

1 Title: Cloud Archiving and Data Mining of High-Resolution Rapid Refresh Forecast Model

2 Output

3 Journal: Computers and Geosciences

4 Authors: Brian Blaylock¹, John Horel¹, Samuel T. Liston²

5 ¹University of Utah, Department of Atmospheric Sciences

6 ²University of Utah, Center for High Performance Computing

7 Corresponding Author: Brian Blaylock, brian.blaylock@utah.edu

8 Address: 135 S 1460 E, Rm 819, Salt Lake City, UT 84112

9 Keywords: object data storage; data stewardship; atmospheric modeling; cloud computing

10

11

12

13

14

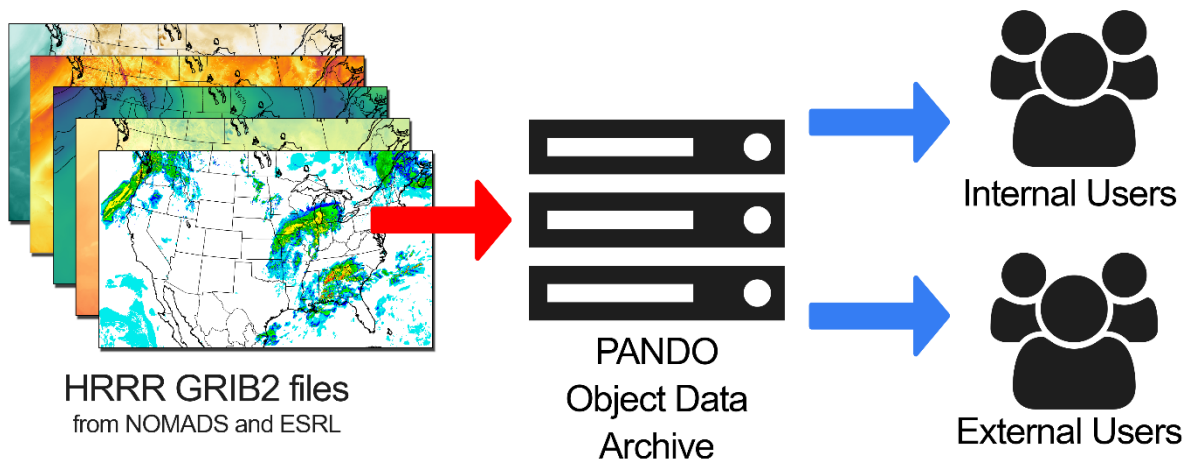
15 **Abstract**

16 Weather-related research often requires synthesizing vast amounts of data that need
17 archival solutions that are both economical and viable during and past the lifetime of the project.
18 Public cloud computing services (e.g., from Amazon, Microsoft, or Google) or private clouds
19 managed by research institutions are providing object data storage systems potentially
20 appropriate for long-term archives of such large geophysical data sets. We illustrate the use of a
21 private cloud object store developed by the Center for High Performance Computing (CHPC) at

22 the University of Utah. Since early 2015, we have been archiving thousands of two-dimensional
23 gridded fields (each one containing over 1.9 million values over the contiguous United States)
24 from the High-Resolution Rapid Refresh (HRRR) data assimilation and forecast modeling
25 system. The archive is being used for retrospective analyses of meteorological conditions during
26 high-impact weather events, assessing the accuracy of the HRRR forecasts, and providing initial
27 and boundary conditions for research simulations. The archive is accessible interactively and
28 through automated download procedures for researchers at other institutions that can be tailored
29 by the user to extract individual two-dimensional grids from within the highly compressed files.
30 Characteristics of the CHPC object storage system are summarized relative to network file
31 system storage or tape storage solutions. The CHPC storage system is proving to be a scalable,
32 reliable, extensible, affordable, and usable archive solution for our research.

33

34 **Graphical Abstract**



35

36 1. Introduction

37 Weather research and operational weather forecasting depends heavily on evaluating the
38 output from high-resolution regional numerical weather prediction models. The Weather
39 Research and Forecasting (WRF) model is the world's most widely-used regional numerical
40 weather prediction model relied upon operationally for life-saving weather forecasts and for
41 aviation, energy, fire prediction, surface transportation, and water resource management
42 applications (Powers et al. 2017). The High-Resolution Rapid Refresh (HRRR) version of the
43 WRF model, developed by the Earth Systems Research Lab (ESRL), is an hourly updating,
44 cloud-resolving, convection-allowing model run operationally by the National Centers for
45 Environmental Prediction's Environmental Modeling Center (EMC) (Benjamin et al. 2016).
46 Output from most U.S. operational weather models run by EMC are available on EMC servers
47 for the current day and then archived by the National Centers for Environmental Information
48 (NCEI). However, the voluminous HRRR model output available each hour for forecast
49 durations from 0-18 h with a grid spacing of 3 km over the contiguous United States (1.9 million
50 grid points) is not yet available from NCEI. To archive in a highly compressed format, a
51 representative sample of the output generated by the operational HRRR model requires over 200
52 TB of disk space per year.

53 Researchers rely heavily on output from regional models such as HRRR and WRF to
54 diagnose the interplay between complex atmospheric processes on spatial scales from $10^2 - 10^6$
55 m and temporal scales from $10^2 - 10^7$ s (Benjamin et al. 2016; Powers et al. 2017). A common
56 research strategy is to focus on case studies of specific weather events as a practical approach to
57 manage the TBs of output generated by the models (e.g., Blaylock et al. 2017, Crosman and
58 Horel 2017). With continued growth in computing capabilities, numerical simulations will

59 continue to transition to finer spatial and temporal resolution over increasingly large regional
60 domains. As these models grow, so does the storage space and monetary cost required to archive
61 model output. Of course, large data storage needs are ubiquitous throughout the atmospheric
62 sciences, for example, to archive satellite imagery (Moody et al. 2016) or multi-decadal
63 numerical simulations of the climate system (Taylor et al. 2012).

64 Molthan et al. (2015) highlight that cloud computing resources (computational services
65 delivered over networks) are providing new capabilities for supporting numerical weather
66 prediction and are a potential solution to archive large volumes of data (Armbrust et al. 2010;
67 Sandholm and Lee 2014). To meet these needs, Sandholm and Lee (2014) described how these
68 services need to be: scalable; fault-tolerant; reliable; high-performance; and easy to use, manage,
69 monitor, and provision efficiently and economically. Public cloud services provided by
70 corporations (e.g. Amazon, Google, or Microsoft) or research consortia (e.g. Open Science Data
71 Cloud, <https://www.opensciencedatacloud.org/>) are increasingly viable options to meet those
72 requirements, although understanding the extent to which they are economical can be difficult
73 (Chou 2015; Amazon Web Services 2017a). Private cloud services are defined as being operated
74 by an organization for which hardware, networking, storage, and other infrastructure are not
75 directly shared with other organizations (Mell and Grance 2011). The Center for High
76 Performance Computing (CHPC) at the University of Utah provides private cloud services
77 through a data center located off campus.

