

Technical note: Using Long Short-term Memory Models to Fill Data Gaps in Hydrological Monitoring Networks

Huiying Ren¹, Erol Cromwell², Ben Kravitz^{3,4}, and Xingyuan Chen⁴

¹Earth Systems Science Division, Pacific Northwest National Laboratory, WA, USA

²Advanced Computing, Mathematics, and Data Division, Pacific Northwest National Laboratory, WA, USA

³Department of Earth and Atmospheric Sciences, Indiana University, Bloomington, IN, USA

⁴Atmospheric Sciences and Global Change Division, Pacific Northwest National Laboratory, WA, USA

Correspondence: Xingyuan Chen (Xingyuan.Chen@pnnl.gov)

Abstract. The spatio-temporal dynamics in subsurface hydrological flows over a long time window are usually quantified through a network of monitoring wells; however, such observations often are spatially sparse and temporal gaps exist due to poor quality or instrument failure. In this study, we explore the ability of recurrent neural networks to fill gaps in a spatially distributed time-series dataset from a well network that monitors the dynamic and heterogeneous hydrologic exchanges between the Columbia River and its adjacent groundwater aquifer at the U.S. Department of Energy's Hanford site. This 10-year-long dataset contains hourly temperature, specific conductance, and groundwater table elevation measurements from 42 wells with various lengths of gaps. We employ a long short-term memory (LSTM) model to capture the temporal variations in the observed system behaviors for gap filling. The performance of the LSTM-based gap filling method was evaluated against a traditional autoregressive integrated moving average (ARIMA) method in terms of both the error statistics and how well they capture the temporal patterns in river corridor wells that exhibit various dynamics signatures. Our study demonstrates that the ARIMA models yield better average error statistics, yet they tend to have larger errors during time windows with abrupt changes or high-frequency (daily and subdaily) variations. The LSTM-based models are found to excel in capturing both the high-frequency and low-frequency (monthly and seasonal) dynamics, although the inclusion of high-frequency fluctuations may also lead to overly dynamic predictions in time windows that lacks such fluctuations. The LSTM is able to take advantage of the spatial information from neighboring wells to improve the gap filling accuracy, especially for long gaps in system states that vary at subdaily scales. Despite the fact that LSTM models require substantial training data and have limited extrapolation power beyond the conditions represented in the training data, they afford the great ~~flexibity~~flexibility to account for the spatial and temporal correlations and nonlinearity in data without ~~a-priori~~a-priori assumptions. Thus, LSTMs provide effective alternatives to fill in data gaps in spatially distributed time-series observations characterized by multiple dominant frequencies of variability, which are essential for advancing our understanding of dynamic complex systems.

Copyright statement. TEXT

1 Introduction

Long-term hydrological monitoring using distributed well networks is of critical importance for understanding how ecosystems respond to chronic or extreme perturbations, as well as for informing policies and decisions related to natural resources and environmental issues (Wett et al., 2002; Taylor and Alley, 2002; Grant and Dietrich, 2017). One of the most common methods to collect hydrologic and chemistry data in groundwater is through wells (Güler and Thyne, 2004; Strobl and Robillard, 2008; Lin et al., 2012); however, most well data have temporal gaps due to instrument failure or poor quality of measurements for numerous reasons. These data gaps degrade the quality of the dataset and increase the uncertainty in the spatial and temporal patterns that are derived from them. While gap filling is essential for developing understanding of dynamic system behaviours and for use in creating continuous, internally consistent boundary conditions for numerical models, one outstanding challenge is to capture the nonstationarity in data.

