

We would like to express our deep appreciation to the two reviewers (Dr. Doherty and the anonymous reviewer) for the thoughtful comments and insightful suggestions. Working to resolve these comments helps us add lots of interesting science and add value to the manuscript.

The primary objective of this study, as pointed out by Reviewer #2, is to build a bridge for linking the parallel PEST and WRF-hydro on the basis of HPC clusters and explore the computational benefits of this bridge. We do not attempt to extensively assess each individual tool or address questions in each individual domain, such as optimizing the objective functions in PEST or calibrating WRF-Hydro to achieve the best set of model parameters. However, we appreciate the opportunity every much during the revision of this manuscript by learning more about PEST especially the method of regularization for calibrating environmental models. We are also very glad that the reviewers found the bridge we built useful for helping WRF-Hydro users with the long and tedious model calibration.

In the revised version of this manuscript, several major changes are made based on both reviewers' comments/suggestions. They are listed below:

1. We re-do the WRF-Hydro calibration using SVD-based regularization method in PEST.
2. We consider prior information for the calibrated parameters.
3. We also consider different weight for the stations that are calibrated, based on their inversed mean of discharges.
4. To test the computational benefits of the bridge, we design five experiments by assigning different amount of computing resource for parallel PEST and for parallel WRF-hydro.
5. To constrain the problem size due to the limits of computing resource, we reduce the number of calibrated parameters to 22 according to the model sensitiveness of this particular study.

Please find our one-on-one response below to each reviewer's comment. A complete list of the changes made for the revised manuscript can be found in the "track changes" version of the manuscript. A clean version of the revised manuscript is also attached at the end.

Sincerely,
Jiali Wang
jjaliwang@anl.gov

Reviewer #1

Interactive comment on “A parallel workflow implementation for PEST version 13.6 in high-performance computing for WRF-Hydro version 5.0: a case study over the Midwestern United States” by Jiali Wang et al.

Doherty (Referee)

johndoherty@ozemail.com.au

Received and published: 6 December 2018

The authors describe modifications that they made to PEST to enhance its use on a HPC. They then describe use of their modified PEST in calibration of a complex surface water model. While I found the paper interesting, I found that it was lacking in information in some respects. For example nothing is said about the interface that they constructed between parallel PEST and the run management software that they employed.

Response:

Thank you for your comment. The interface is the most important thing we built in this study, and testing the operational feasibility and computational benefit of this interface are the main objectives of this manuscript. Hence it definitely should be described as you suggested. We add this paragraph in **Section 3.2 PEST files and settings**:

“The interface we have built between parallel PEST and the management software (SLURM here) is, in general, used for (1) setting the number of workers and the nodes for each worker to conduct a model run (WRF-Hydro here); (2) finding the nodes that are available; (3) setting up the working directory for the workers; (4) identifying the nodes that work for each worker; (5) passing the global files (same for all the working directory) to all the workers (these files include the lookup table files that are not to be calibrated, the namelist files for both LSM and hydrological sector, and restart files that generated by the previous simulations, or spin-up period); and (6) submitting the job for the entire calibration process, including parallel PEST and parallel WRF-hydro. This job can be submitted as a cold-start run or as a restart. The main difference for this interface on different management software is that different management software has its own way to submit jobs and identify available nodes. This difference requires some changes in the script we developed.”

Nor was any reference made to PEST settings.

Response:

Thanks for the comment. In our original version of manuscript, we used estimation mode for PEST, and considered equal weight for all four calibrated stations. There was no singular value decomposition (SVD) nor regularization used.

In our revised manuscript, we conduct the calibration using SVD-based regularization, we assign prior information for all the calibrated parameters, and we also consider different weights for the calibrated stations. We add the PEST setting in **Section 3.2 PEST files and settings**:

“This study presents calibration results from PEST using the SVD-based regularization in regularization mode to ensure numerical stability (Tonkin and Doherty, 2005). We focus on calibrating 22 parameters (see Table 1 and detail description in Sec. 3.3) using 96 observation points and 22 items of prior information for the calibrated parameters. In each item of prior information, a value equal to its default value provided by the WRF-Hydro v5.0 (or the log of its default value) is assigned for each adjustable parameter, assuming

that default values are the preferred values. All prior information equations are assigned a weight of 1.0. We assigned five different regularization groups to the prior information: Manning's roughness coefficients specified by Strahler stream order in CHANPARAM.TBL to one group; the parameters in HYDRO.TBL (Manning's roughness coefficients for overland flow as a function of vegetation types) to another group; and three global parameters for the Noah-MP (xslop1, refdk, and refkdt) in GENPARAM.TBL to the remaining three groups. The 96 observation points are given different weights based on the inversed mean of their observed discharge during the studied period (see the detailed description in Sec. 3.3 and Sec. 4.1). For a detailed description of these settings see the PEST User Manual (Doherty, 2015)."

While I agree with the authors that use of inversion methods that can parallelize model runs and handle the estimation of many parameters employed by a complex model is a much-needed addition to the arsenal of surface water modelling, I think that many more advances could be made than the authors have made. In particular, there was no mention of the use of Tikhonov regularization to accommodate parameter nonuniqueness at the same time as it promulgates uniqueness through obtaining a set of parameters that "make sense" from an expert knowledge point of view. This, I think, is one of the strongest arguments for use of gradient-based, highly parameterized methods in regional surface or land use model calibration, that is the ability to not just accommodate nonuniqueness, but to turn the "wobble room" engendered by nonuniqueness into formulation of an inverse problem that can actually make regionalization and transportability of parameters a reality.

Response:

Thanks for your comment, and we understand this is one of the major concerns for how PEST was used in our original manuscript to calibration a hydrological model and more broadly, environmental models. During the revision of this study, we conduct all the WRF-hydro calibration using SVD-based regularization. We also use prior information for the parameters that we calibrated, as you can see from our previous response about PEST settings. In **Section 5 Summary and discussion**, we also commented that, "In this study, we consider using the prior or regularization information only for the parameters that we calibrate. As is the case with solving inverse problems, prior information is added to improve the smoothness of the solutions. In order to build a more comprehensive calibration, an important aspect that can be considered is to enrich the prior with the available historical data. For example, in this particular case, one can use the historical observation data (e.g., April and May from the past few years) to enrich the prior information for the parameters. Hence, the regularization objective function in PEST will constitute not only the discrepancies between parameters and their "current estimates" but also the discrepancies between WRF-Hydro simulations and preferred values (which is the observed time series of historical discharge). Additionally, one can use the pilot points technique described by Doherty (2005) in conjunction with parameter estimation to add more flexibility to the calibration process. This will be potentially beneficial in improving the predictions"

In addition, to emphasize the importance of regularization for hydrological and environmental model calibration, we mentioned regularization in introduction: "PEST has four modes of operation. One of the modes is regularization mode, which supports the use of Tikhonov regularization and is found better for serving environmental models, because if implemented properly, it supports model predictions of minimum error variance, it is numerically stable and it embraces rather than eschews the heterogeneity of natural systems. Singular value decomposition (SVD) can be used as a regularization device to guarantee numerical stability of the calibration problem".

We also find that regularization mode does generate a better hydrograph shape compared to using estimation mode:

In **Section 4.1 WRF-Hydro calibration and validation**, we mentioned “Compared with the results of calibration using the estimation mode (no regularization) in PEST, the SVD-based regularization generates slightly better hydrograph shape with 24-hour later discharge peaks that are closer to the observations.”

In **Section 4.3 Evaluation of spatial transferability of the calibrated parameters**, we mentioned “We also find that using the SVD-based regularization for the PEST calibration captures the timing of discharge peak better than using the estimation mode, which is one-day earlier than the observations reaching the discharge peak.”

The authors use a simple objective function. This may be ok for some inverse problems. However as they point out, some of the smaller flows (in terms of location in space and location in a single flow time series) are not as well fitted as they could be. Perhaps weights should be a function of flow – and of location. Perhaps other important aspects of the flow time series should be made more visible to PEST through formulation of separate, targeted objective function components to ensure that these aspects of the time series are also well fit.

Response:

We understand this is another major concern about how PEST was setup for this study. We agree that, due to the fact that the calibrated station we chose are on different size in terms of water volume, to better handle the smaller river station (Station 1), considering different weight is fundamental, as also pointed out by Reviewer #2.

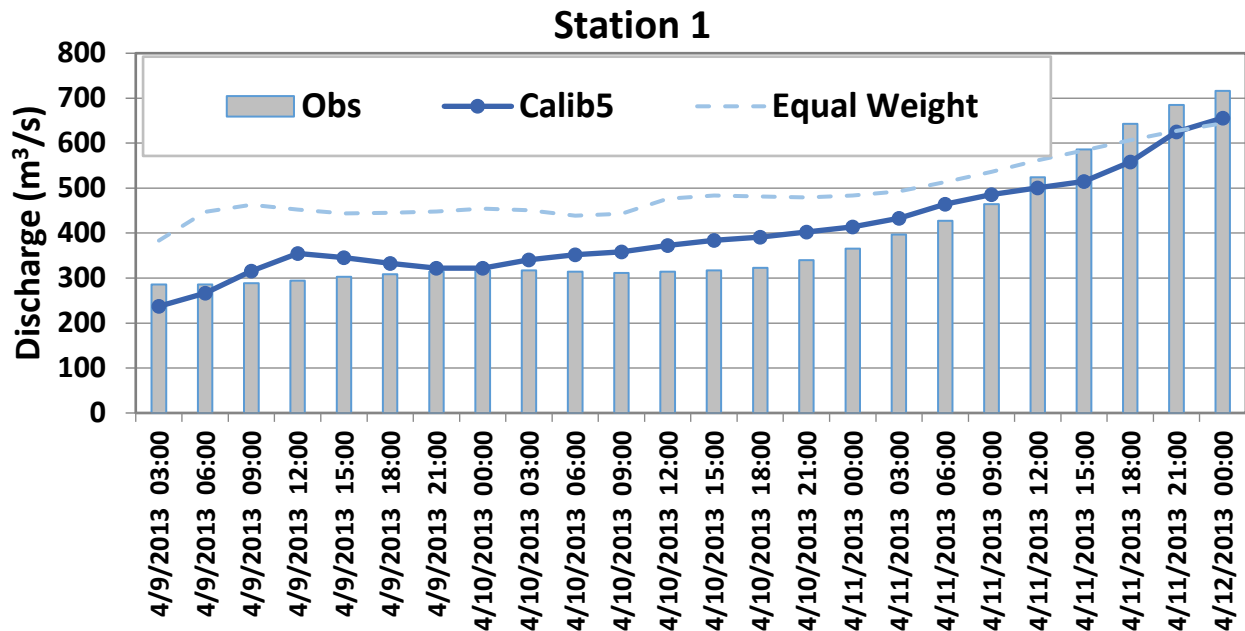
Therefore, during the revision of this study, we conduct all the PEST calibration of WRF-hydro considering a higher weight for Station 1 than for the other three stations. We add description about weight in the revised manuscript in **Section 3.2 PEST files and settings**, as you can find from previous response. We also add description in **Section 3.3 Calibration experiments**:

“We find that by using the same absolute weight for all four stations, the calibration helps three stations (Station 2, 3, and 4) with large water volumes to generate more reasonable results than do the default parameters; however, the results for Station 1, which has a relatively small volume of water, is not always better than the discharge that is modeled by using default parameters. Thus, we assign a weight of 15.0 for Station 1 versus a weight of 1.0 for the other three stations according to the inversed mean of observed discharge over these four stations in April 2013. The ratio of the weights between Station 1 and the other three stations stays similar even if the means are calculated based on different time periods.”

We describe the results in **Section 4.1 WRF-Hydro calibration and validation** by comparing PEST calibration using equal weight and a higher weight for Station 1:

“Compared with the results calibrated by using equal weights for all the stations, by giving a higher weight to Station 1 the model bias over Station 1 is significantly reduced, with a higher NSE (0.87 with higher weight versus 0.14 with equal weight) and lower RMSE (48.1 versus 123.6).”

Here is a figure (Figure I) showing how the weight helps the result of Station 1. “Calib5” represent the results using higher weight for Station 1, while “Equal Weight” the results using equal weight for all the four stations.



The authors make a big deal out of their modifications to parallel PEST so that it is HPC-friendly. Actually, I think that the BEOPEST version of PEST has similar capabilities. The original version of BEOPEST used both MPI and TCP/IP for communication between master and slaves (now called manager and workers). Now only TCP/IP is used. One of the reasons that BEOPEST's capabilities exceed those of parallel PEST in the HPC environment (actually on any network) is that the manager does not need to write model input files and read model output files across the network. This makes run management much faster, more secure, and able to take place in a greater variety of network environments.

Response:

Thanks for the comment. Since our major objective is to build the bridge (or interface) between the parallel PEST and WRF-Hydro on the basis of HPCs, we expect the interface still working if one wants to use BEOPEST instead of parallel PEST to calibrated WRF-hydro, assuming there is also a version of BEOPEST in Linux environment (this is not clear to us by reading the manual). The only main changes for our script would be copying the template files and instruction file into each working directory. The command to execute BEOPEST is also different, which is *beopest* instead of *ppest* for parallel PEST. We add this paragraph in **Section 5 Summary and discussion**:

“While this study ported the parallel PEST to HPC system and linked it to WRF-Hydro, we note that BEOPEST is available in the PEST family. BEOPEST has the same functionality as parallel PEST but uses a different approach for communication between master and workers. Working with HPC-enabled BEOPEST may save total time cost since BEOPEST uses the Transmission Control Protocol (TCP) and the Internet Protocol (IP) instead of message files (reading input and writing output between master and works) for communication. We expect it to be relatively straightforward to use BEOPEST to calibrate WRF-hydro on HPCs since the interface remains similar, except one needs to copy the template and instruction files in addition to the global files (see Section 3.1) into each working folder.”

In summary, I think that what the authors have done is good. However I also think that the potential for regional surface water model calibration and uncertainty analysis in a HPC environment still remains

largely untapped. Some of this potential will be realized with use of singular value decomposition to ensure numerical stability when inverse problems are ill-posed, use of Tikhonov regularisation to ensure parameter sensibility and transportability under the same conditions, and more creative formulation of the objective function than the authors have done.

Response:

We thank Reviewer #1 again for all your valuable comments. Although optimizing the objective functions in PEST is beyond the scope of this study, we do have some thoughts for future studies. We add this in **Section 5 Summary and discussion:**

“In this study, we consider using the prior or regularization information only for the parameters that we calibrate. As is the case with solving inverse problems, prior information is added to improve the smoothness of the solutions. In order to build a more comprehensive calibration, an important aspect that can be considered is to enrich the prior with the available historical data. For example, in this particular case, one can use the historical observation data (e.g., April and May from the past few years) to enrich the prior information for the parameters. Hence, the regularization objective function in PEST will constitute not only the discrepancies between parameters and their “current estimates” but also the discrepancies between WRF-Hydro simulations and preferred values (which is the observed time series of historical discharge). Additionally, one can use the pilot points technique described by Doherty (2005) in conjunction with parameter estimation to add more flexibility to the calibration process. This will be potentially beneficial in improving the predictions.”

Reviewer #2

Interactive comment on “A parallel workflow implementation for PEST version 13.6 in high-performance computing for WRF-Hydro version 5.0: a case study over the Midwestern United States” by Jiali Wang et al. Anonymous Referee #2

Received and published: 16 January 2019

The paper of Wang et al. deals with a potentially interesting implementation of the parallel version of the PEST software. PEST is a powerful and very useful tool for hydrologists, helping them during long and “exhausting” calibration sessions. Therefore, introducing the portability of parallel PEST to HPCs is good news and, specifically for the present paper, the main theme to highlight.

Response:

We are so glad that the reviewer finds this study helpful for hydrologist with the most tedious part of model development and application — calibration. We also wish to thank the reviewer for your great insight about the highlight of this study, which add important value to this manuscript.

Nevertheless, in my opinion the way the paper is structured mainly highlights, instead of the advantages of the novelty, the performances of the PEST calibration, which is something widely and well assessed by the hydrology research community. Almost all figures and tables deal with PEST results. Furthermore, the calibration procedure presented is questionable from different points of view (some of which are exposed later).

Response:

Thanks and we agree your comment. In our revised manuscript, we delete the section of 7-day calibration (figures and tables) considering it is unnecessary to support our main objective. We also delete or shorten the descriptions about WRF-Hydro and PEST, which readers can easily learn from the User Guides/Manuals. We thus focus on only two points: one is to show the operational feasibility, and the other is to explore the computational benefits of the HPC-enabled parallel PEST linked to WRF-Hydro. To demonstrate the first point, we calibrate WRF-Hydro for 3 days using SVD-based regularization method in PEST, and considering different weight for the four calibrated stations. For the second point, we design new experiments, using different computing resources for PEST workers and for WRF-hydro. Details can be found below in our response to your later comments.

The most interesting/innovative Section of the paper is Section 5.1, but the analysis of scale-up capabilities should be described with much more detail. Concerning the main outcomes of the paper highlighted in the summary, points from 2 to 5 are quite obvious (they deal with the recognized skills of the PEST software), while point 1 should be expanded: what does a factor of 30 “with respect to a serial calibration” exactly mean? In my opinion it’s not a rigorous statement. What do the authors exactly mean with “serial”? Even though PEST calibration is serial, WRF-Hydro can run in a parallel fashion, and the speed of the calibration process would depend on the number of nodes used for the hydrological simulation. A possible idea is to provide hints about the trade-off between the number of nodes/CPU’s used for running the parallel model (i.e., WRF-Hydro in this case) and the number of nodes/CPU’s used for running PEST in a parallel fashion. I guess it depends somehow also on the dimensions of the domain (and no information is given here about the number of cells in which the basin is discretized, so the reader has no idea about the actual computational burden).

Response:

Thanks for your comment and suggestions. We re-write **Section 5 Summary and discussion** to only summarize the findings of the study, and then raise some key points that beyond the scope of this study but may inform future studies. In other words, we delete majority of the summary that appears in points 2 to 5 in our original manuscript.

In the revised manuscript we expand the discussion of scale-up capabilities of parallel PEST linked to WRF-hydro by designing more experiments using different computing resource for PEST workers and for WRF-hydro. We add **Section 4.2 computational benefits of parallel PEST on HPCs**. Some key notes about the experiments and our findings are quoted below:

“In this study we test the computational performance of HPC-enabled parallel PEST using different number of workers (6, 12, and 23) for the 22-parameter calibration. As shown in Table 3, we conducted five experiments: Test 1 uses 23 workers, Test 2 uses 12 workers, and Test 3 uses 6 workers. All three tests use two nodes for each worker to run WRF-Hydro in parallel.”

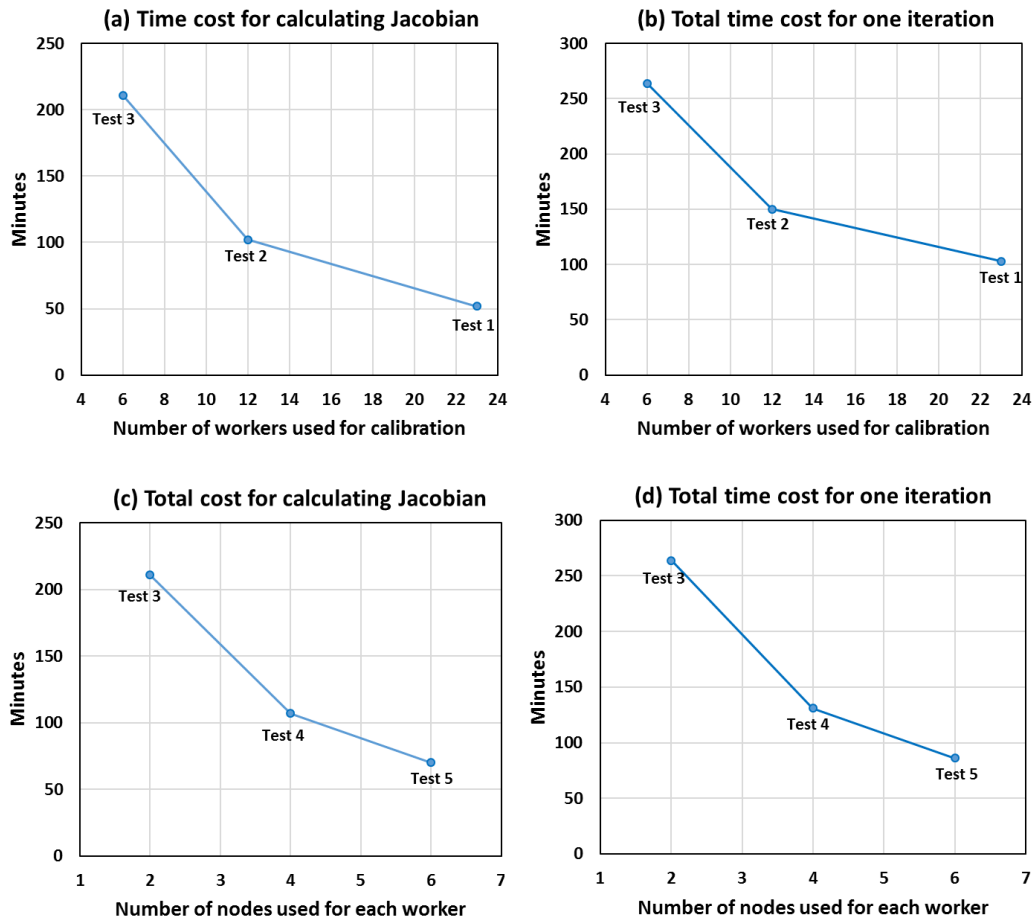
“In order to test the trade-offs between the computing nodes used for running parallel WRF-Hydro and the workers used for running parallel PEST, Tests 4 and 5 use different number of nodes for each worker to run WRF-Hydro in parallel. Explicitly, Test 4 uses four nodes per worker, and Test 5 uses six nodes per worker. Both tests use six workers for running the parallel PEST.”

“when we assign 12 (and 6) workers to parallel PEST, the time cost for calculating the Jacobian matrix is increased by a factor of 2 (and 4) compared with the time cost of using 23 workers. The time cost for the parameter upgrade stays similar for the three experiments because only one cycle of WRF-hydro simulation is conducted to test the Marquardt lambdas. As a result, the total time cost for Test 2 is ~1.5 times more than that for Test 1, and the total time cost for Test 3 is ~1.5 times more than that for Test 2 (Fig. 4b). By extrapolating the speedup curve shown in Fig. 4a and Fig. 4b, we expect the total time cost to be ~1516 minutes when using only one worker (or sequential mode), which is about 15 times slower compared with running the PEST in parallel mode using 23 workers.”

“On the other hand, for the same case study and using the same number of nodes for running parallel WRF-Hydro, we can estimate the computing speedup by assuming an increase in the number of calibrated parameters to 50. This would be the case, for example, to evaluate model sensitiveness to the physics in Noah-MP or the spatial variabilities of certain parameters. We then expect to use 51 workers to achieve the best computing performance for parallel PEST. This would then be 28–30 times faster than running PEST using one worker (or in sequential mode). Similarly, if 100 parameters were used for the calibration for the same case study, a factor of up to 60 speedup in the calibration process would be achieved by running HPC-enabled parallel PEST.”

“In addition, by increasing the number of nodes for each worker to conduct WRF-Hydro (Tests 3, 4, and 5), the time cost for the entire calibration process is significantly reduced (Figs. 4c and 4d). Specifically, the WRF-hydro scales up well when using four and six nodes compared with using two nodes per worker for running the WRF-Hydro. Both the time spent on calculating the Jacobian matrix and the time spent on testing the parameter upgrades are decreased by 49% and 67%, respectively, when using four and six nodes. Therefore, the total time spent is also decreased when using more nodes for each worker (see Table 3). Increasing the number of nodes to eight for each worker will most likely further decrease the time cost by 70–75% compared with using only two nodes per worker. Moreover, if one has a larger study area such as the entire contiguous United States, we expect the WRF-Hydro to have an even better scale-up capability (e.g., on dozens of nodes) than this study.”

