

Evaluation of the ability of the Weather Research and Forecasting model to reproduce a sub-daily extreme rainfall event in Beijing, China using different domain configurations and spin-up times

Qi Chu^{1,2,3}, Zongxue Xu^{1,2}, Yiheng Chen³, and Dawei Han³

¹College of Water Sciences, Beijing Normal University, Beijing, 100875, China

²Beijing Key Laboratory of Urban Hydrological Cycle and Sponge City, Beijing, 100875, China

³Department of Civil Engineering, University of Bristol, Bristol, BS8 1TR, UK

Correspondence to: Zongxue Xu (zongxuexu@vip.sina.com)

Abstract. The rainfall outputs from the latest convection-scale Weather Research and Forecasting (WRF) model are shown to provide an effective means of extending prediction lead times in flood forecasting. In this study, the performance of the WRF model in simulating a regional sub-daily extreme rainfall event centred over Beijing, China is evaluated at high temporal (sub-daily) and spatial (convective-resolving) scales using different domain configurations and spin-up times. Seven objective verification metrics that are calculated against the gridded ground observations and the ERA-Interim reanalysis are analysed jointly using subjective verification methods to identify the likely best WRF configurations. The rainfall simulations are found to be highly sensitive to the choice of domain size and spin-up time at the convective scale. A model run covering northern China with a 1:5:5 horizontal downscaling ratio (1.62 km), 57 vertical layers (less than 0.5 km), and a 60-hour spin-up time exhibits the best performance in terms of the accuracy of rainfall intensity and the spatial correlation coefficient (R'). A comparison of the optimal run and the initial run performed using the most common settings reveals clear improvements in the verification metrics. Specifically, R' increases from 0.226 to 0.67; the relative error of the maximum precipitation at a point rises from -56 % to -11.7 %; and the root mean squared error decreases by 33.65 %. In summary, re-evaluation of the domain configuration options and spin-up times used in WRF is crucial in improving the accuracy and reliability of rainfall outputs used in regional sub-daily heavy rainfall (SDHR)-related applications.

1 **1 Introduction**

2 The possibility that sub-daily heavy rainfall (SDHR) will increase with climate change is of significant societal concern.
3 SDHR-driven flash floods (FFs) are among the most destructive natural hazards that threaten many urban areas in northern
4 and central China and many other parts of the world. In these regions, SDHR is triggered mainly by regional mesoscale
5 circulation systems (MCSs) and occurs with increased intensity and frequency in warm seasons (Yu et al., 2007; Chen et al.,
6 2013). Records from the Emergency Events Database (EM-DAT) indicate that the damages and losses caused by FF events
7 in China have increased significantly over the past several decades. The risks are expected to continue to grow, given the
8 increase in the magnitude of SDHR predicted by most general circulation models (Chen et al., 2012; Willems et al., 2012;
9 Westra et al., 2014). The accelerating pace of urbanization also contributes to the increase in risk; urbanization has already
10 changed the hydrologic characteristics of the land surface considerably, resulting in higher peak flows and shorter flow
11 concentration times (Xu and Zhao, 2016; Gao et al., 2017). In such cases, very short-term (< 6-h) rainfall predictions are not
12 sufficient to provide adequate warning and mobilize emergency response activities. Recently developed statistically-based
13 rainfall generation methods and remote sensing data have been shown to enable the extension of the lead time to 24 hours.
14 However, this lead time is still insufficient to provide effective flood mitigation for medium or large urban areas with very
15 short hydrologic response times (Shih et al., 2014; Li et al., 2017). Therefore, numerical weather prediction (NWP), which
16 represents a means of forecasting heavy rainfall with lead times exceeding 24 h, has come into wide use in flood-related
17 studies and applications (Cuo et al., 2011).

18
19 Precipitation uncertainty accounts for a large proportion of the uncertainty in flood forecasts. Hence, given the large
20 uncertainties of NWP, its use in flood forecasting has long been questioned (Castelli, 1995; Bartholmes and Todini, 2005).
21 The ice wasn't broken until the end of the 20th century; substantial improvements in the predictive skill of NWP were made
22 that resulted from the increases in computational power and storage capacity, which enable parallel processing of
23 high-resolution forcing data and the resolution of convective-scale physical processes (Done et al., 2004; Clark et al., 2016).
24 The NWP models developed during and after this period can perform regional and convective-scale modelling and display
25 good performance in simulating heavy rainfall. Experimental studies have shown that NWP models of this kind, such as the
26 WRF model (Skamarock et al., 2008), tend to capture greater numbers of small-scale processes and the triggers of
27 convective storms (Klemp, 2006; Prein et al., 2015). Increasing numbers of meteorological operational centres and research
28 groups are adopting these new NWP models to carry out simulations of heavy rainfall events or real-time forecasting. The
29 resolutions of the rainfall products have improved from tens of kilometres to less than a kilometre, and the lead times have
30 increased from less than a day to more than a week (WMO, 2013). Meanwhile, case studies have been carried out using

1 regional convective-resolving models to evaluate the local rainfall predictions generated by sophisticated regional nesting
2 techniques or the global smooth grid transition approach on unstructured grids (Hong and Lee, 2009; Soares et al., 2012;
3 Sikder and Hossain, 2016; Heinzeller et al., 2016). The results of these studies demonstrate that, over relatively short periods
4 of time, regional modelling is often superior to large-scale modelling because it better resolves surface heterogeneities,
5 topography and small-scale features in air flow, such as growing instabilities (Miguez et al., 2004; En-Tao et al., 2010; Prein
6 et al., 2015; Brommel et al., 2015).

7
8 Despite the great potential of NWP models to predict heavy rainfall, a number of uncertainties remain that must be
9 considered. The errors induced by the initial and boundary conditions represent one source of these uncertainties; others stem
10 from cognitive errors and the scale effect in the solution of physical models, both of which may be exacerbated by the chaotic
11 nature of NWP. In regional simulations, these uncertainties are expected to be further magnified by downscaling or the use of
12 mesh transition procedures, so re-evaluation and calibration of the related model configurations are commonly required
13 (Warner, 2011; Vrac et al., 2012; Liu et al., 2012). As an example, running the WRF model at convective scales means that
14 convective processes are more likely to be resolved by explicit physical schemes than when sub-grid parameterizations are
15 used, which may incorporate new structural uncertainties related to the model physics (Done et al., 2004; Ruiz et al., 2010;
16 Crétat et al., 2012). In addition to model physics, several other aspects of model configuration, such as the domain size, the
17 spatial resolution and the spin-up time, may also have a substantial impact on the uncertainty of rainfall forecasts through
18 their effects on the initial and boundary conditions (Aligo et al., 2009; Fierro, 2009; Cuo et al., 2011). However, these aspects
19 of model configuration have received less attention in regional case studies because of their insignificant effects on rainfall
20 forecasts in coarse-resolution and long-term model simulations when compared to the physics of the WRF model. Generally,
21 these model configuration aspects are left at the common settings recommended by the official website of the WRF model
22 and by some experimental regional heavy rainfall studies.

23
24 Precipitation is one of the most sensitive variables to NWP model uncertainties. In this study, a re-evaluation of WRF is
25 performed to explore whether the recommended configuration of WRF represents the best choice in reproducing a regional
26 SDHR event happened in Beijing. The WRF model is assessed here because of its superior scalability and computational
27 efficiency; these traits are valued in interdisciplinary studies (Klemp, 2006; Foley et al., 2012; Coen et al., 2013; Yucel et al.,
28 2015). As the latest NWP community model, WRF incorporates up-to-date developments in physics, numerical methods and
29 data assimilation and is thus widely used in theoretical studies and practical applications (Powers et al., 2017). The selected
30 regional SDHR event occurred on July 21st, 2012 and was centred over Beijing, China. Beijing is among the most vulnerable
31 cities to SDHR-induced floods in central China (Yu et al., 2007). The precipitation in this area is caused mainly by monsoon
32 weather systems and enhanced by local orographic effects, and 60 %-80 % of the total annual precipitation occurs during just

1 a few SDHR events (Xu and Chu, 2015). The SDHR event that occurred on 21 July 2012 caused the most disastrous urban
2 flood in Beijing since 1950. The national operational NWP system failed to predict this event, which resulted in 79 deaths
3 and more than 1.6 billion dollars in damage (Wang et al., 2013; Zhou et al., 2013). Thus, several convective-scale studies
4 have been carried out to re-evaluate the optimal combination of the physics options used in the WRF model, such as Di et al.
5 (2015) and Wang et al. (2015). These studies represent the background information that has stimulates this research.

6
7 The second question we attempt to explore is to what extent rainfall simulations could be improved through the use of the
8 likely best set of settings if the recommended model configurations are not the best choices. The aspects of the model
9 configuration that are evaluated in this study are the domain size, vertical resolution, horizontal resolution and spin-up time.
10 These options have been found to have substantial impacts on daily-scale extreme rainfall outputs (Leduc and Laprise, 2009;
11 Aligo et al., 2009; Goswami et al., 2012). A comparative test with four scenarios is designed. Each scenario evaluates one
12 model configuration option to ensure that the simulated disparities can be attributed solely to a single factor each time. In
13 addition, the test is conceived as a progressive process: the optimal setting identified in each scenario will be adopted as the
14 primary choice for the next scenario to help quantify the overall improvement in the accuracy of rainfall outputs. The
15 ‘ground truth’ datasets are gridded observations obtained from Beijing Normal University and the China Meteorological
16 Centre. A coarser-resolution reanalysis called ERA-Interim (Dee et al., 2011) is also employed in identifying departures of
17 the WRF simulations from the driving weather fields as the model setup is varied. Seven objective verification metrics that
18 reflect different features of the model performance are adopted and considered jointly as part of a subjective verification
19 process because no single verification approach has been shown to provide comprehensive information about the quality of
20 rainfall simulations (Sikder and Hossain, 2016). Most of the metrics adopted here are those used to assess the performance of
21 WRF over daily or longer time periods (Liu et al., 2012; Tian et al., 2016). In this research, these metrics are calculated on an
22 hourly basis and averaged over different sub-daily time spans to evaluate the performance of the WRF model using different
23 configurations from a sub-daily and convective-scale perspective.

24 **2 Numerical Model Used to Forecast Heavy Rainfall**

25 The advanced WRF (ARW-WRF) model, version 3.7.1, is utilized as the dynamical downscaling tool. ARW-WRF is a
26 compressible non-hydrostatic and convection-permitting regional NWP model that employs the conservative form of the
27 dynamic Euler equations. As the latest regional NWP community system, WRF is composed of two dynamic cores, a data
28 assimilation system and a platform that facilitates parallel computation and function portability. Observations, model output
29 or assimilated reanalysis output can be used to initialize WRF. In terms of discretization, WRF uses a third-order
30 Runge-Kutta method for temporal separation and an Arakawa C-grid staggering scheme for spatial discretization. The model

1 is capable of conducting either one-way or two-way nested runs for regional downscaling. A detailed introduction to the
2 physics and numerical properties of ARW-WRF can be found in Skamarock et al. (2008). Given its emphasis on efficiency,
3 portability and updates to reflect the state of the art, WRF has been employed in settings ranging from research to
4 applications and has been incorporated into various operational systems, such as the Hurricane-WRF system for hurricane
5 forecasting and the WRF-Hydro system for hydrologic prediction.

6
7 In WRF, the domain size implicitly determines the large-scale dynamics and terrain effects, whereas the vertical and
8 horizontal grid spacings determine the smallest resolvable scale (Goswami et al., 2012). Together, these domain
9 configuration options affect the spectrum of the resolved scales and the nature of scale interactions in the model dynamics
10 (Leduc and Laprise, 2009). Thus, they are responsible for the generation and distribution of precipitation. In regional
11 simulations, small domain sizes are commonly preferred for computational efficiency. Seth and Rojas (2003) demonstrated
12 that simulations with small domain sizes are more likely to benefit from the lateral boundary conditions (LBCs) by
13 dampening the feedback from local perturbations on the large-scale general circulation. However, insufficiently large
14 domains have been shown to prevent the full development of small-scale features over areas of interest. To solve this issue,
15 the official website of WRF provides general guidance (Warner, 2011). This guidance recommends that the ranges of
16 domains should include the major features of the leading MCSs and local surface perturbations, and more than five grid
17 points should exist between adjacent nested domains to allow sufficient relaxation.

