

Application of a New Adjoint Newton Algorithm to the 3D ARPS Storm-Scale Model Using Simulated Data

ZHI WANG

Center for Analysis and Prediction of Storms, University of Oklahoma, Norman, Oklahoma

KELVIN K. DROEGEMEIER

School of Meteorology and Center for Analysis and Prediction of Storms, University of Oklahoma, Norman, Oklahoma

L. WHITE

Department of Mathematics, University of Oklahoma, Norman, Oklahoma

I. M. NAVON

Department of Mathematics and Supercomputer Computations Research Institute, The Florida State University, Tallahassee, Florida

(Manuscript received 16 December 1996, in final form 5 March 1997)

ABSTRACT

The adjoint Newton algorithm (ANA) is based on the first- and second-order adjoint techniques allowing one to obtain the "Newton line search direction" by integrating a "tangent linear model" backward in time (with negative time steps). Moreover, the ANA provides a new technique to find Newton line search direction without using gradient information. The error present in approximating the Hessian (the matrix of second-order derivatives) of the cost function with respect to the control variables in the quasi-Newton-type algorithm is thus completely eliminated, while the storage problem related to storing the Hessian no longer exists since the explicit Hessian is not required in this algorithm. The ANA is applied here, for the first time, in the framework of 4D variational data assimilation to the adiabatic version of the Advanced Regional Prediction System, a three-dimensional, compressible, nonhydrostatic storm-scale model. The purpose is to assess the feasibility and efficiency of the ANA as a large-scale minimization algorithm in the setting of 4D variational data assimilation.

Numerical results using simulated observations indicate that the ANA can efficiently retrieve high quality model initial conditions. It improves upon the efficiency of the usual adjoint method employing the LBFGS algorithm by more than an order of magnitude in terms of both CPU time and number of iterations for test problems presented here. Numerical results also show that the ANA obtains a fast linear convergence rate.

1. Introduction

Four-dimensional variational data assimilation (Zupanski 1993b; Thépaut et al. 1993; Yang et al. 1996; Zupanski and Mesinger 1995, etc.) is widely acknowledged as one of the most promising approaches for implementing real-time analyzed data into an operational weather forecast system. In this approach, the numerical procedure attempts to generate both a close fit to the data and consistency with the dynamic model over a period of time. With present computer power, the only practical way to carry out 4D variational data assimilation is through an appropriate use of the so-called

adjoint of the assimilation model, which is used to calculate the gradient of the cost function with respect to control variables. The computational cost to operationally implement 4D variational data assimilation is still very expensive due to the computer power, the required large storage size of meteorology problems, and the inefficiency of large-scale unconstrained minimization algorithms. Therefore, it is very important to improve the existing algorithms or to develop new ones. While proper scaling and preconditioning (Zupanski 1993a, 1996) may significantly improve the efficiency of a large-scale unconstrained minimization algorithm, they are not the focus points of this work. This paper studies the application of the adjoint Newton algorithm (ANA) to 4D variational data assimilation problems using the Advanced Regional Prediction System (ARPS), a three-dimensional, compressible, nonhydrostatic storm-scale model.

Corresponding author address: Dr. Zhi Wang, 100 E. Boyd, EC 1110, Center for Analysis and Prediction of Storms, University of Oklahoma, Norman, OK 73019.
E-mail: dwang@tornado.gcn.uoknor.edu

The second-order adjoint technique provides information that is related to the Hessian (the matrix of second-order derivatives) of the cost function with respect to the control variables. One integration of the second-order adjoint model yields a Hessian vector product (Wang et al. 1992), which can be used to improve the efficiency of the truncated Newton algorithm (Wang et al. 1993; Wang et al. 1995a). The application of the second-order adjoint technique may also be used to derive the Newton descent direction, which leads to an adjoint Newton algorithm (Wang et al. 1997), which will be applied to the ARPS.

The adjoint Newton algorithm requires integrating a tangent linear model backward in time (Wang et al. 1997). We follow the approach of Pu et al. (1997) and approximate the integration of the tangent linear model backward in time by running the tangent linear model (TLM) with a negative time step, but reversing the sign of friction and diffusion terms, in order to avoid the computational instability that would be associated with these terms if they were run backward. Pu et al. (1996a,b; 1997) also show that the backward integration of a TLM can be used to improve the future forecast skill using past forecast error. They found that the forecast improvement obtained by the quasi-inverse linear method (integrating the TLM backwards in time) is considerably better than that obtained with a single adjoint integration and similar to the one obtained using five iterations of the adjoint method.

The diffusion process and the equation that describes it are irreversible. But it is not true that the solution of the diffusion equation cannot be reversed in time. Just as a film showing the diffusion process may be reversed in time, the equivalent numerical trick may be found to produce the same effect (Smolarkiewicz 1983). His technique (not used in this paper) may be used in the future.

Current popular methods for large-scale unconstrained minimization are

- 1) limited memory conjugate gradient method (Shanno and Phua 1980; Navon and Legler 1987);
- 2) quasi-Newton-type algorithms (Davidon 1959; Gill and Murray 1972; O'Leary 1983);
- 3) limited-memory quasi-Newton methods such as the LBFGS algorithm (Nocedal 1980; Liu and Nocedal 1989); and
- 4) truncated-Newton-type algorithms (Nash 1984a,b, 1985; Nash and Sofer 1989; Schlick 1992a,b; Navon et al. 1992b,c; Wang 1995a).

All those methods find the descent direction using the gradient of the cost function. The adjoint Newton algorithm provides a new approach to find the Newton descent direction by integrating a tangent linear model backward in time (with negative time steps). In section 2b, we will introduce the adjoint Newton algorithm, which is illustrated using a simple example in appendix B. If the tangent linear model can be accurately inte-

grated backward in time, the adjoint Newton algorithm has a quadratic convergence rate. For completeness, the theory of the adjoint Newton algorithm will be provided in appendix A. In the process of the proof of the adjoint Newton algorithm, we derive a relation between the first- and second-order adjoint variables. The ARPS as well as its adjoint will be briefly described in section 3. Also in section 3, different issues of backward integration of tangent linear models will be addressed. Cost function, weights, and scaling factors are addressed in section 4. Numerical results for different experiments are shown in section 5. There the ANA is compared with the usual adjoint method employing the LBFGS algorithm of Liu and Nocedal (1989) in terms of both CPU time and quality of the retrieved fields. Since the backward tangent linear model of the ARPS is not well posed and a modified version of it is used, the adjoint Newton algorithm achieves just a fast linear convergence rate (appendix C) for our test problems. Summary, conclusions, and limitations as well as topics for further research are presented in section 6.

In the general framework of 4D variational data assimilation, there are observation errors and/or model errors. Our first experiment (section 5) uses simulated (model generated) observations without error, while the second experiment uses simulated observations with random errors. Therefore the first experiment is not truly a 4D variational data assimilation problem but is carried out in a setting closely related to 4D variational data assimilation. Since this work represents the first step in the application of the adjoint Newton method using a three-dimensional model, we proceed by assuming that the model is perfect.

2. Method

a. Definitions

Of interest is the large-scale unconstrained minimization of a functional assuming the following form:

$$\min_U J(\mathbf{U}) = \min_U \left\{ \frac{1}{2} \int_{t_0}^{t_f} [\mathbf{C}\mathbf{X}(t) - \mathbf{X}^o(t)]^T \times \mathbf{W}[\mathbf{C}\mathbf{X}(t) - \mathbf{X}^o(t)] dt \right\}. \quad (1)$$

Here $\mathbf{X}: t \mapsto \mathbf{X}(t)$ is the state variable that is a function from $[t_0, t_f]$ into n -dimensional Euclidean space R^n . The operator $\mathbf{C}: R^n \mapsto R^m$ represents a projection from the space (R^n) of the model solution \mathbf{X} to the m -dimensional space (R^m) of observations. The linearity in \mathbf{C} is assumed for simplicity but is not required. The superscript T indicates a transpose. The components of \mathbf{X} are values of the various model fields (wind, temperature, pressure, etc.) at each of the model's grid locations. The number of components of \mathbf{X} is denoted by n , and m represents the number of observations at any given time. In gen-

eral, n is not equal to m . Here, $[t_0, t_f]$ denotes the assimilation window where t_0 and t_f are the initial and final time and \mathbf{W} is a positive definite and symmetric weighting function, which will be defined later. The cost function $J(\mathbf{U})$ is the weighted sum of squares of distance between the model solution and available observations distributed in space and time. The observation variable $\mathbf{X}^o: t \mapsto \mathbf{X}^o(t)$ is a function from $[t_0, t_f]$ into R^m . The control variable $U \in R^n$ and the state vector \mathbf{X} satisfy semidiscrete model equations, that is, a set of ordinary differential equations (usually nonlinear):

$$\frac{d\mathbf{X}}{dt} = F(\mathbf{X}), \tag{2}$$

$$\mathbf{X}(t_0) = \mathbf{U}, \tag{3}$$

where t is the time and F is model vector function of \mathbf{X} . For simplicity, F is assumed to be at least twice differentiable. The control to state mapping (from R^n to R^m) is given by $\mathbf{U} \mapsto \mathbf{X}(t, \mathbf{U})$.

It is important to realize that the control variables belong to a subset of the control space. To determine the admissible set of control, one may use physical information about the control variable. *For simplicity, only initial conditions are taken as control variables in this paper.* However, our method holds when the control variables are initial plus boundary conditions. It is not clear whether parameters (Wang 1993) other than initial and boundary conditions can be a part of the control variables vector for the ANA. This topic is the subject of further investigations.

Let $\mathbf{X}' = [D\mathbf{X}(\mathbf{U})]\mathbf{U}'$, where $\mathbf{U}' = \mathbf{X}'(t_0) \in R^n$ be a perturbation on the initial condition \mathbf{U} in Eq. (3) and $D\mathbf{X}(\mathbf{U}): R^n \mapsto R^m$ be the derivative of $\mathbf{X}(\mathbf{U})$ with respect to \mathbf{U} . Then the tangent linear model is defined as

$$\frac{d\mathbf{X}'}{dt} = DF(\mathbf{X})\mathbf{X}', \tag{4}$$

$$\mathbf{X}'(t_0) = \mathbf{U}'. \tag{5}$$

The first- and second-order adjoint models (Le Dimet and Talagrand 1986; Wang et al. 1992; Wang et al. 1995a) may be defined, respectively, as

$$-\frac{d\mathbf{P}}{dt} = [DF(\mathbf{X})]^T \mathbf{P} + \mathbf{C}^T \mathbf{W}(\mathbf{C}\mathbf{X} - \mathbf{X}^o), \tag{6}$$

$$\mathbf{P}(t_f) = \mathbf{0}, \tag{7}$$

$$-\frac{d\hat{\mathbf{P}}}{dt} = [DF(\mathbf{X})]^T \hat{\mathbf{P}} + [D^2F(\mathbf{X})\mathbf{X}']^T \mathbf{P} + \mathbf{C}^T \mathbf{W} \mathbf{C} \mathbf{X}', \tag{8}$$

$$\hat{\mathbf{P}}(t_f) = \mathbf{0}. \tag{9}$$

Here¹ \mathbf{P} and $\hat{\mathbf{P}} \in R^n$ are the first- and second-order

¹ The definition of adjoint variable \mathbf{P} differs from the definition in some references by a sign. We think that this notation is more consistent with our presentation.

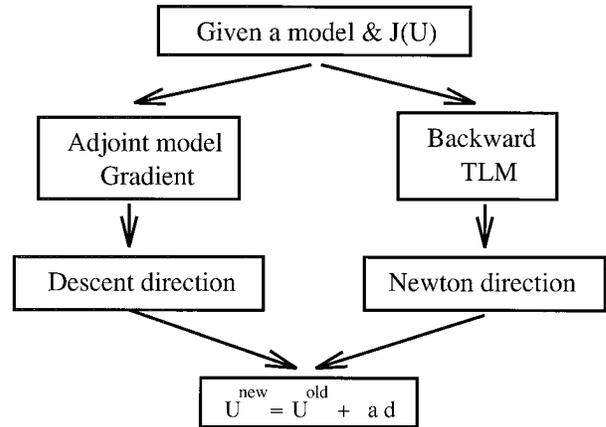


FIG. 1. Schematic presentations of adjoint and tangent linear approaches in finding descent directions. Here U^{new} , U^{old} , a , and d denote the new and old estimates for U , step length, and descent direction, respectively. The update step, $U^{new} = U^{old} + ad$, is not a part of the process of finding descent direction, but for completeness it is shown here.

adjoint variables; $D^2F(\mathbf{X}): R^n \times R^n \mapsto R^n$ denotes the second derivative of $F(\mathbf{X})$ with respect to \mathbf{X} .

It can be proven (Le Dimet and Talagrand 1986; Wang et al. 1992; Wang et al. 1995a) that

$$\mathbf{P}(t_0) = DJ(\mathbf{U}), \tag{10}$$

$$\hat{\mathbf{P}}(t_0) = D^2J(\mathbf{U})\mathbf{U}'. \tag{11}$$

b. Adjoint Newton algorithm

The current algorithms used to solve problem (1) first find its gradient by integrating the adjoint Eqs. (6) and (7) and then use the gradient to find a descent direction. There are other alternatives to calculate the gradient, such as the finite-differencing approach (although this is not practical at all for meteorology problems). Through the study of the first- and second-order adjoint techniques, we found that the solution of a backward integration of a tangent linear model with a “suitable final condition” yields “a descent direction,” which is comparable to the Newton descent direction. We will call this direction the estimated Newton descent direction (Wang et al. 1997). This approach is schematically presented in Fig. 1. Now we introduce the adjoint Newton algorithm here.

