



Post-Processing the National Water Model with Long Short-Term Memory Networks for Streamflow Predictions and Model Diagnostics

Jonathan M. Frame , Frederik Kratzert , Austin Raney II , Mashrekur Rahman , Fernando R. Salas, and Grey S. Nearing

Research Impact Statement: Deep learning improves upon streamflow predictions by the U.S. National Water Model (NWM). Post-processing the NWM does not improve predictive performance.

ABSTRACT: We build three long short-term memory (LSTM) daily streamflow prediction models (deep learning networks) for 531 basins across the contiguous United States (CONUS), and compare their performance: (1) a LSTM post-processor trained on the United States National Water Model (NWM) outputs (LSTM_PP), (2) a LSTM post-processor trained on the NWM outputs and atmospheric forcings (LSTM_PPA), and (3) a LSTM model trained only on atmospheric forcing (LSTM_A). We trained the LSTMs for the period 2004–2014 and evaluated on 1994–2002, and compared several performance metrics to the NWM reanalysis. Overall performance of the three LSTMs is similar, with median NSE scores of 0.73 (LSTM_PP), 0.75 (LSTM_PPA), and 0.74 (LSTM_A), and all three LSTMs outperform the NWM validation scores of 0.62. Additionally, LSTM_A outperforms LSTM_PP and LSTM_PPA in ungauged basins. While LSTM as a post-processor improves NWM predictions substantially, we achieved comparable performance with the LSTM trained without the NWM outputs (LSTM_A). Finally, we performed a sensitivity analysis to diagnose the land surface component of the NWM as the source of mass bias error and the channel router as a source of simulation timing error. This indicates that the NWM channel routing scheme should be considered a priority for NWM improvement.

(KEYWORDS: National Water Model; theory-guided machine learning; long short-term memory; streamflow; model diagnostics.)

INTRODUCTION

The United States (U.S.) National Water Model (NWM), based on WRF-Hydro (Cosgrove et al. 2015), is an emerging large-scale hydrology simulator. Some specific details of the NWM advancements in large-scale hydrology are described by Elmer (2019, p. 11), including increased resolution and number of stream reaches (2.7 million) for a model covering the contiguous United States (CONUS). A purported strength of

WRF-Hydro is simulating hydrologic dynamics, and specifically the timing of hydrologic response (Salas et al. 2018). The predictive performance of the NWM (ability to match streamflow observations) has been shown to vary widely. Hansen et al. (2019) evaluated the performance of the NWM in the Colorado River Basin in terms of drought and low flows; they found better performance in the Upper Colorado River Basin than in the Lower Colorado River Basin, and attributed this discrepancy to the NWM's ability to simulate snowpack. WRF-Hydro has generally poor

Paper No. JAWR-20-0099-P of the *Journal of the American Water Resources Association* (JAWR). Received July 27, 2020; accepted October 13, 2021. © 2021 American Water Resources Association. **Discussions are open until six months from issue publication.**

Department of Geological Sciences (Frame), University of Alabama Tuscaloosa, Alabama, USA; LIT AI Lab and Institute for Machine Learning (Kratzert), Johannes Kepler University Linz, Austria; Department of Geography (Raney), University of Alabama Tuscaloosa, Alabama, USA; Department of Land, Air and Water Resources (Rahman, Nearing), University of California, Davis Davis, California, USA; and NOAA National Water Center (Salas), Tuscaloosa, Alabama, USA (Correspondence to Frame: jmframe@crimson.ua.edu).

Citation: Frame, J.M., F. Kratzert, A. Raney II, M. Rahman, F.R. Salas, and G.S. Nearing. 2021. "Post-Post-Processing the National Water Model with Long Short-Term Memory Networks for Streamflow Predictions and Model Diagnostics." *Journal of the American Water Resources Association* 1–21. <https://doi.org/10.1111/1752-1688.12964>.

performance in the Southwest and Northern Plains (Salas et al. 2018). Salas et al. (2018) hypothesized that error in WRF-hydro might come from lakes, reservoirs, floodplain dynamics, and soil parameter calibration.

NOAA personnel calibrated the NWM (version 2.0) at 1,457 gauged basins within the CONUS domain. As a point of comparison, the U.S. Geological Survey (USGS) records daily streamflow at 28,529 basins (<https://nwis.waterdata.usgs.gov/nwis>, accessed June 2020). Calibrating the model at each stream gauge within the NWM domain (which include all of CONUS and many U.S. territories) is a large computational expense, and while regionalization strategies can be used to improve real-time forecast accuracy without having to calibrate each individual basin, accuracy typically suffers compared to direct calibration. Due to these reasons and others, making accurate hydrologic predictions over large scales is a challenging problem, however, there are promising results in the machine learning (ML) and data science communities that may be directly applicable to improving the NWM.

ML is a powerful tool for hydrologic modeling, and there has been a call to merge ML with traditional hydrologic modeling (Reichstein et al. 2019; Nearing et al. 2020). One example of an ML approach that has been effective for hydrologic prediction is the “long short-term memory” network (LSTM) (Hochreiter 1991; Hochreiter and Schmidhuber 1997). The LSTM is a time-series deep learning method that is particularly well suited to model hydrologic processes because it mimics in certain ways the Markovian input-state-output structure of a dynamical system (Kratzert et al. 2018). LSTMs have been effective at simulating predictions of surface runoff at the daily time scale (Kratzert, Klotz, Herrnegger, et al. 2019), including in ungauged catchments where traditional methods of calibration do not work (Kratzert, Klotz, Shalev, et al. 2019), and also at sub-daily (hourly) time scales (Gauch, Kratzert, et al. 2021). One potential problem with ML, however, is that it lacks a physical basis. While there are emerging efforts in hydrology to merge physical understanding with ML (Karpatne, Watkins, et al. 2017; Pelissier and Frame 2019; Read et al. 2019; Chadalawada and Herath 2020; Daw et al. 2020; Nearing et al. 2020; Tartakovsky et al. 2020; Hoedt et al. 2021), the field of *theory-guided ML* (Karpatne, Atluri, et al. 2017) is still relatively immature in hydrology.

The NWM informs forecasts of many hydrologic conditions, including river ice, snowpack, soil moisture, and inundation, which are used for management applications such as transportation, recreation, agriculture, and fisheries (NOAA 2019). When ML is to be used in the NWM it should not disrupt the delivery of these

hydrologic forecasts, therefore an ML prediction for streamflow that does not also include predictions of the other hydrologic states and variables must be run in parallel with the existing process-based hydrologic model. A natural question arises: does the existing NWM formulation benefit the already highly accurate LSTM predictions of streamflow?

Hydrologic post-processing can remove systematic errors in the model prediction, and has been shown to improve real-time forecast accuracy of both calibrated and uncalibrated basins, particularly in wet basins (Ye et al. 2014). The general methodology of post-processing involves taking the output of a process-based model and feeding it into a data-driven model. In this paper, we applied a LSTM-based post-processor for NWM basin-scale streamflow predictions. This is a straightforward theory-guided ML approach. We tested a LSTM-based post-processor that uses the dynamic NWM model outputs (shown in Table 1 and described below in the methods section) and compared the results against the NWM itself. We also tested a post-processor that included both the NWM outputs and atmospheric forcings as inputs and compared against an LSTM model trained only with atmospheric forcings (no NWM outputs).

We applied the LSTM post-processors to 531 basins across the CONUS. The basins chosen for this large-scale analysis are mostly headwater catchments without engineered control structures, such as dams, canals, and levees. This was a deliberate choice made for the purpose of simulating a close-to-natural rainfall-runoff response. Our goal was to use the post-processor to learn systematic corrections to simulated basin-scale rainfall-runoff processes that can improve forecasts of streamflow, rather than the hydraulic engineering implications resulting from simulated controlled flow, for example a reservoir release. Kim et al. (2020) showed the limitation of the NWM to predict streamflow in a highly engineered watershed and the need for representing controlled releases. Thus, we are using some of the simplest, and top performing, applications of the NWM for these experiments.

METHODS

Data and Models

CAMELS Catchments. This study used the Catchment Attributes and Meteorological dataset for Large Sample Studies (CAMELS) (Newman et al. 2015; Addor et al. 2017). The U.S. National Center for Atmospheric Research curated these data (NCAR;

TABLE 1. National Water Model (NWM) output data.

Feature name	Feature	NWM model component	Resolution
ACCET	Accumulated evapotranspiration	LDAS	1 km
FIRA	Total net long-wave (LW) radiation to atmosphere	LDAS	1 km
FSA	Total absorbed short-wave (SW) radiation	LDAS	1 km
FSNO	Snow cover fraction on the ground	LDAS	1 km
HFX	Total sensible heat to the atmosphere	LDAS	1 km
LH	Latent heat to the atmosphere	LDAS	1 km
SNEQV	Snow water equivalent	LDAS	1 km
SNOWH	Snow depth	LDAS	1 km
SOIL M (4 layers)	Volumetric soil moisture	LDAS	1 km
SOIL W (4 layers)	Liquid volumetric soil moisture	LDAS	1 km
TRAD	Surface radiative temperature	LDAS	1 km
UGDRNOFF	Accumulated underground runoff	LDAS	1 km
streamflow	River Flow	CHRT	Point
q _{lateral}	Runoff into channel reach	CHRT	Point
velocity	River Velocity	CHRT	Point
qSfcLatRunoff	Runoff from terrain routing	CHRT	Point
qBucket	Flux from groundwater bucket	CHRT	Point
qBtmVertRunoff	Runoff from bottom of soil to groundwater bucket	CHRT	Point
Sfcheadsbrt (mean and max)	Ponded water depth	RTOUT	250 km
Zwattablrt (mean and max)	Water table depth	RTOUT	250 km

<https://ral.ucar.edu/solutions/products/camels>, accessed March 2020), and we used the 531 (out of 671) basins that Newman et al. (2015) chose for model benchmarking. Newman et al. (2015) excluded basins with large discrepancies in different methods for measuring basin area and also basins larger than 2,000 km². CAMELS data include corresponding daily streamflow records from USGS gauges, and meteorological forcing data (precipitation, max/min temperature, vapor pressure, and total solar radiation) come from North American Land Data Assimilation System (NLDAS; Xia et al. 2012).

National Water Model. We used the NWM version 2.0 reanalysis, which contains output from a 25-year (January 1993 through December 2019) retrospective simulation (<https://docs.opendata.aws/nwm-archive/readme.html>, accessed June 2020). The NWM retrospective ingests rainfall and other meteorological forcings from atmospheric reanalyses (<https://water.noaa.gov/about/nwm>, accessed June 2020). NWM reanalysis output includes channel outputs (point fluxes: CHRT) and land surface (gridded states and fluxes: LDAS and RT) outputs. The specific features that we used from the NWM reanalysis are shown in Table 1. To be compatible with the LSTM model, which uses a one-day timestep and was trained using all basins simultaneously, we took the mean values of these model outputs across UTC calendar days (midnight–2300) to produce daily records from the hourly NWM when used as input to the LSTM, but for NWM streamflow diagnostics we used the local calendar day (based on U.S. time zone) to be compatible

with the USGS gauge records. We collected channel routing point data (CHRT) at each individual, NWM stream reach that corresponds to the stream gauge associated with each CAMELS catchment. We collected the gridded land surface data (LDAS) from each 1 km² Noah-MP cell (Niu et al. 2011) contained within the boundaries of each CAMELS catchment, and then calculated the averaged to produce a single representative (lumped) value for each catchment. We collected Gridded routing data (RT) from each 250 m² cell, and we included the mean and maximum value within the catchment boundary. We did not include lake input and output fluxes because these would be inconsistent across basins (some basins have zero and some basins have multiple lakes). Note that the units of the NWM outputs are not required for the LSTM post-processor.

LSTM Network. The LSTM is a recurrent neural network that is able to maintain a memory of the system state and dynamics through a period of time (in this case 365 days). This recurrent state space is the main advantage for hydrologic applications over other types of neural networks. We developed our LSTM network from Kratzert et al. (2018), Kratzert, Klotz, Herrnegger, et al. (2019), and Kratzert, Klotz, Shalev, et al. (2019) using a codebase that is now referred to as NeuralHydrology (<https://neuralhydrology.github.io/> accessed March 2021). NeuralHydrology was written in the Python programming language and is based primarily on the Pytorch ML library.

