

Development of the Upgraded Tangent Linear and Adjoint of the Weather Research and Forecasting (WRF) Model

XIN ZHANG AND XIANG-YU HUANG

National Center for Atmospheric Research, Boulder, Colorado*

NING PAN

Fujian Meteorological Bureau, Fuzhou, Fujian, China

(Manuscript received 20 September 2012, in final form 16 January 2013)

ABSTRACT

The authors propose a new technique for parallelizations of tangent linear and adjoint codes, which were applied in the redevelopment for the Weather Research and Forecasting (WRF) model with its Advanced Research WRF dynamic core using the automatic differentiation engine. The tangent linear and adjoint codes of the WRF model (WRFPLUS) now have the following improvements: A complete check interface ensures that developers write accurate tangent linear and adjoint codes with ease and efficiency. A new technique based on the nature of duality that existed among message passing interface communication routines was adopted to parallelize the WRFPLUS model. The registry in the WRF model was extended to automatically generate the tangent linear and adjoint codes of the required communication operations. This approach dramatically speeds up the software development cycle of the parallel tangent linear and adjoint codes and leads to improved parallel efficiency. Module interfaces were constructed for coupling tangent linear and adjoint codes of the WRF model with applications such as four-dimensional variational data assimilation, forecast sensitivity to observation, and others.

1. Introduction

During the past two decades, the use of the adjoint technique in meteorology and oceanography rapidly increased. The adjoint model is a powerful tool in many applications, such as data assimilation, parameter estimation, and sensitivity analysis (Errico and Vukicevic 1992; Errico 1997; Rabier et al. 1996; Langland et al. 1999; Li et al. 1999; Xiao et al. 2002, 2008; Kleist and Morgan 2005a,b).

The Weather Research and Forecasting (WRF) modeling system (Skamarock et al. 2008) is a multiagency effort to provide a next-generation mesoscale forecast model and data assimilation system that will advance

both the understanding and the prediction of mesoscale weather as well as accelerate the transfer of research advances into operations. The WRF model was designed to be an efficient massively parallel computing code to take advantage of advanced high-performance computing systems. The code, which can be configured for both research and operations, also offers numerous physics options. The WRF model is maintained and supported as a community model to facilitate wide use internationally, for research, operations, and teaching. The model is suitable for a broad span of applications across scales ranging from large-eddy to global simulations. Such applications include real-time numerical weather prediction, data assimilation development and studies, parameterized-physics research, regional climate simulations, air quality modeling, atmosphere–ocean coupling, and idealized simulations. As of this writing, the WRF model is in operational and research use around the world, and the number of registered WRF users exceeds 20 000.

The first version of the adiabatic WRF tangent linear model (TLM) and adjoint model (ADM) system

*The National Center for Atmospheric Research is sponsored by the National Science Foundation.

Corresponding author address: Dr. Xin Zhang, MMM, NCAR, P.O. Box 3000, Boulder, CO 80307.
E-mail: xinzhang@ucar.edu

(WAMS) was developed by the National Center for Atmospheric Research (NCAR) around 2007 (Xiao et al. 2008). The WAMS is used in the adjoint sensitivity analysis (Xiao et al. 2008) and the four-dimensional variational data assimilation (4DVAR) (Huang et al. 2009). It took us more than one year to complete the parallelization of the WAMS by following the traditional incremental and iterative parallelization procedures, which include algorithm analysis, identifying the time-consuming hot spot, data dependency analysis, locating the halo exchange, and manually writing the communication routines for data exchanging. In the past few years, because of the tremendous time required for hand-coding parallelization, the WAMS has failed to repeat the rapid development pattern of the WRF model and the WRF data assimilation system (WRFDA) (Barker et al. 2012). The growing gap between the WAMS and the WRF–WRFDA makes the WAMS inconvenient to use with other systems. One example was that the WRF 4DVAR system could not read the WRF boundary conditions of later versions because the WAMS was left behind and it does not include the upgrades of the boundary condition data structure in the WRF. Furthermore, because the WAMS uses the disk input/output (I/O) for storing–reading basic-state trajectory and exchanging data among the WRF full-physics forward model (FWM), TLM, and ADM, the parallel efficiency is unsatisfactory on modern high-performance computers with distributed memory parallelization, especially for 4DVAR applications.

Encouraged by the rapid developments of 4DVAR, cloud analysis, forecast sensitivity to observations, and chemistry data assimilation, we redeveloped the WRF TLM and ADM (called WRFPLUS) based on the latest repository WRF. Compared with the WAMS developed by Xiao et al. (2008) and Huang et al. (2009), the new system is an all-in-one system that includes the FWM, TLM, and ADM; this system also includes the tangent linear check and adjoint check procedures. A set of module interfaces was developed for easily coupling the WRFPLUS model with other systems such as data assimilation and adjoint sensitivity applications. A new approach was applied to develop the parallel code that dramatically reduces the software development cycle of the parallel TLM and ADM, and the derived parallel TLM and ADM have better parallel efficiency compared to the FWM.