Here is a figure to show the computational benefit using parallel PEST to calibrate WRF-hydro.



Another important point, that should be better discussed, is the missed capability of the implemented version of PEST to deal with the calibration of spatially distributed parameters. This is important because it's reasonable to expect parallel PEST executions with WRF-Hydro for wide domains, and wide domains often need spatial differentiation of spatially distributed parameters, like, e.g., OVROUGHRTFAC, RETDEPRTFAC or other spatially distributed parameters available with WRF-Hydro v5.0.

Response:

We do acknowledge the importance of regionalization of parameter calibration, which definitely deserves future studies especially for large domains. We have been trying to add the interface on top of what we have now, to consider the spatial distributed files/parameters. For example, one can add regional OVROUGHRTFACs (e.g., their lower/upper bounds, and default values) in the PEST control file based on catchments/basins/regular regions etc. the potential challenge is that, the selection of the locations and sizes of catchment may introduce significant uncertainties to the calibration results. Thus it requires systematic and comprehensive investigation and understanding of the study area. We add a paragraph in **Section 5 Summary and discussion** about this:

“We only calibrate the parameters in lookup tables. Using a single value to represent a physics may work for a small domain but could be problematic for a large domain, especially we expect the HPC-enabled parallel PEST to execute with WRF-Hydro for large domains, which often need parameter regionalization. For example in WRF-Hydro v5, there are many spatially distributed parameters available such as the overland flow roughness scaling factor (OVROUGHRTFAC), the factor of maximum retention depth (RETDEPRTFAC), and the soil related parameters (when compiled with SPATIAL_SOIL=1). Calibrating these spatial parameters based on grid scale (e.g., catchments) rather than a single value will give the model more flexibility and thus can better fit the observations (Wagener and Wheeler, 2006; Hundecha and Bardossy, 2004). In practice, for example, one can include regional OVROUGHRTFACs (e.g., their lower/upper bounds, and default values) in the PEST control file based on catchments. However, the selection of the locations and sizes of catchment may introduce significant uncertainties to the calibration results, which requires systematic and comprehensive investigation and understanding of the study area.”

By the way, another limitation is that, at least as I understand, the calibration is available only against observed streamflow. Of course, this is the first option but not the unique one (one can decide to calibrate also against, e.g., soil moisture or latent heat flux data).

Response:

For the calibration exercise we did in this study, we use frxst_pts_out.txt as an instruction file which serves the format of output files of each working directory. In this file there is only discharge and water level available, so we calibrate the model using discharge data. It is feasible, however, to calibrate other variables as long as the observation data is available. For example, one can either find the closest point from the gridded dataset to the observation location and then compare that point to observations; or one can change the WRF-Hydro I/O code to output other variables in the frxst_pts_out.txt file, so they can still use the same interface we build here to calibrate other variables in addition to the discharge. We have added this regard in **Section 5 Summary and discussion**.

Finally, another important point is to (at least) discuss the problem of equifinality, which is incidentally (but not explicitly) dealt with in P11 L29 – P12 L5.

Response:

Thanks for your comment. It's actually interesting to think about this together with your other comment that one can also calibrate other variables rather than discharge. Since equifinality is an important source of model uncertainty, to reduce the model uncertainty, one may calibrate the model using multiple variables instead of one variables. This way the calibrate can constrain the model drift and may reduce the model uncertainty of prediction of certain variables (e.g. discharge and soil moisture). We add a paragraph in the **Section 5 Summary and discussion**:

“The optimal parameter set obtained from this study is from the 5th iteration of parallel PEST by testing five Marquardt lambdas. Testing different number of lambdas or calibrating different number of parameters may generate a different set of optimal parameters. These parameter sets can all make physical sense and be equally good for reproducing observed discharges. This problem is named equifinality (Beven and Freer, 2001; Savenije, 2001), which is an important source of model uncertainty. To reduce the model uncertainty through reducing the equifinality, hydrologists carry out additional modelling objective for model evaluation to find more useful parameter sets (Mo and Beven, 2004; Gallart et al., 2007). Alternatively, inspired by No. 3 discussed above, one can calibrate the WRF-hydro model based on more than one variables, such as discharge and soil moisture (or heat flux or water table depth) to reduce the number of optimal parameter sets, and thus reduce the model uncertainty of predictions for these variables.”

Summarizing, though I acknowledge that the research presented is potentially interesting and innovative, I suggest to re-think the paper highlighting much more the computational benefits provided and reviewing the calibration performed in the case study.

Response:

We thank Reviewer #2 again for all your great insights. Our responses and revisions can be found above and below, as well as in the revised manuscript with tracked changes.

Following, a (not comprehensive) list of doubts regarding the calibration procedure and other minor comments and typos. I hope my comments can help improving the research.

Doubts about the calibration procedure:

Even though I acknowledge that authors decided to “focus less on extensively assessing the performance of the WRF-Hydro model”, several aspects of the calibration procedure are very questionable.

1. no information about spin-up. This is extremely important, especially for such a short range calibration (only few days). The model should be run in advance (at least one month, I would say) in order to let several variables (e.g., moisture fields) have a realistic spatial distribution.

Response:

We did run the model for 3 months for spin-up. We apologize for not including it in the model description. Here is what we add in the revised manuscript:

“We start the WRF-Hydro simulation on Jan 1 2013 and run the model for more than 3 months to reach equilibrium. This 3-month period is considered as spin-up time and is excluded from model calibration and evaluation. We calibrate the river discharge calculated by the WRF-Hydro model from 00UTC April 9 to 00UTC April 12 2013, considering it is long enough to achieve our objective. We then evaluate the model performance against U.S. Geological Survey (USGS) observed river discharge from 00UTC April 12 to 00UTC 25, 2013.”

2. the authors state that: April 8-11 moderate rain, April 12-14 no rain, April 15-18 rain, peak flow April 19. 3-day calibration is: April 9-11 (to be precise, April 12 at midnight), then validation is April 13-23 (April 12 is missed). 7-days calibration is April 9-15, validation is April 17-23. To me, it does not make too much sense that 4 more days are added when only the last one is rainy. It would be much better to calibrate the model with respect to a previous flood event, as it is usual. After all, observing graphs in figures 3 and 5 one after another just shows that increasing the number of days used for calibration improves the performances (but this is rather obvious), even though not yet enough.

Response:

As we mentioned earlier, in the revised manuscript we only focus on model calibration during April 9-11, and validation from April 12-24. We delete the 7-day calibration/validation results considering it is not necessary nor helpful to demonstrate our main objective. We add these sentences to emphasize this regard:

“The primary objective of this study is to build a bridge for linking the parallel PEST and WRF-hydro on the basis of HPC clusters and to explore the computational benefits of this bridge. We do not attempt to

extensively assess each individual tool or address questions in each individual domain, such as optimizing the objective functions in PEST or calibrating WRF-Hydro for a long time period considering all the relevant parameters to achieve an optimal parameter set. The calibration period thus is limited to only three days, which we believe long enough to achieve our objective and to understand WRF-Hydro's sensitivity to the calibrated parameters.”

We also add text in **4.1 WRF-Hydro calibration and validation** to explain the reason for the bias in hydrograph shape, such as the early peak and the faster decrease of river discharge:

“For Station 1, the WRF-Hydro almost captures the timing of the peak of discharge, although it still underestimates the water volume by ~25%. The reason is that this study uses a direct pass-through baseflow module, which does not account for slow discharge and long-term storage of the baseflow. Therefore, the largest contribution to river discharge is from precipitation, and groundwater does not contribute much discharge to the channels in a long-term view, as is also true for the other three large river stations. Different from Station 1, for the other three large river stations, the WRF-Hydro modeled discharge increases soon after the peak of precipitation and reaches a peak on April 21, 2013, which is much earlier than the observed peak of river discharge (near April 24). The reason is that the water contributions for these stations are from a larger river basin (Mississippi River) than we included in our current study area. Thus, when a heavy precipitation event occurs over the entire river basin, there will be a significant lag time (especially at the lower part of the basin) between the peak of precipitation amount and the peak of river discharge. For example, the precipitation over the upper part of Mississippi River Basin (MRB) has a peak amount on April 18–19, but the river discharge did not reach its peak until April 24. Because our studied area covers only half of the MRB, the modeled river discharge has a shorter delay period after the peak of precipitation than does the observed river discharge. Enlarging the study area to include the entire MRB may improve this situation. Alternatively, calibrating and validating local rivers that are included in the current study area may also reduce the bias in hydrograph shape compared to calibrating and validating large rivers. On the other hand, the WRF-Hydro simulated river discharge decreases soon after it reaches the peak and much earlier than the observed discharge. The reason is again that the direct pass-through baseflow employed by this study does not account for slow discharge and long-term storage of the baseflow. As a result, the contribution from the baseflow to the river discharge in model simulations does not stay as long as in real situations. In the observations, the river discharge decreases from the peak at a speed of ~500 m³/sec per day, while the modeled river discharge decreases from the peak at a speed of ~1667 m³/sec per day. Using exponential storage-discharge function for the baseflow may improve this situation.”

3. In order to deal with the observed streamflow in Section 1, it is fundamental to work with weights.

Response:

Thanks for your comment. This is also a major concern of Reviewer #1. For all the experiments we present in the revised manuscript, we consider a higher weight for Station 1. We add description about weight in the revised manuscript in **Section 3.2 PEST files and settings**, and in **3.3 Calibration experiments**:

“The 96 observation points are given different weights based on the inversed mean of their observed discharge during the studied period (see the detailed description in Sec. 3.3 and Sec. 4.1).”

“We find that by using the same absolute weight for all four stations, the calibration helps three stations (Station 2, 3, and 4) with large water volumes to generate more reasonable results than do the default parameters; however, the results for Station 1, which has a relatively small volume of water, is not always better than the discharge that is modeled by using default parameters. Thus, we assign a weight of 15.0 for

Station 1 versus a weight of 1.0 for the other three stations according to the inversed mean of observed discharge over these four stations in April 2013. The ratio of the weights between Station 1 and the other three stations stays similar even if the means are calculated based on different time periods.”

We describe the results in **Section 4.1** by comparing PEST calibration using equal weight and a higher weight for Station 1:

“Compared with the results calibrated by using equal weights for all the stations, by giving a higher weight to Station 1 the model bias over Station 1 is significantly reduced, with a higher NSE (0.87 with higher weight versus 0.14 with equal weight) and lower RMSE (48.1 versus 123.6).”

A figure is shown above in the response to Reviewer #1.

Minor comments, grammar and typos

P6 LL6-17: not clear if in this case overland flow is switched on. It should.

Response:

Yes, it is. We change the sentence to emphasize this regard: “Overland flow, saturated subsurface flow, gridded channel routing, and a conceptual baseflow are active in this study”

“If overland flow is active as it is in this study, it passes water directly to the channel model.”

P6 L19: probably “tools”

Response: corrected

P7 L18: GENPARAM.TBL

Response: fixed

P8 LL11-12: master, not mater. The full stop is missing.

Response: fixed

P8 L30: As it is a common problem, it is usually solved ‘simply’ reallocating manually the stations. It’s a pity to miss streamflow data for this reason

Response:

We regret we didn’t do that, and it should be done in future applications. To clarify and to emphasize this regard, we change the sentence to:

“We calibrated WRF-Hydro using four USGS sites (referred to as Station 1, Station 2, Station 3, and Station 4 hereafter), as shown in Fig. 1. More USGS sites could be included if one manually reallocated the stations that were not properly assigned to the desired location on the channel network by the GIS tool.”

P9 L24 and following: I suggest to explicitly declare also the meaning of the ovn parameter

Response: added the vegetation type in Table 1 for each own parameter*

P12 L17: 50%, maybe

Response: Apologize for the wrong number. It should be 66.7%

Figs.2 and 3: April 12 is missing. It should be the first validation day, I guess.

Response: added.

Figs. 4 and 5: the same for April 16

Response: this figure is deleted as it is for the 7-day calibration result

Table 3: the note is incorrect, it refers to information about the 3-day calibration

Response: fixed

Section 4.3: this is a purely “hydrological” analysis that could be skipped, given the numerous limitations of the calibration procedure and the focus on the implementation of the PEST software

Response:

We still keep this section and the figure in the revised manuscript. One of the reasons is that, it is more obvious using these stations to show the benefit (generate better hydrograph shape) of SVD-based regularization compared to non-regularization, which is the method we emphasize that should be applied for calibrating hydrological and environmental models. The other reason is that, these rivers are relatively small and are local rivers that included in the current study area, therefore, (1) the lag time between precipitation peak and discharge peak is much shorter than those for Station 3, 4, 5, and the hydrograph shape is well captured by the optimal parameter set; (2) the slow-discharge effect from baseflow is also relatively small so the discharge decrease faster after the peak than that for Station 3, 4, 5. This is also well captured by the optimal parameter set. Overall, the WRF-Hydro calibration with current configuration actually did a good job to capture the hydrograph features for these stations.

P16 LL9-10: please check the sentence

Response: this sentence is deleted, and changes are made accordingly through the entire manuscript.

P16 L18: to investigate

Response: corrected.

A parallel workflow implementation for PEST version 13.6 in high-performance computing for WRF-Hydro version 5.0: a case study over the ~~Midwestern~~midwestern United States

Jiali¹Jiali Wang, Cheng¹Cheng Wang, Andrew²Vishwas Rao, ¹Andrew Orr, Rao¹Rao Kotamarthi
Argonne¹Argonne National Laboratory, Environmental Science Division, 9700 South Cass Avenue, ArgonneLemont, IL 60439, USA

²Argonne National Laboratory, Mathematics and Computer Science Division, 9700 South Cass Avenue, Lemont, IL 60439, USA

Correspondence to: Jiali Wang (jialiwang@anl.gov)

Abstract. ~~Surface hydrological models must be calibrated for each application region.~~ The Weather Research and Forecasting Hydrological ~~system~~(WRF-Hydro) system is a state-of-the-art numerical model that models the entire hydrological cycle based on physical principles. ~~However,~~ asAs with other hydrological models, WRF-Hydro parameterizes many physical processes. ~~As a result~~Hence, WRF-Hydro needs to be calibrated to optimize its output with respect to observations. ~~However, when~~ for the application region. ~~When~~ applied to a relatively large domain, both WRF-Hydro simulations and calibrations require intensive computing resources and are best performed ~~in parallel on multimode, multicore high-performance computing (HPC) systems.~~ Typically, each physics ~~parameterization-based model~~ requires a calibration process that works specifically with that model, and is not transferrable to a different process or model. ~~Parameter Estimate Tool~~ The parameter estimation tool (PEST) is a flexible and generic calibration tool that can be used in principle to calibrate any ~~numerical code.~~ ~~However, PEST in~~ of these models. ~~In its current existing configuration,~~ however, PEST is not designed to work on the current generation of massively parallel ~~high-performance computing (HPC)~~ clusters. ~~This study~~To address this issue, we ported the parallel PEST to HPCs and adapted it to work with ~~the~~ WRF-Hydro. The porting involved writing scripts to modify the workflow for different workload managers and job schedulers, as well as developing code to connect ~~Parallel-parallel~~ PEST to WRF-Hydro. ~~We~~To test the operational feasibility and the potential computational benefits of this first-of-its-kind HPC-enabled parallel PEST, we developed a case study using a flood in the ~~Midwestern~~midwestern United States in 2013 ~~to test the operational feasibility of the HPC-enabled parallel PEST.~~ ~~We then~~

1 ~~evaluate the WRF Hydro performance in water volume and timing of the flood event. We also~~
2 ~~assess the spatial transferability of the calibrated parameters for the study area. We finally discuss~~
3 ~~the scale-up capability of. Results on a problem involving calibration of 22 parameters show that~~
4 ~~on the same computing resource used for parallel WRF-Hydro, the HPC-enabled parallel PEST ~~to~~~~
5 ~~provide insight for PEST's application to other hydrological models and earth system models on~~
6 ~~current and emerging HPC platforms. We find that, for this particular study, the HPC-enabled~~
7 ~~PEST calibration tool can speed up WRF-Hydro~~the calibration process by a factor of 30~~up to 15~~
8 ~~compared ~~to~~with commonly-used PEST in sequential mode. The speedup factor is expected to be~~
9 greater with a larger calibration approaches problem (e.g., more parameters to be calibrated or a
10 larger size of study area).

11 **1 Introduction**

12 ~~Hydrological models are important tools for research relevant but not limited to, water resource~~
13 ~~management, flood control, and hydrological response to climate change (Zanon et al., 2010;~~
14 ~~Papathanasiou et al., 2015). Conceptual hydrological models express hydrological processes in the~~
15 ~~form of abstract models that come from physical phenomenon and experience. Physically based~~
16 ~~hydrological models contain definite~~Physically based hydrological models contain detailed
17 physical mechanisms to model the hydrological cycle, but many complex physical processes in
18 these models are parameterized. For example, the state-of-the-art Weather Research and
19 Forecasting Hydrological (WRF-Hydro) modeling system (~~WRF-Hydro;~~ (Gochis et al., 2015) has
20 dozens of parameters that can be land- and river-type dependent and are typically specified in
21 lookup tables. ~~Both conceptual hydrological models and physically based~~Therefore, these
22 hydrological models need to be calibrated before they can be applied to research: over different
23 regions. In this context, calibration refers to ~~the hydrologists' need to adjust~~adjusting the values of
24 the model parameters so that the model can closely match the behavior of the real system it
25 represents. In some cases, the appropriate value for a model parameter can be determined through
26 direct measurements conducted on the real system. ~~However, in~~In many situations, however, the
27 model parameters are conceptual representations of abstract watershed characteristics and must be
28 determined through calibration. In fact, model calibration is the most time-consuming step, not
29 only for hydrological models, but also for ~~earth~~Earth system model development, because both
30 parametric estimation and parametric uncertainty analysis require hundreds—if not thousands—

1 of model simulations to understand how perturbations in model parameters affect simulations of
2 dominant physical processes and to find the optimum value of a single parameter.

3
4 WRF-Hydro is a ~~practical physics-based~~ numerical model that can simulate the entire hydrological
5 cycle using advanced high-resolution data such as satellite and radar products. Compared ~~to~~with
6 the traditional land surface model (LSM) used by WRF, WRF-Hydro provides a framework for
7 multiscale representation of surface flow, subsurface flow, channel routing, and baseflow, as well
8 as a simple lake/reservoir routing scheme. As a physics-based model, WRF-Hydro includes many
9 complicated physical processes that are nonlinear and must be parameterized. ~~For example, the~~
10 ~~parameters for channel routing are prescribed as functions of stream order, not space; thus the~~The
11 default parameters given by WRF-Hydro ~~are only~~may be valid ~~over a small~~for one region. ~~Because~~
12 ~~channel routing can affect the accuracy of the model performance, but not for another region.~~
13 Hence calibration of related model parameters is often required in order to use the model in a new
14 domain. In particular, for a large spatial domain such as the entire ~~Contiguous~~contiguous United
15 States (~~CONUS~~),₂ in order to develop the optimal parameter sets in a reasonable amount of time,
16 the calibration must be conducted on ~~HPC~~high-performance computing (HPC) systems in parallel
17 instead of in the traditional sequential mode. To date, ~~there is~~ no such calibration tool ~~that can~~
18 ~~straightforwardly~~efficiently calibrate WRF-Hydro on ~~HPCs~~HPC resources. Typically, each
19 physics-based model needs a calibration code that is custom-designed to work ~~with that particular~~
20 ~~numerical model. These custom-designed calibration codes/tools are highly challenging and do~~
21 ~~not offer flexibility; they are designed to operate~~ with that particular numerical model and its set
22 of physics parameterizations, software architecture, and solvers. These custom-designed
23 calibration codes are highly challenging and do not offer flexibility. Therefore, ~~there is a need for~~
24 a more flexible and generic calibration tool is needed that can calibrate any code that uses Message
25 Passing Interface/Open Multi Processing (MPI/OpenMP) for parallelization on ~~HPCs~~HPC
26 systems.

27
28 ~~There are two general types of calibration methods for hydrological models: manual calibration~~
29 ~~and automatic calibration. Models for individual catchments have traditionally been calibrated by~~
30 ~~manually adjusting key model parameters within established ranges of parameters to obtain a best~~
31 ~~match between observed and simulated discharges. This procedure is time consuming, dependent~~

1 ~~on the skill and experience of the modeler, and therefore prone to inconsistency between modelers.~~
2 ~~Automatic calibration is based on stochastic or mathematical methods and thus is more widely~~One
3 widely used generic and independent calibration tool is the parameter estimation tool (PEST).
4 PEST (Doherty, 2016) conducts calibration automatically based on mathematical methods and
5 thus is applicable for optimizing nonlinear parameters. Compared with manual calibration,
6 automatic calibration is more efficient and effective, because it avoids interference from human
7 factors (Madsen, 2000; Getirana, 2010). ~~One widely used automatic calibration tool~~The
8 uniqueness of PEST is Parameter Estimation Tool (PEST; Doherty 2016), which uniquelythat it
9 operates independent of models.~~There; there~~ is no need to develop additional programs/ or codes
10 for a particular model except preparing the files required by PEST (as described in ~~See~~Sec. 3.2;
11 because). PEST ~~works with that model through~~has four modes of operation. One of the model's
12 own inputmodes is regularization mode, which supports the use of Tikhonov regularization and
13 output files. PEST implements a particularly robust variant of the Gauss-Marquardt-Levenberg
14 method (Levenberg, 1944; Marquardt, 1963) to estimate parameters. This method requires a
15 continuous relationship to exist between model parameters and model outputs, but it can normally
16 find the is found better for serving environmental models because, if implemented properly, it
17 supports model predictions of minimum in the objective function in a fairly shorter time period
18 error variance, is numerically stable, and embraces rather than other parameter estimation methods.
19 This is especially important when model runs are lengthy or when many parameters must eschews
20 the heterogeneity of natural systems. Singular value decomposition (SVD) can be calibrated-used
21 as a regularization device to guarantee numerical stability of the calibration problem. Parallel
22 PEST is able to distribute many runs across many computing nodes using master-~~slave~~worker
23 parallel programming. ~~However, to the~~To our best of our knowledge, however, no approach is
24 available that allows users to submit jobs using PEST parallelization to a typical supercomputing
25 facility that uses job scheduling and workload management using such as Simple Linux Utility for
26 Resource Management (SLURM), Portable Batch System (PBS), and Cobalt. A previous study
27 (Senatore et al., 2015) used PEST to calibrate WRF-Hydro over the Crati River Basin in
28 Southernsouthern Italy. ~~However, because~~Because the study area was relatively small, ~~they~~the
29 authors were able to conduct the calibration using PEST in sequential mode. (Alfonso Senatore,
30 personal communication, 2018).

31

~~In this~~This study, ~~we ported~~ aims to (1) port parallel PEST to HPC clusters operated by the U.S. Department of Energy (DOE) and ~~adapted~~adapt it to work with WRF-Hydro. ~~Porting involved writing additional scripts to modify,~~ (2) evaluate the ~~workflow for SLURM, Cobalt, and PBS and developing code to connect parallel PEST to WRF-Hydro.~~ In particular, we aim to (1) calibrate the parameters of WRF-Hydro to improve model performance with realistic values maintaining their physical meanings; (2) speed up calibration for this particular study case and provide the capability to WRF-Hydro users; and (3) explore the scale-up capability of HPC-enabled parallel PEST linked to WRF-Hydro; by calibrating a flood event, and (3) explore the scale-up capability and computational benefits of HPC-enabled parallel PEST by assigning different computing resource to the entire calibration process.