Various statistical methods have been developed to fill gaps in spatio-temporal datasets, with the most commonly used being the autoregressive integrated moving average (ARIMA) method (Han et al., 2010; Zhang, 2003). For any given spatial location, ARIMA uses temporal autocorrelations to predict unobserved data points in a time series. Spatio-temporal autocorrelations can be considered by using multivariate ARIMA and space-time autoregressive models (Kamarianakis and Prastacos, 2003; Wikle et al., 1998; Kamarianakis and Prastacos, 2005); however, ARIMA cannot capture nonlinear trends because it assumes a linear dependence between adjacent observations (Faruk, 2010; Valenzuela et al., 2008; Ho et al., 2002). In addition, all existing space-time ARIMA models assume fixed global autoregressive and moving average terms, which would fail to capture evolving dynamics in highly dynamic systems (Pfeifer and Deutch, 1980; Griffith, 2010; Cheng et al., 2012, 2014). Spectral-based methods, such as singular spectrum analysis, maximum entropy method, and Lomb-Scargle periodogram, have been used to account for nonlinear trends while filling in gaps in spatio-temporal datasets (Ghil et al., 2002; Hocke and Kämpfer, 2008; Kondrashov and Ghil, 2006). However, these methods use a few optimal spatial or temporal modes occurring at low frequencies to predict the missing values, with the other higher frequency components discarded as noise, which may lead to reduced accuracy of the statistical models in fitting the observations and in predicting missing values (Kondrashov et al., 2010; Wang et al., 2012). Kriging and maximum likelihood estimation used in spatial and spatio-temporal gap filling often face computational challenges in computing the covariance matrix of the data vector, which can be quite large (Katzfuss and Cressie, 2012; Eidsvik et al., 2014). Other nonlinear methods have been explored with some success, including expectation-maximization or Bayesian probabilistic inference including hierarchical models, Gaussian process, and Markov chain Monte Carlo; the spatial and temporal correlations are most effectively captured by using models that build dependencies in different stages or hierarchies (Calculi et al., 2015; Banerjee et al., 2014; Datta et al., 2016; Finley et al., 2013; Stroud et al., 2017). In general, the expectation-maximization algorithm and Bayesian-based methods are sensitive to the choice of initial values and prior distributions in parameter space (Katzfuss and Cressie, 2011, 2012). Moreover, the prior distributions with all the associated parameters in both the spatial and temporal domains need to be specified, which becomes increasingly difficult in more complex systems. Empirical Orthogonal Functions (EOF) related interpolation methods, such as least squares EOF (LSEOF), data interpolation EOF (DINEOF), and recursively subtracted EOF (REEOF), are widely used to fill in missing data

from geophysical fields such as clouds in sea surface temperature datasets or other satellite-based images with regular gridded domains (Beckers and Rixen, 2003; Beckers et al., 2006; Alvera-Azcárate et al., 2016). However, the requirement of gridded data by the EOF methods limits their use in filling data gaps in irregularly spaced monitoring networks.

Deep neural networks (DNNs) (Schmidhuber, 2015) are data-driven tools that, in principle, could provide a powerful way of extracting the nonlinear spatio-temporal patterns hidden in the distributed time-series data without knowing their explicit forms (Långkvist et al., 2014). They are increasingly being used in geoscience domains to extract patterns and insights from the streams of geospatial data and to transform the understanding of complex systems (Reichstein et al., 2019; Shen, 2018; Sun, 2018; Sun et al., 2019; Gentile et al., 2018). The umbrella term of DNN contains numerous categories of architectures, depending on the problem at hand. For the analyses in this paper, which are focused on filling gaps in time-series data, a natural choice of the architecture is recurrent neural networks (RNNs) (JORDAN, 1986). RNNs take sequences (e.g., time series) as input and output single values or sequences that follow. They are designed to use information about previous events to make predictions about future events, essentially by letting the model “remember.” However, for longer sequences of data, RNNs have been shown to lose memory from previously trained data, i.e., they “forget” (Hochreiter et al., 2001). The earlier information becomes exponentially less impactful for the prediction as the size of the sequence increases. Long short-term memory (LSTM) networks are variations of RNNs that are explicitly designed to avoid this problem by using memory cells to retain information about relevant past events (Hochreiter and Schmidhuber, 1997). RNNs and LSTMs have been successfully applied to text prediction (Graves, 2013), text translation (Wu et al., 2016), speech recognition (Graves et al., 2013), and image captioning (You et al., 2016). Recently, the applications of RNNs and LSTMs are also emerging in hydrology. For example, Kratzert et al. (2018) used LSTMs to predict watershed runoff from meteorological observations, Zhang et al. (2018) used LSTMs for predicting sewer overflow events from rainfall intensity and sewer water level measurements, and Fang et al. (2017) used LSTMs to predict soil moisture with high fidelity.