In the usual Newton algorithm (Berger 1977; Stoer and Bulirsch 1976), the function $J(\mathbf{U})$ being minimized is approximated locally by a quadratic function, and this approximated function is minimized exactly. Thus near \mathbf{U} we can approximate $J(\mathbf{U})$ by the truncated Taylor series (Luenberger 1984)

$$J(\mathbf{U} + \mathbf{d}) \approx J(\mathbf{U}) + DJ(\mathbf{U})\mathbf{d} + \frac{1}{2}\mathbf{d}^T D^2J(\mathbf{U})\mathbf{d}, \tag{12}$$

where $\mathbf{d} \in R^n$ is a descent direction. The minimum of

the right-hand side of Eq. (12) will be achieved if \mathbf{d} is the stationary point of the right-hand side of Eq. (12), that is, if \mathbf{d} satisfies Newton's equations

$$D^2J(\mathbf{U})\mathbf{d} = -DJ(\mathbf{U}). \quad (13)$$

Setting $\mathbf{U}'_N = -\mathbf{d}$, then Eq. (13) may be written as

$$D^2J(\mathbf{U})\mathbf{U}'_N = DJ(\mathbf{U}). \quad (14)$$

Our underlying assumption is that the Hessian matrix $D^2J(\mathbf{U})$ is positive definite and Eq. (14) is uniquely solvable. That is, the hypotheses for the convergence of Newton's method hold (Berger 1977; Stoer and Bulirsch 1976). The ANA amounts to a reformulation of Eq. (14) so that the estimated Newton descent direction is obtained using techniques related to the adjoint and there is no need to explicitly calculate or update the Hessian.

According to Eqs. (10) and (11), Eq. (14) implies there exists a unique initial condition \mathbf{U}'_N such that

$$\hat{\mathbf{P}}(t_0) = \mathbf{P}(t_0), \quad (15)$$

where $\hat{\mathbf{P}}(t_0) = D^2J(\mathbf{U})\mathbf{U}'_N$. We seek an approximation \mathbf{U}'_e to \mathbf{U}'_N . It turns out that by integrating a tangent linear model backward from a "suitable final condition," one can obtain just such an initial condition \mathbf{U}'_e . Indeed, if the tangent linear model can be integrated backward in time,² then (proof is given in appendix A) the estimated Newton descent direction is given by $\mathbf{d}_e = -\mathbf{Y}'(t_0) = -\mathbf{U}'_e$, where $\mathbf{Y}'(t_0)$ is the solution of the following backward problem:

$$\frac{d\mathbf{Y}'}{dt} = DF(\mathbf{X})\mathbf{Y}', \quad (16)$$

$$\mathbf{Y}'(t_f) = \mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t_f) - \mathbf{X}^o(t_f)]. \quad (17)$$

If \mathbf{U}_m is the minimum and

$$\mathbf{U}_m = \mathbf{U} - \mathbf{U}', \quad (18)$$

where \mathbf{U} is the current estimate, then the following hold:

$$\|\mathbf{U}' - \mathbf{U}'_m\| = O(\|\mathbf{U}'\|^2), \quad (19)$$

$$\|\mathbf{U}'_e - \mathbf{U}'_m\| = O(\|\mathbf{U}'\|^2), \quad (20)$$

$$\|\mathbf{U}'_e - \mathbf{U}'\| = O(\|\mathbf{U}'\|^2). \quad (21)$$

Here \mathbf{C} is assumed to be invertible and $\|\mathbf{U}\|$ denotes the Euclidean norm of \mathbf{U} .

We present the following remarks.

- 1) From Eq. (18), we know that $-\mathbf{U}'$ is the step to the minimum of $J(\mathbf{U})$ while Newton's direction $-\mathbf{U}'_N$ is the step to the minimum of the right-hand side of Eq. (12).
- 2) Since \mathbf{U}'_e is as good as \mathbf{U}'_N in the sense of Eqs. (19)–(21), we will also call $-\mathbf{U}'_e$ the Newton descent direction.
- 3) Since we only reformulate the Newton's equation,

Newton method's theoretical convergence rate should remain the same. Its convergence proof can be found in Luenberger (1984). A numerical convergence analysis will be provided in appendix C using the ARPS as an example.

- 4) For simplicity it is assumed that the tangent linear model may be integrated backward in time. If the nonlinear model is a mixture of hyperbolic and parabolic types, its corresponding tangent linear model cannot be directly integrated backward in time. Modifications have to be carried out such that the tangent linear model can be integrated backward in time. This issue will be addressed in section 3c.
- 5) The final condition (17) denotes the "forecast error" at t_f . The solution of the backward problem given by Eqs. (16) and (17) at $t < t_f$ is an approximation to the forecast error at t . That is, if the TLM can be accurately integrated backward in time and if there is no observation error, $\mathbf{Y}'(t) \approx \mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t) - \mathbf{X}^o(t)]$. Hence under these conditions, if the observations at $t < t_f$ are available, one may let $\mathbf{Y}'(t) = \mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t) - \mathbf{X}^o(t)]$ and then continue to integrate Eq. (16) to the initial time. Clearly, if the observations at t_0 are available, setting $\mathbf{d}_e = -\mathbf{Y}'(t_0) = -\mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t_0) - \mathbf{X}^o(t_0)]$ the ANA will find the minimum in one step since $\mathbf{d}_e = \mathbf{C}^{-1}\mathbf{X}^o(t_0) - \mathbf{X}(t_0)$ is the step to the minimum; that is, starting from an initial guess $\mathbf{U} = \mathbf{X}(t_0)$, one obtains $\mathbf{X}(t_0) + \mathbf{d}_e = \mathbf{C}^{-1}\mathbf{X}^o(t_0)$, which is the minimum.

In summary, it is assumed that all data are available at a single time t ($t_0 \leq t \leq t_f$) and that the backward integration of the TLM starts from that time.

- 6) The computational cost and storage for integrating the tangent linear model backward in time is similar to that required for integrating the first-order adjoint model.
- 7) The second-order adjoint model is used to derive the ANA but it is not used in the ANA. In the current version of the ANA, the gradient is only used to check the convergence criteria and line search conditions [see (A44), (A42), and (A43) in appendix A]. Convergence criteria and line search conditions without using the gradient information [such as golden section search, etc. (Luenberger 1984)] are under investigation. In other words, the next version of the ANA will not use the gradient of the cost function or the adjoint model.
- 8) The current ANA requires that the operator \mathbf{C} be invertible. In practice, it may not be invertible. For instance, when the observations are incomplete, the operator \mathbf{C} is not invertible in the classic sense. In the case of noninvertible \mathbf{C} , the generalized inverse of \mathbf{C} (see remark 1 following Theorem 2 in appendix A) may have to be used.
- 9) A search parameter α is introduced such that the ANA can be used at points that are relatively remote to the solution [see Eq. (A40) in appendix A and Luenberger (1984)]. Near the solution we expect, on

² The backward integration of a tangent linear model will be addressed in section 3c.

the basis of how Newton's method was derived, that $\alpha_k \approx 1$. Introducing the parameter for general points, however, guards against the possibility that the cost function might increase with $\alpha_k = 1$ (corresponding to the pure ANA), due to nonquadratic terms in the cost function (Luenberger 1984).

A simple example is provided in appendix B to illustrate the process of implementing the ANA.

3. Model descriptions

a. Nonlinear model

The ARPS is a three-dimensional, compressible, non-hydrostatic model developed for storm-scale prediction (Droegemeier et al. 1995; Xue et al. 1995). It is formulated in a curvilinear coordinate system that is orthogonal in the horizontal. The curvilinear coordinates can be defined numerically as well as analytically, making it more flexible than conventional terrain-following coordinates. The governing equations are the result of a direct transformation from the Cartesian system and are expressed in a fully conservative form.

The current adiabatic version of the ARPS includes

the Coriolis force, artificial divergence damping, total buoyancy, and subgrid turbulence mixing. The governing equations are discretized on an Arakawa C grid. Since the model atmosphere described by the governing equations is compressible, the meteorologically unimportant acoustic waves must be handled efficiently to avoid unnecessary restriction on the time step. The ARPS achieves this goal through the use of a splitting time integration technique reported by Klemp and Wilhelmson (1978). This technique divides a leapfrog integration time step into a number of small time steps and updates the acoustically active term every small time step while computing all the other terms only every leapfrog (big) time step. These acoustically active terms are the perturbation pressure gradient terms in the momentum equations and the divergence term in the pressure equation.

In the ARPS, a base state can be defined as either horizontally homogeneous or inhomogeneous. The thermal energy and pressure equations are written as prognostic equations for potential temperature and pressure. The adiabatic version of the ARPS with periodic boundary conditions is used in this study. Its governing equations are

$$\begin{aligned} \frac{\partial(\rho^*u)}{\partial t} = & - \left[(\rho^*u) \frac{\partial u}{\partial \xi} + (\rho^*v) \frac{\partial u}{\partial \eta} + (\rho^*W^c) \frac{\partial u}{\partial \zeta} \right] - \left[\frac{\partial(\bar{p}J_3)}{\partial \xi} + \frac{\partial(\bar{p}J_1)}{\partial \zeta} \right] \\ & - \left\{ \frac{\partial[J_3(p' - \alpha \text{Div}^*)]}{\partial \xi} + \frac{\partial[J_1(p' - \alpha \text{Div}^*)]}{\partial \zeta} \right\} + \rho^*fv - \rho^*\tilde{f}w + G^{1/2}D_u \end{aligned} \quad (22)$$

$$\begin{aligned} \frac{\partial(\rho^*v)}{\partial t} = & - \left[(\rho^*u) \frac{\partial v}{\partial \xi} + (\rho^*v) \frac{\partial v}{\partial \eta} + (\rho^*W^c) \frac{\partial v}{\partial \zeta} \right] - \left[\frac{\partial(\bar{p}J_3)}{\partial \eta} + \frac{\partial(\bar{p}J_2)}{\partial \zeta} \right] \\ & - \left\{ \frac{\partial[J_3(p' - \alpha \text{Div}^*)]}{\partial \eta} + \frac{\partial[J_2(p' - \alpha \text{Div}^*)]}{\partial \zeta} \right\} - \rho^*fu + G^{1/2}D_v \end{aligned} \quad (23)$$

$$\frac{\partial(\rho^*w)}{\partial t} = - \left[(\rho^*u) \frac{\partial w}{\partial \xi} + (\rho^*v) \frac{\partial w}{\partial \eta} + (\rho^*W^c) \frac{\partial w}{\partial \zeta} \right] - \frac{\partial[p' - \alpha \text{Div}^*]}{\partial \zeta} + \rho^*B + \rho^*\tilde{f}u + G^{1/2}D_w \quad (24)$$

$$\frac{\partial(\rho^*\theta)}{\partial t} = - \left[(\rho^*u) \frac{\partial \theta}{\partial \xi} + (\rho^*v) \frac{\partial \theta}{\partial \eta} + (\rho^*W^c) \frac{\partial \theta}{\partial \zeta} \right] + G^{1/2}D_\theta \quad (25)$$

$$\begin{aligned} \frac{\partial(G^{1/2}p')}{\partial t} = & - \left[(G^{1/2}u) \frac{\partial p'}{\partial \xi} + (G^{1/2}v) \frac{\partial p'}{\partial \eta} + (G^{1/2}W^c) \frac{\partial p'}{\partial \zeta} \right] - \left[(G^{1/2}u) \frac{\partial \bar{p}}{\partial \xi} + (G^{1/2}v) \frac{\partial \bar{p}}{\partial \eta} \right] + G^{1/2}W^c \bar{\rho}g \\ & - \frac{1}{\bar{\rho}c^2} \left[\frac{\partial(G^{1/2}u)}{\partial \xi} + \frac{\partial(G^{1/2}v)}{\partial \eta} + \frac{\partial(G^{1/2}W^c)}{\partial \zeta} \right], \end{aligned} \quad (26)$$

where

- variables with an overbar denote base-state quantities that are functions only of height;

- u , v , w , and W^c are two components of horizontal velocity, vertical velocity, and vertical contravariant velocity, respectively;

- $G^{1/2}$ (Droegemeier et al. 1995; Xue et al. 1995) is the determinant of the Jacobian matrix of the transformation from (ξ, η, ζ) system to (x, y, z) system; J_1 , J_2 , and J_3 are the nonconstant Jacobian of the transformation from (x, y, z) to (ξ, η, ζ) ;
- $\rho^* = G^{1/2}\bar{\rho}$, and $\bar{\rho}$ is the base-state density;
- p' and \bar{p} denote the perturbation and base-state pressure, respectively;
- αDiv^* is an artificial divergence damping term designed to attenuate acoustic waves;
- $f = 2\Omega \sin(\phi)$ and $f = 2\Omega \cos(\phi)$, where Ω is the angular velocity of the earth and ϕ is latitude;
- D_u, D_v, D_w , and D_θ are subgrid-scale turbulence mixing;
- B is the total buoyancy;
- θ is potential temperature; and
- \bar{c} is the sound speed.