The LSTM in previous studies used two types of inputs: daily meteorological forcings and static

catchment attributes. Again, note that the units of the forcing data are irrelevant when used as inputs for the LSTM, which does not include a mass or energy balance. We normalized all inputs to the LSTM, including static and dynamic inputs by subtracting the mean and dividing by the standard deviation of the training data. We used 18 catchment attributes from the CAMELS dataset related to climate, vegetation, topography, geology, and soils. These are described in more detail by Addor et al. (2017) and listed here in Table 2. Catchment attributes are static for each basin (do not change in time). LSTMs are trained to make predictions that are appropriate for individual basins according to their static attributes (Kratzert, Klotz, Shalev, et al. 2019), allowing us to train a single model that can be applied on any basin (we tested them on 531 CAMELS basins). The static attributes position a particular basin within an input space that is suitable for a particular hydrologic response (Nearing et al. 2021). For instance, the geologic permeability may influence the mass difference between total rainfall and runoff in a particular basin, as it would as a parameter in a process-based model. For the post-processing runs, we added the NWM model output predictions from version 2.0 of the NWM shown in Table 1.

We trained the LSTM models to make predictions at all 531 CAMELS catchments used in the analysis. We split the data temporally into a training period and testing period, and we present no results from the training period as these results are unrepresentative of the out-of-sample predictions. We trained the LSTMs on water years 2004 through 2014 and tested on water years 1994 through 2002. We included no spatial splits in the training procedure. The LSTMs used a 365-day LSTM look-back period, so a full year gap was left between training and testing to prevent bleedover (i.e., information exchange) between the two periods. We trained separate LSTMs with 10 unique random seeds for initializing weights and biases, and calculated benchmarking statistics using the ensemble mean hydrograph. The LSTMs make predictions representing runoff in units [mm], reflecting an area normalized volume of water that moves through a stream at each model time step. USGS gauge records (and the NWM predictions) are in streamflow units [L^3/T]. We used the geospatial fabric estimate of the catchment area provided in the CAMELS dataset to convert all streamflow to units [L] for our diagnostic comparison. We trained the LSTMs with the protocol and features described in appendix B of Kratzert, Klotz, Shalev, et al. (2019): this includes 30 epochs, a hyperbolic tangent activation function, a hidden layer size of 256 cell states, a look-back of 365 days, variable learning rates set at

TABLE 2. North American Land Data Assimilation System forcings and static catchment attributes.

Meteorological forcing data (used only in models denoted with an “A”)	
Maximum air temp (T_{Max})	2-m daily maximum air temperature
Minimum air temp (T_{Min})	2-mr daily minimum air temperature
Precipitation (PRCP)	Average daily precipitation
Radiation (SRAD)	Surface-incident solar radiation
Vapor pressure (V_p)	Near-surface daily average
Static catchment attributes (used in each of the LSTM models)	
Precipitation mean	Mean daily precipitation
PET mean	Mean daily potential evapotranspiration
Aridity index	Ratio of mean PET to mean precipitation
Precipitation seasonality	Estimated by representing annual precipitation and temperature as sin waves positive (negative) values indicate precipitation peaks during the summer (winter). Values of approx. 0 indicate uniform precipitation throughout the year
Snow fraction	Fraction of precipitation falling on days with temp [C]
High precipitation frequency	Frequency of days with $\leq 5 \times$ mean daily precipitation. Average duration of high precipitation events (number of consecutive days with $\leq 5 \times$ mean daily precipitation)
Low precipitation frequency	Frequency of dry days (<1 mm/day)
Low precipitation duration	Average duration of dry periods (number of consecutive days with precipitation <1 mm/day)
Elevation	Catchment mean elevation
Slope	Catchment mean slope
Area	Catchment area
Forest fraction	Fraction of catchment covered by forest
LAI max	Maximum monthly mean of leaf area index
LAI difference	Difference between the max. and min. mean of the leaf area index
GVF max	Maximum monthly mean of green vegetation fraction
GVF difference	Difference between the maximum and minimum monthly mean of the green vegetation fraction
Soil depth (pelletier)	Depth to bedrock (maximum 50 m)
Soil depth (STATSGO)	Soil depth (maximum 1.5 m)
Soil porosity	Volumetric porosity
Soil conductivity	Saturated hydraulic conductivity
Max water content	Maximum water content of the soil
Sand fraction	Fraction of sand in the soil
Silt fraction	Fraction of silt in the soil
Clay fraction	Fraction of clay in the soil
Carbonate rocks fraction	Fraction of the catchment area characterized as “carbonate sedimentary rocks”
Geological permeability	Surface permeability (log10)

Note: LSTM, long short-term memory.

epoch 0 to 0.001, epoch 11 to 0.005 and epoch 21 to 0.0001, dropout rate of 0.4 and an input sequence length: 270.

Overfitting of deep learning models can lead to poor performance when the models make predictions on data that is not part of the training set. The methods described above to ensure that information in the testing set (water years 1994 through 2002) is not part of the training set helps build confidence in our modeling results. In addition, the dropout rate is an important hyper-parameter for preventing overfitting. The dropout probabilistically removed some connections from the LSTM network during training, in our case with a probability of 0.4. This avoids the network relying too heavily on specific connections. Model runs during testing did not include dropout.

Experimental Design

We tested the results from LSTM post-processing against the NWM and also against a LSTM trained with atmospheric forcings as dynamic inputs to the model, with no inputs from the NWM model outputs (referred to as LSTM_A, in which the A stands for atmospheric forcing). Table 3 will guide the reader through the setup of each model.

Simple schematics of the LSTMs used in this study are shown in Figure 1. The LSTM post-processors (LSTM_PP and LSTM_PPA) used NWM outputs as LSTM inputs, and the process-based NWM predictions influenced the LSTM-based streamflow predictions. This is a straightforward method of theory-guided (or physics-informed) ML, but is commonly referred to as post-processing (Han 2021).

As a quality check, we compared the results from each LSTM ensemble member, and found a relative standard error of the mean streamflow about 1%, and relative standard error of the Nash–Sutcliffe efficiency (NSE) value of about 0.01%. This means that all LSTM solutions are similar between random initialization seeds. Gauch, Mai, et al. (2021) attributed

TABLE 3. Models.

Model label	Number of dynamic LSTM inputs	Model description
NWM	N/A	NWM mean daily streamflow predictions
LSTM_PP	28	LSTM trained with NWM output for post-processing
LSTM_PPA	33	LSTM trained with NWM output and atmospheric forcings for post-processing
LSTM_A	5	LSTM trained with atmospheric forcing conditions

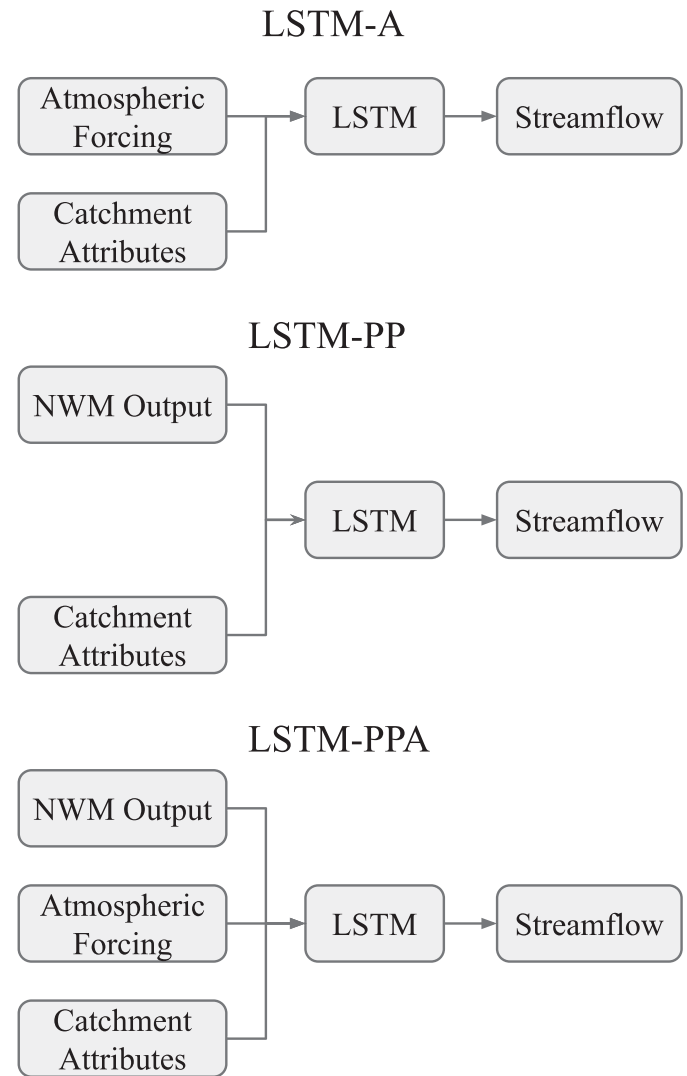


FIGURE 1. Flowchart showing the LSTM_A and the LSTM post-processors with NWM data as inputs (LSTM_PP and LSTM_PPA). LSTM_PP is the post-processor which used only NWM outputs as input to an LSTM, and LSTM_PPA used both the NWM outputs and atmospheric forcings.

a 0.01 discrepancy in NSE values of the LSTM predictions to nondeterminism of the loss function minimization. In our experiments discrepancies in the loss function occur between different random seed initializations, but running the training procedure twice with the same random seed gives an identical solution, satisfying the definition of determinism.

Model comparisons. We tested/evaluated all models (NWM and all LSTMs) on the same daily data and the same time period (years 1994–2002). We trained the LSTMs on data from years 2004–2014 and evaluated them on years 1994–2002. The NWM was calibrated by NOAA on the time period 2007–2013 (https://ral.ucar.edu/sites/default/files/public/9_Rafieei_Nasab_CalibOverview_CUAHSI_Fall019_0.pdf, accessed

August 2021), though no journal publications thoroughly describe the details of this calibration. For this study, we tested the performance of the NWM reanalysis only on the time period 1994–2002 (the same time period as the LSTM).

Performance Metrics. We calculated several metrics to evaluate predictive performance, including the NSE and Kling–Gupta efficiency (KGE) values (Gupta et al. 2009). We calculated the variance, bias, and Pearson correlation metrics separately as components of the NSE (Gupta et al. 2009); these tell us about relative variability, mass conservation, and linear correlation between the modeled/observed streamflow values, respectively. Observed streamflow values are from the USGS streamflow gauges associated with each of the CAMELS basins. We calculated the metrics in two ways: (1) at each basin and then averaged together, and (2) using all of the flows from all basins combined.

Our graphical results focus on three performance metrics: (1) NSE measures the overall predictive performance as a correlation coefficient for the 1:1 linear fit between simulations and observations, (2) Peak timing error measures the absolute value of differences (in units days) between simulated and observed peak flows for a given event, and (3) total (absolute) bias measures the overall bias of the simulated hydrograph relative to observations and represents how well the model matches the total volume of partitioned rainfall that passes through the stream gauge at each basin.

We also calculated performance metrics on different flow regimes. Rising limbs and falling limbs were characterized by a one-day derivative, where positive derivatives were categorized as rising limb, and negative derivatives as falling limb. High flows were characterized as all flow above the 80th percentile in a given basin, and low flows as below the 20th percentile in a given basin.

We tested the performance of the LSTM post-processors in different regions. We split the basins by USGS designated “water resource regions” (<https://water.usgs.gov/GIS/regions.html>, accessed July 2020). To analyze the regions individually we averaged the NSE, bias, and timing error of the CAMELS basins within each region.

We set an alpha value for statistical significance to $\alpha = 0.05$. To control for multiple comparisons we adjusted the alpha values using family-wise error rate equal to $1 - (1 - \alpha)^m$, with m being the number of significance tests (86 in total), which brought our effective alpha value down to 0.049. We tested for statistical significance with a Wilcoxon signed-rank test against the null hypothesis that our test models (LSTM post-processors) performance across basins

came from the same distribution as our base models (NWM and LSTM_A).

