

The June 2012 North American Derecho: A testbed for evaluating regional and global climate modeling systems at cloud-resolving scales

Weiran Liu¹, Paul Ullrich¹, Jianfeng Li², Colin M. Zarzycki³, Peter Martin Caldwell⁴, L. Ruby Leung⁵, and Yun Qian⁶

¹University of California Davis

²Pacific Northwest National Laboratory

³Pennsylvania State University

⁴Lawrence Livermore National Laboratory (DOE)

⁵PNNL

⁶Pacific Northwest National Laboratory (DOE)

November 24, 2022

Abstract

In this paper, we introduce a testbed for evaluating and comparing climate modeling systems at cloud resolving scales using hindcasts of the June 2012 North American derecho. The testbed is applied to two models: the regionally-refined Simple Cloud-Resolving E3SM Atmosphere Model (SCREAM) at horizontal resolutions ranging from 6.5 to 1.625 km and the Weather Research and Forecasting (WRF) model with 4 km grid spacing. We find the simulation results to be highly sensitive to the initial conditions, initialization time, and model configurations, with initial conditions from the Rapid Refresh (RAP) producing the best simulation. Significant improvement is identified in the SCREAM simulations as horizontal grid spacing is refined. While a propagation delay of approximately 2 hours is found in both models, SCREAM at 1.625 km simulates the observed bow echo structure of the derecho well and predicts strong surface gusts that exceed 30 m/s. In comparison, WRF hardly produces surface wind over 25 m/s, and the derecho wind gust in WRF is 42-46% lower than in SCREAM. Moreover, WRF has a lower bias in simulating cold clouds but overestimates the precipitation intensity. Both models well reproduce the observed outgoing longwave radiation spatial patterns (Pearson correlation > 0.88) while they simulate larger areas of composite radar reflectivity > 40 dBZ by up to 4 times and underestimate the precipitating area by $\sim 70\%$ in WRF and 47% in SCREAM compared to observations.

1 **The June 2012 North American Derecho: A testbed for**
2 **evaluating regional and global climate modeling**
3 **systems at cloud-resolving scales**

4 **W. Liu¹, P.A. Ullrich¹, J. Li², C. Zarzycki³, P. M. Caldwell⁴, L. R. Leung²,**
5 **Y. Qian²**

6 ¹Department of Land, Air, and Water Resources, University of California-Davis, Davis, CA, USA

7 ²Pacific Northwest National Laboratory, Richland, WA, USA

8 ³Department of Meteorology and Atmospheric Science, Pennsylvania State University, University Park,
9 PA, USA

10 ⁴Lawrence Livermore National Lab, Livermore, CA, USA

11 **Key Points:**

- 12 • A testbed of observational products, diagnostics, and metrics is constructed to eval-
13 uate hindcasts of the June 2012 North American derecho.
14 • Hindcast results are sensitive to initial conditions, initialization time, horizontal
15 resolutions, and convective and microphysics schemes.
16 • The Simple Cloud-Resolving E3SM Atmosphere Model at 1.625km successfully
17 reproduces severe surface gusts with wind speeds above 30 m/s.

Corresponding author: Weiran Liu, wraliu@ucdavis.edu

Abstract

In this paper, we introduce a testbed for evaluating and comparing climate modeling systems at cloud resolving scales using hindcasts of the June 2012 North American derecho. The testbed is applied to two models: the regionally-refined Simple Cloud-Resolving E3SM Atmosphere Model (SCREAM) at horizontal resolutions ranging from 6.5 to 1.625 km and the Weather Research and Forecasting (WRF) model with 4 km grid spacing. We find the simulation results to be highly sensitive to the initial conditions, initialization time, and model configurations, with initial conditions from the Rapid Refresh (RAP) producing the best simulation. Significant improvement is identified in the SCREAM simulations as horizontal grid spacing is refined. While a propagation delay of approximately 2 hours is found in both models, SCREAM at 1.625 km simulates the observed bow echo structure of the derecho well and predicts strong surface gusts that exceed 30 m/s. In comparison, WRF hardly produces surface wind over 25 m/s, and the derecho wind gust in WRF is 42-46% lower than in SCREAM. Moreover, WRF has a lower bias in simulating cold clouds but overestimates the precipitation intensity. Both models well reproduce the observed outgoing longwave radiation spatial patterns (Pearson correlation > 0.88) while they simulate larger areas of composite radar reflectivity > 40 dBZ by up to 4 times and underestimate the precipitating area by $\sim 70\%$ in WRF and 47% in SCREAM compared to observations.

Plain Language Summary

This paper describes a testbed for evaluating model performance on a particular high-impact weather event – the June 2012 North American derecho, a storm event associated with extreme winds and precipitation. The testbed is applied to evaluate the Simple Cloud-Resolving E3SM Atmosphere Model (SCREAM) and the Weather Research and Forecasting (WRF) model at resolutions that resolve cloud systems. The performance of both models is shown to be sensitive to the dataset used for model initialization. Finer grid resolution generally leads to better model performance. All simulations show a 2-hour delay in predicting the evolution of the derecho and produce more intense rainfall. SCREAM generates a more realistic convective front than WRF and produces stronger surface winds. The evaluation protocol can be used to better understand the credibility of model simulations of extreme events and guide model development.

1 Introduction

Climate modeling systems are among our best tools for understanding the climate system and future impacts of climate change (Kharin et al., 2007). In the pursuit of models of the highest quality, modeling groups cycle between developing new functionality, testing that functionality in isolation, integrating it into comprehensive modeling systems, and evaluating those combined systems. In the case of climate models, evaluation has generally focused on average behavior over large regions or long time periods (Gleckler et al., 2008; Eyring et al., 2019). However, this form of generalized analysis does not address whether climate models are able to simulate the most extreme and high-impact weather phenomena, such as extreme mesoscale convective systems (MCSs), with high fidelity. To this end, in this paper we propose one such extreme event testbed for evaluating climate modeling systems that operate at cloud resolving scales. The testbed focuses on historical simulation of a single event, in this case the June 2012 North American derecho and accompanying MCS, with the intention of providing a standard and comprehensive suite of metrics for model assessment and intercomparison.

MCSs are responsible for a variety of severe atmospheric hazards, such as flood-producing heavy rainfall events (Schumacher & Johnson, 2006; Stevenson & Schumacher, 2014; Hu et al., 2020a), lightning (Carey et al., 2005), and damaging winds (Bernardet & Cotton, 1998; Schoen & Ashley, 2011). Due to their important role in the hydrolog-

ical cycle (Hu et al., 2020b) and land-atmosphere interactions (Hu et al., 2021), there is considerable demand for evaluating their representation in one model or as part of an intercomparison across different models (Van Weverberg et al., 2013; Schumacher & Clark, 2014; A. F. Prein et al., 2020; Feng, Song, et al., 2021; Na et al., 2022). However, most previous studies focus on the long-term climatological metrics averaged over multiple years, which does not take specific events into consideration (Demaria et al., 2011; Pinto et al., 2015). A few studies evaluated the performance of climate models in the hindcast of individual extreme MCS events qualitatively but without a uniform set of metrics (Toll et al., 2015; Grunzke & Evans, 2017). Even in the short-term, for case studies of extreme MCSs that last for 1-2 days, evaluations are mainly conducted through qualitative analysis of spatial patterns (Toll et al., 2015; Grunzke & Evans, 2017) rather than quantifying the model performance using a uniform set of metrics. Previous studies (Davis et al., 2006; N. Roberts, 2008; N. M. Roberts & Lean, 2008; Mittermaier & Roberts, 2010) proposed and discussed fractions skill score (FSS) as a variation of the Brier skill score to assess a common dataset that consisted of WRF model precipitation forecasts in geometric cases. However, as indicated in Davis et al. (2006), this skill score only considered the precipitation and was highly dependent upon the threshold values and the domain sizes. While the FSS provided a measure of the spatial accuracy of precipitation forecasts, additional techniques are needed to determine behaviors of other features to gain a comprehensive understanding of the convective systems.

It is challenging to simulate individual convective storms accurately, due to the need to adequately resolve complex physical interactions between dynamical and microphysical processes over a wide range of scales (Stensrud et al., 2013; Houze Jr, 2004; Weisman & Rotunno, 2004; A. F. Prein et al., 2015; Feng et al., 2018). Previous studies have demonstrated that the performance of MCS simulations is greatly influenced by a number of factors, such as horizontal grid spacing (Tao & Chern, 2017; Squitieri & Gallus, 2020), initial conditions (ICs) (Vié et al., 2011; Brousseau et al., 2016; Weyn & Durran, 2017), model configuration (Schumacher & Clark, 2014), and choice of parameterizations (Elliott et al., 2016; Wheatley et al., 2014; Feng et al., 2018). As a result, sensitivity tests and simulation ensembles are often carried out in MCS studies to determine optimal model configurations. However, different MCS tracking algorithms and evaluation criteria are employed in these studies (Fiolleau & Roca, 2013; Haberlie & Ashley, 2019; Feng, Leung, et al., 2021), leading to possible inconsistencies in the reported results and a lack of clarity regarding the strengths and weaknesses of various models. Storm evaluation is also subject to uncertainties due to the observations or reanalyses selected as reference datasets and the selected thresholds (Kolios & Feidas, 2010; Huang et al., 2018). Therefore, a comprehensive and robust evaluation process and a uniform suite of metrics and diagnostics are much needed to streamline the process and provide greater comparability across climate modeling studies for understanding MCS features and impacts, particularly in the context of large ensembles.

The testbed proposed herein can be used to evaluate the representation of multiple storm characteristics in regional and global climate models at cloud system resolving scales. The proposed evaluation protocol is subsequently applied to compare and contrast regional and regionally-refined global climate models for a specific severe storm event, which we recommend as a standard for broader intercomparison. In this study, we limit our investigation to the Weather Research and Forecasting (WRF) model at 4km and the regionally refined model (RRM) approach using the Simple Cloud-Resolving E3SM Atmosphere Model (SCREAM) with different model configurations. Sensitivity tests address RRM grid spacing (6.5 - 1.625 km), differences between hydrostatic and nonhydrostatic dynamical cores, low-resolution and high-resolution model configurations, initialization time, and data source for the ICs.

This paper is organized as follows: section 2 presents the proposed testbed, including a brief introduction of the severe weather event we selected, observations employed,

121 and metrics selected; section 3 describes the SCREAM-RRM and WRF models in de-
122 tail; section 4 evaluates the simulation results of both models; finally, section 5 provides
123 a summary of our findings and conclusions.

124 **2 The June 2012 North American Derecho Testbed**

125 In this section we describe the proposed testbed, based on a 24-hour hindcast of
126 the June 2012 North American derecho, and designed for evaluation and intercompar-
127 ison of climate modeling systems at cloud resolving scales. The testbed consists of a sim-
128 ulation protocol, a set of observational products, and a comprehensive set of diagnos-
129 tics and statistical metrics that leverage those observations. Section 2.1 provides a me-
130 teorological overview of the derecho and previous relevant studies. Section 2.2 describes
131 the selected observational datasets. Section 2.3 presents four essential storm features that
132 are examined in this framework, and section 2.4 explains the calculations of the metrics.

133 **2.1 Meteorology**

134 Johns and Hirt (1987) categorized derechos as meteorological events with severe
135 wind gusts and precipitation lasting for several hours, in conjunction with a linear MCS.
136 An extensive study (Corfidi et al., 2016) more recently defined a derecho as an event with:
137 1) convectively induced wind damage and/or gusts of > 25.7 m/s over an area with a
138 major axis of 400 km, 2) geographically-consistent reports, and 3) presence of 3 or more
139 reports of gusts > 33.4 m/s within the affected area. Among all historical derechos in
140 North America, the June 2012 North American derecho (or June 2012 Mid-Atlantic and
141 Midwest derecho) is one of the most infamous – a progressive derecho event that became
142 one of the most destructive and fastest-moving derechos in US history.

143 The June 2012 North American derecho was characterized by an intense bow-echo
144 MCS causing widespread severe wind damage across the upper Midwest and the Ohio
145 River valley, as well as the mid-Atlantic states, during the afternoon and evening of 29
146 June and early morning of 30 June in 2012 (Shourd, 2017; Shourd & Kaplan, 2021). This
147 particular event was selected because of the significant socioeconomically-hazardous im-
148 pact and the high forecast difficulty. At initiation, a relatively small cluster of storm cells
149 began to form as embryonic convection in eastern Iowa around 14:00 UTC on 29 June.
150 Around 16:00 UTC, the small storm cluster began rapidly forming a well-defined MCS
151 before crossing through Chicago, Illinois. Afterward, the MCS expanded into an asym-
152 metric bow echo over Indiana as it accelerated southeastward at about 25 m/s slightly
153 to the north of the frontal boundary. The MCS intensified further as it crossed Indiana
154 and Ohio, transforming into a derecho MCS. The MCS continued along its destructive
155 path until reaching the Atlantic coast of Virginia and Maryland about 06:00 UTC on
156 30 June. As estimated by the Storm Prediction Center (SPC), a damaging wind swath
157 of about 1000 km in length resulted from this event, with over 800 wind damage reports
158 during the 10-hour lifetime. Severe wind gust reports ranging between 25–33 m/s were
159 widespread with peak gusts in excess of 40 m/s reported over eastern Indiana and west-
160 ern Ohio.

161 As indicated in Johns and Hirt (1987), progressive derechos are frequently challeng-
162 ing for operational meteorologists to forecast due to their weakly forced nature. The June
163 2012 North American derecho was underforecasted days and hours ahead of time, as well
164 as throughout much of the duration of the storm. Most numerical weather prediction
165 models showed no indication that any convective cells would develop, illustrating the fore-
166 cast difficulty (Halverson, 2014; Guastini & Bosart, 2016; Schumacher & Rasmussen, 2020).

167 This forecast difficulty serves as the motivation for the following studies. Fierro et
168 al. (2014) evaluated the short-term forecast (≤ 6 h) of the derecho event from the re-
169 gional WRF model at 3 km resolution to compare two distinct data assimilation tech-

170 niques. Shourd and Kaplan (2021) simulated the derecho using the WRF model in a nested
 171 domain with the inner domain at 2 km resolution and reproduced the super derecho. How-
 172 ever, no quantitative evaluation metrics were used in these two analyses, resulting in no
 173 clear conclusions drawn as to the quality of the reproduction, especially when compared
 174 with other studies. Shepherd et al. (2021) performed an 11-member ensemble of convection-
 175 permitting regional simulations using WRF and tested the sensitivities to model con-
 176 figuration including microphysics parameterizations, lateral boundary conditions, start
 177 dates, and use of nudging. All 11 members had difficulty capturing the realistic evolu-
 178 tion of the derecho, exhibiting a time delay (ranging from 2 - 8 hours) in simulating the
 179 derecho intensification and passage.

180 Previous studies that focused on simulating the June 2012 North American dere-
 181 cho (Fierro et al., 2014; Shourd, 2017; Shourd & Kaplan, 2021; Schumacher & Rasmussen,
 182 2020) have emphasized the analysis and evaluation of composite radar reflectivity. Nev-
 183 ertheless, wind gusts are an integral component of the definition of derecho. In order to
 184 provide a more thorough evaluation of the event, our study has an additional empha-
 185 sis on evaluating the wind speed, along with precipitation and composite radar reflec-
 186 tivity.

187 2.2 Observations

188 It is well known that precipitation products diverge considerably across regions,
 189 even in the regional means at daily to seasonal timescales, and particularly across in-situ,
 190 reanalysis and satellite products (Miao et al., 2015; Beck et al., 2017, 2019; Sadeghi et
 191 al., 2021). Our testbed requires a detailed comparison of hourly precipitation pattern
 192 and magnitude at fine horizontal resolution, where the differences between these prod-
 193 ucts are particularly large. Therefore, three high-resolution gauge-based precipitation
 194 datasets are used to evaluate the simulated precipitation:

- 195 1. The National Centers for Environmental Prediction (NCEP) 4km Gridded Stage
 196 IV Data (Lin & Mitchell, 2005; Du, 2011), which is a merged ground-based and
 197 radar-derived hourly rainfall accumulation dataset from 140 radars and ~ 5500
 198 gauges over the continental United State (CONUS). The NCEP Stage IV dataset
 199 provides highly accurate precipitation estimates and is, therefore, widely used as
 200 a reference for the evaluation of precipitation (Hong et al., 2004; AghaKouchak
 201 et al., 2011, 2012; Nelson et al., 2016; X. Zhang et al., 2018).
- 202 2. NASA Integrated Multi-satellite Retrievals for Global Precipitation Measurement
 203 (IMERG) V06B final run (Huffman et al., 2015), which intercalibrates, merges,
 204 and interpolates all estimates of the Global Precipitation Measurement (GPM)
 205 constellation, infrared (IR) estimates, gauge observations, and other potential sen-
 206 sors' data with a $0.1^\circ \times 0.1^\circ$ spatial resolution and 30 minute temporal resolution.
- 207 3. NOAA Climate Prediction Center Morphing technique (CMORPH) bias-corrected
 208 V1.0 (Joyce et al., 2004; Xie et al., 2017, 2019) – this 8 km resolution dataset pro-
 209 duces 30 minute estimates of rainfall derived from passive microwave observations
 210 and extrapolates them backwards and forwards in time via spatial propagation
 211 information obtained from geostationary IR satellite data.

212 Following previous efforts (Beck et al., 2019; Feng et al., 2018), we use the NCEP Stage
 213 IV dataset as our primary precipitation reference, but also provide supplementary re-
 214 sults from IMERG and CMORPH. While the NCEP Stage IV precipitation dataset is
 215 of high quality, it is available only over the US. IMERG and CMORPH are included to
 216 generalize our framework for testbed cases worldwide, considering their broader cover-
 217 age. An intercomparison of different precipitation datasets is out of the scope of this pa-
 218 per.

219 Outgoing longwave radiation (OLR) is evaluated using the brightness temperature
 220 (Tb) from the NCEP half-hourly 4 km IR V1 dataset (Janowiak et al., 2017), which con-
 221 tains globally-merged geostationary satellites with parallax correction and viewing an-
 222 gles correction. Tb is converted to OLR following the empirical formulation provided by
 223 Yang and Slingo (2001).

224 For observations of radar reflectivity, we use the hourly three-dimensional high-resolution
 225 Next-Generation Radar (NEXRAD) (Bowman & Homeyer, 2017), which covers most of
 226 the contiguous U.S merged from 125 National Weather Service WSR-88D weather radars.
 227 The raw spatial resolution of NEXRAD is $0.02^\circ \times 0.02^\circ$ and a vertical resolution of 1 km.
 228 Composite reflectivity (cREF) is calculated as the maximum reflectivity for each column
 229 and time step in both NEXRAD and the simulations.

230 For wind speed evaluation, we use station records from the National Weather Ser-
 231 vice Automated Surface Observation System (ASOS) (Nadolski, 1998). There are 90 ASOS
 232 stations in the analysis domain ($76^\circ\text{-}88^\circ\text{W}$, $36.5^\circ\text{-}42^\circ\text{N}$), shown as black circles in Figure
 233 1c. The temporal frequency of the ASOS record is 5 minutes, although several records
 234 are missing. Two wind-related parameters from the ASOS are used:

- 235 1. **Wind speeds:** ASOS stations measure wind direction and speed once every sec-
 236 ond using meteorological equipment at a height of 10 meters. Five-second wind
 237 direction and wind speed averages are computed from the 1-second measurements.
 238 These 5-second averages are rounded to the nearest knot and retained for 2 min-
 239 utes. The resolution of the wind speed is 1 knot and converted from knots to m/s
 240 in all analyses of this study.
- 241 2. **Gust wind speeds:** The gust wind speeds represent the maximum five-second
 242 wind speed measured in each five-minute period when gust criteria are met (Nadolski,
 243 1998). Gusts are rounded up to the nearest whole knot and converted from knots
 244 to m/s. Gust wind speed is not a standard parameter and only reported when:
 - 245 (a) Gust wind speed is at least 3 knots (1.54 m/s) above the current running 2-minute
 246 mean wind speed.
 - 247 (b) Gust wind speed exceeds the minimum five-second average in the last 10 min-
 248 utes by at least 10 knots (5.14 m/s).
 - 249 (c) The current 2-minute average wind speed is at least 3 knots (1.54 m/s).