18
19 As for grid spacing, it appears plausible that WRF model runs performed with relatively small grid spacings would provide
20 more accurate outputs because such runs would resolve more small-scale phenomena of interest that are not present in the
21 LBCs. This statement is generally accepted as true when a relatively coarse-resolution run (>10 km horizontally or >1 km
22 vertically) is compared with a relatively finely resolved run at the convective scale (1-5 km horizontally or <1 km vertically)
23 in representing a convective storm. However, this conclusion is controversial when the comparison is conducted among
24 convective-scale model runs. Taking the horizontal resolution as an example, although there is evidence to show that WRF
25 runs performed at relatively high resolution capture more convective-scale features, the accuracy of rainfall outputs either
26 shows considerable or no statistical improvement (Roberts and Lean, 2008; Kain et al., 2008; Schwartz et al., 2009). In one
27 study, Fierro (2009) suggested that some features detected in convective-scale runs with too small horizontal grid spacings
28 tend to weaken the kinetic structures that favour torrential rainfall. A similar conclusion was drawn by Aligo et al. (2009) in
29 evaluating the impact of the vertical grid spacing on simulations of summer rainfall performed using WRF. Thus, horizontal
30 and vertical grid spacings of approximately 4 km and 1 km, respectively, have been employed as a reasonable compromise
31 between accuracy and computational efficiency in several regional studies.

32

1 In regional modelling, a spin-up period is often required to balance the inconsistencies between the results simulated by the
2 model physics and the initial and boundary conditions provided by the forcing data (Luna et al., 2013). The proper spin-up
3 time depends on the time needed for initialization, which can be affected by the size of the domain and the local boundary
4 perturbations (Warner, 1995; Kleczek et al., 2014). Moreover, the presence of chaotic behaviour, which causes reductions in
5 the predictive skill of models over time, imposes an upper bound on the spin-up time. Therefore, in cases where short
6 forecast lead times are expected, e.g., real-time rainfall forecasting, the spin-up time is mainly determined by the domain size
7 and the regional initial and boundary conditions. However, in cases where long forecast lead times are needed, e.g., warnings
8 of extreme rainfall, the effects of chaotic behaviour should be relatively evident. In practice, this issue is commonly
9 addressed by regularly updating the lateral boundary information derived from the latest forecasts or analyses to maintain
10 consistency between the regional model solutions and the atmospheric forcing conditions. In such cases, the best-fit
11 performance may occur for model runs with long spin-up times. Based on most previous studies, a spin-up time of 12 hours
12 is recommended to obtain an initial state; however, this spin-up time is often regarded as the suitable choice in many regional
13 case studies without further verification.

14 **3 Studied Event and Experimental Design**

15 As mentioned above, one aim of this study is to re-evaluate whether the recommended WRF domain configuration options
16 and spin-up time represent the optimal model configuration for reproducing a regional SDHR event when evaluated at a
17 sub-daily scale. Here, the SDHR event that occurred on 21 July 2012 and was centred on Beijing, China is selected as a case
18 study. The reasons why this event is selected, the synoptic and physical features that drove this event, and the model physics
19 adopted in this study are presented before the entire procedure of the experimental design is introduced.

20 **3.1 Study Event Selection and WRF Physical Schemes**

21 Beijing is selected as the study area because it is one of the most vulnerable cities to SDHR-induced FF hazards in China.
22 Beijing is located in central China. It has an area of 16 411 km², and its weather is mainly affected by the semi-humid warm
23 continental monsoon climate. The flows of air that favour local precipitation are cold, dry flows of air from high-latitude
24 areas to the north and hot, wet flows of air from the ocean to the south. The interactions between these two flows of air lead
25 to clear divergence in the temporal distribution of rainfall amount; 60 %-80 % of the annual precipitation occurs during just
26 a few heavy rainfall events during the warm season. Of all of the heavy rainfall events, the intensity and frequency of SDHR
27 events have been shown to display the greatest increasing tendencies over the past several decades. Meanwhile, Beijing, as
28 the capital of China, has experienced a significant expansion of its urban area and rapid increases in its population and
29 economic development. The negative effects of this expansion, such as losses of natural water bodies, increases in land cover
30 with low permeability and increases in urban drainage pipe networks, have led to continuous decreases in the hydrologic

1 response time. In addition, most of the population lives in the southwestern plain area. This region is downstream of
2 mountainous areas with steep terrain that varies in elevation from 60 m to 2 300 m (**Fig. 1**). All of these factors contribute to
3 the continuing increase in the exposure of this city to the high risks of flooding and waterlogging caused by SDHR events
4 (Xu and Chu, 2015).

5
6 **[Figure 1]**

7
8 The case study examines the largest heavy rainfall event that has occurred in Beijing in the past 65 years. The rainfall event
9 lasted for 16 hours (from 2 am to 6 pm) on 21 July 2012 (UTC), and the highest hourly rainfall intensity (100 mm/h) was
10 experienced in the southwestern part of the plain area. The associated FF hazard led to 79 deaths and damages totalling 1.6
11 billion US dollars, and more than 1.6 million people were affected. In addition to Beijing, the adjacent provinces, including
12 Hubei and Liaoning, were all significantly affected by this event and experienced severe FF hazards. The synoptic features
13 that triggered the rainfall were an eastward-moving vortex in the middle to high troposphere, a northward-moving zone of
14 subtropical high pressure and sharp vertical wind shear (Sun et al., 2013). The rainfall event as a whole can be divided into
15 two phases. From 2 am to 2 pm, the convective rain was dominated and enhanced by the orographic effect. The frontal rain
16 was then followed by the arrival of a cold front moving from the northwest until 6 pm (Guo et al., 2015). The rainfall
17 intensity in the second phase was relatively low compared to that in the first phase, due to the lack of strong kinetic forcing to
18 maintain the occurrence of precipitation.

19
20 The ERA-Interim reanalysis and 30-second static geographical data are employed to initialize the surface and meteorological
21 fields of the WRF. ERA-Interim is produced by an integrated forecasting system (IFS) used by the European Centre for
22 Medium-Range Weather Forecasts (ECMWF). The IFS is an Earth system model that incorporates a data assimilation
23 system and an atmospheric model that is fully coupled with land-surface and oceanic processes. The atmospheric model
24 provides output every 30 min at a spectral resolution of T255 (approximately 81 km over Beijing). This output is then
25 employed as prior information and combined with available observations twice a day to produce the reanalysis output using
26 the four-dimensional variation (4D-Var) assimilation system. The final reanalysis product, ERA-Interim, is a global gridded
27 dataset that is available at a spectral resolution of T255 and at both the 60 levels used in the model and 38 interpolated
28 pressure levels for all dates beginning on 1 January 1979 (Berrisford et al., 2009; Dee et al., 2011). Here, the ERA-Interim
29 pressure-level data are selected as the initial forcing. One reason is that, as is necessary, the vertical grid spacing between the
30 adjacent pressure layers is less than 1 km in the free troposphere, where the convective processes mainly occurred during the
31 Beijing SDHR event. In addition, the NWP models used by the Chinese Meteorological Centre mainly employ 31 vertical
32 levels in regional forecasting (WMO, 2013).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

As shown in **Fig. 2**, ERA-Interim captures the vortex and the subtropical high pressure well that occurred at the beginning of the rainfall event. In addition, the patterns of the leading MCSs and the primary synoptic features shown in this figure also correspond well to those described in previous studies (Zhou et al., 2013). The setup of the model physics is based mainly on the results of sensitive, high-resolution studies on the physics of the WRF model in simulating the same event (Wang et al., 2015; Di et al., 2015). The ‘resolved rain’ is driven by the single-moment 6-class microphysics scheme (Hong and Lim., 2006), whereas the ‘convective rain’ is resolved using the Grell-Devenyi cumulus parameterization scheme (Grell and Devenyi, 2002). The Noah land-surface model (Chen and Dudhia, 2001) is used and coupled with the Monin-Obukhov surface layer model (Ek et al., 2003). The radiation processes are represented by the RRTMG shortwave radiation and the RRTMG longwave radiation schemes (Iacono et al., 2008). For the planetary boundary layer scheme, the Yonsei University method (Hong et al., 2006) is adopted.

[Figure 2]

3.2 Experimental Design: Domain Configuration Options and Spin-Up Time

The comparative test is designed as a progressive process to help quantify the overall improvement in the performance of WRF after re-evaluating the WRF experiments performed using different domain configuration options and spin-up times. The test is classified into four successive scenarios. The first three scenarios investigate the domain configuration options, including the domain size, vertical resolution, and horizontal resolution; the fourth scenario concerns the spin-up time. During the entire procedure, the optimum configuration identified in each scenario is then adopted as the primary choice for the corresponding configuration in the following scenario. The initial datasets and the model physics are the same for all of the domains throughout the entire comparative procedure. Because the area of interest is located in the middle latitudes, the Lambert conformal projection is employed in all of the experiments, which is centred on the same latitude (42.25° N) and longitude (114.0° E). Moreover, sigma vertical coordinates with a top level of 50 hPa are used in all of the experiments.

Initially, the WRF domain configuration options and the spin-up time are set to the recommended values described in Sect. 2. Three levels of two-way nested domains are adopted so that the horizontal resolution in the smallest domain is sufficiently high to explicitly resolve convective-scale processes (**Fig. 1**). An odd downscaling ratio (1:3:3) is selected to reduce the initial error introduced by interpolating the initial fields to the assigned Arakawa grid. For the same reason, the boundaries of each domain are set along specific grid lines of the ERA-Interim dataset. Of the three nested domains, the outermost domain (D01) has the largest horizontal grid spacing of 40.5 km over north-central China, where the main perturbed synoptic

1 features occur. The innermost domain (D03) has the smallest horizontal grid spacing of nearly 4.5 km over the area of
2 interest, Beijing. The second domain (D02) is the child of D01 and the parent of D03 and has a horizontal grid spacing of
3 13.5 km. The distance between D01 and D02 is similar to that between D02 and D03, both of which exceed five grid points.
4 The grid numbers of D01, D02, and D03 are 40×40 , 72×72 and 90×90 , respectively. The eta values utilized in the initial
5 run are set based on the pressure values at the 29 vertical layers of the ERA-Interim pressure-level data. A spin-up time of
6 twelve hours (12 h) is selected; the outputs are saved every three hours in D03 and every hour in D02. The LBCs are updated
7 every six hours using ERA-Interim.

8
9 As shown in **Table 1**, the first experiment (C0) adopts the model configuration options mentioned above. To determine
10 whether the domain configuration options and the spin-up time used in C0 are the likely best set, four scenarios are designed.
11 The first scenario (S1) focuses on evaluating the effect of the WRF domain size. For computational efficiency, the MCS
12 systems that drive the local synoptic features are not completely contained within the outermost domain of C0, the
13 information of which is compensated by the updated LBCs from ERA-Interim. Two comparative experiments, C1 and C2,
14 are devised to verify that the domain size assigned to C0 is large enough to enable the full development of small-scale
15 features. Of the three experiments, C2 has the largest outermost domain size, which incorporates the leading MCS systems
16 over the entire Northeastern Hemisphere. The intermediate domain, which is centred between the outermost and innermost
17 domains, is then adopted as the outermost domain of C1. The purpose of scenario two (S2) is to evaluate whether the use of
18 a higher vertical resolution in a WRF model run results in better performance. In this scenario, the starting experiment is the
19 optimal experiment identified in S1 (OS1), forced by the ERA-Interim pressure-level data with 29 vertical levels. This
20 starting experiment is then followed by two experiments, C3 and C4, which incorporate one and two times more vertical
21 levels than OS1 (57 and 85 vertical levels), respectively. In the Beijing SDHR event, the pressure-level data meet the
22 requirement of a grid spacing of less than 1 km in the troposphere; however, this condition is not necessarily satisfied in other
23 regions. Thus, an experiment forced by the ERA-Interim model-level data with 38 vertical levels (C5) is also designed for
24 comparison. The three experiments (OS2, C6, and C7) in scenario three (S3) differ in terms of their horizontal resolutions
25 and nesting ratios, with increased nesting ratio of 1:3:3 (4.5 km grid spacing in D03), 1:5:5 (1.62 km in D03) and 1:7:7
26 (0.826 km in D03). The last scenario (S4) is designed to identify a reasonable optimal model run with the maximum spin-up
27 time after minimizing the uncertainties introduced by inappropriate domain configuration options. It contains one starting
28 experiment (OS3) and twelve comparative experiments (C8-C19). Except for C8, which includes no spin-up time, the
29 remaining experiments (C9-C19) include spin-up times that increase from 24 hours to 144 hours by every twelve hours.