Simulated Hydrograph Representation of Hydrologic Signatures. Hydrologic signatures help us understand how well a model represents important aspects of real-world streamflow, and where improvement should be made to the model’s conceptualization (Gupta and Wagener 2008). We analyzed the hydrologic signatures described by Addor et al. (2018), and these are listed below in Table 4. We calculated the true signatures with USGS streamflow observations, and calculated model representations with predicted values of daily streamflow. We compared true values and predicted values with a correlation coefficient (r^2) across basins (one value of the observed and predicted hydrologic signatures were calculated per basin), higher values indicate a better representation of hydrologic signature across basins by the model. We used the Steiger method to test for statistically significant changes between the LSTM_A, NWM, and the LSTM post-processor (Steiger and Browne 1984).

Identifying Basins Best Suited for Post-Processing with Multi-Linear Regression. The LSTM post-processors did not improve performance at every basin. It therefore would be valuable to know if a LSTM post-processor will work in any particular basin before implementation. We trained a multi-linear regression, using the Scikit-learn library in Python, to predict the performance changes between the NWM and the LSTM post-processors (LSTM_PP and LSTM_PPA) at each individual basin. The multi-linear regression analysis included performance scores of the NWM streamflow predictions, hydrologic signatures, and catchment characteristics as inputs. These regressors are useful to help interpret what basins might benefit most from an LSTM

TABLE 4. Hydrologic signatures (adapted from Addor et al. 2018).

Signature description	Signature name
Average duration of low-flow events	low_q_dur
Frequency of days with zero flow	zero_q_freq
Average duration of high-flow events	high_q_dur
Streamflow precipitation elasticity	stream_elas
Frequency of high-flow days	high_q_freq
Slope of the flow duration curve	slope_fdc
Frequency of low-flow days	low_q_freq
Baseflow index	baseflow_index
Runoff ratio	runoff_ratio
Mean half-flow date	hfd_mean
5% flow quantile	q5
95% flow quantile	q95
Mean daily discharge	q_mean

post-processor. We trained and tested multi-linear regression models using k -fold cross-validation with 20 splits ($k = 20$) over the 531 basins. We report the correlation (r^2) of out-of-sample regression predictions of post-processing changes vs. actual post-processing changes.

Interpretation of LSTM with Integrated Gradients. We aim to explain the relationship between a model's predictions in terms of its features. This will help us understand feature importance, identifying data issues, and inform NWM process diagnostics from the post-processors. We calculated integrated gradients (Sundararajan and Taly 2017) to attribute the LSTM inputs (both atmospheric forcings and NWM outputs) to the total prediction of streamflow. Integrate gradients are a type of sensitivity analysis that are relatively insensitive to low gradients (e.g., at the extremes of neural network activation functions). We calculated integrated gradients separately for each input, at each timestep, for each lookback timestep, in each basin. This means that for nine years of test data with a 365-day lookback there were about 1.2 million integrated gradients per input, per basin. The unit of the integrated gradient is technically normalized streamflow, but we were mostly interested in the relative values of integrated gradients of each individual LSTM input.

Interpretation of LSTM with Correlations between Performance and NWM Inputs. We made a direct connection between LSTM post-processor improvements with the NWM outputs using correlation. We calculated Pearson R values between the basin average value of each NWM input feature and the total performance change (NSE, bias, and peak timing). We calculated these correlations for different flow regimes (all flows using the whole hydrograph, rising/falling limbs using the single day differentials, and high/low flows using the top 80% and bottom 20%). The strengths of these correlations (positive or negative) indicated which types of basins (via NWM features) are benefiting most from a LSTM post-processor. Results for rising limbs and falling limbs of the hydrograph were qualitatively similar to this figure, and were therefore omitted.

Splitting the CAMELS Catchments by Calibrated/Uncalibrated. Of the NWM calibrated basins, 480 overlap with the 531 CAMELS catchments used in this study. In a separate set of experiments, we trained the LSTM_A and the LSTM post-processors LSMT_PP and LSTM_PPA) on only the 480 calibrated basins. We then used the full set of 531 catchments to test the performance out-of-sample.

We analyzed the 480 in-sample basins and 51 out-of-sample basins separately using the NSE, bias, and timing error metrics. This allowed us to determine if the LSTM is a suitable post-processing method to use in uncalibrated basins. If the post-processors trained only on calibrated basins can improve streamflow predictions at uncalibrated basins, then they would be considered suitable, particularly if there is no statistical difference between the post-processor's performance improvement over the NWM and/or LSTM_A.

Sensitivity Analysis and NWM Process Diagnostics. We trained a set of LSTM post-processors using different combinations of NWM outputs as input to the LSTM, as described in Table 5. To test the sensitivity to the NWM streamflow prediction itself, we trained an LSTM with only streamflow (LSTM_Q_only), and excluded it from another (LSTM_PP_noQ). We tested the sensitivity to the channel routing (LSTM_chrt) and land surface (LSTM_ldas) components of the NWM by training LSTMs with only these dynamic inputs. We trained these models with the same specifications as the LSTM_A, LSTM_PPA, and LSTM_PP.

Each of these models (Table 5), in addition to the main post-processing models presented in Table 3, have a distinct flow of information that we can use to diagnose NWM model processes. Figure 2 shows the information flow of each of the model subcomponents. We used the performance results of the different post-processing models to assess how much information passes between the model components. Nearing et al. (2018) described the method to quantify the information exchange down a modeling chain (i.e., integrating over the expected effect of the conditional probability), but since we used limited outputs from the NWM reanalysis, rather than the full state space, we examined the NWM only qualitatively for information loss between the major NWM subcomponents (land surface runoff, overland router, and channel

TABLE 5. Additional models for sensitivity analysis and NWM diagnostics.

Model label	Number of dynamic LSTM parameters	Model description
LSTM_PP_noQ	26	LSTM post-processor (LSTM_PP) but without streamflow or velocity
LSTM_Q_only	1	LSTM trained with NWM streamflow only
LSTM_chrt	6	LSTM trained with NWM channel routing outputs only
LSTM_ldas	18	LSTM trained with NWM land surface outputs only

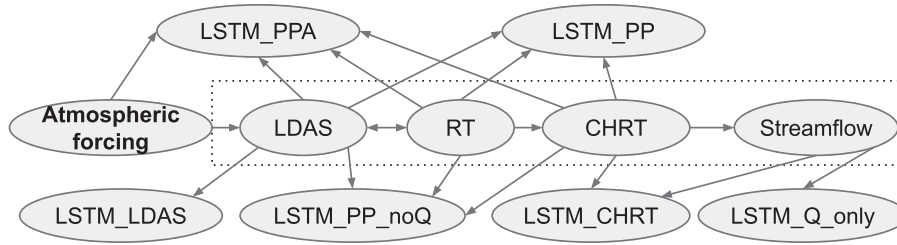


FIGURE 2. Process network diagram showing the information flow of each of these models. Arrows indicate the information flow from one component of the model to another. The NWM components are outlined with the dashed box. This is also a good guide for understanding the inputs to the different post-processing models.

router). The LSTM extracts information from its input to make predictions about its target, in our case streamflow, and we assumed higher streamflow prediction accuracy indicated more information is available in the NWM components used as input. If a post-processor made less accurate streamflow predictions than the LSTM_A, then this indicates that the NWM modeling chain lost information from the atmospheric forcings.

RESULTS

Overall Model Performance

Post-processing the NWM with LSTMs significantly improved predictive performance, both with or without including the atmospheric forcings as inputs into the model. The LSTM_A, however, is the overall better performing model. Figure 3 shows the cumulative distributions of three performance metrics (NSE, peak timing error, and total bias).

The LSTM_PP improved the NSE score of the NWM mean daily streamflow at a total of 465 (88%) and reduced accuracy in 66 basins (12%) of the total 531 CAMELS basins, improved the total bias of the NWM mean daily streamflow at a total of 325 (61%) of basins and improved the peak timing error at a total of 488 (92%) of basins. The LSTM_PPA post-processor improved the NSE score of the NWM mean daily streamflow at a total of 488 (92%) and reduced accuracy in 43 basins (8%) of the total 531 CAMELS basins. The LSTM_PPA post-processor improved the total bias of the NWM mean daily streamflow at a total of 331 (62%) of basins and improved the peak timing error at a total of 494 (93%) of basins. The LSTM_A (without NWM model output) outperformed the NWM at a total of 473 (89%) and reduced accuracy in 58 basins (11%), improved the total bias of the NWM mean daily streamflow at a total of 339 basins (64%) and improved the peak timing error at a

total of 484 basins (91%). The LSTM_PPA improved the greatest number of basins in terms of NSE and peak timing error and the LSTM_A was the best performing model in terms of total bias. Figure 4 shows scatter plots of the post-processor performance at individual basins against the performance of the NWM and LSTM_A.

The post-processing models (LSTM_PP and LSTM_PPA) improved relative to the NWM in similar basins. The improvements of the two post-processing methods are correlated across all basins ($r^2 = 0.995$). Performance comparisons between the LSTM models and the NWM for each basin are plotted spatially in Figure 5. Notice that some of the highest NSE improvements between the LSTM_PP and the NWM are the worst NSE detriments between the LSTM_PPA and the LSTM_A, particularly in the northern plains. This indicates that although the post-processor greatly improves the NWM, the information from the NWM at bad basins hinders the performance of the LSTM, or in other words, the NWM passes bad information to the LSTM.

Performance by Flow Regime

The LSTM post-processors improved the predictive performance of the NWM according to the NSE and KGE metrics, as well as their components (variance and correlation). A full set of performance metrics broken down by flow regime are shown in Table 6. The left side of the table shows the average of metrics calculated individually at each basin, and the right side of the table shows the metrics as calculated by combining the flows from all basins. The NSE includes both mean and median averages, but the rest of the metrics are only averaged by the median.

In general Table 6 shows that the LSTM post-processors improved over the NWM in nearly all flow regimes according to most metrics. The LSTM_PPA also improved upon the LSTM_A in more than half the basins, and by most metrics, though not significantly. The prediction of rising limb and high flow

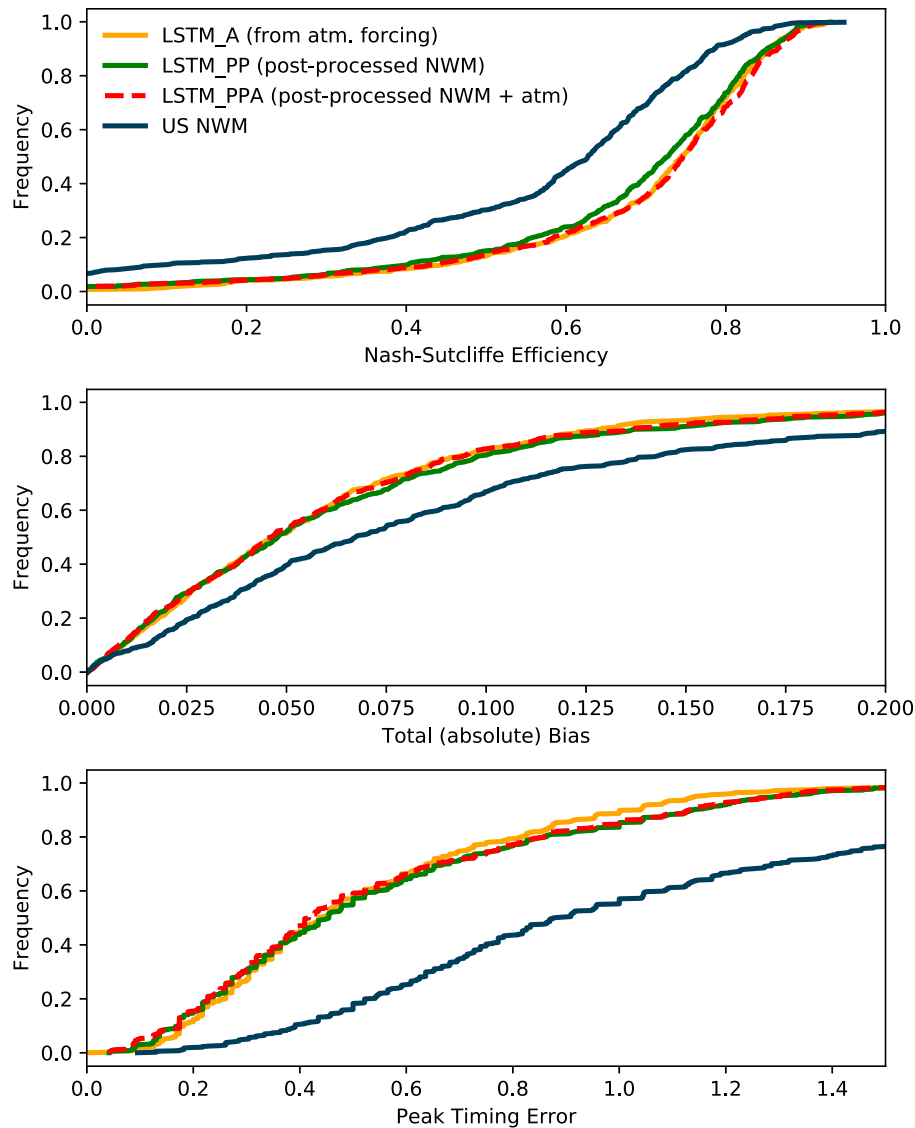


FIGURE 3. Results showing the cumulative distributions of model performance calculated as Nash–Sutcliffe efficiency (NSE), total bias, and peak timing error over a 10-year test period in 531 CAMELS catchments. The NWM reanalysis streamflow was averaged daily, LSTM networks shown used (1) the original atmospheric inputs (LSTM_A), (2) NWM states and fluxes only (LSTM_PP), and (3) both atmospheric forcings and NWM states and fluxes (LSTM_PPA). These figures omit the distribution tails for clarity.

regimes was improved upon by the LSTM post-processors according to every metric.