The purpose of this paper is to describe the technical aspects of the newly developed WRFPLUS model. A brief introduction of the development of the WRFPLUS model is presented in section 2, followed by a demonstration of the linearity and adjoint tests in section 3. A detailed description of the parallelization strategy

for tangent linear and adjoint codes, with the demonstrated parallel performance are in section 4. Section 5 introduces the module interfaces constructed in the WRFPLUS model for coupling purposes, such as in WRF 4DVAR. Concluding remarks appear in section 6.

2. Description of the WRF tangent linear and adjoint models

After the release of the WRF model version 3.2, we started to use TAPENADE (Hascoët and Pascual 2004) to redevelop the TLM and the ADM of the Advanced Research WRF (ARW-WRF) core based on the latest WRF model. The development of the WRFPLUS model follows the same three phases proposed by Xiao et al. (2008). First, numerical experiments were conducted to make sure that the adiabatic version of the WRF model with simplified physics parameterization routines can produce the major features that the full-physics model does. Second, the TLM and its ADM were generated by TAPENADE and modified manually whenever necessary. Third, the TLM and the ADM were verified for correctness.

TAPENADE is a source-to-source automatic differentiation (AD) tool for programs written in FORTRAN 77–95; that is, TAPENADE generates a TLM or an ADM from the source code of a given model. Like other AD tools, TAPENADE struggles with such complicated codes as WRF's third-order Runge–Kutta large time steps and small acoustic time steps (Xiao et al. 2008). Checking and improving TAPENADE-generated code can require manual intervention (i) to represent first-order physical effects on the model evolution; (ii) to minimize code length, so that developing its adjoint is simple; and (iii) to allow the code to run quickly in the iterations, which is required to lower costs. We also adopted three simplified physics packages that have maximum impact on a forecast compared to a no-physics model: surface friction, cumulus parameterization, and large-scale condensation. These were developed by J. Dudhia (2012, personal communication) and are similar to the physics packages in the WAMS (Xiao et al. 2008).

3. Linearity test and adjoint test

Testing TLM consistency with FWM as well as ADM consistency with TLM before either is used in any real application is essential (Vukicevic 1991; Errico and Vukicevic 1992; Gilmour et al. 2001). We developed the tangent linear and adjoint check procedures following Navon et al. (1992).

Let $f(\mathbf{x})$, $g_{-}f(\mathbf{x}, g_{-}\mathbf{x})$, and $a_{-}f(\mathbf{x}, a_{-}\mathbf{x})$ denote an FWM, a TLM, and an ADM, respectively, where \mathbf{x} , $g_{-}\mathbf{x}$, and $a_{-}\mathbf{x}$ are the column vectors of model-state variables, perturbations of state variables, and the adjoint of state variables, respectively. The correctness of the TLM can be tested as

$$\Phi(\lambda) = \frac{\|f(\mathbf{x} + \lambda g_{-}\mathbf{x}) - f(\mathbf{x})\|}{\|g_{-}f(\mathbf{x}, \lambda g_{-}\mathbf{x})\|}, \quad \lim_{\lambda \rightarrow 0} \Phi(\lambda) = 1, \quad (1)$$

where $\|\cdot\|$ denotes the norm of the vector. The adjoint relation is tested by

$$\langle g_{-}f(\mathbf{x}, g_{-}\mathbf{x}), g_{-}f(\mathbf{x}, g_{-}\mathbf{x}) \rangle = \langle a_{-}f[\mathbf{x}, g_{-}f(\mathbf{x}, g_{-}\mathbf{x})], g_{-}\mathbf{x} \rangle. \quad (2)$$

If the tangent linear and adjoint codes are correct, then the above-mentioned two relations should hold up to the machine accuracy. Because different model variables have different magnitudes, we also designed the capability to perform checks on individual variables separately. We performed the tangent linear and adjoint checks with the test case being integrated up to 24 h. We sequentially reduced the initial perturbations by a factor (λ) of 10 and repeatedly calculated $\Phi(\lambda)$ in Eq. (1) on NCAR's IBM machine with 64-bit precision. Table 1 shows the value of $\Phi(\lambda)$ from the tangent linear forecasts and the differences between two nonlinear forecasts over the whole domain. This indicated that the tangent linear forecast approximates the difference between two nonlinear forecasts as the initial perturbations decrease and approach zero.

In the adjoint relationship, the left-hand side (lhs) involves only the tangent linear code, while the right-hand side (rhs) involves the adjoint code. If lhs and rhs have the same value with machine accuracy, then the adjoint code is correct compared with the tangent linear code. Using the same test case with 24-h integration, lhs and rhs for the test case are $0.141\,827\,207\,298\,78 \times 10^{14}$ and $0.141\,827\,207\,298\,83 \times 10^{14}$, respectively. This indicates that the adjoint code is correct.