The control variables of the cost function for our numerical tests are the perturbation variables at the initial time, defined as

$$\begin{aligned}
 u'(x, y, z, t_0) &= u(x, y, z, t_0) - \bar{u}(z) \\
 v'(x, y, z, t_0) &= v(x, y, z, t_0) - \bar{v}(z) \\
 w'(x, y, z, t_0) &= w(x, y, z, t_0) \\
 \theta'(x, y, z, t_0) &= \theta(x, y, z, t_0) - \bar{\theta}(z) \\
 p'(x, y, z, t_0) &= p(x, y, z, t_0) - \bar{p}(z). \quad (27)
 \end{aligned}$$

For simplicity, we drop the prime in subsequent discussion.

b. Model initialization

In this paper, the base state used follows that described by Weisman and Klemp (1982) and is horizontally homogeneous, hydrostatic, and time invariant. The computational domain extends 10 km in both the east-west and north-south directions, and 5 km in the vertical, with horizontal and vertical grid spacings of 1 km and 0.5 km, respectively. Convection is initiated with five buoyant thermal disturbances (“bubbles”) placed in the boundary layer. The maximum amplitude of these “bubbles” is 3 K. The centroid locations of the bubbles are $(x_c, y_c, z_c) = (5 \text{ km}, 5 \text{ km}, 1.5 \text{ km})$, $(3 \text{ km}, 3 \text{ km}, 1.5 \text{ km})$, $(7 \text{ km}, 3 \text{ km}, 1.5 \text{ km})$, $(7 \text{ km}, 7 \text{ km}, 1.5 \text{ km})$, and $(2 \text{ km}, 7 \text{ km}, 1.5 \text{ km})$. Above the ground, horizontal and vertical radius are $(x_r, y_r, z_r) = (3 \text{ km}, 3 \text{ km}, 1.5 \text{ km})$. The bubbles overlap. Starting with this initial condition, we integrate the ARPS for 15 min and use the output at this time as the initial conditions to obtain the model-simulated “observations” at the end of assimilation window (17 min). The initial fields at $z = 3.5 \text{ km}$ of v and p are shown in Fig. 2. All fields of v and p and difference fields between reference and retrieved fields will be shown at $z = 3.5 \text{ km}$ and we will mention this no more. The assimilation window extends from 15

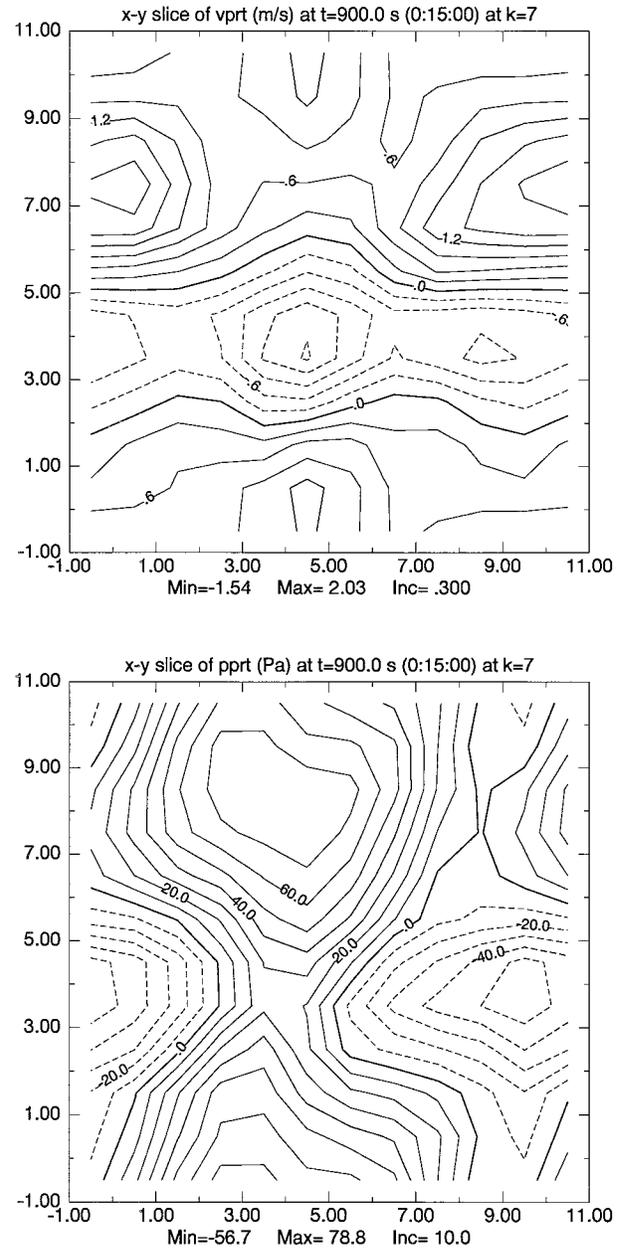


FIG. 2. Initial perturbation fields (reference fields) of v (top panel) and p (bottom panel) at $z = 3.5 \text{ km}$. Zero contour lines are kept since they clearly separate positive and negative contour lines.

to 17 min. The simulated observations are available only at the end of assimilation window. The simulated observations at the initial time will be used only to check the retrieved fields. The length of assimilation window is 2 min for our test problems since it is about the optimal assimilation length, which is obtained by looking at the graph (not shown) of rms errors between the retrieved and reference fields for different assimilation lengths. Sun (1992) discussed in detail the optimal assimilation length for a model similar to ours.

c. Tangent linear model and its backward integration

A TLM describes the evolution of perturbations in a forecast model. It uses approximations that are linear with respect to perturbations of the model fields. It is called a tangent model because the linearization is around a time-evolving solution, and, therefore, the coefficients of the linear model are defined by slopes of tangents to the nonlinear model trajectory in phase space (Errico et al. 1993). For the validity of a TLM, readers may consults Lacarra and Talagrand (1988), Vukićević (1991), and Errico et al. (1993). We will focus on the backward integration of a TLM for the purpose of the ANA application.

If a TLM can be integrated backward in time, given a perturbation (forecast error) at the end of forecast period, a backward integration of a TLM determines what perturbation (initial error), say at the initial time, causes the perturbation at the end of forecast period. At the end of the integration, the initial error (sensitivity) of the forecast error is obtained.

Unfortunately, in general, most TLMs may not be integrated backward in time due to the presence of irreversible processes such as mixing (diffusion) and divergence damping. However, simple modifications to TLMs lead to modified TLMs that can be integrated backward in time and yield approximations to the initial errors (sensitivities) of the forecast errors. In the following, we describe the backward integration of TLMs.

Most numerical prediction systems such as ARPS are a mixture of hyperbolic, elliptic, and parabolic systems although hyperbolic terms (advection terms) are dominant. Because of irreversible mixing and damping terms (parabolic in nature), their corresponding TLMs cannot be integrated backward in time. Pu et al. (1996a,b; 1997) suggested that the inverse TLM can be approximated by reversing the signs of mixing and damping terms, thus making the problem well posed. They denote the backward integration of such a modified TLM as the quasi inverse of the TLM. Reversing signs of mixing terms in a TLM leads to mixings backward in time in the quasi inverse of the TLM. The difference between the initial condition of a well-posed forward TLM and the solution at the initial time of the quasi inverse of the TLM starting from a solution of the well-posed TLM at a future time is hard to quantify since it is the difference of fields from two different models. Empirical results show that they are “close.” We will demonstrate this using the ARPS model. Pu et al. (1996a,b; 1997) introduced a similar approach of integrating the model backward in time to calculate the sensitivity of forecast error to changes in the initial conditions.

In the quasi inverse of the TLM of the ARPS, the artificial divergence damping terms and subgrid-scale turbulence mixing terms have opposite signs as those in the ARPS. Here, αDiv^* and $-G^{1/2}D_u$, $-G^{1/2}D_v$, $-G^{1/2}D_w$, and $-G^{1/2}D_\theta$ are used in the quasi inverse of the TLM instead of $-\alpha \text{Div}^*$, $G^{1/2}D_u$, $G^{1/2}D_v$, $G^{1/2}D_w$, and

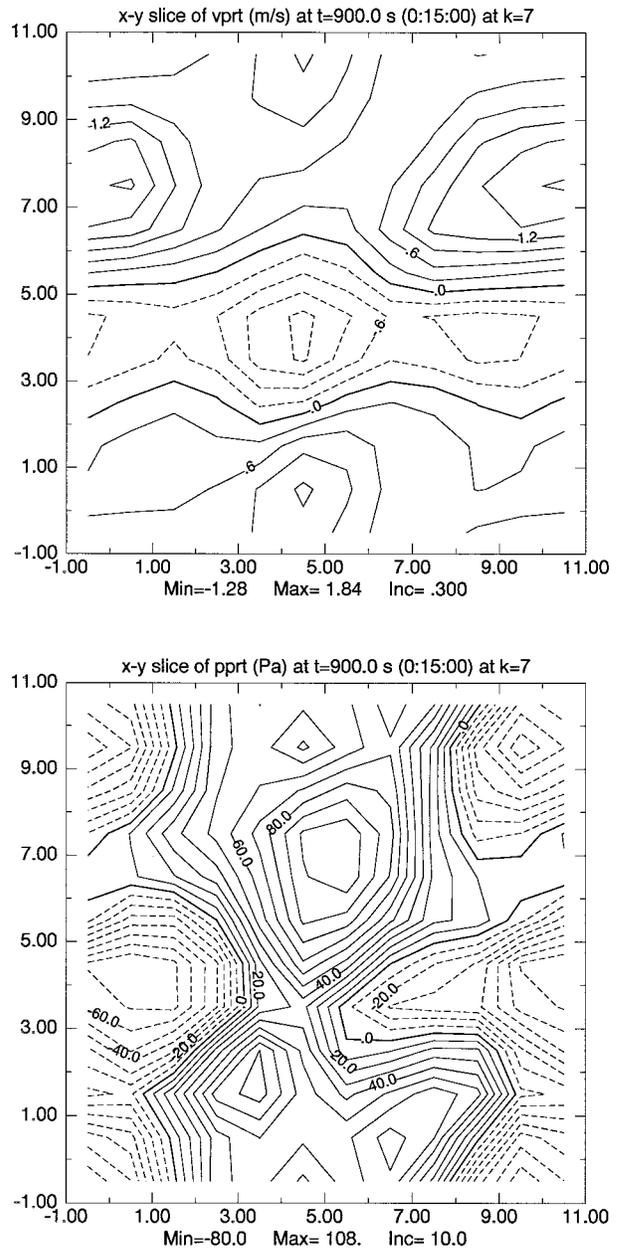


FIG. 3. Initial sensitivity fields of v (top panel) and p (bottom panel) at $z = 3.5$ km obtained by integrating the modified TLM backward in time.

$G^{1/2}D_\theta$ in the TLM of the ARPS. Let U' be the same initial fields (Fig. 2) used to obtain the “simulated observations.” The TLM solution fields at 17 min, $\mathbf{X}'(t_p)$, are used as the final condition for the quasi inverse of the TLM. The solution fields of v and p of the quasi inverse of the TLM at 15 min are shown in Fig. 3. Comparing Fig. 2 with Fig. 3, one concludes that the solution of the quasi inverse of the TLM of the ARPS approximately yields the initial perturbation fields. Although the initial perturbation fields thus obtained are

different from the exact fields corresponding to U' , their overall eddy structures, magnitudes, and signs are “similar” to those in the exact fields. The correlation coefficients of anomaly between the exact u , v , w , θ , and p at time t_0 and those obtained by the quasi inverse of the TLM are 0.9997, 0.9996, 0.9729, 0.9306, and 0.8275, respectively. This indicates that the quasi-inverse of the TLM of the ARPS yields meaningful and useful results “direction wise.”

The validity of the quasi inverse of a TLM may depend on the validity of the TLM itself. The study of the time length for which a quasi inverse of a TLM yields “meaningful” results is an important issue, which is under investigation.

d. Adjoint model development

The adjoint of the ARPS was developed (Wang et al. 1995b). The verification of the correctness of the gradient is conducted as follows. A Taylor series expansion applied in the direction $\mathbf{U}' = \mathbf{W}[\mathbf{X}(t_f) - \mathbf{X}^o(t_f)]$ yields (Navon et al. 1992a)

$$\phi(\alpha) = \frac{J(\mathbf{U} + \alpha\mathbf{U}') - J(\mathbf{U})}{\alpha\mathbf{U}'^T \cdot DJ(\mathbf{U})} = 1 + O(\alpha\|\mathbf{U}'\|). \quad (28)$$

For small α but not too close to machine accuracy, the value of $\phi(\alpha)$ should be close to unity if the gradient $DJ(\mathbf{U})$ is correctly calculated. Numerical results (not shown) indicate that this is exactly the case.

4. Cost function, weighting, and scaling

If the “simulated observations” are available at the beginning of the assimilation window without error, the ANA would find the minimum in one iteration since there is no need to integrate the quasi inverse of the TLM, whose solution, $\mathbf{X}(t_0) - \mathbf{X}^o(t_0)$, is known in this case and $\mathbf{X}^o(t_0) - \mathbf{X}(t_0)$ is the step to the minimum (see remark 5 in section 2b). Therefore, in order to test the robustness of the ANA, simulated observations are made available at the end of the assimilation window. For simplicity, we assume that the simulated observations are of the same dimension as the model state. It is worth noting that although the ANA is presented in a semidiscrete form, to carry out numerical calculations only the cost function needs to be defined as a sum instead of an integral. For instance, in our experiments, the cost function is defined as

$$J(\mathbf{U}) = \frac{1}{2}[\mathbf{X}(t_f) - \mathbf{X}^o(t_f)]^T \mathbf{W}[\mathbf{X}(t_f) - \mathbf{X}^o(t_f)], \quad (29)$$

where $\mathbf{X} = (u, v, w, p, \theta)$ is the state vector including all gridpoint values of model solution at a time and \mathbf{W} is weighting matrix. Since there are no statistics available for the ARPS, a constant diagonal weighting matrix $\mathbf{W} = \text{diag}(W_u, W_v, W_w, W_p, W_\theta)$ is used where

$$W_u^{-1} = \frac{1}{N} \sum_{i=1}^N \left\{ u_i(t_f) - u_i^o(t_f) - \frac{1}{N} \sum_{i=1}^N [u_i(t_f) - u_i^o(t_f)] \right\}^2, \quad (30)$$

and N is the number of u -gridpoint values. Model solution $\mathbf{X} = (u, v, w, p, \theta)$ is obtained by integrating the ARPS from zero initial guess. The weights thus defined nondimensionalize the cost function.