Bias was the only metric that was reduced due to post-processing, and the difference was highest in low flow regimes. All models poorly predicted flows below the 20th percentile. This is likely due to the fact that all models tend to have difficulty predicting zero streamflow, and the 101 basins with periods of zero streamflow affected the average performance metrics. This will be discussed further in terms of hydrologic signatures.

The right side of the table has better performance values than the average of metrics calculated

individually at each basin. This is a result of some of the better performing basins compensating for poorer performing basins, or from a different perspective, some basins have a relatively poor performance which weighs down the average.

Performance by Region

Results from a regional analysis of performance are shown below in Figure 6. The LSTM post-processors significantly improved the NSE over the NWM in 15 of the 18 regions, the peak timing error

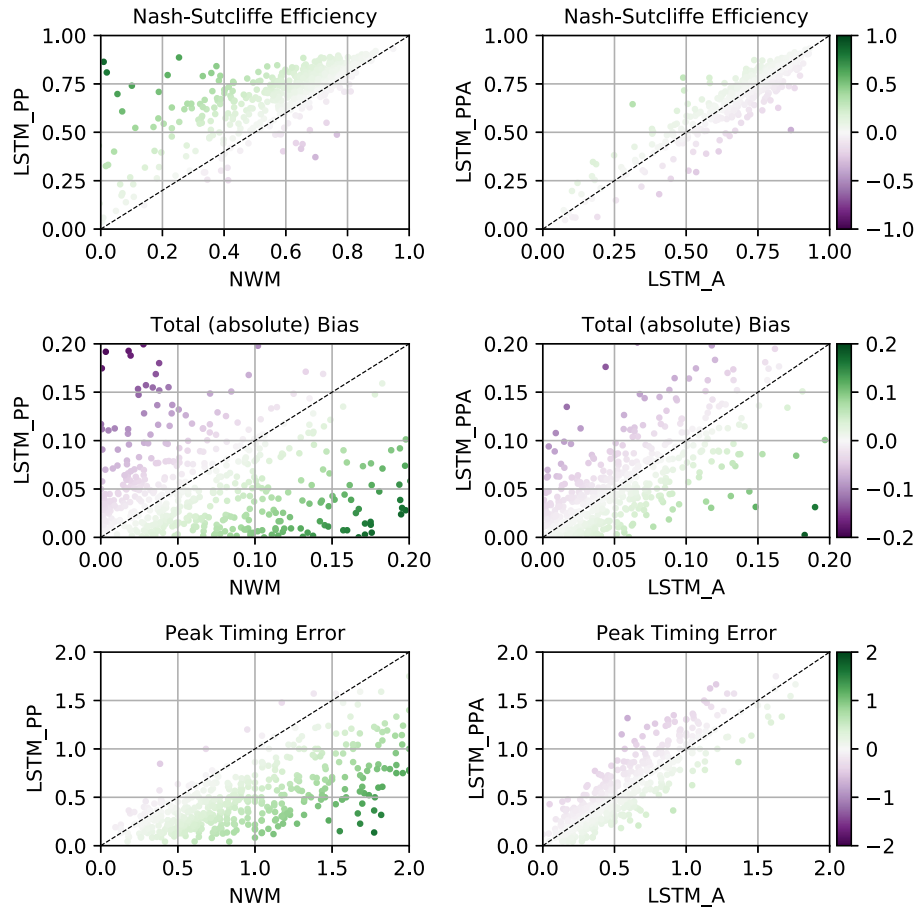


FIGURE 4. Performance differences of the post-processors against the NWM and LSTM_A in 531 CAMELS basins across contiguous United States (CONUS). Green indicates basins where post-processing improved performance over the NWM and LSTM_A (darker indicates larger relative improvement), and purple indicates basins where there was a decrease in performance (darker indicating worse relative detriment). The first column shows the performance difference between the LSTM_PP and the NWM. The second column shows the performance difference between the LSTM_PPA and the LSTM_A.

in 16 regions (all regions with enough basins for a statistical evaluation) and significantly improved bias in only one region. Note that region 9 was represented by only two CAMELS basins, which is not sufficient for statistical evaluation. The bias was better represented by the NWM than the post-processor in five of the 18 regions, including the entire East Coast (regions 1, 2, and 3), the Pacific Northwest (17), and the Lower-Colorado River (15).

The regional performance of the LSTM post-processors and the regional performance of the LSTM_A were correlated with the regional performance of the NWM in terms of NSE ($r^2 = 0.78$ for post-processors and 0.63 for LSTM_A) and peak timing error ($r^2 = 0.96$ for post-processors and 0.92 for LSTM_A), but not in terms of bias ($r^2 = 0.24$, calculated on bias although absolute bias is plotted for clarity). The post-processors and the LSTM_A are correlated in terms of their bias ($r^2 = 0.91$). A better model has a higher NSE, bias closer to zero, and a lower timing error.

Regression to Predict Post-Processing Performance Improvement

The performance of the LSTM_A was more predictable than the post-processors. We performed a multi-linear regression on the target of performance improvement over the NWM, with inputs being the catchment attributes and hydrologic signatures, as well as the NWM performance itself. Figure 7 shows the results predicting the LSTM improvement over the NWM at each basin with an r^2 value of 0.97 , 0.88 , and 0.89 for the LSTM_A, LSTM_PPA, and LSTM_PP, respectively. The high r^2 value is due in part to the outlier basins with abnormally large performance improvements from the LSTM models (LSTM_A, LSTM_PPA, and LSTM_PP). This means that the magnitude of the LSTM_A and post-processors improvement is directly related to the performance of the NWM.

The aim of these results is to understand whether it is possible to predict where post-processing might

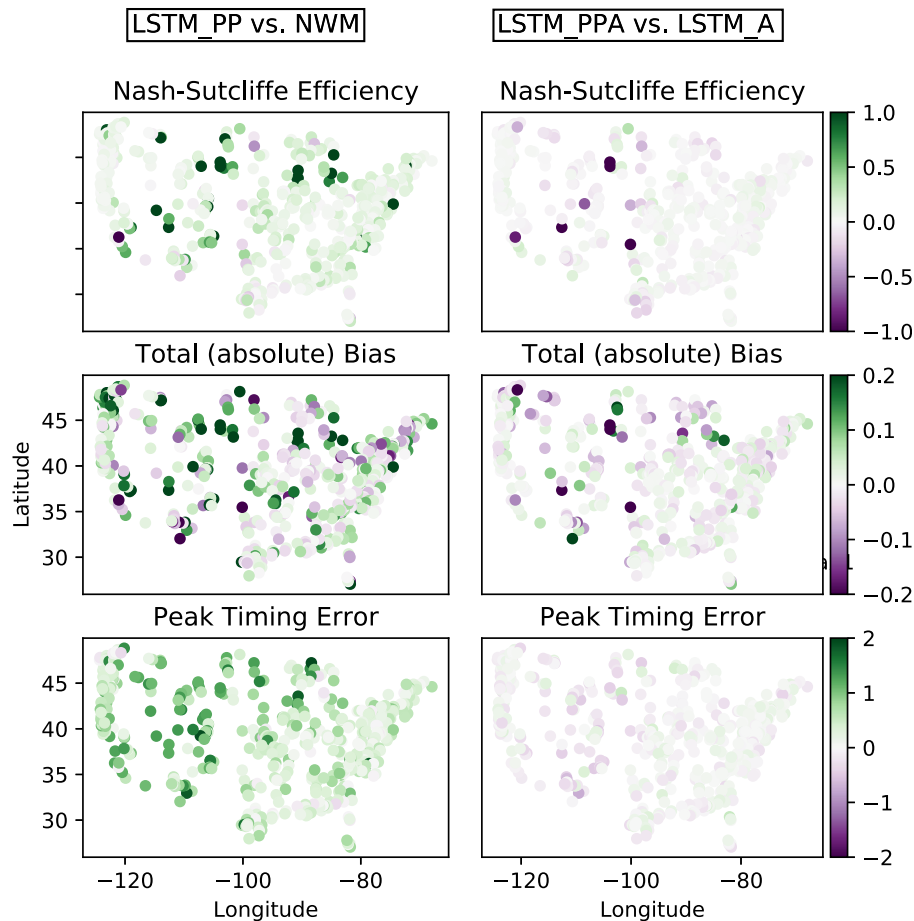


FIGURE 5. Per-basin performance change between the post-processors and NWM and LSTM_A in 531 CAMELS basins across CONUS. Green indicates basins where post-processing improved performance over the NWM and LSTM_A (darker indicates larger relative improvement), and purple indicates basins where there was a decrease in performance (darker indicating worse relative detriment). The first column shows the performance change between the LSTM_PP and the NWM. The second column shows the performance change between the LSTM_PPA and the LSTM_A.

be beneficial (remember that post-processing helped in most basins). Although we found relatively high predictability in the improvement expected from post-processing, a problem is that this requires knowing ahead of time the NWM performance. This prevents us from predicting post-processing improvement in *ungauged* basins, since calculating the NWM performance requires streamflow observations. The correlation analysis below may help inform future efforts to learn general patterns of post-processor improvement over both the NWM and the LSTM_A.

Correlations between NWM Inputs and Improvements

Figure 8 shows correlations (over 531 basins) between the time-averaged NWM inputs and changes in performance metric scores of the post-processor relative to the NWM and LSTM_A. The LSTM_PP was compared against the NWM and the LSTM_PPA

was compared against the LSTM_A, although qualitatively both post-processor models were similar. The rows of this figure show that correlation was weaker for differences in NSE score than total bias and peak timing error. Performance differences between the NWM and the post-processor were most strongly (anti)correlated with stream velocity from the channel router and accumulated underground runoff from the land surface model component: basins with lower stream velocity (velocity) and less underground runoff (UGDRNOFF) saw greater performance improvement from (daily) post-processing. This means that in basins with high underground runoff and/or high stream velocity the post-processor improvements were smaller. In contrast, basins with higher total radiation (TRAD) and higher latent heat flux (LH) saw greater improvement due to post-processing. This means that in basins with more radiation and heat flux the post-processor improvements were larger. A direct interpretation of this could be that a flat

TABLE 6. Predictive performance for NWM, LSTM_A, and the LSTM post-processors during various flow regimes. The NSE and Kling-Gupta efficiency (KGE) are overall performance metrics of prediction quality. Variance, bias, and correlation (R) are the components of the NSE. We calculated these in two ways: (1) at each basin and averaged across all basins, and (2) once using the observed and predicted streamflow values from all basins combined. Note that calculations done once across all basins do not include a test of significance.