2 Model description

2.1 Study area

The case presented here is one of the worst floods experienced by greater Chicago area in the ~~lastpast~~ three decades, which; the storm occurred on April 18, 2013 (Campos and Wang, 2015). According to the National Weather Service (NWS), the heaviest 24-hour accumulated rainfall during this storm reached 201.4, 171.1, and 136.4 mm across Illinois, Iowa, and Missouri, respectively. The Mississippi River crested at 10.8 m (1.7 m above flood stage), and the Illinois River crested in Peoria, Illinois, at 8.95 m; ~~this~~these river cresting broke the previous record of 8.78 m, set in 1943, and was 4.55 m above the historical normal river stage (NWS, 2013). Campos and Wang (2015) conducted three-domain nested WRF simulations to understand the dynamical and microphysical mechanisms of the event. Our study builds on the smallest domain of that study, which covers the majority of Illinois, Iowa, and Missouri at a spatial resolution of 3 km ~~for the atmospheric and land surface model (Fig. 1).~~(Fig. 1). The domain size is 750 km from west to east and 660 km from south to north.

~~During the 10-day period of this studied case, light to moderate rain occurred on April 8 through 11, 2013, followed by a relatively dry period from April 12 to 14. Then a heavy rain event began on April 15 and peaked on April 18. The heaviest rain band moved east of the study area on April 19. The rainy event ended over the study area on April 20.~~

2.2 WRF-Hydro configuration

~~This study employs WRF-Hydro version 5 with a basic configuration. This configuration does not use nudging techniques or spatially distributed soil-related parameters as used in the National Water Model configuration. WRF-Hydro has been tested in several different cases that focused on different hydrometeorological forecasting and simulation problems (e.g., Gochis et al., 2018; Yucl et al., 2015; Senatore et al., 2015; Arnault et al., 2016), and it shows reasonable accuracy in simulated streamflow after being carefully calibrated. For details of the WRF-Hydro modeling system, see Gochis et al. (2018). WRF-Hydro employs a multiscale modeling approach to handle the local landscape gradient features. Specifically, WRF-Hydro uses a subgrid disaggregation-aggregation procedure. For each time step at which forcing data are available, the column moisture stays within the LSM and is disaggregated from the LSM grid to a high-resolution routing grid (Gochis and Chen 2003). After disaggregation, the routing schemes are executed using the high-resolution grid values. After execution of the routing schemes, the high-resolution grid values are aggregated back to the native LSM grid. For details of each routing component, see Gochis et al. (2015), Yucl et al. (2015), and Senatore et al. (2015).~~

Currently, two LSMs are available in WRF-Hydro for representing land-surface column physics: Noah (Chen and Dudhia, 2001) and Noah Multi-parameterization (~~NoahMP~~Noah-MP; Niu et al. 2011). We utilize NoahMPNoah-MP LSM because compared ~~to~~with Noah LSM it shows obvious improvements in reproducing surface fluxes, skin temperature over dry periods, snow water equivalent, snow depth, and runoff (Niu et al. 2011). ~~Compared to LSM, one major advantage of WRF-Hydro system~~The Noah-MP is that, WRF-Hydro system can keep the infiltration capacity exceedance as ponded water within the model domain. This ponded water is subsequently available for lateral redistribution, which combine the ponded water with new precipitation for calculating the infiltration amount in the next time step~~configured. WRF-Hydro has been tested in several different cases that focused on different hydrometeorological forecasting and simulation problems (e.g., Gochis et al., 2015; Yucl et al., 2015; Senatore et al., 2015; Arnault et al., 2016), and it shows reasonable accuracy in simulated streamflow.~~

~~This study employs WRF-Hydro version 5 with a basic configuration. This configuration does not use nudging technique as used in the National Water Model configuration and spatially distributed~~

1 ~~soil-related parameters. The LSM is~~ at a grid spacing of 3 km, and the aggregation factor is 15;
2 that is, starting from a 3-km LSM resolution in the domain shown in Fig. 1, hydrological routing
3 is performed at a ~~spatial grid~~ resolution of 200 m, ~~with 3285 south-north × 3735 west-east grid~~
4 ~~cells~~. We use a time step of 10 seconds for the routing grid ~~in order~~ to maintain model stability
5 and prevent numerical dispersion of overland flood waves. The time step also meets the Courant
6 condition criteria for diffusive wave routing on a 200-m resolution grid. The WRF-Hydro is
7 configured to be ~~in~~ offline or uncoupled mode—~~—~~there is no online interaction ~~withbetween the~~
8 ~~WRF-Hydro hydrological model and the~~ WRF atmospheric model. ~~SurfaceOverland~~ flow,
9 saturated subsurface flow, gridded channel routing, and a conceptual baseflow are active in this
10 study. The gridded channel network uses an explicit, one-dimensional, variable time-stepping
11 diffusive wave. A direct output-equals-input “pass-through” relationship is adopted ~~here~~ to
12 estimate the baseflow. Although the baseflow module is not physically explicit, it is ~~very~~ important
13 because the water flow in the channel routing ~~areis~~ contributed by both ~~the~~ overland flow and
14 baseflow. If ~~the~~ overland flow is active ~~as it is in this study~~, it passes water directly to the channel
15 model. In this case the soil drainage is the only water resource flowing into the baseflow buckets.
16 ~~IfHowever, if the~~ overland flow is deactivated but channel routing is still active, then WRF-Hydro
17 collects excess surface infiltration water from the land model; and passes this water into the
18 baseflow bucket. This bucket then contributes the water from both overland and soil drainage to
19 the channel flow. Therefore, the baseflow must be active if the overland flow is switched off. This
20 study does not consider lakes and reservoirs.

21
22 We use the geographic information system (GIS) tool ~~that are~~(Sampson and Gochis, 2018)
23 developed by the WRF-Hydro team to delineate the stream channel network, open water (i.e., lake,
24 reservoir, and ocean) grid cells, and groundwater/baseflow basins. Meteorological input for ~~the~~
25 WRF-Hydro model system includes hourly precipitation; near-surface air temperature, humidity,
26 ~~and~~ wind speed; incoming shortwave and longwave radiation; and surface pressure. In this study,
27 the hourly precipitation is from the National Centers for Environmental Prediction (NCEP) Stage
28 IV analysis at a spatial resolution of 4 km. The Stage IV data is based on combined radar and
29 gauge data (Lin and Mitchell, 2005; Prat and Nelson, 2015), and has been shown to be temporally
30 well correlated with high-quality measurements from individual gauges (see, e.g., Sapiano and
31 Arkin, 2009; Prat and Nelson, 2015). The other hourly meteorological ~~inputinputs~~ are from the

1 second phase of the multi-institution North American Land Data Assimilation System project,
2 phase 2 (NLDAS-2) (Xia et al., 2012a,b), at a spatial resolution of 12 km. NLDAS-2 is an offline
3 data assimilation system featuring uncoupled LSMs ~~that are~~ driven by observation-based
4 atmospheric forcing.

5
6 During the 15-day period of this studied case, light to moderate rain occurred on April 8 through
7 11, 2013, followed by a relatively dry period from April 12 to 15. Then a heavy rain event began
8 on April 16 and peaked on April 18. The heaviest rain band moved east of the study area on April
9 19. The rainy event ended over the study area on April 20 (see Fig. S1 in Supporting Information).
10 We start the WRF-Hydro simulation on Jan. 1, 2013, and run the model for more than three months
11 to reach equilibrium. This 3-month period is considered as spin-up time and is excluded from
12 model calibration and evaluation. We calibrate the river discharge calculated by the WRF-Hydro
13 model from 00UTC April 9 to 00UTC April 12, 2013, considering it long enough to achieve our
14 objective. We then evaluate the model performance against U.S. Geological Survey (USGS)
15 observed river discharge from 00UTC April 12 to 00UTC April 25, 2013.

16 **3 Calibration**

17 **3.1 Platforms**

18 We customized parallel PEST to work on ~~two~~three different workload managers and job
19 schedulers: SLURM at the National Energy Research Scientific Computing Center (NERSC), PBS
20 at the Argonne National Laboratory Computing Resource Center, and Cobalt at the Argonne
21 Leadership Computing Facility. The tests presented here are conducted on Edison ~~of~~at NERSC,
22 which uses the SLURM workload manager and job scheduler. Edison is a Cray XC30 with a peak
23 performance of 2.57 petaflops per second, 133,824 compute cores, 357 terabytes of memory, and
24 7.56 petabytes of disk storage. It has 5,586 nodes and 24 cores per node.

25 **3.2 PEST files**

26 ~~Parallel PEST requires four types of input file:~~

27 ~~Template files, which~~

1 The interface we have built between parallel PEST and the management software (SLURM here)
2 is, in general, used for (1) setting the number of workers and the nodes for each worker to conduct
3 a model run (WRF-Hydro here); (2) finding the nodes that are available; (3) setting up the working
4 directory for the workers; (4) identifying the nodes that work for each worker; (5) passing the
5 global files (same for all the working directory) to all the workers (these files include the lookup
6 table files that are not to be calibrated, the namelist files for both LSM and hydrological sector,
7 and restart files that generated by the previous simulations, or spin-up period); and (6) submitting
8 the job for the entire calibration process, including parallel PEST and parallel WRF-hydro. This
9 job can be submitted as a cold-start run or as a restart. The main difference for this interface on
10 different management software is that different management software has its own way to submit
11 jobs and identify available nodes. This difference requires some changes in the script we
12 developed.

13 3.2 PEST files and settings

14 ~~1. PEST requires three file types in both sequential and parallel mode. They are template files~~
15 ~~to define the parameters to be calibrated. For example, we generated CHANNEL.TPL,~~
16 ~~HYDRO.TPL, and GENPARAM.TPL based on the format of their corresponding lookup~~
17 ~~tables, which are CHANNEL.TBL, HYDRO.TBL, and GENPARAM.TBL, respectively.~~
18 ~~CHANNEL.TBL describes the features of a channel, such as bottom width, channel side~~
19 ~~slope, and Manning's roughness coefficients. HYDRO.TBL contains Manning's~~
20 ~~roughness coefficients for land use types. GENPARAM.TPL describes the parameters~~
21 ~~used in the Noah MP LSM.~~

22 ~~2. An, an instruction file, which defines to define the format of model-generated output files.~~
23 ~~For example, WRF Hydro can output time series of streamflow over the forecast points~~
24 ~~(frxst_pts_out.txt) specified during model configuration. The instruction file follows the~~
25 ~~format of frxst_pts_out.txt and specifies the line number of each calibrated forecast point~~
26 ~~in frxst_pts_out.txt.~~

27 ~~3. A, and a control file, which supplies to supply PEST with the size of the problem (e.g.,~~
28 ~~how many parameters to be calibrated; how many observational points); initial parameter~~
29 ~~values and their lower and upper bounds; the increment of each parameter for forward-~~
30 ~~calculation; the names of all template and instruction files; observational values, and~~

1 weight for each parameter to be calibrated. PEST requires all these three file types in both
2 sequential and parallel mode.

3 4. and the settings for the calibration method. Parallel PEST uses a “master-worker” paradigm
4 that starts model runs simultaneously by different workers (or in different folders). The
5 master of parallel PEST communicates with each of its workers many times during a
6 calibration. To run PEST in parallel mode, one also needs a management file to inform
7 PEST where ~~each slave’s~~the working folder is, ~~as well as for each worker and what~~
8 names and paths ~~of~~are for each model input file that PEST must write (i.e., lookup tables
9 that come from template files) and each model output file that PEST must read (such as
10 frsxt_pts_out.txt).

11
12 ~~Parallel PEST uses a “master-slave” paradigm that starts model runs simultaneously in different~~
13 ~~folders (or by different “slaves”). The master of parallel PEST communicates with each of its~~
14 ~~slaves many times during the course of a calibration. When PEST needs to run a~~ The management
15 file also set the maximum running time for each worker. For those workers that take longer than
16 the maximum running time, PEST will stop the model in arun by that particular folder, the master
17 notifies the slave to start the worker and assign that model in that folder. Each slave starts the model
18 execution accordingly, and informs the master that the model starts running. Once the simulation
19 is completed in a particular folder, the slave signals the master, so the mater can read the particular
20 output ~~However, to the best of our knowledge~~run to another worker if there is one with nothing
21 else to do.

22
23 To the best of our knowledge, however, parallel PEST is not designed to run on HPCs directly.
24 We developed scripts and an interface to enable parallel PEST to run on HPCs using SLURM,
25 PBS, or Cobalt workload managers and job schedulers. ~~This enables~~The development involved
26 writing scripts to modify the workflow for different workload managers and job schedulers, as
27 well as developing code to connect parallel PEST to WRF-Hydro. These developments enable
28 parallel PEST to run many ~~slaves on~~workers at the ~~HPC~~same time; each ~~slave~~worker runs a
29 parallel code (~~such as~~here WRF-Hydro) that uses more than one node, which could significantly
30 ~~increase~~reduce the ~~computational performance~~wall-clock time of model calibrations. Although
31 this master-~~slave~~worker parallelism may not be as efficient as a fully MPI approach, it is sufficient

1 for model calibration and requires the least effort for the current parallel PEST to run on HPCs
2 systems.

3
4 This study presents calibration results from PEST using the SVD-based regularization in
5 regularization mode to ensure numerical stability (Tonkin and Doherty, 2005). We focus on
6 calibrating 22 parameters (see Table 1 and detail description in Sec. 3.3) using 96 observation
7 points and 22 items of prior information for the calibrated parameters. In each item of prior
8 information, a value equal to its default value provided by the WRF-Hydro v5.0 (or the log of its
9 default value) is assigned for each adjustable parameter, assuming that default values are the
10 preferred values. All prior information equations are assigned a weight of 1.0. We assigned five
11 different regularization groups to the prior information: Manning’s roughness coefficients
12 specified by Strahler stream order in CHANPARAM.TBL to one group; the parameters in
13 HYDRO.TBL (Manning’s roughness coefficients for overland flow as a function of vegetation
14 types) to another group; and three global parameters for the Noah-MP (xslop1, refdk, and refkdt)
15 in GENPARAM.TBL to the remaining three groups. The 96 observation points are given different
16 weights based on the inversed mean of their observed discharge during the studied period (see the
17 detailed description in Sec. 3.3 and Sec. 4.1). For a detailed description of these settings see the
18 PEST User Manual (Doherty, 2016).

20 **3.3 Calibrated experiments**

21 The primary objective of this study is to ~~present~~build a bridge for linking the ~~operational-parallel~~
22 PEST and WRF-hydro on the scale-up capability ~~basis of HPC clusters and to explore the HPC-~~
23 ~~enabled parallel PEST for use with WRF-Hydro calibration over a relatively large~~
24 ~~domain-computational benefits of this bridge.~~ We ~~focus less on~~do not attempt to extensively
25 ~~assessing the performance of~~assess each individual tool or address questions in each individual
26 domain, such as optimizing the objective functions in PEST or calibrating WRF-Hydro for a long
27 time period considering all the WRF-Hydro model-relevant parameters to achieve an optimal
28 parameter set. The calibration ~~and validation~~period thus is limited to only ~~7~~three days, ~~considering~~
29 ~~it is~~which we believe long enough to achieve our objective and to understand WRF-Hydro’s
30 sensitivity to ~~multiple parameters.~~ ~~The calibration compares WRF-Hydro modeled river discharge~~

1 ~~to U.S. Geological Survey (USGS) surface water observations.~~ the calibrated parameters. We
2 ~~originally choose 11 USGS sites across the study area. However, because of inaccuracies~~
3 ~~introduced when projecting geospatial data from one coordinate system to another by the ArcGIS~~
4 ~~tool, three of the observational sites were not properly assigned to the desired location on the~~
5 ~~channel network. This situation is common in hydrographic data processing and well known to~~
6 ~~hydrologists (Sampson and Gochis, 2018). Among the remaining eight sites, four have~~
7 ~~discontinuous or missing data over the calibration period. Therefore, we~~ calibrated WRF-Hydro
8 using four USGS sites (referred to as Station 1, Station 2, Station 3, and Station 4 hereafter), as
9 shown in Fig. 1 ~~with their site number.~~ (More USGS sites could be included if one manually
10 reallocated the stations that were not properly assigned to the desired location on the channel
11 network by the GIS tool.) We then transfer the calibrated parameters to other ~~sub-basins~~subbasins
12 in the study area to assess the transferability of the calibrated parameters. Although ~~there are~~ many
13 parameters, including spatially distributed parameters and constant parameters in the lookup
14 tables, ~~that~~ affect the model performance, we ~~only~~ calibrate only the parameters in lookup tables
15 and do not consider the spatial variability of ~~each parameter~~other parameters or their scaling
16 factors. We acknowledge that ~~there are some~~ studies ~~that~~ calibrate a single scaling factor (without
17 considering its spatial variability, however) of overland roughness coefficients
18 (OVROUGHRTFAC) rather than the actual value of each land type in the lookup table (e.g.,
19 Kerandi et al., 2018). Although this approach reduces the number of calibrated parameters,
20 however, it has less flexibility because changing one factor will change all the parameters that use
21 the same proportion. ~~In addition, a single scaling factor holds the same for the entire domain,~~
22 ~~which may work well for a small domain, but could be problematic for a large domain. Thus, we~~
23 ~~suggest the calibration of spatially distributed parameters requires more knowledge and~~
24 ~~understanding of the study area and deserves future studies. In this study, we calibrate the~~
25 ~~roughness coefficients for each land type rather than calibrating a single scaling factor.~~

26
27 For ~~most~~the calibration exercises we ~~document~~conduct here, the retention depth factor
28 (RETDEPRTFAC) is fixed at 0.001. This value is reasonable because the ~~modelled~~modeled
29 discharge of our particular configuration (~~See~~Sec. 2.2) using default parameters is ~~much~~ lower
30 than observed discharge. Reducing this factor from 1 to 0.001 keeps less water in water ponds and
31 more water on the surface so it can contribute to river discharge. First, we calibrate 48 parameters

1 based on a 3-day simulation from April 9 to ~~12~~April 11, 2013 (Table S1 in Supporting
2 Information). This calibration uses the estimation mode in the PEST tool and considers equal
3 weight for all four USGS stations. We calibrate~~the~~ Manning's roughness coefficients for both
4 channels and land-use types, the deep drainage (SLOPE), infiltration-scaling parameter
5 (REFKDT), and saturated soil lateral conductivity (REFDK). ~~The~~ Manning's roughness
6 coefficients control the hydrograph shape and the timing of the peaks; the ~~infiltration factor,~~
7 ~~saturated hydraulic conductivity~~SLOPE, REFKDT, and ~~deep drainage~~REFDK control the total
8 water volume. Second, based on the knowledge we learn from the ~~3-day~~48-parameter calibration
9 (see details in ~~Seet~~Sec. 4.1), for the same 3-day period, we ~~redefine~~reduce the number of calibrated
10 parameters from 48 to calibrate22 according to the sensitiveness of the WRF-Hydro model to the
11 adjustable parameters. For example, during the calibration we find that Manning's roughness
12 coefficients for several land types barely change because these land types (e.g., tundra, snow/ice)
13 are not present in the study area. We also learn that even though the calibrated WRF-Hydro
14 parameters can generate discharge results that closely resemble observations, the physical meaning
15 of several parameters are not appropriate because of the wide range of those parameters that we
16 set in the PEST control file. For example, Manning's roughness coefficient for stream order 1
17 (0.199) is calibrated smaller than that for stream order 2 (0.218); the overland roughness
18 coefficients for evergreen needleleaf forest (0.043) and mixed forest (0.023) are calibrated smaller
19 than for cropland/woodland (0.046). Neither of these is true in the real world. We therefore adjust
20 the range of many parameters according to the literature (Soong et al., 2012) to maintain their
21 physical meanings (Table 1). ~~We also extend our calibration period to 7 days to include a heavy~~
22 ~~precipitation period. Although a period of 7 days is still very short compared to the traditional~~
23 ~~calibration period of at least 1 year, we find it provides more appropriate parameter estimation—~~
24 ~~as well as better results of simulated hydrograph shape and the total water volume—than does the~~
25 ~~3-day calibration.~~We find that by using the same absolute weight for all four stations, the
26 calibration helps three stations (Station 2, 3, and 4) with large water volumes to generate more
27 reasonable results than do the default parameters; however, the results for Station 1, which has a
28 relatively small volume of water, is not always better than the discharge that is modeled by using
29 default parameters. Thus, we assign a weight of 15.0 for Station 1 versus a weight of 1.0 for the
30 other three stations according to the inversed mean of observed discharge over these four stations

in April 2013. The ratio of the weights between Station 1 and the other three stations stays similar even if the means are calculated based on different time periods.

3.4 Statistics

This study employs three statistical criteria: Nash–Sutcliffe efficiency (NSE; Nash and Sutcliffe, 1970; Moriasi et al., 2007), root-mean-square error (RMSE), and Pearson correlation coefficient (PCC). RMSE and PCC evaluate model performance in terms of bias and temporal variation. NSE quantitatively describes the accuracy of ~~modelled~~modeled discharge compared ~~to~~with the mean of the observed data. Equation (1) calculates the NSE with defined variables:

$$NSE = 1 - \frac{\sum_{t=0}^n (Y_t^{obs} - Y_t^{sim})^2}{\sum_{t=0}^n (Y_t^{obs} - Y_{mean}^{sim})^2}, \quad (1)$$

where Y_t^{obs} is the t th observed value from USGS sites for river discharge, Y_t^{sim} is the t th simulated value from the WRF-Hydro output, Y_{mean}^{obs} is the temporal average of USGS observed discharge, and n is the total number of observation time points. An efficiency of 1 (NSE = 1) corresponds to a perfect match between modeled discharge and observed data. An efficiency of 0 (NSE = 0) indicates that the model predictions are as accurate as the mean of the observed data. An efficiency below zero (NSE < 0) occurs when the model is worse than the observed mean. Essentially, the closer the NSE is to 1, the more accurate the model is.

4 Results

4.1 ~~Three-day~~WRF-Hydro calibration and validation

Based on the knowledge we gained from the 48-parameter 3-day calibration, we adjust the range of critical parameters in the PEST control file to main their physical meanings. For example, we set Manning’s roughness coefficient larger for stream order 1 than for stream order 2. We also adjust the parameter range of the overland roughness coefficient for multiple land covers, such as forests. We exclude the parameters that WRF-Hydro is not sensitive to for this study, in order to constrain the problem size considering the availability of computational resources. However, if the studied area is much larger with more land types than the study area here, then there would be more parameters to calibrate. Also, hundreds of constant parameters in the Noah-MP model could

1 affect the WRF-Hydro results (Cuntz et al. 2016) and can be calibrated. Both these situations
2 would increase the burden of WRF-Hydro calibration. We perform the same 3-day calibration
3 from April 9 to April 11, 2013. Figure 2 shows the results of the 3-day modeled discharge (in cubic
4 meters) using default and calibrated parameters after five iterations, as well as observed discharge
5 from April 9 to 12-. The four stations are calibrated by considering different weights. Compared
6 with the results calibrated by using equal weights for all the stations, by giving a higher weight to
7 Station 1 the model bias over Station 1 is significantly reduced, with a higher NSE (0.87 with
8 higher weight versus 0.14 with equal weight) and lower RMSE (48.1 versus 123.6). Over Stations
9 2, 3, and 4, which sit on rivers with relatively large water volumes, the modeled discharge using
10 the default parameter underestimates the streamflow by more than 65%. PEST detects this
11 underestimation and immediately adjusts the parameters and increases the modeled discharge
12 during the first iteration. After the ~~fifth~~third iteration, the difference in calibrated results between
13 different iterations is relatively small, ~~and~~. We allow the PEST ~~performed 12~~to conduct five
14 iterations ~~before finding the~~and use the parameters obtained from the fifth iteration as our optimum
15 parameters. ~~Here we only show results generated by default parameters and by parameters~~
16 ~~calibrated from the 1st, 5th, and 12th iterations. Over Stations 2, 3, and 4, which sit on rivers with~~
17 ~~relatively large water volumes, the discharge modeled by the default parameters is much lower~~
18 ~~than discharge seen in observations. PEST detects this underestimation. It immediately adjusts the~~
19 ~~parameters and increases the modeled discharge during the first iteration. After adjusting the~~
20 ~~parameters for several iterations~~As shown in Table 2, when the optimum parameters are used, the
21 modeled ~~discharge gets~~discharges are much closer to the observations compared ~~to the modeled~~
22 ~~results that used the~~with the modeled results when the default parameters were used. The NSEs
23 for the four stations increased from 0.73 (Station 1), -54.4 (Station 2), 157.3 (Station 3) and -
24 1316.9 (Station 4) to 0.87, 0.64, 0.05, and -58.78, respectively, being closer to 1. The RMSEs
25 decreased from 69.3, 3925.2, 3981.3, and 4391.3 m³/sec to 48.1, 318.2, 308.7, and 934.6 m³/sec,
26 respectively. Giving a lower weight for the three large river stations does not change the calibration
27 results much.