30
31 [Table 1]
32

1 4 Verification Schemes

2 Both objective and subjective verification methods are applied to the innermost domain (D03) at a sub-daily scale. D03 is
3 selected because it covers the area of interest, Beijing, and the convective processes in this domain can be explicitly resolved
4 in all of the experiments. The rainfall data used for comparison in D03 are 3-hourly 0.05-degree data that were produced by
5 fusing rain gauge observations and the CMORPH data (Huang et al., 2013). The ERA-Interim reanalysis is utilized as well to
6 monitor the possible departures of the model simulations from the driving fields. Because the sub-daily rainfall is not
7 available from the reanalysis, the atmospheric precipitable water vapour (PW), which determines the possible maximum
8 precipitation, is instead compared with the model outputs every six hours. In addition, the model outputs that cover a larger
9 domain (D02) are compared with an hourly 0.1-degree gridded dataset obtained from the China Meteorological Centre (http://data.cma.cn/data/cdcdetail/dataCode/SEVP_CLI_CHN_MERGE_CMP_PRE_HOUR_GRID_0.10.html). The comparison
10 over domain two is used only as an auxiliary method for subjective verification, based on the assumption that an experiment
11 with good performance in the inner domain should also capture the large-scale features in the outer domain, as the
12 appropriate representation of these large-scale features will result in more accurate boundary conditions.
13

14
15 Seven error metrics that describe different features of precipitation are selected for use as objective verification metrics. Five
16 are rainfall-related and compared by bilinear interpolation of the output of the simulations to the grid of the ground truth data.
17 The accumulated areal rainfall is assessed using the relative error of the total precipitation (RE_{TP}). The percentage of correct
18 rainfall hits is measured using the probability of detection (POD) with a threshold of 0.1 mm. The root mean squared error
19 ($RMSE$) represents the amount of continuous error in the predicted precipitation. Detailed illustrations of these three metrics
20 can be found in Liu et al. (2012) and Tian et al. (2016). The other two rainfall-related metrics are the relative error of the
21 maximum grid precipitation (RE_{PMAX}) and the Pearson correlation coefficient (R), which describe the spatial association
22 between the simulations and the ground truth data (**Eq. (1) and Eq. (2)**). The two metrics selected for the verification of PW
23 (PW-related metrics) are the root mean squared error ($WRMSE$) and the Pearson correlation coefficient (WR). For
24 comparison, the PW fields of the reanalysis are remapped to the grids of the model outputs using the WRF Preprocessing
25 System (WPS). In this study, all of the metrics are calculated between the simulations and the reference data on the same grid
26 at each time step (3 h in D03). The values of these metrics are then averaged over four different time periods (6 h, 12 h, 18 h,
27 and 24 h) counted from 12 am on 21 July 2012. Different time periods are selected with the purpose of determining whether
28 the performance of WRF differs when the evaluation is conducted using different durations.
29

$$30 \quad RE = \frac{1}{N} \sum_{i=1}^N \left[\frac{f-r}{r} \times 100 \% \right] \quad (1)$$

$$R = \frac{1}{N} \sum_{i=1}^N \left(\frac{\sum_{j=1}^M (f_j - \bar{f})(r_j - \bar{r})}{\sqrt{\sum_{j=1}^M (f_j - \bar{f})^2 \sum_{j=1}^M (r_j - \bar{r})^2}} \right) \quad (2)$$

Here, R is the empirical spatial correlation coefficient; M is the total number of grid points within the evaluated domain of the starting experiment; f_j is the value of the j th grid point in the tested field at time step i ; r_j is the value of the reference field; N is the total number of time steps, depending on the time period considered; and RE is the relative error. For the maximum precipitation, f is the tested value of the maximum gridded precipitation over the area of interest, and r is the reference value of the maximum gridded precipitation over the same area.

To facilitate evaluation, the metrics are further adjusted to ensure that the ideal value of all of the metrics is 1. In this study, $RMSE$ and $WRMSE$ are first divided by a rescaling factor to fall into the range of 0-1 and then subtracted from 1 to provide an indication of good performance. The rescaled metrics, $RMSE'$ and $WRMSE'$, have the value 1 representing the lowest accumulated error (highest accuracy). The factor used for rescaling is determined by the largest values of each error metric in all of the experiments and is kept at the same value for all of the evaluated time periods (Sikder and Hossain, 2016). RE_{PMAX} and RE_{TP} are added by 1 to have the ideal value of 1. The rescaled metrics are $PMAX'$ and TP' , respectively. The other metrics are not rescaled because they already have ideal values of 1, but they are assigned a new set of symbols to distinguish them from the original metrics used before rescaling. For example, POD is replaced with POD' , and R is replaced with R' . **Table 2** shows the correlations between the original metrics and the rescaled metrics. Given that the metrics describe different features of the rainfall simulations, the values of these metrics are checked and considered together in subjective verification to determine the likely best set of domain configuration options and to search for the longest reasonable spin-up time.

[Table 2]

5 Results and Analyses

In each scenario, the metrics are compared among the experiments that consider different durations and cover the same domain (D03). The results are presented in four sub-graphs; each sub-graph shows the values of the metrics calculated for individual evaluated time periods. The spatial distribution of rainfall is also presented over domain two (D02) when evident discrepancies are noted in the results obtained for the inner domain (D03) and the outer domain (D02). **Table 1** shows the categories of the scenarios and the model configurations adopted in each experiment. In the following section, the domain

1 size scenario (S1) is evaluated first, followed by the vertical resolution scenario (S2) and the horizontal resolution scenario
2 (S3).

3

4 **5.1 Results of the Domain Size Scenario**

5 **Figure 3** shows the spatial values of the verification metrics for the WRF domain size experiments. The performance of the
6 experiments clearly worsens as the evaluated temporal duration increases from 6 h to 24 h. The most evident deteriorations
7 are detected in the point-to-point accuracy of the rainfall; the reversed root mean squared error ($RMSE'$) decreases by 0.8,
8 which represents a six-fold increase in the cumulative spatial error. The spatial association between the simulations and the
9 gridded observations also declines; the Pearson correlation coefficient (R') decreases by 0.3 on average. Although a slight
10 increase is observed in the percentage of correct hits (POD') during the first 18 hours, this increase is followed by a rapid
11 decrease of nearly 14 % during the last stage of the rainfall event. The relative bias in the accumulated areal rainfall (TP')
12 indicates that the total rainfall amount is underestimated throughout the entire evaluated temporal period. The maximum
13 gridded precipitation ($PMAX'$) is also underestimated; the largest negative bias occurs during the heavy convective rainfall
14 stage. For PW, a slight decrease is found in the reversed accumulated error ($WRMSE'$), whereas an increase of 5 %-9 % is
15 detected in the spatial correlation coefficient (WR'). Such variations may be attributable to the role of the updated boundary
16 conditions in adjusting the local model solutions to approach the large-scale atmospheric circulation conditions.

17

18 **[Figure 3]**

19

20 Comparison of the four sub-graphs shows that the values of the metrics do not point to a single perfect experiment in a given
21 period, and their ranked predictive skills determined using a given metric differ when evaluated over different time periods.
22 During the early stage of the rainfall event (6 h), C0 yields better performance than C1 and C2 in terms of $RMSE'$, R' and
23 $PMAX'$; it simultaneously displays the lowest value of POD' and the largest bias in estimating the total precipitation.
24 Although the superiority of C0 is more evident in the second period, a sharp deterioration is then observed in capturing the
25 point-to-point accuracy of precipitation for the 18-h duration, where the lowest R' is obtained. Meanwhile, C1, which
26 employs a domain of moderate size, displays greater skill than C0 in capturing the correct hits and the spatial pattern of the
27 simulated rainfall. C2 employs the largest domain. Although it shows the best fit to the rainfall observations on the daily
28 scale (24 h), it displays the worst performance over the three shorter time periods. For the PW fields, the highest similarity
29 with the ERA-Interim reanalysis is found for C0, whereas the lowest similarity is found for C2. These results demonstrate
30 indirectly that small domains are more likely to be influenced by updated boundary conditions.

1

2 In this scenario, if the experiments are merely evaluated in D03, the conclusion that C0 displays the best performance during
3 most of the evaluated time periods may be reached. However, when evaluated in D02, clear differences between C0 and the
4 ground truth in both the spatial characteristics of the rainfall and the magnitude of the maximum precipitation are detected.
5 **Figure 4** shows the spatial distribution of the accumulated six-hour precipitation over the domain two area of C0. Note that
6 the speed of movement of the belt of heavy rain simulated in C0 is a few kilometres per hour faster than those in C1 and C2,
7 leading to an early end of the heavy rainfall event. This difference may explain why the modelling skill of C0 declines
8 significantly as the end of the rainfall event approaches. The belt of heavy rain in C0 displays an orientation that is shifted
9 nearly ten degrees northward from those simulated in C1 and C2 during the first six hours, and the storm centre in C0
10 displays the smallest range; it is nearly half of the area in C2. The results indicate that the domain size of C0 is not broad
11 enough to allow the model physics to fully develop the small-scale features that favour heavy rainfall. The spatial
12 characteristics of precipitation are relatively similar in the other two experiments, but C1 outperforms C2 in both the
13 rainfall-related and the PW-related features over domain two. It may be that C2 does not yield better performance than C1
14 because of its inefficient use of boundary conditions to adjust the false perturbations generated by the local model run.
15 Therefore, C1 is verified as reasonable from both statistical and physical perspectives and is chosen as the optimal
16 experiment in the domain size scenario (OS1).

17

18

[Figure 4]

19

20 **5.2 Results of the Vertical Resolution Scenario**

21 Based on the analysed results, C1 is selected as the starting experiment in the vertical resolution scenario. As mentioned
22 above, C1 is forced with the ERA-Interim pressure-level data with 29 vertical levels. C3 and C4 are forced with the same
23 pressure-level data with 57 and 85 vertical levels, respectively, whereas C5 is forced with the model-level data with 38
24 vertical levels. As shown in **Fig. 5**, a decline in model performance is also obtained for all of the vertical resolution
25 experiments as the evaluated time period increases in length. Moreover, the largest deterioration in $RMSE'$ is also observed;
26 it decreases by 0.82 on average. The values of TP' and $PMAX'$ derived from the simulations are slightly higher than those
27 predicted in S1 but are still less than those calculated for the actual precipitation over the entire rainfall event. POD' displays
28 an evident decrease during the end stage of the rainfall event, and its magnitude decreases 50 % less relative to that shown in
29 C1. The most obvious difference from the domain size scenario is that the values of R' calculated between the simulations
30 and the ground truth vary slightly and remain almost the same between the different time periods. In addition, the

1 performance of the vertical resolution experiments seems to be less sensitive to the boundary conditions because they result
2 in relatively small variations in $WRMSE'$ and WR' .

4 [Figure 5]

5
6 Unlike the apparent discrepancies noted in the metrics obtained for the domain size experiments, the differences in the
7 rainfall-related metrics among the experiments with different numbers of vertical levels are not evident, especially during the
8 less rainy period (6 h) and the period when convective rainfall dominates (12 h). During the first 12 hours, C4 displays better
9 agreement with the gridded observations than the other three experiments in terms of the accuracy and spatial correlation of
10 the rainfall amount. However, over the longer time periods, C3 displays the greatest skill, according to most of the
11 verification metrics. Comparing C3 and C1 shows that increases in the vertical resolution may increase WRF's ability to
12 explicitly resolve small-scale physical processes and improve the accuracy of the amount and distribution of the simulated
13 rainfall. Comparing C3 and C4 shows that, although C4 include further refinement of the vertical resolution, its performance
14 is worse than that of C3 when the evaluated time period increases to more than 12 hours. This result may occur because
15 progressive reductions in the vertical grid spacing magnify the propagation of surface perturbations through the vertical grid
16 columns, potentially weakening the kinetic energy that favours precipitation. Examining the values of $WRMSE'$ and WR'
17 shows that the differences between the simulations and the reanalysis are more distinct in C3 and C4 than in C1. This
18 discrepancy may occur due to the exaggeration of the initial errors introduced by the interpolation process and the
19 incorporation of false surface perturbations introduced by the limited accuracy and resolution of the initial forcing data. C5
20 shows either better or worse performance than C1 in each period but produces less accurate rainfall simulations than C3 over
21 most of the evaluated durations. As such, C3 is identified as yielding the best performance in the vertical resolution scenario.