Flow categories	Calculated per-basin						All basins			
	NSE (mean)	NSE (median)	KGE	Variance	Bias	R	NSE	Variance	Bias	R
All flows										
NWM	0.46	0.62	0.64	0.82	-0.01^1	0.82	0.75	0.85	-0.02	0.87
LSTM_PP	0.65^2	0.73^2	0.74^2	0.86	0.02	0.87^2	0.81	0.92	0.02	0.90
LSTM_A	0.69	0.74	0.74	0.83	0.02	0.88	0.82	0.89	0.01	0.90
LSTM_PPA	0.67	0.75	0.76	0.87	0.02	0.88	0.82	0.93	0.02	0.91
Rising limbs										
NWM	0.47	0.60	0.60	0.77	-0.07	0.81	0.73	0.82	-0.05	0.85
LSTM_PP	0.64^2	0.70^2	0.72^2	0.83^2	0.00^2	0.86^2	0.78	0.88	0.00	0.88
LSTM_A	0.66	0.71	0.72	0.80	-0.01	0.86	0.78	0.85	-0.01	0.88
LSTM_PPA	0.65	0.72	0.74	0.85	0.00	0.87	0.79	0.89	0.00	0.89
Falling limbs										
NWM	0.29	0.62	0.64	0.94	0.03	0.83	0.78	0.90	0.00	0.88
LSTM_PP	0.62^2	0.75^2	0.76^2	0.95^2	0.07	0.90^2	0.87	0.99	0.04	0.93
LSTM_A	0.69	0.78	0.77	0.92	0.05	0.90	0.87	0.96	0.03	0.93
LSTM_PPA	0.65	0.77	0.77	0.94	0.05	0.90	0.87	0.98	0.03	0.93
Above 80th percentile										
NWM	0.17	0.41	0.54	0.80	-0.13	0.73	0.69	0.83	-0.10	0.84
LSTM_PP	0.47^2	0.57^2	0.64^2	0.82	-0.08^2	0.80^2	0.76	0.89	-0.04	0.90
LSTM_A	0.53	0.58	0.67	0.81	-0.08	0.81	0.78	0.86	-0.06	0.88
LSTM_PPA	0.50	0.60	0.69	0.84	-0.07	0.81	0.79	0.90	-0.04	0.89
Below 20th percentile										
NWM	$-18,384.37$	-17.47	-1.96	3.79	1.89^1	0.36	0.37	1.31	0.22	0.81
LSTM_PP	$-6,941.62^2$	-15.66^2	-1.28^2	2.84^2	3.21	0.43^2	0.53	1.30	0.33	0.90
LSTM_A	$-4,749.68$	-16.35	-1.31	2.85	3.27	0.43	0.56	1.26	0.33	0.89
LSTM_PPA	$-5,147.62$	-14.66	-1.24	2.85	2.87	0.43	0.58	1.28	0.30	0.90

¹Post-processing significantly hurts the NWM.

²Post-processing significantly helps the NWM.

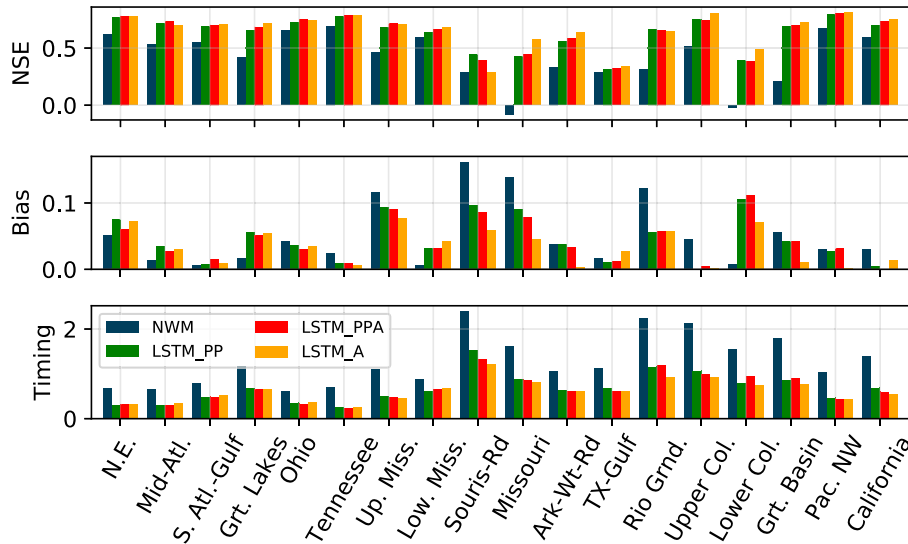


FIGURE 6. Regionally averaged performance metrics for NWM, LSTM_A, and the LSTM post-processors (LSTM_PP and LSTM_PPA) in different United States Geological Survey water resources regions.

meandering stream in the Southwest will benefit from post-processing, which is consistent with the findings of Salas et al. (2018) that WRF-Hydro's

performance is generally poor in the Southwest. Performance differences between the LSTM_A and the post-processor were most strongly correlated with

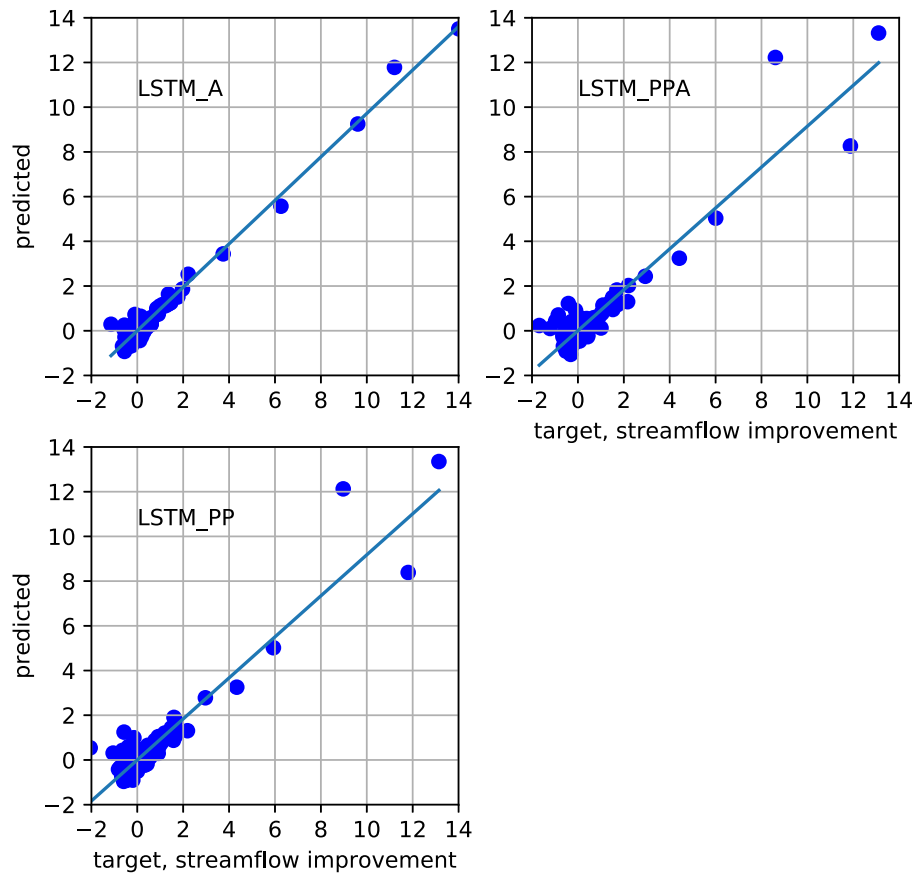


FIGURE 7. Predicting LSTM_A, LSTM_PP, and LSTM_PPA performance over the NWM at each basin using a linear regression with NWM performance and hydrologic signatures as inputs. Scatter plots with all of the 531 basins.

snow water equivalent and snow depth. This is consistent with the findings of Hansen et al. (2019) that the NWM represents snowpack hydrology well.

Integrated Gradients

Figure 9 shows the relative strength of the total attribution of the dynamic inputs to the LSTM_PPA averaged across the entire validation period and across basins. The ordered magnitudes of the integrated gradients can be interpreted as corresponding to the order of importance of inputs. The most important dynamic features for the LSTM_PPA were: (1) precipitation from NLDAS, and (2) routed streamflow from the NWM point data. Precipitation inputs were weighted higher than the NWM streamflow output itself, which means that even when NWM streamflow data were available, the LSTM_PPA generally learned to get information directly from forcings rather than from the NWM streamflow output. This indicates that the LSTM_PPA generated a new rain-fall-runoff relationship rather than relying on the

NWM, which is consistent with the overall results (Figure 2) that showed similar performance between the LSTM_A and LSTM_PPA.

Figure 10 shows the relative strength of the total attribution of the dynamic inputs to the LSTM_PP. Without the atmospheric forcings included in the post-processor inputs, the NWM streamflow output was by far the highest contributing dynamic input feature to the LSTM_PP. The static permeability of the catchment was the next highest.

Representations of Hydrologic Signatures

Results of the analysis of hydrologic signature representation are shown in Figure 11, which also shows that the hydrologic signatures best represented by the NWM were similarly those best represented by the LSTM_PPA. The same was true for the most poorly represented hydrologic signatures in both models.

The LSTM post-processors hurt the representation of the frequency of days with zero flow. There were

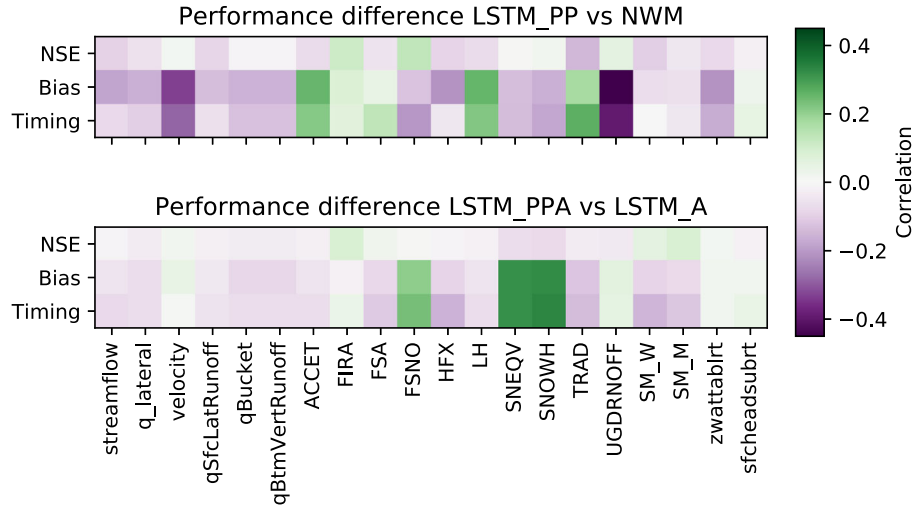


FIGURE 8. Correlations between the time-averaged NWM related inputs vs. performance metric differences between the LSTM post-processors (LSTM_PP and LSTM_PPA) and NWM and LSTM_A.

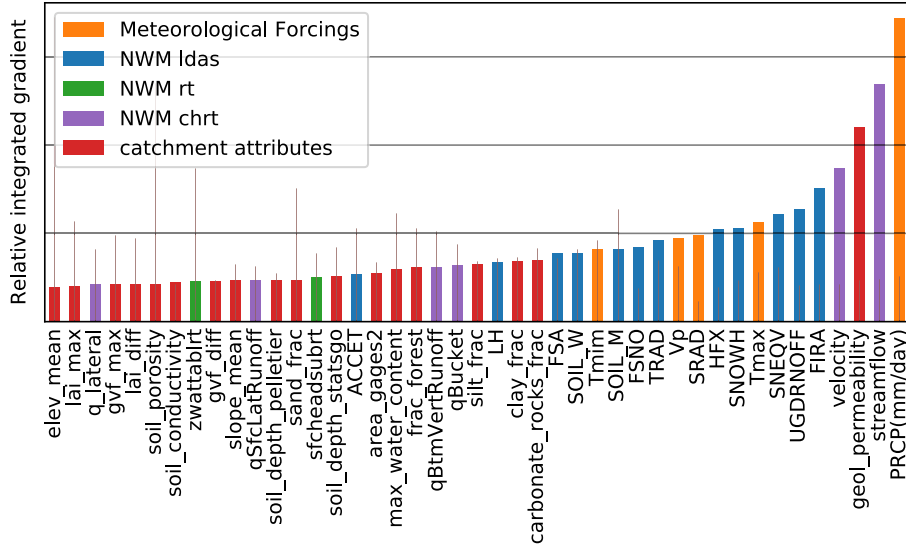


FIGURE 9. Attributions to the LSTM_PPA predictions. The vertical axis shows the relative magnitude of attribution (importance) for each input, with precipitation (PRCP) as the top contributor and NWM-predicted runoff into channel reach (q_lateral) contributing the least.