Scaling is a crucial issue in the success of nonlinear unconstrained optimization problems, and considerable research has been carried out on scaling nonlinear problems. It is well known that a badly scaled nonlinear programming problem can be almost impossible to solve (Navon and de Villiers 1983; Navon et al. 1992a). An effective automatic scaling procedure would ease these difficulties and could also render problems that are well scaled easier to minimize by improving condition number of their Hessian matrix (Thacker 1989).

Scaling by variable transformation converts variables from units that reflect physical properties to units that display desirable properties for the minimization process. Given a diagonal scaling matrix, $\mathbf{S} = \text{diag}(S_u, S_v, S_w, S_p, S_\theta)$, where S_u, S_v, S_w, S_p , and S_θ are constant diagonal submatrices, the general scaling procedure may be written as

$$\mathbf{X} = \mathbf{S}\mathbf{X}^s, \quad (31)$$

$$\mathbf{g}^s = \mathbf{S}\mathbf{g}, \quad (32)$$

$$\mathbf{H}^s = \mathbf{S}^T \mathbf{H} \mathbf{S}, \quad (33)$$

where \mathbf{H} is the Hessian matrix. The constant diagonal elements of submatrix S_u will be calculated by

$$S_u = \left\langle \frac{1}{N} \sum_{i=1}^N \left\{ u_i(t_f) - u_i^o(t_f) - \frac{1}{N} \sum_{i=1}^N [u_i(t_f) - u_i^o(t_f)] \right\}^2 \right\rangle^{1/2}, \quad (34)$$

and similarly for S_v, S_w, S_p , and S_θ .

In summary, the cost function is defined in terms of the original variables scaled by their standard deviation.

For complicated functions, difficulties may be encountered in choosing suitable scaling factors. Good scaling is problem dependent. A basic rule is that the variables of the scaled problem should be of similar magnitude and of order unity because within optimization routines convergence tolerances and other criteria are necessarily based on an implicit definition of “small” and “large,” and, thus, variables with widely varying orders of magnitude may cause difficulties of convergence for minimization algorithms (Gill and Murray 1981). One simple direct way to determine the scaling factor is to use the typical values for different fields.

This issue is closely related to the issue of preconditioning (Yang et al. 1996).

5. Numerical results

a. Experiment description

Two experiments are carried out to assess the feasibility and efficiency of the ANA as a large-scale minimization algorithm in the framework of 4D variational data assimilation. In all experiments, the length of the assimilation window is 2 min, and the big and small time steps are 6 and 1 s, respectively. The two experiments are as follows.

- 1) Simulated observations without error: The model-generated fields of u , v , w , θ , and p at 17 min (the end of the assimilation window) are used as simulated observations. *The initial guess is zero for all variables*, that is, for the three components of the 3D wind, pressure, and potential temperature fields. *It is emphasized that the simulated observations are available in all fields and at final time t_f only.*
- 2) Simulated observations with random errors: Same as experiment 1 except that at most 20% random error of each point is added to the simulated observations at the *final time*. Since random observation errors are just random noise, their distributions are not shown.

The ANA will be compared with the usual adjoint method employing the LBFGS algorithm of Liu and Nocedal [(1989), which will be called the LBFGS method for simplicity]. The full description of the LBFGS method may be found in Liu and Nocedal (1989). The same convergence criterion,

$$\|\mathbf{g}_{k+1}\| \leq 10^{-2}\|\mathbf{g}_0\|, \quad (35)$$

will be used for both the ANA and LBFGS methods. This convergence criterion yields a reasonable quality of the retrieved fields in terms of rms errors and correlation coefficients of anomaly relative to the reference fields (Fig. 2) of u , v , w , θ , and p at the initial time.

The gradient of the cost function with respect to control variables is calculated by integrating the adjoint of the ARPS. In the LBFGS method, the gradient is used to find the descent direction and to check the convergence criterion (35) and line search conditions (A42) and (A43) in appendix A, while it is only used to check the convergence criterion (35) and line search conditions (A42) and (A43) in the ANA. In order to be consistent, for all experiments the ANA and the LBFGS method of Liu and Nocedal (1989) are fixed. No fine-tuning is allowed. The number of updates in the LBFGS method is set to 5 at where the LBFGS algorithm performs best for our test problems. Computations were performed using the Cray C90 at the Pittsburgh Supercomputing Center.

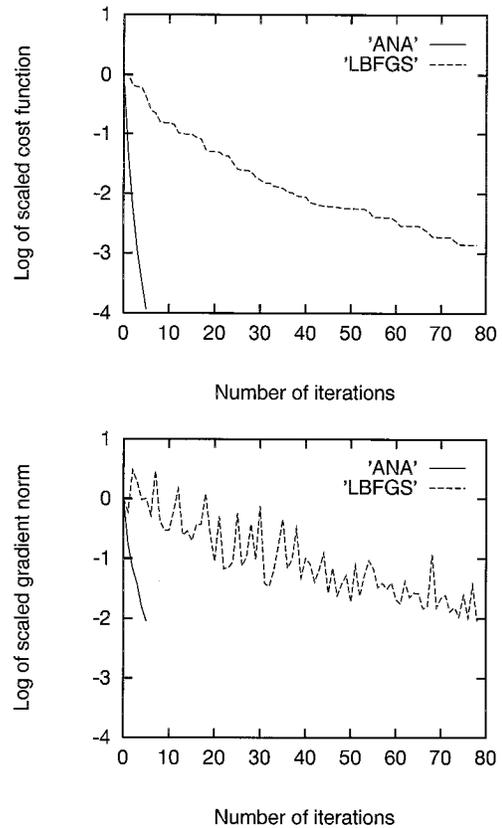


FIG. 4. Variations of the log of the scaled cost function (J/J_0 , top panel) and scaled gradient norm ($\|\mathbf{g}\|/\|\mathbf{g}_0\|$, bottom panel) with the number of iterations using algorithms: adjoint Newton algorithm (ANA) and LBFGS algorithm, respectively, in experiment 1.

b. Results and comparison of the ANA and LBFGS algorithms

For experiment 1, ANA requires 5 iterations, 11 function calls, and 26.847 s of CPU time to satisfy the convergence criterion (35), while the LBFGS method requires 78 iterations, 224 function calls, and 431.131 s to satisfy the same convergence criterion (35). Therefore the ANA is more than an order of magnitude faster than the LBFGS method in terms of both number of iterations and CPU time in this experiment. The variation of the cost function scaled by its initial value (J/J_0) as well as those of the norm of the gradient also scaled by its initial value ($\|\mathbf{g}\|/\|\mathbf{g}_0\|$) as functions of the number of iterations are displayed in Fig. 4, which indicates that the cost function obtained using ANA decreases one order of magnitude more than that obtained by using LBFGS method. The quality of the retrieved fields obtained by using ANA is also much better than that obtained by using the LBFGS method, which can be seen by comparing (given below) the retrieved fields with reference fields.

Figures 5 and 6 show that the retrieved fields obtained using ANA are almost identical to the reference fields (Fig. 2), while those obtained by using the LBFGS meth-

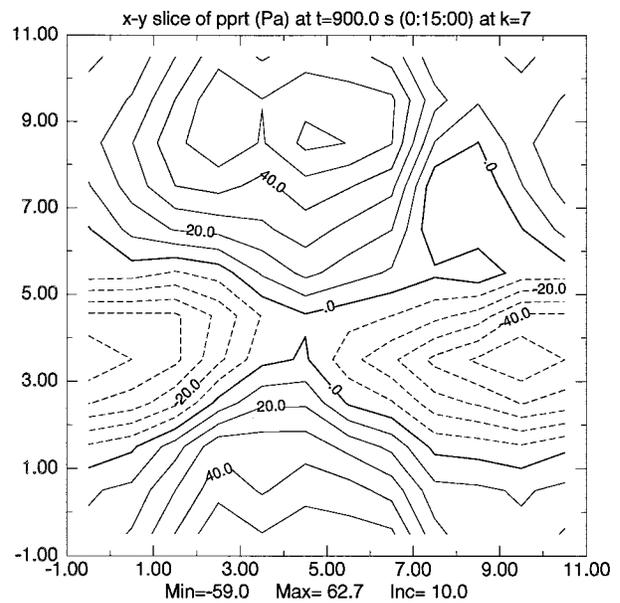
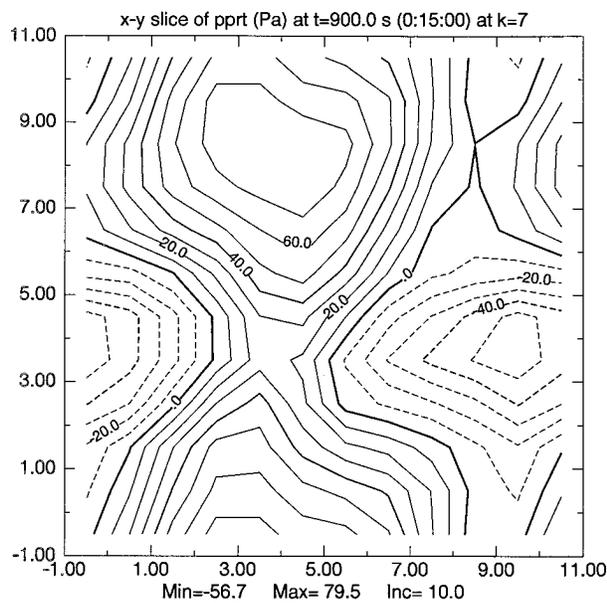
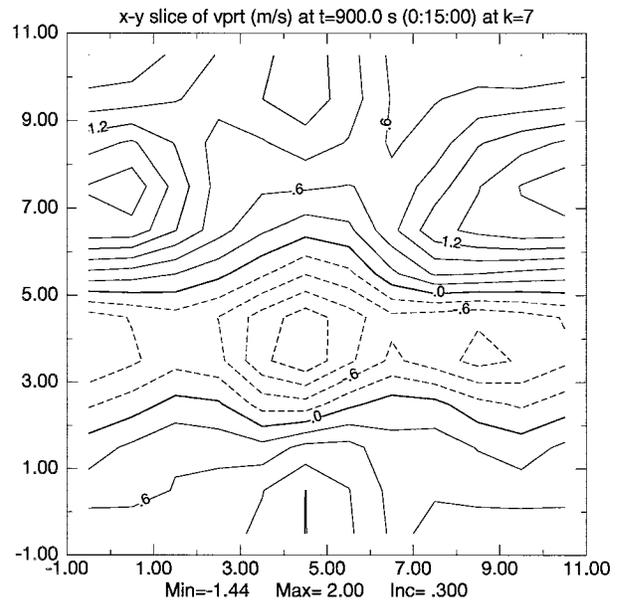
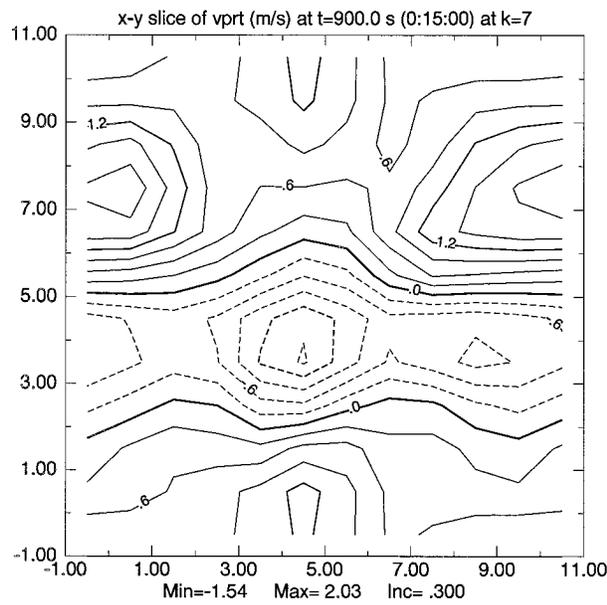


FIG. 5. Retrieved fields of v (top panel) and p (bottom panel) at the initial time and $z = 3.5$ km using the ANA in experiment 1.

FIG. 6. Same as Fig. 5 except LBFGS is used.

od display clear differences with the reference fields (Fig. 2). Difference fields of v and p between the retrieved and reference fields are displayed in Figs. 7 and 8, respectively. After the minimization, the maximum difference values of v and p obtained using the ANA are about an order of magnitude smaller than those obtained using the LBFGS method. The difference fields of u , w , and θ are similar (not shown).

Table 1 shows the rms errors between the retrieved u , v , w , θ , and p and their corresponding reference fields at the end of assimilation using both ANA and LBFGS

methods, while Table 2 shows the corresponding correlation coefficients of anomaly. Table 1 clearly indicates that the rms errors in all fields obtained using ANA are an order of magnitude smaller than those obtained using the LBFGS method. The correlation coefficients of the anomaly (Table 2) obtained using ANA are larger than those obtained using the LBFGS method in all fields.