78 The objective of this paper is to illustrate the utility and cost effectiveness of a PB disk-
79 based object storage data system managed by the CHPC for archiving large data sets. The
80 capabilities of object data storage systems for geoscience applications will be illustrated in terms
81 of an archive of operational and experimental forecasts from the HRRR model in the contiguous

82 United States and Alaska from early 2015 to the present. While we have relied extensively over
83 the years on other CHPC storage media (such as a robotic tape archive system and over 100 TB
84 of network file system disk storage), the object data storage system is meeting several of our
85 interwoven needs that are less practical using other traditional data archival approaches: (1)
86 efficient expandable storage for thousands of large data files; (2) data analysis using fast retrieval
87 of user selectable byte-ranges within those data files; and (3) the ability to have the data publicly
88 accessible to the atmospheric science research community.

89 The remainder of the paper describes how the archive is built and how users can access
90 the data (section 2), followed by applications for which data from the HRRR archive have been
91 used (section 3), and concludes with a discussion of the growing need for large archives and
92 some limitations that should be resolved in the future (section 4).

93

94 **2. Methods**

95 *2.1 Pando Object Storage System*

96 The CHPC has dramatically increased its network file system data storage capabilities
97 over the past 10 years from ~400TB to ~14PB due to decreased hardware costs and development
98 of cost-effective storage solutions (Center for High Performance Computing 2017). However,
99 archival storage capacity primarily in terms of a robotic tape system has not increased as rapidly,
100 leaving a large fraction of the data without backup. To help mitigate this shortcoming, CHPC
101 developed a disk-based object storage solution referred to as Pando (named for a vast stand of
102 aspen trees in Utah that is thought to be the largest and oldest single living organism). Currently
103 at 1 PB in usable capacity, Pando was developed at lower cost than other archival options and
104 has greater resiliency, accessibility, and expandability. Researchers lease dedicated amounts of

105 archival space over a 5-year span to help recover some of the costs for Pando. They then manage
106 their own space, which helps reduce CHPC's administrative burden to manage the archive.

107 The CHPC took into consideration that an improved archival system needed to scale to a
108 much larger size than what might be affordable initially. Large network file systems or
109 Redundant Array of Independent Disks (RAID) sets do not scale well as the number and size of
110 drives increase, particularly since recovering and repairing after an error or disk corruption may
111 require disks to be offline for many days. The CHPC selected Red Hat's Ceph object-based open
112 source storage system (Maltzahn et al. 2010) to address the shortcomings of both RAID and file
113 systems based on published performance comparisons (e.g., Poat et al. 2015) and testing over
114 several years. Low-level operations, such as block or file level I/O, are managed by a software
115 layer that manipulates objects for the user or administrator such that expensive RAID controllers
116 are not necessary and archived objects can be replicated or made redundant according to
117 configurable parameters.

118 Pando was formatted using the 6+3 erasure coding, i.e., all objects are broken into 9
119 pieces—6 data pieces and 3 redundancy pieces necessary for data protection and reconstruction.
120 The initial 1 PB Pando archive consists of 9 storage servers each with sixteen 8 TB drives that
121 are coordinated by 3 monitor nodes that efficiently maintain the map of the objects in the system
122 (Fig. 1). If the file system on a single drive becomes corrupt, then: (1) that drive is logically
123 removed by the system administrator; (2) the administrator recreates the file system and logically
124 adds it back in; and (3) the objects are redistributed within the new file system automatically by
125 the Ceph software to maintain the configured level of redundancy. The 6+3 erasure coding
126 ensures no data loss even if every disk fails on three servers. The Pando system has the capacity
127 to contain 44 servers before additional network infrastructure must be purchased making it

128 expandable to approximately 5PB with current drive capacities. To ensure that Pando is in
129 production past disk warranty periods, Ceph can transparently migrate the data to new hardware
130 when old hardware is retired.

131 The Amazon Simple Storage Service (S3) has been implemented on Pando through a
132 Reliable Autonomic Distributed Object Store (RADOS) Gateway node to focus on
133 long-term storage needs separate from the other mounted file systems available to CHPC users
134 (Nawaz et al. 2016). The RADOS Gateway node (Fig. 1) serves as an interface between client
135 computers and objects managed by the RADOS software layer. Present usage suggests that
136 additional RADOS Gateway nodes will be necessary in the future to avoid throughput
137 bottlenecks (speeds of only 5 GB s^{-1} during high loads) that limit optimal utilization of the Pando
138 system. Objects are most efficiently uploaded to Pando from the CHPC local file systems using
139 rclone (Wood 2017), which is open source software commonly used to download or upload files
140 between hard disk and cloud storage systems.

141 *2.2 HRRR Data Archive*

142 Several implementations of the HRRR modeling system have been developed by ESRL
143 researchers with staff at EMC maintaining its operational version for the contiguous United
144 States (Benjamin et al. 2016). To support air quality research at the University of Utah (Horel et
145 al. 2016; Blaylock et al. 2017), we started archiving operational HRRR analysis (forecast hour 0)
146 output files beginning April 2015 on local network file system disks obtained from the NOAA
147 Operational Model Archive and Distribution System (NOMADS). Other research projects led us
148 to download selected meteorological fields from the operational HRRR 1-18 h forecast files
149 beginning in summer 2017 and analysis and forecast fields from experimental versions of the
150 HRRR for the contiguous United States and Alaska. The thousands of 2-dimensional

151 meteorological fields available from the HRRR are stored as gridded binary-2 (GRIB2) files, a
152 highly efficient binary format that relies on Joint Photographic Experts Group (JPEG) 2000
153 image compression (Silver and Zender 2017).

154 By early 2017, local file system storage for the HRRR products grew to over 20 TB with
155 the expectation that by later in 2017, over 100 GB of model grids would be added per day. That
156 storage approach was becoming unwieldy to manage across multiple file server partitions and
157 not practical to facilitate access to the archive for an increasing number of atmospheric science
158 researchers external to the University of Utah, who became aware of it through online searches
159 for HRRR model output. After initial testing of the Pando system, all the locally-archived HRRR
160 files were transferred to it and removed from the local file system.

161 Since EMC and ESRL provide efficient access for anyone interested in HRRR model
162 output for the current and previous day (Bowman and Lees 2015), we prefer external users to not
163 overwhelm our archival system by requesting what is already easily available from those
164 sources. We execute download scripts after 00 UTC to retrieve files for the previous day to our
165 local CHPC network file storage, a process that can take several hours to complete even with
166 multithreading. The files are then copied to the Pando archive using the open source rclone
167 utility. The s3cmd utility is used to change permissions for each file from private to public so
168 they can be accessed by other researchers at the University of Utah and elsewhere.