250 2.3 Storm Characteristics

251 To provide a near comprehensive evaluation of the relevant meteorological char-
 252 acteristics of the derecho, the proposed testbed focuses on four essential parameters: pre-
 253 cipitation, cREF, OLR, and wind speed. We define three features based on the commonly
 254 used thresholds in the previous MCS analyses to locate and track the derecho: the cold
 255 cloud shield, the precipitation feature, and the cREF feature. The cold cloud shield is
 256 defined as the contiguous area with Tb lower than 241 K (Maddox, 1980; Feng et al.,
 257 2018; Feng, Leung, et al., 2021). Following the empirical formulation provided by Yang
 258 and Slingo (2001), this Tb threshold is instead applied to the OLR (which is output di-
 259 rectly from the models), using a threshold of 163.44 W/m^2 . The precipitation feature
 260 is defined as the contiguous area with precipitation rate higher than 1 mm/hour (Peters
 261 et al., 2009; Yuan & Houze, 2010; Feng et al., 2018). The cREF feature is defined as a
 262 continuous area with composite radar reflectivity greater than 40 dBZ (Dye et al., 1989;
 263 Zipser & Lutz, 1994; Haberlie & Ashley, 2019).

264 The latitude and longitude of the midpoint of a certain feature is calculated as the
 265 mean of the maximum and minimum of the latitude and longitude of the object. While
 266 the centroid of the feature polygon could have been similarly employed (Pinto et al., 2015;
 267 Davis et al., 2006), we observed similar results to those obtained via the simple midpoint.

268 Therefore, this study only uses the midpoint instead of the centroid because of the sim-
 269 plicity of computation. The features are further isolated and tracked using TempestEx-
 270 tremes (Ullrich & Zarzycki, 2017; Ullrich et al., 2021), as shown in the appendix. The
 271 area of an isolated feature is calculated as the sum of areas of grid points that are de-
 272 tected in the TempestExtremes.

273 Note that the definitions of MCSs and thresholds are diverse in the past studies
 274 (Schumacher & Johnson, 2005; Yuan & Houze, 2010). While we choose the most widely
 275 used thresholds, the involvement of thresholds and tracking algorithms would still in-
 276 duce a certain degree of uncertainty, as mentioned in section 1. Therefore, we will use
 277 metrics without incorporating the storm detection and tracking if possible besides the
 278 features described above.

279 2.4 Evaluation Metrics

280 Quantitative evaluation of the SCREAM and WRF experiments is performed via
 281 a variety of statistical techniques over the analysis region shown in Figure 1c. To facil-
 282 itate comparison, the simulated variables are interpolated onto the coarse observation
 283 grid (i.e., 0.05° for OLR and cREF and 0.1° for precipitation). Here we use precipita-
 284 tion as an example, but similar calculations are applied to other variables (OLR and cREF).

285 The model bias is measured by the mean error (ME),

$$ME = \frac{1}{N} \sum_{i=1}^N (p_i - o_i), \quad (1)$$

286 where N is the total number of verification grid point, and p and o are the simulated and
 287 observed values, respectively. Mean absolute error (MAE), is calculated as

$$MAE = \frac{1}{N} \sum_{i=1}^N |p_i - o_i|. \quad (2)$$

288 The root-mean-square of the error (RMSE) (Anthes, 1983) is defined as

$$RMSE = \left[\frac{1}{N} \sum_{i=1}^N (p_i - o_i)^2 \right]^{1/2}. \quad (3)$$

289 The pattern correlation between simulations and observations are represented by the Pear-
 290 son product-moment correlation coefficient, calculated as

$$r = \frac{\sum_{i=1}^N (p_i - \bar{p})(o_i - \bar{o})}{\sqrt{\sum_{i=1}^N (p_i - \bar{p})^2} \sqrt{\sum_{i=1}^N (o_i - \bar{o})^2}}. \quad (4)$$

291 We choose the centered form, which measures the similarity of the pattern after remov-
 292 ing the regional mean (Santer et al., 1993), because it provides additional information
 293 independent of the mean bias. We also use the Spearman rank correlation coefficient as
 294 a robust and resistant alternatives to the Pearson product-moment correlation coefficient.
 295 The Spearman correlation is simply the Pearson correlation coefficient computed using
 296 the ranks of the data,

$$r_s = 1 - \frac{6 \sum_{i=1}^N D_i^2}{N(N^2 - 1)}, \quad (5)$$

297 where D_i is the difference in ranks between the i th pair of values.

298 It is important to stress that a particular simulation can exhibit a bias close to zero,
 299 along with poor correlation (e.g., the regionally-averaged precipitation rate is similar to
 300 the reference dataset but the precipitation patterns are distorted), or a high correlation,

301 but with a high bias (e.g., a consistent spatial distribution of precipitation but with in-
 302 tensified rainfall rate relative to that of the reference dataset). As such, the conclusions
 303 derived from single metrics could be misleading, suggesting a need to incorporate mul-
 304 tiple measures in such an analysis.

305 Our evaluation metrics also include two scores normally used in the assessment of
 306 accuracy of weather prediction. The first is the bias score (BS), which indicates whether
 307 the model over or under predicts the fractional areal coverage of precipitation for a cer-
 308 tain threshold. On the other hand, the threat score (TS), ranging from 0 (worst) to 1
 309 (best), is used to measure the skill of predicting the area of precipitation exceeding a cer-
 310 tain intensity threshold. The BS and TS are defined as

$$BS = \frac{P}{O}, \quad (6)$$

311 and

$$TS = \frac{H}{P + O - H}, \quad (7)$$

312 where P and O are the number of grid points with values higher/lower (i.e., higher for
 313 precipitation and cREF; lower for OLR) than the threshold in the simulation and ref-
 314 erence dataset, respectively. H is the number of grid points higher/lower than the thresh-
 315 old in both the simulation and the observation.

316 3 Models

317 3.1 The SCREAM Regional Refined Model

318 A series of RRM simulations are conducted using SCREAM (Caldwell et al., 2021;
 319 Liu et al., 2022), configured with a high-resolution (HR) grid located in the northeast-
 320 ern US, a low-resolution (LR) grid covering the remaining globe, and a transition area
 321 between them (Figure 1b). Figure 1a shows the SCREAM RRM grid in the global or-
 322 thographic projection. The grid is based on the unstructured cubed-sphere finite-element
 323 grid with 4 Gauss-Lobatto-Legendre (GLL) nodes per element's edge (`np4`). Our LR grid
 324 uses 128×128 spectral elements on each face, denoted `ne128`, corresponding to a hor-
 325 izontal grid spacing of 0.25° (~ 28 km). Using the offline software tool SquadGen (Ullrich,
 326 2014), three RRM grids were constructed using the same low base resolution (`ne128`) and
 327 various high resolutions: `ne512` (6.5 km), `ne1024` (3.25 km), and `ne2048` (1.625 km). The
 328 HR portion of the grid is large enough to comprise the region where the derecho initi-
 329 ated, as well as its propagation path. While the derecho eventually migrated to the At-
 330 lantic in its decay phase, our analysis only focuses on processes on land where the dam-
 331 age occurred and, therefore, does not cover broad oceanic area in the HR portion.

332 The RRM approach has been validated in other models over many regions of in-
 333 terest (Zarzycki & Jablonowski, 2014; Sakaguchi et al., 2015, 2016; Rhoades et al., 2018;
 334 Wu et al., 2017; Xu et al., 2018) and demonstrated to be effective for regional climate
 335 studies at a reduced computational cost compared to uniform GCMs. For example, Zarzycki
 336 and Jablonowski (2014, 2015) demonstrated improved skill in simulating tropical cyclones
 337 in the Community Atmosphere Model with a refined mesh (0.25°) over the North At-
 338 lantic at multidecadal timescale. Huang and Ullrich (2017) reproduced the geographic
 339 patterns of 26-year historical precipitation climatology over the western US with the variable-
 340 resolution Community Earth System Model with a fine grid resolution of 0.25° . Two
 341 studies (Sakaguchi et al., 2015; Tang et al., 2019) demonstrated that RRM reproduced
 342 the seasonal precipitation of the high-resolution model over the CONUS.

343 However, previous RRM studies were performed with the highest horizontal res-
 344 olution of around 0.25° and seasonal or longer timescales. Our study adopt the RRM ap-
 345 proach in SCREAM with finer horizontal resolutions from 6.5 - 1.625 km and a timescale
 346 of 1-day. Since no optimal grid spacing has been identified for MCS simulations (Weisman

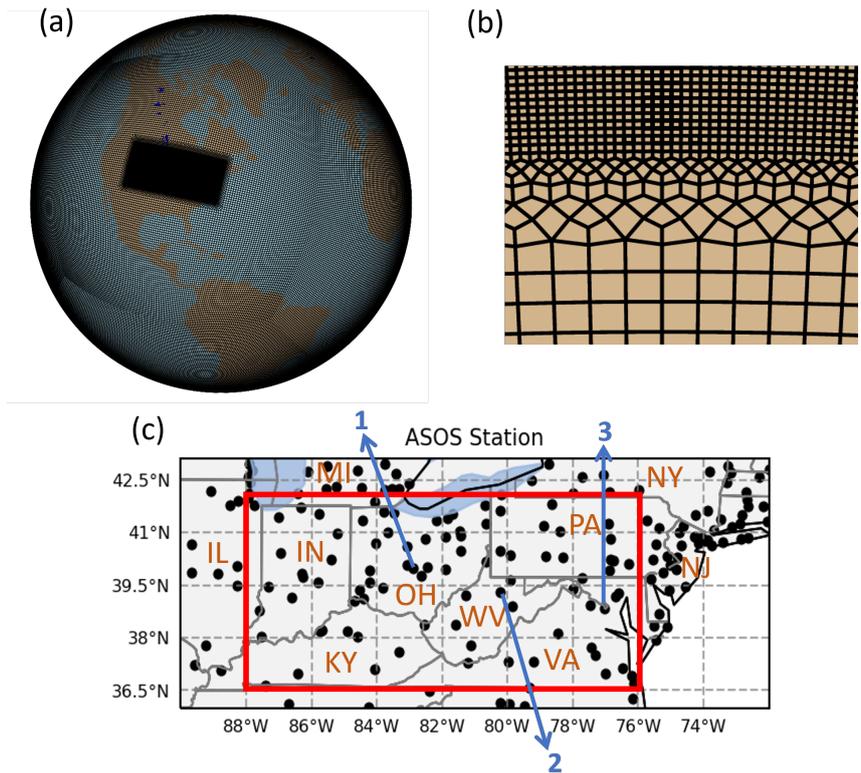


Figure 1. (a) The SCREAM RRM grid, shown using a global orthographic projection. (b) The transition region in the RRM grid from LR to HR resolution. (c) Locations of ASOS stations in black circles. The red box shows the analysis region (76°-88°W, 36.5°-42°N).

347 et al., 1997; Squitieri & Gallus, 2020), we investigate a variety of grid spacings between
 348 6.5 km and 1.625 km to examine the impact of horizontal resolution on the derecho simu-
 349 lation in SCREAM.

350 Prescribed SST and sea ice extent are used for all SCREAM simulations. The land
 351 initial file is generated from a 12-month spinup land simulation prior to the initial date.
 352 Native output is saved every 15 minutes and later remapped using the TempestRemap
 353 software suite (Ullrich & Taylor, 2015; Ullrich et al., 2016) before calculating derived vari-
 354 ables or performing analyses.

355 A summary of the SCREAM RRM simulations conducted in this study is provided
 356 in Table 1, including horizontal resolutions, ICs, initialization time, LR/HR configura-
 357 tions, and dynamical cores. Specifically, simulations SCREAM_6.5km, SCREAM_3.25km,
 358 and SCREAM_1.625km are designed to examine the sensitivity of model performance
 359 to grid spacing. Simulations SCREAM_ERA5 and SCREAM_ERAI serve as sensitivity
 360 tests of the model to the IC source. All simulations are initialized at 12:00 UTC 29 June
 361 2012 and end at 12:00 UTC 30 June 2012, except for SCREAM_06Z initialized at 06:00
 362 UTC 29 June 2012. The SCREAM HR configuration, where deep convection is turned
 363 off and includes a 128 layer vertical grid with a model top at 40 km (2.25 hPa), is em-
 364 ployed in most simulations. SCREAM_LR uses the LR configuration; the model is run
 365 with 72 vertical levels with a top at 60 km, and the Zhang-McFarlane deep convection
 366 scheme (G. J. Zhang & McFarlane, 1995) is applied. These two configurations follow the
 367 vertical levels, model tops, and application of the deep convection scheme as described

Table 1. A summary of the SCREAM RRM simulations conducted and compared in this study.

Simulation Abbreviation	Fine Resolution	IC	Initialization Time (UTC)	LR/HR Configuration	Dynamical Core
SCREAM.6.5km	ne512 (6.5km)	ERA5+RAP	12:00 29 June 2012	HR	NH
SCREAM.3.25km	ne1024 (3.25km)	ERA5+RAP	12:00 29 June 2012	HR	NH
SCREAM.1.625km	ne2048 (1.625km)	ERA5+RAP	12:00 29 June 2012	HR	NH
SCREAM.ERA5	ne1024 (3.25km)	ERA5	12:00 29 June 2012	HR	NH
SCREAM.ERA1	ne1024 (3.25km)	ERA1	12:00 29 June 2012	HR	NH
SCREAM.06Z	ne1024 (3.25km)	ERA5+RAP	06:00 29 June 2012	HR	NH
SCREAM.LR	ne1024 (3.25km)	ERA5+RAP	12:00 29 June 2012	LR	NH
SCREAM.H	ne1024 (3.25km)	ERA5+RAP	12:00 29 June 2012	HR	H

Table 2. Timesteps of the SCREAM RRM simulations.

Simulation Name(s)	Fine Resolution	Dynamics Timestep (s)	Physics Timestep (s)
SCREAM.6.5km	ne512 (6.5km)	$16 \frac{2}{3}$	300
SCREAM.3.25km, SCREAM.ERA5, SCREAM.ERA1, SCREAM.06Z SCREAM.LR, SCREAM.H	ne1024 (3.25km)	$8 \frac{1}{3}$	100
SCREAM.1.625km	ne2048 (1.625km)	$4 \frac{1}{6}$	50

368 in Caldwell et al. (2021) and Caldwell et al. (2019), respectively. The dynamical equa-
369 tions are solved using the High Order Method Modeling Environment (HOMME) (J. Den-
370 nis et al., 2005; J. M. Dennis et al., 2012; Evans et al., 2013). The simulations mostly
371 use the HOMME nonhydrostatic (NH) spectral element dynamical core (Taylor et al.,
372 2020; Bertagna et al., 2020; Liu et al., 2022) with one addition sensitivity test (SCREAM.H)
373 utilizing the HOMME hydrostatic (H) dynamical core (Golaz et al., 2019; Caldwell et
374 al., 2019; S. Zhang et al., 2020). Both the dynamics and physics timesteps are scaled across
375 different RRM grids, controlled by the fine resolution, as shown in Table 2. Because of
376 the horizontal resolution differences among three RRM grids, the topography is repre-
377 sented differently in these configurations.

378 The IC files are derived from three datasets: the Rapid Refresh (RAP) (Benjamin
379 et al., 2016), the fifth generation of atmospheric reanalysis (ERA5) (Hersbach et al., 2018)
380 and ERA-Interim (ERA1) (Dee et al., 2011), with details summarized in Table 3. A hind-
381 cast initialization suite (Betacast, Zarzycki and Jablonowski (2015)) is used to generate
382 the IC files for the model from the above datasets. Since the RRM is a global model, ERA5
383 and ERA1 data are directly mapped from the reanalysis grids to the model grid. The
384 RAP analysis only covers North America, so for the simulations initialized using RAP,
385 a two step approach is applied where a global ‘base’ IC is first generated using ERA5
386 and then the RAP analysis is used to overwrite the model state fields over the valid RAP
387 region, displayed as ERA5 + RAP in Table 1. To eliminate noise associated with map-
388 ping the analyses across different grids, a hydrostatic correction is applied at each grid
389 point to correct the hydrostatic surface pressure field between the analysis and model
390 orographies, following the method described in Trenberth et al. (1993). Finally, all prog-
391 nostic state variables in the vertical column are then reinterpolated based on the adjusted
392 surface pressure since SCREAM uses a terrain-following coordinate.

393 3.2 WRF Model

394 WRF v4.3.3 (Skamarock et al., 2019) at 4 km is employed for intercomparison with
395 the SCREAM RRM simulations. The WRF domain extends from 30.69°N to 48.04°N
396 and from 102.78°W to 62.01°W. Four WRF simulations are run using different setups

Table 3. A summary of datasets used to generate IC files.

Dataset Name	Coverage	Temporal Resolution	Grid Spacing	Reference
RAP	North America	Hourly	13 km	Benjamin et al. (2016)
ERA5	Global	Hourly	0.25°	Hersbach et al. (2018)
ERA-Interim	Global	6-hourly	0.75°	Dee et al. (2011)

Table 4. A summary of the WRF simulations.

Simulation Abbreviation	IC	Number of Vertical Levels	Microphysical Scheme
WRF_RAP	ERA5+RAP	45	Thompson
WRF_NARR	NARR	45	Thompson
WRF_HR	ERA5+RAP	72	Thompson
WRF_HR_P3	ERA5+RAP	72	P3

(Table 4) with the same simulation period, initialized on 12:00 UTC 29 June 2012, and output frequency as the SCREAM RRM simulations. The time step for integration is 10 seconds in the WRF simulations. The baseline simulation (WRF_RAP in Table 4), has 45 vertical layers with a thickness of ~ 50 m for the lowest layer and a top at 100 hPa. Physics schemes used in WRF_RAP include the Thompson microphysics scheme (Thompson et al., 2008), the Rapid Radiative Transfer Model for General Circulation Models (RRTMG) shortwave and longwave radiation schemes (Iacono et al., 2008), the Mellor-Yamada-Janjic (MYJ) planetary boundary layer scheme (Janjić, 1994), the Eta similarity surface layer scheme, the Noah Land Surface Model (Chen & Dudhia, 2001), and the Building Energy Model coupled with the Building Environment Parameterization (BEP + BEM) for urban physics (Salamanca et al., 2010). Initial and lateral boundary conditions are from ERA5 and RAP, where ERA5 provides soil conditions while RAP provides atmospheric and land surface conditions.

Besides the baseline simulation, three sensitivity tests are performed (WRF_NARR, WRF_HR, and WRF_HR_P3; Table 4) to examine the impacts of different initial and boundary conditions, vertical resolutions, and microphysical schemes. The configuration of WRF_NARR is the same as WRF_RAP, except using initial and boundary conditions from the NCEP North American Regional Reanalysis (NARR) product (Mesinger et al., 2006). Compared to WRF_RAP, WRF_HR has 72 vertical layers with a vertical resolution of ~ 20 m near the surface. The difference between WRF_HR and WRF_HR_P3 is that WRF_HR_P3 uses the Predicted Particle Property (P3) microphysics scheme with 3-moment ice (Morrison & Milbrandt, 2015).

4 Evaluation

Our discussion begins with a snapshot of the mature stage of the derecho at 00:00 UTC 30 June 2012 in section 4.1. The temporal evolution of the derecho is investigated in section 4.2 followed by section 4.3, which presents the metrics to quantify the fidelity of the models. We also display how to interpret the quantified metrics to understand the derecho characteristics in section 4.3. Section 4.4 provides additional discussion about the simulated 10-m wind speed.