In experiment 2 we test the feasibility and efficiency of the ANA in the case of "simulated observations with random errors." As in experiment 1, the ANA is very robust compared to the LBFGS method. It requires 7

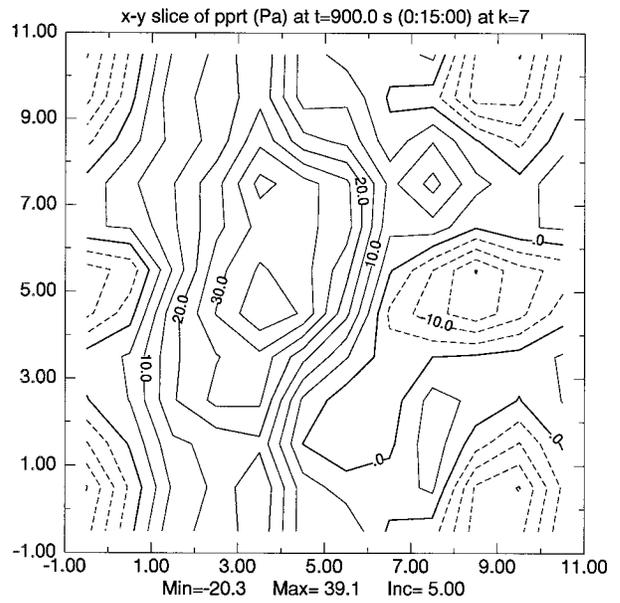
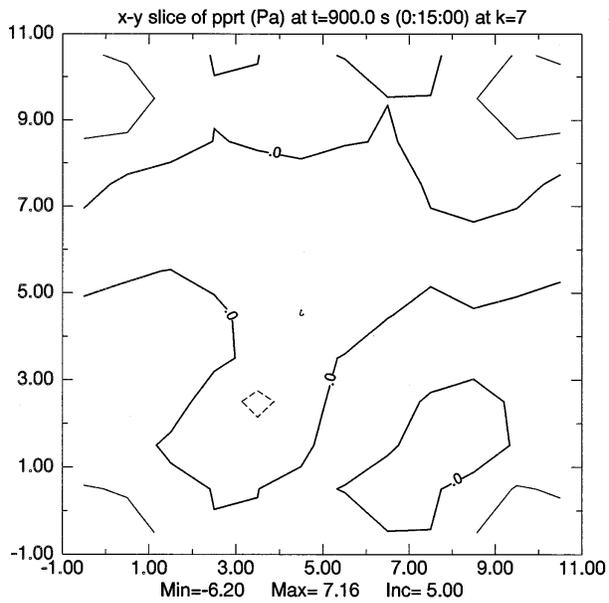
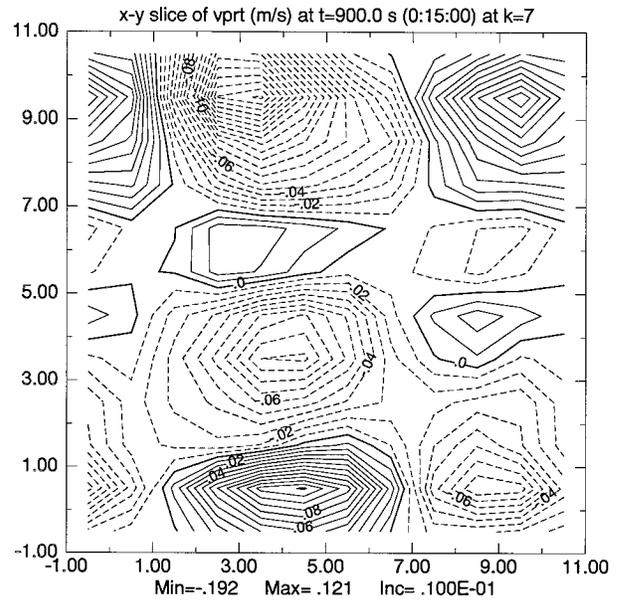
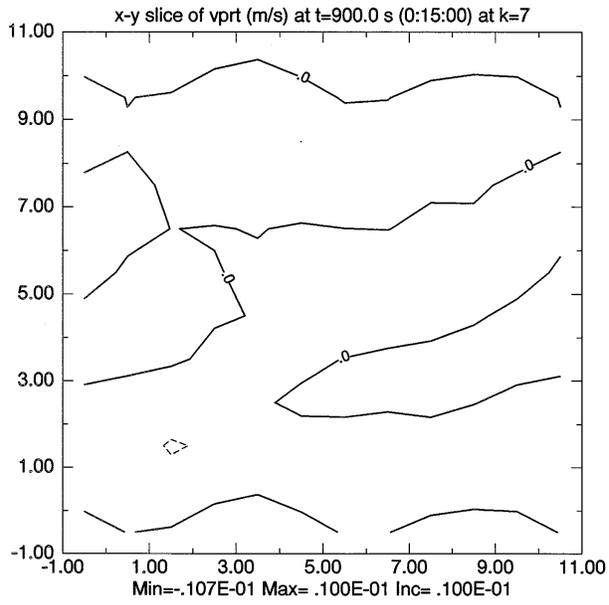


FIG. 7. Difference fields between the retrieved fields of v (top panel) and p (bottom panel) at $z = 3.5$ km using the ANA method and their corresponding reference fields in experiments 1, respectively.

FIG. 8. Same as Fig. 7 except LBFSGS is used.

TABLE 1. The rms errors at the initial time between the retrieved u , v , w , θ , and p fields and corresponding reference fields normalized by their initial values in experiments 1 and 2, respectively. "Iter." denotes number of iterations.

Experiments	Algorithms	rms in u	v	w	θ	p	Iter.	CPU (s)
1	ANA	0.0054	0.0045	0.0090	0.0026	0.0463	5	26.847
1	LBFSGS	0.1340	0.0619	0.0760	0.0314	0.2344	78	431.131
2	ANA	0.0560	0.0421	0.0973	0.0485	0.2076	7	36.556
2	LBFSGS	0.1125	0.0644	0.1133	0.0530	0.2451	108	568.006

TABLE 2. Correlation coefficients of anomaly at the initial time for u , v , w , θ , and p between the retrieved and corresponding reference fields for experiments 1 and 2, respectively. "Corr." denotes correlation coefficient.

Experiments	Algorithms	Corr. in u	v	w	θ	p
1	ANA	0.99998	0.99999	0.99996	1	0.99892
1	LBFGS	0.99292	0.99792	0.99756	0.99966	0.97411
2	ANA	0.99825	0.99903	0.99529	0.99876	0.97966
2	LBFGS	0.99406	0.99776	0.99413	0.99857	0.97256

iterations, 15 function calls, and 36.556 s of CPU time to satisfy the convergence criterion (35), while the LBFGS method requires 108 iterations, 295 function calls, and 568.006 s to satisfy the same condition. At the end of assimilation, the cost function obtained using ANA decreases by one order of magnitude more than that obtained by using the LBFGS method (Fig. 9).

The rms errors and correlation coefficients of anomaly in Tables 1 and 2 indicate that the ANA yields more accurate results than the LBFGS method does, that is, smaller rms errors and larger correlation coefficients of anomaly in all fields, but the differences are smaller comparing with those in experiment 1.

Figures 10 and 11 show that the retrieved fields obtained using the ANA are a little closer to the reference fields (Fig. 2) than those obtained by using the LBFGS method. After the minimization, the difference between the maximum and minimum values of v and p in the

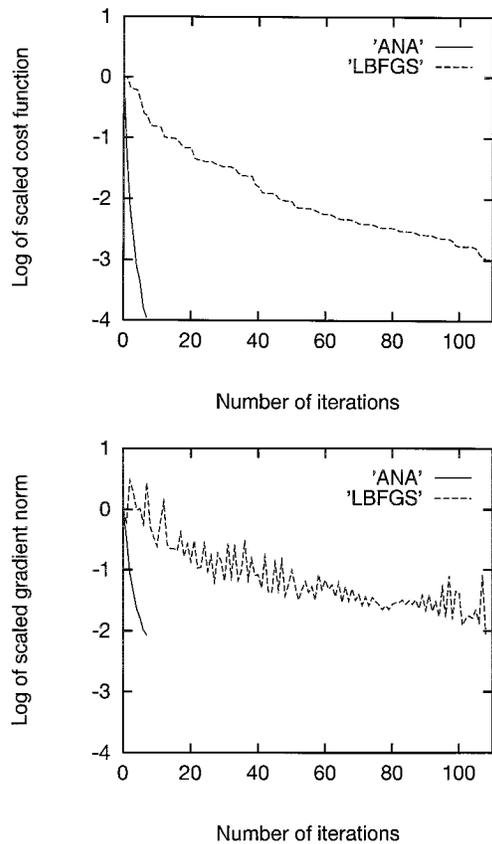


FIG. 9. Same as Fig. 4 except they are for experiment 2.

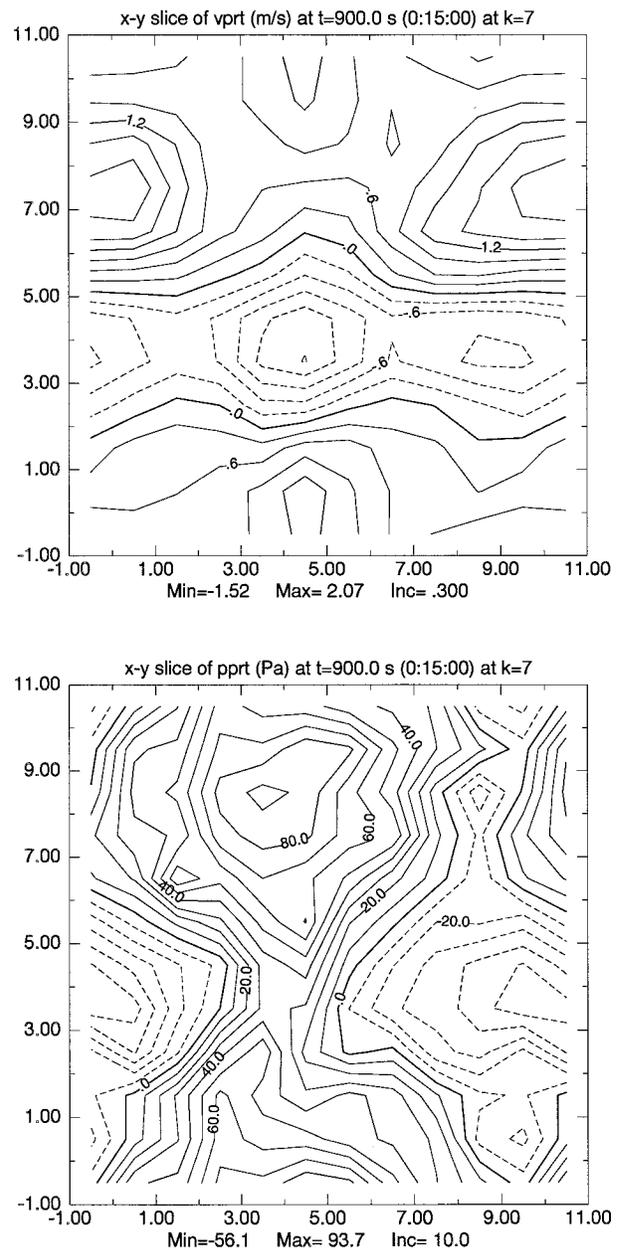


FIG. 10. Same as Fig. 5 except they are for experiment 2.

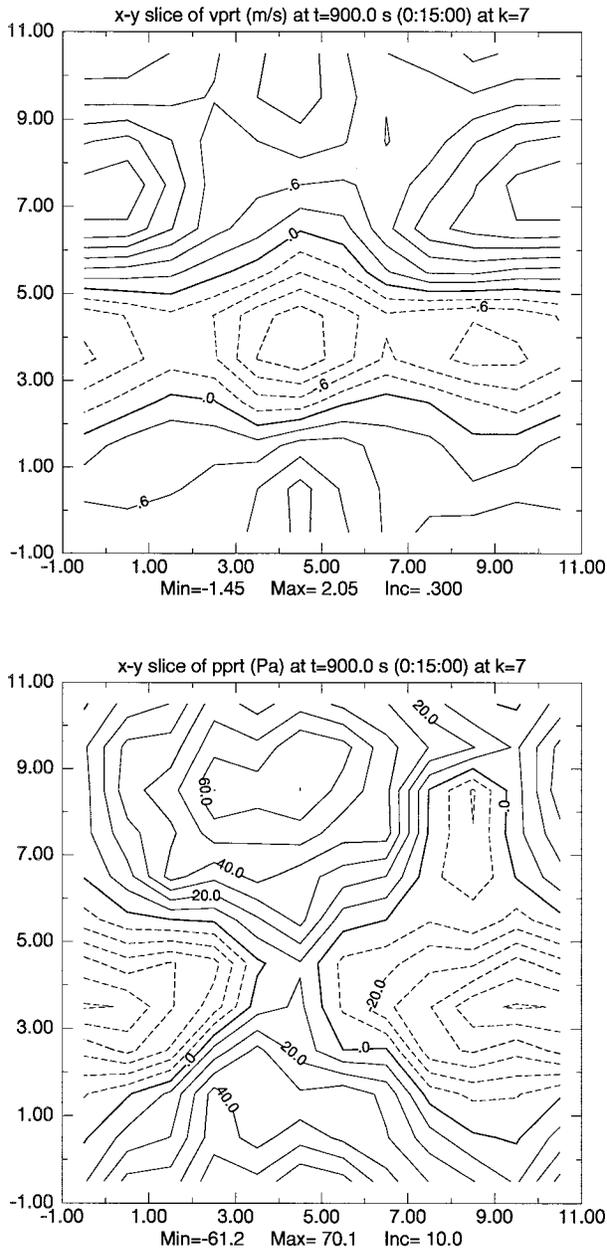


FIG. 11. Same as Fig. 6 except they are for experiment 2.

difference fields of v and p between the retrieved and reference fields are 0.266 m s^{-1} and 59.8 Pa when the ANA is used and 0.354 m s^{-1} and 62.9 Pa when the LBFGS method is used. The quality of the retrieved fields in this case is similar, although the ANA yields slightly better results (Tables 1 and 2).