169 The present implementation of Ceph on Pando limits the ability to view the contents or
170 manipulate the data object files. Rather, each file has a unique URL that can be used to download
171 it via HTTPS. While anyone can attempt to directly download such files from the archive, web
172 pages have been developed for identifying which files are available to simplify interactive
173 downloads (<https://hrrr.chpc.utah.edu>; Fig. 2). Users are encouraged to avoid excessive reliance

174 on the interactive pages and create automated download procedures using wget or cURL with
175 example code provided on the aforementioned web page.

176 Since most users prefer to access a relatively small number of the meteorological fields
177 contained within each of the large HRRR GRIB2 files, it is cumbersome to retrieve the entire file
178 and then process it to extract the fields of interest. To facilitate access to specific 2-dimensional
179 fields, we use the wgrib2 tool (Climate Prediction Center 2017) to create a metadata file for each
180 GRIB2 file and provide that information on a local web server since there is no need to store
181 them as objects in Pando. These index files contain for each field its abbreviated variable name,
182 vertical level, beginning byte, time of the model run, and forecast hour. Hence, it is
183 straightforward to derive the corresponding byte range for a variable and retrieve using cURL its
184 2-dimensional field. Unfortunately, it is not currently possible to retrieve a byte range within a
185 GRIB2 formatted file for a subsection of the two-dimensional grid (e.g., for a state or regional
186 area of interest). This is a present limitation of object storage and GRIB2 file formats that may be
187 solved through continued development of object storage systems or archiving the gridded data in
188 a different file format. Hence, the smallest granule that can be retrieved from a GRIB2 HRRR
189 file is a single field over that entire domain (~1 MB). Multiprocessing and multithreading
190 techniques such as those available using Python's multiprocessing module can be leveraged to
191 spread the work across multiple cores and reduce download time and greatly increase the data
192 processing speed when fields from multiple files are needed. We have developed Python multi-
193 processor procedures that rely on basic cURL commands to efficiently access the HRRR files
194 from a single dedicated CHPC server. For example, computing the minimum, mean, and
195 maximum wind speed from nearly 17,000 hourly analyses at the 1.9 million grid points in the
196 operational HRRR model was done in less than 15 minutes using 30 processors.

197 The current HRRR archive directory tree for both the Pando and metadata archive is
198 branched by model type (operational HRRR, experimental HRRR, and experimental HRRR
199 Alaska), by file type (sfc files contain a selection of 2-dimensional fields while many more 2-
200 dimensional fields at fixed pressure levels in the vertical as well as other levels are available in
201 the prs files), and by date (year, month, and day).

```
202       HRRR/  
203       ├── oper/  
204       │   ├── sfc/  
205       │   │   ├── YYYYMMDD/  
206       │   ├── prs/  
207       │   │   ├── YYYYMMDD/  
208       ├── alaska/  
209       │   ├── sfc/  
210       │   │   ├── YYYYMMDD/  
211       │   ├── prs/  
212       │   │   ├── YYYYMMDD/  
213       ├── exp/  
214       │   ├── sfc/  
215       │   │   ├── YYYYMMDD/
```

216 Each file within the daily directories follow the same naming convention used by NOMADS
217 when the file is first downloaded (files from ESRL are renamed to match the NOMADS naming
218 convention). The files are named by the model type, the initialization hour, variable field, and the
219 forecast hour ([hrrr/hrrrAK/hrrrX].t[hour]z.wrf[sfc/prs]f[forecast].grib2). For example, the
220 following request will download the full surface field file from the operational HRRR analysis
221 for 14:00 UTC 5 April 2017:

222 <https://pando-rgw01.chpc.utah.edu/HRRR/oper/sfc/20170405/hrrr.t14z.wrfsfcf00.grib2>.

223 Metadata for the corresponding HRRR file can be found in the GRIB2 index file located here:

224 <https://api.mesowest.utah.edu/archive/HRRR/oper/sfc/20170405/hrrr.t14z.wrfsfcf00.grib2.idx>.

225 The index file can be used to request specific variables within a byte range. If a user was only
226 interested in 10 m gusts, then the index file indicates that the byte range for the gusts variable for
227 that file is between 3478099 and 4879421. Using cURL, a user can download the gust variable
228 from the larger file as follows:

```
229 curl -o downloaded_file.grib2 --range 2757386-4110515 https://pando-  
230 rgw01.chpc.utah.edu/HRRR/oper/sfc/20170405/hrrr.t14z.wrfsfcf00.grib2.
```

231

232 **3. Applications**

233 *3.1 High-Impact Weather Events*

234 While voluminous sets of graphics of analysis and forecasts fields from the HRRR model
235 runs are generated routinely by ESRL, EMC, academic institutions, and commercial sources of
236 weather information, those usually depict only conditions within the past few days and only
237 show a small fraction of the information contained in the HRRR GRIB2 files. The HRRR Pando
238 archive provides users access to all the fields contained in the HRRR grib2 files. These files can
239 be used to create customized graphics of high impact weather events or other features of interest
240 to the user. For example, the major New England snowstorm on 14 March 2017 is depicted by
241 the HRRR mean sea level pressure analysis valid at 1700 UTC 14 March 2017 (Fig. 3).

242 Hourly changes in atmospheric conditions at specific locales can be examined by
243 downloading the requisite grids each hour, which can be easily retrieved from the Pando archive
244 using the procedures described above. Figure 4 illustrates the conditions analyzed by the HRRR
245 centered on 2100 UTC 27 April 2017 at which time a wildfire near O'Donnell Texas traversed
246 across the site of a West Texas Mesonet station (Schroeder et al. 2005) as evident by the 58°C
247 observed 2-m air temperature at that time. The HRRR hourly analyses closely track observations

248 (albeit not the temperature spike associated with the fire) as well as provide additional diagnostic
249 variables, such as winds at 80 m above ground level and estimates of the boundary layer depth.

250 Since the primary purpose of the operational HRRR model is to provide short-term (0-18
251 h) weather forecast guidance updated every hour to predict severe weather (Benjamin et al.
252 2016), assessing the model's ability to properly forecast such conditions is of high interest. For
253 example, 30 tornadoes and hundreds of reports of hail and high winds were received on 4-5 April
254 2017 from Missouri to Ohio extending southward to Alabama and Georgia (Storm Prediction
255 Center 2017). Airline operations in Atlanta were severely affected on 5 April causing thousands
256 of delayed or canceled flights. Figure 5 contrasts the simulated composite reflectivity and gust
257 analyses from the HRRR model at 1400 UTC 5 April 2017 to the 16 h forecast from the HRRR
258 run initialized 2200 UTC 4 April 2017. The model forecast at 16 h highlights many of the
259 locations that later received heavy precipitation and strong winds.