Our study aims to evaluate the potential of using the LSTM models for filling gaps in spatio-temporal time series data collected from a distributed network. The LSTM-based gap filling method is tested using datasets collected to understand the interactions between a regulated river and a contaminated groundwater aquifer. We treat the gap filling as a forecasting problem, i.e., we use the historical data as input to predict the missing values in the data gaps. The performance of the LSTM-based gap filling method is compared with traditional time series approaches (i.e., ARIMA) to identify situations in which LSTM models outperform or under-perform the ARIMA models.

2 Study Site and Data Description

A 10-year (2008–2018) spatio-temporal dataset was collected from a network of groundwater wells that monitor temperature (Water conductivity and temperature probe CS547A by Campbell Scientific), specific conductance (SpC) (Water conductivity and temperature probe CS547A by Campbell Scientific), and water-table elevation (stainless-steel pressure transducer CS451 by the Campbell Scientific) at the 300 Area of the U.S. Department of Energy Hanford site, located in southeastern Washington State. The groundwater well network was originally built to monitor the attenuation of legacy contaminants. The groundwater

aquifer at our study site is composed of two distinct geologic formations: a highly permeable formation (Hanford formation, consisting of coarse gravelly sand and sandy gravel) underlain by a much less permeable formation (the Ringold Formation, consisting of silt and fine sand). The dominant hydrogeologic features of the aquifer are defined by the interface between the Hanford and Ringold formations and the heterogeneity within the Hanford formation (Chen et al., 2012, 2013).

5 The intrusion of river water into the adjacent groundwater aquifer causes mixing of two water bodies with distinct geo-chemistry and stimulates biogeochemical reactions at the interface. The river water has lower SpC ($0.1\text{--}0.12\text{ mS/cm}$) than the groundwater (averaging $\sim 0.4\text{ mS/cm}$). Groundwater has a nearly constant temperature ($16\text{--}17^\circ\text{C}$) as opposed to seasonally varying river temperature ($3\text{--}22^\circ\text{C}$). The highly heterogeneous coarse-textured aquifer (Zachara et al., 2013) interacts with dynamic river stages to create complex river intrusion and retreat pathways and dynamics. The time series of multi-year SpC
10 and temperature observations at the selected set of wells in the network have demonstrated these complicated processes of river water intrusion into our study site (Figure 1). Wells near the river shoreline (e.g., wells 1-1, 1-10A, 2-2, and 2-3) tend to be strongly affected by river water intrusion in spring and summer. As such, the dynamic patterns of SpC and temperature correspond well with river stage fluctuations, specifically that SpC decreases and temperature increases with increasing river stage. Fluctuations of SpC in well 2-2 appear to be stronger and at higher frequency than in other wells, likely indicating its
15 higher connectivity with the river. For wells that are farther inland (e.g., well 1-15), on the other hand, temperatures remain consistently within the groundwater temperature range and SpC has three noticeable dips (dropping from 0.5 to 0.4 mS/cm range), coinciding with the high river stages in years 2011, 2012, and 2017, which are featured with higher peak river stages than other years such that the river water was able to intrude further into the groundwater aquifer. In wells located at an intermediate distance from the river, such as Well 2-5, the intrusion of river water is evident in most of the years except in low-flow
20 years such as 2009 and 2015, during which both SpC and temperature remain nearly unchanged.

The understanding we developed from earlier studies is that the physical heterogeneity contributes to the different response behaviors at different locations while the river stage dynamics lead to multi-frequency dynamics in those responses. The seasonal and annual variations are driven by natural climatic forcing (Amaranto et al., 2019, 2018), whereas the higher-frequency (i.e., daily and sub-daily) fluctuations are primarily induced by operations of the upstream hydroelectric dam operations to
25 meet various demands of human society (Song et al., 2018). Our system is representative of many dam-regulated gravel-bed rivers across the world, where dam operations, as a typical anthropogenic activity, have significantly altered the hydrologic exchanges between river water and groundwater, as well as the associated thermal and biogeochemical processes (Song et al., 2018; Shuai et al., 2019; Zachara et al., 2020). Note that the multi-frequency variations in data characterize the dynamic features of data, which could exist in both short-term and long-term time series data as a result of short-term or long-term monitoring
30 effort.