22

23 5.3 Results of the Horizontal Resolution Scenario

24 Based on the results obtained for scenario S2, C3 is selected as the starting experiment in the horizontal resolution scenario.
25 The modelling skill of the S3 experiments shows similar temporal trends as that of the S2 experiments (**Fig. 5** and **Fig. 6**).
26 However, the sensitivity of the metrics to the variation of the horizontal resolution is more evident than that with different
27 vertical resolutions. Over most of the evaluated time periods, C6, which has a grid spacing of 1.62 km, displays better
28 performance than C3 and C7 having grid spacings of 4.5 km and 0.826 km, respectively. Comparison of C3 and C6 shows
29 that C6 tends to produce more accurate spatial patterns of rainfall throughout the heavy rainfall event in Beijing. Higher
30 values of P_{MAX}' and TP' are also detected in C6 when compared to C3. This result stems in part from the explicit
31 resolution of the convective processes by the WRF microphysics scheme, which may explain why the P_{MAX}' of C7 is

1 higher than C6 over most of the tested durations. Note that the modelling skill of C7 deteriorates rapidly after the heavy rain
2 begins (12 h); the lowest POD' and R' values of the three experiments are obtained for this simulation and time period.
3 Analysis of the $WRMSE'$ values suggests that simulation C7 displays significant departures from the coarser-scale PW
4 fields that are used to force the model. Thus, model simulations with excessively high horizontal resolutions may also
5 display poor performance. Theoretically, this deterioration may be attributed to the accumulated errors introduced by the
6 imperfect model physics or biases in the initial and boundary conditions, which can be exaggerated by the chaotic nature of
7 NWP systems. According to the above analysis, C6 yields the best agreement with the ground truth data among the
8 horizontal resolution experiments.

9
10 **[Figure 6]**
11

12 **5.4 Searching for the Likely Ideal Spin-up Time**

13 To limit the effects of the chaotic nature of NWP on the model simulations and extend the lead time, the scenario in which the
14 spin-up time used in WRF is varied is placed at the end of the experimental design, after the possible errors introduced by
15 inappropriate domain configuration options have been reduced. In S4, C6 is adopted as the starting experiment (OS3).
16 Unlike the previous scenarios, the ranks of the spin-up time experiments, as sorted by the metrics, are nearly the same across
17 the different time periods. Hence, **Fig. 7** presents only the modelling skill of the spin-up time experiments over the time
18 period of 18 h. The model performance of WRF in simulating heavy rainfall clearly varies with the spin-up time. For most of
19 the metrics, an obvious diurnal tendency is found from 0 h to 60 h, followed by a short-term decrease until 72 h; random
20 fluctuations occur after 72 h. Before 72 h, the variations in the rainfall and PW metrics are almost consistent; thus, the good
21 fits of the simulations produced by the model runs with longer spin-up times are also physically reasonable within this period.
22 The discrepancies among these experiments may be due to differences in the initial conditions (e.g., the water vapour
23 amounts and the times of day when the simulations begin).

24
25 **[Figure 7]**
26

27 From TP' , it is found that all of the spin-up time experiments underestimate the total rainfall amount during the heavy
28 rainfall event. Of all of the rainfall-related metrics, POD' is found to display the least sensitivity to the spin-up time;
29 however, it displays similar variations over time as $PMAX'$, R' , and $RMSE'$ before 72 h, with the highest values shown in
30 the experiment with a spin-up time of 48 h (C11). Positive biases are detected in $PMAX'$ in C9 (which is run 24 hours ahead)
31 and C11, in which the largest positive biases are detected in the simulated amount of water vapour across the analysed

1 periods and earlier (during the initialization period). This result may occur because the atmospheric water vapour content
2 determines the maximum possible rainfall amount. C12, which includes a spin-up time of 60 h, is ranked third in terms of
3 $PMAX'$, whereas it displays better performance than C9 and C11 in terms of TP' , WR' , and $WRMSE'$. As seen in **Fig. 8**, C9,
4 C11 and C12 also rank in the top three, based on the values of the rainfall-related metrics calculated over domain two.
5 However, larger departures from the forcing PW fields are seen in C9 and C11 than in C12. The difference is that C12 shows
6 the best agreement with the ground truth data in terms of both the rainfall- and PW-related fields. Overall, C12 is regarded as
7 the experiment that best reproduces the Beijing SDHR event with the optimal set of domain configuration options and the
8 longest spin-up time.

9
10 **[Figure 8]**
11

12 **6 Discussion**

13 The results reveal that the initial experiment with the most commonly employed WRF domain settings does not yield the best
14 performance in reproducing the temporal and spatial characteristics of SDHR on the convective scale. In S1, the assigned
15 domain size of C0 is not sufficiently broad to allow the model physics to fully develop local small-scale features, resulting in
16 obvious reductions in modelling skill as the evaluated time duration increases from 12 h to 24 h. Further refinement of the
17 grid spacing of C0 in S2 and S3 is shown to enable more explicit resolution of convective processes, leading to more accurate
18 rainfall simulations. The comparison made in S4 suggests that the proper spin-up time is determined by both the time needed
19 for model initialization and the accuracy of the initial conditions fed into the model run. Moreover, experiments with too
20 large domains, too high spatial resolutions, or too long spin-up times also yield poor performance in rainfall simulations.
21 Therefore, the reasonableness of these WRF settings should be checked before the model is utilized in regional NWP
22 systems for flood forecasting or as a reference for the design of flood mitigation strategies.

23
24 In addition to exploring whether the recommended WRF domain configuration options and spin-up time are optimal for
25 application in SDHR-prone urban areas, the performance of the model is quantified, and its total improvement is evaluated
26 by comparing the values of the verification metrics yielded by the experiments. **Table 3** compares the values of the
27 verification metrics obtained for the optimal experiments in each scenario with the values obtained for the initial experiment.
28 Here, the 18 h time duration is selected for evaluation because it covers most of the heavy rainfall event, and the metrics
29 calculated over this period display a greater range and thus greater ability in identifying the simulation with the best
30 performance. One exception is the domain size scenario, in which C0 presents the most obvious reduction in performance

1 during the last stage of the rainfall event (24 h). Therefore, the improvement in $C1$ relative to $C0$ is mainly represented by R'
2 and POD' across D03 over the 18-h time period. The improvement produced by refining the vertical resolution is indicated
3 by all of the rainfall-related metrics but is accompanied by a decrease in $WRMSE'$ that stems in part from the reduction in
4 kinetic energy, which promotes rainfall. $C6$ yields higher values of POD' , $RMSE'$, R' , and $PMAX'$ when compared with $C3$,
5 indicating that appropriate increases in the horizontal resolution can increase the accuracy of rainfall simulations. The largest
6 differences in the metrics between $C6$ and $C12$ occurs for $PMAX'$, which may relate to the different initial weather
7 conditions at the different starting times of the model runs.

8
9 **[Table 3]**

10

11 Overall, although the magnitudes of the increases in the rainfall metrics differ, they all reflect an increase in model skill after
12 the re-evaluation process has been conducted. Specifically, R' increases from 0.226 in $C0$ to 0.67 in $C12$; $RMSE'$ increases
13 from 0.098 to 0.402; and $PMAX'$ increases from 0.44 to 0.883. As the complete assessment is based on objective
14 verification metrics and checked by subjective verification methods, it can be concluded that the domain configuration
15 options and the spin-up time have significant effects on regional simulations of SDHR. Therefore, re-evaluating the values of
16 those settings used in high-resolution regional studies is certainly worthwhile, and the accuracy of predictions of heavy rain
17 clearly benefit from these analyses. For the evaluated metrics, evaluations based on a single type of metric or a single time
18 period may clearly result in partially accurate conclusions. The use of datasets from multiple sources in verification can help
19 increase the comprehensiveness of the analyses, such as the use of $WRMSE'$ and WR' in this study. The use of different
20 time periods helps to determine the optimal configurations with higher physical rationality, such as the selection of the
21 proper domain size. In addition, the verification results may also depend on the fields and temporal-spatial scales of interest.
22 To further understand the effects of WRF model configuration options on regional simulations of sub-daily heavy rainfall,
23 more objective verification metrics for SDHR should be developed, and more case studies of SDHR events are also needed.
24 Given that the uncertainties in the regional NWP studies result mainly from the inaccurate boundary conditions associated
25 with grid nesting techniques, methods that can serve as alternate schemes to reduce these uncertainties are also worth
26 studying. Examples include the mesh transitions approach used on irregular grids. In addition, more accurate simulations are
27 expected when the model is driven with forcing data with higher temporal or spatial resolutions than those of the
28 ERA-Interim reanalysis because the uncertainties and errors introduced by the input data could be further reduced.

29

1 7 Conclusions

2 In this study, a comparative test is designed to evaluate the effects of WRF domain configuration options and the spin-up
3 time on simulations of the precipitation during the SDHR event that occurred on July 21st, 2012 in Beijing, China. Three
4 nested domains are established; D01 is the largest, has the coarsest resolution, and covers the leading synoptic features, and
5 D03 is the smallest and covers the area of interest, Beijing. The initial conditions of the three domains are provided by the
6 ERA-Interim reanalysis and the 30-second static geographical datasets. For the LBCs, D01 is forced by the ERA-Interim
7 reanalysis, whereas D02 is forced by D01, and D03 is forced by D02. The reference ground truth data used for verification is
8 3-hourly 0.05 gridded rainfall observations and the coarser-scale ERA-Interim reanalysis. Five rainfall-related error metrics
9 and two PW-related indices that monitor the departure of the model simulations from the driving fields are calculated at the
10 convective-resolving scale over different sub-daily time spans. These metrics are then checked and considered together as
11 part of a subjective verification process that is intended to pinpoint the likely best combination of the domain configuration
12 options and spin-up time and to help quantify the possible improvements in the model performance of WRF in reproducing
13 severe SDHR events after carrying out the entire re-evaluation process.

14
15 Precipitation simulations are sensitive to changes in domain size, vertical resolution, horizontal resolution and spin-up time.
16 Of all of the configurations, the most obvious variations are found when adjusting the domain size and the spin-up time. This
17 analysis shows that domains that cover only the area of interest may be insufficiently broad to permit full development of
18 small-scale features, resulting in poor performance in capturing the spatial pattern of heavy rainfall, especially in the early
19 stages of rainfall events. Despite the dominant role of chaotic processes, it is still possible that model runs with longer
20 spin-up times may result in better rainfall simulations, given favourable initial weather conditions. The effects of the vertical
21 and horizontal resolutions are smaller, but the accuracy of the rainfall amount and the correct hits exhibit evident increases in
22 runs with slightly higher spatial resolutions. A comparison of C12, which uses the evaluated optimum configurations, and C0,
23 which uses the recommended settings, shows that the metrics clearly increase. Specifically; R' increases from 0.226 to 0.67;
24 $RE_{P_{MAX}}$ rises from -56 % to -11.7 %; and $RMSE$ decreases by 33.65 %. Thus, substantial benefits may result from
25 re-evaluating the WRF domain configuration options and spin-up times used in regional studies of SDHR.