28
29 During the validation period, compared with the modeled discharge using default parameters-, as
30 shown in Table 2-~~shows~~, the ~~statistics of model performance using default and calibrated~~
31 parametersNSEs for all four stations ~~during~~are increased to be closer to 1; RMSEs are decreased

1 by 50% or more; and the calibration correlation coefficients between the observed and validation
2 period. Compared to the discharge that was modeled using the default WRF-Hydro parameters,
3 overall, the calibrated modeled discharge matches observations better are increased from 0.8, 0.76,
4 0.21, and 0.72 to 0.98, 0.82, 0.80, and 0.75. Compared with the results of calibration using the
5 estimation mode (no regularization) in PEST, the SVD-based regularization generates slightly
6 better hydrograph shape with 24-hour later discharge peaks that are closer to the observations.
7 However, a problem remains with the hydrograph shapes of the modeled discharge, especially at
8 the three stations with large volumes of with the modeled peak of discharge. For Station 1, the
9 WRF-Hydro almost captures the timing of the peak of discharge, although it still underestimates
10 the water. Note volume by ~25%. The reason is that this 3-day period only experienced a light
11 rain over the study area. The streamflow in the rivers is, therefore, mostly from groundwater and
12 overland flow from upstream or from previous precipitation events. The contribution of overland
13 flow is small for this period because the amount of precipitation was also small, so the main
14 contributor to river discharge in the real situation was from the groundwater. However, in this
15 study WRF-Hydro study uses a direct pass-through groundwater model baseflow module, which
16 does not account for slow discharge and long-term storage of the baseflow. Therefore, Therefore,
17 the largest contribution to river discharge is from precipitation, and groundwater does not
18 contribute much discharge to the channels. This situation causes the model to greatly
19 underestimate discharge, so the calibration adjusts critical parameters aggressively to increase the
20 streamflow to match the observations. When we apply these calibrated parameters to the following
21 days, the calibrated discharge is much higher than the observed discharge during the heavy
22 precipitation period, as shown by Fig. 3 and the RMSEs for the validation period in Table 2.
23 However, after the in a long-term view, as is also true for the other three large river stations.
24 Different from Station 1, for the other three large river stations, the WRF-Hydro modeled
25 discharge increases soon after the peak of precipitation and reaches a peak on April 21, 2013,
26 which is much earlier than the observed peak of river discharge (near April 24). The reason is that
27 the water contributions for these stations are from a larger river basin (Mississippi River) than we
28 included in our current study area. Thus, when a heavy precipitation event, the modeled discharge
29 occurs over the entire river basin, there will be a significant lag time (especially at the lower part
30 of the basin) between the peak of precipitation amount and the peak of river discharge. For
31 example, the precipitation over the upper part of Mississippi River Basin (MRB) has a peak amount

1 on April 18–19, but the river discharge did not reach its peak until April 24. Because our studied
2 area covers only half of the MRB, the modeled river discharge has a shorter delay period after the
3 peak of precipitation than does the observed river discharge. Enlarging the study area to include
4 the entire MRB may improve this situation. Alternatively, calibrating and validating local rivers
5 that are included in the current study area may also reduce the bias in hydrograph shape compared
6 to calibrating and validating large rivers. On the other hand, the WRF-Hydro simulated river
7 discharge decreases soon after it reaches the peak and much faster earlier than in the observed
8 situation (Fig. 3). This discharge. The reason is again might be due to that the direct pass-through
9 groundwater model we adopt in baseflow employed by this study, which uses an output equals
10 input relationship between soil drainage and the discharge into river channels. This model does
11 not allow does not account for slow discharge and long-term storage of baseflow in each
12 conceptual bucket, and thereby not be able to fully represent the baseflow. As a result, the
13 contribution of groundwater to streamflow from the baseflow to the river discharge in model
14 simulations does not stay as long as in real situations. In the observations, the river discharge
15 decreases from the peak at a speed of ~500 m³/sec per day, while the modeled river discharge
16 decreases from the peak at a speed of ~1667 m³/sec per day. Using exponential storage-discharge
17 function for the baseflow may improve this situation.

18
19 ~~From this 3-day calibration experiment, we learn that the WRF-Hydro output is not sensitive to~~
20 ~~several parameters we calibrated in this particular study. For example, Manning's roughness~~
21 ~~coefficients for several land types barely change during the calibration because these land types~~
22 ~~(e.g., tundra, snow/ice) are not present in the study period and area. We also learn that even though~~
23 ~~the calibrated WRF-Hydro parameters can generate discharge results that closely resemble~~
24 ~~observations, the physical meaning of several parameters are not appropriate due to the wide range~~
25 ~~of those parameters that we set in the PEST control file. For example, as shown in Table S1, the~~
26 ~~Manning's roughness coefficient for stream order 1 (0.199) is calibrated smaller than that for~~
27 ~~stream order 2 (0.218); the overland roughness coefficients for evergreen needleleaf forest (0.043)~~
28 ~~and mixed forest (0.023) are calibrated smaller than cropland/woodland (0.046). Neither of these~~
29 ~~is true in the real world.~~

4.2 Seven-day calibration and validation

Based on the knowledge we gained from the 3-day calibration, we adjust the range of critical parameters in the PEST control file. For example, we set the Manning's roughness coefficient larger for stream order 1 than for stream order 2. We also adjust the parameter range of the overland roughness coefficient for multiple land covers, such as forests. With the adjusted range of parameters, we perform 7-day calibration from 00:00 UTC on April 9, 2013, to 00:00 UTC on April 16, 2013, when there is an increased streamflow that the simulation does not capture using 3-day calibrated parameters. The entire calibration takes 12 iterations. Figure 4 shows the results of modeled discharge (in cubic meters) using default and calibrated parameters (from the first, fifth, and 12th iterations), as well as observed discharge from April 9 to 16. Over Stations 2, 3, and 4, the modeled discharge using the default parameter underestimates the streamflow by more than 100%. PEST detects this underestimation and starts adjusting parameters to increase the discharge to match the observations. Compared to the modeled discharge using default parameters during the validation period, as shown in Table 3, the RMSE decreased from 624.9 (Station 1), 5162.9 (Station 2), 4990.0 (Station 3), and 5098.3 m³/sec (Station 4) to 283.1, 637.9, 666.8, and 1202.8 m³/sec, respectively. The correlation coefficient between observed and modeled discharge increased from 0.71, 0.90, 0.87, and 0.82 to 0.97, 0.99, 0.96, and 0.86. Note that, although the calibration helps three stations (Station 2, 3, and 4) with large water volumes to generate more reasonable results than the default parameters, the results for Station 1, which has a relatively small volume of water, is not always better than the discharge that is modeled using default parameters (Tables 2 and 3). This might be because we use the same absolute weight for all the stations when we perform the calibration. Using a higher weight for Station 1 may help improve this situation and generate better results for this station.

Although a period of 7 days is still very short for calibration compared to traditional calibration period of at least 1 year, we find that the 7-day period provides better and more appropriate parameter estimation than does the 3-day calibration, and it does a better job of capturing the hydrograph shape and the total water volume. Comparing the validation statistics between Tables 2 and 3 as well as Figs. 3 and 5, we find the 7-day calibrated parameters generate better results than do the 3-day calibrated parameters for both water volume and the hydrograph shape over the validation period. Compared to the discharge modeled using 3-day calibrated optimum parameters,

1 there is a 17–33% increase in the simulated streamflow (1,400–2,800 m³/sec) calculated over
2 Station 2, 3, and 4 using the 7-day calibrated parameters from April 12 to 16. The RMSE is 3,400–
3 3,600 m³/sec when calculated using the 3-day calibration, but this decreases to 600–1,200 m³/sec
4 when calculated using the 7-day calibration. The correlation coefficient between observed and
5 modeled discharge using the 7-day calibration is 0.8–0.99, but that using the 3-day calibration is
6 only 0.7–0.8. However, there is still a problem with the temporal variability of the modeled
7 discharge, especially over the rivers with large discharge. When there is precipitation, the
8 discharge immediately increases and is higher than the observed discharge. After the precipitation
9 period, when the observed discharge still stays high, the modeled discharge decreases sooner and
10 thus is smaller than the observed discharge. This might be the direct pass-through approach
11 simplified groundwater flow, and does not represent the interaction between stream flow and
12 groundwater properly in this case study.

13 ~~4.3 Evaluation of spatial transferability of the modeling system~~

14 In this section, we apply the calibrated parameters for the four stations (black circles) in Fig. 1 to
15 other 13 stations in the study area. As mentioned before, because of inaccuracies in the spatial
16 location of station data and digital elevation models, and because of small errors introduced when
17 projecting geospatial data from one coordinate system to another using the ArcGIS tool, only four
18 stations are mapped on the river systems (crosses in Fig. 1); others are slightly shifted out of their
19 closest grid cell. One of these four sites (Station 5) is located on a relatively small river, and others
20 are located on larger rivers. The following analyses assess the transferability of the calibrated
21 parameters from particular sites to other sites that are in the study area but not calibrated. ~~The~~
22 ~~assessment compares the observed discharge with the closest grid cells from the discharge output~~
23 ~~of WRF-Hydro.~~ Figure 6 shows the observed and modeled discharge using default and calibrated
24 parameters. Overall, WRF-Hydro's default parameters underestimate the discharge. WRF-Hydro
25 also generates an earlier discharge peak compared to observations over the four stations (Stations
26 5, 6, 7, and 8) in this particular study. The calibrated model results increase the discharge and
27 generate a hydrograph shape that is closer to the observations than the default model results do.
28 The absolute error of simulated discharge decreases by 12.8%, 22.8%, 46.8%, and 49.9%,
29 respectively, over Stations 5 through 8, compared to the default simulated discharge. However,

1 ~~because we did not specifically calibrate these stations based on observations, there are still~~
2 ~~differences between the calibrated results and observations.~~

3 ~~5 Discussion and summary~~

4 ~~5.1 Scale-up capabilities~~

5 4.2 Computational benefits of parallel PEST on HPCs

6 The ability to scale up the calibration of WRF-Hydro by using parallel PEST on ~~HPCs~~HPC
7 systems is determined by two factors: the scale-up capability of ~~WRF-Hydro, parallel PEST~~ and
8 the scale-up capability of ~~PEST~~WRF-Hydro. In ~~the course of~~ calibrating WRF-Hydro, PEST ~~must~~
9 ~~run the WRF-Hydro model many times.~~ PEST~~first~~ makes ~~some~~as many model runs as there are
10 adjustable parameters to calculate Jacobian matrix (Doherty, 2016). The Jacobian matrix has a
11 column for each calibrated parameter and a row for each observation and each item of prior
12 information that set in the PEST control file. These model runs are independent between
13 ~~slaves, workers and can be easily parallelized.~~ Each ~~slave run~~worker runs the model using with
14 temporarily incremented parameters that are defined in the template and control files. ~~These model~~
15 ~~runs can be easily parallelized. However~~Then, PEST ~~also need~~needs to make ~~some other~~additional
16 model runs to test parameter upgrades. ~~These runs are calculated based on updates.~~ Different from
17 the Jacobian runs, these additional runs are performed by using different Marquardt lambdas. ~~The,~~
18 and the search for a Marquardt lambda that achieves the best set of parameters is a serial procedure
19 ~~—what iterative process.~~ The lambda to use for the next run depends on the outcome of the model
20 run conducted using the previously chosen lambda. ~~This in fact is the major bottleneck of~~
21 ~~parallelization of the PEST code.~~ Although serial testing of Marquardt lambdas may quickly find
22 the optimal Marquardt lambda in the first or second series of model runs, it is an inefficient use of
23 computing resources because other processors are idle while only one process is searching the
24 lambdas. This is especially true when the model domain is large and requires extensive computing
25 resources.

26
27 This study employs “partial parallelization” for the lambda-testing procedure (Doherty, 2016), so
28 all the processorsmultiple workers can be used to calculate parameter upgrades based on a series
29 of lambda values that are related to each other by a factor of RLAMFAC set in the PEST control

1 file. This partial parallelization makes the scale up challenging when more processors are in use,
2 because generating many Marquardt lambdas does not always guarantee that the best Marquardt
3 lambdas were the ones generated. As a result, the calibration process may converge more slowly
4 when using more slaves than it does when using less slaves. We tested different numbers of slaves
5 (35, 50, 70, and 105) for the 32-parameter calibration experiment. In total, each of these tests uses
6 71, 101, 141, and 211 nodes; two nodes for each slave run WRF-Hydro, and one node runs PEST
7 master to coordinate jobs and communicate with the slaves. The results shown in Figs. 2-6 and
8 Tables 2-3 are from a calibration using 35 slaves; PEST conducted 12 iterations before finding
9 the optimum parameters. We find that using different numbers of slaves generates slightly different
10 parameter values and involves different numbers of iterations. For example, using 70 slaves only
11 takes eight iterations and 41% of the wall-clock time 35 slaves used to find the optimum
12 parameters. The calibrated parameters are slightly different from those generated by 35 slaves
13 (Table 1, last column), and they generate slightly better results than does the 35-slave test
14 compared to observed discharge. We finish the 12 iterations using 35 slaves (71 nodes) within 73
15 hours, and the eight iterations using 70 slaves (141 nodes) within 30 hours. More than 800 model
16 runs were conducted for entire calibration process including calculating the Jacobian Matrix as
17 well as testing the parameter upgrades. In fact if more nodes were used by each slave for the
18 calculation, the wall-clock time can be further reduced. If these calibration were conducted
19 sequentially on personal computers, the same calibration process would have taken 60-80 days for
20 a 7-day calibration. We also set the value of PARLAM to -9999 in the management file so only
21 one cycle of parallel WRF-hydro runs is devoted to testing Marquardt lambdas. For additional
22 details on these parameters and their settings see the PEST User Manual (Doherty, 2016).

23
24 **Our**

25 In this study finds that, depending on we test the computational performance of HPC-enabled
26 parallel PEST using different number of parameters being calibrated (e.g., 32 parameters in this
27 particular study), using 32 to 64 slaves shows fairly good scale-up capability; most of the time
28 consumed by PEST is for running WRF-Hydroworkers (6, 12, and the number of slaves can be
29 used to carry out model runs to generate the Jacobian matrix. However, using 105 slaves (211
30 nodes) does not result in fewer iterations or a shorter wall-clock time than using 70 slaves. In fact,
31 using more than 64 slaves may not be necessary because generating many 23) for the 22-parameter

1 calibration. As shown in Table 3, we conducted five experiments: Test 1 uses 23 workers, Test 2
2 uses 12 workers, and Test 3 uses 6 workers. All three tests use two nodes for each worker to run
3 WRF-Hydro in parallel. The maximum number of lambda-testing runs undertaken per iteration is
4 set to 15, 10, and 5 for Test 1, 2, and 3, respectively, to make sure that only one cycle of WRF-
5 hydro runs is devoted (using 15, 10 and 5 workers from Tests 1, 2, and 3, respectively) to testing
6 Marquardt lambdas. Note that the maximum number of lambda-testing runs should be set equal to
7 or less than the workers available. Otherwise, another cycle of WRF-hydro runs needs to be
8 conducted. In fact, generating more Marquardt lambdas does not always guarantee ~~generating~~that
9 the best Marquardt lambdas. ~~In addition, at least for~~ are generated. In contrast, it may make the
10 calibrations conducted in this study, in each iteration model convergence slower (here, PEST ~~runs~~
11 the) or even model ~~either 32 or 64~~failure.

12
13 In order to test the trade-offs between the computing nodes used for running parallel WRF-Hydro
14 and the workers used for running parallel PEST, Tests 4 and 5 use different number of nodes for
15 each worker to run WRF-Hydro in parallel. Explicitly, Test 4 uses four nodes per worker, and Test
16 5 uses six nodes per worker. Both tests use six workers for running the parallel PEST. The
17 maximum number of lambda-testing runs undertaken per iteration is set to five for both Tests 4
18 and 5. Note that the time costs in Table 3 are limited to only one iteration. Conducting more
19 iterations will increase the cost of wall-clock time and computing, but will not change the
20 conclusion for the scale-up capability and computational benefits for HPC-enabled parallel PEST
21 linked to WRF-hydro.

22
23 PEST needs to run the WRF-Hydro model at least as many times as the number of calibrated
24 parameters (22 here). In fact, PEST runs the model 23 times in the first round (or the first iteration)
25 with initial parameter values and for the first Jacobian matrix. From the second iteration, it runs
26 the model 22 times to calculate ~~the Jacobin~~Jacobian matrix. Therefore, ~~having~~if there are fewer
27 than 23 workers, the time cost for the first round of Jacobian matrix calculation will increase
28 accordingly. For example, as shown in Fig. 4a, when we assign 12 (and 6) workers to parallel
29 PEST, the time cost for calculating the Jacobian matrix is increased by a factor of 2 (and 4)
30 compared with the time cost of using 23 workers. The time cost for the parameter upgrade stays
31 similar for the three experiments because only one cycle of WRF-hydro simulation is conducted

1 to test the Marquardt lambdas. As a result, the total time cost for Test 2 is ~1.5 times more than 64
2 slaves that for Test 1, and the total time cost for Test 3 is ~1.5 times more than that for Test 2 (Fig.
3 4b). By extrapolating the speedup curve shown in Fig. 4a and Fig. 4b, we expect the total time cost
4 to be ~1516 minutes when using only one worker (or sequential mode), which is about 15 times
5 slower compared with running the PEST in parallel mode using 23 workers. For this particular
6 study with 22 adjustable parameters, we expect the time cost most likely to stay the same even if
7 one increases the number of workers to more than 23, because PEST runs WRF-Hydro only 23 or
8 22 times for each iteration. Assigning more workers for this particular study would most likely
9 render some slavesworkers idle and is not an efficient use of computing resources. PEST may run
10 WRF-Hydro more than 22 times (e.g., 44 times) if higher-order finite differences are employed.
11 In this case, assigning more workers (e.g. 45 workers) may further speed up the calibration process.
12 On the other hand, for the same case study and using the same number of nodes for running parallel
13 WRF-Hydro, we can estimate the computing speedup by assuming an increase in the number of
14 calibrated parameters to 50. This would be the case, for example, to evaluate model sensitiveness
15 to the physics in Noah-MP or the spatial variabilities of certain parameters. We then expect to use
16 51 workers to achieve the best computing performance for parallel PEST. This would then be 28–
17 30 times faster than running PEST using one worker (or in sequential mode). Similarly, if 100
18 parameters were used for the calibration for the same case study, a factor of up to 60 speedup in
19 the calibration process would be achieved by running HPC-enabled parallel PEST.

20
21 In addition, by increasing the number of nodes for each worker to conduct WRF-Hydro (Tests 3,
22 4, and 5), the time cost for the entire calibration process is significantly reduced (Figs. 4c and 4d).
23 Specifically, the WRF-hydro scales up well when using four and six nodes compared with using
24 two nodes per worker for running the WRF-Hydro. Both the time spent on calculating the Jacobian
25 matrix and the time spent on testing the parameter upgrades are decreased by 49% and 67%,
26 respectively, when using four and six nodes. Therefore, the total time spent is also decreased when
27 using more nodes for each worker (see Table 3). Increasing the number of nodes to eight for each
28 worker will most likely further decrease the time cost by 70–75% compared with using only two
29 nodes per worker. Moreover, if one has a larger study area such as the entire contiguous United
30 States, we expect the WRF-Hydro to have an even better scale-up capability (e.g., on dozens of
31 nodes) than this study. Overall, based on the experiments we conduct here, using 23 workers for

1 parallel PEST and six nodes for each worker to run parallel WRF-Hydro would cost the least wall-
2 clock time—about 32 min for one iteration for this particular study.

4 **4.3 Evaluation of spatial transferability of the calibrated parameters**

5 To assess the transferability of the calibrated parameters, we apply the optimum parameters
6 obtained from the calibration for the four stations (black circles) in Fig. 1 to another set of four
7 stations (crosses in Fig. 1) in the study area. All four sites are located on relatively small rivers, so
8 the lag time between precipitation peak and the discharge peak are much shorter than that for the
9 stations on the lower part of MRB (e.g., Stations 2, 3, and 4). The assessment compares the
10 observed discharge with the closest grid cells from the discharge output of WRF-Hydro. 5-2 Figure
11 5 shows the observed and modeled discharge using default and the optimum parameters. Overall,
12 WRF-Hydro’s default parameters underestimate the discharge and misrepresent the timing of
13 discharge peaks compared with observations over the four assessed stations (Stations 5, 6, 7, and
14 8). By using the calibrated parameters from other sites over the area, the model results increase the
15 discharge and shift the hydrograph shape so they are much closer to the observations than model
16 results using default parameters. The absolute error of simulated discharge decreases by 13.1%,
17 38.3%, and 71.6%, respectively, over Stations 6 through 8 (Station 5 shows a 6% increase of
18 absolute error), compared with the default simulated discharge. We also find that using the SVD-
19 based regularization for the PEST calibration captures the timing of discharge peak better than
20 using the estimation mode, which is one-day earlier than the observations reaching the discharge
21 peak.

22 **5 Summary and discussion**

23 WRF-Hydro is a new, and perhaps the first practical, computer code that can run on ~~HPCs~~HPC
24 systems and can model the entire hydrological cycle using physics-based ~~sub-models~~submodels
25 and-very high-resolution input datasets (e.g., radar). The hydrological community has desired this
26 capability for decades, although it requires intensive computing resources. Thus, the calibration of
27 this model would ideally be conducted on HPCs in parallel as well, especially when the model
28 covers a large domain rather than the basin scale. This study ports an independent model
29 calibration tool, parallel PEST, to HPC clusters and links it to WRF-Hydro to help WRF-Hydro

1 users calibrate the model within a much shorter wall-clock time period. ~~This tool's uniqueness lies~~
2 ~~in its flexibility~~The bridge we build here (between parallel PEST and ~~robustness to calibration any~~
3 ~~parameters in~~WRF-Hydro. ~~It is also unique in its use~~ on the basis of ~~two levels of parallelization~~
4 ~~across many slaves running PEST, with each slave running a simulation of WRF-Hydro. The~~
5 ~~calibration tool presented in this study also applies~~HPC systems) can be applied to any other
6 hydrological models and ~~similar earth~~Earth system models that use parameterization
7 ~~to parameterizations to represent~~ model physics. We present the operational feasibility of the HPC-
8 enabled parallel PEST by evaluating the performance of calibrated WRF-Hydro against
9 observation in hydrograph features such as volume and timing of flood events. We examine the
10 scale-up capability and computational benefits of the tool by assigning different computing
11 resource for PEST and for WRF-Hydro. While this study presents the optimum parameters
12 identified from the calibration of ~~this~~the particular ~~study case and area, but the calibrated~~flood
13 event, the parameters can be significantly different if one uses different physics, such as
14 exponential storage-discharge function for a groundwater model, ~~or reach-based channel routing.~~
15 or reach-based channel routing. Our preliminary testing shows that using exponential storage-
16 discharge function with the default parameters provided by WRF-Hydro, the modeled discharge
17 was larger than that of observations. Thus, the calibration will need to adjust the parameters to
18 reduce the discharge. Our study finds that for calibrating 22 parameters, using the same computing
19 resource for running WRF-hydro, the HPC-enabled PEST calibration tool can speed up WRF-
20 Hydro calibration by a factor of 15, compared with running PEST in sequential mode. The speedup
21 factor can be larger when the number of parameters needing calibration is higher (e.g., 50 or 100).