4.1 00:00 UTC snapshot

Figures 2-3 show the instantaneous simulation outputs of OLR and cREF at 00:00 UTC 30 June 2012 in eight SCREAM RRM simulations (Table 1), two WRF simulations

(WRF_RAP and WRF_NARR, Table 4), and observations at 0.05° resolution. The observation panel is marked with red title in all figures. Unlike Figures 2-3, which show instantaneous outputs, Figure 4 shows the precipitation amount in the simulations and reference datasets accumulated from 00:00 to 01:00 UTC (since the NCEP Stage IV precipitation dataset is accumulated hourly). Some spatial displacement is clear between the cREF and precipitation due to the propagating nature of the derecho. We will discuss the two WRF sensitivity tests (WRF_HR and WRF_HR_P3; see Table 4) in section 4.4 and, therefore, not display their results in this section.

Figures 4a-c clearly show that the precipitation patterns in different products are divergent. While all three products include gauge corrections, IMERG shows significantly higher rainfall rates than NCEP Stage IV dataset and CMORPH, especially in the southern part of the derecho near (80°W, 39°N). However, the NCEP Stage IV rainfall is the most widely used reference dataset (Beck et al., 2019; Feng et al., 2018) and has the best agreement with the ASOS station records among all three products (not shown). Accordingly, NCEP Stage IV dataset is used as the primary precipitation reference dataset in the following analysis.

Based on the comparison of SCREAM RRM simulations at three different grid spacings (SCREAM_6.5km, SCREAM_3.25km, and SCREAM_1.625km) with observations, it is clear that simulations at higher horizontal resolutions appear to better represent the derecho. Specifically, the simulated derecho at 6.5 km resolution is underdeveloped, producing the smallest cold cloud shield (by $\sim 65\%$) and most compact cREF/precipitation feature. While the derechos at all three horizontal resolutions are all located upstream (northwest side) of the observed feature, the discrepancy between the simulation and the observation decreases as the resolution becomes finer. The bow-shape echo and the axis angle of the convective core are more qualitatively similar to observations in the 1.625 km simulation.

The simulation performance exhibits substantial sensitivity to the IC sources (SCREAM_3.25km, SCREAM_ERA5, and SCREAM_ERAI). This dependency has been pointed out in past research examining convection simulations, as summarized in section 1. Consistent with the results in Shepherd et al. (2021), despite the higher resolution and larger data assimilation volume of ERA5, the simulation initialized with ERA5 does not show significantly better performance than the one with ERAI. However, simulation initialization with RAP shows significantly improved performance compared to both ERA5 and ERAI. Notably, Figurski et al. (2017) also showed that simulations using ERA5 produce scattered reflectivity fields that are very different from those observed. WRF simulations (WRF_RAP and WRF_NARR) are also sensitive to the IC source, with better performance apparent in WRF_RAP than WRF_NARR. Interestingly, despite the good performance of SCREAM_3.25km, the simulation initialized 6 hours earlier at 06:00 UTC (SCREAM_06UTC) shows little precipitation and cREF, along with a weaker cold cloud shield, indicating the high sensitivity to the IC source even when applying the same dataset at different initialization times (Figurski et al., 2017).

The SCREAM simulation with LR model configuration (SCREAM_LR) is not able to reproduce the derecho successfully: namely, convective clouds do not form when the deep convective scheme is active. This is perhaps unsurprising, as previous studies have indicated better simulation of individual convective events in a convection-permitting model without a convective parameterization scheme than those in GCMs and RCMs (A. F. Prein et al., 2015; Fosser et al., 2015). This suggests a significant benefit comes from resolving convection explicitly, as the use of a convective parameterization scheme leads to common errors such as misrepresentation of the diurnal cycle of convective precipitation (Dai et al., 1999; Brockhaus et al., 2008) and the underestimation of hourly precipitation intensity (A. Prein et al., 2013; Fosser et al., 2015; Ban et al., 2014; Gao et al., 2017).

481 Although Liu et al. (2022) demonstrated the discrepancy between nonhydrostatic
 482 and hydrostatic simulations is significant over certain hotspots in the seasonal simula-
 483 tion ensembles, the simulation with hydrostatic dynamical core (SCREAM_H) is not sig-
 484 nificantly different from its nonhydrostatic counterpart in the snapshots of this short-
 485 term hindcast, producing remarkably similar result to SCREAM_3.25km. This suggests
 486 that even a hydrostatic dynamical core can simulate MCSs with comparable fidelity to
 487 a nonhydrostatic dynamical core, even far into the classical nonhydrostatic regime, po-
 488 tentially because the physics parameterizations dominate the model behaviors.

489 The detailed structure of the derecho in observations is particularly well simulated
 490 in the SCREAM_1.625km. Specifically, the bow-shape echo of the cREF core is tilted
 491 in a northeast-southwest direction, forming a classic bow echo described in Fujita (1978).
 492 The precipitation feature (Figure 4f) displays a similar tilting shape along with a larger
 493 precipitating area and higher rainfall intensity in the northeast tail. A secondary cluster
 494 is found in the southwest tail with relatively low rainfall rate in the center part of
 495 the derecho. In contrast, the shapes of precipitating and cREF features in the WRF_RAP
 496 simulation are aligned in a more east-west direction with the most intense rainfall show-
 497 ing in the northwest part of the derecho. The meridional spread in WRF_RAP ($\sim 1.5^\circ$)
 498 is about half that in the observation and SCREAM_1.625km.

499 Even in the simulations with relatively better representation of the derecho (e.g.,
 500 SCREAM_3.25km, SCREAM_1.625km, and WRF_RAP), the simulated precipitation rate
 501 is higher than the observed (i.e., NCEP Stage IV). These moist biases are consistent with
 502 past studies, such as a study of daily WRF hindcasts of monsoon convections in Moker Jr
 503 et al. (2018). While observational precipitation bias may be a factor here, there is no ev-
 504 idence to suggest this is the case.

505 In Figures 2-4, it is obvious that some simulations (i.e, SCREAM_ERA5, SCREAM_ERAI,
 506 SCREAM_06UTC, SCREAM_LR, and WRF_NARR) are not able to capture the dere-
 507 cho accurately and are simply not comparable to the other simulations. Therefore, in
 508 the following discussions, we will only show results in the better simulations (i.e, SCREAM_6.5km,
 509 SCREAM_3.25km, SCREAM_1.625km, and WRF_RAP). Given the clear similarity of SCREAM_3.25km
 510 and SCREAM_H, the result from SCREAM_H will also not be displayed, except when
 511 there is a noteworthy result.

512 As mentioned in sections 1 and 2.3, our study emphasizes the assessment of the sim-
 513 ulated 10-m wind speeds because of its relevance to storm damage (Shourd, 2017; Shep-
 514 herd et al., 2021). Figure 5 shows the simulated 10-m wind speed maximum (m/s; shaded)
 515 in the period of 00:00 - 01:00 UTC 30 June 2012, calculated as the maximum of 15-minute
 516 instantaneous wind speed. The region of high wind speed in Figure 5 is wider/larger than
 517 the instantaneous gust front because it includes the wind swaths over an hour, captur-
 518 ing the movement of the derecho. The dot markers indicate the gust wind maximum (m/s;
 519 left panels) and wind speed maximum (m/s; right panels) calculated from 5-minute ASOS
 520 stations' records. To simplify the figure, only ASOS stations with gust reports are shown
 521 in the left panels and right panels only show stations with wind speed maximum higher
 522 than 5 m/s. All simulation results are shown at native grid points without regridding
 523 to minimize interpolation error. The ASOS stations have the caveat that they possibly
 524 do not capture the highest wind speed due to their limited spatial and temporal cover-
 525 age. Note that the ASOS gust wind speed is generally higher than regular wind speed
 526 by 2-10 m/s (see section 2.2 for details).

527 Compared with ASOS, SCREAM RRM performs well in simulating the observed
 528 10-m wind speed. The bow-shaped convective feature that produces extended swaths of
 529 damaging surface winds is one of the most important feature of the derecho, which is clearly
 530 shown in SCREAM_1.625km as a curved wind front, related to either a very strong rear-
 531 inflow jet or a strong downdraft (Fujita, 1978). It is obvious that the area with high wind
 532 speed at 3.25 and 1.625 km resolutions is significantly larger than that at 6.5 km res-

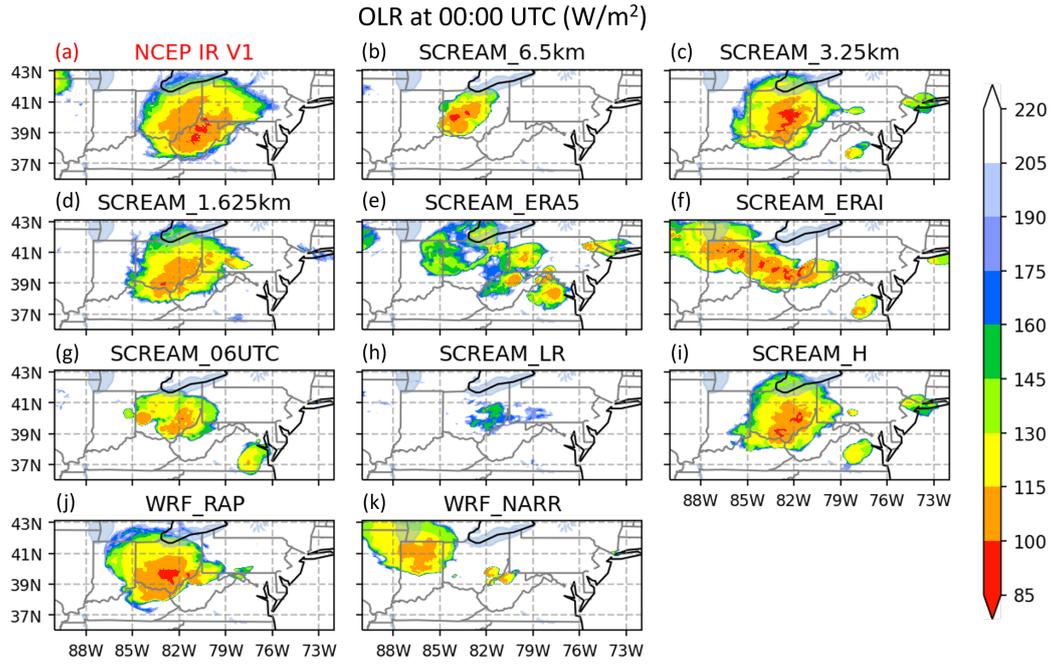


Figure 2. OLR (W/m^2) at 00:00 UTC 30 June 2012 in (a) NCEP IR V1, (b) SCREAM_6.5km, (c) SCREAM_3.25km, (d) SCREAM_1.625km, (e) SCREAM_ERA5, (f) SCREAM_ERAI, (g) SCREAM_06UTC, (h) SCREAM_LR, (i) SCREAM_H, (j) WRF_RAP, and (k) WRF_NARR. All datasets are remapped to 0.05° resolution. The panel with red title denotes the reference dataset.

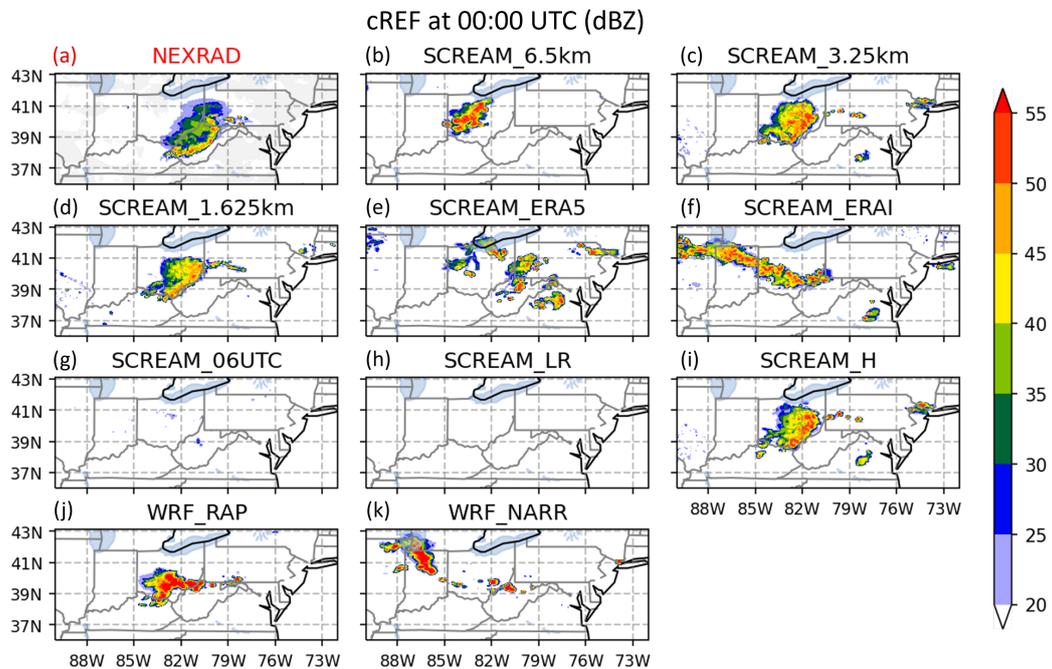


Figure 3. Same as Figure 2 but for cREF (dBZ). Panel (a) shows cREF in NEXRAD dataset.

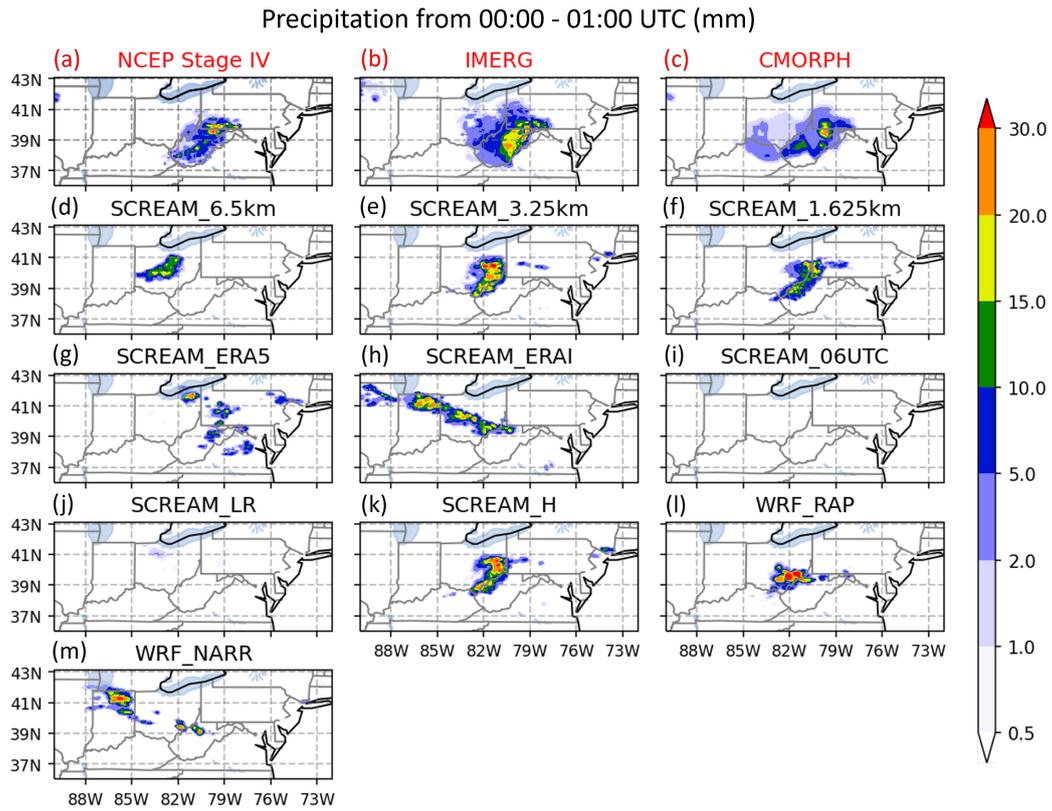


Figure 4. Same as Figure 2 but for accumulated precipitation (mm) from 00:00 - 01:00 UTC 30 June 2012. Panels (a-c) show NCEP Stage IV, IMERG, and CMORPH precipitation, respectively. All datasets are remapped to 0.1° .

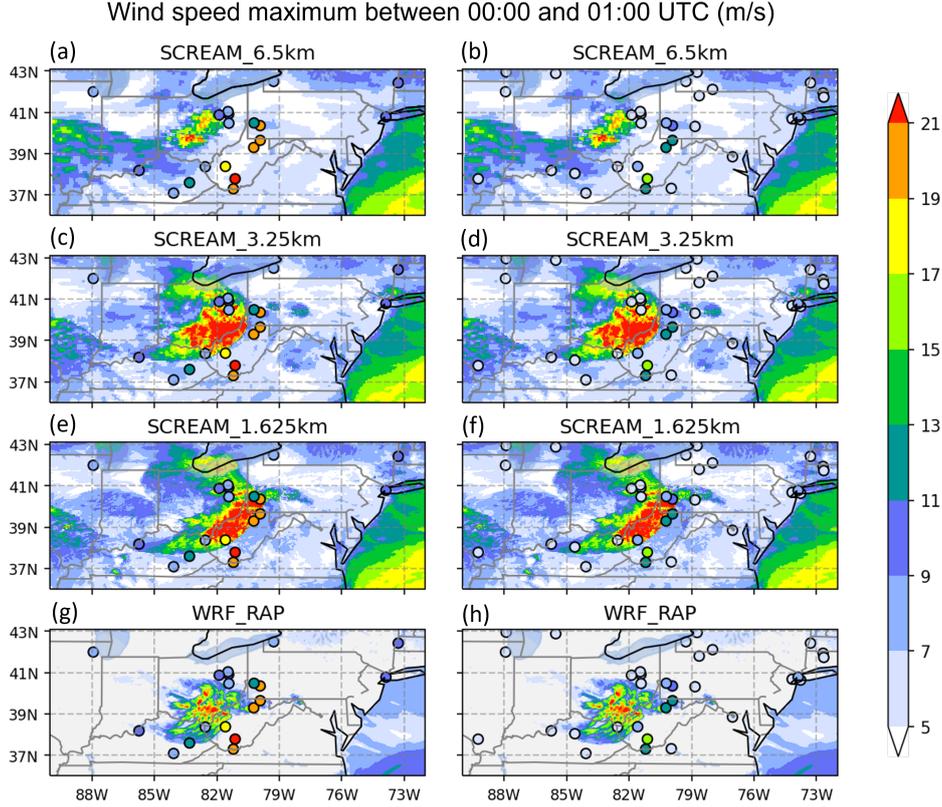


Figure 5. Wind speed maximum (m/s; shaded) between 00:00 and 01:00 UTC 30 June 2012 in (a-b) SCREAM_6.5km, (c-d) SCREAM_3.25km, (e-f) SCREAM_1.625km, and (g-h) WRF_RAP. The dot markers represent the ASOS gust wind speed maximum (m/s) in the left panels and wind speed maximum (m/s) in the right panels. The ASOS stations with wind speed maximum lower than 5 m/s are not shown in the right panels. All simulation results are displayed at raw grids.

533 olution, consistent with previous figures where enhanced fidelity is found at finer resolu-
 534 tion. Additionally, the wind front forward of the derecho is closest to the ASOS sites
 535 with wind gust reports in the SCREAM_1.625km as the derecho location is simulated
 536 best at the 1.625 km resolution (Figure 4f). The WRF simulation shows lower wind speeds
 537 than SCREAM not only in the derecho-covered area but also over the entire analysis do-
 538 main in general. More discussions about the simulated wind speeds and the different be-
 539 haviors between WRF and SCREAM will be presented in section 4.4.

540 4.2 Time Evolution

541 Figure 6 shows 2-hourly evolution of cREF (dBZ) from 18:00 UTC 29 June to 06:00
 542 UTC 30 June 2012 in the NEXRAD, SCREAM_6.5km, SCREAM_3.25km, SCREAM_1.625km,
 543 and WRF_RAP at 0.05° resolution. We only show cREF feature associated with the dere-
 544 cho identified using the TempestExtremes described in section 2.3 and remove other small
 545 clusters. Some time slots are not displayed (e.g., 18, 20, 22 UTC in SCREAM_6.5km)
 546 because the cREF feature does not qualified to be identified at that time (either too weak
 547 or too small using the defined thresholds).