260

261 *3.2 HRRR Model Composites*

262 Statistics derived over long-time intervals from model output can provide useful
263 information, such as availability of wind and solar energy resources (James et al. 2017) or
264 identifying model performance characteristics (Katona et al. 2016, Ikeda et al. 2017).
265 Preliminary basic statistics (minimum, mean, maximum, and percentiles) of meteorological
266 variables (temperature, wind speed, snow cover, lightning, etc.) have been derived from the 2-
267 year archive of HRRR analysis grids. Multiprocessing techniques were used to speed up
268 downloading the files from the archive and processing the grids for each of the 1.9 million grid
269 points. Figure 6 shows the 95th percentile of the 10 m gusts analyzed by the operational HRRR
270 at 2300 UTC during all days between 18 April 2015 and 30 March 2017. Such statistics are

271 intended to be used to provide realistic bounds for observations of wind and other variables at
272 over 25,000 locations in the United States that are available within the past 20 years as well as
273 received continuously as part of the MesoWest and SynopticLabs projects (Horel et al. 2002;
274 SynopticLabs 2017). Simultaneous calculations that require less memory (e.g., extreme and
275 mean values) were completed in about 15 minutes for one variable over the entire contiguous
276 United States. Brute-force approaches to calculate multiple percentile values (e.g., 1st, 5th, 10th,
277 90th, 95th, and 99th) for each hour of the day necessary to generate Figure 6 required storing more
278 values in memory and required roughly an hour for a single variable. Improved approaches using
279 approximation techniques are possible to efficiently compute percentiles and other statistics and
280 avoid excessive memory consumption on our compute nodes.

281

282 *3.3 Initializing WRF Simulations*

283 The original impetus for our archive of the HRRR output was to obtain the best possible
284 high-resolution WRF simulations over northern Utah to understand a poor air quality episode in
285 the vicinity of Salt Lake City during 17-18 June 2015. Blaylock et al. (2017) ran a 1 km WRF
286 simulation for northern Utah with initial and boundary conditions obtained from the HRRR
287 hourly analyses beginning at 0000 UTC 14 June 2015 and continuing until 0700 UTC 19 June
288 2015.

289 While many researchers initialize high-resolution model simulations from operational and
290 reanalysis modeling systems (e.g., Foster et al. 2017; Li et al. 2017), the HRRR provides
291 significant advantages in terms of its 3 km grid spacing, hourly output files, and advanced data
292 assimilation techniques. To the best of our knowledge, the study by Blaylock et al. (2017) was
293 the first one to use HRRR analyses to initialize and provide the requisite lateral boundary

294 conditions for WRF research simulations. While ESRL maintains an internal tape archive of
295 HRRR model output, the HRRR archive on Pando is currently the only readily available resource
296 for other researchers to initialize high-resolution WRF simulations with HRRR boundary
297 conditions. While it is recommended to initialize WRF simulations with native or model-level
298 HRRR files, we don't archive the native level files at this time due to its large file sizes (> 600
299 GB per file). However, WRF can be initialized with the HRRR pressure-level analysis files
300 available on Pando. The steps required to initialize WRF with HRRR boundary conditions have
301 been documented by Blaylock (2017).

302

303 **4. Discussion and Conclusions**

304 The management and distribution of large geoscience data sets have received increasing
305 attention, particularly given the explosion in public and private cloud-based resources. For
306 example, an Amazon Web Service (AWS) S3 object store hosts the level 2 retrospective and
307 real-time archive of Next Generation Weather Radar (NEXRAD) data (Amazon Web Services
308 2017b). Our research group in the Department of Atmospheric Sciences uses Amazon AWS
309 including its S3 object store for other applications that require uninterruptible computational
310 resources and require a relatively fixed small amount of disk storage (SynopticLabs 2017). The
311 complexity and volatility in the egress costs to upload or download data depending on the
312 policies of each public cloud storage facility precluded our use of one of them for the HRRR
313 archive.

314 The private cloud CHPC Pando object storage archive has made it possible to efficiently
315 archive, access, and analyze the HRRR model output. Pando is also being used by other
316 atmospheric scientists, anthropologists, geneticists, and cancer researchers at the University of

317 Utah. Our HRRR archive has many of the properties of an ideal data archive described by
318 Kruger et al. (2006)—it is scalable, extensible, inexpensive, and usable. Having fixed leasing
319 costs over a 5-year period allows us to plan as our archival needs grow. The private cloud Pando
320 system provides faster access to our long-term data archive for our needs as well as provide
321 reasonable access times for the several dozen researchers outside the University of Utah that
322 have already discovered its utility in the short time that the archive has been available.

323 The major limitation of the present Pando object storage systems is that Ceph constrains
324 how the objects can be managed and accessed. Red Hat now supports Ceph File System (Ceph
325 FS, Red Hat 2017) as a Portable Operating System Interface (POSIX) compliant file system that
326 is more flexible to handle the objects in the storage cluster. However, S3-type objects still must
327 be downloaded to a local disk before the data contained within them can be processed. To avoid
328 excessive downloading of data not of interest to a user, the highly efficient GRIB2 format of the
329 HRRR model output allows selecting by byte range and returning only the fields of interest from
330 the many two-dimensional fields contained within an object. Other file formats, such as
331 Hierarchical Data Format Version 5 or Network Common Data Format, may eventually allow
332 subsetting of S3 objects by variable, region, single grid point, all vertical levels at a point, etc.,
333 but that capability is not presently available.

334 We expect that NCEI or other government or institutional repositories will begin to
335 archive operational HRRR model output at some point. Although long-term archives of evolving
336 experimental versions of models are seldom undertaken, having the ability as we do to compare
337 output from experimental and operational versions of the same model makes it possible to assess
338 model improvements more efficiently. Research agencies such as the National Science
339 Foundation now require data management plans that describe what will happen to the data and

340 metadata that led to the research results. While a small number of geoscience data repositories
341 exist (e.g., the National Center for Atmospheric Research), those entities have strict standards for
342 accepting large data sets that are often difficult to meet. At the present time, geoscience data
343 journals require that data sets be in such data repositories prior to publication such as that by
344 Jacques et al. (2016). Academic institutions will increasingly need to consider having facilities
345 like the Pando archive to effectively meet those data stewardship requirements. However, it
346 remains unclear whether those institutions are willing to subsidize the cost of maintaining large
347 archives that are necessary to store results once research projects have been completed and funds
348 are no longer available from the granting agencies.

349

350 **Acknowledgements**

351 We appreciate the support and resources made available by the Center for High Performance
352 Computing at the University of Utah. We would like to thank the model developers at ESRL and
353 EMC for the ongoing development of the HRRR as well as providing the model output. We also
354 would like to thank Chris Galli for his suggestions about this research and comments on the
355 manuscript. Funding: This research has been supported by the National Science Foundation
356 (NSF Grant 1443046) and the NOAA Collaborative Science, Technology, and Applied Research
357 (CSTAR) Program (NOAA Grant NA13NWS4680003).