4. Parallelization of the WRFPLUS model

TAPENADE has few problems in generating tangent linear and adjoint codes from the sequential forward codes. However, it cannot derive such codes of parallel communication routines inside a parallel forward model. In most of the atmospheric and oceanic models, the communication routines to parallelize finite difference algorithms are linear operators; hence, as matrices the

TABLE 1. Ratio of norms between the tangent linear forecasts and the differences of the two nonlinear model forecasts at 24 h. Norm is defined as the summation of the squares of all variables (perturbations of tangent linear model and difference of two nonlinear models) over the whole domain at 24 h. Here, λ is the perturbation scaling factors of the initial perturbation.

$\lambda = 0.1000 \times$	Ratio
10	$0.101\,142\,811\,984\,81 \times 10^{-1}$
10^{-1}	$0.100\,085\,452\,404\,48 \times 10^{-1}$
10^{-2}	$0.100\,008\,324\,840\,54 \times 10^{-1}$
10^{-3}	$0.100\,000\,958\,062\,29 \times 10^{-1}$
10^{-4}	$0.100\,000\,075\,039\,57 \times 10^{-1}$
10^{-5}	$0.100\,000\,017\,434\,69 \times 10^{-1}$
10^{-6}	$0.100\,000\,003\,442\,15 \times 10^{-1}$
10^{-7}	$0.999\,999\,985\,519\,13 \times 10^{-1}$
10^{-8}	$0.100\,000\,014\,534\,68 \times 10^{-1}$
10^{-9}	$0.100\,000\,073\,020\,81 \times 10^{-1}$
10^{-10}	$0.100\,007\,756\,313\,70 \times 10^{-1}$

adjoint is simply the transpose, which is the dual operator, and the tangent linear is the original linear operator acting on the perturbations (see Cheng 2006; Utke et al. 2009). With the proper coding structure and available parallel communication routines in forward codes, it is straightforward to write communication routines for tangent linear models. We need only use the same parallel communication templates in forward codes and add the corresponding perturbation variables. In the adjoint code, the data flow of the original program needs to be reversed and any communication needs to be reversed as well. Because of the duality between MPI_SEND and MPI_RECV calls, in transforming FWM to ADM, we replace MPI_SEND calls with MPI_RECV, and vice versa where MPI is a message passing interface. Please note that it is the adjoint variable that is subject to being exchanged in the ADM. Figures 1a and 1b could be helpful to understanding the data flow of communication in the FWM and the ADM, respectively. In the FWM, the variable U in the ghost region of processor P1 will be overwritten by the value of U received from processor P0. In the ADM, the communication needs to be reversed. The adjoint variable $a_{-}U$ in the ghost region of P1 will be sent to P0 and added to the value of $a_{-}U$ of P0 and then the $a_{-}U$ in the ghost region of P1 will be set to zero.

In the WRF model, hundreds of thousands of lines of code are automatically generated from a user-edited table called the registry (Michalakes and Schaffer 2004). The registry provides a high-level single point of control over the fundamental structure of the model data. It contains lists describing state data fields and their attributes: dimensionality, binding to particular solvers, association with WRF I/O streams, communication operations, and runtime configuration options (namelist