548 The modelled track of the convective line broadly matches the observed one. The
 549 derecho-producing system proceeds southeastward from northern Indiana across central
 550 and southern Ohio with a strengthening convective core, reaching western West Virginia
 551 by 00:00 UTC. Over Ohio, the derecho system attains its greatest organization and strength.
 552 A rear-inflow notch at the back edge of system, which indicates an evaporatively cooled
 553 strong rear-inflow jet (Grim et al., 2009; Alliss & Hoffman, 2010), is evident before and
 554 during the leading line’s transformation into a bow echo over Ohio. The mature bow echo
 555 contains two bookend vortices, generally marking a region of enhanced downdraft and
 556 an increased probability of stronger winds at the surface. The signature progressive bow-
 557 ing presentation is evident in the SCREAM simulation at 3.25 and 1.626 km resolutions.
 558 For a sufficiently persistent MCS, the Coriolis force eventually leads to a strengthening
 559 of the cyclonic (or poleward) bookend vortex and a weakening of the anticyclonic (or equa-
 560 torward) vortex (Przybylinski, 1995; Schenkman & Xue, 2016). Accordingly, relatively
 561 fast eastward propagation is favored north of the front, with slower speed to its south
 562 in the observation and simulations. The storm system weakens as it moves into the south-
 563 ern New Jersey.

564 The development of the derecho from 18:00 to 22:00 UTC is significantly under-
 565 estimated in the SCREAM_6.5km, which is corrected at finer resolutions. The observed
 566 weakening (displayed as a discontinuity in the track) around 04:00 UTC and the north-
 567 ward jump around 06:00 UTC near New Jersey are also well reproduced by the SCREAM_1.625km.
 568 The location of the derecho shows roughly 2-hour delay in SCREAM_3.25km and WRF_RAP,
 569 and larger (~ 3 -hour) delay in SCREAM_6.5km. Comparing to SCREAM_3.25km, SCREAM
 570 at 1.625 km resolution reduces the delay by ~ 0.5 -1 hour. Longer postponements rang-
 571 ing from 3-8 hours were found in Shepherd et al. (2021), dependent on the model con-
 572 figurations.

573 All SCREAM and WRF simulations show larger cREF feature coverage than the
 574 NEXRAD. With that said, SCREAM_1.625km is the best ensemble simulating a nar-
 575 row linear core with cREF > 50 dBZ, most similar to NEXRAD, with extended spread
 576 in 40-50 dBZ. WRF shows significantly higher cREF than the SCREAM and NEXRAD
 577 by ~ 10 dBZ, consistent with the overestimated rainfall intensity in Figure 4.

578 To compare the location of the simulated derecho to the observation accurately,
 579 Figure 7 shows the time series of longitude (left panels) and latitude (right panels) of
 580 the cold cloud shield (top), cREF feature (middle), and precipitation feature (bottom)
 581 from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in SCREAM_6.5km (dark blue),
 582 SCREAM_3.25km (yellow), SCREAM_1.625km (red), WRF_RAP (green), and observa-
 583 tion (black). The solid line represents the center of the derecho at 15-minute frequency
 584 for all simulations and hourly frequency for the observation. Circle and triangle mark-
 585 ers denote the 2-hourly maximum and minimum of the longitude/latitude, respectively.

586 All simulations show western and southern biases ranging from 0 - 3° , associated with
 587 the time delay of the migration. SCREAM_1.625km provides the best simulated posi-
 588 tion in the zonal direction among all simulations with the eastern progressive edge of the
 589 derecho following the observed one. WRF_RAP simulates the best location in the merid-
 590 ional direction while SCREAM_1.625km exhibits a more northern position by 0 - 0.5° .

591 Figure 8 shows the time series of the cold cloud shield, cREF feature, and precip-
 592 itation feature areas. The dashed lines represent the raw simulation result frequency (15
 593 minutes) while results averaged to the observation frequency (hourly for cREF and pre-
 594 cipitation, and half hourly for OLR) are shown in the solid lines. The observed tempo-
 595 ral evolution of cold cloud shield is reproduced best by WRF_RAP with the largest cold
 596 cloud shield present around 03:00 UTC. WRF_RAP shows a similar cold cloud size to
 597 SCREAM_3.25km before 00:00 UTC and grows up to twice the size of SCREAM after
 598 00:00 UTC. SCREAM_1.625km captures the extending cold cloud shield before 23:00 UTC

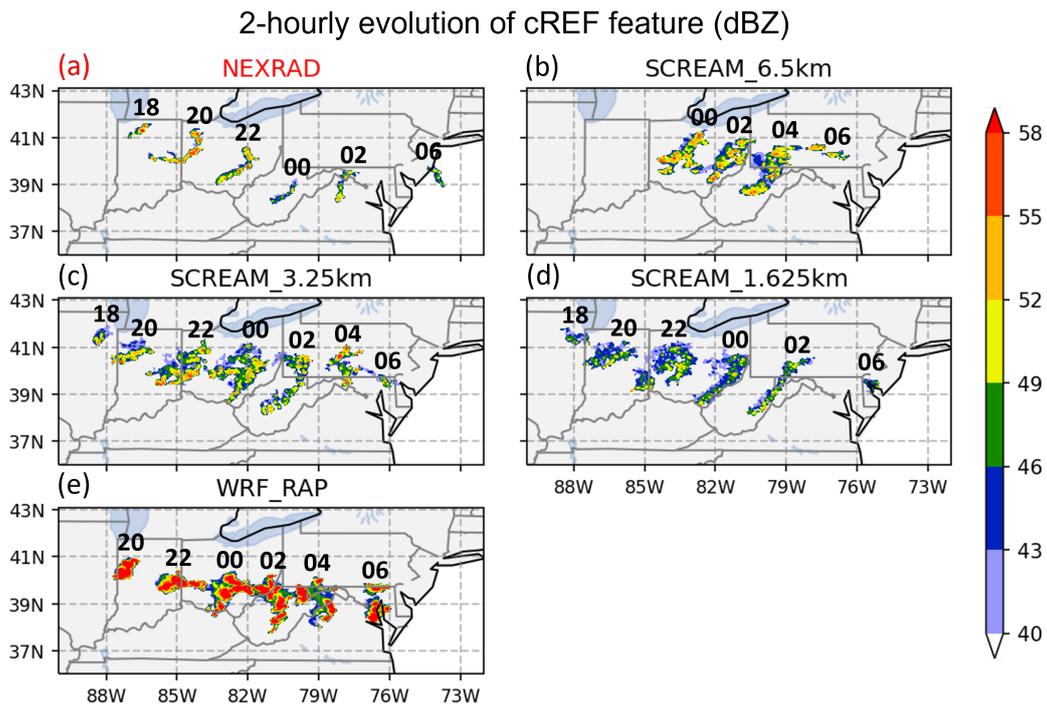


Figure 6. 2-hourly evolution of cREF feature (dBZ) from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in (a) NEXRAD, (b) SCREAM.6.5km, (c) SCREAM.3.25km, (d) SCREAM.1.625km, and (e) WRF_RAP at 0.05° resolution. The black bold numbers mark the hours in UTC.

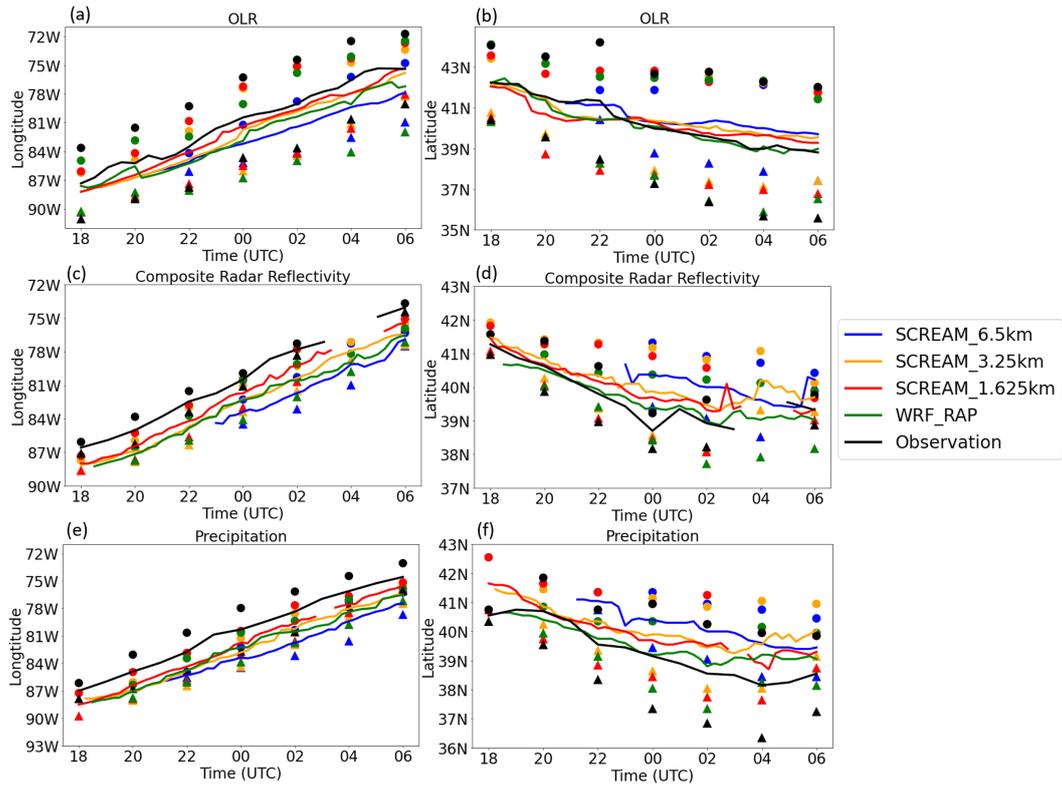


Figure 7. Time series of longitude (left) and latitude (right) of the cold cloud shield (top), cREF feature (middle), and precipitation features (bottom) from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in the SCREAM.6.5km (dark blue), SCREAM.3.25km (yellow), SCREAM.1.625km (red), WRF_RAP (green), and observation (black). The solid line represents the center of the derecho at 15-minute frequency for all simulations and hourly frequency for the observation. Circle and triangle markers denote the 2-hourly longitude/latitude maximum and minimum, respectively.

599 but shows a smaller cold cloud shield after 00:00 UTC by up to 50% than the observa-
600 tion.

601 All simulations overestimate the observed cREF feature size by up to four times
602 while they underestimate the precipitation feature size by $\sim 50\%$. The precipitation fea-
603 ture is significantly larger than the cREF feature by up to 10 times in the observation
604 (black lines in Figures 8b-c) indicating only approximately 10% of the precipitation fea-
605 ture is associated with high cREF (> 40 dBZ) whereas the rest of it has relatively low
606 cREF. However, precipitation and cREF features show comparable sizes in WRF_RAP
607 (greens lines in Figures 8b-c) suggesting almost the entire precipitation feature is asso-
608 ciated with high cREF. In the SCREAM simulations, about 50% of the precipitation fea-
609 ture is associated with high cREF. The results are consistent with Figures 3 and 6 where
610 the observation shows the most linear cREF area while WRF_RAP shows the widest cov-
611 erage of the high cREF.

612 The observed cREF feature develops strongly between 18:00-20:00 UTC reaching
613 its maximum coverage at 20:00 UTC, and persists until around 23:00 UTC. The precip-
614 itating area keeps expanding until 00:00 UTC when it starts to shrink while the cold cloud
615 shield remains extending for three more hours. Similarly, in SCREAM simulations, the
616 cREF and precipitation features show their coverage maxima 1-2 hours earlier than the
617 cold cloud. Despite the propagation delay (Figure 6), the largest precipitating area of
618 SCREAM occurs two hours earlier than observed, associated with the early decay of the
619 cold cloud shield (Figure 8a). The peak time of precipitation, cREF, and cold cloud fea-
620 ture is almost simultaneous in the WRF simulation with a delayed cREF/precipitation
621 feature area maximum by approximately 2.5-3 hours than the observation.

622 To evaluate the precipitation intensity, Figure 9a shows time series of regional-averaged
623 precipitation rate in the analysis domain (76° - 88° W, 36.5° - 42° N), shown as the red box
624 in Figure 1c, in the simulations and NCEP Stage IV precipitation dataset. Figure 9b is
625 similar to Figure 9a but averaged only over precipitating grid points with rainfall rate
626 > 1 mm/day. Additionally, averaged precipitation over the derecho identified by Tem-
627 pestExtremes is also examined (not shown) and implies similar results to Figure 9b. It
628 is not shown considering the calculation processes in Figure 9b are much simpler and achiev-
629 able for the broad research community without additional steps using TempestExtremes.

630 The observed regional-averaged precipitation peak time is captured by WRF ac-
631 curately in both Figures 9a and b, but the averaged precipitation magnitude over the
632 precipitating grid points is approximately twice as great as the observed (Figure 9b). WRF
633 has wet bias over the precipitating grids but the precipitating area (Figure 8c) is reduced
634 resulting in the similar magnitudes of precipitation peaks in the regional means (Figure
635 9a). While the maximum of averaged precipitation over the precipitating grid points is
636 similar at three SCREAM resolutions, the time delay in the peak time is greatest in SCREAM_6.5km
637 and declines at higher resolutions. SCREAM RRM simulations also show higher precip-
638 itation peaks than the observation by roughly 45%.

639 Figure 9c shows frequency distribution of the hourly precipitation rates at all grid
640 points within the analysis domain (76° - 88° W, 36.5° - 42° N) during 18:00 UTC 29 June to
641 06:00 UTC 30 June 2012 in solid lines. The dashed lines are the same as the solid lines
642 except for applying a 2-hour forward shift in the simulations, resulting in a period of 20:00
643 UTC 29 June to 08:00 UTC 30 June 2012; however, the conclusions are not sensitive to
644 the 2-hour shift. WRF_RAP shows the strongest wet bias, strongly overestimating the
645 observed precipitation rates higher than ~ 30 mm/day. However, WRF displays signif-
646 icantly lower frequency of precipitation rates below 30 mm/day and higher frequency for
647 precipitation above 50 mm/day, consistent with a smaller coverage of relatively shallow
648 precipitation and an overwhelming intense precipitation core in Figure 4. The SCREAM
649 RRM simulations also produce an excess of extremely high rainfall rates (> 350 mm/day)
650 but show lower frequency for the rainfall rates between 100 and 250 mm/day. SCREAM

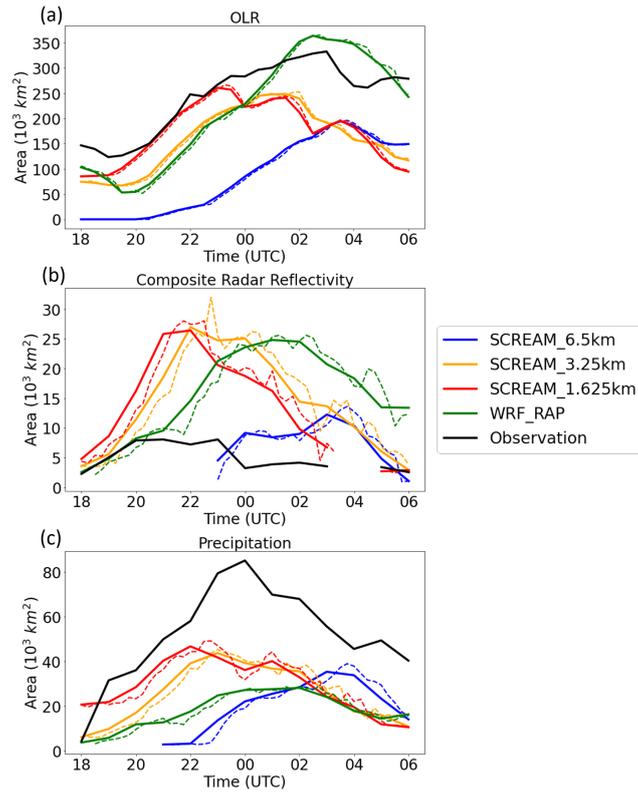


Figure 8. Time series of the area (10^3 km^2) of the identified features using (a) OLR, (b) cREF, and (c) precipitation from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in SCREAM_6.5km (dark blue), SCREAM_3.25km (yellow), SCREAM_1.625km (red), WRF_RAP (green), and observation (black). The dashed lines represent the results at 15-minute frequency for all simulations. The solid lines denote hourly frequency for cREF and precipitation and half hourly frequency for OLR.

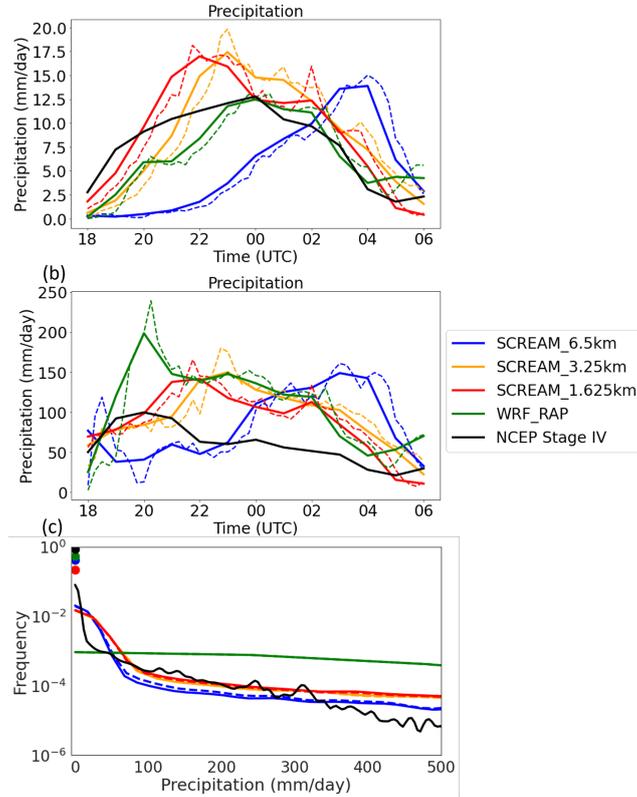


Figure 9. (a) Time series of regional-averaged precipitation rate (mm/day) in the analysis domain (76° - 88° W, 36.5° - 42° N), shown as the red box in the Figure 1c, from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in SCREAM.6.5km (dark blue), SCREAM.3.25km (yellow), SCREAM.1.625km (red), WRF_RAP (green), and NCEP Stage IV dataset (black). The dashed lines represent the simulation results at 15-minute frequency and the solid lines represent the hourly frequency. (b) is the same (a) but averaged only over precipitating grid points with rainfall rate > 1 mm/day. (c) Frequency distribution of hourly precipitation rates of all grid points within the analysis domain in the period of 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in the solid lines. The dashed lines are the same as the solid lines but apply a 2-hour forward shift for the simulations resulting in a period as 20:00 UTC 29 June to 08:00 UTC 30 June 2012. The dots on the y axis denote the frequencies of zero precipitation.

Table 5. Locations of three ASOS stations.

Number in Figure 1c	Station ID	Station Full Name	State	Longitude	Latitude
1	KCMH	COLUMBUS PORT COLUMBUS INTL AP	OH	39.99139°N	82.88083°W
2	KCKB	CLARKSBURG BENEDUM AP	WV	39.29556°N	80.22889°W
3	KDCA	WASHINGTON REAGAN AP	VA	38.8483°N	77.0341°W

Table 6. Metrics derived from the OLR averaged from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in the analysis region (76°-88°W, 36.5°-42°N) in each simulation using NCEP IR V1 dataset as the reference. The calculations of the metrics are present in section 2.4. Scores in parentheses are calculated by applying a two-hour forward shift to the simulation results (i.e., the averaging period for the simulations changes to 20:00 UTC 29 June to 08:00 UTC 30 June 2012). The red numbers denote the best scores in each category. BS and TS are calculated using the threshold of 230 W/m^2 .