358 **References**

- 359 Amazon Web Services, 2017a. Amazon EC2 pricing. URL:
360 <https://aws.amazon.com/ec2/pricing/>.
- 361 Amazon Web Services, 2017b. NEXRAD data archive. URL: [https://aws.amazon.com/noaa-big-](https://aws.amazon.com/noaa-big-data/nexrad/)
362 [data/nexrad/](https://aws.amazon.com/noaa-big-data/nexrad/).
- 363 Armbrust, M., and coauthors, 2010. A view of cloud computing. *Communications of the ACM*,
364 53, 50-58. doi: 10.1145/1721654.1721672.
- 365 Benjamin, S., and coauthors, 2016. A North American Hourly Assimilation and Model Forecast
366 Cycle: The Rapid Refresh. *Monthly Weather Review*, 144, 1669-1694,
367 doi:10.1175/MWR-D-15-0242.1.
- 368 Blaylock, B., 2017. How to initialize WRF with HRRR boundary conditions. URL:
369 http://home.chpc.utah.edu/~u0553130/Brian_Blaylock/hrrr.html.
- 370 Blaylock, B., J. Horel, E. Crosman, 2017. Impact of Lake Breezes on Summer Ozone
371 Concentrations in the Salt Lake Valley. *Journal of Applied Meteorology and Climatology*,
372 56, 353-370, doi: 10.1175/JAMC-D-16-0216.1.
- 373 Bowman, D., J. Lees, 2015. Near real time weather and ocean model data access with
374 rNOMADS. *Computers & Geosciences*, 78, 88-95. doi: 10.1016/j.cageo.2015.02.013.
- 375 Center for High Performance Computing, 2017. Storage services at CHPC. URL:
376 https://www.chpc.utah.edu/resources/storage_services.php.
- 377 Chou, D., 2015. Cloud computing: A value creation model. *Computer Standards & Interfaces*,
378 38, 72-77. doi: 10.1016/j.csi.2014.10.001.
- 379 Climate Prediction Center, 2017. WGRIB2: Utility to read and write grib2 files. URL:
380 <http://www.cpc.ncep.noaa.gov/products/wesley/wgrib2/>.

381 Crosman, E., J. Horel, 2017. Large-eddy simulations of a Salt Lake Valley cold-air pool.
382 *Atmospheric Research*, 193, 10-25. doi: 10.1016/j.atmosres.2017.04.010.

383 Foster, C., E. Crosman, J. Horel, 2017. Simulations of a Cold-Air Pool in Utah's Salt Lake
384 Valley: Sensitivity to Land Use and Snow Cover. *Boundary-Layer Meteorology*, 164, 63-87.
385 doi: 10.1007/s10546-017-0240-7.

386 Horel, J., and Coauthors, 2002. Mesowest: Cooperative Mesonets in the Western United States.
387 *Bulletin of the American Meteorological Society*, 83, 211–225, doi: 10.1175/1520-
388 0477(2002)083<0211:MCMITW>2.3.CO;2.

389 Horel, J., E. Crosman, A. Jacques, B. Blaylock, S. Arens, A. Long, J. Sohl, R. Martin, 2016.
390 Influence of the Great Salt Lake on summer air quality over nearby urban areas. *Atmospheric*
391 *Science Letters*, 17, 480-486. doi: 10.1002/asl.680.

392 Ikeda, K., M. Steiner, G. Thompson, 2017. Examination of mixed-phase precipitation forecasts
393 from the High-Resolution Rapid Refresh model using surface observations and sounding
394 data. *Weather and Forecasting*, 32, 949-967. doi: 10.1175/WAF-D-16-0171.1.

395 Jacques, A., J. Horel, E. Crosman, F. Vernon, J. Tytell, 2016. The Earthscope US Transportable
396 Array 1 Hz Surface Pressure Dataset. *Geoscience Data Journal*, 3, 29–36. doi:
397 10.1002/gdj3.37.

398 James, E., S. Benjamin, M. Marquis, 2017. A unified high-resolution wind and solar dataset from
399 a rapidly updating numerical weather prediction model. *Renewable Energy*, 102, 390-
400 405. doi: 10.1016/j.renene.2016.10.059.

401 Katona, B., P. Markowski, C. Alexander, S. Benjamin, 2016. The Influence of Topography on
402 Convective Storm Environments in the Eastern United States as Deduced from the
403 HRRR. *Weather and Forecasting*, 31, 1481-1490. doi: 10.1175/WAF-D-16-0038.1.

404 Kruger, A., R. Lawrence, E. Dragut, 2006. Building a terabyte NEXRAD radar database for
405 hydrometeorology research. *Computers and Geosciences*, 32, 247-258. doi:
406 10.1016/j.cageo.2005.06.001.

407 Li, Y. and coauthors, 2017. A Numerical Study of the June 2013 Flood-Producing Extreme
408 Rainstorm over Southern Alberta. *J. Hydrometeor*, 18, 2057-2078, doi: 10.1175/JHM-D-
409 15-0176.1.

410 Mell, P., T. Grance, 2011. The NIST Definition of Cloud Computing: Recommendations of the
411 National Institute of Standards and Technology. URL:
412 <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>. doi:
413 10.6028/NIST.SP.800-145.

414 Maltzahn, C., E. Molina-Estolano, A. Khurana, A. Nelson, S. Brandt, S. Weil, 2010. Ceph as a
415 Scalable Alternative to the Hadoop Distributed File System. *USENIX Magazine*, 35, 38-
416 49. URL: <http://static.usenix.org/publications/login/2010-08/openpdfs/maltzahn.pdf>.

417 Molthan, A., J. Case, J. Venner, R. Schroeder, M. Checchi, B. Zavodsky, A. Limaye, R.
418 O'Brien, 2015. Clouds in the cloud: Weather forecasts and applications within cloud
419 computing environments. *Bulletin of the American Meteorological Society*, 96, 1369-
420 1379. doi: 10.1175/BAMS-D-14-00013.1.

421 Moody D., M. Warren, S. Skillman, R. Chartrand, S. Brumby, R. Keisler, T. Kelton, M. Mathis,
422 2016. Building a living atlas of the Earth in the cloud. *50th Asilomar Conference on*
423 *Signals, Systems and Computers*. 1273-1277. doi: 10.1109/ACSSC.2016.7869578.

424 Nawaz, H., G. Juve, R. da Silva, E. Deelman, 2016. Performance Analysis of an I/O-Intensive
425 Workflow executing on Google Cloud and Amazon Web Services. *Parallel and*

426 *Distributed Processing Symposium Workshops, 2016 IEEE International*. IEEE, 2016.
427 535-544. doi: 10.1109/IPDPSW.2016.90.

428 Poat, M., J. Lauret, W. Betts, 2015. POSIX and Object Distributed Storage Systems Performance
429 Comparison Studies with Real-Life Scenarios in an Experimental Data Taking Context
430 Leveraging OpenStack Swift & Ceph. *Journal of Physics: Conference Series*. 664, 1-9.
431 doi: 10.1088/1742-6596/664/4/042031.

432 Powers, J. and coauthors, 2017. The Weather Research and Forecasting (WRF) Model:
433 Overview, System Efforts, and Future Directions. *Bulletin of the American*
434 *Meteorological Society*, In Press. doi: 10.1175/BAMS-D-15-00308.1.