26
27 Given the intensification of SDHR and the increased risks posed by SDHR-induced hazards, the demands of the operational
28 flood management community for more accurate rainfall predictions with longer lead times, especially over highly affected
29 areas with very short hydrologic response times, are increasing. One method that has now been proven to be effective is to
30 dynamically downscale freely available global NWP products to areas of interest using high-resolution regional NWP
31 models (e.g., WRF). Therefore, the uncertainties associated with the downscaling process, such as errors in boundary

1 conditions and the issues associated with grid nesting, should be carefully evaluated to ensure that the rainfall simulations
2 produced are both statistically accurate and physically reasonable before they are employed in flood forecasting systems.
3 This study illustrates the importance of re-evaluating the domain configuration options and spin-up times used in WRF in
4 improving regional rainfall simulations. Comparisons of the metrics indicate that evaluations based on just one category of
5 metrics or values of metrics calculated over only one time period (e.g., 24 h) do not result in comprehensive comparisons and
6 may lead to partially accurate conclusions. The use of PW fields calculated against reanalysis output is verified to be helpful
7 in determining the optimal set of model configurations when analyses of rainfall-related metrics do not yield uniform
8 conclusions. In addition, evaluations conducted over larger-scale domains are demonstrated to have utility in establishing the
9 reasonableness of the evaluated results. Overall, the evaluation process is partly subjective. To simplify the assessment
10 process, verification methods that can replace this subjective verification procedure should be developed. More regional case
11 studies are also needed to further investigate the effects of configuration options in simulations of regional SDHR and to
12 explore methods of reducing the uncertainties in regional NWP modelling associated with the scale-variation procedures. In
13 addition, the use of more accurate forcing data with higher temporal and spatial resolutions is also expected to reduce the
14 errors in the initial and boundary conditions and could thus be helpful in further improving the accuracy of rainfall
15 simulations and extending the lead times of forecasts.

16

17 *Data availability:* The ERA-Interim reanalysis dataset used as the initial forcing in the text is freely available at
18 <https://www.ecmwf.int/en/forecasts/datasets/archive-datasets/reanalysis-datasets/era-interim>.

19

20 *Competing interests:* The authors declare that they have no conflict of interest.

21

22 *Acknowledgement:* This study is supported by the key research projects “Sponge city construction and urban
23 flooding/waterlogging disaster in the sub-center of Beijing City” (Z171100002217080), Beijing Municipal Science and
24 Technology Commission, and “Urban storm flooding/waterlogging disasters under changing environment”
25 (2017YFC1502701), national key research and development plan, Ministry of Science and Technology, China. Support is
26 also received from the Resilient Economy and Society by Integrated Systems modeling (RESIST), Newton Fund via Natural
27 Environment Research Council (NERC) and Economic and Social Research Council (ESRC) (NE/N012143/1), and the
28 National Natural Science Foundation of China (No: 4151101234). The China Scholarship Council supports the first author
29 for her academic visit to the University of Bristol, UK.

30

31 Edited by: Uwe Ehret

32 Reviewed by: two anonymous referees

1 **References**

- 2 Aligo, E. A., Gallus Jr., W. A., and Segal, M.: On the impact of WRF model vertical grid resolution on Midwest summer
3 rainfall forecasts. *Weather and Forecasting*, **24**, 575-594, 2009.
- 4 Bartholmes, J., and Todini, E.: Coupling meteorological and hydrological models for flood forecasting. *Hydrol. Earth Syst.*
5 *Sci.: Discussions*, **9(4)**, 333-346, 2005.
- 6 Berrisford, P., Dee, D. P., Fielding, K., Fuentes, M., Kallberg, P., Kobayashi, S., and Uppala, S. M.: The ERA-Interim
7 Archive. *ERA Report Series*, **1**, 1-16, 2009.
- 8 Castelli, F.: Atmosphere modeling and hydrologic-prediction uncertainty. U.S. - Italy Research Workshop on the
9 Hydrometeorology, impacts and management of extreme floods, Perugia, 1995.
- 10 Chen, F., and Dudhia, J.: Coupling an advanced land surface-hydrology model with the Penn State-NCAR MM5 modeling
11 system. Part I: Model implementation and sensitivity. *Mon. Weather Rev.*, **129(4)**, 569-585, 2001.
- 12 Chen, H., Sun, J., Chen, X., and Zhou, W.: CGCM projections of heavy rainfall events in China. *Int. J. Climatol.*, **32(3)**,
13 441-450, 2012.
- 14 Clark, P., Roberts, N., Lean, H., Ballard, S. P., and Charlton-Perez, C.: Convection-permitting models: a step-change in
15 rainfall forecasting. *Meteor. Appl.*, **23(2)**, 165-181, 2016.
- 16 Coen, J. L., Cameron, M., Michalakes, J., Patton, E. G., Riggan, P. J., and Yedinak, K. M.: WRF-Fire: coupled
17 weather-wildland fire modeling with the weather research and forecasting model. *J. Appl. Meteor. Climatol.*, **52(1)**, 16-38,
18 2013.
- 19 Crétat, J., Pohl, B., Richard, Y., and Drobinski, P.: Uncertainties in simulating regional climate of Southern Africa: sensitivity
20 to physical parameterizations using WRF. *Clim. Dyn.*, **38(3-4)**, 613-634, 2012.
- 21 Cuo, L., Pagano, T. C., and Wang, Q. J.: A review of quantitative precipitation forecasts and their use in short-to
22 medium-range streamflow forecasting. *J. Hydrometeor.*, **12(5)**, 713-728, 2011.
- 23 Dee, D. P., and Coauthors: The ERA-Interim reanalysis: configuration and performance of the data assimilation system. *Q. J.*
24 *R. Meteorol. Soc.*, **137(656)**, 553-597, 2011.
- 25 Di, Z. H., and Coauthors: Assessing WRF model parameter sensitivity: A case study with five-day summer precipitation
26 forecasting in the Greater Beijing Area. *Geophys. Res. Lett.*, **42**, 579-587, 2015.
- 27 Done, J., Davis, C. A., and Weisman, M.: The next generation of NWP: Explicit forecasts of convection using the Weather
28 Research and Forecasting (WRF) model. *Atmos. Sci. Lett.*, **5(6)**, 110-117, 2004.
- 29 En-Tao, Y. U., Hui-Jun, W. A. N. G., and Jian-Qi, S. U. N.: A quick report on a dynamical downscaling simulation over
30 China using the nested model. *Atmos. Oceanic Sci. Lett.*, **3(6)**, 325-329, 2010.

1 Ek, M. B., Mitchell, K. E., Lin, Y., Rogers, E., Grunmann, P., Koren, V., Gayno, G., and Tarpley, J. D.: Implementation of
2 Noah land surface model advances in the National Centers for Environmental Prediction operational mesoscale Eta model.
3 *J. Geophys. Res.: Atmos.*, **108(D22)**, 2003.

4 Fierro, A. O., Rogers, R. F., Marks, F. D., and Nolan, D. S.: The impact of horizontal grid spacing on the microphysical and
5 kinematic structures of strong tropical cyclones simulated with the WRF-ARW model. *Mon. Weather Rev.*, **137(11)**,
6 3717-3743, 2009.

7 Foley, A. M., Leahy, P. G., Marvuglia, A., and McKeogh, E. J.: Current methods and advances in forecasting of wind power
8 generation. *Renewable Energy*, **37(1)**, 1-8, 2012.

9 Gao, Y., Yuan, Y., Wang, H., Schmidt, A. R., Wang, K., and Ye, L.: Examining the effects of urban agglomeration polders on
10 flood events in Qinhuai River basin, China with HEC-HMS model. *Water Sci. Technol.*, **75(9)**, 2130-2138, 2017.

11 Goswami, P., Shivappa, H., and Goud, S.: Comparative analysis of the role of domain size, horizontal resolution and initial
12 conditions in the simulation of tropical heavy rainfall events. *Meteor. Appl.*, **19(2)**, 170-178, 2012.

13 Grell, G. A., and Dévényi, D.: A generalized approach to parameterizing convection combining ensemble and data
14 assimilation techniques. *Geophys. Res. Lett.*, **29(14)**, 38-31, 2002.

15 Guo, C., Xiao, H., Yang, H., and Tang, Q.: Observation and modeling analyses of the macro-and microphysical
16 characteristics of a heavy rain storm in Beijing. *Atmospheric Research*, **156**, 125-141, 2015.

17 Heinzeller, D., Duda, M. G. and Kunstmann, H.: Towards convection-resolving, global atmospheric simulations with the
18 Model for Prediction Across Scales (MPAS) v3. 1: an extreme scaling experiment. *Geosci. Model Dev.*, **9(1)**, 77, 2016.

19 Hong, S. Y., and Lee, J. W.: Assessment of the WRF model in reproducing a flash-flood heavy rainfall event over Korea.
20 *Atmos. Res.*, **93(4)**, 818-831, 2009.

21 Hong, S. Y., and Lim, J. O. J.: The WRF single-moment 6-class microphysics scheme (WSM6). *J. Korean Meteor. Soc.*,
22 **42(2)**, 129-151, 2006.

23 Hong, S. Y., Noh, Y., and Dudhia, J.: A new vertical diffusion package with an explicit treatment of entrainment processes.
24 *Mon. Weather Rev.*, **134(9)**, 2318-2341, 2006.

25 Huang, C., Zheng, X., Tait, A., Dai, Y., Yang, C., Chen, Z., Li, T., and Wang Z.: On using smoothing spline and residual
26 correction to fuse rain gauge observations and remote sensing data. *J. Hydrol.*, **508**, 410-417, 2013.

27 Kain, J. S., and Coauthors: Some practical considerations regarding horizontal resolution in the first generation of
28 operational convection-allowing NWP. *Weather and Forecasting*, **23(5)**, 931-952, 2008.

29 Kleczek, M. A., Steeneveld, G. J., and Holtslag, A. A.: Evaluation of the weather research and forecasting mesoscale model
30 for GABLS3: impact of boundary-layer schemes, boundary conditions and spin-up. *Boundary-layer meteor.*, **152(2)**,
31 213-243, 2014.

32 Klemp, J. B.: Advances in the WRF model for convection-resolving forecasting. *Adv. Geosci.*, **7**, 25-29, 2006.

1 Leduc, M., and Laprise, R.: Regional climate model sensitivity to domain size. *Clim. Dyn.*, **32(6)**, 833-854, 2009.

2 Liu, J., Bray, M., and Han, D.: Sensitivity of the Weather Research and Forecasting (WRF) model to downscaling ratios
3 and storm types in rainfall simulation. *Hydrol. Processes*, 26(20), 3012-3031, 2012.

4 Li, J., Chen, Y., Wang, H., Qin, J., Li, J., and Chiao, S.: Extending flood forecasting lead time in a large watershed by
5 coupling WRF QPF with a distributed hydrological model. *Hydrol. Earth Syst. Sci.*, **21(2)**, 1279, 2017.

6 Luna, T., Castanheira, M., and Rocha, A.: Assessment of WRF-ARW forecasts using warm initializations. 2013. [Available
7 online at http://climetua.fis.ua.pt/publicacoes/APMG_extended_abstract_2013_Luna_et_al.pdf]

8 Miguez-Macho, G., Stenchikov, G. L., and Robock, A.: Spectral nudging to eliminate the effects of domain position and
9 geometry in regional climate model simulations. *J. Geophys. Res.: Atmos.*, **109(D13)**, 2004.

10 Mlawer, E. J., and Clough, S. A.: Shortwave and longwave enhancements in the rapid radiative transfer model. *Proceedings*
11 *of the 7th Atmospheric Radiation Measurement (ARM) Science Team Meeting*, 499-504, 1998.

12 Mlawer, E. J., Taubman, S. J., Brown, P. D., Iacono, M. J., and Clough, S. A.: Radiative transfer for inhomogeneous
13 atmospheres: RRTM, a validated correlated-k model for the longwave. *J. Geophys. Res.: Atmos.*, **102(D14)**, 16663-16682,
14 1997.

15 Prein, A. F., and Coauthors: A review on regional convection-permitting climate modeling: Demonstrations, prospects, and
16 challenges. *Rev. Geophys.*, **53(2)**, 323-361, 2015.

17 Powers, J. G., and Coauthors: The Weather Research and Forecasting (WRF) Model: Overview, System Efforts, and Future
18 Directions. *Bull. Amer. Meteor. Soc.*, 2017.

19 Roberts, N. M., and Lean, H. W.: Scale-selective verification of rainfall accumulations from high-resolution forecasts of
20 convective events. *Mon. Weather Rev.*, **136(1)**, 78-97, 2008.

21 Ruiz, J. J., Saulo, C., and Nogués-Paegle, J.: WRF model sensitivity to choice of parameterization over South America:
22 validation against surface variables. *Mon. Weather Rev.*, **138(8)**, 3342-3355, 2010.

23 Schwartz, C. S., and Coauthors: Next-day convection-allowing WRF model guidance: A second look at 2-km versus 4-km
24 grid spacing. *Mon. Weather Rev.*, **137(10)**, 3351-3372, 2009.