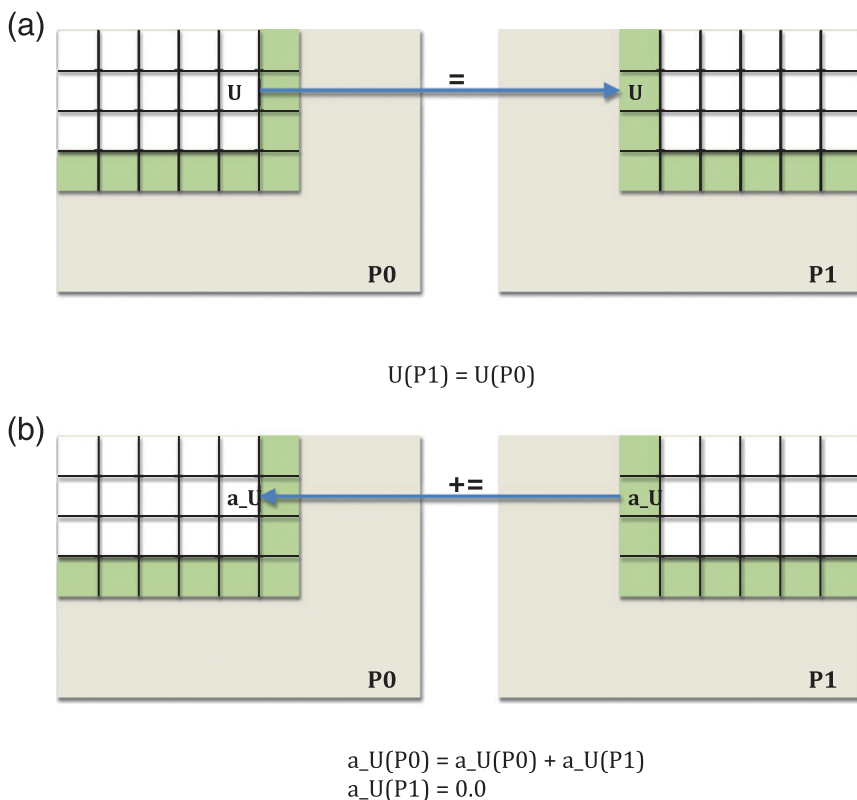


FIG. 1. Schematic diagram for (a) halo exchange between two neighbor processors in FWM and (b) adjoint of halo exchange between two neighbor processors in ADM. Gray area denotes the entire model domain. Green zone denotes the ghost area. White zone is the computational patch for each processor. Basic-state variable in the FWM is U , and adjoint variable in the ADM is a_U .

elements and their bindings to model control structures). Adding or modifying a state variable to WRF involves modifying a single line of a single file; this single change is then automatically propagated to scores of locations in the source code the next time the code is compiled.

Halo entries in the registry define communication operations in the model. Halo entries specify halo updates around a patch horizontally. A typical halo entry is

```
halo HALO_EM_C dyn_em 4:u_2,v_2
```

The first field is the keyword `halo`. The second entry, `HALO_EM_C`, which is the given name of the halo exchange template, is used in the model to refer to the communication operation being defined. The third entry denotes the associated solver, and the fourth entry is a list of information about the operation. This example specifies that four points (one cell each in the north, south, east, and west directions, respectively) of the stencil are used in updating the state arrays for fields u_2 and v_2 across the processors. During compilation,

the WRF registry automatically generates a code segment based on this halo entry:

```
...
CALL HALO_EM_C_sub (grid, local_communicator,
... )
...

```

The code snippet is included in the places where the communications are needed. At the same time, the registry also generates the subroutine `HALO_EM_C_sub`:

```
...
! for Y direction
CALL RSL_LITE_PACK (local_communicator,
grid%u_2, ..., 0, ...
CALL RSL_LITE_PACK (local_communicator,
grid%v_2, ..., 0, ...
CALL RSL_LITE_EXCH_Y
CALL RSL_LITE_PACK (local_communicator,
grid%u_2, ..., 1, ...
CALL RSL_LITE_PACK (local_communicator,
grid%v_2, ..., 1, ...
...

```

```

! for X direction
CALL RSL_LITE_PACK (local_communicator,
  grid%u_2, ...,0, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%v_2, ...,0, ...
CALL RSL_LITE_EXCH_X
CALL RSL_LITE_PACK (local_communicator,
  grid%u_2, ...,1, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%v_2, ...,1, ...
...

```

In the above-mentioned code segment, the outgoing slices of u_2 and v_2 for y direction (south–north) exchanging are packed by `RSL_LITE_PACK` into a local contiguous memory region. Therefore, one call of `RSL_LITE_EXCH_Y` is able to complete the data exchanges in the south–north direction. Once every processor receives the incoming data, `RSL_LITE_PACK` will be called again with different arguments to unpack the data to the ghost area position. The similar operations are performed in the x direction (east–west) as well.

The forward communication codes show an efficient chain of established relationships. Therefore, perturbing and adjoining the code is simple. However, it is an error-prone and time-consuming task to manually write all tangent linear and adjoint codes of communication subroutines. Since the WRF registry is able to generate halo exchange routines, there is the possibility of letting the registry generate the tangent linear and adjoint codes of halo exchanges too. To enable the WRF registry to generate the corresponding tangent linear and adjoint communication codes automatically, we modified the registry and added a new entry `halo_nta` as

```
halo_nta HALO_EM_C dyn_em 4:u_2,v_2
```

With this new entry, the registry will not only generate `HALO_EM_C_sub` but also `HALO_EM_C_TL_sub` for tangent linear codes and `HALO_EM_C_AD_sub` for adjoint codes. The subroutine `HALO_EM_C_TL_sub` looks like this:

```

...
! for Y direction
CALL RSL_LITE_PACK (local_communicator,
  grid%u_2, ...,0, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%v_2, ...,0, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%g_u_2, ...,0, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%g_v_2, ...,0, ...
CALL RSL_LITE_EXCH_Y

```

```

CALL RSL_LITE_PACK (local_communicator,
  grid%u_2, ...,1, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%v_2, ...,1, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%g_u_2, ...,1, ...
CALL RSL_LITE_PACK (local_communicator,
  grid%g_v_2, ...,1, ...
...
! for X direction

```

The TLM has exactly the same exchange stencil and data flow as the FWM; however, in addition to the basic-state fields (u_2 and v_2), the perturbation fields (g_{u_2} and g_{v_2}) exchanged as well, which is required by the TLM.