Simulation Name	SCREAM.6.5km	SCREAM.3.25km	SCREAM.1.625km	SCREAM_H	WRF_RAP
RMSE	47.77 (39.35)	29.13 (26.43)	27.21 (25.29)	29.39 (27.69)	20.48 (17.65)
MAE	43.81 (34.61)	26.18 (22.61)	23.98 (21.46)	26.16 (23.89)	16.70 (14.15)
ME	43.81 (34.49)	22.31 (17.85)	18.86 (15.95)	21.66 (18.09)	8.18 (-0.11)
Pearson Correlation	0.89 (0.88)	0.89 (0.88)	0.88 (0.88)	0.87 (0.86)	0.88 (0.90)
Spearman Correlation	0.90 (0.86)	0.85 (0.80)	0.78 (0.78)	0.79 (0.74)	0.84 (0.82)
BS	0.20 (0.39)	0.58 (0.71)	0.73 (0.79)	0.61 (0.70)	0.83 (1.05)
TS	0.20 (0.39)	0.54 (0.64)	0.65 (0.71)	0.55 (0.61)	0.75 (0.87)

651 at 6.5 km exhibits a lower frequency in rainfall rates above 30 mm/day than simulations
652 at 3.25 and 1.625 km, associated with the smaller precipitation feature (Figure 8c).

653 Figure 10 shows the time series of wind speeds at three ASOS stations, marked by
654 the blue arrows in Figure 1c with details in Table 5. The three stations are selected to
655 be airport stations in the derecho propagation path, spread over three states to capture
656 various stages of the derecho life cycle, and using Figure 6 as a reference. In addition,
657 the three stations are confirmed to not have missing wind speed records during the anal-
658 ysis period (18:00-06:00 UTC). The simulation results are shown at the closest single grid
659 point to the specific ASOS station. While displaying time series at all ASOS stations is
660 not feasible on a single plot, the three stations selected provide insights into the timing
661 of the simulated gust fronts. Note that the ASOS station records have a high time fre-
662 quency of 5 minutes and the simulation results are derived at individual grid points, caus-
663 ing the high-frequency fluctuations in the time series.

664 The delayed wind speed peaks representing the gust fronts are found in all simu-
665 lations with reduced timing biases at finer resolution, consistent with the observed im-
666 provement in timing at 1.625 km resolution (Figure 6). The timing biases are approx-
667 imately 1-1.5 hours at 1.625 km resolution, 2-3 hours at 3.25 km resolution, and 3-4 hours
668 at 6.5 km resolution. The magnitudes of wind speed peaks in SCREAM_1.625km and
669 SCREAM_3.25km are either comparable to or larger than ASOS winds (black lines) by
670 0-30%, and lower than ASOS gust speed (purple lines) by $\sim 30\%$. On the other hand,
671 WRF_RAP shows the wind speed peak lower than ASOS wind speed by 27-70% and ASOS
672 gust by 56-85%.

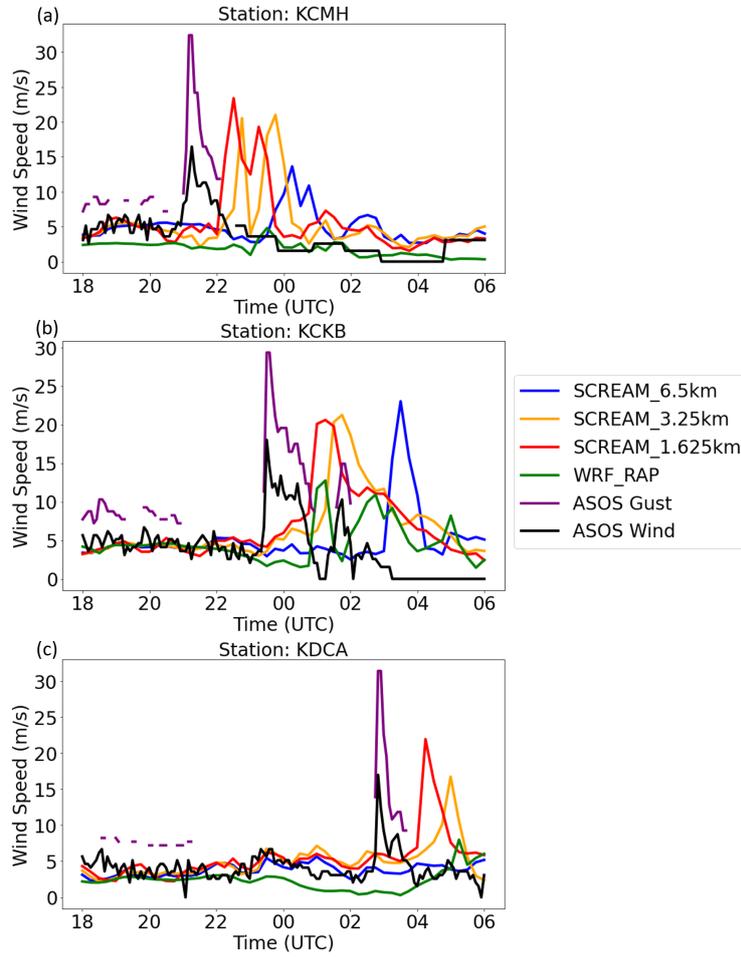


Figure 10. Time series of wind speeds (m/s) in the SCREAM.6.5km (dark blue), SCREAM.3.25km (yellow), SCREAM.1.625km (red), WRF_RAP (green), ASOS gust (purple), and ASOS wind (black) from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 at ASOS station (a) KCKB, (b) KCMH, and (c) KDCA. The time intervals of ASOS records and simulation outputs are 5 minutes and 15 minutes, respectively. The simulation results are shown at the closest grid point to the specific ASOS station.

Table 7. As Table 6 but derived from the precipitation accumulated from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 using NCEP Stage IV dataset as the reference. The last two columns show metrics calculated using the CMORPH and IMERG precipitation datasets comparing to the NCEP Stage IV dataset. BS and TS are calculated using the threshold of 15 mm.

Simulation/Observation Name	SCREAM_6.5km	SCREAM_3.25km	SCREAM_1.625km	SCREAM_H	WRF_RAP	CMORPH	IMERG
RMSE	6.99 (7.14)	8.50 (8.48)	7.73 (7.52)	8.17 (8.07)	9.26 (9.05)	5.78	8.65
MAE	4.32 (4.37)	4.66 (4.65)	4.49 (4.34)	4.81 (4.71)	4.87 (4.78)	3.88	5.47
ME	-1.73 (-1.46)	-0.70 (0.67)	1.03 (0.71)	1.16 (0.96)	-0.67 (-0.71)	2.89	4.56
Pearson Correlation	0.54 (0.53)	0.50 (0.51)	0.54 (0.55)	0.54 (0.54)	0.52 (0.53)	0.77	0.72
Spearman Correlation	0.70 (0.72)	0.73 (0.75)	0.71 (0.73)	0.71 (0.73)	0.73 (0.74)	0.87	0.86
BS	1.01 (1.07)	1.46 (1.48)	1.49 (1.41)	1.72 (1.67)	1.07 (1.08)	2.17	2.54
TS	0.28 (0.27)	0.23 (0.23)	0.23 (0.24)	0.24 (0.25)	0.25 (0.25)	0.30	0.28

673

4.3 Metrics

674

675

676

677

678

Table 6 displays metrics derived from the OLR in each simulation (see section 2.4 for the calculations of the metrics) with NCEP IR V1 dataset as the reference. The red number marks the best score in each category. The scores in parentheses are calculated by applying two-hour forward shift to the simulation results, providing better results (smaller biases) in all metrics except for two correlation scores.

679

680

681

682

683

684

685

686

687

688

689

690

691

WRF_RAP produces smaller biases in OLR than SCREAM when compared to observations. RMSE, MAE, and ME are lowest for WRF_RAP, particularly in the two-hour shifted ones, indicating WRF_RAP simulates the OLR field better than SCREAM. It is notable that SCREAM at finer resolutions shows better performance (in RMSE, MAE, and ME) than at coarser resolutions. The positive-biased OLR indicates a lower cloud top along with a smaller cold cloud shield (Figure 8a) in the simulations than the observed. The differences of Pearson correlations among all simulations are minor (< 0.03). Interestingly, despite previous analyses and other metrics showing better performance at finer resolution, Spearman correlation is highest in the coarsest simulation (SCREAM_6.5km), possibly caused by the underestimation of the cold cloud area at finer resolutions after 23:00 UTC (Figure 8a). BS showing values < 1 also indicates an underestimated cold cloud area in the simulations. WRF_RAP has the best representation of the cold cloud shield area, as indicated by the BS closest to 1, especially in the two-hour shifted BS.

692

693

694

695

696

697

698

The simulation using the H dynamical core (SCREAM_H) is also listed in Table 6. While it is not significantly different from the NH simulation in the snapshots of section 4.1 (Figures 2-5) and the time series in section 4.2 (not shown), we investigate whether the difference is more pronounced as the period is prolonged here. The H simulation shows slightly higher biases ($< 6\%$) in the two-hour shifted RMSE, MAE, and ME than SCREAM_3.25km, but the difference is much smaller than that among other simulation ensembles. As such, we attribute this difference to simulation variability rather than structural uncertainty.

699

700

701

702

703

704

705

706

707

708

709

Table 7 is the same as Table 6 but derived from precipitation accumulated from 18:00 UTC 29 June to 06:00 UTC 30 June 2012. Figure 11 shows the accumulated precipitation patterns along with two-hour shifted patterns shown in Figure S1. WRF_RAP shows the largest biases in RMSE and MAE, related to the overestimates of precipitation (Figures 4, 9, and 11). However, the ME in WRF_RAP is smallest in magnitude and becomes the best score among all simulations because the wet bias from enhanced precipitation intensity is offset by the reduced precipitating area (Figures 8c and 11e). The ME changes from negative to positive when SCREAM resolution becomes finer and the precipitating area of the derecho increases (Figure 8c). The ranges of RMSE and MAE in the simulations are comparable to those in CMORPH and IMERG, suggesting reasonable model performance in line with observational uncertainty. The Pearson and Spear-

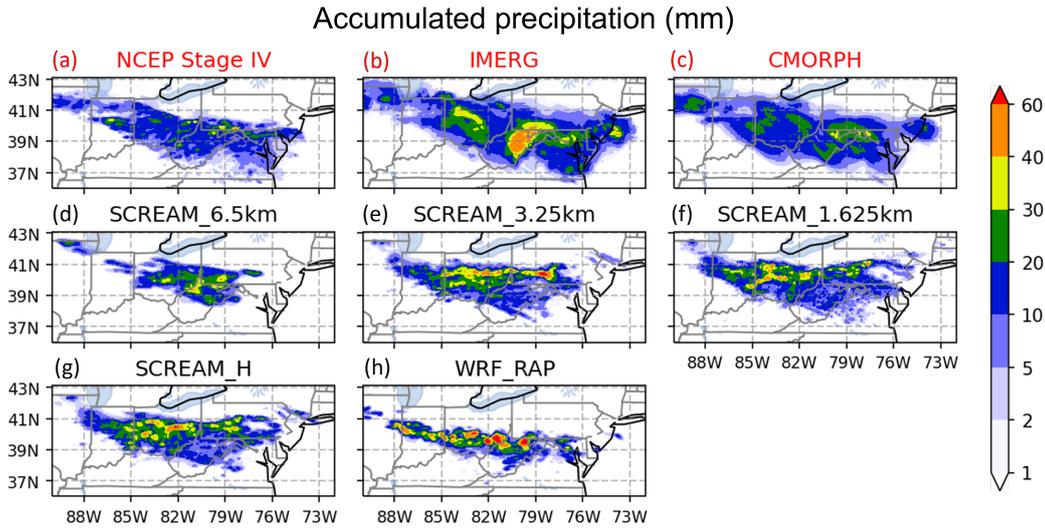


Figure 11. Accumulated precipitation (mm) from 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in (a) NCEP Stage IV, (b) CMORPH, (c) IMERG, (d) SCREAM_6.5km, (e) SCREAM_3.25km, (f) SCREAM_1.625km, (g) SCREAM_H, and (h) WRF_RAP.

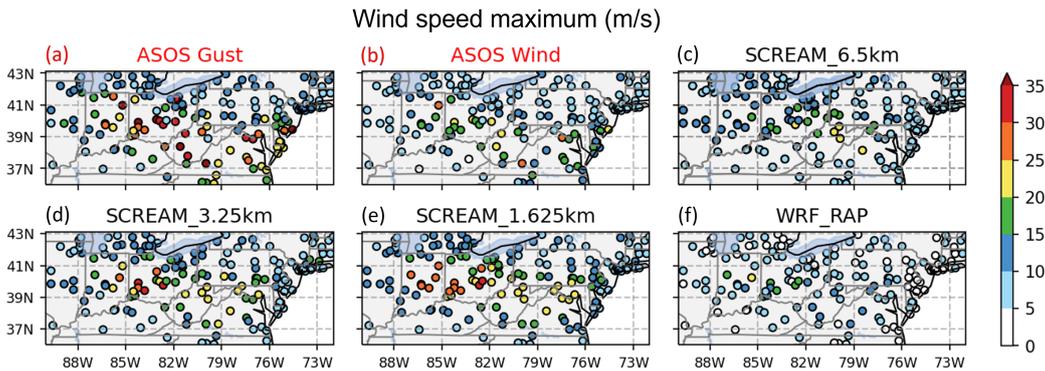


Figure 12. Wind speed maximum (m/s) between 18:00 UTC 29 June and 06:00 UTC 30 June 2012 in (a) ASOS gust, (b) ASOS wind, (c) SCREAM_6.5km, (d) SCREAM_3.25km, (e) SCREAM_1.625km, and (f) WRF_RAP. Only ASOS sites with gust reports during the period are displayed. Panels (c-f) show results at the closest grid points to the ASOS locations.

710 man correlations with 2-hour shift do not display significant differences among simulations
 711 (with difference < 0.03). The simulations have larger areas with higher accumu-
 712 lated precipitation amount (> 15 mm) than observations, illustrated by BS larger than
 713 1, and consistent with Figure 9c. Comparing SCREAM_H to SCREAM_3.25km, the H
 714 simulation shows worse BS, indicating the area with greater accumulated precipitation
 715 amount (> 15 mm) is enhanced when employing the H dynamical core.

716 Figure 12 shows 10-m wind speed maximum between 18:00 UTC 29 June and 06:00
 717 UTC 30 June 2012 in ASOS records and simulations. Only ASOS sites with gust reports
 718 during the period are displayed. Figure 13 shows the histogram of wind speed maximum
 719 to quantify the number of stations with wind speed maximum in each 5 m/s interval.

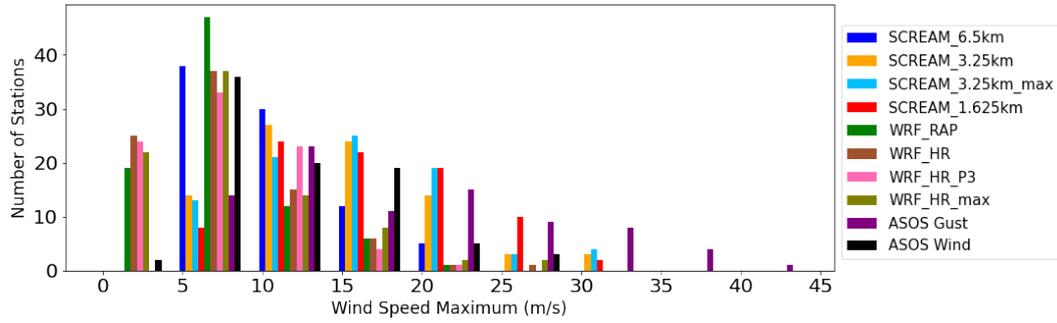


Figure 13. Histogram of wind speed maximum between 18:00 UTC 29 June and 06:00 UTC 30 June 2012 in SCREAM_6.5km (dark blue), SCREAM_3.25km (yellow), SCREAM_3.25km_max (light blue), SCREAM_1.625km (red), WRF_RAP (green), WRF_HR (brown), WRF_HR_P3 (pink), WRF_HR_max (olive), ASOS gust (purple), and ASOS wind (black) in the analysis region.

720 Figures S2-S3 are the same as Figures 12-13 but with the two-hour shift, producing similar results.
721

722 The wind speed maximum in SCREAM is generally between the values reported
723 by ASOS gust and ASOS wind. SCREAM produces wind speeds exceeding 25 m/s at
724 3.25 km and 1.625 km resolutions during approximately the first half of the derecho life
725 cycle, over central Indiana and Ohio, while the wind speed maximum is lower (20-25 m/s)
726 during the second half of the derecho life cycle over northern West Virginia when the dere-
727 cho size and intensity decline (Figures 8 and 9). SCREAM_6.5km displays an opposite
728 pattern with higher wind speeds in the second half of the derecho path, associated the
729 under-development of the system before 00:00 UTC at coarser resolution (Figure 6b).
730 WRF_RAP shows lower wind speeds than both ASOS gust and ASOS wind speeds, with
731 maximum wind speed lower than 25 m/s. This result is consistent with the WRF result
732 in Shepherd et al. (2021) (see their Figure 9) where the wind speed maximum is under
733 25 m/s. Given that wind damage generally occurs above 25.7 m/s and is proportional
734 to the cube of the wind speed, the underestimation of 10-m wind speeds may make it
735 difficult to leverage the WRF simulations to estimate wind damage during such storms.

736 In Figure 13, the distribution of wind speed maximum in SCREAM_1.625km agrees
737 best with the ASOS gust showing most of the stations in the range of 10-25 m/s as well
738 as ~ 15 stations with extremely high speed (> 25 m/s). None of the SCREAM simu-
739 lations display the wind maximum lower than 5 m/s. As the SCREAM grid spacing de-
740 creases, the histogram shifts towards being right-skewed, representing generally higher
741 wind speeds. For WRF_RAP, wind speed maximum higher 25 m/s is not represented.
742 Further, it has twice as many stations with wind speed maximum lower than 10 m/s as
743 as in ASOS wind, but only 35-70% of the ASOS stations fall into categories greater than
744 10 m/s.

745 4.4 Sensitivity Tests of Wind Speeds

746 To better understand the discrepancy between SCREAM and WRF simulations,
747 in term of 10-m wind speed (Figures 5, 10, 12 and 13), we consider three factors that
748 might affect the simulated wind speed and conduct additional sensitivity tests to pro-
749 vide more insights into the wind speed in diverse model configurations.

Table 8. Regional-averaged 10-m wind speed maximum (m/s) between 18:00 UTC 29 June to 06:00 UTC 30 June 2012 in the analysis region (76°-88°W, 36.5°-42°N). The reference simulation in the calculation of percentage is shown in the parenthesis of the last column.

Simulation Name	Regional-averaged 10-m Wind Speed Maximum (m/s)	Increase of Regional-averaged 10-m Wind Speed Maximum (%)
WRF_RAP	8.99	-
WRF_HR	9.36	4.12 (WRF_RAP)
WRF_HR_P3	9.46	1.07 (WRF_HR)
WRF_HR_max	10.28	9.83 (WRF_HR)
SCREAM_6.5km	10.71	-
SCREAM_3.25km	15.52	-
SCREAM_3.25km_max	16.47	6.12 (SCREAM_3.25km)
SCREAM_1.625km	16.62	-

750 Firstly, the SCREAM HR configuration (used in all SCREAM simulations except
751 for SCREAM_LR; Table 1) uses 128 vertical levels (92 levels below 100 hPa), while WRF_RAP
752 uses 45 vertical levels (all below 100 hPa), which may cause higher 10-m wind speed in
753 SCREAM than WRF. A new WRF simulation run with 72 vertical levels (WRF_HR;
754 Table 4) is performed and shown in Figure 14a. Comparing WRF_HR with Figure 12f,
755 increasing the number of vertical levels does lead to higher 10-m wind speed, especially
756 in central Indiana and southern Ohio. Table 8 quantifies the regional-averaged 10-m wind
757 speed maximum. WRF_HR has a similar wind maximum pattern as WRF_RAP but with
758 an increased regional-averaged wind maximum by 4.12%. The histogram of wind speed
759 maximum in the new WRF_HR is also shown in Figure 13. More vertical levels reduce
760 the wind speed maximum in the 5-10 m/s bin by 30% and slightly increases frequency
761 of high wind speed (> 25 m/s). However, the WRF_HR wind speed maximum distribu-
762 tion still exhibits a low bias when compared with the ASOS wind in all categories above
763 10 m/s.