435 RedHat 2017: CephFS: Ceph File System. URL: <http://docs.ceph.com/docs/master/cephfs/>.

436 Sandholm, T., D. Lee, 2014. Notes on Cloud computing principles. *Journal of Cloud Computing*.
437 3:21, 1-10. doi: 10.1186/s13677-014-0021-5.

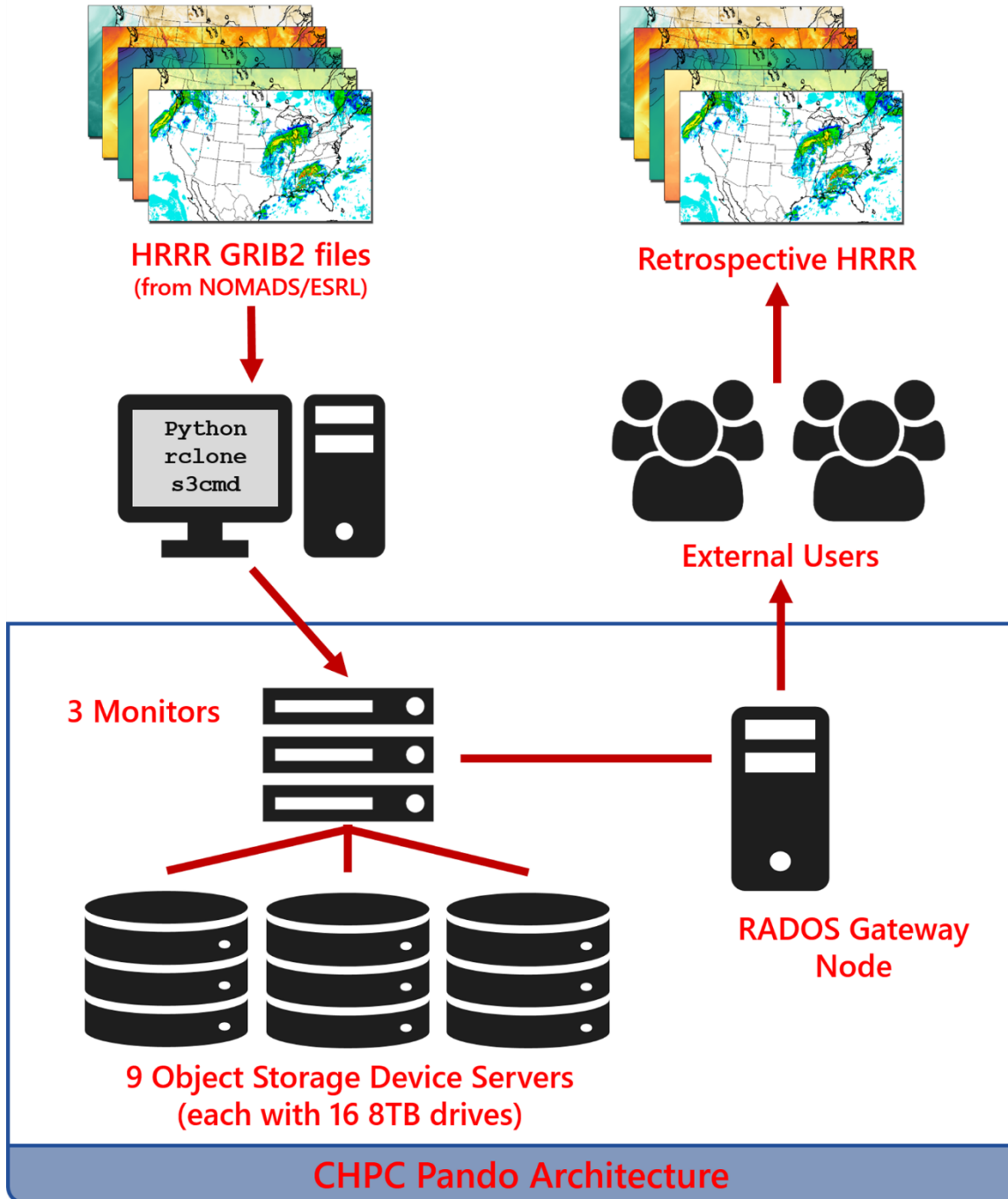
438 Schroeder, J., W. Burgett, K. Haynie, I. Sonmez, G. Skwira, A. Doggett, J. Lipe, 2005. The West
439 Texas mesonet: a technical overview. *Journal of Atmospheric and Oceanic*
440 *Technology*, 22, 211-222. doi: 10.1175/JTECH-1690.1.

441 Silver, J., C. Zender, 2017. The compression-error trade-off for large gridded data
442 sets. *Geoscientific Model Development*, 10, 413-423. doi: 10.5194/gmd-10-413-2017.

443 Storm Prediction Center, 2017. Storm Reports for 4-5 April 2017. URL:
444 http://www.spc.noaa.gov/climo/reports/170404_rpts.html and
445 http://www.spc.noaa.gov/climo/reports/170405_rpts.html.

446 SynopticLabs, 2017: MesoWest & SynopticLabs - Fostering collaboration within the weather
447 observing community. URL: <https://synopticlabs.org/>.

- 448 Taylor, K., R. Stouffer, G. Meehl. 2012: An Overview of CMIP5 and the Experiment Design.
449 *Bulletin of the American Meteorological Society*, 93, 485-498. doi: 10.1175/BAMS-D-
450 11-00094.1.
- 451 Wood, N., 2017: RCLONE- rsync for cloud storage. URL: <https://rclone.org/>.



452

453 Fig. 1. Present architecture of the Pando archive system.

HRRR Download Page

Have you Registered?

Best Practices

HRRR FAQ

Scripting Tips

Web Download Instructions

Model Type: HRRR (operational)

Variables Field: Surface (sfc, 2D fields)