22
23 ~~We apply the HPC-enabled parallel PEST to WRF-Hydro to investigated a major flood event that~~
24 ~~occurred over the Midwestern United States in April 2013. Our precipitation inputs were derived~~
25 ~~from a radar, gauge, and satellite rainfall product named Stage IV. The calibrated parameters~~
26 ~~include Manning's roughness coefficients for both channels and land-use types, deep drainage, the~~
27 ~~infiltration scaling parameter, and saturated soil lateral conductivity. We evaluated the~~
28 ~~performance of WRF-Hydro in hydrograph features such as volume and timing of flood events.~~
29 ~~We also assessed the spatial transferability of the calibrated parameters in the study area.~~

1 The following are ~~the primary findings of several key points that we would like to mention to~~
2 inform future studies:

3 1. In this study, we consider using the prior or regularization information only for the
4 parameters that we calibrate. As is the case with solving inverse problems, prior
5 information is added to improve the smoothness of the solutions. In order to build a more
6 comprehensive calibration, an important aspect that can be considered is to enrich the prior
7 with the available historical data. For example, in this particular case, one can use the
8 historical observation data (e.g., April and May from the past few years) to enrich the prior
9 information for the parameters. Hence, the regularization objective function in PEST will
10 constitute not only the discrepancies between parameters and their “current estimates” but
11 also the discrepancies between WRF-Hydro simulations and preferred values (which is the
12 observed time series of historical discharge). Additionally, one can use the pilot points
13 technique described by Doherty (2005) in conjunction with parameter estimation to add
14 more flexibility to the calibration process. This will be potentially beneficial in improving
15 the predictions.

16 1. ~~For this particular study, the HPC-enabled PEST calibration tool can speed up WRF-Hydro~~
17 ~~calibration by a factor of 30, compared to a serial calibration procedure.~~

18 2. ~~Calibrated WRF-Hydro improves the modeled hydrographs compared to the default model~~
19 ~~results. The RMSE of discharge over the three large stations are reduced by 76–86% with~~
20 ~~calibration for the validation period.~~

21 3. ~~Although the calibration period in this study is relatively short, we found that the longer~~
22 ~~the calibration period is, the better the model results are when compared to observations.~~
23 ~~It is difficult to precisely define what length of data is sufficient to identify model~~
24 ~~parameters so that they can also be used for other periods, because different models have~~
25 ~~different levels of complexity and different catchments have different information content~~
26 ~~in each year of hydrological record. Because of the heavy computation load for the~~
27 ~~calibration of WRF-Hydro, it would be challenging to calibrate yearlong time series.~~

28 4. ~~Although there are inaccuracies when mapping the USGS stations onto the river systems,~~
29 ~~we found the WRF-Hydro calibrated parameters are helpful for nearby locations. The~~
30 ~~absolute bias over the four assessed stations decreased by 12–50% as a result of using the~~

1 calibrated parameters, compared to using the default parameters for their simulations. The
2 mapping issue are often random, or non-systematic; there is no generalizable way to
3 automate the correction procedure with a high degree of fidelity. Manual manipulation or
4 specification of the data is often required (Sampson and Gochis 2018).

5 ~~5. Using different groundwater models can generate very different results and will require a~~
6 ~~completely different set of parameters for WRF Hydro to model the observed discharge.~~
7 ~~Our preliminary testing shows that, using exponential storage discharge function with the~~
8 ~~default parameters provided by WRF Hydro, the modeled discharge was larger than~~
9 ~~observations. Thus, the calibration will need to adjust the parameters to reduce the~~
10 ~~discharge.~~

11
12 2. To focus on our main goal, we calibrate only the parameters in lookup tables. However,
13 we acknowledge that using a single value to represent a physics for a large domain could
14 be problematic, especially we expect the HPC-enabled parallel PEST to execute with
15 WRF-Hydro for large domains. This situation often needs parameter regionalization. For
16 example, WRF-Hydro v5 has many spatially distributed parameters available, such as the
17 overland flow roughness scaling factor (OVROUGHRTFAC), the factor of maximum
18 retention depth (RETDEPRTFAC), and the soil-related parameters (when compiled with
19 SPATIAL_SOIL=1). Calibrating these spatial parameters based on grid scale (e.g.,
20 catchments) rather than a single value will give the model more flexibility and thus better
21 fit the observations (Hundecha and Bardossy, 2004; Wagener and Wheeler, 2006). In
22 practice, for example, one can include regional OVROUGHRTFACs (e.g., their
23 lower/upper bounds, and default values) in the PEST control file based on catchments.
24 However, the selection of the locations and sizes of catchment may introduce significant
25 uncertainties to the calibration results, which require systematic and comprehensive
26 investigation and understanding of the study area.

27 3. This study is limited to calibrating the observed streamflow only based on the format of
28 one of WRF-Hydro model outputs for individual station or point (frxst_pts_out.txt). It is
29 feasible, however, to calibrate other variables as long as the observation data is available.
30 For example, one can either find the closest point from the gridded dataset to the
31 observation location and then compare that model grid to observations; or one can change

1 the WRF-Hydro input/output code to output other variables in the frxst_pts_out.txt file, so
2 they can still use the same interface we developed here to calibrate other variables instead
3 in addition to the discharge.

4 4. The optimal parameter set obtained from this study is from the 5th iteration of parallel
5 PEST by testing five Marquardt lambdas. Testing different number of lambdas or
6 calibrating different number of parameters may generate a different set of optimal
7 parameters. These parameter sets can all make physical sense and be equally good for
8 reproducing observed discharges. This problem is named equifinality (Beven and Freer,
9 2001; Savenije, 2001), which is an important source of model uncertainty. To reduce the
10 model uncertainty through reducing the equifinality, hydrologists carry out additional
11 modelling objective for model evaluation to find more useful parameter sets (Mo and
12 Beven, 2004; Gallart et al., 2007). Alternatively, inspired by No. 3 discussed above, one
13 can calibrate the WRF-hydro model based on more than one variables, such as discharge
14 and soil moisture (or heat flux or water table depth) to reduce the number of optimal
15 parameter sets, and thus reduce the model uncertainty of predictions for these variables.

16 5. While this study ported the parallel PEST to HPC system and linked it to WRF-Hydro, we
17 note that BEOPEST is available in the PEST family. BEOPEST has the same functionality
18 as parallel PEST but uses a different approach for communication between master and
19 workers. Working with HPC-enabled BEOPEST may save total time cost since BEOPEST
20 uses the Transmission Control Protocol (TCP) and the Internet Protocol (IP) instead of
21 message files (reading input and writing output between master and works) for
22 communication. We expect it to be relatively straightforward to use BEOPEST to calibrate
23 WRF-hydro on HPCs since the interface remains similar, except one needs to copy the
24 template and instruction files in addition to the global files (see Section 3.1) into each
25 working folder.

26
27 *Data and Code availability.* The observed river discharge is downloaded from the USGS Surface-
28 Water Data website, available at <https://waterdata.usgs.gov/nwis/sw>. The Stage IV precipitation
29 data were downloaded from <https://data.eol.ucar.edu/dataset/21.093>. PEST was downloaded from
30 <http://www.pesthomepage.org/Downloads.php>. We use the Unix PEST version 13.6. The scripts
31 and files that are developed in this study and required by PEST for calibrating WRF-Hydro are

1 available at https://www.zenodo.org/record/1490230#.W_XI6TFRdhE. DOI:
2 [10.5281/zenodo.1490230](https://doi.org/10.5281/zenodo.1490230). [http://doi.org/10.5281/zenodo.2588506](https://doi.org/10.5281/zenodo.2588506).

3
4 *Author contributions.* JW proposed the project and developed the study case in WRF and WRF-
5 Hydro. CW developed the scripts/code to port the parallel PEST to DOE supercomputers and adapt
6 it to work with WRF-Hydro. [VR provided important input for the regularization calibration](#)
7 [method](#). AO operated the ArcGIS tool to delineate the high-resolution grid cells to include stream
8 channel network, open water, and groundwater/baseflow basins. RK provide high-level guidance
9 and insight for the entire project. All authors commented on this manuscript.

10
11 *Competing interests.* The authors declare that they have no conflict of interest

12
13 ~~*Acknowledgements*~~*Acknowledgments.* This work is supported under a Laboratory Directed
14 Research and Development (LDRD) Program at Argonne National Laboratory, through U.S.
15 Department of Energy (DOE) contract DE-AC02-06CH11357. Computational resources are
16 provided by the DOE-supported National Energy Research Scientific Computing Center, Argonne
17 National Laboratory Computing Resource Center, and Argonne Leadership Computing Facility.
18 Our special thanks to [the PEST developer John Doherty](#) and ~~the~~ entire WRF-Hydro
19 team, especially Kevin Sampson, for his guidance on the ArcGIS tool. [We gratefully thank the two](#)
20 [reviewers for their valuable comments and suggestions, which tremendously improved this](#)
21 [manuscript](#).

22 **References**

23 Arnault, J., Wagner, S., Rumlmer, T., Fersch, B., Bliefernicht, J., Andresen, S., and Kunstmann,
24 H.: Role of runoff-infiltration partitioning and resolved overland flow on land-atmosphere
25 feedbacks: A case study with the WRF-Hydro coupled modeling system for West Africa, J.
26 Hydrometeorol., 17, 1489–1516, 2016.

27
28 Campos, E., and Wang, J.: Numerical simulation and analysis of the April 2013 Chicago Floods,
29 J. Hydrol., 531, 454–474, 2015.

1 Chen, F. and Dudhia, J.: Coupling an advanced land surface-hydrology model with the Penn State-
2 NCAR MM5 modeling system. Part I: Model implementation and sensitivity, *Mon. Weather Rev.*,
3 129, 569–585, 2001.

4

5 [Cuntz, M., Mai, J., Samaniego, L., Clark, M., Wulfmeyer, V., Branch, O., Attinger, S., and Thober,](#)
6 [S.: The impact of standard and hard-coded parameters on the hydrologic fluxes in the Noah-MP](#)
7 [land surface model, *J. Geophys. Res. Atmos.*, 121, 10,676–10,700, doi:10.1002/2016JD025097,](#)
8 [2016.](#)

9

10 Doherty, J.: PEST: Model Independent Parameter Estimation, User Manual, 6th ed., Watermark
11 Numerical Computing, Brisbane, Queensland, Australia, 2016.

12

13 [Doherty, J.: Ground water model calibration using pilot points and regularization, *Groundwater*,](#)
14 [41\(2\), 170–177, 2005.](#)

15

16 [Gallart, F., Latron, J., Llorens, P., and Beven, K. J.: Using internal catchment information to reduce](#)
17 [the uncertainty of discharge and baseflow predictions. *Adv. Water Resour.* 30\(4\), 808–823, 2007.](#)

18

19 Getirana, A. C. V.: Integrating spatial altimetry data into the automatic calibration of hydrological
20 models, *J. Hydrol.*, 387 (3-4), 244–255, doi: 10.1016/j.jhydrol.2010.04.013, 2010.

21

22 ~~Gochis, D. J. and Chen, F.: Hydrological enhancements to the community Noah land surface~~
23 ~~model, NCAR Technical Note NCAR/TN-454+STR, doi: 10.5065/D60P0X00, 2003.~~

24

25 ~~Gochis, D.J., M.,~~ Barlage, ~~A.M.,~~ Dugger, ~~K.A.,~~ FitzGerald, ~~L.K.,~~ Karsten, ~~M.L.,~~ McAllister, ~~J.M.,~~
26 McCreight, J., Mills, ~~A.J.,~~ RafieeiNasab, ~~L.A.,~~ Read, ~~K.L.,~~ Sampson, ~~D.K.,~~ Yates, ~~W.D.,~~ and Yu,
27 ~~W.: (2018).~~ The WRF-Hydro modeling system technical description, (Version 5.0). NCAR
28 Technical Note. 107 pages. Available online at:
29 <https://ral.ucar.edu/sites/default/files/public/WRFHydroV5TechnicalDescription.pdf>, 2018.

30

1 [Hundecha, Y., and Bárdossy, A.: Modeling of the effect of land use changes on the runoff](#)
2 [generation of a river basin through parameter regionalization of a watershed model, J. Hydrol.,](#)
3 [292, 281–295, 2004.](#)

4
5 Kerandi, N., Arnault, J., Laux, P., Wagner, S., Kitheka, J., and Kunstmann, H.: Joint atmospheric-
6 terrestrial water balances for East Africa: A WRF-Hydro case study for the upper Tana River basin,
7 Theor. Appl. Climatol., 131, 1337–1355, doi: 10.1007/s00704-017-2050-8, 2018.

8
9 ~~Levenberg, K.: A method for the solution of certain non-linear problems in least squares, Q. Appl.~~
10 ~~Math., 2(2), 164–168, 1944.~~

11
12 Lin, Y., and Mitchell, K. E.: The NCEP stage II/IV hourly precipitation analyses: Development
13 and applications, Preprints, 19th Conf. on Hydrology, San Diego, CA, Amer. Meteor. Soc., 1.2.,
14 2005.

15
16 Madsen, H.: Automatic calibration of a conceptual rainfall–runoff model using multiple
17 objectives, J. Hydrol., 235, 276–288, 2000.

18
19 ~~Marquardt, D. W.: An algorithm for least squares estimation of non-linear parameters, J. Soc.~~
20 ~~Indust. Appl. Math., 11(2), 431–441, 1963.~~

21
22 [Mo, X., and Beven, K.: Multi-objective parameter conditioning of a three-source wheat canopy](#)
23 [model. Agricultural & Forest Meteorol. 122\(1–2\), 39–63, 2004.](#)

24
25 Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D., and Veith, T. L.:
26 Model evaluation guidelines for systematic quantification of accuracy in watershed simulations,
27 Transactions of the ASABE, 50 (3), 885–900, 2007.

28
29 Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models, part I — A
30 discussion of principles, J. Hydrol., 10(3), 282–290, doi: 10.1016/0022-1694(70)90255-6, 1970.

1 Niu, G.-Y., Yang, Z.-L., Mitchell, K. E., Chen, F., Ek, M. B., Barlage, M., Kumar, A., Manning,
2 K., Niyogi, D., Rosero, E., Tewari, M., and Xia, Y.: The community Noah land surface model with
3 multiparameterization options (Noah-MP): 1. Model description and evaluation with local-scale
4 measurements, *J. Geophys. Res.*, 116, D12109, doi: 10.1029/2010JD015139, 2011.

5

6 NWS (National Weather Service): Record ~~River Flooding~~river flooding of April 2013,
7 <https://www.weather.gov/ilx/apr2013flooding>, 2013.

8

9 ~~Papathanasiou, C., Makropoulos, C., and Mimikou, M.: Hydrological modelling for flood~~
10 ~~forecasting: Calibrating the post fire initial conditions, *J. Hydrol.*, 529, 1838–1850, doi:~~
11 ~~10.1016/j.jhydrol.2015.07.038, 2015.~~

12

13 Prat, O. P., and Nelson, B. R.: Evaluation of precipitation estimates over CONUS derived from
14 satellite, radar, and rain gauge data sets at daily to annual scales (2002-2012), *Hydrol. Earth Syst.*
15 *Sci.*, 19, 2037–2056, doi: 10.5194/hess-19-2037-2015, 2015.

16

17 Sampson, K., and Gochis, D.: WRF Hydro GIS Pre-~~Processing Tools~~processing tools, Version
18 5.0 Documentation, 2018.

19

20 Sapiano, M. R. P., and Arkin, P.A.: An intercomparison and validation of high-resolution satellite
21 precipitation estimates with 3-hourly gauge data, *J. Hydrometeorol.*, 10, 149–166, doi:
22 10.1175/2008JHM1052.1, 2009.

23

24 Senatore, A., Mendicino, G., Gochis, D. J., Yu, W., Yates, D. N., and Kunstmann, H.: Fully
25 coupled atmosphere-hydrology simulations for the central Mediterranean: Impact of enhanced
26 hydrological parameterization for short and long time scales, *J. Adv. Model. Earth Syst.*, 7(4),
27 1693–1715, doi: 10.1002/2015MS000510, 2015.

28

29 Soong, D. T., Prater, C. D., Halfar, T. M., and Wobig, L. A.: Manning’s roughness coefficients for
30 Illinois streams, U.S. Geological Survey Data Series 668, 2012.

31

1 [Tonkin, M. J., and Doherty, J.: A hybrid regularized inversion methodology for highly](#)
2 [parameterized environmental models, *Water Resource Research*, 41, W10412,](#)
3 [doi:10.1029/2005WR003995, 2005.](#)

4
5 [Wagener, T., and Wheater, H. S.: Parameter estimation and regionalization for continuous](#)
6 [rainfall-runoff models including uncertainty, *J. Hydrol.*, 320, 132–154, 2006.](#)

7
8 Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H.,
9 Meng, J., Livneh, B., Lettenmaier, D., Koren, V., Duan, Q., Mo, K., Fan, Y., and Mocko, D.:
10 Continental-scale water and energy flux analysis and validation for the North American Land Data
11 Assimilation System project phase 2 (NLDAS-2)-), 1: Intercomparison and application of model
12 products, *J. Geophys. Res.*, 117, D03109, doi: 10.1029/2011JD016048, 2012a.

13
14 Xia, Y., Mitchell, K., Ek, M., Cosgrove, B., Sheffield, J., Luo, L., Alonge, C., Wei, H., Meng, J.,
15 Livneh, B., Duan, Q., and Lohmann, D.: Continental-scale water and energy flux analysis and
16 validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2). 2.
17 Validation of model-simulated streamflow, *J. Geophys. Res.*, 117, D03110, doi:
18 10.1029/2011JD016051, 2012b.

19
20 Yucel, I., Onen, A. Yilmaz, K. K., and Gochis, D. J.: Calibration and evaluation of a flood
21 forecasting system: Utility of numerical weather prediction model, data assimilation and satellite-
22 based rainfall, *J. Hydrol.*, 523, 49–66, 2015.

23
24 ~~Zanon, F., Borga, M., Zoccatelli, D., Marchi, L., Gaume, E., Bonnifait, L., Delrieu, G.:~~
25 ~~Hydrological analysis of a flash flood across a climatic and geologic gradient: The September 18,~~
26 ~~2007 event in Western Slovenia, *J. Hydrol.*, 394, 182–197, doi: 10.1016/j.jhydrol.2010.08.020,~~
27 ~~2010.~~

1 **Table 1: Calibrated 22 parameters and the optimum parameters found after five iterations.**

Calibrated Parameter	Default	Lower Bound	Upper Bound	Optimum Parameter
mannn1	0.55	0.35	0.6	0.517599
mannn2	0.35	0.15	0.35	0.153894
mannn3	0.15	0.08	0.15	8.00E-02
mannn4	0.1	0.05	0.15	5.00E-02
mannn5	7.00E-02	0.02	0.1	6.677379E-02
mannn6	5.00E-02	0.015	0.1	1.628244E-02
mannn7	4.00E-02	0.01	0.08	1.298054E-02
mannn8	3.00E-02	0.005	0.06	5.00E-03
xslope1	0.1	1.00E-04	1	0.496680
refdk	2.00E-06	1.00E-08	1.00E-05	2.899043E-07
refkdt	1	0.01	5	1.66664
ovn1 (urban)	2.50E-02	0.005	0.06	6.00E-02
ovn2 (dry crop)	3.50E-02	0.015	0.06	1.50E-02
ovn3 (irrigated crop)	3.50E-02	0.015	0.06	1.50E-02
ovn5 (crop/grass)	3.50E-02	0.015	0.06	2.822497E-02
ovn6 (crop/wood)	6.80E-02	0.035	0.25	4.568903E-02
ovn7 (grass)	5.50E-02	0.015	0.25	1.50E-02
ovn10 (savanna)	5.50E-02	0.015	0.3	1.50E-02
ovn11 (deciduous forest)	0.2	0.1	0.3	0.30
ovn14 (evergreen forest)	0.2	0.1	0.3	0.164557
ovn15 (mixed forest)	0.2	0.1	0.3	0.112134
ovn16 (water)	5.00E-03	0.001	0.01	1.00E-02

1 **Table 2: Statistics of model performance using optimum and default (in parentheses)**
 2 **parameters for Stations 1–4 during the calibration and validation period.^a**

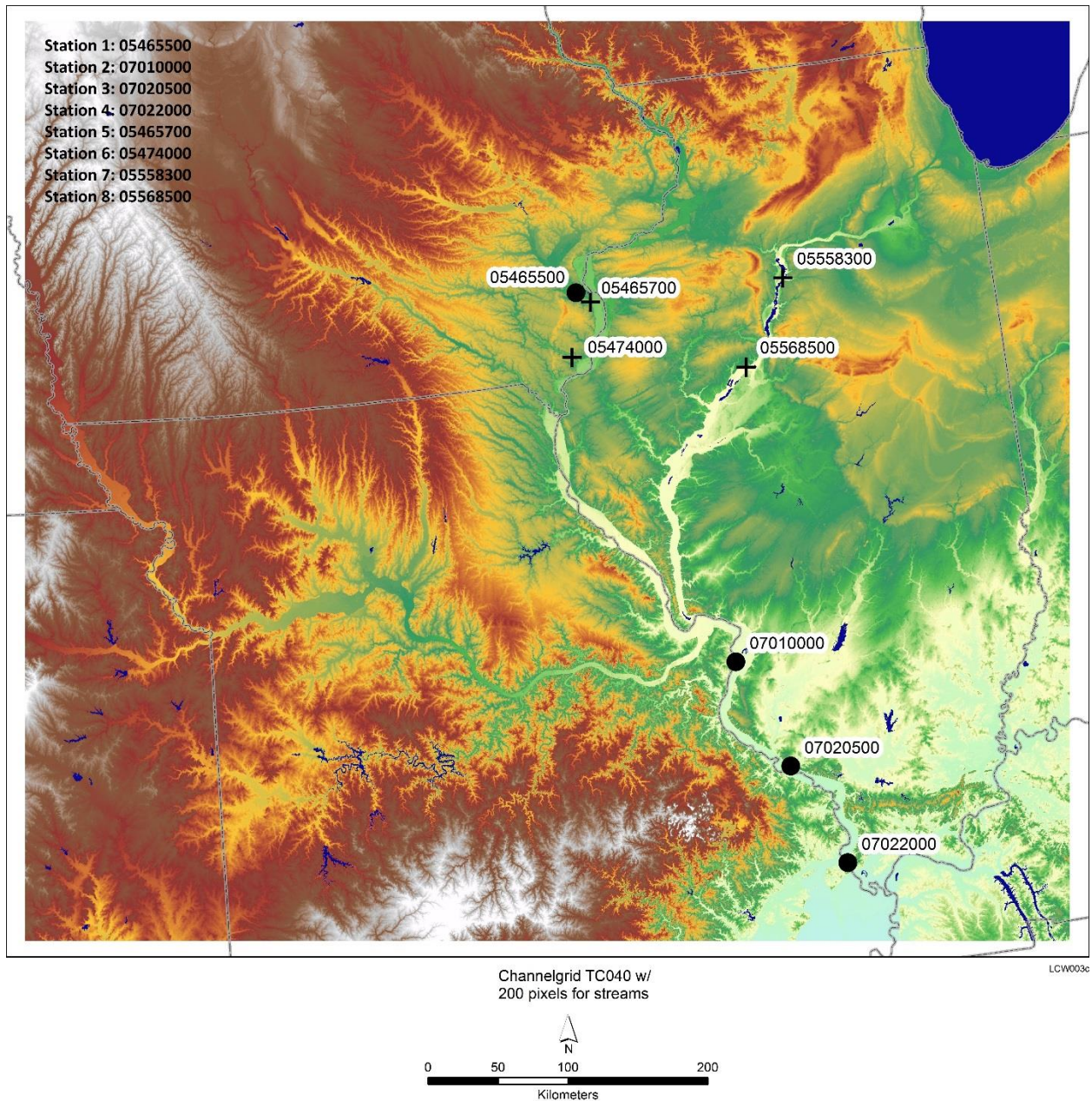
Statistics	Station 1	Station 2	Station 3	Station 4
Calibration				
NSE	0.87 (0.73)	0.64 (-54.4)	0.05 (-157.3)	-58.78 (-1316.9)
RMSE	48.1 (69.3)	318.2 (3925.2)	308.7 (3981.3)	934.6 (4391.3)
PCC	0.95 (0.91)	0.87 (0.92)	0.91 (0.87)	0.53 (0.66)
Validation				
NSE	0.83 (0.41)	-0.08 (-3.5)	-0.08 (-27.4)	-0.12 (-3.33)
RMSE	259.9 (487.3)	3264.3 (6670.1)	3170.1 (16305.7)	3283.9 (6854.3)
PCC	0.83 (0.8)	0.98 (0.69)	0.29 (0.19)	0.94 (0.64)

3 ^a The calibration period is 3 days (April 9–11) and includes 22 parameters. The validation period
 4 is April 12–24. Bold typeface indicates the calibrated model results are closer to observations
 5 compared with the default model results. NSE and PCC are unitless; RMSE is in m³/sec.