764 Secondly, the simulated 10-m wind speed is also related to the microphysical scheme
765 applied since the microphysical scheme affects the convective structure and the cold pool
766 associated with it. A new simulation WRF_HR_P3 (Figure 14b, Table 4) is conducted
767 by replacing the microphysical scheme in WRF_HR with P3, which produces an insignif-
768 icant increase in regionally-averaged wind maximum (1%). The P3 scheme's impact is
769 more noticeable in shifting the derecho propagation path southward by 1-3° rather than
770 modifying the 10-m wind speed magnitude exclusively.

771 Thirdly, both SCREAM and WRF results in the previous analyses are instantane-
772 ous outputs at 15-minute frequency, which possibly do not capture the highest wind
773 speed during the 15-minute period. Therefore, we further output the wind speed max-
774 imum during each 15-minute period in SCREAM_3.25km and WRF_HR, labeled as SCREAM_3.25km_max
775 (Figure 14d) and WRF_HR_max (Figure 14c), respectively. This change causes an in-
776 crease of regional-averaged wind maximum by 9.83% in WRF and 6.12% in SCREAM
777 (Table 8), becoming the most influential factor in this section to the wind speed. This
778 suggests that future work involving the assessment of 10-m wind hindcast against high-
779 frequency observations (such as 5-minute ASOS) should consider a higher output fre-
780 quency of the wind speed than other variables. The highest wind speed in SCREAM_3.25km_max
781 then exceeds 30 m/s. WRF_HR_max display the most stations with wind speed max-
782 imum above 25 m/s among all WRF simulations. Although the wind speed increases (by
783 9.83%) in WRF_HR_max than WRF_RAP, there is an underestimation of wind speed
784 in every category above 10 m/s comparing WRF_HR_max to ASOS gust or wind (Fig-
785 ure 13). Specifically, the total number of stations with wind speeds > 10 m/s is 47 in
786 ASOS wind, 71 in ASOS gust, and 26 in WRF_HR_max (less than ASOS wind and gust
787 by 45% and 63%, respectively).

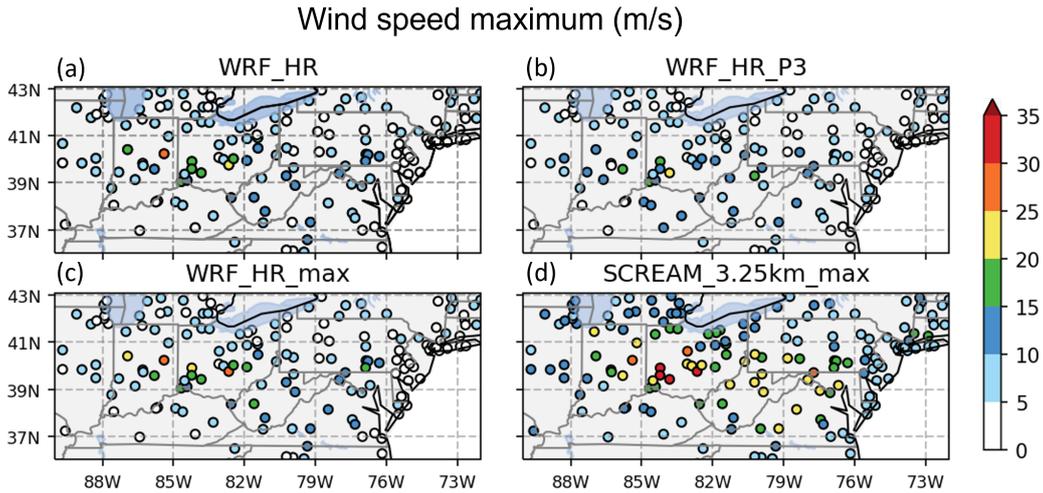


Figure 14. Same as Figure 12 but in (a) WRF_HR, (b) WRF_HR_P3, (c) WRF_HR_max, and (d) SCREAM_3.25km_max.

5 Conclusions

Climate models have been evaluated primarily using average behavior over large areas or long time periods (Gleckler et al., 2008; Eyring et al., 2019). However, it is rare for evaluations to consider the fidelity of simulating the most extreme and high-impact weather phenomena. Evaluating the capability of models to reproduce poorly predicted, but severe historic events will facilitate further model development and comparison, enable optimization of model configuration, and provide context for examining future changes in such events. In this work, we present one such extreme event testbed for evaluating climate modeling systems that operate at cloud resolving scales. The testbed focuses on the hindcast of the June 2012 North American derecho, which is one of the most devastating and strongly under-forecasted events in the US, with the intention of providing a largely comprehensive suite of diagnostics and statistical metrics for model assessment and intercomparison.

The metrics aim to assess the spatiotemporal characteristics of the derecho using RRM approach in SCREAM at various resolutions (Table 1) and regional WRF model at 4km (Table 4) against observations. Sensitivity tests address RRM grid spacing ranging from 6.5 - 1.625 km, differences between hydrostatic and nonhydrostatic dynamical cores, low-resolution and high-resolution model configurations, initialization time, and source for the ICs (i.e., RAP, ERA5, and ERAI). Besides OLR, precipitation, and composite radar reflectivity fields, this study places additional emphasis on the 10-m wind speed evaluation, which has not been thoroughly covered in previous studies. It is worth mentioning that the metrics package evaluated here is independent of the tracking method.

The representation of the derecho is shown to benefit from the finer horizontal resolutions in SCREAM, particularly at 1.625 km grid spacing, in a variety of ways including: cold cloud temperature and its coverage, radar reflectivity structure of the MCS, derecho position during the propagation, and simulated surface gust wind speed. The derecho-associated cold cloud in the SCREAM simulation at the coarsest resolution (6.5 km) is significantly underdeveloped with a smaller coverage maximum by $\sim 23\%$ and a longer lag (~ 2 hours) in the peak time compared to simulations at finer resolution. These results reinforce the need for higher resolution in operational convection-permitting models.

819 The simulations exhibit high dependence on the IC source and the initialization
 820 time, revealing the initial environment to be one of the most important factors for the
 821 simulation quality. Although it is impossible to determine the most superior product,
 822 as results may vary on a case by case basis, in this study RAP is found to be the best
 823 choice of IC source. Simulations initialized with RAP provide significantly improved per-
 824 formance in both SCREAM and WRF, compared with ERA5, ERAI, and NARR (Fig-
 825 ures 2-4) for this specific event. In particular, SCREAM simulations with ERA5 and ERAI
 826 are not able to generate a realistic organized convection pattern.

827 The SCREAM HR model configuration (no deep convection scheme and more ver-
 828 tical model levels) produces a significantly better storm than the LR model configura-
 829 tion, which fails to develop an organized precipitating system over the affected region
 830 (Figures 2-4). The simulation with the hydrostatic dynamical core is similar to the non-
 831 hydrostatic one when examining individual snapshots (Figures 2-4) but shows greater
 832 biases ($< 6\%$) in the averaged OLR over the 12-hour period (Table 6).

833 While both SCREAM and WRF models show high pattern correlations (> 0.88)
 834 between the simulated OLR and the observation (Table 6), SCREAM is characterized
 835 by lower cloud top (indicated by 33-42% more biases in OLR RMSE; Table 6) and smaller
 836 cold cloud coverage by up to 50% than WRF (Figure 8), especially in the second half
 837 of the derecho life cycle.

838 SCREAM and WRF simulations both capture the observed derecho track, but both
 839 produce a delay of approximately 2 hours in feature location and associated gust front
 840 timing (Figures 6 and 10). Among all simulations, SCREAM at 1.625 km resolution dis-
 841 plays the smallest time lag with an difference of ~ 0.5 -1 hour from the 3.25 km simu-
 842 lation. Both models overestimate the precipitation intensity over the precipitating grid
 843 points (up to 100% in WRF and 45% in SCREAM; Figure 9b) and the areas with com-
 844 posite radar reflectivity > 40 dBZ (up to 4 times in both models; Figure 8b), and un-
 845 derestimate the precipitating area ($\sim 70\%$ in WRF and 47% in SCREAM; Figure 8c).
 846 WRF yields higher wet biases (up to 20% higher in accumulated precipitation RMSE;
 847 Table 7 and Figure 9c) but over smaller precipitation feature by $\sim 45\%$ than SCREAM
 848 (Figure 8c). The overall bias magnitudes of 12-hour accumulated precipitation in the mod-
 849 els fall in the range of CMORPH and IMERG compared to the NCEP Stage IV precip-
 850 itation, except for a higher RMSE in WRF (Table 7). Our results highlight the impor-
 851 tance of using multiple metrics to reveal different aspects of the simulations and errors.

852 SCREAM captures the bow-shape echo with a tilted axis more realistically than
 853 WRF (Figure 3). Moreover, the largest discrepancies between SCREAM and WRF are
 854 apparent in the 10-m wind speed. SCREAM simulates a 10-m wind speed maximum in
 855 between ASOS wind and ASOS gust speeds and a highest wind speed above 30 m/s, sig-
 856 nificantly higher than WRF by $\sim 73\%$ (Table 8). WRF underestimates the wind speed
 857 maximum compared to either ASOS wind (by 27-70%) or gust speeds (by 56-85%; Fig-
 858 ure 10) and does not produce damaging wind speeds > 25 m/s (Figure 13). Further in-
 859 vestigation shows that this underestimation of the 10-m wind speed in WRF could be
 860 partly reduced by finer vertical resolution (4.12%) or changing the analyzed output from
 861 the 15-minute instantaneous model result to the maximum during each 15-minute in-
 862 terval (9.83%; Table 8).

863 Last but not least, we suggest some potential applications for future studies. SCREAM
 864 RRM demonstrates competitive utility for studying individual high-impact weather events
 865 when compared to a high-resolution regional climate model (WRF), and so could be em-
 866 ployed for future regional climate model simulations. We argue that it could be useful
 867 for assessing and tuning resolution-dependent configurations in global models and for
 868 short-term weather prediction at fine scales (Zarzycki & Jablonowski, 2015). We further
 869 expect the extreme weather testbed described here is useful for future cloud-resolving
 870 model intercomparisons, such as to models from the DYnamics of the Atmospheric gen-

871 eral circulation On Non-hydrostatic Domains (DYAMOND) project (Stevens et al., 2019),
 872 performed in similar hindcast mode. This suite of assessment will be useful in objectively
 873 evaluating model design choices related to extreme weather phenomenon, building cred-
 874 ibility for extreme event attribution, and developing physical climate storylines to ex-
 875 plore plausible changes of extreme events in the future.

876 **Appendix A Derecho tracking with TempestExtremes**

877 For feature tracking in the simulations and observations, we use TempestExtremes
 878 2.2.1 (Ullrich & Zarzycki, 2017; Ullrich et al., 2021). The exact commands employed in
 879 this analysis are provided here for reference.

```
880 $STEMPESTEXTREMESDIR/DetectBlobs --in_data FLUT.nc --out DetectBlobs.FLUT.nc
881 --thresholdcmd "FLUT,<,163,0" --geofiltercmd "area,>=,5000km2"
882 --lonname lon --latname lat --regional
```

```
883 $STEMPESTEXTREMESDIR/StitchBlobs --in DetectBlobs.FLUT.nc
884 --out StitchBlobs.FLUT.nc --var "binary_tag" --outvar "id" --mintime "6h"
885 --min_overlap_prev 50 --regional --lonname lon --latname lat
```

886 **Data Availability Statement**

887 SCREAM is available online (E3SM Project, 2022). Simulation results (including
 888 SCREAM and WRF) and scripts used to plot figures could be archived at Zenodo ([https://](https://doi.org/10.5281/zenodo.6617206)
 889 doi.org/10.5281/zenodo.6617206).

890 **Acknowledgments**

891 This research was supported as part of the Energy Exascale Earth System Model (E3SM)
 892 project “The Simple Cloud-Resolving E3SM Atmosphere Model (SCREAM)”, funded
 893 by the U.S. Department of Energy, Office of Science, Office of Biological and Environ-
 894 mental Research. This work was performed under the auspices of the U.S. Department
 895 of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

896 We wish to thank the editor and anonymous reviewers for their suggestions to im-
 897 prove this manuscript. The authors would also like to thank Ben Hillman for his instruc-
 898 tions about building new RRM grids in SCREAM.

899 **References**

- 900 AghaKouchak, A., Behrangi, A., Sorooshian, S., Hsu, K., & Amitai, E. (2011).
 901 Evaluation of satellite-retrieved extreme precipitation rates across the central
 902 united states. *Journal of Geophysical Research: Atmospheres*, *116*(D2).
 903 AghaKouchak, A., Mehran, A., Norouzi, H., & Behrangi, A. (2012). Systematic and
 904 random error components in satellite precipitation data sets. *Geophysical Re-*
 905 *search Letters*, *39*(9).
 906 Alliss, R., & Hoffman, M. (2010). Quasi-linear convective system mesovorticies and
 907 tornadoes. *Meteorology Program, Iowa State University, Ames*.
 908 Anthes, R. A. (1983). Regional models of the atmosphere in middle latitudes.
 909 *Monthly weather review*, *111*(6), 1306–1335.
 910 Ban, N., Schmidli, J., & Schär, C. (2014). Evaluation of the convection-resolving re-
 911 gional climate modeling approach in decade-long simulations. *Journal of Geo-*
 912 *physical Research: Atmospheres*, *119*(13), 7889–7907.
 913 Beck, H. E., Pan, M., Roy, T., Weedon, G. P., Pappenberger, F., Van Dijk, A. I.,
 914 ... Wood, E. F. (2019). Daily evaluation of 26 precipitation datasets using

- 915 stage-iv gauge-radar data for the conus. *Hydrology and Earth System Sciences*,
 916 *23*(1), 207–224.
- 917 Beck, H. E., Vergopolan, N., Pan, M., Levizzani, V., Van Dijk, A. I., Weedon, G. P.,
 918 ... Wood, E. F. (2017). Global-scale evaluation of 22 precipitation datasets
 919 using gauge observations and hydrological modeling. *Hydrology and Earth*
 920 *System Sciences*, *21*(12), 6201–6217.
- 921 Benjamin, S. G., Weygandt, S. S., Brown, J. M., Hu, M., Alexander, C. R.,
 922 Smirnova, T. G., ... others (2016). A north american hourly assimilation
 923 and model forecast cycle: The rapid refresh. *Monthly Weather Review*, *144*(4),
 924 1669–1694.
- 925 Bernardet, L. R., & Cotton, W. R. (1998). Multiscale evolution of a derecho-
 926 producing mesoscale convective system. *Monthly weather review*, *126*(11),
 927 2991–3015.
- 928 Bertagna, L., Guba, O., Taylor, M. A., Foucar, J. G., Larkin, J., Bradley, A. M.,
 929 ... Salinger, A. G. (2020). A Performance-Portable Nonhydrostatic At-
 930 mospheric Dycore for the Energy Exascale Earth System Model Running at
 931 Cloud-Resolving Resolutions. In *Sc20: International conference for high per-*
 932 *formance computing, networking, storage and analysis* (pp. 1–14).
- 933 Bowman, K. P., & Homeyer, C. R. (2017). *Gridrad - three-dimensional gridded*
 934 *nextrad wsr-88d radar data*. Boulder CO: Research Data Archive at the Na-
 935 tional Center for Atmospheric Research, Computational and Information
 936 Systems Laboratory. Retrieved from <https://doi.org/10.5065/D6NK3CR7>
- 937 Brockhaus, P., Luthi, D., & Schar, C. (2008). Aspects of the diurnal cycle in a re-
 938 gional climate model. *Meteorologische Zeitschrift*, *17*(4), 433–444.
- 939 Brousseau, P., Seity, Y., Ricard, D., & Léger, J. (2016). Improvement of the fore-
 940 cast of convective activity from the arome-france system. *Quarterly Journal of*
 941 *the Royal Meteorological Society*, *142*(699), 2231–2243.
- 942 Caldwell, P. M., Mametjanov, A., Tang, Q., Van Roekel, L. P., Golaz, J.-C., Lin,
 943 W., ... others (2019). The DOE E3SM coupled model version 1: Descrip-
 944 tion and results at high resolution. *Journal of Advances in Modeling Earth*
 945 *Systems*, *11*(12), 4095–4146.
- 946 Caldwell, P. M., Terai, C. R., Hillman, B., Keen, N. D., Bogenschutz, P., Lin, W.,
 947 ... others (2021). Convection-permitting simulations with the e3sm global
 948 atmosphere model. *Journal of Advances in Modeling Earth Systems*, *13*(11),
 949 e2021MS002544.
- 950 Carey, L. D., Murphy, M. J., McCormick, T. L., & Demetriades, N. W. (2005).
 951 Lightning location relative to storm structure in a leading-line, trailing-
 952 stratiform mesoscale convective system. *Journal of Geophysical Research:*
 953 *Atmospheres*, *110*(D3).
- 954 Chen, F., & Dudhia, J. (2001). Coupling an advanced land surface–hydrology model
 955 with the penn state–ncar mm5 modeling system. part i: Model implementation
 956 and sensitivity. *Monthly weather review*, *129*(4), 569–585.
- 957 Corfidi, S. F., Coniglio, M. C., Cohen, A. E., & Mead, C. M. (2016). A proposed
 958 revision to the definition of “derecho”. *Bulletin of the American Meteorological*
 959 *Society*, *97*(6), 935–949.
- 960 Dai, A., Giorgi, F., & Trenberth, K. E. (1999). Observed and model-simulated di-
 961 urnal cycles of precipitation over the contiguous united states. *Journal of Geo-*
 962 *physical Research: Atmospheres*, *104*(D6), 6377–6402.
- 963 Davis, C., Brown, B., & Bullock, R. (2006). Object-based verification of precipi-
 964 tation forecasts. part i: Methodology and application to mesoscale rain areas.
 965 *Monthly Weather Review*, *134*(7), 1772–1784.
- 966 Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S.,
 967 ... others (2011). The era-interim reanalysis: Configuration and performance
 968 of the data assimilation system. *Quarterly Journal of the royal meteorological*
 969 *society*, *137*(656), 553–597.