Date: 4/30/2017

Get this: GRIB2 Metadata Sample

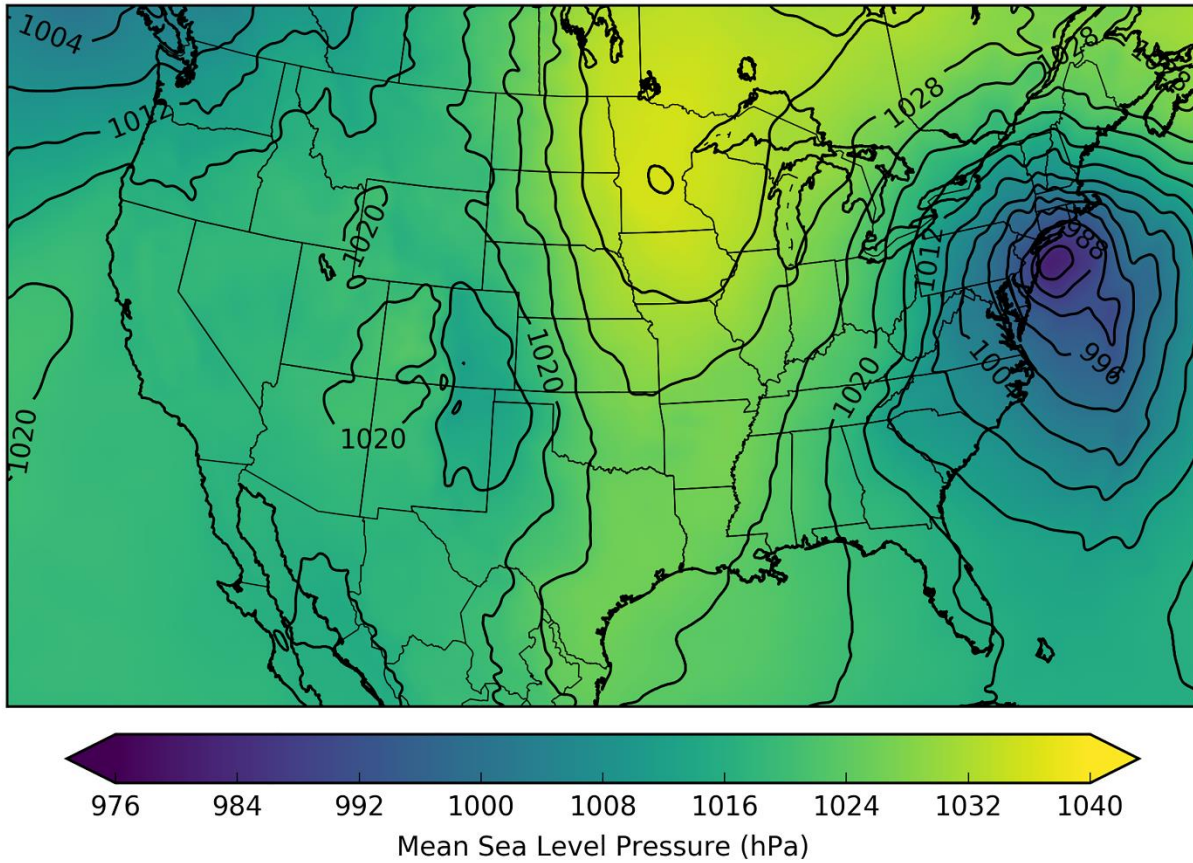
Submit

Tap to download **grib2** from 2017-04-30:

Hour 00	f00	f01	f02	f03	f04	f05	f06	f07	f08	f09	f10	f11	f12	f13	f14	f15	f16	f17	f18
Hour 01	f00	f01	f02	f03	f04	f05	f06	f07	f08	f09	f10	f11	f12	f13	f14	f15	f16	f17	f18
Hour 02	f00	f01	f02	f03	f04	f05	f06	f07	f08	f09	f10	f11	f12	f13	f14	f15	f16	f17	f18
Hour 03	f00	f01	f02	f03	f04	f05	f06	f07	f08	f09	f10	f11	f12	f13	f14	f15	f16	f17	f18
Hour 04	f00	f01	f02	f03	f04	f05	f06	f07	f08	f09	f10	f11	f12	f13	f14	f15	f16	f17	f18
Hour 05	f00	f01	f02	f03	f04	f05	f06	f07	f08	f09	f10	f11	f12	f13	f14	f15	f16	f17	f18

454

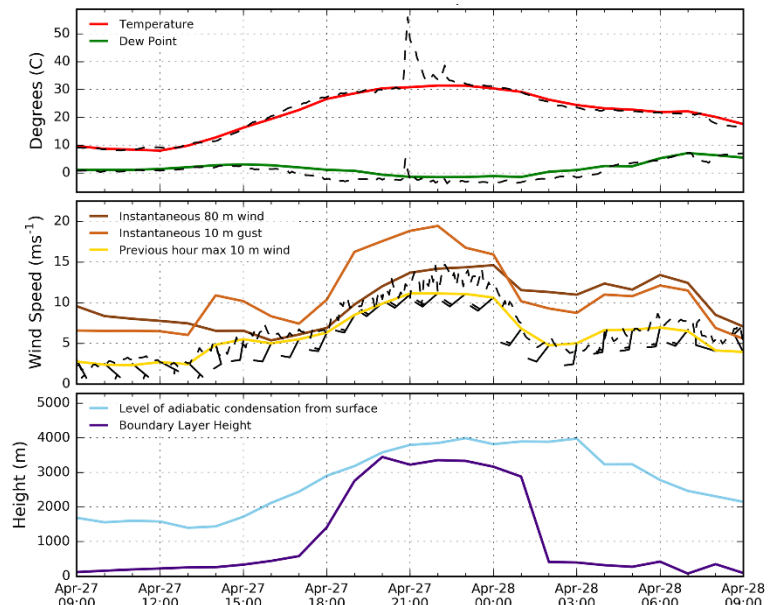
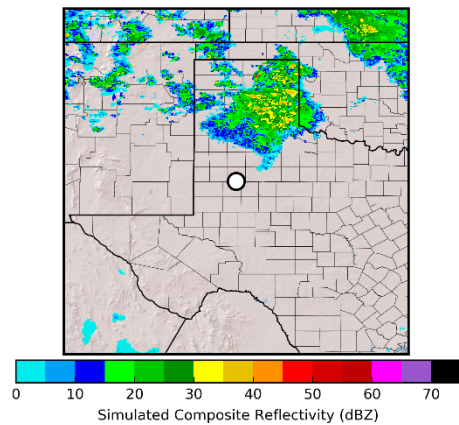
455 Fig. 2. Web interface to interactively access HRRR model output at <http://hrrr.chpc.utah.edu>.