101 basins with any periods of zero flow. None of these models do well simulating zero flow, but the NWM is better at handling this situation, predicting zero flow periods in 56 of the 101 basins. The LSTM_A, LSTM_PPA, and LSTM_PP only predicted periods of zero flows at 35, 29, and 25 basins, respectively. This is an important characteristic in basins in the Southwest, where the NWM could use the benefit of a LSTM post-processor, so this would be a good place to focus future research of theory-guided ML for hydrology.

The LSTM post-processor made a significant improvement over the NWM for several signatures. The improvement to runoff ratio, which is the

fraction of precipitation that makes it through the stream gauge at the surface, could be a compensation for the uncalibrated soil parameters in the NWM mentioned by Salas et al. (2018). The LSTM post-processor improved both high and low flow predictions (5% and 95% flow quantiles), which are important for natural resources management. The mean daily discharge was the best represented hydrologic signature by all models.

The LSTM_PPA post-processor made significant improvements over the LSTM for baseflow index. This is the only sign that an LSTM post-processor improved over both the NWM and the LSTM_A. This signature estimates the contribution of baseflow to

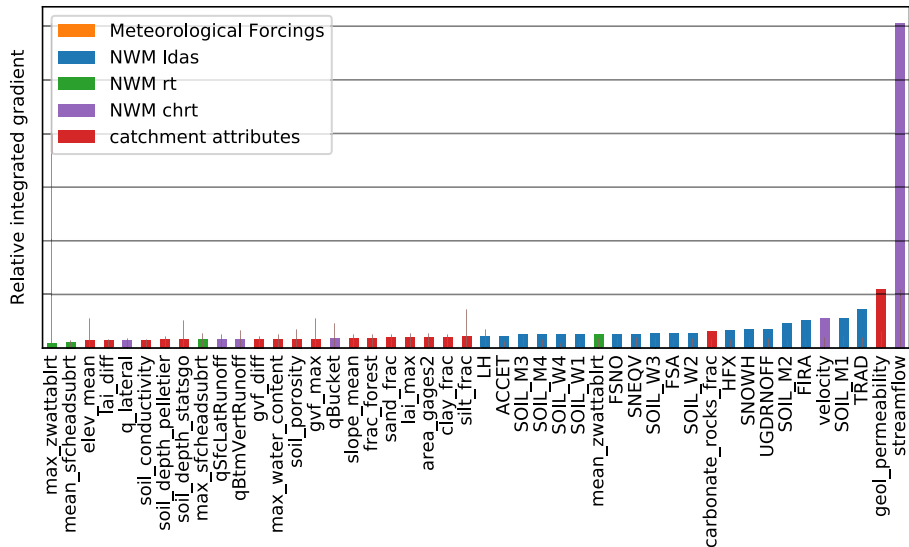


FIGURE 10. Attributions for the LSTM_PP model. Color coded by LSTM input source. The streamflow is overwhelmingly the highest contributor to the post-processed streamflow prediction.

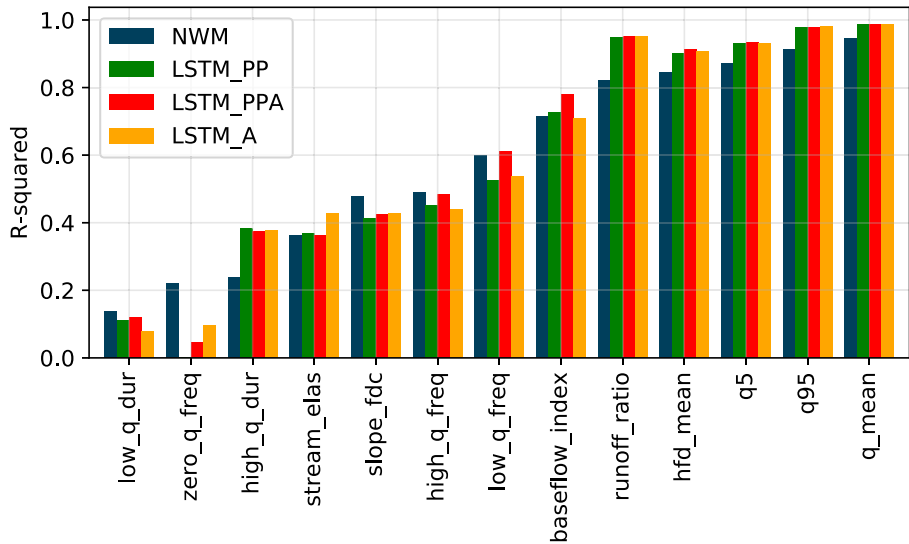


FIGURE 11. Correlation between simulated and observed per-basin hydrologic signatures from the NWM (blue), LSTM_A (orange), LSTM_PPA (green), and LSTM_PP (red). Larger values indicate better performance.

the total discharge, which is computed by hydrograph separation. Klemeš (1986) (summarizing Lindsly’s Applied Hydrology) cautioned strongly against using hydrograph separation, because there is no real basis for distinguishing the source of flow in a stream.

Results Comparing Gauged Basins vs. Ungauged Basins

Results in Table 7 summarize an analysis designed to replicate prediction in ungauged basins. The table has metrics from predictions by the NWM, LSTM_A

and the LSTM post-processors (LSTM_PP and LSTM_PPA) calculated only at basins that were either calibrated or uncalibrated, but not both. There was no statistical difference between the calibrated and uncalibrated samples. This indicates that the LSTM post-processor works in uncalibrated basins. When post-processors were trained only in calibrated basins (denoted with a “C” in Table 7), however, the performance in uncalibrated basins significantly deteriorated. But this is true for the LSTM_A as well, so it is not a result of the calibration (as calibration would not influence the LSTM_A), but a result of prediction at ungauged type basins. However, the

	Calibrated basins				Uncalibrated basins			
	Mean	Median	Max	Min	Mean	Median	Max	Min
NSE								
NWM	0.49	0.64	0.95	-10.81	0.18	0.48	0.79	-7.10
LSTM_PP	0.65	0.73	0.93	-3.32	0.69	0.71	0.89	0.38
LSTM_A	0.68	0.74	0.93	-0.64	0.73	0.75	0.89	0.43
LSTM_PPA	0.66	0.75	0.93	-3.61	0.71	0.73	0.89	0.42
LSTM_PP(C)	0.65	0.73	0.93	-1.86	0.21	0.57	0.75	-8.12
LSTM_A(C)	0.67	0.74	0.93	-1.13	0.51	0.67	0.84	-2.54
LSTM_PPA(C)	0.67	0.75	0.94	-2.71	0.13	0.58	0.84	-14.07
Total bias								
NWM	0.01	-0.01	2.57	-0.63	0.00	-0.06	1.84	-0.58
LSTM_PP	0.04	0.02	1.05	-0.24	0.02	0.01	0.27	-0.12
LSTM_A	0.02	0.02	0.56	-0.22	0.02	0.01	0.20	-0.11
LSTM_PPA	0.03	0.02	0.98	-0.21	0.01	0.00	0.22	-0.11
LSTM_PP(C)	0.01	-0.01	0.92	-0.25	0.06	-0.04	2.15	-0.51
LSTM_A(C)	0.02	0.02	0.62	-0.21	0.09	0.04	0.99	-0.20
LSTM_PPA(C)	0.01	0.00	0.95	-0.22	0.06	-0.05	2.89	-0.41
Peak timing error								
NWM	1.06	0.91	3.00	0.10	1.04	0.77	2.70	0.25
LSTM_PP	0.55	0.45	1.95	0.04	0.52	0.35	1.59	0.04
LSTM_A	0.53	0.43	1.76	0.00	0.51	0.41	1.50	0.04
LSTM_PPA	0.54	0.42	1.75	0.04	0.51	0.36	1.45	0.05
LSTM_PP(C)	0.55	0.45	2.10	0.00	0.59	0.41	1.76	0.09
LSTM_A(C)	0.52	0.43	1.77	0.00	0.57	0.50	1.50	0.08
LSTM_PPA(C)	0.54	0.41	1.83	0.04	0.57	0.41	1.65	0.13

TABLE 7. Performance of the LSTM and the LSTM post-processor split between basins where the NWM was calibrated vs. uncalibrated. The “C” in the model name denotes that the model training set only included calibrated basins.

median performance of the post-processor predictions at ungauged type basins when trained at only calibrated basins was still significantly better than the NWM in the uncalibrated basins.

The NWM, LSTM_A, and the LSTM_PPA had higher NSE scores in calibrated basins than the uncalibrated basins. Note that these results are from the LSTMs (with and without NWM model outputs) trained on only basins where the NWM was calibrated. In the case of the LSTM post-processors the mean NSE scores in uncalibrated basins were very low for NSE. This is a result of two outlier basins (1466500, MCDONALDS BRANCH, Lat: 39.9, Lon: -74.5, Area: 5.7 km; and 01484100 BEAVERDAM BRANCH, Lat: 38.9, Lon: -75.5, Area: 7.8 km). Both of those outlier basins are much smaller, and have lower flows, than the average of the training set. Without these basins the mean NSE scores were 0.32, 0.51, 0.56 and 0.56 for the NWM, LSTM_PP, LSTM_A, and LSTM_PPA, respectively. Table 7 also shows that the median value of the LSTM_PPA was higher than the NWM, as was the maximum NSE value, but the minimum value was exceptionally low.

The total bias in calibrated basins was generally better (lower) than the uncalibrated basins. The timing error of the NWM was actually better in the uncalibrated basins, but the LSTM_A and LSTM post-processors had better performance in the calibrated basins. The NSE values for the NWM,

LSTM_A, and the LSTM post-processors (LSTM_PP and LSTM_PPA) were significantly different in the calibrated basins vs. the uncalibrated basins, as were the differences between the LSTM_A and LSTM post-processors (LSTM_PP and LSTM_PPA) compared to the NWM. The bias values were significantly different between the two samples (calibrated vs. uncalibrated), but the differences between LSTM_A and LSTM post-processors vs. the NWM were not statistically different. This means that the LSTM models were successful at predicting streamflow at basins outside of the calibration set.