25 Seth, A., and Rojas, M.: Simulation and sensitivity in a nested modeling system for South America. Part I: Reanalyses
26 boundary forcing. *J. Clim.*, **16(15)**, 2437-2453, 2003.

27 Shih, D. S., Chen, C. H., and Yeh, G. T.: Improving our understanding of flood forecasting using earlier
28 hydro-meteorological intelligence. *J. Hydrol.*, **512**, 470-481, 2014.

29 Sikder, S., and Hossain, F.: Assessment of the weather research and forecasting model generalized parameterization schemes
30 for advancement of precipitation forecasting in monsoon-driven river basins. *J. Adv. Modeling Earth Syst.*, **8(3)**,
31 1210-1228, 2016.

1 Skamarock, W. C., and Coauthors: A description of the advanced research WRF Ver. 30, NCAR Technical Note.
2 NCAR/TN-475, 2008.

3 Soares, P. M., Cardoso, R. M., Miranda, P. M., de Medeiros, J., Belo-Pereira, M., and Espirito-Santo, F.: WRF high
4 resolution dynamical downscaling of ERA-Interim for Portugal. *Clim. Dyn.*, **39(9-10)**, 2497-2522, 2012.

5 Sun M. S., Yang L. Q., Yin Q., Niu Z. Y., and Gao L. M.: Analysis of the cause of a torrential rain occurring in Beijing on 21
6 July 2012(II). *Torrential Rain and Disasters (in Chinese)*, **32(3)**, 218-223, 2013.

7 Swinbank, R. and James Purser, R.: Fibonacci grids: A novel approach to global modeling. *Q. J. R. Meteorol. Soc.*,
8 **132(619)**, 1769-1793, 2006.

9 Tian, J. Y., Liu, J., Li, C. Z., and Yu, F. L.: Numerical rainfall simulation with different spatial and temporal evenness by
10 using WRF multi-physics ensembles. *Nat. Hazards Earth Syst. Sci.*, **17(4)**, 563-579, 2017.

11 Vrac, M., Drobinski, P., Merlo, A., Herrmann, M., Lavaysse, C., Li, L., and Somot, S.: Dynamical and statistical
12 downscaling of the French Mediterranean climate: uncertainty assessment. *Nat. Hazards Earth Syst. Sci.*, **12(9)**, 2769,
13 2012.

14 Wang, K., Wang, L., Wei, Y. M., and Ye, M.: Beijing storm of July 21, 2012: observations and reflections. *Nat. hazards*,
15 **67(2)**, 969-974, 2013.

16 Wang S. L., Kang H. W., Gu X. Q., and Ni Y. Q.: Numerical Simulation of Mesoscale Convective System in the Warm
17 Sector of Beijing '7.21' Severe Rainstorm. *Meteor. Mon.*, **41(5)**, 544-553, 2015.

18 Warner, T. T., Peterson, R. A., and Treadon, R. E.: A tutorial on lateral boundary conditions as a basic and potentially serious
19 limitation to regional numerical weather prediction. *Bull. Amer. Meteor. Soc.*, **78(11)**, 2599, 1997.

20 Warner, T. T.: Quality assurance in atmospheric modeling. *Bull. Amer. Meteor. Soc.*, **92(12)**, 1601-1610, 2011.

21 Westra, S., and Coauthors: Future changes to the intensity and frequency of short-duration extreme rainfall. *Rev. Geophys.*,
22 **52(3)**, 522-555, 2014.

23 Willems, P., and Coauthors: Climate change impact assessment on urban rainfall extremes and urban drainage: methods and
24 shortcomings. *Atmos. Res.*, **103**, 106-118, 2012.

25 WMO: Anticipated advances in numerical weather prediction, and the growing technology gap in weather forecast. 2013.
26 [Available online at https://www.wmo.int/pages/prog/www/swfdp/Meetings/documents/Advances_NWP.pdf]

27 Xu, Z.X., and Chu, Q.: Climatological features and trends of extreme precipitation during 1979–2012 in Beijing, China.
28 *Proceedings of the International Association of Hydrological Sciences*, **369**, 97-102, 2015.

29 Xu, Z. X., and Zhao, G.: Impact of urbanization on rainfall-runoff processes: case study in the Liangshui River Basin in
30 Beijing, China. *Proceedings of the International Association of Hydrological Sciences*, **373**, 7-12, 2016.

31 Yu, R., Xu, Y., Zhou, T., and Li, J.: Relation between rainfall duration and diurnal variation in the warm season precipitation
32 over central eastern China. *Geophys. Res. Lett.*, **34(13)**, 2007.

- 1 Yu, W., Nakakita, E., Kim, S., and Yamaguchi, K.: Impact Assessment of Uncertainty Propagation of Ensemble NWP
- 2 Rainfall to Flood Forecasting with Catchment Scale. *Adv. Meteor.*, 2016.
- 3 Yucel, I., Onen, A., Yilmaz, K. K., and Gochis, D. J.: Calibration and evaluation of a flood forecasting system: Utility of
- 4 numerical weather prediction model, data assimilation and satellite-based rainfall. *J. Hydrol.*, **523**, 49-66, 2015.
- 5 Zhou Y. S., Liu L., Zhu K. F., and Li J. T.: Simulation and evolution characteristics of mesoscale systems occurring in
- 6 Beijing on 21 July 2012. *Chinese J. Atmos. Sci. (in Chinese)*, **38 (5)**, 885-896, 2014.
- 7

Table 1: Categories of experiments with different domain sizes, vertical resolutions, horizontal resolutions and spin-up times.

Scenario	Experiment Number	Domain Size	Vertical Levels	Horizontal Resolution (nesting ratio)	Spin-up Time
Domain Size (S1)	Case 0 (C0)	D01 40×40 D02 72×72 D03 90×90	29 (pressure level)	D01 40.5km; D02 13.5km; D03 4.5km 1:3:3	12 h
	Case 1 (C1)	D01 80×64 D02 120×120	As C0	As C0	As C0
	Case 2 (C2)	D01 160×128 D02 240×192	As C0	As C0	As C0
Vertical Resolution (S2)	Optimal Case in S1 (OS1)	As OS1	29	As C0	As C0
	Case 3 (C3)	As OS1	57	As C0	As C0
	Case 4 (C4)	As OS1	85	As C0	As C0
	Case 5 (C5)	As OS1	38 (model level)	As C0	As C0
Horizontal Resolution (S3)	Optimal Case in S2 (OS2)	As OS1	As OS2	1:3:3	As C0
	Case 6 (C6)	As OS1	As OS2	D01 40.5km; D02 8.1km; D03 1.62km 1:5:5	As C0
	Case 7 (C7)	As OS1	As OS2	D01 40.5km; D02 5.785km; D03 0.826km 1:7:7	As C0
Spin-up Time (S4)	Optimal Case in S3 (OS3)	As OS1	As OS2	As OS3	12 h
	Case 8 (C8)	As OS1	As OS2	As OS3	0 h
	Case 9-Case 19 (C9 - C19)	As OS1	As OS2	As OS3	24 h – 144 h per 12 h

Table 2: Correlations between the original and rescaled objective verification metrics.

Original Metrics	Representative Meaning	Rescaled Metrics	Threshold Value
POD	Probability of Detection	$POD' = POD$	N/A
$RMSE$	Root Mean Squared Error	$RMSE' = 1 - RMSE/RMSE_{max}$	+ 62.5 max
R	Pearson Correlation Coefficients	$R' = R$	N/A
$WRMSE$	RMSE of the Precipitable Water	$WRMSE' = 1 - WRMSE/WRMSE_{max}$	+ 8.3 max
WR	R of the Precipitable Water	$WR' = WR$	N/A
$RE_{P_{MAX}}$	Relative Error of the Maximum Precipitation	$P_{MAX}' = RE_{P_{MAX}} + 1$	N/A
RE_{TP}	Relative Error of the Total Precipitation	$TP' = RE_{TP} + 1$	N/A

Table 3: Comparison of the values of the error metrics in the initial experiment and the optimum experiments identified for each scenario.

Experiment Number	POD'	$RMSE'$	R'	$WRMSE'$	WR'	$PMAX'$	TP'
Case 0 (C0)	0.950	0.098	0.226	0.789	0.980	0.440	0.478
Case 1 (C1)	0.960	0.064	0.376	0.622	0.967	0.436	0.471
Case 3 (C3)	0.969	0.110	0.373	0.610	0.967	0.515	0.496
Case 6 (C6)	0.963	0.205	0.375	0.600	0.956	0.582	0.592
Case 12 (C12)	0.959	0.402	0.670	0.807	0.977	0.883	0.920

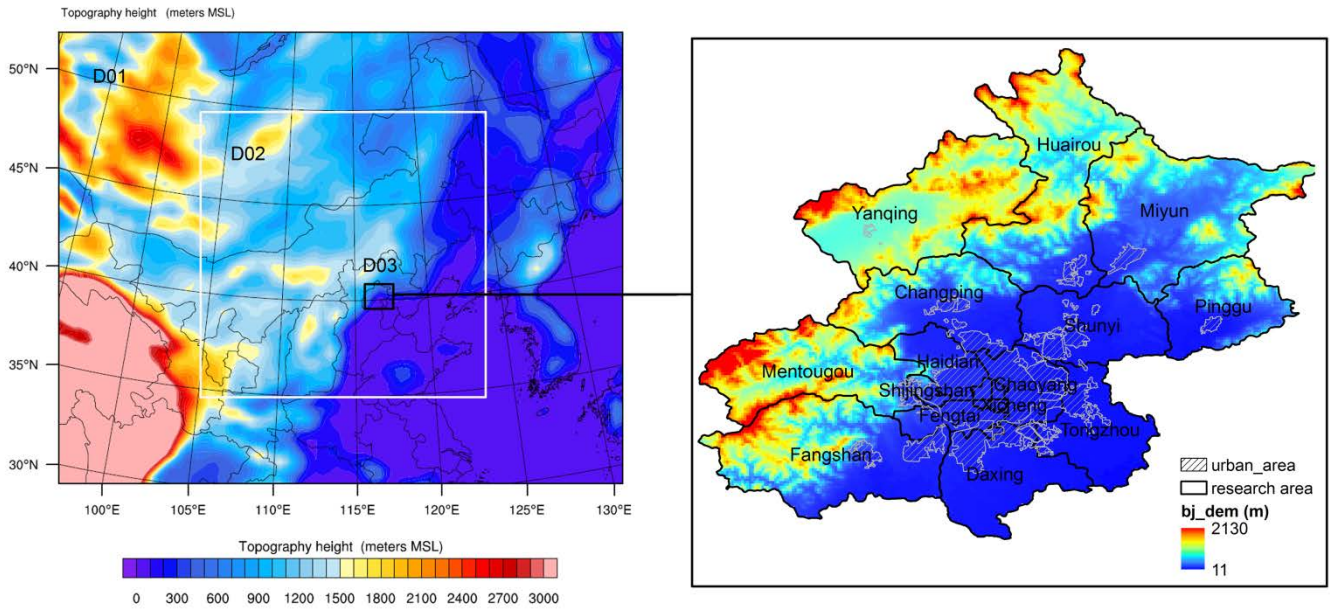


Figure 1: Relative location of the study area. The left panel shows the three nested domains adopted in most of the experiments, of which domain three (D03) covers the entire Beijing area; the right panel depicts the geographic features of the Beijing area.