The adjoint code `HALO_EM_C_AD_sub` looks like this:

```

...
! for Y direction
CALL RSL_LITE_PACK_AD (local_communicator,
  grid%a_u_2, ...,0, ...
CALL RSL_LITE_PACK_AD (local_communicator,
  grid%a_v_2, ...,0, ...
CALL RSL_LITE_EXCH_Y
CALL RSL_LITE_PACK_AD (local_communicator,
  grid%a_u_2, ...,1, ...
CALL RSL_LITE_PACK_AD (local_communicator,
  grid%a_v_2, ...,1, ...
...
! for X direction

```

The ADM also has the same exchange stencil as the FWM and the TLM, but the data flow is reversed. On each processor, the entire patch of basic-state variables (including halo) is stored in the memory stack during the forward recomputation stage and will be restored from the memory stack during the adjoint calculation. Therefore, in adjoint communication codes, only the adjoint variables (a_{u_2} and a_{v_2}) need to be exchanged and the new subroutines `RSL_LITE_PACK_AD` will pack the data slices from where we receive data in the FWM and the TLM for sending and unpack the received data to the slices where we send out data in the FWM and the TLM.

With the upgraded WRF registry, all tangent linear and adjoint communication routines are generated automatically during compilation. To manually insert the tangent linear and adjoint communication interfaces into the TLM and the ADM is straightforward. In the TLM, the tangent linear communication routines are inserted into the same locations where the forward communication routines reside in the FWM. In the ADM,

the calling sequence of the adjoint communication routines is the reverse of that in the FWM and the TLM. From the code snippets presented above, we found that because both basic-state variables and perturbation variables are packed together, although the amount of data to communicate is doubled in the TLM, the communication latency is kept the same as in the FWM. This is highly desirable for a modern distributed memory parallel computer system with high bandwidth network. In the ADM, in general, there are two stages: the first is the forward recomputation part, which propagates the basic states within one time step and has the same communication latency and amount of data to communicate with the FWM; the second is the adjoint backward part, which has the same communication latency and amount as the first. Therefore, the total communication latency and amount of data to communicate are doubled in the ADM. However, because the adjoint code may have 3–4 times the amount of computation code than the forward code and the computations to communications ratio is still large, the communication overhead is not substantial and should not impact the parallel scalability of the ADM.

5. Parallel performance

To demonstrate the parallel efficiency of the WRFPLUS model, we prepared the initial condition and boundary conditions for a 15-km-resolution domain (not shown) with a domain center at 32.6°N, 110°E. There are 665 × 363 grid points in the horizontal direction and 45 levels in the vertical direction; the time step is 90 s. We tested this case on NCAR’s two supercomputers: Lynx and Bluefire. Lynx is a single-cabinet Cray XT5m supercomputer composed of 76 compute nodes, each with 12 processors on two hex-core AMD 2.2-GHz Opteron chips, with a total of 912 processors. Each Lynx computer node has 16 GB of memory, for 1.33 GB per processor, and 1.216 TB of memory in the system. Bluefire is an IBM clustered symmetric multiprocessing (SMP) system composed of 4096 Power 6 processors. The 4096 processors are deployed in 128 nodes, each containing 32 processors. Nodes are interconnected with an InfiniBand switch for parallel processing using MPI. We used the default compilation options and the default processors topology calculated by the WRF model. We did not do any further optimization to get the best performance. Therefore, the following results do not reflect the best performance possible on a specific supercomputer.

We ran the WRFPLUS model on 16, 32, 64, 128, 256, 512, 1024, and 2048 processors of Bluefire and measured the parallel performance. Figure 2 shows the results for

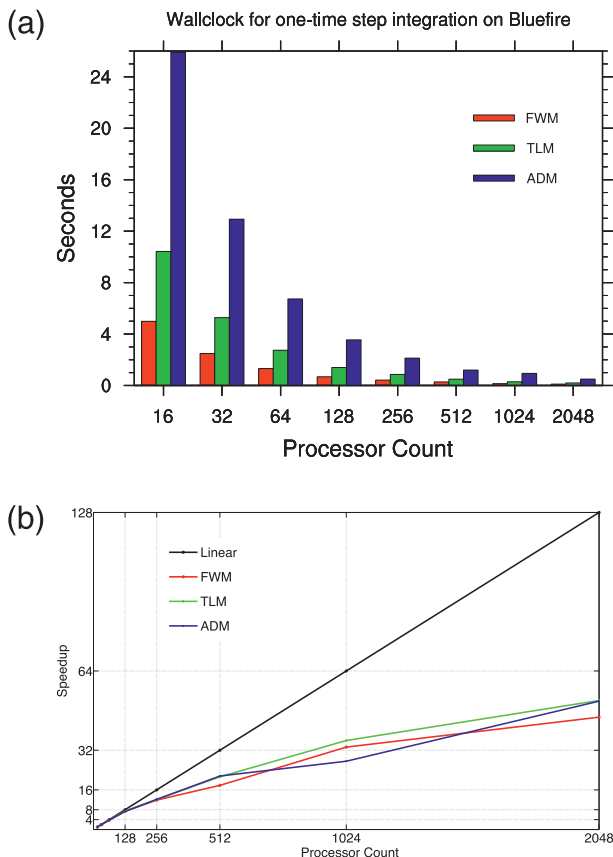


FIG. 2. Parallel performance for 1-time-step integration on Bluefire in terms of (a) wall-clock time and (b) parallel speedup.