LSTM Post-Processor Sensitivity to Inputs and Application for Process Representation Diagnostics

Figure 12 shows results from the LSTM models with inputs from different parts of the NWM (land surface model only, channel router only, predicted streamflow only, and all states and fluxes). The best performing LSTM models (LSTM_A and LSTM_PPA) were the ones trained with inputs that included the five atmospheric forcing variables with (LSTM_PPA) and without (LSTM_A) the NWM output (these are the same models discussed in previous sections above). This implies that LSTM in general was able to extract more information from the atmospheric forcings than the NWM. Each of the LSTM post-

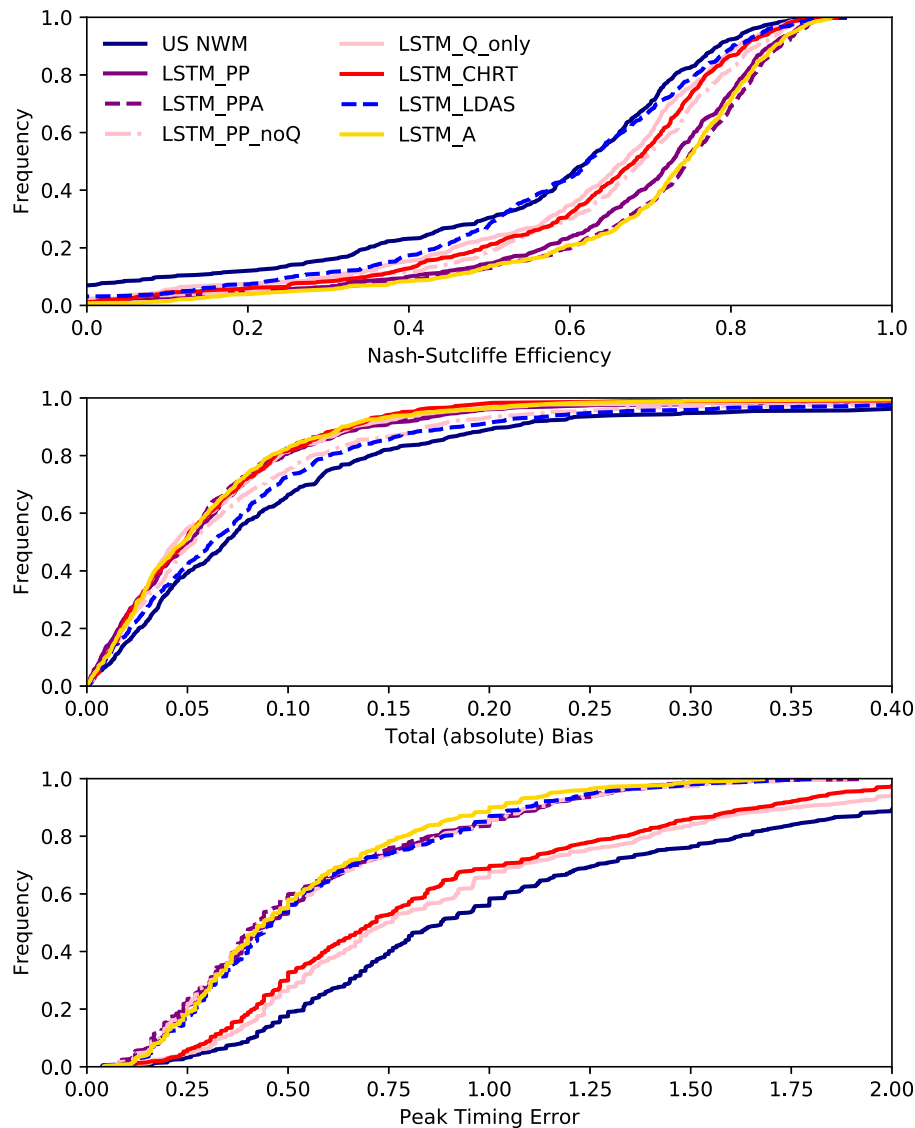


FIGURE 12. Performance of the LSTM post-processor trained with different sets of NWM output. Each of these post-processors outperform the NWM. LSTM_A is the LSTM trained with atmospheric forcings as dynamic inputs. LSTM_PP is the NWM post-processor trained with the outputs of the NWM as dynamic inputs. LSTM_PPA used both the NWM outputs and atmospheric forcings as inputs. LSTM_PP_noQ used all the NWM outputs except for streamflow and velocity from the channel router. LSTM_Q_only used only streamflow from the NWM output. LSTM_chrt used only the NWM channel router outputs. LSTM_ldas used only the land surface fluxes as inputs.

processors made better average daily streamflow predictions than the NWM itself, indicating that information from the atmospheric forcings is lost in the NWM model structure before the streamflow prediction is made. For example, the LSTM that took as inputs only the LDAS model output from the NWM made better predictions than the NWM itself, indicating that there is more information in the LDAS states and fluxes than the NWM is able to translate into streamflow predictions. The same was true for the states and fluxes of the CHRT component of the NWM, meaning that information is also lost in the CHRT component of the NWM model structure.

DISCUSSION

Comparison between the LSTM_A and the Post-Processors (LSTM_PP and LSTM_PPA)

The LSTM_A, trained only on atmospheric forcings as dynamic inputs, was better at extrapolating hydrologic conditions outside the training set than the LSTM post-processors (LSTM_PP and LSTM_PPA), and thus LSTM_A is the better performing model. This is shown in the analysis of prediction in ungauged basins, specifically Table 7. The post-processors both

failed to make reasonable predictions at two basins that were much smaller than any basins included in the training set. The LSTM_A was able to make good predictions in these basins. Including the NWM output as dynamic inputs to the LSTM constrained the model and prevented it from learning general hydrologic relationships that can be extracted to basins with characteristics that might be unrecognizable.

Potential for Improving the Performance of Both the NWM and ML

Results presented here show that the LSTM post-processors are unreliable for improving predictions of the NWM. The LSTM post-processors did provide significant benefit to the NWM streamflow predictions at almost all (88% and 92% for LSTM_PP and LSTM_PPA, respectively) of the 531 basins analyzed here, but was severely detrimental to two basins in our tests of ungauged basins. In the basins where this was not the case, it may be possible to use fine tuning a version of the post-processor that is specific to each gauge location (as would be done in traditional model calibration); however, the LSTM_A did not have this problem and is more reliable. We trained the LSTMs on headwater basins, so further work would be needed to include reservoirs, urban areas, and other management practices. It is worth noting that these LSTM models can be trained on a laptop computer in a few hours, a relatively minor computational cost, and the computational cost of forward prediction is negligible. By comparison the computational cost of calibrating the NWM is much higher — typically requiring HPC or cloud systems.

The NWM performance and the performance improvement from the LSTM post-processors (LSTM_PP and LSTM_PPA) were negatively correlated: basins with a low performance by the NWM have the highest performance change from the LSTM post-processors. This means that post-processing can be expected to correct situations where the NWM gives bad predictions. Conversely, the performance of the NWM and the LSTM_A (the LSTM trained without NWM model outputs) were minimally correlated ($r^2 = 0.42, 0.30, \text{ and } 0.67$ for NSE, bias, and timing, respectively). Considering also that the overall performance of the LSTM_A changed only minimally from the addition of the NWM inputs (as shown in Figures 3–5; Table 6) and that the LSTM_PPA still preferred to extract more information from precipitation forcings (shown in Figure 9), we might conclude that the LSTM post-processors learned new patterns of the rainfall–runoff response, which are not fully represented by the NWM. But this relationship is also

learned by LSTM_A, without the influence of the NWM. The overall improvement in the representation of hydrologic signatures indicates the post-processor may be a better representation of physical flow patterns than either the NWM or the LSTM_A, though not significantly. The interpretation of the integrated gradient (Figures 9 and 10) and the correlations between improvement and NWM features (Figure 8) indicate that this improvement of flow patterns comes from information in the NWM representation of streamflow and snow states.

Application to Real-Time Forecasting

The NWM is not simply a rainfall–runoff simulator; it simulates flow through 2.7 million river reaches around CONUS, dam operations, land surface processes, hydraulics, and other complications of large domain hydrology. The nature of the CAMELS catchments selected in these experiments is such that they have few engineered control structures and are under 20,000 km². The results presented in this paper show that the LSTMs improved streamflow predictions in the catchments studied here, which all had limited human disturbance (e.g., dams, reservoirs, etc.). Kratzert, Klotz, Herrnegger, et al. (2019) showed that LSTM_A predictions extend into ungauged basins, and this is consistent with our results. Our results (section “*Results comparing calibrated basins vs. uncalibrated basins*”) show that the LSTM_A is a much better choice than the post-processors in ungauged basins, which is the majority of the NWM domain. The immediate potential for improving real-time forecasting could be deploying an LSTM_A for streamflow prediction in undisturbed catchments, and undisturbed subcatchments upstream of unnatural hydrologic conditions such as dams, agriculture lands, and urban centers. This would allow for retaining conceptual representations of lakes and reservoirs that already exist in the NWM.

Diagnosing Process-Based Models, Physical Processes, and Data Concerns

The sensitivity analysis reported in Figure 12 showed that some components of the NWM caused poor predictions. Specifically, information was lost in channel router (CHRT) component of the model. This diagnostic method could be used to compare different schemes for future versions of the NWM. For instance, changing the routing function might conserve timing information from the land surface fluxes, or modifying the evapotranspiration options in Noah-MP may conserve mass bias information from

the NWM forcing engine. Such improvements could be quantified with this post-processing method.

Each of the post-processing models tested for sensitivity (Figure 12) fall, roughly and inclusively, between the NWM and the LSTM_A. Based on the relative positions between those bounding curves, we can identify sources of information loss through the NWM modeling chain:

1. The channel routing outputs contain more information of simulation bias than timing, meaning the channel router moves with poor timing, but conserves mass well.
2. The land surface outputs contain more information of simulation timing than bias, meaning the land surface component does not conserve mass well, but delivers water to the channel at appropriate times.
3. Information is lost during channel routing after the mass is delivered, indicating the channel router is not functioning properly.

There is potential to expand this analysis, breaking down the NWM components even further. Quantification can be done with the full state space from the NWM. Retrospective runs using new versions of the NWM should output the full state space for these types of analysis. This diagnostics analysis using ML post-processing is possible with any physics-based, conceptual or process-based dynamics model.

Moving Forward with Theory-Guided ML

The post-processing procedure presented here is one of the cruder techniques currently available for combining process-based and data-driven models. Several other methods of combining the benefits of ML (predictability) with the benefits of physically realistic hydrologic theory (robustness) are in development. For example, Pelissier et al. (2019) integrated a trained Gaussian Processes into the state-space dynamics of a process-based land surface model for predicting soil moisture time series. Another example is using physical principles to constrain the loss function of an ML model during training — for example, Hoedt et al. (2021) integrated mass balance constraints into an LSTM and applied this model to the same 531 basins used in this study. Implementing post-processing is relatively straightforward compared to other techniques such as adding physics into ML code or using ML to dynamically update the state variables, but is unreliable when the process-based models used as input is uncalibrated.

Using ML for post-processing has the potential for advancing the explainability of data-driven models.

We showed that the LSTM model representation of hydrologic signatures (with and without NWM model outputs) is highly correlated with the NWM. This indicates that the “learned” functions mapping inputs to streamflow are actually quite similar. We might have trouble expressing the “learned” LSTM with compact formulas (e.g., PDEs), given the high number of trained model weights, but we can use them with confidence knowing their structural similarities with process-based models like the NWM.

CONCLUSION

The LSTM post-processors (LSTM_PPA and LSTM_PP) significantly outperformed the NWM, but did not consistently, nor significantly, outperformed the LSTM_A (the LSTM model trained without the NWM model outputs as LSTM inputs). LSTMs, in general, are capable of learning the dynamics of rainfall-runoff processes, gaining little additional information from the conceptualizations coded within the NWM. The “pure” post-processing model (LSTM_PP) outperformed the NWM in terms of bias, and significantly outperformed the NWM in terms of NSE and timing. A decision to use the LSTM as a post-processor for the NWM should be made with professional judgment, considering the comparison of the NWM, LSTM, and LSTM post-processor’s performance. In locations where the NWM is not calibrated, or the hydrologic conditions are not well understood, it would be best to use the LSTM without the influence from the NWM.

The results indicate that there is more information in the atmospheric forcings about streamflow observations than in the NWM outputs, including the NWM streamflow prediction. The NWM loses information between the atmospheric forcing inputs and the outputs. The NWM land surface component (LDAS) loses information about mass conservation (shown from the bias error), and the channel router (CHRT) loses information about streamflow timing. The NWM routing scheme should be considered as a priority for improving the NWM.

DATA AVAILABILITY

All data and code used in this paper are publicly available in the following locations: U.S. National Water Model: <https://docs.opendata.aws/nwm-archive/readme.html>. CAMELS data: <https://ral.ucar.edu/solutions/products/camels>. Data processing code: <https://github.com/jmframe/nwm-reanalysis-model-data-processing>; <https://doi.org/10.5281/zenodo.4642605>.

LSTM code: https://github.com/kratzert/ealstm_regional_modeling. Post-processing and analysis code: <https://github.com/jmframe/nwm-post-processing-with-lstm>; <https://doi.org/10.5281/zenodo.4642603>.

ACKNOWLEDGMENTS

Frederik Kratzert from Johannes Kepler University was partially supported by a Google Faculty Research Award. Jonathan Frame from the University of Alabama was supported by the NASA Terrestrial Hydrology Program. Grey Nearing (then at the University of Alabama) was partially supported by the NCAR COMET Program on a cooperative award with the National Water Center. The LSTM models presented here were trained using computational resources from the NASA Center for Climate Simulation. We thank the reviewers for their valuable feedback and constructive criticism throughout the review of this paper.