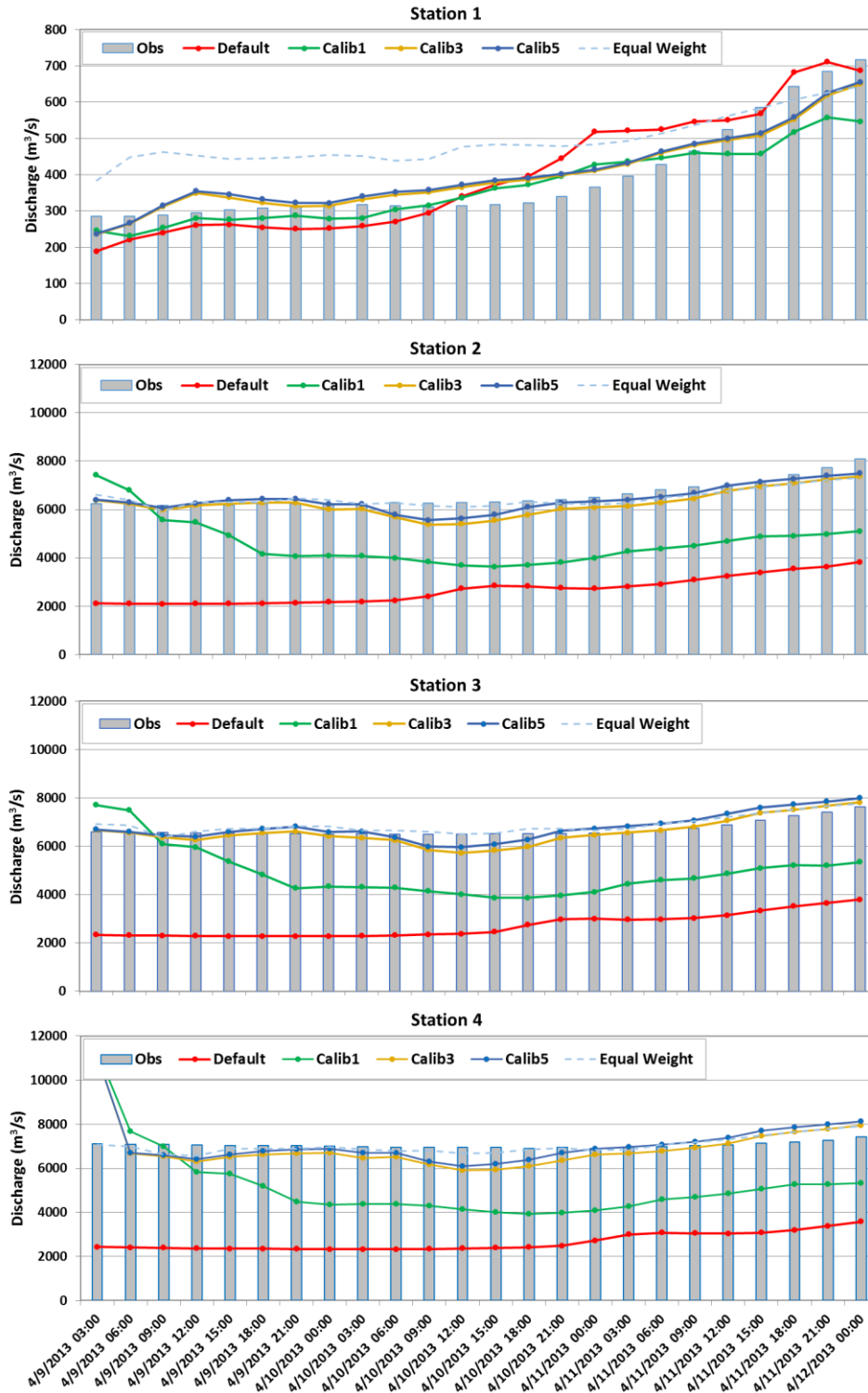
1 **Table 3. Experiments designed to test the scale-up capability and computational benefits of**
 2 **HPC-enabled parallel PEST linked to WRF-Hydro.**

Test	No. of Workers	No. of Lamdas	No. of Nodes for Each Worker	Total Time Cost (min)	Time Cost for Calculating Jacobian Matrix	Time Cost for Testing Parameter Upgrades
Test 1	23	15	2	103	52	51
Test 2	12	10	2	150	102	48
Test 3	6	5	2	264	211	53
Test 4	6	5	4	131	107	24
Test 5	6	5	6	86	70	16

3

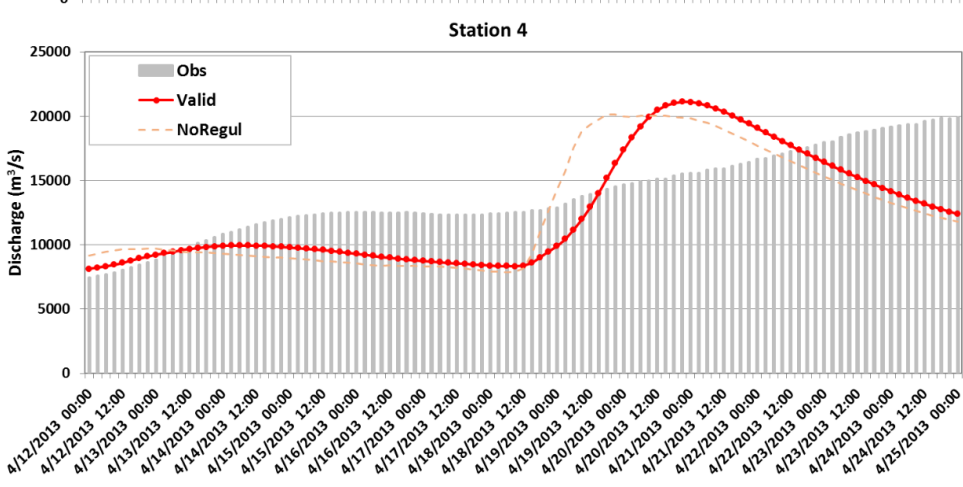
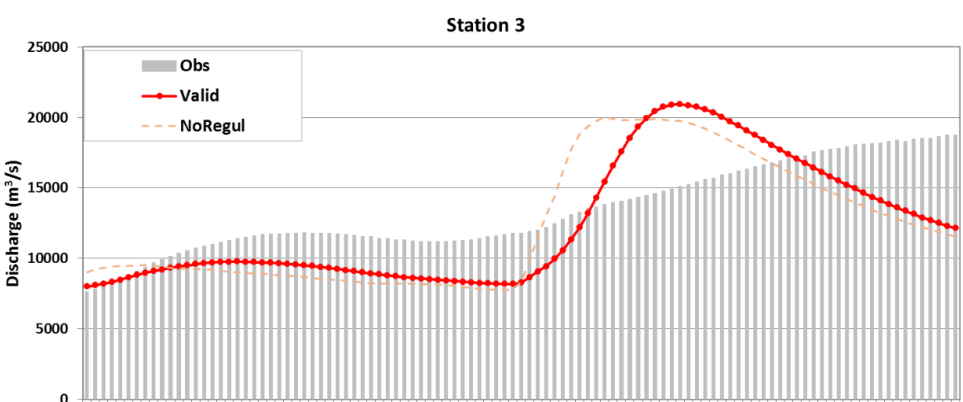
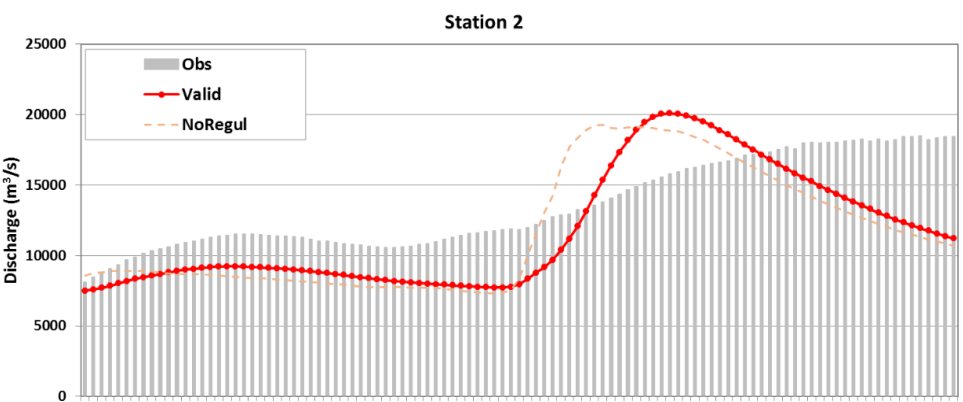
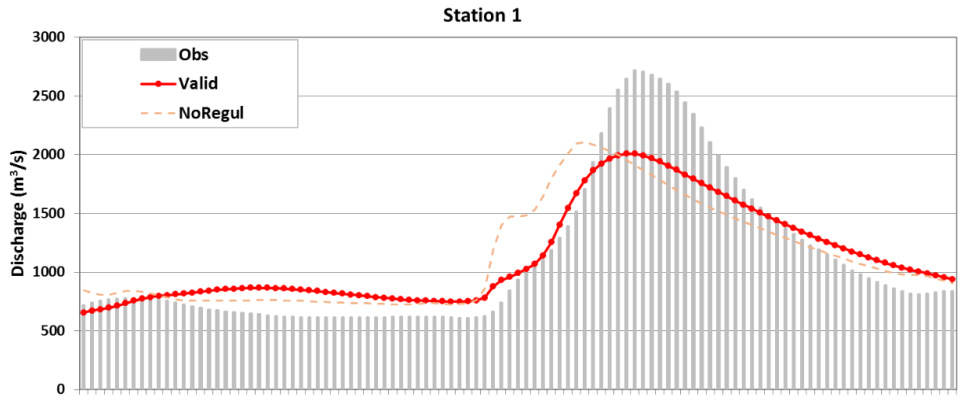


1
 2 **Figure 1: Eight USGS sites over the study area (750 km x 660 km). The four circles are sites**
 3 **that are used for calibrations; the four crosses are sites that are used for transferability**
 4 **assessment. USGS site numbers corresponding to the site index used in this study are listed**
 5 **on the top left corner of the map.**

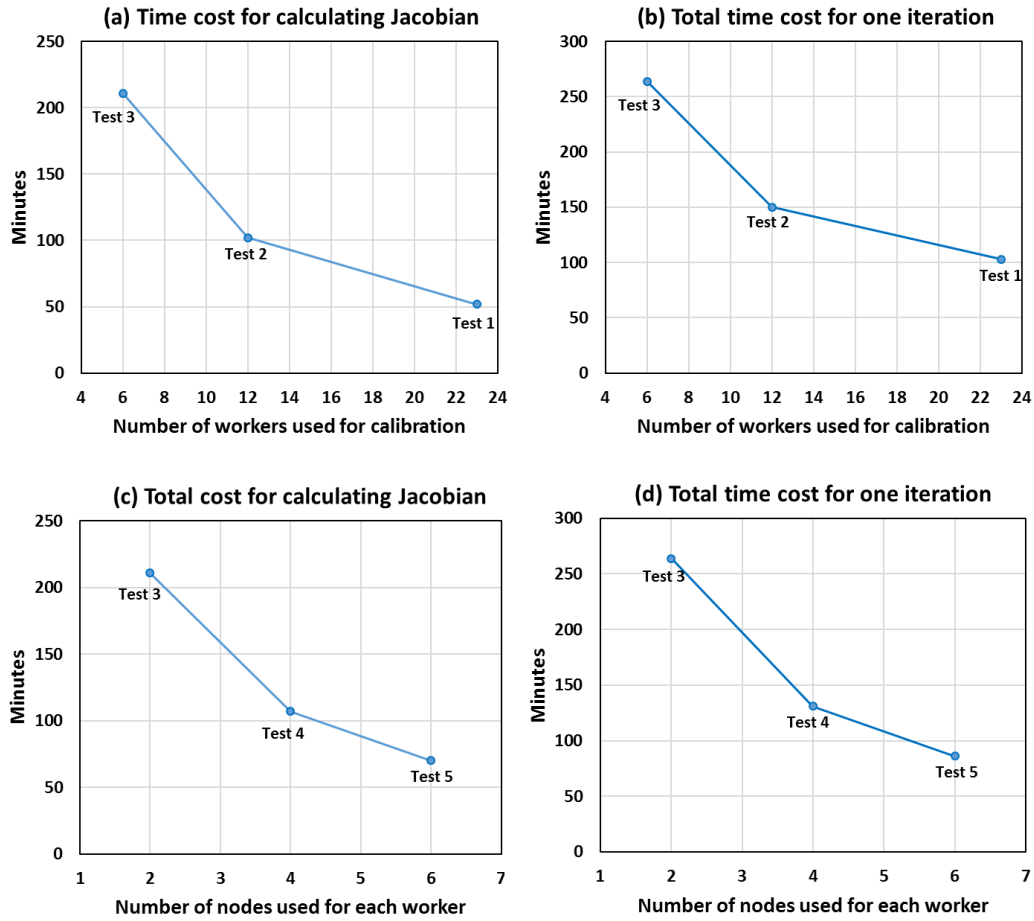


1
 2 **Figure 2: Observed and modeled discharge (m^3/sec) using default and calibrated parameters**
 3 **during a 3-day calibration period (April 9–11, 2013) over the four stations indicated by the**
 4 **black circles in Fig. 1. The calibrations shown in solid lines are conducted by using SVD-**

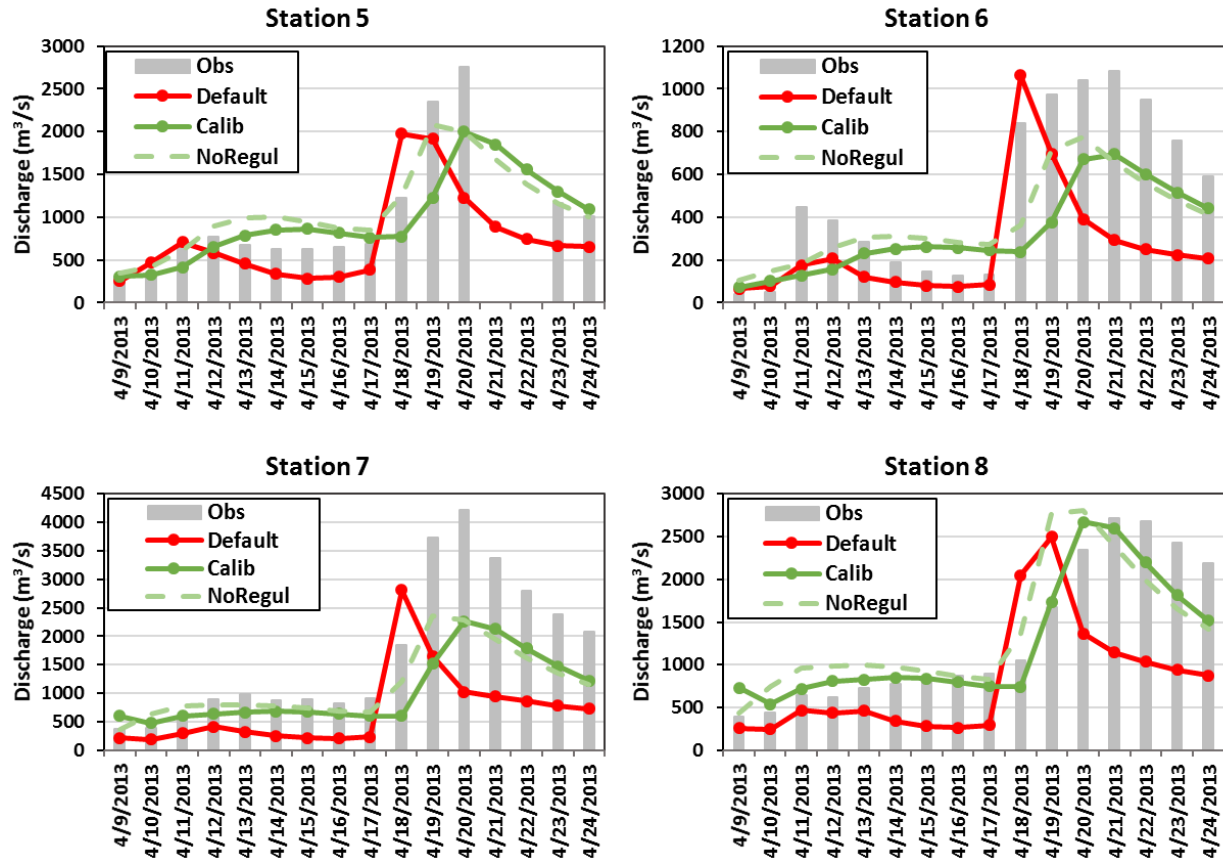
- 1 **based regularization and a higher weight for Station 1. The dashed line is the optimum result**
- 2 **calibrated by using equal weight for all four sites.**



1 **Figure 3: Observed and modeled discharge (m^3/sec) during a validation period (April 12–24,**
2 **2013) using optimum parameters identified from a 3-day calibration over the four stations**
3 **indicated by black circles in Fig. 1. The solid line uses the optimum parameters that**
4 **identified by PEST with SVD-based regularization and a higher weight for Station 1. The**
5 **dashed line uses the optimum result calibrated by using estimation mode (no regularization).**



1
 2 **Figure 4. Time cost for calculating Jacobian matrix and total time cost for one iteration for**
 3 **the five experiments (Table 3) using different number of workers to conduct PEST (a, b) and**
 4 **different number of nodes for each worker (c, d) to conduct WRF-Hydro.**



1
 2 **Figure 5: Observed and modeled daily averaged discharge (m^3/sec) from April 9–24 using**
 3 **default and the optimum parameters (shown in Table 1) identified by the 3-day calibration**
 4 **over four stations that are in the study area (indicated by crosses in Fig. 1). The calibrations**
 5 **shown in solid lines are conducted by using SVD-based regularization and a higher weight**
 6 **for Station 1. The dashed line is the optimum result calibrated by using estimation mode (no**
 7 **regularization).**

A parallel workflow implementation for PEST version 13.6 in high-performance computing for WRF-Hydro version 5.0: a case study over the midwestern United States

¹Jiali Wang, ¹Cheng Wang, ²Vishwas Rao, ¹Andrew Orr, ¹Rao Kotamarthi

¹Argonne National Laboratory, Environmental Science Division, 9700 South Cass Avenue, Lemont, IL 60439, USA

²Argonne National Laboratory, Mathematics and Computer Science Division, 9700 South Cass Avenue, Lemont, IL 60439, USA

Correspondence to: Jiali Wang (jialiawang@anl.gov)

Abstract. The Weather Research and Forecasting Hydrological (WRF-Hydro) system is a state-of-the-art numerical model that models the entire hydrological cycle based on physical principles. As with other hydrological models, WRF-Hydro parameterizes many physical processes. Hence, WRF-Hydro needs to be calibrated to optimize its output with respect to observations for the application region. When applied to a relatively large domain, both WRF-Hydro simulations and calibrations require intensive computing resources and are best performed on multimode, multicore high-performance computing (HPC) systems. Typically, each physics-based model requires a calibration process that works specifically with that model and is not transferrable to a different process or model. The parameter estimation tool (PEST) is a flexible and generic calibration tool that can be used in principle to calibrate any of these models. In its existing configuration, however, PEST is not designed to work on the current generation of massively parallel HPC clusters. To address this issue, we ported the parallel PEST to HPCs and adapted it to work with WRF-Hydro. The porting involved writing scripts to modify the workflow for different workload managers and job schedulers, as well as developing code to connect parallel PEST to WRF-Hydro. To test the operational feasibility and the potential computational benefits of this first-of-its-kind HPC-enabled parallel PEST, we developed a case study using a flood in the midwestern United States in 2013. Results on a problem involving calibration of 22 parameters show that on the same computing resource used for parallel WRF-Hydro, the HPC-enabled parallel PEST can speed the calibration process by a factor of up to 15 compared with commonly used

1 PEST in sequential mode. The speedup factor is expected to be greater with a larger calibration
2 problem (e.g., more parameters to be calibrated or a larger size of study area).

3 **1 Introduction**

4 Physically based hydrological models contain detailed physical mechanisms to model the
5 hydrological cycle, but many complex physical processes in these models are parameterized. For
6 example, the state-of-the-art Weather Research and Forecasting Hydrological (WRF-Hydro)
7 modeling system (Gochis et al., 2015) has dozens of parameters that can be land- and river-type
8 dependent and are typically specified in lookup tables. Therefore, these hydrological models need
9 to be calibrated before they can be applied to research over different regions. In this context,
10 calibration refers to adjusting the values of the model parameters so that the model can closely
11 match the behavior of the real system it represents. In some cases, the appropriate value for a
12 model parameter can be determined through direct measurements conducted on the real system. In
13 many situations, however, the model parameters are conceptual representations of abstract
14 watershed characteristics and must be determined through calibration. In fact, model calibration is
15 the most time-consuming step, not only for hydrological models, but also for Earth system model
16 development, because both parametric estimation and parametric uncertainty analysis require
17 hundreds—if not thousands—of model simulations to understand how perturbations in model
18 parameters affect simulations of dominant physical processes and to find the optimum value of a
19 single parameter.