6. Summary and conclusions

A new ANA suitable for 4D variational data assimilation was applied to the ARPS. The new ANA finds the Newton descent direction by integrating the quasi

inverse of a TLM backward in time. Since most TLMs of numerical weather prediction models are not well posed when they are integrated backward in time due to mixing and damping effects, they have to be modified by reversing the signs of the mixing and damping terms such that the quasi inverses of the TLMs are well posed when they are integrated backward in time (Pu et al. 1996a, b; Pu et al. 1997). The solution of a quasi inverse of a TLM at the initial time is only an approximation to Newton descent direction, which leads to a fast linear convergence rate of the ANA for our experiments. Numerical experiments using the ARPS indicate that the ANA is suitable for 4D variational data assimilation in the setting where an analysis is used as observations. It is very efficient in terms of both the number of iterations and CPU times for our test problems. Therefore it may be a very promising and efficient alternative to large-scale unconstrained minimization algorithms for meteorology problems. Compared to the LBFGS method of Liu and Nocedal (1989), the ANA is a clear winner for our test problems both in terms of efficiency as well as in terms of quality of retrieved fields.

The current ANA has not considered how to deal with

- 1) parameters other than initial and boundary conditions as part of control variables;
- 2) noninvertible operator \mathbf{C} —this limitation can be avoided by reformulating the ANA, as pointed out in the first remark following Theorem 2 in appendix A;
- 3) how the validity of the backward integration of a TLM may depend on the validity of the TLM itself;
- 4) the backward integrations of the TLMs with physical processes; and
- 5) model errors.

How to formulate the ANA in these settings is under investigation.

The current version of ANA requires the operator \mathbf{C} to be invertible. In the case of incomplete observations, the operator \mathbf{C} is not invertible in the classic sense, and one has to use the generalized inverse of \mathbf{C} (Bennett 1992; Caradus 1974; Golub and van Loan 1989; Rao and Mitra 1971). Based on the paper by Zou et al. (1992), incomplete data in a 4D variational data assimilation may affect several issues, namely, conditioning of the Hessian matrix, uniqueness of the solution, convergence of the minimization process, quality of the retrieved fields, and finally the quality of the ensuing forecast. A detailed study addressing these issues is crucial in order to demonstrate that this method has potential for application to real data 4D variational data assimilation. The first remark following Theorem 2 in appendix A is a good starting point. This important issue will be addressed in a separate paper in which real observations will be mimicked in two possible ways: 1) assume data are completely absent for one or more model variables (to represent small-scale data assimilation using radar data) and 2) assume data are available for

all variables but absent in some large areas (to represent large-scale data assimilation using conventional data).

Convergence criteria and line search conditions in the current version of the ANA can be carried out without the information of the gradient using techniques such as the golden section search, etc. (Luenberger 1984). Namely, the next version of the ANA to be tested in the near future will not use the gradient and the adjoint model. In this paper, the ANA is only tested using an adiabatic model and simulated observations. Our next paper will also deal with real data cases using models that include additional physical processes. The backward integration of TLMs with physical processes will also be discussed in a follow-up paper.

Acknowledgments. The first author gratefully acknowledges the insightful comments and useful discussions with Qin Xu at Naval Research Laboratory and thanks Prof. S. Lakshmiarahan and Prof. Alan Shapiro at the University of Oklahoma for their interests and insightful comments. The first author is also thankful to John Lewis at National Severe Storms Laboratory and Andrew Crook and T. Vukićević at NCAR for their interest and encouragement in this project. This research was supported by NSF Grant ATM-912009. Supercomputing was provided by the Pittsburgh Supercomputing Center. Professor Navon was partially supported by NSF Grant ATM-910-2851, managed by Dr. Pamela Stephens. Additional support was provided by the Supercomputer Computations Research Institute, which is partially funded by the Department of Energy.

APPENDIX A

The Theory of the Adjoint Newton Algorithm

For completeness, in this appendix we provide the mathematical foundations of the ANA (Wang et al. 1996). Notations similar to those in the text will be used. For simplicity, the proof is presented under the best, and even unrealistic, scenario by assuming the model is perfect, the TLM can be integrated backward in time, and the operator C is invertible. However, these assumptions may not limit the applications of the ANA as shown in the previous sections and the remarks below. We also assume that $F: R^n \mapsto R^n$ is sufficiently differentiable to carry out our arguments. Discontinuous physical processes are very important. However, we do not deal with them at this time.

THEOREM 1. *Let:*

- (a) Variables \mathbf{X} , \mathbf{X}' , \mathbf{P} , and $\hat{\mathbf{P}} \in R^n$ satisfy Eqs. (2), (4), (6), and (8), respectively.
- (b) The cost function defined by (1) and function $F(\mathbf{X})$ in Eq. (2) be sufficiently differentiable and the operator $\mathbf{C}: R^n \mapsto R^n$ in (1) be invertible.

- (c) There are no observation errors and the cost function defined by (1) has a unique minimum, \mathbf{U}_m , such that $\mathbf{CX}(\mathbf{U}_m)(t) = \mathbf{X}^o(t)$ for any time $t \in [t_0, t_f]$, where $\mathbf{X}(\mathbf{U}_m)(t)$ is the model solution with initial condition \mathbf{U}_m .

Then for a given \mathbf{U} near the global minimum³ \mathbf{U}_m with

$$\mathbf{U}' = \mathbf{U} - \mathbf{U}_m, \quad (\text{A1})$$

and for any time $t \in [t_0, t_f]$, we have

- 1) The forcing terms in the first- and second-order adjoint models satisfy

$$\begin{aligned} \mathbf{C}^T \mathbf{W}[\mathbf{CX}(t) - \mathbf{X}^o(t)] &= [D^2 F(\mathbf{X})\mathbf{X}'(t)]^T \mathbf{P}(t) \\ &+ \mathbf{C}^T \mathbf{W}\mathbf{C}\mathbf{X}'(t) + O(\|\mathbf{U}'\|^2); \end{aligned} \quad (\text{A2})$$

- 2) The first- and second-order adjoint variables satisfy

$$\mathbf{P}(t) = \hat{\mathbf{P}}(t) + O(\|\mathbf{U}'\|^2). \quad (\text{A3})$$

PROOF

1) The nonlinear model and tangent linear model solutions, \mathbf{X} and \mathbf{X}' , obtained with initial conditions \mathbf{U} and \mathbf{U}' , satisfy Eqs. (2) and (4), respectively. Using Taylor expansion (Berger 1977), one has

$$\mathbf{X}(\mathbf{U}) = \mathbf{X}(\mathbf{U}_m) + D\mathbf{X}(\mathbf{U})\mathbf{U}' + O(\|\mathbf{U}'\|^2); \quad (\text{A4})$$

that is,

$$\mathbf{X}(\mathbf{U}) = \mathbf{X}(\mathbf{U}_m) + \mathbf{X}' + O(\|\mathbf{U}'\|^2). \quad (\text{A5})$$

Using Eq. (A5) and $\mathbf{CX}(\mathbf{U}_m)(t) = \mathbf{X}^o(t)$, one obtains

$$\begin{aligned} \mathbf{CX}(U)(t) - \mathbf{X}^o(t) &= \mathbf{C}[\mathbf{X}(U_m)(t) + \mathbf{X}'(t) \\ &+ O(\|\mathbf{U}'\|^2)] - \mathbf{CX}(U_m)(t) \\ &= \mathbf{CX}'(t) + O(\|\mathbf{U}'\|^2). \end{aligned} \quad (\text{A6})$$

Hence $\mathbf{C}^T \mathbf{W}[\mathbf{CX}(t) - \mathbf{X}^o(t)] = O(\|\mathbf{U}'\|)$. From Eq. (6), one obtains (Dieudonne 1960)

$$\mathbf{P}(t) = O(\|\mathbf{U}'\|). \quad (\text{A7})$$

Since F is sufficiently differentiable, $\|D^2 F(\mathbf{X})\|$ is bounded. Hence

$$[D^2 F(\mathbf{X})\mathbf{X}'(t)]^T \mathbf{P}(t) = O(\|\mathbf{U}'\|^2),$$

which together with Eq. (A6) leads to

$$\begin{aligned} \mathbf{C}^T \mathbf{W}[\mathbf{CX}(t) - \mathbf{X}^o(t)] &= \langle \{D^2 F[\mathbf{X}(t)]\mathbf{X}'(t)\}^T \mathbf{P} + \mathbf{C}^T \mathbf{W}\mathbf{C}\mathbf{X}'(t) \rangle \\ &= \mathbf{C}^T \mathbf{W}[\mathbf{CX}'(t) + O(\|\mathbf{U}'\|^2)] - O(\|\mathbf{U}'\|^2) \\ &= \mathbf{C}^T \mathbf{W}\mathbf{C}\mathbf{X}'(t) + O(\|\mathbf{U}'\|^2). \end{aligned} \quad (\text{A8})$$

Equation (A8) yields Eq. (A2).

2) Subtracting Eq. (6) from Eq. (8), one can see that $\hat{\mathbf{P}}(t) - \mathbf{P}(t)$ satisfies

³ See remark 1 after the proof of Theorem 1 for the meaning of "near the global minimum \mathbf{U}_m ."

$$\begin{aligned}
 -\frac{d[\hat{\mathbf{P}}(t) - \mathbf{P}(t)]}{dt} &= [DF(\mathbf{X})]^T[\hat{\mathbf{P}}(t) - \mathbf{P}(t)] \\
 &+ [D^2F(\mathbf{X})\mathbf{X}']^T\mathbf{P}(t) + \mathbf{C}^T\mathbf{W}\mathbf{C}\mathbf{X}' \\
 &- \mathbf{C}^T\mathbf{W}(\mathbf{C}\mathbf{X} - \mathbf{X}^o), \tag{A9}
 \end{aligned}$$

with the final condition

$$\mathbf{P}(t_f) - \mathbf{P}(t_f) = \mathbf{0}. \tag{A10}$$

Equation (A9) is a linear equation with a forcing term satisfying Eq. (A8). Hence Eqs. (A9) and (A10) yield Eq. (A3) (Dieudonne 1960).

REMARKS:

1) Under the assumption of sufficient conditions for the convergence of Newton’s method (Berger 1977; Stoer and Bulirsch 1976), we know that there exists a neighborhood, B , of \mathbf{U}_m such that if $\mathbf{U} \in B$, the Newton’s method converges. “For a given U near the global minimum $\mathbf{U}_m \dots$ ” in Theorem 1 means $\mathbf{U} \in B$.

2) In most cases, Eqs. (A2) and (A3) will be approximately satisfied. In appendix B, we will show that Eqs. (A2) and (A3) can be exactly satisfied using a linear example. For nonlinear model equations, Eqs. (A2) and (A3) are satisfied with indicated accuracy. Interested readers may use

$$\frac{d\mathbf{X}}{dt} = -\mathbf{X}^2, \tag{A11}$$

$$\mathbf{X}(0) = \mathbf{U}, \tag{A12}$$

where time $t \in [0, 1]$ to verify Theorem 1.

COROLLARY 1. *If the same assumptions in Theorem 1 hold, then*

$$\mathbf{X}'(t_f) = \mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t_f) - \mathbf{X}^o(t_f)] + O(\|\mathbf{U}'\|^2), \tag{A13}$$

where \mathbf{C}^{-1} is the inverse of \mathbf{C} .

PROOF. From Theorem 1, we know that Eq. (A2) is true at any time. Hence it is true at the time t_f in particular; that is,

$$\begin{aligned}
 &\mathbf{C}^T\mathbf{W}[\mathbf{C}\mathbf{X}(t_f) - \mathbf{X}^o(t_f)] \\
 &= [D^2F(\mathbf{X})\mathbf{X}'(t_f)]^T\mathbf{P}(t_f) + \mathbf{C}^T\mathbf{W}\mathbf{C}\mathbf{X}'(t_f) \\
 &+ O(\|\mathbf{U}'\|^2). \tag{A14}
 \end{aligned}$$

Noticing that $\mathbf{P}(t_f) = \mathbf{0}$ given by Eq. (7) and \mathbf{C} is invertible, solving Eq. (A14) for $\mathbf{X}'(t_f)$, one obtains Eq. (A13) as a result.