the average wall-clock time for one-time-step integration (Fig. 2a) and speedup for the FWM, TLM, and ADM, respectively (Fig. 2b). Please note that the timing results for the FWM are different from the standard WRF run, which is compiled with 4-bytes-long real size. In the WRFPLUS, the FWM is compiled with 8-bytes-long real size. In general, because of the higher precision requirement in the WRFPLUS, the actual computational performance of the FWM is slower than the standard WRF. Speedup for N processors was calculated as the wall-clock time using 16 processors divided by the wall-clock time using N processors. The computing times for all models are considerably reduced with an increased number of processors up to 2048. In general, the TLM has a better speedup than the FWM and the ADM has a slightly better or comparable speedup than the FWM. The results confirm the successful implementation of the new parallel approach.

We also ran the WRFPLUS model on 16, 32, 64, 128, 256, and 512 processors of Lynx and measured the parallel performance. Figure 3 shows the results for the average wall-clock time for one-time-step integration (Fig. 3a) and speedup (Fig. 3b). We drew similar

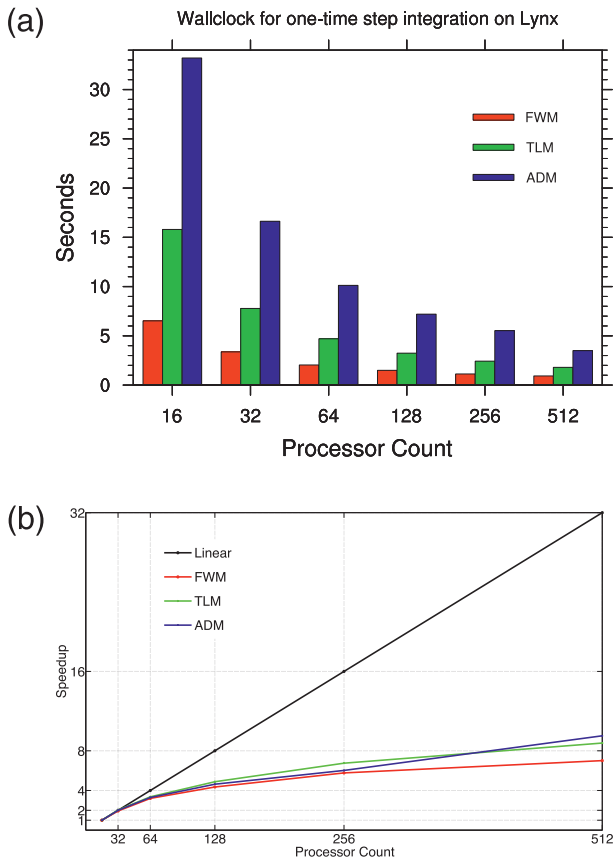


FIG. 3. As in Fig. 2, but on Lynx.

conclusions and again confirmed the high efficiency of the parallelization strategy of the WRFPLUS.

6. Implementation of an inline coupling interface in the WRFPLUS model

One of the motivations to upgrade the WRF TLM and ADM is to improve the computational performance of the WRF 4DVAR. The old WRF 4DVAR system contains a two-way coupling between the WRFDA and the WAMS (Huang et al. 2009), exchanges information via disk files, and is a multiple program multiple data (MPMD) system. During the coupling, the exchanged data are written to disk and a signal file is prepared to inform the other component that data are ready to be read. Exchange of a field between the WRFDA and the WAMS consists of gathering and scattering operations across the processors via disk files, which is very inefficient on modern distributed-memory supercomputers with a large number of processors. This limits the number of processors that can be used for high-resolution modeling.

Since the WRFDA and the WRFPLUS share the same software infrastructure including parallelization,

field definition, I/O, registry etc., it is straightforward to couple the WRFDA and the WRFPLUS to a single executable 4DVAR system, in which all information (basic states, perturbation, and adjoint forcing) from the WRFPLUS is passed as arguments to the coupling interfaces and the WRFDA fetches the data from the coupling interfaces instead of disk files.

For this purpose, three major developments were needed:

- 1) Enable the WRFPLUS model to be callable from the WRFDA with a simple application programming interface consisting of the following:
 - (i) initializing the WRFPLUS model;
 - (ii) advancing one of the WRFPLUS model components (FWM, TLM, and ADM);
 - (iii) finalizing the WRFPLUS model.
- 2) Develop a set of regridding routines that can interpolate data on the WRFPLUS grid to the WRFDA grid (and vice versa), which can be called by the WRFDA in full MPI parallel mode.
- 3) Modify the WRFDA to allow it to call the WRFPLUS with forcing data from the WRFDA and retrieve field data from the WRFPLUS (e.g., gradients).