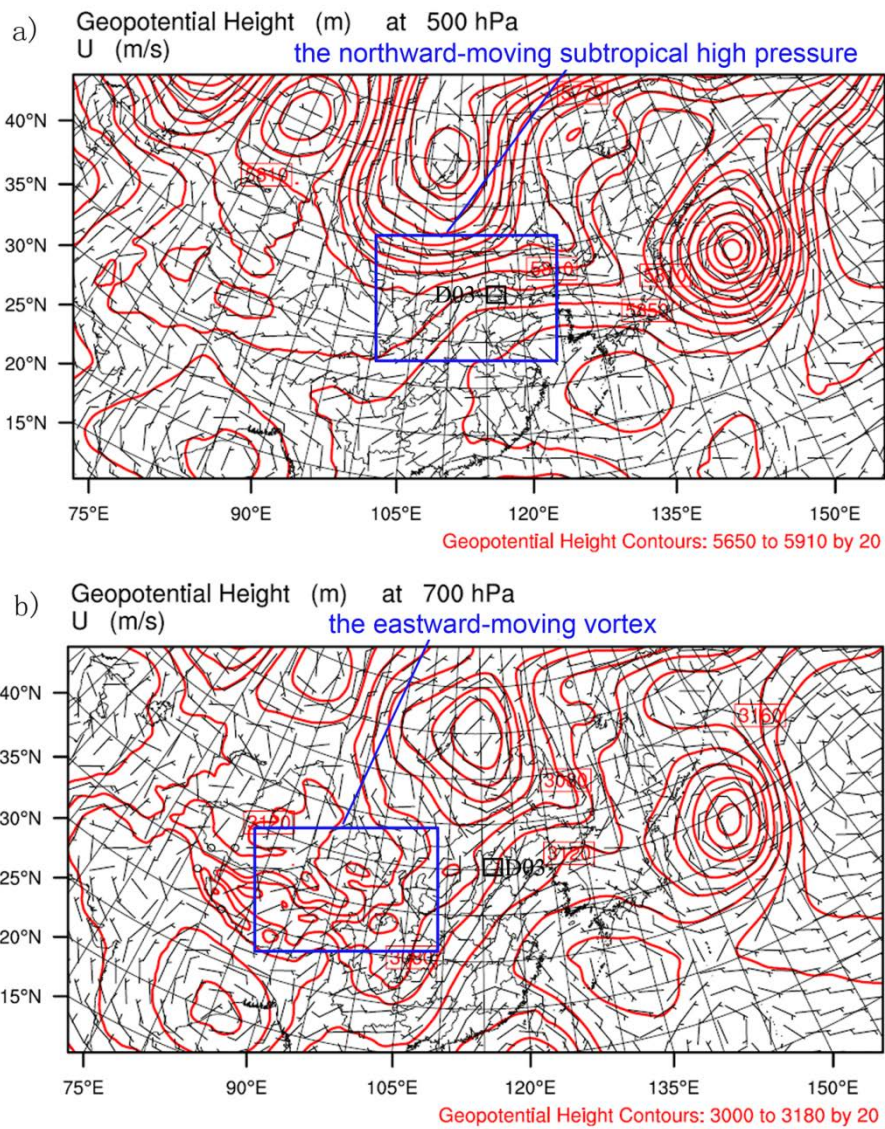


Figure 2: Initial wind field and geopotential height field at 12 pm on 20 July 2012 over the Northeastern Hemisphere obtained from the ERA-Interim reanalysis. (a) The fields at 500 hPa; (b) the fields at 700 hPa.

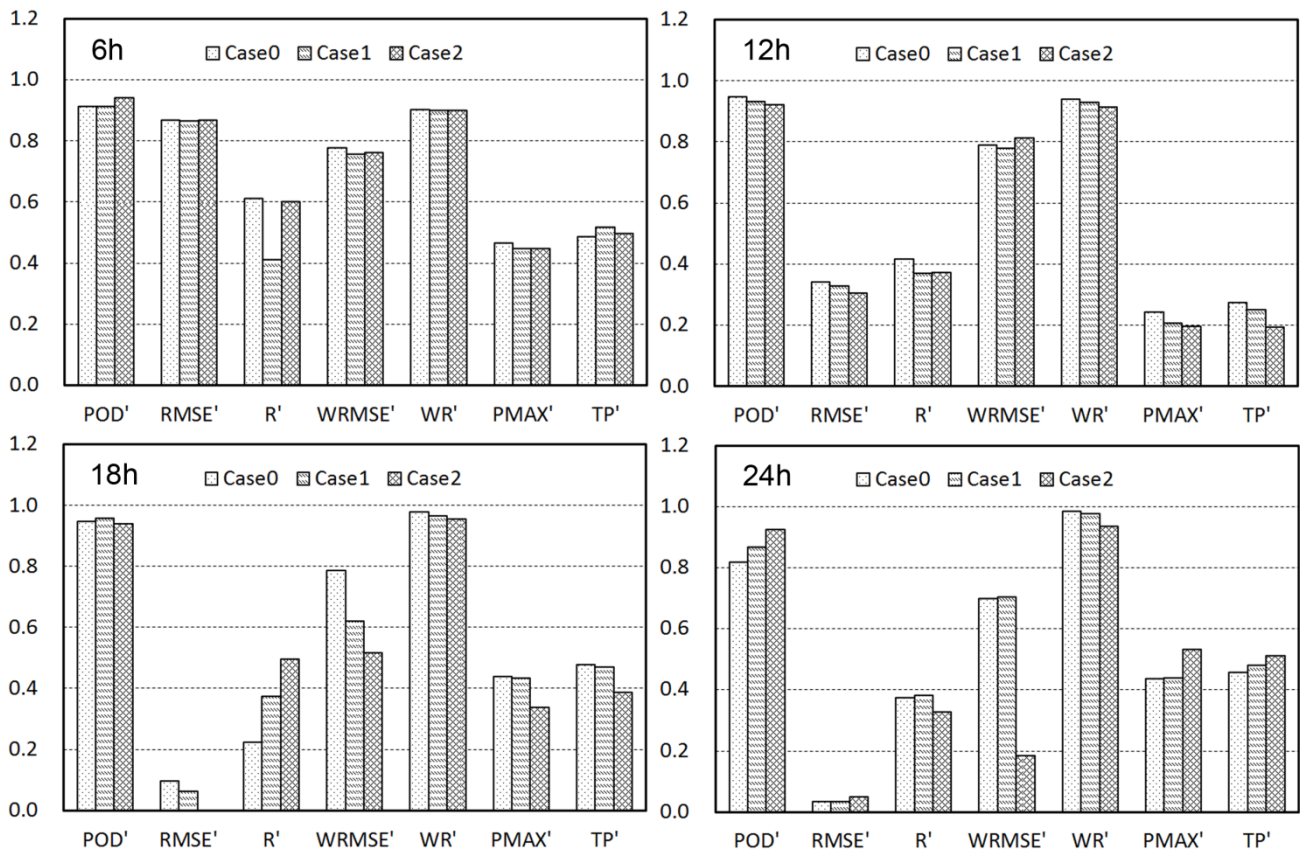


Figure 3: Spatial values of the verification metrics for the WRF domain size experiments, calculated over different temporal durations and over domain three. Case 0 (C0) incorporates the smallest domain, which covers north-central China; Case 1 (C1) incorporates a domain of intermediate size that covers northern China and part of Mongolia; and Case 2 (C2) incorporates the largest domain, which covers the Northeastern Hemisphere. The metrics are calculated over time periods of 6 h, 12 h, 18 h, and 24 h that begin at 12 am on 21 July 2012.

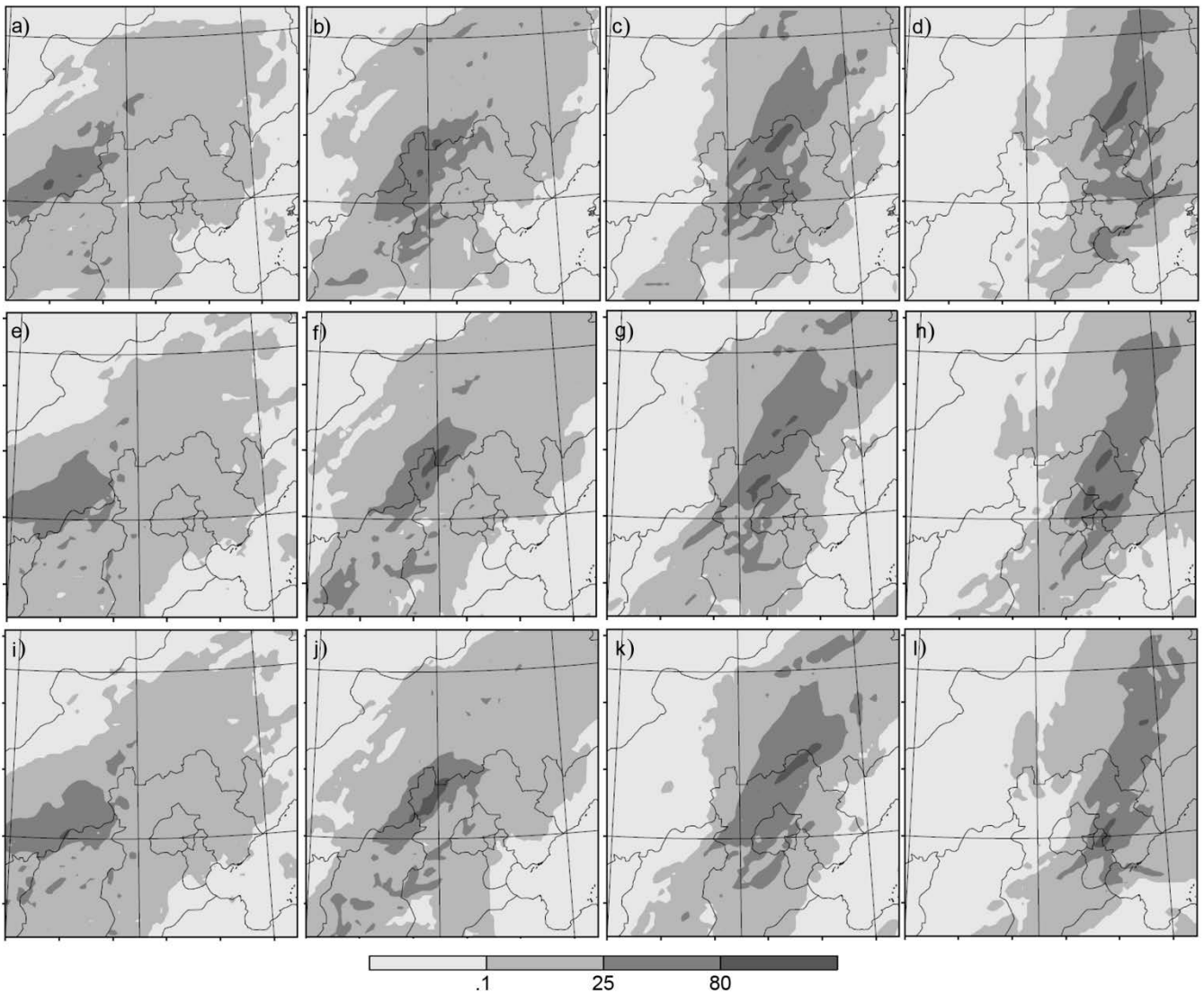


Figure 4: Spatial distribution of 6-h accumulated precipitation for the domain size experiments over the domain two area of C0 during the Beijing heavy rainfall beginning at 12 am on July 21, 2012. (a) Accumulated precipitation (AP) in C0 during the first 6-h period (0 h-6 h); (b) AP in C0 during the second 6-h period (6 h-12 h); (c) AP in C0 during the third 6-h period (12 h-18 h); (d) AP in C0 during the fourth 6-h period (18 h-24 h); (e) AP in C1 during the first 6-h period; (f) AP in C1 during the second 6-h period; (g) AP in C1 during the third 6-h period; (h) AP in C1 during the fourth 6-h period; (i) AP in C2 during the first 6-h period; (j) AP in C2 during the second 6-h period; (k) AP in C2 during the third 6-h period; and (l) AP in C2 during the fourth 6-h period.