THEOREM 2. *If the same assumptions in Theorem 1 hold and the tangent linear model can be integrated backward in time, then the estimated Newton descent direction is given by $\mathbf{d}_e = -\mathbf{Y}(t_0) = -\mathbf{U}'_e$, where $\mathbf{Y}'(t_0)$ is the solution of the following backward problem:*

$$\frac{d\mathbf{Y}'}{dt} = DF(\mathbf{X})\mathbf{Y}', \tag{A15}$$

$$\mathbf{Y}'(t_f) = \mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t_f) - \mathbf{X}^o(t_f)] \tag{A16}$$

such that

$$\|\mathbf{U}' - \mathbf{U}'_N\| = O(\|\mathbf{U}'\|^2), \tag{A17}$$

$$\|\mathbf{U}'_e - \mathbf{U}'_N\| = O(\|\mathbf{U}'\|^2), \tag{A18}$$

$$\|\mathbf{U}'_e - \mathbf{U}'\| = O(\|\mathbf{U}'\|^2). \tag{A19}$$

PROOF. Since the tangent linear model can be integrated backward in time, Eqs. (A15) and (A16) indicate that

$$\mathbf{Y}' = D\mathbf{X}(\mathbf{U})\mathbf{U}'_e. \tag{A20}$$

According to Eq. (A6) and noticing $\mathbf{X}' = O(\|\mathbf{U}'\|)$, one has

$$\mathbf{Y}'(t_f) = \mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t_f) - \mathbf{X}^o(t_f)] = \mathbf{X}'(t_f) + O(\|\mathbf{U}'\|^2). \tag{A21}$$

Hence

$$\mathbf{Y}'(t) = O(\|\mathbf{U}'_e\|) = O(\|\mathbf{U}'\|). \tag{A22}$$

Since the tangent linear model can be integrated backward in time, the solution of the following equations,

$$\frac{d\mathbf{Z}'}{dt} = DF(\mathbf{X})\mathbf{Z}', \tag{A23}$$

$$\mathbf{Z}'(t_f) = \mathbf{X}'(t_f) = [D\mathbf{X}(\mathbf{U})\mathbf{U}']_{t_f}, \tag{A24}$$

is the same as that of Eqs. (4) and (5). Namely,

$$\mathbf{Z}' = D\mathbf{X}(\mathbf{U})\mathbf{U}' = \mathbf{X}'. \tag{A25}$$

Equation (A6) at t_f yields

$$\mathbf{C}^{-1}[\mathbf{C}\mathbf{X}(t_f) - \mathbf{X}^o(t_f)] - \mathbf{X}'(t_f) = O(\|\mathbf{U}'\|^2), \tag{A26}$$

which is a relationship between final conditions (A16) and (A24). Hence,

$$\begin{aligned}
 &\mathbf{Y}'(t) - \mathbf{X}'(t) \\
 &= \mathbf{Y}'(t) - \mathbf{Z}'(t) = D\mathbf{X}(\mathbf{U})\mathbf{U}'_e - D\mathbf{X}(\mathbf{U})\mathbf{U}' \\
 &= O(\|\mathbf{U}'\|^2). \tag{A27}
 \end{aligned}$$

At the minimum, the gradient of the cost function with respect to control variables is zero; that is,

$$\begin{aligned}
 0 &= DJ(\mathbf{U}_m) = DJ(\mathbf{U} - \mathbf{U}') \\
 &= DJ(\mathbf{U}) - D^2J(\mathbf{U})\mathbf{U}' + O(\|\mathbf{U}'\|^2) \tag{A28}
 \end{aligned}$$

and

$$D^2J(\mathbf{U})\mathbf{U}' = DJ(\mathbf{U}) + O(\|\mathbf{U}'\|^2). \tag{A29}$$

Subtracting Eq. (14) from Eq. (A29), one obtains

$$D^2J(\mathbf{U})(\mathbf{U}' - \mathbf{U}'_N) = O(\|\mathbf{U}'\|^2). \tag{A30}$$

Since by assumption $D^2J(\mathbf{U})$ is positive definite, Eq. (A30) yields Eq. (A17). Now consider the forcing terms of Eq. (A9) with \mathbf{X}' replaced by \mathbf{Y}' and apply Eqs. (A22), (A7), (A27), and (A6); one obtains

$$\begin{aligned} & \left\| \{ [D^2F(\mathbf{X})\mathbf{Y}']^T \mathbf{P}(t) + \mathbf{C}^T \mathbf{W} \mathbf{C} \mathbf{Y}' \} - \mathbf{C}^T \mathbf{W} (\mathbf{C} \mathbf{X} - \mathbf{X}^o) \right\| \\ &= \left\| [D^2F(\mathbf{X})\mathbf{Y}']^T \mathbf{P}(t) + \mathbf{C}^T \mathbf{W} \mathbf{C} (\mathbf{Y}' - \mathbf{X}') + \mathbf{C}^T \mathbf{W} (\mathbf{C} \mathbf{X}' - (\mathbf{C} \mathbf{X} - \mathbf{X}^o)) \right\| \\ &= O(\|\mathbf{U}\|^2) + O(\|\mathbf{U}'\|^2) + O(\|\mathbf{U}'\|^2) = O(\|\mathbf{U}'\|^2). \end{aligned} \tag{A31}$$

That is

$$\hat{\mathbf{P}}(t) - \mathbf{P}(t) = D^2J(\mathbf{U})\mathbf{U}'_e - DJ(\mathbf{U}) = O(\|\mathbf{U}'\|^2). \tag{A32}$$

Using Newton's equation Eq. (14), Eq. (A32) becomes

$$D^2J(\mathbf{U})(\mathbf{U}'_e - \mathbf{U}'_N) = O(\|\mathbf{U}'\|^2), \tag{A33}$$

which yields Eq. (A18). From Eqs. (A17) and (A18), one obtains

$$\|\mathbf{U}'_e - \mathbf{U}'\| \leq \|\mathbf{U}'_e - \mathbf{U}'_N\| + \|\mathbf{U}'_N - \mathbf{U}'\| + O(\|\mathbf{U}'\|^2), \tag{A34}$$

which yields Eq. (A19). If \mathbf{U}' is sufficiently small, \mathbf{U}'_e is a descent direction since

$$\begin{aligned} \mathbf{U}'^T \mathbf{U}'_e &= \mathbf{U}'^T (\mathbf{U}'_e - \mathbf{U}' + \mathbf{U}') \\ &\geq \|\mathbf{U}'\|^2 - \|\mathbf{U}'\| \|\mathbf{U}'_e - \mathbf{U}'\| \\ &= \|\mathbf{U}'\| (\|\mathbf{U}'\| - \|\mathbf{U}'_e - \mathbf{U}'\|) \\ &= \|\mathbf{U}'\| (\|\mathbf{U}'\| - O(\|\mathbf{U}'\|^2)) > 0. \end{aligned} \tag{A35}$$

REMARKS:

1) For simplicity, the operator \mathbf{C} is assumed to be invertible. However, \mathbf{C} is not necessarily invertible except in case it is applied to certain discrete systems. In this case, one may reformulate the Newton's equation

$$D^2J(\mathbf{U})\mathbf{U}'_N = DJ(\mathbf{U}) \tag{A36}$$

as a minimization problem of the following form

$$\min_{\mathbf{U}'} \frac{1}{2} \|D^2J(\mathbf{U})\mathbf{U}' - DJ(\mathbf{U})\|^2, \tag{A37}$$

where \mathbf{U}' is the initial condition of tangent linear model given by Eqs. (4) and (5). Minimization problem ((A37) is quadratic and is similar to problem (1). All techniques such as penalty, model error techniques, etc., suitable for problem (1) are also suitable for the quadratic problem (A37). In this case, the tangent linear model is integrated forward. This is related to generalized inverse problem (Bennett 1992; Caradus 1974; Golub and van Loan 1989; Rao and Mitra 1971) and is successfully applied to the ARPS. Our next paper will deal with this general case.

2) For completeness, the algorithm form of the ANA is listed here:

(a) Choose an initial guess \mathbf{U}_0 and set $k = 0$.

(b) Calculate the gradient of the cost function with respect to initial conditions

$$\mathbf{g}_k = \mathbf{g}(\mathbf{U}_k) = DJ(\mathbf{U}_k), \tag{A38}$$

by integrating the first-order adjoint model equations.

(c) Obtain a Newton line search direction,

$$\mathbf{d}_k = -\mathbf{Y}'(t_0), \tag{A39}$$

by integrating the backward tangent linear model given by Eqs. (16) and (17).

(d) Set

$$\mathbf{U}_{k+1} = \mathbf{U}_k + \alpha_k \mathbf{d}_k, \tag{A40}$$

where α_k is the step size obtained by conducting a line search,

$$J(\mathbf{U}_k + \alpha_k \mathbf{d}_k) = \min_{\alpha} J(\mathbf{U}_k + \alpha \mathbf{d}_k), \tag{A41}$$

using Davidon's cubic interpolation method for the line search of the step size and that satisfies the following Wolfe conditions (Liu and Nocedal 1989):

$$J(\mathbf{U}_k + \alpha_k \mathbf{d}_k) \leq J(\mathbf{U}_k) + \beta' \alpha_k \mathbf{g}_k^T \mathbf{d}_k \tag{A42}$$

and

$$\frac{\mathbf{g}_k(\mathbf{U}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k}{\mathbf{g}_k^T \mathbf{d}_k} \leq \beta, \tag{A43}$$

where $\beta' = 0.0001$, $\beta = 0.9$.

(e) Check the convergence condition. Given a tolerance criterion ϵ : $10^{-14} \leq \epsilon \leq 10^{-1}$, and if

$$\|\mathbf{g}_{k+1}\| \leq \epsilon \|\mathbf{g}_0\|, \tag{A44}$$

stop. Otherwise, set $k = k + 1$ and go to step (b).

3) The adjoint Newton method requires the user to provide a subroutine to calculate the cost function and a subroutine to calculate the Newton direction by integrating the backward tangent linear model. In the current version of the ANA, the gradient is only used to check the convergence criterion (A44) and line search conditions (A42) and (A43). This criterion and conditions could be carried out without the gradient information (Luenberger 1984). In addition to the above calculations, each iteration of the adjoint Newton iteration requires two stages: setup (performed once per iteration) and line search iteration. The numerical costs of the two stages are approximately $2n$ flops and $4n$ flops (Nash and Nocedal 1989), respectively, where "flops" denotes

additions, subtractions, multiplications, or divisions. The ANA is a large-scale unconstrained minimization Newton algorithm that does not have the limitation of the storage problem since it does not require the calculation of either the Hessian or its inverse. However, it may be sensitive to the initial guess and may only be applied to the problems where their corresponding backward problems could be solved with “reasonable accuracy.”

APPENDIX B

A Simple Example

We first show that Eqs. (A2) and (A3) are exactly satisfied by using a very simple linear example and then we illustrate the ANA by using the same example.

Let us consider the following one-dimensional model equation:

$$\frac{dX}{dt} = -X, \tag{B1}$$

$$X(0) = U, \tag{B2}$$

where time $t \in [0, 1]$. The solution of the model is

$$X = Ue^{-t}. \tag{B3}$$

Suppose that the simulated observations are model generated with the initial condition

$$X(0) = 1, \tag{B4}$$

then from Eq. (B3) the simulated observation may be written as

$$X^o = e^{-t}. \tag{B5}$$

Let us define the cost function as

$$J(U) = \frac{1}{2} \int_0^1 (X - X^o)^2 dt, \tag{B6}$$

then

$$J(U) = \frac{(U - 1)^2}{4}(1 - e^{-2}). \tag{B7}$$

The first-order adjoint model of Eqs. (B1) and (B2) may be written as

$$-\frac{dP}{dt} = -P + (X - X^o); \tag{B8}$$

that is,

$$-\frac{dP}{dt} = -P + (U - 1)e^{-t}, \tag{B9}$$

$$P(1) = 0, \tag{B10}$$

where P is the adjoint variable. The gradient of the cost function with respect to the initial conditions is given by

$$DJ(U) = P(0). \tag{B11}$$

Equation (B9) has an analytic solution of the following form:

$$\begin{aligned} P(t) &= \bar{P}e^t - e^t \int_0^t e^{-\tau}(U - 1)e^{-\tau} d\tau \\ &= \bar{P}e^t + \frac{(U - 1)}{2}(e^{-t} - e^t), \end{aligned} \tag{B12}$$

where \bar{P} is a constant to be determined by the final condition (B10). Therefore $\bar{P} = [(U - 1)/2](1 - e^{-2})$, and Eq. (B12) yields

$$P(t) = \frac{(U - 1)}{2}(1 - e^{-2})e^t + \frac{(U - 1)}{2}(e^{-t} - e^t) \tag{B13}$$

and

$$P(0) = \frac{(U - 1)}{2}(1 - e^{-2}). \tag{B14}$$

Equation (B14) is exactly the gradient of the cost function [Eq. (B7)] with respect to the initial condition U .

Let us now consider a perturbation, U' , on the initial condition U for X . Then the tangent linear and second-order adjoint models are

$$\frac{dX'}{dt} = -X', \tag{B15}$$

$$X'(0) = U', \tag{B16}$$

and

$$-\frac{d\hat{P}}{dt} = -\hat{P} + X', \tag{B17}$$

$$\hat{P}(1) = 0, \tag{B18}$$

respectively. They have the following exact solutions

$$X' = U'e^{-t} \tag{B19}$$

and

$$\hat{P}(t) = \frac{U'}{2}(1 - e^{-2})e^t + \frac{U'}{2}(e^{-t} - e^t), \tag{B20}$$

respectively. Here, $\hat{P}(0) = U'(1 - e^{-2})/2$ is exactly Hessian vector product $D^2J(U)U'$.

It can be shown that

$$\hat{P}(0) = U' \frac{(1 - e^{-2})}{2} = \frac{(U - 1)}{2}(1 - e^{-2}) = P(0)$$

if $U' = U - 1$; that is, Eq. (A3) can be exactly satisfied in this linear case. In nonlinear cases, Eq. (A3) is only an approximation. When $U' = U - 1$, the forcing terms of Eqs. (B9) and (B17) are exactly equal, $(U - 1)e^{-t} = U'e^{-t}$, which indicates that Eq. (A2) is exactly satisfied for this problem.

The variational data assimilation aims to find the best initial condition, $X(0) = 1$, which minimizes the cost function given by (B7). It should be realized that since

the model itself is linear, its tangent linear model is identical to itself. Since we can solve it exactly, an exact Newton descent direction could be obtained and the minimum point is found in one step.