In this paper, we discuss only the first development; the other two developments will be introduced in a separate paper (Zhang et al. 2013, manuscript submitted *J. Atmos. Oceanic Technol.*).

The WRF model already has a well-defined routine that advances the integration, which makes calling it from an external model straightforward. New initialization and finalization routines had to be coded mostly to deal with the TLM and the ADM. A namelist option `dyn_opt` was borrowed to allow the WRFPLUS to decide which model will be advanced: `dyn_opt = 2` will activate the FWM, `dyn_opt = 202` will activate the TLM, and `dyn_opt = 302` will activate the ADM.

The fully implemented interfaces as seen from the WRFDA point of view looks like this:

- Components routines: The interfaces to activate forward, tangent linear, and adjoint models
 - `wrf_run`: Interface to run forward (nonlinear) model
 - `wrf_run_tl`: Interface to run tangent linear model
 - `wrf_run_ad`: Interface to run adjoint model
- Data exchange routines: The interfaces to exchange data between the WRFDA and the WRFPLUS
 - `read_xtraj`: Read trajectories from FWM integration
 - `save_xtraj`: Save trajectories from FWM integration
 - `read_tl_pert`: Read initial perturbation for TLM integration
 - `save_tl_pert`: Save trajectories of perturbation from TLM integration

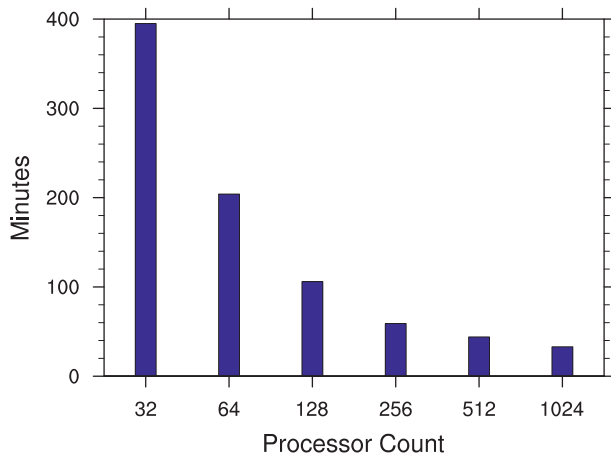


FIG. 4. Parallel performance for 4DVAR with five iterations on Bluefire.

- read_ad_forcing: Read adjoint forcing for ADM integration
- save_ad_forcing: Save initial adjoint forcing from ADM integration

These interfaces are written not only for the coupling with the WRFDA but also designed for general coupling purposes. In addition to the coupling of the WRFDA and the WRFPLUS to construct the WRF 4DVAR, we successfully coupled the WRFPLUS with the community gridpoint statistical interpolation (GSI) to construct a GSI-based WRF 4DVAR (Zhang and Huang 2013). Figure 4 shows the parallel performance of WRF 4DVAR run with 15-km-resolution continental U.S. (CONUS) domain (not shown). This domain has 450×450 horizontal grids and 51 vertical levels. The assimilation window is 6 h and the integration time step is 90 s. Only Global Telecommunication System (GTS) conventional data are assimilated. The WRF 4DVAR shows impressive scalability; that is, with the addition of more processors, the total performance of the WRF 4DVAR increased accordingly. The implementation and coupling work quite well; WRFPLUS already replaced the old WRF 4DVAR system since the version 3.4 release of WRF 4DVAR. Beside the performance benefit, there are other advantages compared to the old WRF 4DVAR system. The execution of the new 4DVAR is simpler than before because it is unnecessary to launch an MPMD collection of different executables.

7. Conclusions

In this paper we describe the implemented technique of the new tangent linear and adjoint codes of the WRF ARW core and demonstrated how, using TAPANADE (a free AD tool), the WRFPLUS was developed and

carried forward into later versions. Compared to the WAMS, the new WRFPLUS has the following improvements: 1) the tangent linear and adjoint codes are maintained to be consistent with the latest WRF changes, 2) parallelization strategy and efficiency are enhanced, 3) complete tangent linear and adjoint checks ensure the accuracy of the existing codes and newly developed codes, and 4) module interfaces for coupling were constructed in the WRFPLUS that led to a single executable WRF 4DVAR system with efficient parallel performance and scalability.

Acknowledgments. This work is supported by the Air Force Weather Agency. The WRFPLUS system benefited greatly from close collaboration with Dr. Qiang Chen and his adjoint code generator. We thank Dong-Kyou Lee and Gyu-Ho Lim of Seoul National University for their comments on the manuscript and generous support through the Korea and U.S. weather and climate centers. We thank Fuqing Zhang and Jon Poterjoy of The Pennsylvania State University for their valuable suggestions and comments on the manuscript. Mary Golden helped edit the manuscript.