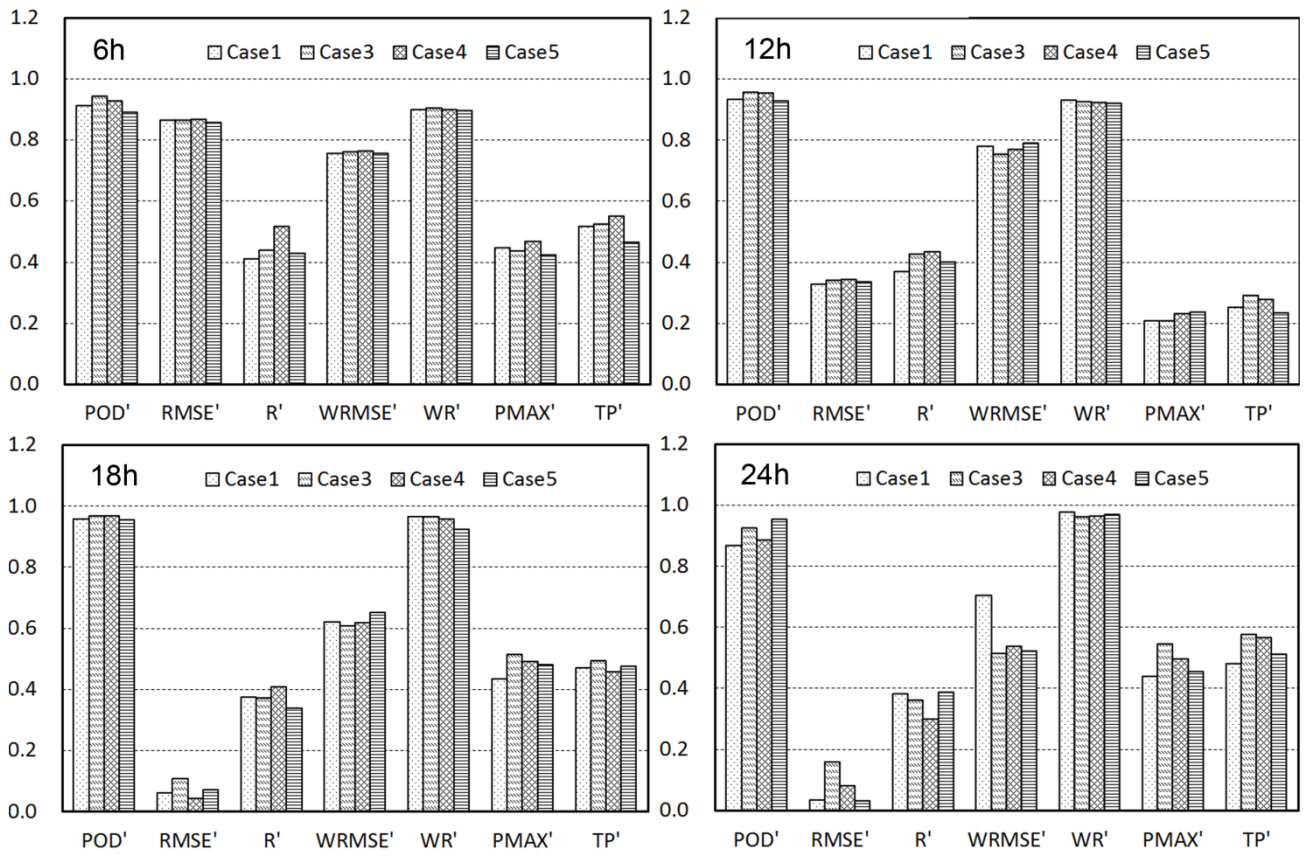


Figure 5: As in Fig. 3, but for the experiments in scenario two with different vertical resolutions. Case 1 is forced by the ERA-Interim pressure-level data with 29 vertical levels; Cases 3 and 4 are forced by the same data but include double and triple the number of vertical levels, respectively; Case 5 is forced by the ERA-Interim model-level data with 38 vertical levels.

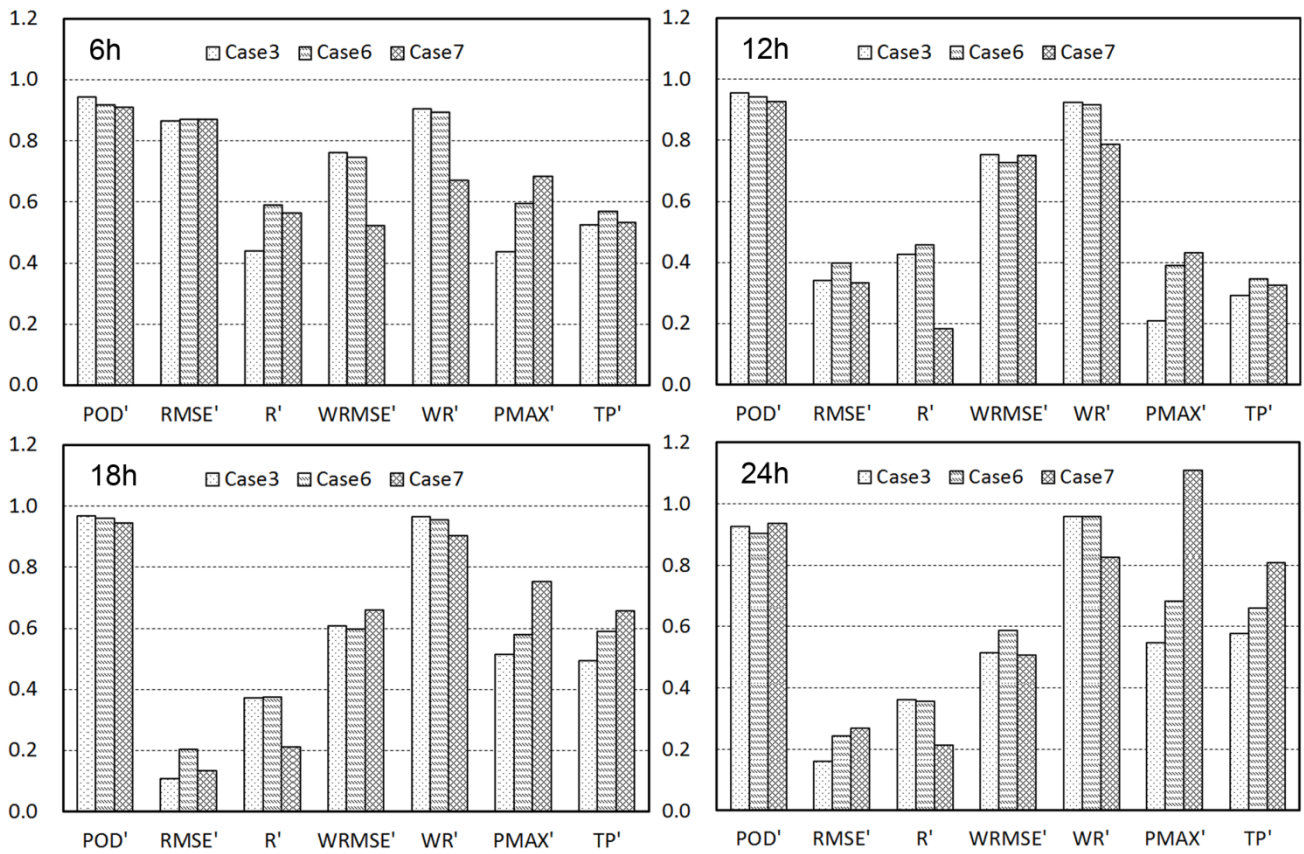


Figure 6: As in Fig. 3, but for the experiments in scenario three with different horizontal resolutions. Case 3 has an initial downscaling ratio of 1:3:3 with horizontal grid spacing of 40.5 km, 13.5 km and 4.5 km, whereas Cases 6 and 7 have the same large horizontal grid spacing with nesting ratios of 1:5:5 and 1:7:7, respectively. The innermost grid spacing is 1.62 km in Case 6 and 0.826 km in Case 7.

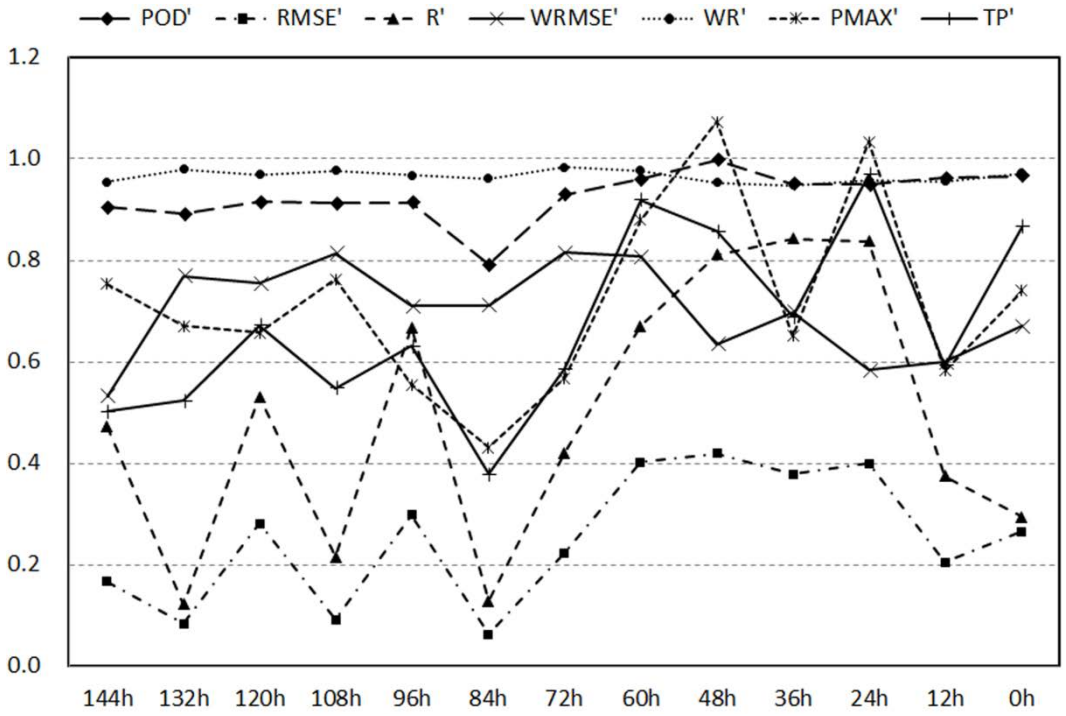


Figure 7: Spatial values of the verification metrics for the WRF spin-up experiments, calculated over 18-h periods and over domain three. Case 6 employs an initial spin-up time of 12 h; Case 8 employs a spin-up time of 0 h; and from Case 9 to Case 19, the spin-up time is increased from 24 h to 144 h by every twelve hours.

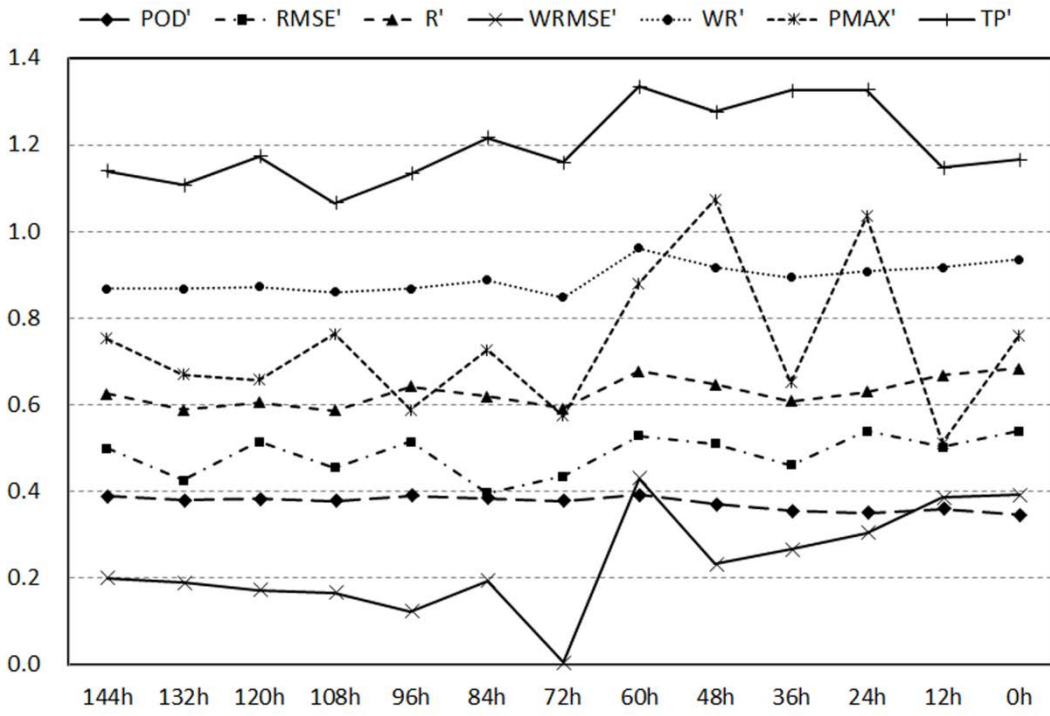


Figure 8: As in Fig. 7, but the metrics are calculated over 18-h periods and over domain two in Case 6.