To understand the multi-frequency variations of the river water and groundwater mixing in each well at the study site, we perform spectral analysis on multi-year SpC observations at each selected well using a discrete wavelet transform (DWT). The DWT is widely used for time–frequency analysis of time series and relies on a "mother wavelet", which is chosen to be the Morlet wavelet (Grossmann and Morlet, 1984) to deal with the time-varying frequency and amplitude in time-series data at
35 this site (Stockwell et al., 1996; Grinsted et al., 2004). We illustrate the Wavelet Power Spectrum (WPS) in log scale and its

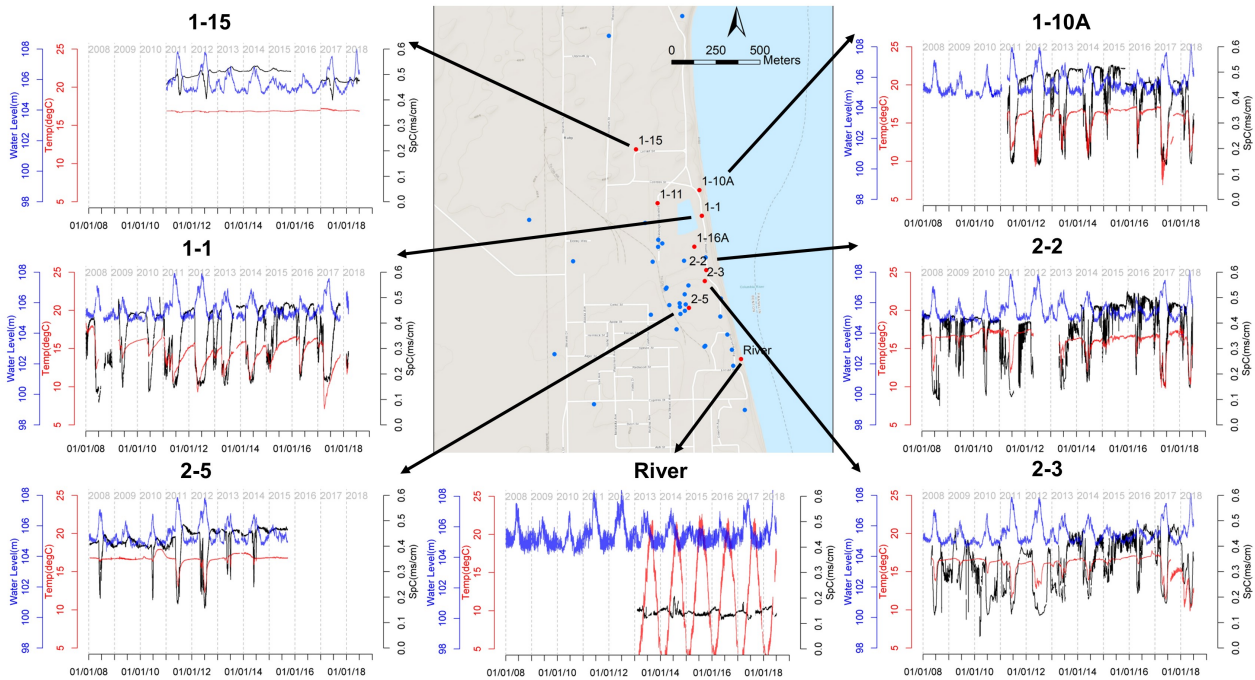


Figure 1. Groundwater monitoring well network at the 300 Area of the Hanford site and the monitoring data at select wells. Each well represented by a dot is instrumented to measure groundwater elevation, temperature, and SpC. The wells selected for this study are marked with red dots with well names. The three variables monitored in wells and in the river are shown in time-series plots with blue (water elevation), black (SpC), and red (temperature) lines. Base map @Google Maps.

normalized global power spectrum (average WPS over the time domain) for the multi-year SpC time series in the first two columns of Figure 2. Data gaps are shown as blank regions in Figure 2; examples include early year 2009 at well 1-1, the beginning of year 2011 at well 1-10A, and the later part of year 2012 at well 2-2. The amplitude of WPS represents the relative importance of variation at a given frequency compared to the variations at other frequencies across the spectrum. At wells 1-1, 1-10A, 2-3, 2-5, and 2-2, the strong intensities of SpC signals appear at the half-year and yearly frequencies; however, well 1-15 has a different pattern in that most of its high intensities are below the 256-hour frequency. The averaged WPS further shows the contrast in behaviors: wells 1-1, 1-10A, 2-3, and 2-5 have a dominant frequency at half a year; well 2-2 has multiple dominant frequencies at daily, monthly, and seasonal scales; while well 1-15 has similar intensities at the half-year and hourly scales. Applying this information to the task at hand, we hypothesize that gap filling at well 2-2 could be more challenging due to such mixture of dynamics signatures.