REFERENCES

- Barker, D., and Coauthors, 2012: The Weather Research and Forecasting model's community variational/ensemble data assimilation system: WRFDA. *Bull. Amer. Meteor. Soc.*, **93**, 831–843.
- Cheng, B., 2006: A duality between forward and adjoint MPI communication routines. *Comput. Methods Sci. Technol.*, **12**, 23–24. [Available online at <http://www.man.poznan.pl/cmst/2005/las/03-Cheng-Gra.pdf>.]
- Errico, R. M., 1997: What is an adjoint model? *Bull. Amer. Meteor. Soc.*, **78**, 2577–2591.
- , and T. Vukicevic, 1992: Sensitivity analysis using an adjoint of the PSU–NCAR mesoscale model. *Mon. Wea. Rev.*, **120**, 1644–1660.
- Gilmour, I., L. Smith, and R. Buizza, 2001: Linear regime duration: Is 24 hours a long time in synoptic weather forecasting? *J. Atmos. Sci.*, **58**, 3525–3539.
- Hascoët, L., and V. Pascual, 2004: TAPENADE 2.1 user's guide. INRIA Tech. Rep. 0300, 78 pp. [Available online at <http://hal.inria.fr/docs/00/06/98/80/PDF/RT-0300.pdf>.]
- Huang, X.-Y., and Coauthors, 2009: Four-dimensional variational data assimilation for WRF: Formulation and preliminary results. *Mon. Wea. Rev.*, **137**, 299–314.
- Kleist, D. T., and M. C. Morgan, 2005a: Application of adjoint-derived forecast sensitivities to the 24–25 January 2000 U.S. East Coast snowstorm. *Mon. Wea. Rev.*, **133**, 3148–3175.
- , and —, 2005b: Interpretation of the structure and evolution of adjoint-derived forecast sensitivity gradients. *Mon. Wea. Rev.*, **133**, 466–484.
- Langland, R. H., R. Gelaro, G. D. Rohaly, and M. A. Shapiro, 1999: Target observations in FASTEX: Adjoint-based targeting procedures and data impact experiments in IOP17 and IOP18. *Quart. J. Roy. Meteor. Soc.*, **125**, 3241–3270.

- Li, Z., A. Barcilon, and I. M. Navon, 1999: Study of block onset using sensitivity perturbations in climatological flows. *Mon. Wea. Rev.*, **127**, 879–900.
- Michalakes, J., and D. Schaffer, 2004: WRF tiger team documentation: The registry. UCAR Tech. Rep., 17 pp. [Available online at http://www.mmm.ucar.edu/wrf/WG2/software_2.0/registry_schaffer.pdf.]
- Navon, I. M., X. Zou, J. Derber, and J. Sela, 1992: Variational data assimilation with an adiabatic version of the NMC spectral model. *Mon. Wea. Rev.*, **120**, 1433–1446.
- Rabier, F., E. Klinker, P. Courtier, and A. Hollingsworth, 1996: Sensitivity of forecast errors to initial conditions. *Quart. J. Roy. Meteor. Soc.*, **122**, 121–150.
- Skamarock, W. C., and Coauthors, 2008: A description of the advanced research WRF version 3. NCAR Tech. Note NCAR/TN-475+STR, 113 pp. [Available from UCAR at 360 Communications, P.O. Box 3000, Boulder, CO 80307-3000.]
- Utke, J., L. Hascoët, P. Heimbach, C. Hill, P. Hovland, and U. Naumann, 2009: Toward adjoinable MPI. *Proc. 10th IEEE Int. Workshop on Parallel and Distributed Scientific and Engineering (PDSEC'09)*, Rome, Italy, IEEE, 8 pp. [Available online at <http://www-sop.inria.fr/tropics/papers/Utke2009TAM.pdf>.]
- Vukicevic, T., 1991: Nonlinear and linear evolution of initial forecast errors. *Mon. Wea. Rev.*, **119**, 1602–1611.
- Xiao, Q., X. Zou, M. Pondecà, M. A. Shapiro, and C. S. Velden, 2002: Impact of *GMS-5* and *GOES-9* satellite-derived winds on the prediction of a NORPEX extratropical cyclone. *Mon. Wea. Rev.*, **130**, 507–528.
- , and Coauthors, 2008: Application of an adiabatic WRF adjoint to the investigation of the May 2004 McMurdo, Antarctica, severe wind event. *Mon. Wea. Rev.*, **136**, 3696–3713.
- Zhang, X., and X. Y. Huang, 2013: Development of GSI-based WRF 4DVAR system. Preprints, *17th Conf. on Integrated Observing and Assimilation Systems for the Atmosphere, Oceans, and Land Surface*, Austin, TX, Amer. Meteor. Soc., 4A.2. [Available online at <https://ams.confex.com/ams/93Annual/webprogram/Paper217182.html>.]
- , —, J. Liu, and J. Poterjoy, 2013: Development of an efficient regional four-dimensional variational data assimilation system for WRF. *J. Atmos. Oceanic Technol.*, submitted.