- 970 Demaria, E., Rodriguez, D., Ebert, E., Salio, P., Su, F., & Valdes, J. (2011). Evalu-
 971 ation of mesoscale convective systems in south america using multiple satellite
 972 products and an object-based approach. *Journal of Geophysical Research:*
 973 *Atmospheres*, 116(D8).
- 974 Dennis, J., Fournier, A., Spatz, W. F., St-Cyr, A., Taylor, M. A., Thomas, S. J., &
 975 Tufo, H. (2005). High-resolution mesh convergence properties and parallel
 976 efficiency of a spectral element atmospheric dynamical core. *The International*
 977 *Journal of High Performance Computing Applications*, 19(3), 225–235.
- 978 Dennis, J. M., Edwards, J., Evans, K. J., Guba, O., Lauritzen, P. H., Mirin, A. A.,
 979 ... Worley, P. H. (2012). CAM-SE: A scalable spectral element dynamical
 980 core for the Community Atmosphere Model. *The International Journal of High*
 981 *Performance Computing Applications*, 26(1), 74–89.
- 982 Du, J. (2011). *Ncep/emc u.s. stage iv imagery. version 1.0*. UCAR/NCAR
 983 - Earth Observing Laboratory. Retrieved from [https://doi.org/10.26023/](https://doi.org/10.26023/VOKW-KGOV-9JOW)
 984 [VOKW-KGOV-9JOW](https://doi.org/10.26023/VOKW-KGOV-9JOW)
- 985 Dye, J., Winn, W., Jones, J., & Breed, D. (1989). The electrification of new mexico
 986 thunderstorms: 1. relationship between precipitation development and the on-
 987 set of electrification. *Journal of Geophysical Research: Atmospheres*, 94(D6),
 988 8643–8656.
- 989 E3SM Project. (2022). *Simple cloud-resolving e3sm atmosphere model (scream)*.
 990 Retrieved from [https://github.com/E3SM-Project/scream/releases/tag/](https://github.com/E3SM-Project/scream/releases/tag/SCREAM_Derecho_RRM_JAMES2022)
 991 [SCREAM_Derecho_RRM_JAMES2022](https://github.com/E3SM-Project/scream/releases/tag/SCREAM_Derecho_RRM_JAMES2022)
- 992 Elliott, E. J., Yu, S., Kooperman, G. J., Morrison, H., Wang, M., & Pritchard, M. S.
 993 (2016). Sensitivity of summer ensembles of fledgling superparameterized us
 994 mesoscale convective systems to cloud resolving model microphysics and grid
 995 configuration. *Journal of Advances in Modeling Earth Systems*, 8(2), 634–649.
- 996 Evans, K. J., Lauritzen, P., Mishra, S., Neale, R., Taylor, M., & Tribbia, J. (2013).
 997 Amip simulation with the cam4 spectral element dynamical core. *Journal of*
 998 *Climate*, 26(3), 689–709.
- 999 Eyering, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P.,
 1000 ... others (2019). Taking climate model evaluation to the next level. *Nature*
 1001 *Climate Change*, 9(2), 102–110.
- 1002 Feng, Z., Leung, L. R., Houze Jr, R. A., Hagos, S., Hardin, J., Yang, Q., ... Fan, J.
 1003 (2018). Structure and evolution of mesoscale convective systems: Sensitivity to
 1004 cloud microphysics in convection-permitting simulations over the united states.
 1005 *Journal of Advances in Modeling Earth Systems*, 10(7), 1470–1494.
- 1006 Feng, Z., Leung, L. R., Liu, N., Wang, J., Houze Jr, R. A., Li, J., ... Guo, J. (2021).
 1007 A global high-resolution mesoscale convective system database using satellite-
 1008 derived cloud tops, surface precipitation, and tracking. *Journal of Geophysical*
 1009 *Research: Atmospheres*, 126(8), e2020JD034202.
- 1010 Feng, Z., Song, F., Sakaguchi, K., & Leung, L. R. (2021). Evaluation of mesoscale
 1011 convective systems in climate simulations: Methodological development and
 1012 results from mpas-cam over the united states. *Journal of Climate*, 34(7),
 1013 2611–2633.
- 1014 Fierro, A. O., Gao, J., Ziegler, C. L., Mansell, E. R., MacGorman, D. R., & Dem-
 1015 bek, S. R. (2014). Evaluation of a cloud-scale lightning data assimilation
 1016 technique and a 3dvar method for the analysis and short-term forecast of the
 1017 29 june 2012 derecho event. *Monthly Weather Review*, 142(1), 183–202.
- 1018 Figurski, M., Nykiel, G., Jaczewski, A., Baldysz, Z., & Wdowikowski, M. (2017).
 1019 The impact of initial and boundary conditions on severe weather event simu-
 1020 lations using a high-resolution wrf model. case study of the derecho event in
 1021 poland on 11 august 2017. *Meteorology Hydrology and Water Management*.
- 1022 Fiolleau, T., & Roca, R. (2013). An algorithm for the detection and tracking of
 1023 tropical mesoscale convective systems using infrared images from geostation-
 1024 ary satellite. *IEEE transactions on Geoscience and Remote Sensing*, 51(7),

- 1025 4302–4315.
- 1026 Fosser, G., Khodayar, S., & Berg, P. (2015). Benefit of convection permitting cli-
 1027 mate model simulations in the representation of convective precipitation. *Cli-*
 1028 *mate Dynamics*, 44(1), 45–60.
- 1029 Fujita, T. T. (1978). Manual of downburst identification for project nim-
 1030 rod(atmospheric circulation). *Satellite and Mesometeorology Research Paper*,
 1031 104.
- 1032 Gao, Y., Leung, L. R., Zhao, C., & Hagos, S. (2017). Sensitivity of us summer
 1033 precipitation to model resolution and convective parameterizations across gray
 1034 zone resolutions. *Journal of Geophysical Research: Atmospheres*, 122(5),
 1035 2714–2733.
- 1036 Gleckler, P. J., Taylor, K. E., & Doutriaux, C. (2008). Performance metrics for cli-
 1037 mate models. *Journal of Geophysical Research: Atmospheres*, 113(D6).
- 1038 Golaz, J.-C., Caldwell, P. M., Van Roekel, L. P., Petersen, M. R., Tang, Q., Wolfe,
 1039 J. D., ... others (2019). The DOE E3SM coupled model version 1: Overview
 1040 and evaluation at standard resolution. *Journal of Advances in Modeling Earth*
 1041 *Systems*, 11(7), 2089–2129.
- 1042 Grim, J. A., Rauber, R. M., McFarquhar, G. M., Jewett, B. F., & Jorgensen, D. P.
 1043 (2009). Development and forcing of the rear inflow jet in a rapidly develop-
 1044 ing and decaying squall line during bamex. *Monthly weather review*, 137(4),
 1045 1206–1229.
- 1046 Grunzke, C. T., & Evans, C. (2017). Predictability and dynamics of warm-core
 1047 mesoscale vortex formation with the 8 may 2009 “super derecho” event.
 1048 *Monthly Weather Review*, 145(3), 811–832.
- 1049 Guastini, C. T., & Bosart, L. F. (2016). Analysis of a progressive derecho climatol-
 1050 ogy and associated formation environments. *Monthly Weather Review*, 144(4),
 1051 1363–1382.
- 1052 Haberlie, A. M., & Ashley, W. S. (2019). A radar-based climatology of mesoscale
 1053 convective systems in the united states. *Journal of Climate*, 32(5), 1591–1606.
- 1054 Halverson, J. B. (2014). A mighty wind: The derecho of june 29, 2012. *Weatherwise*,
 1055 67(4), 24–31.
- 1056 Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J.,
 1057 ... others (2018). *ERA5 hourly data on pressure levels from 1979 to present*.
 1058 Copernicus Climate Change Service (C3S) Climate Data Store (CDS). Re-
 1059 trieved from [https://cds.climate.copernicus.eu/cdsapp#!/dataset/
 1060 reanalysis-era5-pressure-levels?tab=overview](https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels?tab=overview)
- 1061 Hong, Y., Hsu, K.-L., Sorooshian, S., & Gao, X. (2004). Precipitation estimation
 1062 from remotely sensed imagery using an artificial neural network cloud classifi-
 1063 cation system. *Journal of Applied Meteorology*, 43(12), 1834–1853.
- 1064 Houze Jr, R. A. (2004). Mesoscale convective systems. *Reviews of Geophysics*,
 1065 42(4).
- 1066 Hu, H., Leung, L. R., & Feng, Z. (2020a). Observed warm-season characteristics of
 1067 mcs and non-mcs rainfall and their recent changes in the central united states.
 1068 *Geophysical Research Letters*, 47(6), e2019GL086783.
- 1069 Hu, H., Leung, L. R., & Feng, Z. (2020b). Understanding the distinct impacts of
 1070 mcs and non-mcs rainfall on the surface water balance in the central united
 1071 states using a numerical water-tagging technique. *Journal of Hydrometeorol-*
 1072 *ogy*, 21(10), 2343–2357.
- 1073 Hu, H., Leung, L. R., & Feng, Z. (2021). Early warm-season mesoscale convective
 1074 systems dominate soil moisture–precipitation feedback for summer rainfall
 1075 in central united states. *Proceedings of the National Academy of Sciences*,
 1076 118(43).
- 1077 Huang, X., Hu, C., Huang, X., Chu, Y., Tseng, Y.-h., Zhang, G. J., & Lin, Y.
 1078 (2018). A long-term tropical mesoscale convective systems dataset based
 1079 on a novel objective automatic tracking algorithm. *Climate dynamics*, 51(7),

- 1080 3145–3159.
- 1081 Huang, X., & Ullrich, P. A. (2017). The changing character of twenty-first-century
1082 precipitation over the western united states in the variable-resolution cesm.
1083 *Journal of Climate*, *30*(18), 7555–7575.
- 1084 Huffman, G. J., Bolvin, D. T., Braithwaite, D., Hsu, K., Joyce, R., Xie, P., & Yoo,
1085 S.-H. (2015). Nasa global precipitation measurement (gpm) integrated multi-
1086 satellite retrievals for gpm (imerg). *Algorithm Theoretical Basis Document*
1087 *(ATBD) Version, 4*, 26. Retrieved from [https://disc.gsfc.nasa.gov/
1088 datasets/GPM_3IMERGDF_06/summary?keywords=%22IMERG%20final%22](https://disc.gsfc.nasa.gov/datasets/GPM_3IMERGDF_06/summary?keywords=%22IMERG%20final%22)
- 1089 Iacono, M. J., Delamere, J. S., Mlawer, E. J., Shephard, M. W., Clough, S. A., &
1090 Collins, W. D. (2008). Radiative forcing by long-lived greenhouse gases:
1091 Calculations with the aer radiative transfer models. *Journal of Geophysical*
1092 *Research: Atmospheres*, *113*(D13).
- 1093 Janjić, Z. I. (1994). The step-mountain eta coordinate model: Further developments
1094 of the convection, viscous sublayer, and turbulence closure schemes. *Monthly*
1095 *weather review*, *122*(5), 927–945.
- 1096 Janowiak, J., Joyce, B., & Xie, P. (2017). *Ncep/cpc l3 half hourly 4km global (60s*
1097 *- 60n) merged ir v1 (gpm_mergir)*. Goddard Earth Sciences Data and Informa-
1098 tion Services Center (GES DISC). Retrieved from [https://disc.gsfc.nasa
1099 .gov/datasets/GPM_MERGIR_1/summary](https://disc.gsfc.nasa.gov/datasets/GPM_MERGIR_1/summary)
- 1100 Johns, R. H., & Hirt, W. D. (1987). Derechos: Widespread convectively induced
1101 windstorms. *Weather and Forecasting*, *2*(1), 32–49.
- 1102 Joyce, R. J., Janowiak, J. E., Arkin, P. A., & Xie, P. (2004). Cmorph: A method
1103 that produces global precipitation estimates from passive microwave and in-
1104 frared data at high spatial and temporal resolution. *Journal of hydrometeorol-*
1105 *ogy*, *5*(3), 487–503.
- 1106 Kharin, V. V., Zwiers, F. W., Zhang, X., & Hegerl, G. C. (2007). Changes in tem-
1107 perature and precipitation extremes in the ipcc ensemble of global coupled
1108 model simulations. *Journal of Climate*, *20*(8), 1419–1444.
- 1109 Kolios, S., & Feidas, H. (2010). A warm season climatology of mesoscale convec-
1110 tive systems in the mediterranean basin using satellite data. *Theoretical and*
1111 *applied climatology*, *102*(1), 29–42.
- 1112 Lin, Y., & Mitchell, K. E. (2005). 1.2 the ncep stage ii/iv hourly precipitation anal-
1113 yses: Development and applications. In *Proceedings of the 19th conference hy-*
1114 *drology, american meteorological society, san diego, ca, usa* (Vol. 10).
- 1115 Liu, W., Ullrich, P. A., Guba, O., Caldwell, P. M., & Keen, N. D. (2022). An as-
1116 sessment of nonhydrostatic and hydrostatic dynamical cores at seasonal time
1117 scales in the energy exascale earth system model (e3sm). *Journal of Advances*
1118 *in Modeling Earth Systems*, *14*(2), e2021MS002805.
- 1119 Maddox, R. A. (1980). Mesoscale convective complexes. *Bulletin of the American*
1120 *Meteorological Society*, 1374–1387.
- 1121 Mesinger, F., DiMego, G., Kalnay, E., Mitchell, K., Shafran, P. C., Ebisuzaki, W.,
1122 ... others (2006). North american regional reanalysis. *Bulletin of the Ameri-*
1123 *can Meteorological Society*, *87*(3), 343–360.
- 1124 Miao, C., Ashouri, H., Hsu, K.-L., Sorooshian, S., & Duan, Q. (2015). Evaluation
1125 of the persiann-cdr daily rainfall estimates in capturing the behavior of ex-
1126 treme precipitation events over china. *Journal of Hydrometeorology*, *16*(3),
1127 1387–1396.
- 1128 Mittermaier, M., & Roberts, N. (2010). Intercomparison of spatial forecast verifica-
1129 tion methods: Identifying skillful spatial scales using the fractions skill score.
1130 *Weather and Forecasting*, *25*(1), 343–354.
- 1131 Moker Jr, J. M., Castro, C. L., Arellano Jr, A. F., Serra, Y. L., & Adams, D. K.
1132 (2018). Convective-permitting hindcast simulations during the north ameri-
1133 can monsoon gps transect experiment 2013: Establishing baseline model
1134 performance without data assimilation. *Journal of Applied Meteorology and*

- 1135 *Climatology*, 57(8), 1683–1710.
- 1136 Morrison, H., & Milbrandt, J. A. (2015). Parameterization of cloud microphysics
1137 based on the prediction of bulk ice particle properties. part i: Scheme descrip-
1138 tion and idealized tests. *Journal of the Atmospheric Sciences*, 72(1), 287–311.
- 1139 Na, Y., Fu, Q., Leung, L. R., Kodama, C., & Lu, R. (2022). Mesoscale convective
1140 systems simulated by a high-resolution global nonhydrostatic model over
1141 the united states and china. *Journal of Geophysical Research: Atmospheres*,
1142 127(7), e2021JD035916.
- 1143 Nadolski, V. (1998). Automated surface observing system (asos) user’s guide.
1144 *National Oceanic and Atmospheric Administration, Department of Defense,*
1145 *Federal Aviation Administration, United States Navy*, 20. Retrieved from
1146 <https://www.ncei.noaa.gov/pub/data/asos-fivemin/6401-2012/>
- 1147 Nelson, B. R., Prat, O. P., Seo, D.-J., & Habib, E. (2016). Assessment and im-
1148 plications of ncep stage iv quantitative precipitation estimates for product
1149 intercomparisons. *Weather and Forecasting*, 31(2), 371–394.
- 1150 Peters, O., Neelin, J. D., & Nesbitt, S. W. (2009). Mesoscale convective systems and
1151 critical clusters. *Journal of the atmospheric sciences*, 66(9), 2913–2924.
- 1152 Pinto, J. O., Grim, J. A., & Steiner, M. (2015). Assessment of the high-resolution
1153 rapid refresh model’s ability to predict mesoscale convective systems using
1154 object-based evaluation. *Weather and Forecasting*, 30(4), 892–913.
- 1155 Prein, A., Gobiet, A., Suklitsch, M., Truhetz, H., Awan, N., Keuler, K., &
1156 Georgievski, G. (2013). Added value of convection permitting seasonal simula-
1157 tions. *Climate Dynamics*, 41(9), 2655–2677.
- 1158 Prein, A. F., Langhans, W., Fosser, G., Ferrone, A., Ban, N., Goergen, K., ... oth-
1159 ers (2015). A review on regional convection-permitting climate modeling:
1160 Demonstrations, prospects, and challenges. *Reviews of geophysics*, 53(2),
1161 323–361.
- 1162 Prein, A. F., Liu, C., Ikeda, K., Bullock, R., Rasmussen, R. M., Holland, G. J., &
1163 Clark, M. (2020). Simulating north american mesoscale convective systems
1164 with a convection-permitting climate model. *Climate Dynamics*, 55(1), 95–
1165 110.
- 1166 Przybylinski, R. W. (1995). The bow echo: Observations, numerical simulations, and
1167 severe weather detection methods. *Weather and Forecasting*, 10(2), 203–218.
- 1168 Rhoades, A. M., Ullrich, P. A., & Zarzycki, C. M. (2018). Projecting 21st century
1169 snowpack trends in western usa mountains using variable-resolution cesm. *Cli-*
1170 *mate Dynamics*, 50(1), 261–288.
- 1171 Roberts, N. (2008). Assessing the spatial and temporal variation in the skill of pre-
1172 cipitation forecasts from an nwp model. *Meteorological Applications: A journal*
1173 *of forecasting, practical applications, training techniques and modelling*, 15(1),
1174 163–169.
- 1175 Roberts, N. M., & Lean, H. W. (2008). Scale-selective verification of rainfall accu-
1176 mulations from high-resolution forecasts of convective events. *Monthly Weather*
1177 *Review*, 136(1), 78–97.
- 1178 Sadeghi, M., Nguyen, P., Naeini, M. R., Hsu, K., Braithwaite, D., & Sorooshian, S.
1179 (2021). Persiann-ccs-cdr, a 3-hourly 0.04° global precipitation climate data
1180 record for heavy precipitation studies. *Scientific Data*, 8(1), 1–11.
- 1181 Sakaguchi, K., Leung, L. R., Zhao, C., Yang, Q., Lu, J., Hagos, S., ... Lauritzen,
1182 P. H. (2015). Exploring a multiresolution approach using amip simulations.
1183 *Journal of Climate*, 28(14), 5549–5574.
- 1184 Sakaguchi, K., Lu, J., Leung, L. R., Zhao, C., Li, Y., & Hagos, S. (2016). Sources
1185 and pathways of the upscale effects on the southern hemisphere jet in mpas-
1186 cam4 variable-resolution simulations. *Journal of Advances in Modeling Earth*
1187 *Systems*, 8(4), 1786–1805.
- 1188 Salamanca, F., Krpo, A., Martilli, A., & Clappier, A. (2010). A new building energy
1189 model coupled with an urban canopy parameterization for urban climate simu-

- 1190 lations—part i. formulation, verification, and sensitivity analysis of the model.
1191 *Theoretical and applied climatology*, 99(3), 331–344.
- 1192 Santer, B., Wigley, T., & Jones, P. (1993). Correlation methods in fingerprint detec-
1193 tion studies. *Climate Dynamics*, 8(6), 265–276.
- 1194 Schenkman, A. D., & Xue, M. (2016). Bow-echo mesovortices: A review. *Atmo-
1195 spheric Research*, 170, 1–13.
- 1196 Schoen, J. M., & Ashley, W. S. (2011). A climatology of fatal convective wind events
1197 by storm type. *Weather and forecasting*, 26(1), 109–121.
- 1198 Schumacher, R. S., & Clark, A. J. (2014). Evaluation of ensemble configurations
1199 for the analysis and prediction of heavy-rain-producing mesoscale convective
1200 systems. *Monthly Weather Review*, 142(11), 4108–4138.
- 1201 Schumacher, R. S., & Johnson, R. H. (2005). Organization and environmental
1202 properties of extreme-rain-producing mesoscale convective systems. *Monthly
1203 weather review*, 133(4), 961–976.
- 1204 Schumacher, R. S., & Johnson, R. H. (2006). Characteristics of us extreme rain
1205 events during 1999–2003. *Weather and Forecasting*, 21(1), 69–85.
- 1206 Schumacher, R. S., & Rasmussen, K. L. (2020). The formation, character and chang-
1207 ing nature of mesoscale convective systems. *Nature Reviews Earth & Environ-
1208 ment*, 1(6), 300–314.
- 1209 Shepherd, T. J., Letson, F. L., Barthelmie, R. J., & Pryor, S. C. (2021). How well
1210 are hazards associated with derechos reproduced in regional climate simula-
1211 tions? *Natural Hazards and Earth System Sciences Discussions*, 1–42.
- 1212 Shourd, K. N. (2017). *The multi-scale dynamics of the 29-30 june 2012” super dere-
1213 cho”* (Unpublished doctoral dissertation). University of Nevada, Reno.
- 1214 Shourd, K. N., & Kaplan, M. L. (2021). The multiscale dynamics of the 29 june
1215 2012 super derecho. *Climate*, 9(11), 155.
- 1216 Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., ...
1217 others (2019). A description of the advanced research wrf model version 4.
1218 *National Center for Atmospheric Research: Boulder, CO, USA*, 145, 145.
- 1219 Squitieri, B. J., & Gallus, W. A. (2020). On the forecast sensitivity of mcs cold
1220 pools and related features to horizontal grid spacing in convection-allowing wrf
1221 simulations. *Weather and Forecasting*, 35(2), 325–346.
- 1222 Stensrud, D. J., Wicker, L. J., Xue, M., Dawson II, D. T., Yussouf, N., Wheatley,
1223 D. M., ... others (2013). Progress and challenges with warn-on-forecast.
1224 *Atmospheric Research*, 123, 2–16.
- 1225 Stevens, B., Satoh, M., Auger, L., Biercamp, J., Bretherton, C. S., Chen, X., ...
1226 others (2019). Dyamond: the dynamics of the atmospheric general circula-
1227 tion modeled on non-hydrostatic domains. *Progress in Earth and Planetary
1228 Science*, 6(1), 1–17.
- 1229 Stevenson, S. N., & Schumacher, R. S. (2014). A 10-year survey of extreme rain-
1230 fall events in the central and eastern united states using gridded multisensor
1231 precipitation analyses. *Monthly Weather Review*, 142(9), 3147–3162.
- 1232 Tang, Q., Klein, S. A., Xie, S., Lin, W., Golaz, J.-C., Roesler, E. L., ... others
1233 (2019). Regionally refined test bed in e3sm atmosphere model version 1
1234 (eamv1) and applications for high-resolution modeling. *Geoscientific Model
1235 Development*, 12(7), 2679–2706.
- 1236 Tao, W.-K., & Chern, J.-D. (2017). The impact of simulated mesoscale convective
1237 systems on global precipitation: A multiscale modeling study. *Journal of Ad-
1238 vances in Modeling Earth Systems*, 9(2), 790–809.
- 1239 Taylor, M. A., Guba, O., Steyer, A., Ullrich, P. A., Hall, D. M., & Eldrid, C. (2020).
1240 An energy consistent discretization of the nonhydrostatic equations in primi-
1241 tive variables. *Journal of Advances in Modeling Earth Systems*, 12(1).
- 1242 Thompson, G., Field, P. R., Rasmussen, R. M., & Hall, W. D. (2008). Explicit fore-
1243 casts of winter precipitation using an improved bulk microphysics scheme. part
1244 ii: Implementation of a new snow parameterization. *Monthly Weather Review*,