456

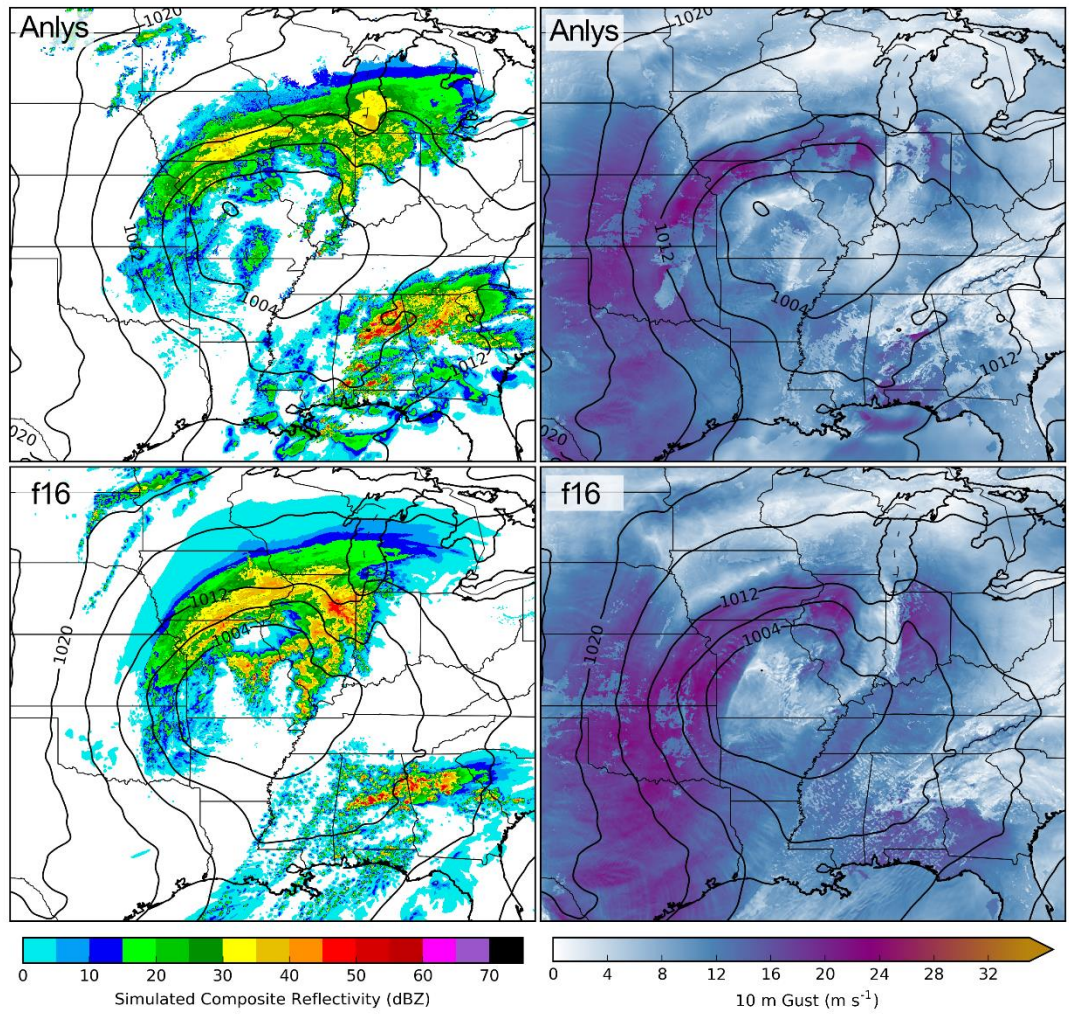
457 Fig. 3. Mean sea level pressure (hPa) from HRRR analysis at 1700 UTC 14 March 2017 during a

458 high impact New England snowstorm.

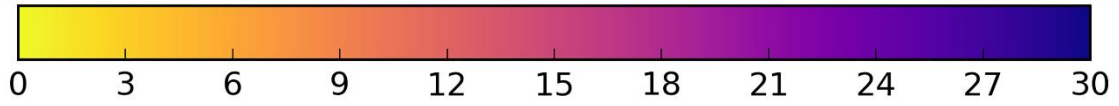
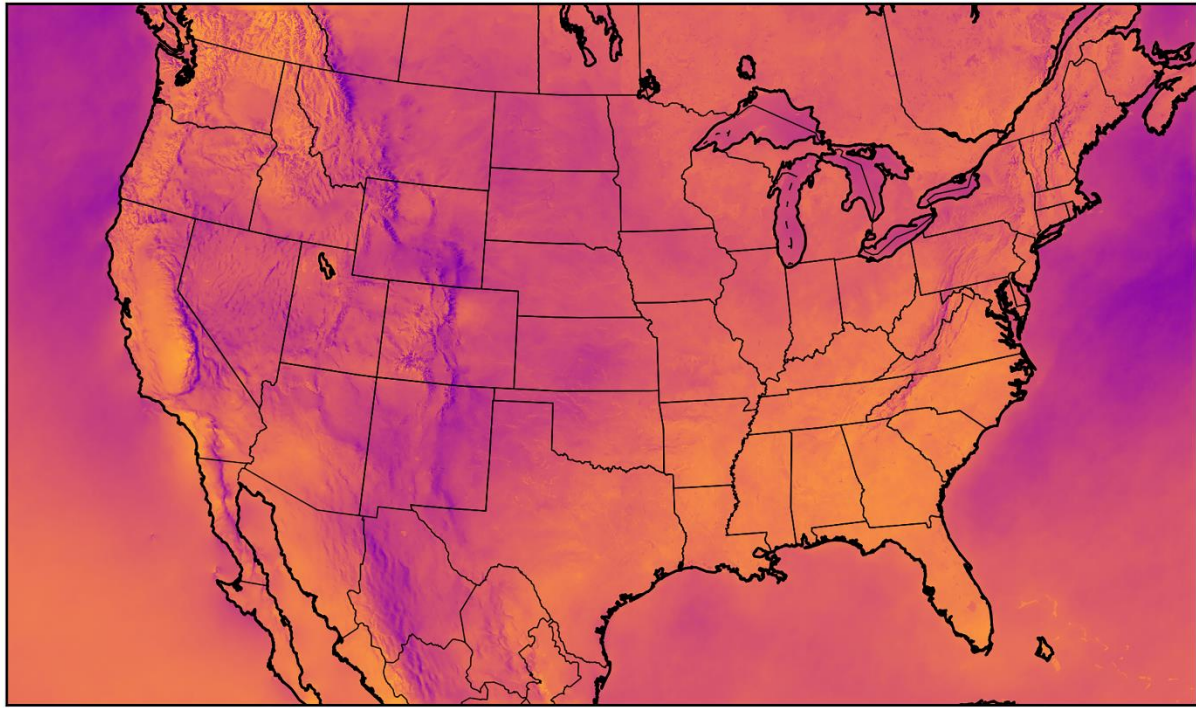


459

460 Fig. 4. (Left) HRRR simulated radar reflectivity (dBZ) at 2100 UTC 27 April 2017 at the time of
 461 a wildfire near O'Donnell, Texas (white circle). (Right) HRRR analysis of temperature ($^{\circ}\text{C}$), dew
 462 point temperature ($^{\circ}\text{C}$), 80 m wind speed (m s^{-1}), 10 m gust (m s^{-1}), 10 m maximum wind speed
 463 (m s^{-1}), 10 m wind speed and direction (half and full barbs denote 2.5 and 5 m s^{-1} , respectively
 464 and direction from which the wind blows denoted by the shaft), boundary layer height (m), and
 465 level of adiabatic condensation (m) between 0900 UTC 27 April 2017 and 900 UTC 28 April
 466 2017 near O'Donnell, Texas (white circle on the left). Observed temperature, dew point
 467 temperature, and wind speed from the O'Donnell West Texas mesonet site are shown by dashed
 468 black lines in the upper two panels.



470 Fig. 5. HRRR analyses (top panels) and HRRR 16 h forecasts (bottom panels) of mean sea level
 471 pressure (contours at intervals of 4 hPa) valid 1400 UTC 5 April 2017 with simulated composite
 472 radar reflectivity (left panels in dBZ) and 10 m gusts (right panels in m s^{-1}).



473

474 Fig. 6. 95th percentile 10 m gusts (m s^{-1}) from HRRR analyses at 2300 UTC for all days between
475 18 April 2015 and 30 March 2017.

476