20
21 WRF-Hydro is a numerical model that can simulate the entire hydrological cycle using advanced
22 high-resolution data such as satellite and radar products. Compared with the traditional land
23 surface model (LSM) used by WRF, WRF-Hydro provides a framework for multiscale
24 representation of surface flow, subsurface flow, channel routing, and baseflow, as well as a simple
25 lake/reservoir routing scheme. As a physics-based model, WRF-Hydro includes many complicated
26 physical processes that are nonlinear and must be parameterized. The default parameters given by
27 WRF-Hydro may be valid for one region but not for another region. Hence calibration of related
28 model parameters is often required in order to use the model in a new domain. In particular, for a
29 large spatial domain such as the entire contiguous United States, in order to develop the optimal
30 parameter sets in a reasonable amount of time, the calibration must be conducted on high-

1 performance computing (HPC) systems in parallel instead of in the traditional sequential mode.
2 To date, no such calibration tool can efficiently calibrate WRF-Hydro on HPC resources.
3 Typically, each physics-based model needs a calibration code that is custom designed to work with
4 that particular numerical model and its set of physics parameterizations, software architecture, and
5 solvers. These custom-designed calibration codes are highly challenging and do not offer
6 flexibility. Therefore, a more flexible and generic calibration tool is needed that can calibrate any
7 code that uses Message Passing Interface/Open Multi Processing (MPI/OpenMP) for
8 parallelization on HPC systems.

9

10 One widely used generic and independent calibration tool is the parameter estimation tool (PEST).
11 PEST (Doherty, 2016) conducts calibration automatically based on mathematical methods and
12 thus is applicable for optimizing nonlinear parameters. Compared with manual calibration,
13 automatic calibration is more efficient and effective because it avoids interference from human
14 factors (Madsen, 2000; Getirana, 2010). The uniqueness of PEST is that it operates independent
15 of models: there is no need to develop additional programs or codes for a particular model except
16 preparing the files required by PEST (as described in Sec. 3.2). PEST has four modes of operation.
17 One of the modes is regularization mode, which supports the use of Tikhonov regularization and
18 is found better for serving environmental models because, if implemented properly, it supports
19 model predictions of minimum error variance, is numerically stable, and embraces rather than
20 eschews the heterogeneity of natural systems. Singular value decomposition (SVD) can be used as
21 a regularization device to guarantee numerical stability of the calibration problem. Parallel PEST
22 is able to distribute many runs across many computing nodes using master-worker parallel
23 programming. To our best knowledge, however, no approach is available that allows users to submit
24 jobs using PEST parallelization to a typical supercomputing facility that uses job scheduling and
25 workload management such as Simple Linux Utility for Resource Management (SLURM),
26 Portable Batch System (PBS), and Cobalt. A previous study (Senatore et al., 2015) used PEST to
27 calibrate WRF-Hydro over the Crati River Basin in southern Italy. Because the study area was
28 relatively small, the authors were able to conduct the calibration using PEST in sequential mode
29 (Alfonso Senatore, personal communication, 2018).

30

1 This study aims to (1) port parallel PEST to HPC clusters operated by the U.S. Department of
2 Energy (DOE) and adapt it to work with WRF-Hydro, (2) evaluate the performance of HPC-
3 enabled parallel PEST linked to WRF-Hydro by calibrating a flood event, and (3) explore the
4 scale-up capability and computational benefits of HPC-enabled parallel PEST by assigning
5 different computing resource to the entire calibration process.

6 **2 Model description**

7 **2.1 Study area**

8 The case presented here is one of the worst floods experienced by greater Chicago area in the past
9 three decades; the storm occurred on April 18, 2013 (Campos and Wang, 2015). According to the
10 National Weather Service (NWS), the heaviest 24-hour accumulated rainfall during this storm
11 reached 201.4, 171.1, and 136.4 mm across Illinois, Iowa, and Missouri, respectively. The
12 Mississippi River crested at 10.8 m (1.7 m above flood stage), and the Illinois River crested in
13 Peoria, Illinois, at 8.95 m; these river cresting broke the previous record of 8.78 m, set in 1943,
14 and was 4.55 m above the historical normal river stage (NWS, 2013). Campos and Wang (2015)
15 conducted three-domain nested WRF simulations to understand the dynamical and microphysical
16 mechanisms of the event. Our study builds on the smallest domain of that study, which covers the
17 majority of Illinois, Iowa, and Missouri at a spatial resolution of 3 km (Fig. 1). The domain size is
18 750 km from west to east and 660 km from south to north.

19

20 **2.2 WRF-Hydro configuration**

21 This study employs WRF-Hydro version 5 with a basic configuration. This configuration does not
22 use nudging techniques or spatially distributed soil-related parameters as used in the National
23 Water Model configuration. WRF-Hydro has been tested in several different cases that focused on
24 different hydrometeorological forecasting and simulation problems (e.g., Gochis et al., 2018;
25 Yucel et al., 2015; Senatore et al., 2015; Arnault et al., 2016), and it shows reasonable accuracy in
26 simulated streamflow after being carefully calibrated. For details of the WRF-Hydro modeling
27 system, see Gochis et al. (2018). Currently, two LSMs are available in WRF-Hydro for
28 representing land-surface column physics: Noah (Chen and Dudhia, 2001) and Noah Multi-

1 parameterization (Noah-MP; Niu et al. 2011). We utilize Noah-MP LSM because compared with
2 Noah LSM it shows obvious improvements in reproducing surface fluxes, skin temperature over
3 dry periods, snow water equivalent, snow depth, and runoff (Niu et al. 2011). The Noah-MP is
4 configured at a grid spacing of 3 km, and the aggregation factor is 15; that is, starting from a 3 km
5 LSM resolution in the domain shown in Fig. 1, hydrological routing is performed at a grid
6 resolution of 200 m, with 3285 south-north \times 3735 west-east grid cells. We use a time step of 10
7 seconds for the routing grid in order to maintain model stability and prevent numerical dispersion
8 of overland flood waves. The time step also meets the Courant condition criteria for diffusive wave
9 routing on a 200 m resolution grid. The WRF-Hydro is configured to be in offline or uncoupled
10 mode—there is no online interaction between the WRF-Hydro hydrological model and the WRF
11 atmospheric model. Overland flow, saturated subsurface flow, gridded channel routing, and a
12 conceptual baseflow are active in this study. The gridded channel network uses an explicit, one-
13 dimensional, variable time-stepping diffusive wave. A direct output-equals-input “pass-through”
14 relationship is adopted to estimate the baseflow. Although the baseflow module is not physically
15 explicit, it is important because the water flow in the channel routing is contributed by both the
16 overland flow and baseflow. If the overland flow is active as it is in this study, it passes water
17 directly to the channel model. In this case the soil drainage is the only water resource flowing into
18 the baseflow buckets. However, if the overland flow is deactivated but channel routing is still
19 active, then WRF-Hydro collects excess surface infiltration water from the land model and passes
20 this water into the baseflow bucket. This bucket then contributes the water from both overland and
21 soil drainage to the channel flow. Therefore, the baseflow must be active if the overland flow is
22 switched off. This study does not consider lakes and reservoirs.

23

24 We use the geographic information system (GIS) tool (Sampson and Gochis, 2018) developed by
25 the WRF-Hydro team to delineate the stream channel network, open water (i.e., lake, reservoir,
26 and ocean) grid cells, and groundwater/baseflow basins. Meteorological input for the WRF-Hydro
27 model system includes hourly precipitation; near-surface air temperature, humidity, and wind
28 speed; incoming shortwave and longwave radiation; and surface pressure. In this study, the hourly
29 precipitation is from the National Centers for Environmental Prediction (NCEP) Stage IV analysis
30 at a spatial resolution of 4 km. The Stage IV data is based on combined radar and gauge data (Lin
31 and Mitchell, 2005; Prat and Nelson, 2015), and has been shown to be temporally well correlated

1 with high-quality measurements from individual gauges (see, e.g., Sapiano and Arkin, 2009; Prat
2 and Nelson, 2015). The other hourly meteorological inputs are from the second phase of the multi-
3 institution North American Land Data Assimilation System project, phase 2 (NLDAS-2) (Xia et
4 al., 2012a,b), at a spatial resolution of 12 km. NLDAS-2 is an offline data assimilation system
5 featuring uncoupled LSMs driven by observation-based atmospheric forcing.

6
7 During the 15-day period of this studied case, light to moderate rain occurred on April 8 through
8 11, 2013, followed by a relatively dry period from April 12 to 15. Then a heavy rain event began
9 on April 16 and peaked on April 18. The heaviest rain band moved east of the study area on April
10 19. The rainy event ended over the study area on April 20 (see Fig. S1 in Supporting Information).
11 We start the WRF-Hydro simulation on Jan. 1, 2013, and run the model for more than three months
12 to reach equilibrium. This 3-month period is considered as spin-up time and is excluded from
13 model calibration and evaluation. We calibrate the river discharge calculated by the WRF-Hydro
14 model from 00UTC April 9 to 00UTC April 12, 2013, considering it long enough to achieve our
15 objective. We then evaluate the model performance against U.S. Geological Survey (USGS)
16 observed river discharge from 00UTC April 12 to 00UTC April 25, 2013.

17 **3 Calibration**

18 **3.1 Platforms**

19 We customized parallel PEST to work on three different workload managers and job schedulers:
20 SLURM at the National Energy Research Scientific Computing Center (NERSC), PBS at the
21 Argonne National Laboratory Computing Resource Center, and Cobalt at the Argonne Leadership
22 Computing Facility. The tests presented here are conducted on Edison at NERSC, which uses the
23 SLURM workload manager and job scheduler. Edison is a Cray XC30 with a peak performance
24 of 2.57 petaflops per second, 133,824 compute cores, 357 terabytes of memory, and 7.56 petabytes
25 of disk storage. It has 5,586 nodes and 24 cores per node.

26
27 The interface we have built between parallel PEST and the management software (SLURM here)
28 is, in general, used for (1) setting the number of workers and the nodes for each worker to conduct
29 a model run (WRF-Hydro here); (2) finding the nodes that are available; (3) setting up the working

1 directory for the workers; (4) identifying the nodes that work for each worker; (5) passing the
2 global files (same for all the working directory) to all the workers (these files include the lookup
3 table files that are not to be calibrated, the namelist files for both LSM and hydrological sector,
4 and restart files that generated by the previous simulations, or spin-up period); and (6) submitting
5 the job for the entire calibration process, including parallel PEST and parallel WRF-hydro. This
6 job can be submitted as a cold-start run or as a restart. The main difference for this interface on
7 different management software is that different management software has its own way to submit
8 jobs and identify available nodes. This difference requires some changes in the script we
9 developed.

10 **3.2 PEST files and settings**

11 PEST requires three file types in both sequential and parallel mode. They are template files to
12 define the parameters to be calibrated, an instruction file to define the format of model-generated
13 output files, and a control file to supply PEST with the size of the problem and the settings for the
14 calibration method. Parallel PEST uses a “master-worker” paradigm that starts model runs
15 simultaneously by different workers (or in different folders). The master of parallel PEST
16 communicates with each of its workers many times during a calibration. To run PEST in parallel
17 mode, one also needs a management file to inform PEST where the working folder is for each
18 worker and what the names and paths are for each model input file that PEST must write (i.e.,
19 lookup tables that come from template files) and each model output file that PEST must read (such
20 as frsxt_pts_out.txt). The management file also set the maximum running time for each worker.
21 For those workers that take longer than the maximum running time, PEST will stop the model run
22 by that particular worker and assign that model run to another worker if there is one with nothing
23 else to do.

24
25 To the best of our knowledge, however, parallel PEST is not designed to run on HPCs directly.
26 We developed scripts and an interface to enable parallel PEST to run on HPCs using SLURM,
27 PBS, or Cobalt workload managers and job schedulers. The development involved writing scripts
28 to modify the workflow for different workload managers and job schedulers, as well as developing
29 code to connect parallel PEST to WRF-Hydro. These developments enable parallel PEST to run
30 many workers at the same time; each worker runs a parallel code (here WRF-Hydro) that uses

1 more than one node, which could significantly reduce the wall-clock time of model calibrations.
2 Although this master-worker parallelism may not be as efficient as a fully MPI approach, it is
3 sufficient for model calibration and requires the least effort for the current parallel PEST to run on
4 HPC systems.

5
6 This study presents calibration results from PEST using the SVD-based regularization in
7 regularization mode to ensure numerical stability (Tonkin and Doherty, 2005). We focus on
8 calibrating 22 parameters (see Table 1 and detail description in Sec. 3.3) using 96 observation
9 points and 22 items of prior information for the calibrated parameters. In each item of prior
10 information, a value equal to its default value provided by the WRF-Hydro v5.0 (or the log of its
11 default value) is assigned for each adjustable parameter, assuming that default values are the
12 preferred values. All prior information equations are assigned a weight of 1.0. We assigned five
13 different regularization groups to the prior information: Manning’s roughness coefficients
14 specified by Strahler stream order in CHANPARAM.TBL to one group; the parameters in
15 HYDRO.TBL (Manning’s roughness coefficients for overland flow as a function of vegetation
16 types) to another group; and three global parameters for the Noah-MP (xslop1, refdk, and refkdt)
17 in GENPARAM.TBL to the remaining three groups. The 96 observation points are given different
18 weights based on the inversed mean of their observed discharge during the studied period (see the
19 detailed description in Sec. 3.3 and Sec. 4.1). For a detailed description of these settings see the
20 PEST User Manual (Doherty, 2016).

21

22 **3.3 Calibrated experiments**

23 The primary objective of this study is to build a bridge for linking the parallel PEST and WRF-
24 hydro on the basis of HPC clusters and to explore the computational benefits of this bridge. We
25 do not attempt to extensively assess each individual tool or address questions in each individual
26 domain, such as optimizing the objective functions in PEST or calibrating WRF-Hydro for a long
27 time period considering all the relevant parameters to achieve an optimal parameter set. The
28 calibration period thus is limited to only three days, which we believe long enough to achieve our
29 objective and to understand WRF-Hydro’s sensitivity to the calibrated parameters. We calibrated
30 WRF-Hydro using four USGS sites (referred to as Station 1, Station 2, Station 3, and Station 4

1 hereafter), as shown in Fig. 1. (More USGS sites could be included if one manually reallocated
2 the stations that were not properly assigned to the desired location on the channel network by the
3 GIS tool.) We then transfer the calibrated parameters to other subbasins in the study area to assess
4 the transferability of the calibrated parameters. Although many parameters, including spatially
5 distributed parameters and constant parameters in the lookup tables, affect the model performance,
6 we calibrate only the parameters in lookup tables and do not consider the spatial variability of
7 other parameters or their scaling factors. We acknowledge that some studies calibrate a single
8 scaling factor (without considering its spatial variability, however) of overland roughness
9 coefficients (OVROUGHRTFAC) rather than the actual value of each land type in the lookup table
10 (e.g., Kerandi et al., 2018). Although this approach reduces the number of calibrated parameters,
11 however, it has less flexibility because changing one factor will change all the parameters that use
12 the same proportion.

13

14 For the calibration exercises we conduct here, the retention depth factor (RETDEPRTFAC) is
15 fixed at 0.001. This value is reasonable because the modeled discharge of our particular
16 configuration (Sec. 2.2) using default parameters is lower than observed discharge. Reducing this
17 factor from 1 to 0.001 keeps less water in water ponds and more water on the surface so it can
18 contribute to river discharge. First, we calibrate 48 parameters based on a 3-day simulation from
19 April 9 to April 11, 2013 (Table S1 in Supporting Information). This calibration uses the
20 estimation mode in the PEST tool and considers equal weight for all four USGS stations. We
21 calibrate Manning's roughness coefficients for both channels and land-use types, the deep drainage
22 (SLOPE), infiltration-scaling parameter (REFKDT), and saturated soil lateral conductivity
23 (REFDK). Manning's roughness coefficients control the hydrograph shape and the timing of the
24 peaks; the SLOPE, REFKDT, and REFDK control the total water volume. Second, based on the
25 knowledge we learn from the 48-parameter calibration (see details in Sec. 4.1), for the same 3-day
26 period, we reduce the number of calibrated parameters from 48 to 22 according to the sensitiveness
27 of the WRF-Hydro model to the adjustable parameters. For example, during the calibration we
28 find that Manning's roughness coefficients for several land types barely change because these land
29 types (e.g., tundra, snow/ice) are not present in the study area. We also learn that even though the
30 calibrated WRF-Hydro parameters can generate discharge results that closely resemble
31 observations, the physical meaning of several parameters are not appropriate because of the wide

1 range of those parameters that we set in the PEST control file. For example, Manning’s roughness
 2 coefficient for stream order 1 (0.199) is calibrated smaller than that for stream order 2 (0.218); the
 3 overland roughness coefficients for evergreen needleleaf forest (0.043) and mixed forest (0.023)
 4 are calibrated smaller than for cropland/woodland (0.046). Neither of these is true in the real world.
 5 We therefore adjust the range of many parameters according to the literature (Soong et al., 2012)
 6 to maintain their physical meanings (Table 1). We find that by using the same absolute weight for
 7 all four stations, the calibration helps three stations (Station 2, 3, and 4) with large water volumes
 8 to generate more reasonable results than do the default parameters; however, the results for Station
 9 1, which has a relatively small volume of water, is not always better than the discharge that is
 10 modeled by using default parameters. Thus, we assign a weight of 15.0 for Station 1 versus a
 11 weight of 1.0 for the other three stations according to the inversed mean of observed discharge
 12 over these four stations in April 2013. The ratio of the weights between Station 1 and the other
 13 three stations stays similar even if the means are calculated based on different time periods.

14

15 **3.4 Statistics**

16 This study employs three statistical criteria: Nash–Sutcliffe efficiency (NSE; Nash and Sutcliffe,
 17 1970; Moriasi et al., 2007), root-mean-square error (RMSE), and Pearson correlation coefficient
 18 (PCC). RMSE and PCC evaluate model performance in terms of bias and temporal variation. NSE
 19 quantitatively describes the accuracy of modeled discharge compared with the mean of the
 20 observed data. Equation (1) calculates the NSE with defined variables:

$$21 \quad NSE = 1 - \frac{\sum_{t=0}^n (Y_t^{obs} - Y_t^{sim})^2}{\sum_{t=0}^n (Y_t^{obs} - Y_{mean}^{obs})^2}, \quad (1)$$

22 where Y_t^{obs} is the t th observed value from USGS sites for river discharge, Y_t^{sim} is the t th
 23 simulated value from the WRF-Hydro output, Y_{mean}^{obs} is the temporal average of USGS observed
 24 discharge, and n is the total number of observation time points. An efficiency of 1 (NSE = 1)
 25 corresponds to a perfect match between modeled discharge and observed data. An efficiency of 0
 26 (NSE = 0) indicates that the model predictions are as accurate as the mean of the observed data.
 27 An efficiency below zero (NSE < 0) occurs when the model is worse than the observed mean.
 28 Essentially, the closer the NSE is to 1, the more accurate the model is.

1 **4 Results**

2 **4.1 WRF-Hydro calibration and validation**

3 Based on the knowledge we gained from the 48-parameter 3-day calibration, we adjust the range
4 of critical parameters in the PEST control file to main their physical meanings. For example, we
5 set Manning’s roughness coefficient larger for stream order 1 than for stream order 2. We also
6 adjust the parameter range of the overland roughness coefficient for multiple land covers, such as
7 forests. We exclude the parameters that WRF-Hydro is not sensitive to for this study, in order to
8 constrain the problem size considering the availability of computational resources. However, if the
9 studied area is much larger with more land types than the study area here, then there would be
10 more parameters to calibrate. Also, hundreds of constant parameters in the Noah-MP model could
11 affect the WRF-Hydro results (Cuntz et al. 2016) and can be calibrated. Both these situations
12 would increase the burden of WRF-Hydro calibration. We perform the same 3-day calibration
13 from April 9 to April 11, 2013. Figure 2 shows the results of the 3-day modeled discharge (in cubic
14 meters) using default and calibrated parameters after five iterations, as well as observed discharge.
15 The four stations are calibrated by considering different weights. Compared with the results
16 calibrated by using equal weights for all the stations, by giving a higher weight to Station 1 the
17 model bias over Station 1 is significantly reduced, with a higher NSE (0.87 with higher weight
18 versus 0.14 with equal weight) and lower RMSE (48.1 versus 123.6). Over Stations 2, 3, and 4,
19 which sit on rivers with relatively large water volumes, the modeled discharge using the default
20 parameter underestimates the streamflow by more than 65%. PEST detects this underestimation
21 and immediately adjusts the parameters and increases the modeled discharge during the first
22 iteration. After the third iteration, the difference in calibrated results between different iterations
23 is relatively small. We allow the PEST to conduct five iterations and use the parameters obtained
24 from the fifth iteration as our optimum parameters. As shown in Table 2, when the optimum
25 parameters are used, the modeled discharges are much closer to the observations compared with
26 the modeled results when the default parameters were used. The NSEs for the four stations
27 increased from 0.73 (Station 1), -54.4 (Station 2), 157.3 (Station 3) and -1316.9 (Station 4) to 0.87,
28 0.64, 0.05, and -58.78, respectively, being closer to 1. The RMSEs decreased from 69.3, 3925.2,
29 3981.3, and 4391.3 m³/sec to 48.1, 318.2, 308.7, and 934.6 m³/sec, respectively. Giving a lower
30 weight for the three large river stations does not change the calibration results much.

1
2 During the validation period, compared with the modeled discharge using default parameters, as
3 shown in Table 2, the NSEs for all four stations are increased to be closer to 1; RMSEs are
4 decreased by 50% or more; and the correlation coefficients between the observed and modeled
5 discharge are increased from 0.8, 0.76, 0.21, and 0.72 to 0.98, 0.82, 0.80, and 0.75. Compared with
6 the results of calibration using the estimation mode (no regularization) in PEST, the SVD-based
7 regularization generates slightly better hydrograph shape with 24-hour later discharge peaks that
8 are closer to the observations. However, a problem remains with the hydrograph shapes of the
9 modeled discharge, especially with the modeled peak of discharge. For Station 1, the WRF-Hydro
10 almost captures the timing of the peak of discharge, although it still underestimates the water
11 volume by ~25%. The reason is that this study uses a direct pass-through baseflow module, which
12 does not account for slow discharge and long-term storage of the baseflow. Therefore, the largest
13 contribution to river discharge is from precipitation, and groundwater does not contribute much
14 discharge to the channels in a long-term view, as is also true for the other three large river stations.
15 Different from Station 1, for the other three large river stations, the WRF-Hydro modeled
16 discharge increases soon after the peak of precipitation and reaches a peak on April 21, 2013,
17 which is much earlier than the observed peak of river discharge (near April 24). The reason is that
18 the water contributions for these stations are from a larger river basin (Mississippi River) than we
19 included in our current study area. Thus, when a heavy precipitation event occurs over the entire
20 river basin, there will be a significant lag time (especially at the lower part of the basin) between
21 the peak of precipitation amount and the peak of river discharge. For example, the precipitation
22 over the upper part of Mississippi River Basin (MRB) has a peak amount on April 18–19, but the
23 river discharge did not reach its peak until April 24. Because our studied area covers only half of
24 the MRB, the modeled river discharge has a shorter delay period after the peak of precipitation
25 than does the observed river discharge. Enlarging the study area to include the entire MRB may
26 improve this situation. Alternatively, calibrating and validating local rivers that are included in the
27 current study area may also reduce the bias in hydrograph shape compared to calibrating and
28 validating large rivers. On the other hand, the WRF-Hydro simulated river discharge decreases
29 soon after it reaches the peak and much earlier than the observed discharge. The reason is again
30 that the direct pass-through baseflow employed by this study does not account for slow discharge
31 and long-term storage of the baseflow. As a result, the contribution from the baseflow to the river

1 discharge in model simulations does not stay as long as in real situations. In the observations, the
2 river discharge decreases from the peak at a speed of $\sim 500 \text{ m}^3/\text{sec}$ per day, while the modeled river
3 discharge decreases from the peak at a speed of $\sim 1667 \text{ m}^3/\text{sec}$ per day. Using exponential storage-
4 discharge function for the baseflow may improve this situation.