AUTHOR CONTRIBUTIONS

Jonathan M. Frame: Conceptualization; Data curation; Formal analysis; Investigation; Methodology; Project administration; Resources; Software; Supervision; Validation; Visualization; Writing-original draft; Writing-review & editing. **Frederik Kratzert:** Data curation; Software. **Austin Raney:** Data curation; Methodology; Software. **Mashrekur Rahman:** Writing-review & editing. **Fernando Salas:** Conceptualization. **Grey S. Nearing:** Conceptualization; Data curation; Formal analysis; Investigation; Methodology; Project administration; Resources; Software; Supervision; Validation; Writing-review & editing.

LITERATURE CITED

Addor, N., G. Nearing, C. Prieto, A.J. Newman, N. Le Vine, and M.P. Clark. 2018. "A Ranking of Hydrological Signatures Based on Their Predictability in Space." *Water Resources Research* 54 (11): 8792–812. <https://doi.org/10.1029/2018WR022606>.

Addor, N., A.J. Newman, N. Mizukami, and M.P. Clark. 2017. "The CAMELS Data Set: Catchment Attributes and Meteorology for Large-Sample Studies." *Hydrology and Earth System Sciences* 21: 5293–313. <https://doi.org/10.5194/hess-21-5293-2017>.

Chadalawada, J., H.M.V.V. Herath, and V. Babovic. 2020. "Hydrologically Informed Machine Learning for Rainfall-Runoff Modeling: A Genetic Programming-Based Toolkit for Automatic Model Induction." *Water Resources Research* 56 (4): 1–23. <https://doi.org/10.1029/2019WR026933>.

Cosgrove, B., D. Gochis, E.P. Clark, Z. Cui, A.L. Dugger, G.M. Fall, X. Feng *et al.* 2015. "Hydrologic Modeling at the National Water Center: Operational Implementation of the WRF-Hydro Model to Support National Weather Service Hydrology." In AGU Fall Meeting Abstracts.

Daw, A., R.Q. Thomas, C.C. Carey, J.S. Read, A.P. Appling, and A. Karpatne. 2020. "Physics-Guided Architecture (PGA) of Neural Networks for Quantifying Uncertainty in Lake Temperature Modeling." Proceedings of the 2020 SIAM International Conference on Data Mining, 532–40. <https://doi.org/10.1137/1.9781611976236.60>.

Elmer, N.J. 2019. "Using Satellite Observations of River Height and Vegetation to Improve National Water Model Initialization and Streamflow Prediction." PhD diss., The University of Alabama in Huntsville.

Gauch, M., F. Kratzert, D. Klotz, G. Nearing, J. Lin, and S. Hochreiter. 2021. "Rainfall-Runoff Prediction at Multiple Timescales with a Single Long Short-Term Memory Network." *Hydrology and Earth System Sciences* 25 (4): 2045–62. <https://doi.org/10.5194/hess-25-2045-2021>.

Gauch, M., J. Mai, and J. Lin. 2021. "The Proper Care and Feeding of CAMELS: How Limited Training Data Affects Streamflow Prediction." *Environmental Modelling & Software* 135: 104926. <https://doi.org/10.1016/j.envsoft.2020.104926>.

Gupta, H.V., H. Kling, K.K. Yilmaz, and G.F. Martinez. 2009. "Decomposition of the Mean Squared Error and NSE Performance Criteria: Implications for Improving Hydrological Modeling." *Journal of Hydrology* 377 (1–2): 80–91. <https://doi.org/10.1016/j.jhydrol.2009.08.003>.

Gupta, H.V., T. Wagener, and Y. Liu. 2008. "Reconciling Theory with Observations: Elements of a Diagnostic Approach to Model Evaluation." *Hydrological Processes* 22 (18): 3802–13. <https://doi.org/10.1002/hyp.6989>.

Han, H. 2021. "Improving Hydrologic Modeling of Runoff Processes Using Data-Driven Models." Doctoral diss., Colorado State University. https://mountainscholar.org/bitstream/handle/10217/232583/Han_colostate_0053A_16465.pdf?sequence=1&isAllowed=y.

Hansen, C., J.S. Shiva, S. McDonald, and A. Nabors. 2019. "Assessing Retrospective National Water Model Streamflow with Respect to Droughts and Low Flows in the Colorado River Basin." *Journal of the American Water Resources Association* 55 (4): 964–75. <https://doi.org/10.1111/1752-1688.12784>.

Hochreiter, S. 1991. "Untersuchungen Zu Dynamischen Neuronalen Netzen." Doctoral diss., Institut Für Informatik, Technische Universität, München. <http://people.idsia.ch/~juergen/SeppHochreiter1991ThesisAdvisorSchmidhuber.pdf>.

Hochreiter, S., and J. Schmidhuber. 1997. "Long Short-Term Memory." *Neural Computation* 9 (8): 1735–80. <https://doi.org/10.1162/neco.1997.9.8.1735>.

Hoedt, P.J., F. Kratzert, D. Klotz, C. Halmich, M. Holzleitner, G. Nearing, S. Hochreiter & G. Klambauer. 2021. "MC-LSTM: Mass-Conserving LSTM." Proceedings of the 38th International Conference on Machine Learning, PMLR 139. <http://arxiv.org/abs/2101.05186>.

Karpatne, A., G. Atluri, J.H. Faghmous, M. Steinbach, A. Banerjee, A. Ganguly, S. Shekhar, N. Samatova, and V. Kumar. 2017. "Theory-Guided Data Science: A New Paradigm for Scientific Discovery from Data." *IEEE Transactions on Knowledge and Data Engineering* 29 (10): 2318–31. <https://doi.org/10.1109/TKDE.2017.2720168>.

Karpatne, A., W. Watkins, J. Read, and V. Kumar. 2017. "Physics-Guided Neural Networks (PGNN): An Application in Lake Temperature Modeling." <http://arxiv.org/abs/1710.11431>.

Kim, J., L. Read, L.E. Johnson, D. Gochis, R. Cifelli, and H. Han. 2020. "An Experiment on Reservoir Representation Schemes to Improve Hydrologic Prediction: Coupling the National Water Model with the HEC-ResSim." *Hydrological Sciences Journal* 65 (10): 1652–66. <https://doi.org/10.1080/02626667.2020.1757677>.

Klemeš, V. 1986. "Dilettantism in Hydrology: Transition or Destiny?" *Water Resources Research* 22 (9S): 177S–88. <https://doi.org/10.1029/WR022i09Sp0177S>.

- Kratzert, F., D. Klotz, C. Brenner, K. Schulz, and M. Herrnegger. 2018. "Rainfall-Runoff Modelling Using Long Short-Term Memory (LSTM) Networks." *Hydrology and Earth System Sciences* 22 (11): 6005–22. <https://doi.org/10.5194/hess-22-6005-2018>.
- Kratzert, F., D. Klotz, M. Herrnegger, A.K. Sampson, S. Hochreiter, and G.S. Nearing. 2019. "Towards Improved Predictions in Ungauged Basins: Exploiting the Power of Machine Learning." *Water Resources Research*: 2019WR026065. <https://doi.org/10.1029/2019WR026065>.
- Kratzert, F., D. Klotz, G. Shalev, G. Klambauer, S. Hochreiter, and G.S. Nearing. 2019. "Towards Learning Universal, Regional, and Local Hydrological Behaviors via Machine Learning Applied to Large-Sample Datasets." *Hydrology and Earth System Sciences* 23 (12): 5089–110. <https://doi.org/10.5194/hess-23-5089-2019>.
- Nearing, G.S., F. Kratzert, D. Klotz, P.J. Hoedt, G. Klambauer, S. Hochreiter, H. Gupta, S. Nevo, and Y. Matias. 2020. "A Deep Learning Architecture for Conservative Dynamical Systems: Application to Rainfall-Runoff Modeling." AI for Earth Sciences Workshop at NEURIPS 2020.
- Nearing, G.S., F. Kratzert, A.K. Sampson, C.S. Pelissier, D. Klotz, J.M. Frame, C. Prieto, and H.V. Gupta. 2021. "What Role Does Hydrological Science Play in the Age of Machine Learning?" *Water Resources Research* 57: e2020WR028091. <https://doi.org/10.1029/2020WR028091>.
- Nearing, G.S., B.L. Ruddell, M.P. Clark, B. Nijssen, and C. Peters-Lidard. 2018. "Benchmarking and Process Diagnostics of Land Models." *Journal of Hydrometeorology* 19 (11): 1835–52. <https://doi.org/10.1175/JHM-D-17-0209.1>.
- Newman, A.J., M.P. Clark, K. Sampson, A. Wood, L.E. Hay, A. Bock, R.J. Viger *et al.* 2015. "Development of a Large-Sample Watershed-Scale Hydrometeorological Data Set for the Contiguous USA: Data Set Characteristics and Assessment of Regional Variability in Hydrologic Model Performance." *Hydrology and Earth System Sciences* 19 (1): 209–23. <https://doi.org/10.5194/hess-19-209-2015>.
- Niu, G.-Y., Z.-L. Yang, K.E. Mitchell, F. Chen, M.B. Ek, M. Barlage, A. Kumar *et al.* 2011. "The Community Noah Land Surface Model with Multiparameterization Options (Noah-MP): 1. Model Description and Evaluation with Local-Scale Measurements." *Journal of Geophysical Research* 116: D12109. <https://doi.org/10.1029/2010JD015139>.
- NOAA (National Oceanic and Atmospheric Administration). 2019. "Stakeholder Engagement to Inform National Weather Service Hydrologic Products and Services to Meet User Needs NOAA National Weather Service Water Resources Services Branch and Office of Water Prediction." <https://www.weather.gov/media/water/StakeholderEngagementtoInformNWMProductsServices2019.pdf>.
- Pelissier, C., J. Frame, and G. Nearing. 2019. "Combining Parametric Land Surface Models with Machine Learning." arXiv preprint arXiv:2002.06141.
- Read, J.S., X. Jia, J. Willard, A.P. Appling, J.A. Zwart, S.K. Oliver, A. Karpatne *et al.* 2019. "Process-Guided Deep Learning Predictions of Lake Water Temperature." *Water Resources Research* 55: 9173–90. <https://doi.org/10.1029/2019WR024922>.
- Reichstein, M., G. Camps-Valls, B. Stevens, M. Jung, J. Denzler, and N. Carvalhais. 2019. "Deep Learning and Process Understanding for Data-Driven Earth System Science." *Nature* 566: 195–204. <https://doi.org/10.1038/s41586-019-0912-1>.
- Salas, F.R., M.A. Somos-Valenzuela, A. Dugger, D.R. Maidment, D.J. Gochis, C.H. David, W. Yu, D. Ding, E.P. Clark, and N. Noman. 2018. "Towards Real-Time Continental Scale Streamflow Simulation in Continuous and Discrete Space." *Journal of the American Water Resources Association* 54 (1): 7–27. <https://doi.org/10.1111/1752-1688.12586>.
- Steiger, J., and M. Browne. 1984. "The Comparison of Interdependent Correlations between Optimal Linear Composites." *Psychometrika* 49 (1): 11–24. <https://doi.org/10.1017/CBO9781107415324.004>.
- Sundararajan, M., A. Taly, and Q. Yan. 2017. "Axiomatic Attribution for Deep Networks." 34th International Conference on Machine Learning, ICML 2017 7: 5109–18.
- Tartakovsky, A.M., C.O. Marrero, P. Perdikaris, G.D. Tartakovsky, and D. Barajas-Solano. 2020. "Physics-Informed Deep Neural Networks for Learning Parameters and Constitutive Relationships in Subsurface Flow Problems." *Water Resources Research* 56 (5): 1–16. <https://doi.org/10.1029/2019WR026731>.
- Xia, Y., K. Mitchell, M. Ek, J. Sheffield, B. Cosgrove, E. Wood, L. Luo *et al.* 2012. "Continental-Scale Water and Energy Flux Analysis and Validation for the North American Land Data Assimilation System Project Phase 2 (NLDAS-2): 1. Intercomparison and Application of Model Products." *Journal of Geophysical Research: Atmospheres* 117 (D3). <https://doi.org/10.1029/2011JD016048>.
- Ye, A., Q. Duan, X. Yuan, E.F. Wood, and J. Schaake. 2014. "Hydrologic Post-Processing of MOPEX Streamflow Simulations." *Journal of Hydrology* 508: 147–56. <https://doi.org/10.1016/j.jhydrol.2013.10.055>.