- 1245 136(12), 5095–5115.
- 1246 Toll, V., Männik, A., Luhamaa, A., & Rõõm, R. (2015). Hindcast experiments of
1247 the derecho in estonia on 08 august, 2010: Modelling derecho with nwp model
1248 harmonie. *Atmospheric Research*, 158, 179–191.
- 1249 Trenberth, K. E., Berry, J. C., & Buja, L. E. (1993). *Vertical interpolation and trun-*
1250 *cation of model-coordinate data*. Citeseer.
- 1251 Ullrich, P. A. (2014). *Squadgen: Spherical quadrilateral grid generator*. Retrieved
1252 from <https://climate.ucdavis.edu/squadgen.php>
- 1253 Ullrich, P. A., Devendran, D., & Johansen, H. (2016). Arbitrary-order conserva-
1254 tive and consistent remapping and a theory of linear maps: Part II. *Monthly*
1255 *Weather Review*, 144(4), 1529–1549.
- 1256 Ullrich, P. A., & Taylor, M. A. (2015). Arbitrary-order conservative and consis-
1257 tent remapping and a theory of linear maps: Part I. *Monthly Weather Review*,
1258 143(6), 2419–2440.
- 1259 Ullrich, P. A., & Zarzycki, C. M. (2017). TempestExtremes: A framework for scale-
1260 insensitive pointwise feature tracking on unstructured grids. *Geoscientific*
1261 *Model Development*, 10(3), 1069–1090.
- 1262 Ullrich, P. A., Zarzycki, C. M., McClenny, E. E., Pinheiro, M. C., Stansfield, A. M.,
1263 & Reed, K. A. (2021). TempestExtremes v2.1: a community framework
1264 for feature detection, tracking and analysis in large datasets. *Geoscientific*
1265 *Model Development Discussions*, 1–37. Retrieved from [https://github.com/](https://github.com/ClimateGlobalChange/tempestextremes)
1266 [ClimateGlobalChange/tempestextremes](https://github.com/ClimateGlobalChange/tempestextremes)
- 1267 Van Weverberg, K., Vogelmann, A., Lin, W., Luke, E., Cialella, A., Minnis, P., ...
1268 Jensen, M. (2013). The role of cloud microphysics parameterization in the sim-
1269 ulation of mesoscale convective system clouds and precipitation in the tropical
1270 western pacific. *Journal of the Atmospheric Sciences*, 70(4), 1104–1128.
- 1271 Vié, B., Nuissier, O., & Ducrocq, V. (2011). Cloud-resolving ensemble simulations
1272 of mediterranean heavy precipitating events: Uncertainty on initial conditions
1273 and lateral boundary conditions. *Monthly Weather Review*, 139(2), 403–423.
- 1274 Weisman, M. L., & Rotunno, R. (2004). “a theory for strong long-lived squall lines”
1275 revisited. *Journal of the Atmospheric Sciences*, 61(4), 361–382.
- 1276 Weisman, M. L., Skamarock, W. C., & Klemp, J. B. (1997). The resolution de-
1277 pendence of explicitly modeled convective systems. *Monthly Weather Review*,
1278 125(4), 527–548.
- 1279 Weyn, J. A., & Durran, D. R. (2017). The dependence of the predictability of
1280 mesoscale convective systems on the horizontal scale and amplitude of initial
1281 errors in idealized simulations. *Journal of the Atmospheric Sciences*, 74(7),
1282 2191–2210.
- 1283 Wheatley, D. M., Yussouf, N., & Stensrud, D. J. (2014). Ensemble kalman filter
1284 analyses and forecasts of a severe mesoscale convective system using different
1285 choices of microphysics schemes. *Monthly Weather Review*, 142(9), 3243–
1286 3263.
- 1287 Wu, C., Liu, X., Lin, Z., Rhoades, A. M., Ullrich, P. A., Zarzycki, C. M., ...
1288 Rahimi-Esfarjani, S. R. (2017). Exploring a variable-resolution approach
1289 for simulating regional climate in the rocky mountain region using the vr-cesm.
1290 *Journal of Geophysical Research: Atmospheres*, 122(20), 10–939.
- 1291 Xie, P., Joyce, R., Wu, S., Yoo, S.-H., Yarosh, Y., Sun, F., & Lin, R. (2017). Re-
1292 processed, bias-corrected cmorph global high-resolution precipitation estimates
1293 from 1998. *Journal of Hydrometeorology*, 18(6), 1617–1641.
- 1294 Xie, P., Joyce, R., Wu, S., Yoo, S.-H., Yarosh, Y., Sun, F., & Lin, R. (2019). *Noaa*
1295 *climate data record (cdr) of cpc morphing technique (cmorph) high resolution*
1296 *global precipitation estimates, version 1*. NOAA National Centers for Environ-
1297 mental Information. Retrieved from <https://doi.org/10.25921/w9va-q159>
- 1298 Xu, Z., Rhoades, A. M., Johansen, H., Ullrich, P. A., & Collins, W. D. (2018). An
1299 intercomparison of gcm and rcm dynamical downscaling for characterizing

- 1300 the hydroclimatology of california and nevada. *Journal of Hydrometeorology*,
1301 19(9), 1485–1506.
- 1302 Yang, G.-Y., & Slingo, J. (2001). The diurnal cycle in the tropics. *Monthly Weather*
1303 *Review*, 129(4), 784–801.
- 1304 Yuan, J., & Houze, R. A. (2010). Global variability of mesoscale convective system
1305 anvil structure from a-train satellite data. *Journal of Climate*, 23(21), 5864–
1306 5888.
- 1307 Zarzycki, C. M., & Jablonowski, C. (2014). A multidecadal simulation of atlantic
1308 tropical cyclones using a variable-resolution global atmospheric general circula-
1309 tion model. *Journal of Advances in Modeling Earth Systems*, 6(3), 805–828.
- 1310 Zarzycki, C. M., & Jablonowski, C. (2015). Experimental tropical cyclone forecasts
1311 using a variable-resolution global model. *Monthly Weather Review*, 143(10),
1312 4012–4037.
- 1313 Zhang, G. J., & McFarlane, N. A. (1995). Sensitivity of climate simulations to the
1314 parameterization of cumulus convection in the canadian climate centre general
1315 circulation model. *Atmosphere-ocean*, 33(3), 407–446.
- 1316 Zhang, S., Fu, H., Wu, L., Li, Y., Wang, H., Zeng, Y., ... others (2020). Optimizing
1317 high-resolution community earth system model on a heterogeneous many-core
1318 supercomputing platform. *Geoscientific Model Development*, 13(10), 4809–
1319 4829.
- 1320 Zhang, X., Anagnostou, E. N., & Schwartz, C. S. (2018). Nwp-based adjustment of
1321 imerg precipitation for flood-inducing complex terrain storms: Evaluation over
1322 conus. *Remote Sensing*, 10(4), 642.
- 1323 Zipser, E. J., & Lutz, K. R. (1994). The vertical profile of radar reflectivity of con-
1324 vective cells: A strong indicator of storm intensity and lightning probability?
1325 *Monthly Weather Review*, 122(8), 1751–1759.

Supporting Information for “The June 2012 North American Derecho: A testbed for evaluating regional and global climate modeling systems at cloud-resolving scales”

W. Liu¹, P.A. Ullrich¹, J. Li ², C. Zarzycki ³, P. M. Caldwell⁴, L. R. Leung ²,
Y. Qian ²

¹Department of Land, Air, and Water Resources, University of California-Davis, Davis, CA, USA

²Pacific Northwest National Laboratory, Richland, WA, USA

³Department of Meteorology and Atmospheric Science, Pennsylvania State University, University Park, PA, USA

⁴Lawrence Livermore National Lab, Livermore, CA, USA

Contents of this file

1. Figure S1
2. Figure S2
3. Figure S3
4. Figure S4

Corresponding author: Weiran Liu, Department of Land, Air, and Water Resources, University of California-Davis, Davis, CA, USA. (wraliu@ucdavis.edu)

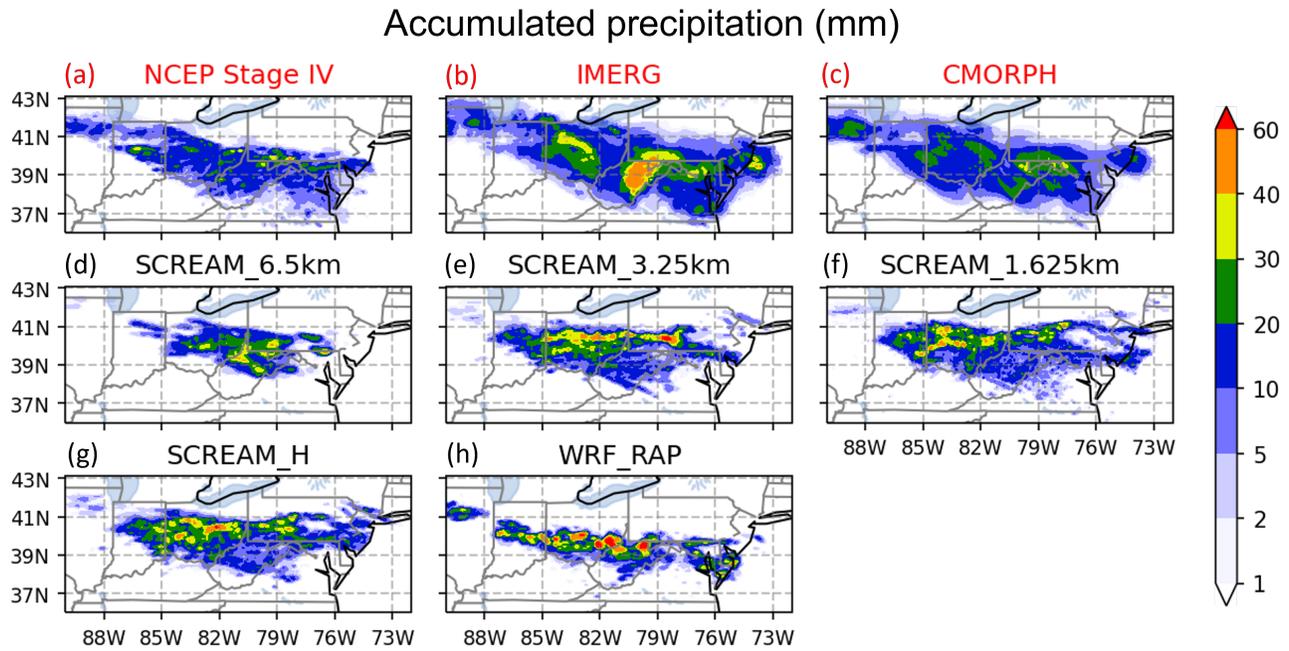


Figure S1. Same as Figure 11 but have a 2-hour shift for simulations resulting in the period changed to 20:00 UTC 29 June - 08:00 UTC 30 June 2012. Note that the time shift is applied only for simulation (i.e., panels a-e) and not applied for precipitation products (i.e., panels f-h).

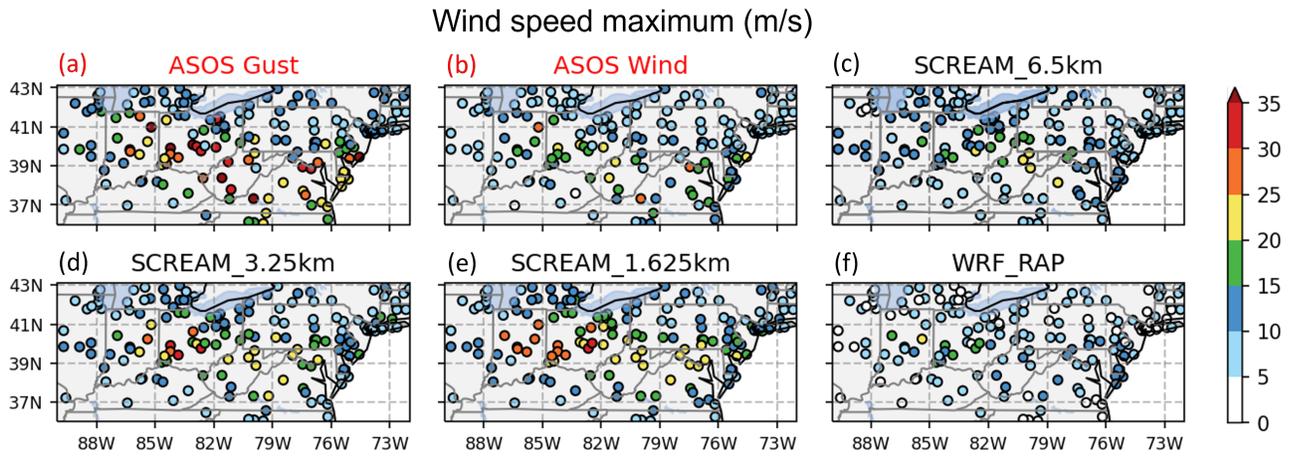


Figure S2. Same as Figure 12 but have a 2-hour shift for simulations resulting in the period changed to 20:00 UTC 29 June - 08:00 UTC 30 June 2012. Note that the time shift is applied only for simulation (i.e, panels a-e) and not applied for precipitation products (i.e., panels f-h).

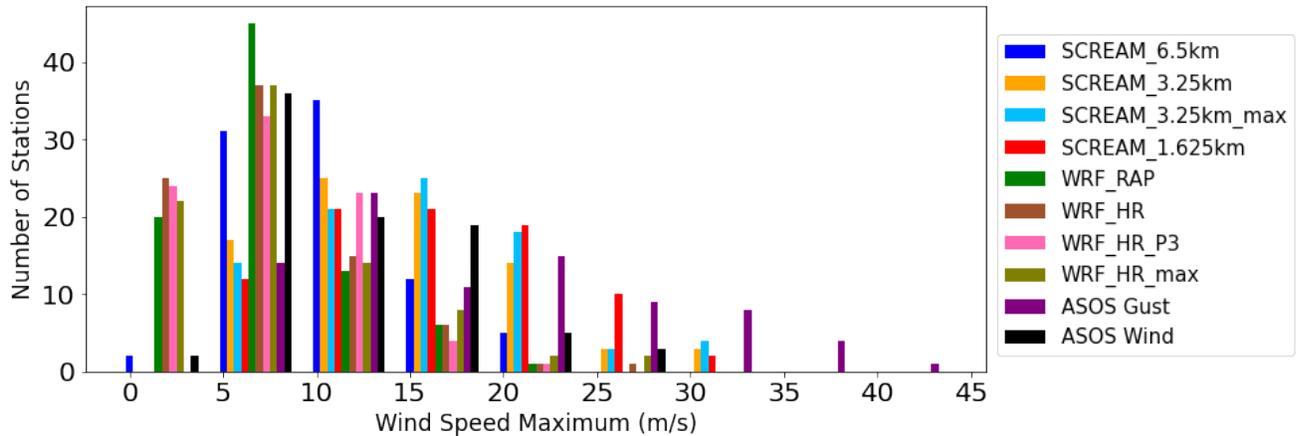


Figure S3. Same as Figure 13 but have a 2-hour shift for simulations resulting in the period changed to 20:00 UTC 29 June - 08:00 UTC 30 June 2012.

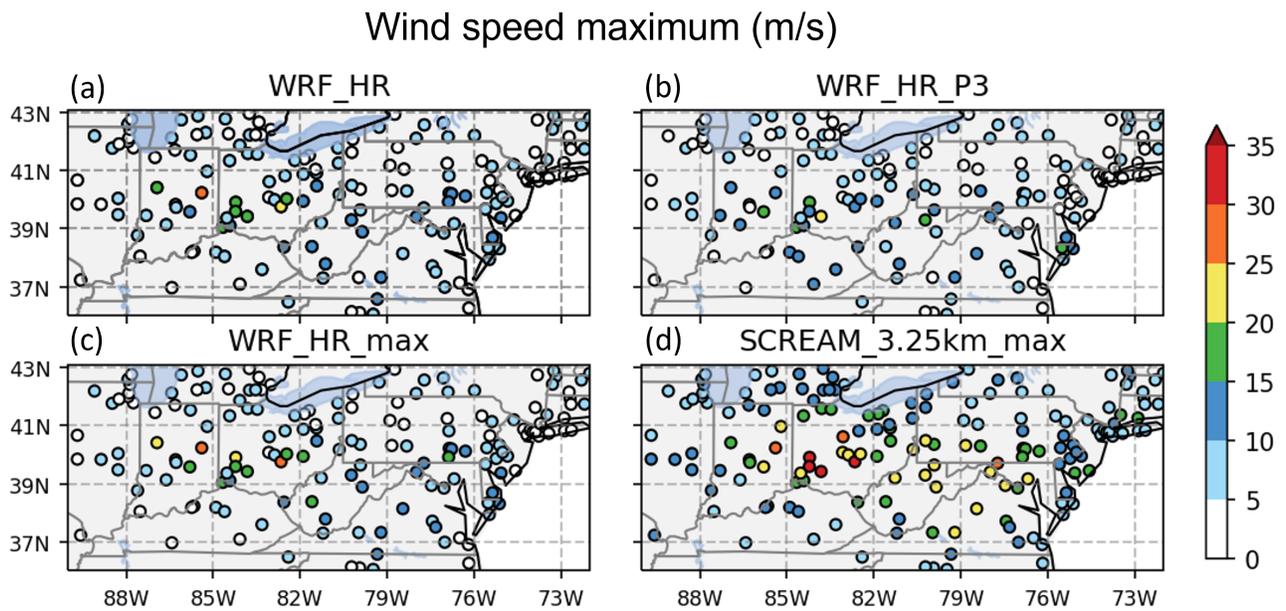


Figure S4. Same as Figure 14 but have a 2-hour shift for simulations resulting in the period changed to 20:00 UTC 29 June - 08:00 UTC 30 June 2012.