5

6 **4.2 Computational benefits of parallel PEST on HPCs**

7 The ability to scale up the calibration of WRF-Hydro by using parallel PEST on HPC systems is
8 determined by two factors: the scale-up capability of parallel PEST and the scale-up capability of
9 WRF-Hydro. In calibrating WRF-Hydro, PEST first makes as many model runs as there are
10 adjustable parameters to calculate Jacobian matrix (Doherty, 2016). The Jacobian matrix has a
11 column for each calibrated parameter and a row for each observation and each item of prior
12 information that set in the PEST control file. These model runs are independent between workers
13 and can be easily parallelized. Each worker runs the model with temporarily incremented
14 parameters that are defined in the template and control files. Then, PEST needs to make additional
15 model runs to test parameter updates. Different from the Jacobian runs, these additional runs are
16 performed by using different Marquardt lambdas, and the search for a Marquardt lambda that
17 achieves the best set of parameters is a serial iterative process. The lambda to use for the next run
18 depends on the outcome of the model run conducted using the previously chosen lambda. Although
19 serial testing of Marquardt lambdas may quickly find the optimal Marquardt lambda in the first or
20 second series of model runs, it is an inefficient use of computing resources because other
21 processors are idle while only one process is searching the lambdas. This is especially true when
22 the model domain is large and requires extensive computing resources. This study employs “partial
23 parallelization” for the lambda-testing procedure (Doherty, 2016), so multiple workers can be used
24 to calculate parameter upgrades based on a series of lambda values that are related to each other
25 by a factor of RLAMFAC set in the PEST control file. We also set the value of PARLAM to -9999
26 in the management file so only one cycle of parallel WRF-hydro runs is devoted to testing
27 Marquardt lambdas. For additional details on these parameters and their settings see the PEST
28 User Manual (Doherty, 2016).

29

1 In this study we test the computational performance of HPC-enabled parallel PEST using different
2 number of workers (6, 12, and 23) for the 22-parameter calibration. As shown in Table 3, we
3 conducted five experiments: Test 1 uses 23 workers, Test 2 uses 12 workers, and Test 3 uses 6
4 workers. All three tests use two nodes for each worker to run WRF-Hydro in parallel. The
5 maximum number of lambda-testing runs undertaken per iteration is set to 15, 10, and 5 for Test
6 1, 2, and 3, respectively, to make sure that only one cycle of WRF-hydro runs is devoted (using
7 15, 10 and 5 workers from Tests 1, 2, and 3, respectively) to testing Marquardt lambdas. Note that
8 the maximum number of lambda-testing runs should be set equal to or less than the workers
9 available. Otherwise, another cycle of WRF-hydro runs needs to be conducted. In fact, generating
10 more Marquardt lambdas does not always guarantee that the best Marquardt lambdas are
11 generated. In contrast, it may make the model convergence slower (here, PEST) or even model
12 failure.

13

14 In order to test the trade-offs between the computing nodes used for running parallel WRF-Hydro
15 and the workers used for running parallel PEST, Tests 4 and 5 use different number of nodes for
16 each worker to run WRF-Hydro in parallel. Explicitly, Test 4 uses four nodes per worker, and Test
17 5 uses six nodes per worker. Both tests use six workers for running the parallel PEST. The
18 maximum number of lambda-testing runs undertaken per iteration is set to five for both Tests 4
19 and 5. Note that the time costs in Table 3 are limited to only one iteration. Conducting more
20 iterations will increase the cost of wall-clock time and computing, but will not change the
21 conclusion for the scale-up capability and computational benefits for HPC-enabled parallel PEST
22 linked to WRF-hydro.

23

24 PEST needs to run the WRF-Hydro model at least as many times as the number of calibrated
25 parameters (22 here). In fact, PEST runs the model 23 times in the first round (or the first iteration)
26 with initial parameter values and for the first Jacobian matrix. From the second iteration, it runs
27 the model 22 times to calculate Jacobian matrix. Therefore, if there are fewer than 23 workers, the
28 time cost for the first round of Jacobian matrix calculation will increase accordingly. For example,
29 as shown in Fig. 4a, when we assign 12 (and 6) workers to parallel PEST, the time cost for
30 calculating the Jacobian matrix is increased by a factor of 2 (and 4) compared with the time cost
31 of using 23 workers. The time cost for the parameter upgrade stays similar for the three

1 experiments because only one cycle of WRF-hydro simulation is conducted to test the Marquardt
2 lambdas. As a result, the total time cost for Test 2 is ~1.5 times more than that for Test 1, and the
3 total time cost for Test 3 is ~1.5 times more than that for Test 2 (Fig. 4b). By extrapolating the
4 speedup curve shown in Fig. 4a and Fig. 4b, we expect the total time cost to be ~1516 minutes
5 when using only one worker (or sequential mode), which is about 15 times slower compared with
6 running the PEST in parallel mode using 23 workers. For this particular study with 22 adjustable
7 parameters, we expect the time cost most likely to stay the same even if one increases the number
8 of workers to more than 23, because PEST runs WRF-Hydro only 23 or 22 times for each iteration.
9 Assigning more workers for this particular study would most likely render some workers idle and
10 is not an efficient use of computing resources. PEST may run WRF-Hydro more than 22 times
11 (e.g., 44 times) if higher-order finite differences are employed. In this case, assigning more
12 workers (e.g. 45 workers) may further speed up the calibration process. On the other hand, for the
13 same case study and using the same number of nodes for running parallel WRF-Hydro, we can
14 estimate the computing speedup by assuming an increase in the number of calibrated parameters
15 to 50. This would be the case, for example, to evaluate model sensitiveness to the physics in Noah-
16 MP or the spatial variabilities of certain parameters. We then expect to use 51 workers to achieve
17 the best computing performance for parallel PEST. This would then be 28–30 times faster than
18 running PEST using one worker (or in sequential mode). Similarly, if 100 parameters were used
19 for the calibration for the same case study, a factor of up to 60 speedup in the calibration process
20 would be achieved by running HPC-enabled parallel PEST.

21
22 In addition, by increasing the number of nodes for each worker to conduct WRF-Hydro (Tests 3,
23 4, and 5), the time cost for the entire calibration process is significantly reduced (Figs. 4c and 4d).
24 Specifically, the WRF-hydro scales up well when using four and six nodes compared with using
25 two nodes per worker for running the WRF-Hydro. Both the time spent on calculating the Jacobian
26 matrix and the time spent on testing the parameter upgrades are decreased by 49% and 67%,
27 respectively, when using four and six nodes. Therefore, the total time spent is also decreased when
28 using more nodes for each worker (see Table 3). Increasing the number of nodes to eight for each
29 worker will most likely further decrease the time cost by 70–75% compared with using only two
30 nodes per worker. Moreover, if one has a larger study area such as the entire contiguous United
31 States, we expect the WRF-Hydro to have an even better scale-up capability (e.g., on dozens of

1 nodes) than this study. Overall, based on the experiments we conduct here, using 23 workers for
2 parallel PEST and six nodes for each worker to run parallel WRF-Hydro would cost the least wall-
3 clock time—about 32 min for one iteration for this particular study.
4

5 **4.3 Evaluation of spatial transferability of the calibrated parameters**

6 To assess the transferability of the calibrated parameters, we apply the optimum parameters
7 obtained from the calibration for the four stations (black circles) in Fig. 1 to another set of four
8 stations (crosses in Fig. 1) in the study area. All four sites are located on relatively small rivers, so
9 the lag time between precipitation peak and the discharge peak are much shorter than that for the
10 stations on the lower part of MRB (e.g., Stations 2, 3, and 4). The assessment compares the
11 observed discharge with the closest grid cells from the discharge output of WRF-Hydro. Figure 5
12 shows the observed and modeled discharge using default and the optimum parameters. Overall,
13 WRF-Hydro’s default parameters underestimate the discharge and misrepresent the timing of
14 discharge peaks compared with observations over the four assessed stations (Stations 5, 6, 7, and
15 8). By using the calibrated parameters from other sites over the area, the model results increase the
16 discharge and shift the hydrograph shape so they are much closer to the observations than model
17 results using default parameters. The absolute error of simulated discharge decreases by 13.1%,
18 38.3%, and 71.6%, respectively, over Stations 6 through 8 (Station 5 shows a 6% increase of
19 absolute error), compared with the default simulated discharge. We also find that using the SVD-
20 based regularization for the PEST calibration captures the timing of discharge peak better than
21 using the estimation mode, which is one-day earlier than the observations reaching the discharge
22 peak.

23 **5 Summary and discussion**

24 WRF-Hydro is a new, and perhaps the first practical, computer code that can run on HPC systems
25 and can model the entire hydrological cycle using physics-based submodels and high-resolution
26 input datasets (e.g., radar). The hydrological community has desired this capability for decades,
27 although it requires intensive computing resources. Thus, the calibration of this model would
28 ideally be conducted on HPCs in parallel as well, especially when the model covers a large domain
29 rather than the basin scale. This study ports an independent model calibration tool, parallel PEST,

1 to HPC clusters and links it to WRF-Hydro to help WRF-Hydro users calibrate the model within
2 a much shorter wall-clock time period. The bridge we build here (between parallel PEST and
3 WRF-Hydro on the basis of HPC systems) can be applied to any other hydrological models and
4 Earth system models that use parameterizations to represent model physics. We present the
5 operational feasibility of the HPC-enabled parallel PEST by evaluating the performance of
6 calibrated WRF-Hydro against observation in hydrograph features such as volume and timing of
7 flood events. We examine the scale-up capability and computational benefits of the tool by
8 assigning different computing resource for PEST and for WRF-Hydro. While this study presents
9 the optimum parameters identified from the calibration of the particular flood event, the parameters
10 can be significantly different if one uses different physics, such as exponential storage-discharge
11 function for a groundwater model or reach-based channel routing. Our preliminary testing shows
12 that using exponential storage-discharge function with the default parameters provided by WRF-
13 Hydro, the modeled discharge was larger than that of observations. Thus, the calibration will need
14 to adjust the parameters to reduce the discharge. Our study finds that for calibrating 22 parameters,
15 using the same computing resource for running WRF-hydro, the HPC-enabled PEST calibration
16 tool can speed up WRF-Hydro calibration by a factor of 15, compared with running PEST in
17 sequential mode. The speedup factor can be larger when the number of parameters needing
18 calibration is higher (e.g., 50 or 100).

19

20 The following are several key points that we would like to mention to inform future studies:

- 21 1. In this study, we consider using the prior or regularization information only for the
22 parameters that we calibrate. As is the case with solving inverse problems, prior
23 information is added to improve the smoothness of the solutions. In order to build a more
24 comprehensive calibration, an important aspect that can be considered is to enrich the prior
25 with the available historical data. For example, in this particular case, one can use the
26 historical observation data (e.g., April and May from the past few years) to enrich the prior
27 information for the parameters. Hence, the regularization objective function in PEST will
28 constitute not only the discrepancies between parameters and their “current estimates” but
29 also the discrepancies between WRF-Hydro simulations and preferred values (which is the
30 observed time series of historical discharge). Additionally, one can use the pilot points
31 technique described by Doherty (2005) in conjunction with parameter estimation to add

1 more flexibility to the calibration process. This will be potentially beneficial in improving
2 the predictions.

- 3 2. To focus on our main goal, we calibrate only the parameters in lookup tables. However,
4 we acknowledge that using a single value to represent a physics for a large domain could
5 be problematic, especially we expect the HPC-enabled parallel PEST to execute with
6 WRF-Hydro for large domains. This situation often needs parameter regionalization. For
7 example, WRF-Hydro v5 has many spatially distributed parameters available, such as the
8 overland flow roughness scaling factor (OVROUGHRTFAC), the factor of maximum
9 retention depth (RETDEPRTFAC), and the soil-related parameters (when compiled with
10 SPATIAL_SOIL=1). Calibrating these spatial parameters based on grid scale (e.g.,
11 catchments) rather than a single value will give the model more flexibility and thus better
12 fit the observations (Hundecha and Bardossy, 2004; Wagener and Wheater, 2006). In
13 practice, for example, one can include regional OVROUGHRTFACs (e.g., their
14 lower/upper bounds, and default values) in the PEST control file based on catchments.
15 However, the selection of the locations and sizes of catchment may introduce significant
16 uncertainties to the calibration results, which require systematic and comprehensive
17 investigation and understanding of the study area.
- 18 3. This study is limited to calibrating the observed streamflow only based on the format of
19 one of WRF-Hydro model outputs for individual station or point (frxst_pts_out.txt). It is
20 feasible, however, to calibrate other variables as long as the observation data is available.
21 For example, one can either find the closest point from the gridded dataset to the
22 observation location and then compare that model grid to observations; or one can change
23 the WRF-Hydro input/output code to output other variables in the frxst_pts_out.txt file, so
24 they can still use the same interface we developed here to calibrate other variables instead
25 in addition to the discharge.
- 26 4. The optimal parameter set obtained from this study is from the 5th iteration of parallel
27 PEST by testing five Marquardt lambdas. Testing different number of lambdas or
28 calibrating different number of parameters may generate a different set of optimal
29 parameters. These parameter sets can all make physical sense and be equally good for
30 reproducing observed discharges. This problem is named equifinality (Beven and Freer,
31 2001; Savenije, 2001), which is an important source of model uncertainty. To reduce the

1 model uncertainty through reducing the equifinality, hydrologists carry out additional
2 modelling objective for model evaluation to find more useful parameter sets (Mo and
3 Beven, 2004; Gallart et al., 2007). Alternatively, inspired by No. 3 discussed above, one
4 can calibrate the WRF-hydro model based on more than one variables, such as discharge
5 and soil moisture (or heat flux or water table depth) to reduce the number of optimal
6 parameter sets, and thus reduce the model uncertainty of predictions for these variables.

- 7 5. While this study ported the parallel PEST to HPC system and linked it to WRF-Hydro, we
8 note that BEOPEST is available in the PEST family. BEOPEST has the same functionality
9 as parallel PEST but uses a different approach for communication between master and
10 workers. Working with HPC-enabled BEOPEST may save total time cost since BEOPEST
11 uses the Transmission Control Protocol (TCP) and the Internet Protocol (IP) instead of
12 message files (reading input and writing output between master and works) for
13 communication. We expect it to be relatively straightforward to use BEOPEST to calibrate
14 WRF-hydro on HPCs since the interface remains similar, except one needs to copy the
15 template and instruction files in addition to the global files (see Section 3.1) into each
16 working folder.

17
18 *Data and Code availability.* The observed river discharge is downloaded from the USGS Surface-
19 Water Data website, available at <https://waterdata.usgs.gov/nwis/sw>. The Stage IV precipitation
20 data were downloaded from <https://data.eol.ucar.edu/dataset/21.093>. PEST was downloaded from
21 <http://www.pesthomepage.org/Downloads.php>. We use the Unix PEST version 13.6. The scripts
22 and files that are developed in this study and required by PEST for calibrating WRF-Hydro are
23 available at <http://doi.org/10.5281/zenodo.2588506>.

24
25 *Author contributions.* JW proposed the project and developed the study case in WRF and WRF-
26 Hydro. CW developed the scripts/code to port the parallel PEST to DOE supercomputers and adapt
27 it to work with WRF-Hydro. VR provided important input for the regularization calibration
28 method. AO operated the ArcGIS tool to delineate the high-resolution grid cells to include stream
29 channel network, open water, and groundwater/baseflow basins. RK provide high-level guidance
30 and insight for the entire project. All authors commented on this manuscript.

31

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

2

3 *Acknowledgments.* This work is supported under a Laboratory Directed Research and
4 Development (LDRD) Program at Argonne National Laboratory, through U.S. Department of
5 Energy (DOE) contract DE-AC02-06CH11357. Computational resources are provided by the
6 DOE-supported National Energy Research Scientific Computing Center, Argonne National
7 Laboratory Computing Resource Center, and Argonne Leadership Computing Facility. Our special
8 thanks to the PEST developers and entire WRF-Hydro team, especially Kevin Sampson for his
9 guidance on the ArcGIS tool. We gratefully thank the two reviewers for their valuable comments
10 and suggestions, which tremendously improved this manuscript.

11 **References**

12 Arnault, J., Wagner, S., Rumlmer, T., Fersch, B., Bliefernicht, J., Andresen, S., and Kunstmann,
13 H.: Role of runoff–infiltration partitioning and resolved overland flow on land–atmosphere
14 feedbacks: A case study with the WRF-Hydro coupled modeling system for West Africa, *J.*
15 *Hydrometeorol.*, 17, 1489–1516, 2016.

16

17 Campos, E., and Wang, J.: Numerical simulation and analysis of the April 2013 Chicago Floods,
18 *J. Hydrol.*, 531, 454–474, 2015.

19

20 Chen, F. and Dudhia, J.: Coupling an advanced land surface-hydrology model with the Penn State-
21 NCAR MM5 modeling system, Part I: Model implementation and sensitivity, *Mon. Weather Rev.*,
22 129, 569–585, 2001.

23

24 Cuntz, M., Mai, J., Samaniego, L., Clark, M., Wulfmeyer, V., Branch, O., Attinger, S., and Thober,
25 S.: The impact of standard and hard-coded parameters on the hydrologic fluxes in the Noah-MP
26 land surface model, *J. Geophys. Res. Atmos.*, 121, 10,676–10,700, doi:10.1002/2016JD025097,
27 2016.

28

29 Doherty, J.: PEST: Model Independent Parameter Estimation, User Manual, 6th ed., Watermark
30 Numerical Computing, Brisbane, Queensland, Australia, 2016.

1
2 Doherty, J.: Ground water model calibration using pilot points and regularization, *Groundwater*,
3 41(2), 170–177, 2005.
4
5 Gallart, F., Latron, J., Llorens, P., and Beven, K. J.: Using internal catchment information to reduce
6 the uncertainty of discharge and baseflow predictions. *Adv. Water Resour.* 30(4), 808–823, 2007.
7
8 Getirana, A. C. V.: Integrating spatial altimetry data into the automatic calibration of hydrological
9 models, *J. Hydrol.*, 387 (3-4), 244–255, doi: 10.1016/j.jhydrol.2010.04.013, 2010.
10
11 Gochis, D. J., Barlage, M., Dugger, A., FitzGerald, K., Karsten, L., McAllister, M., McCreight, J.,
12 Mills, J., RafieeiNasab, A., Read, L., Sampson, K., Yates, D., and Yu, W.: The WRF-Hydro
13 modeling system technical description, (Version 5.0). NCAR Technical Note. 107 pages.
14 Available online at:
15 <https://ral.ucar.edu/sites/default/files/public/WRFHydroV5TechnicalDescription.pdf>, 2018.
16
17 Hundecha, Y., and Bárdossy, A.: Modeling of the effect of land use changes on the runoff
18 generation of a river basin through parameter regionalization of a watershed model, *J. Hydrol.*,
19 292, 281–295, 2004.
20
21 Kerandi, N., Arnault, J., Laux, P., Wagner, S., Kitheka, J., and Kunstmann, H.: Joint atmospheric-
22 terrestrial water balances for East Africa: A WRF-Hydro case study for the upper Tana River basin,
23 *Theor. Appl. Climatol.*, 131, 1337–1355, doi: 10.1007/s00704-017-2050-8, 2018.
24
25 Lin, Y., and Mitchell, K. E.: The NCEP stage II/IV hourly precipitation analyses: Development
26 and applications, Preprints, 19th Conf. on Hydrology, San Diego, CA, Amer. Meteor. Soc., 1.2.,
27 2005.
28
29 Madsen, H.: Automatic calibration of a conceptual rainfall–runoff model using multiple
30 objectives, *J. Hydrol.*, 235, 276–288, 2000.
31

1
2 Mo, X., and Beven, K.: Multi-objective parameter conditioning of a three-source wheat canopy
3 model. *Agricultural & Forest Meteorol.* 122(1–2), 39–63, 2004.
4 Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D., and Veith, T. L.:
5 Model evaluation guidelines for systematic quantification of accuracy in watershed simulations,
6 *Transactions of the ASABE*, 50 (3), 885–900, 2007.
7
8 Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models, part I – A
9 discussion of principles, *J. Hydrol.*, 10(3), 282–290, doi: 10.1016/0022-1694(70)90255-6, 1970.
10
11 Niu, G.-Y., Yang, Z.-L., Mitchell, K. E., Chen, F., Ek, M. B., Barlage, M., Kumar, A., Manning,
12 K., Niyogi, D., Rosero, E., Tewari, M., and Xia, Y.: The community Noah land surface model with
13 multiparameterization options (Noah-MP): 1. Model description and evaluation with local-scale
14 measurements, *J. Geophys. Res.*, 116, D12109, doi: 10.1029/2010JD015139, 2011.
15
16 NWS (National Weather Service): Record river flooding of April 2013,
17 <https://www.weather.gov/ilx/apr2013flooding>, 2013.
18
19 Prat, O. P., and Nelson, B. R.: Evaluation of precipitation estimates over CONUS derived from
20 satellite, radar, and rain gauge data sets at daily to annual scales (2002-2012), *Hydrol. Earth Syst.*
21 *Sci.*, 19, 2037–2056, doi: 10.5194/hess-19-2037-2015, 2015.
22
23 Sampson, K., and Gochis, D.: WRF Hydro GIS Pre-processing tools, Version 5.0 Documentation,
24 2018.
25
26 Sapiano, M. R. P., and Arkin, P.A.: An intercomparison and validation of high-resolution satellite
27 precipitation estimates with 3-hourly gauge data, *J. Hydrometeorol.*, 10, 149–166, doi:
28 10.1175/2008JHM1052.1, 2009.
29
30 Senatore, A., Mendicino, G., Gochis, D. J., Yu, W., Yates, D. N., and Kunstmann, H.: Fully
31 coupled atmosphere-hydrology simulations for the central Mediterranean: Impact of enhanced

1 hydrological parameterization for short and long time scales, *J. Adv. Model. Earth Syst.*, 7(4),
2 1693–1715, doi: 10.1002/2015MS000510, 2015.

3

4 Soong, D. T., Prater, C. D., Halfar, T. M., and Wobig, L. A.: Manning’s roughness coefficients for
5 Illinois streams, U.S. Geological Survey Data Series 668, 2012.

6

7 Tonkin, M. J., and Doherty, J.: A hybrid regularized inversion methodology for highly
8 parameterized environmental models, *Water Resource Research*, 41, W10412,
9 doi:10.1029/2005WR003995, 2005.

10

11 Wagener, T., and Wheater, H. S.: Parameter estimation and regionalization for continuous
12 rainfall-runoff models including uncertainty, *J. Hydrol.*, 320, 132–154, 2006.

13

14 Xia, Y., Mitchell, K., Ek, M., Sheffield, J., Cosgrove, B., Wood, E., Luo, L., Alonge, C., Wei, H.,
15 Meng, J., Livneh, B., Lettenmaier, D., Koren, V., Duan, Q., Mo, K., Fan, Y., and Mocko, D.:
16 Continental-scale water and energy flux analysis and validation for the North American Land Data
17 Assimilation System project phase 2 (NLDAS-2), 1: Intercomparison and application of model
18 products, *J. Geophys. Res.*, 117, D03109, doi: 10.1029/2011JD016048, 2012a.

19

20 Xia, Y., Mitchell, K., Ek, M., Cosgrove, B., Sheffield, J., Luo, L., Alonge, C., Wei, H., Meng, J.,
21 Livneh, B., Duan, Q., and Lohmann, D.: Continental-scale water and energy flux analysis and
22 validation for the North American Land Data Assimilation System project phase 2 (NLDAS-2). 2.
23 Validation of model-simulated streamflow, *J. Geophys. Res.*, 117, D03110, doi:
24 10.1029/2011JD016051, 2012b.

25

26 Yucel, I., Onen, A., Yilmaz, K. K., and Gochis, D. J.: Calibration and evaluation of a flood
27 forecasting system: Utility of numerical weather prediction model, data assimilation and satellite-
28 based rainfall, *J. Hydrol.*, 523, 49–66, 2015.