5**, 1–43.
- Saujani, S., and T. G. Shepherd, 2002: Comments on “Balance and the slow quasimanifold: Some explicit results.” *J. Atmos. Sci.*, **59**, 2874–2877.
- Schneider, T., I. M. Held, and S. T. Garner, 2003: Boundary effects in potential vorticity dynamics. *J. Atmos. Sci.*, **60**, 1024–1040.
- Vanneste, J., and I. Yavneh, 2004: Exponentially small inertia-gravity waves and the breakdown of quasigeostrophic balance. *J. Atmos. Sci.*, **61**, 211–223.
- Warn, T., 1997: Nonlinear balance and quasigeostrophic sets. *Atmos.–Ocean*, **35**, 135–145.
- , and R. Ménard, 1986: Nonlinear balance and gravity-inertial wave saturation in a simple atmospheric model. *Tellus*, **38A**, 285–294.
- , O. Bokhove, T. G. Shepherd, and G. K. Vallis, 1995: Rossby number expansions, slaving principles, and balance dynamics. *Quart. J. Roy. Meteor. Soc.*, **121**, 723–739.
- Whitaker, J. S., 1993: A comparison of primitive and balance equation simulation of baroclinic waves. *J. Atmos. Sci.*, **50**, 1519–1530.