- 1) Assume an arbitrary initial guess U_0 .
- 2) Solving the backward problem

$$\frac{dY'}{dt} = -Y', \tag{B21}$$

$$Y'(1) = X(t_f) - X^o(t_f) = (U_0 - 1)e^{-1}, \tag{B22}$$

one obtains

$$d_0 = -Y'(0) = -(U_0 - 1). \tag{B23}$$

- 3) Let

$$U_1 = U_0 + \alpha_0 d_0, \tag{B24}$$

where α_0 is the step size obtained by the following line search:

$$J(U_0 + \alpha_0 d_0) = \min_{\alpha} \left\{ \frac{[U_0 - \alpha(U_0 - 1) - 1]^2}{2} (1 - e^{-2}) \right\}, \tag{B25}$$

which has an exact solution $\alpha_0 = 1$.

- 4) Using the newly found α_0 , we can see that

$$U_1 = U_0 + \alpha_0 d_0 = U_0 - (U_0 - 1) = 1, \tag{B26}$$

which is the best initial condition that minimizes the cost function given by Eq. (B7), that is,

$$J(U_1) = 0. \tag{B27}$$

So the adjoint Newton algorithm algorithm has found the minimum in a single step.

APPENDIX C

Convergence Analysis

For a quick assessment of the rate of convergence, a table can be constructed of

$$\xi_k = J_{k-1} - J_k$$

for the last few values of k (Gill and Murray 1972). In our experiment, J_k denotes the value of the cost function at the k th iteration scaled by J_0 . Superlinear convergence would be indicated if $\xi_{k+1} \approx r \xi_k$, where $r > 1$. Fast linear convergence would be indicated if $\xi_{k+1} \approx \xi_k/M$, where $M > 2$.

In forming the sequence $\{\xi_k\}$, the user needs to be aware that eventually all algorithms either fail to make further progress or display slow linear convergence near the limiting accuracy of the solution. What may occur with a superlinearly convergent algorithm is that the sequence $\{\xi_k\}$ will demonstrate superlinear convergence

TABLE C1. Convergence analysis information for adjoint Newton algorithm in experiment 1.

k	J_k	ξ_k	$\xi_k/\xi_{k+1} = M$
1	4.847×10^{-2}	0.952	22.50
2	6.183×10^{-3}	4.229×10^{-2}	8.438
3	1.170×10^{-3}	5.012×10^{-3}	6.176
4	3.590×10^{-4}	8.116×10^{-4}	3.847
5	1.165×10^{-4}	2.424×10^{-4}	3.781
14	1.060×10^{-7}	8.804×10^{-8}	1.809
15	5.734×10^{-8}	4.867×10^{-8}	2.069
16	3.382×10^{-8}	2.352×10^{-8}	1.585
17	1.897×10^{-8}	1.484×10^{-8}	2.001
18	1.155×10^{-8}	7.417×10^{-9}	1.460

for a few iterations only and then lapse into slow linear convergence when limiting accuracy has been achieved. Therefore, it is important, especially if a failure has been indicated, to examine the sequence $\{\xi_k\}$ at iterations that sufficiently precede the final stage (Gill and Murray 1972).

After 18 iterations, the cost function decreased eight orders of magnitude when the ANA is used. We construct Table C1 for iterations 1–5 and 14–18. This table indicates that the ANA has a fast linear convergence rate for the first five iterations with average $M = 8.77$, while its convergence rate slows down for the last five iterations with average $M = 1.79$. Similar studies are carried out for the LBFGS method. It is found out that LBFGS has a linear convergence rate with average $M = 2.38$ for the first 5 iterations and $M = 0.63$ for iterations 74–78.

As pointed out in the third remark in section 2b, ANA has a theoretical quadratic convergence rate. It is a Newton method. However, its numerical convergence rate in this example is only fast linear. This is due to the fact that once modifications are introduced in the backward integration of the tangent linear model, its solution at the initial time will not be as good as the Newton descent direction. Hence the corresponding ANA will not display a quadratic convergence rate.

REFERENCES

Bennett, A., 1992: *Inverse Problem in Physical Oceanography*. Cambridge University Press, 346 pp.
 Berger, M. S., 1977: *Nonlinearity and Functional Analysis*. Academic Press, 417 pp.
 Caradus, S. R., 1974: Operator theory of the pseudo-inverse. *Queen's Papers in Pure and Applied Mathematics* 38, 67 pp. [Available from Queen's University, Kingston, ON K7L-5C4, Canada.]
 Davidon, W. C., 1959: Variable metric method for minimization. A. E. C. Research and Development Rep. ANL-5990. [Available from Argonne National Laboratory, Argonne, IL 60439.]
 Dieudonne, J., 1960: *Foundations of Modern Analysis*. Academic Press, 361 pp.
 Droegemeier, K., and Coauthors, 1995: Weather prediction: A scalable storm-scale model. *High Performance Computing*, G. W. Sabot, Ed., Addison Wesley, 45–92.
 Errico, R. M., T. Vukićević, and K. Raeder, 1993: Examination of the accuracy of a tangent linear model. *Tellus*, **45A**, 462–477.

- Gill, P. E., and W. Murray, 1972: Quasi-Newton methods for unconstrained optimization. *J. Inst. Math. Its Appl.*, **9**, 91–108.
- , —, and M. H. Wright, 1981: *Practical Optimization*. Academic Press, 401 pp.
- Golub, G. H., and C. F. van Loan, 1989: *Matrix Computations*. 2d ed. The Johns Hopkins University Press, 642 pp.
- Klemp, J. B., and R. B. Wilhelmson, 1978: The simulation of three-dimensional convective storm dynamics. *J. Atmos. Sci.*, **35**, 1070–1096.
- Lacarra, J. F., and O. Talagrand, 1988: Short-range evolution of small perturbations in a barotropic model. *Tellus*, **40A**, 81–95.
- Le Dimet, F. X., and O. Talagrand, 1986: Variational algorithms for analysis and assimilation of meteorological observations: Theoretical aspects. *Tellus*, **38A**, 97–110.
- Liu, D. C., and J. Nocedal, 1989: On the limited memory BFGS method for large scale minimization. *Math. Prog.*, **45**, 503–528.
- Luenberger, D. G., 1984: *Linear and Nonlinear Programming*. 2d ed. Addison-Wesley, 491 pp.
- Nash, S. G., 1984a: Newton-type minimization via the Lanczos method. *SIAM J. Numer. Anal.*, **21** (4), 770–788.
- , 1984b: Truncated-Newton methods for large-scale function minimization. *Applications of Nonlinear Programming to Optimization and Control*, Pergamon Press, H. E. Rauch, Ed., 91–100.
- , 1985: Preconditioning of truncated-Newton methods. *SIAM J. Sci. Stat. Comput.*, **6** (3), 599–616.
- , and J. Nocedal, 1989: A numerical study of the limited memory BFGS method and the truncated-Newton method for large scale optimization. Tech. Rep. NAM, 02, Dept. of Electrical Engineering and Computer Science, Northwestern University, 19 pp. [Available from Dept. of Electrical Engineering and Computer Science, Northwestern University, Evanston, IL 60208.]
- , and A. Sofer, 1989: Block truncated-Newton methods for parallel optimization. *Math. Prog.*, **45**, 529–546.
- Navon, I. M., and R. de Villiers, 1983: Combined penalty multiplier optimization methods to enforce integral invariants conservation. *Mon. Wea. Rev.*, **111**, 1228–1243.
- , and D. M. Legler, 1987: Conjugate gradient methods for large scale minimization in meteorology. *Mon. Wea. Rev.*, **115**, 1479–1502.
- , X. L. Zou, J. Derber, and J. Sela, 1992a: Variational data assimilation with an adiabatic version of the NMC Spectral Model. *Mon. Wea. Rev.*, **120**, 1433–1446.
- , —, M. Berger, P. K. H. Phua, T. Schlick, and F. X. Le Dimet, 1992b: Numerical experience with limited memory quasi-Newton and truncated Newton methods. *Optimization Techniques and Applications*, K. H. Phua et al., Eds., Vol. 1, World Scientific Publishing, 33–48.
- , —, —, —, —, and —, 1992c: Testing for reliability and robustness of optimization codes for large scale optimization problems. *Optimization Techniques and Applications*, K. H. Phua et al., Eds., Vol. 1, World Scientific Publishing, 445–480.
- Nocedal, J., 1980: Updating quasi-Newton matrices with limited storage. *Math. Comput.*, **35**, 773–782.
- O’Leary, D. P., 1983: A discrete Newton algorithm for minimizing a function of many variables. *Math. Prog.*, **23**, 20–23.
- Pu, Z., E. Kalnay, J. Derber, and J. Sela, 1996a: Using past forecast error to improve the future forecast skill. Preprints, *11th Conf. on Numerical Weather Prediction*, Norfolk, VA, Amer. Meteor. Soc., 142–143.
- , —, and —, 1996b: Use of tangent linear model in the study of the sensitivity of forecast error to initial conditions. WMO Weather Prediction Research Programs. CAS/JSC Working Group on Numerical Experimentation, Research Activities in Atmospheric Oceanic Modeling, Rep. 23, WMO/ID-No. 734, 626–627. [Available from WMO, CP No. 2300, CH-1211, Geneva 2, Switzerland.]
- , —, and J. Sela, 1997: Sensitivity of forecast errors to initial conditions with a quasi-inverse linear method. *Mon. Wea. Rev.*, **125**, 2479–2503.
- Rao, C. R., and S. K. Mitra, 1971: *Generalized Inverse of Matrices and Its Applications to Statistics*. John Wiley and Sons, 240 pp.
- Schlick, T., and A. Fogelson, 1992a: TNPACK—A truncated Newton minimization package for large-scale problems: I. Algorithm and usage. *ACMTOMS*, **18** (1), 46–70.
- , and —, 1992b: TNPACK—A truncated Newton minimization package for large-scale problems: II. Implementation examples. *ACMTOMS*, **18** (1), 71–111.
- Shanno, D. F., and K. H. Phua, 1980: Remark on algorithm 500—A variable method sub-routine for unconstrained nonlinear minimization. *ACMTOMS*, **6**, 618–622.
- Smolarkiewicz, P. K., 1983: A simple positive definite advection scheme with small implicit diffusion. *Mon. Wea. Rev.*, **111**, 479–486.
- Stoer, J., and R. Bulirsch, 1976: *Introduction to Numerical Analysis*. 2d ed. Springer-Verlag, 659 pp.
- Sun, J., 1992: Convective scale 4-D data assimilation using simulated single Doppler radar observations. Ph.D. thesis, University of Oklahoma, Norman, OK, 172 pp. [Available from 100 East Boyd Street, EC1310, School of Meteorology, University of Oklahoma, Norman, OK 73019.]
- Thacker, W. C., 1989: The role of Hessian matrix in fitting models to measurements. *J. Geophys. Res.*, **94**, 6177–6196.
- Thépaut, J.-N., D. Vasiljevic, and P. Courtier, 1993: Variational assimilation of conventional meteorological observations with a multilevel primitive-equation model. *Quart. J. Roy. Meteor. Soc.*, **119**, 153–186.
- Vukićević, T., 1991: Nonlinear and linear evolution of initial forecast error. *Mon. Wea. Rev.*, **119**, 1602–1611.
- Wang, Z., 1993: Variational data assimilation with 2-D shallow water equations and 3-D FSU global spectral models. Tech. Rep. FSU-SCRI-93T-149, The Florida State University, Tallahassee, FL, 235 pp. [Available from SCRI, FSU, Tallahassee, FL 32306-4052.]
- , I. M. Navon, F. X. Le Dimet, and X. Zou, 1992: The second order adjoint analysis: Theory and application. *Meteor. Atmos. Phys.*, **50**, 3–20.
- , —, X. Zou, and F. X. Le Dimet, 1995a: The adjoint truncated Newton algorithm for large-scale unconstrained optimization. *Comput. Optimization Appl.*, **4** (3), 241–262.
- , K. Droegemeier, M. Xue, and S. Park, 1995b: The sensitivity of a 3-D compressible storm-scale model to input parameters. Preprints, *Int. Symp. on Assimilation of Observations in Meteorology and Oceanography*, Tokyo, Japan, Japan Meteorological Agency, 437–443.
- , —, and L. White, 1997: The adjoint Newton algorithm for large-scale unconstrained optimization in meteorology applications. *Comput. Optimization Appl.*, in press.
- Weisman, M. L., and J. B. Klemp, 1982: The dependence of numerically simulated convective storms on vertical wind shear and buoyancy. *Mon. Wea. Rev.*, **110**, 504–520.
- Xue, M., K. Droegemeier, V. Vong, A. Shapiro, and K. Brewster, 1995: ARPS Version 4.0 User’s Guide. Center for the Analysis and Prediction of Storms, University of Oklahoma, 380 pp. [Available from Center for the Analysis and Prediction of Storms, University of Oklahoma, Norman, OK 73019.]
- Yang, W., I. M. Navon, and P. Courtier, 1996: A new Hessian preconditioning method applied to variational data assimilation experiments using NASA general circulation models. *Mon. Wea. Rev.*, **124**, 1000–1017.
- Zou, X., I. M. Navon, and F. X. LeDimet, 1992: Incomplete observations and control of gravity waves in variational data assimilation. *Tellus*, **44A** (4), 273–296.
- Zupanski, D., and F. Mesinger, 1995: Four-dimensional variational data assimilation of precipitation data. *Mon. Wea. Rev.*, **123**, 1112–1127.
- Zupanski, M., 1993a: A preconditioning algorithm for large scale minimization problems. *Tellus*, **45A**, 578–592.
- , 1993b: Regional four-dimensional variational data assimilation in a quasi-operational forecasting environment. *Mon. Wea. Rev.*, **121**, 2396–2408.
- , 1996: A preconditioning algorithm suitable for a general four-dimensional data assimilation. Preprints, *11th Conf. on Numerical Weather Prediction*, Norfolk, VA, Amer. Meteor. Soc., 241–242.