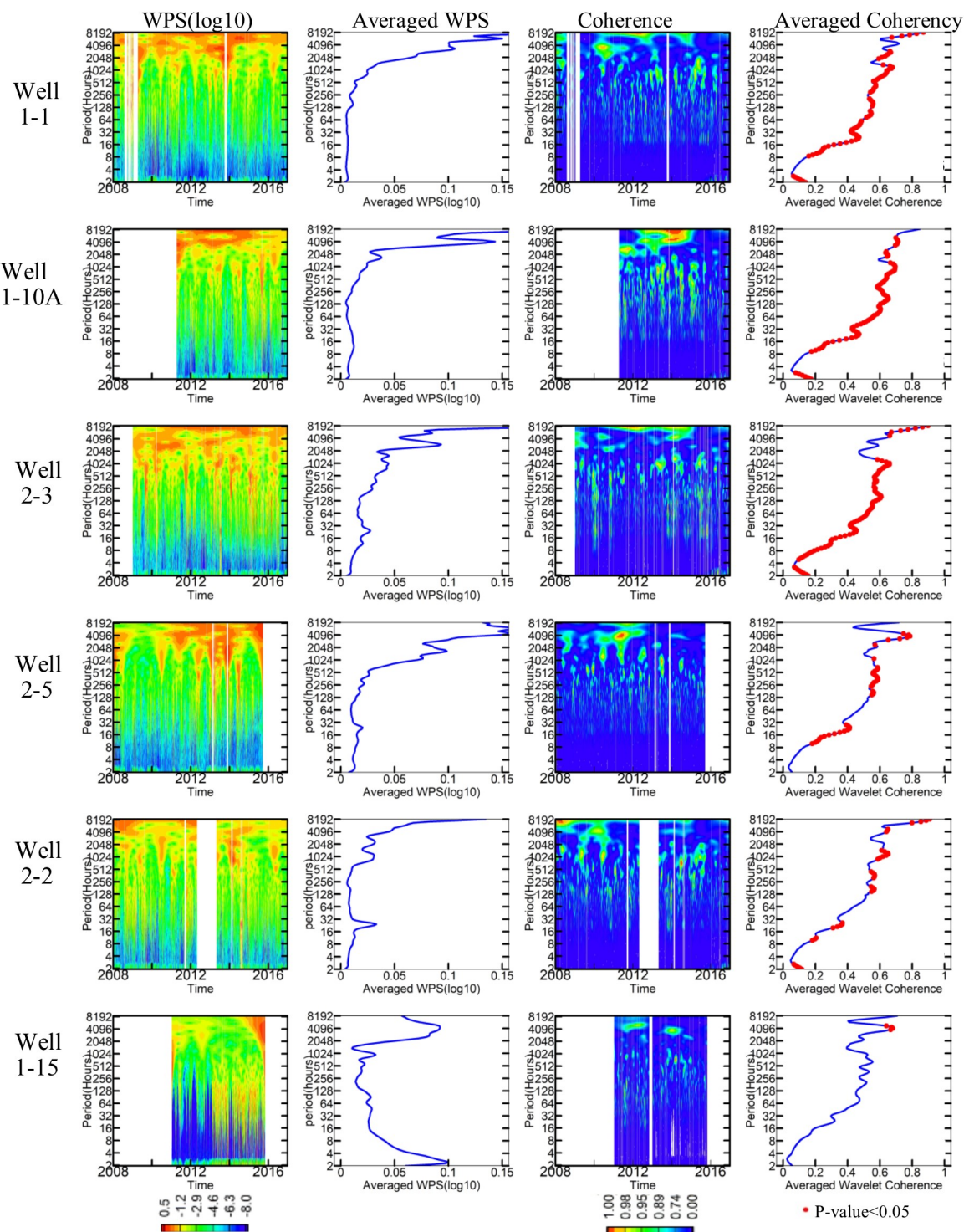


Figure 2. WPS analysis of SpC at each well from 2008 to 2018. The first column is the spectrogram (in log10 scale) of SpC in each well; the second column is the averaged WPS; the third column is the coherence between SpC in each well and the river stage; and the fourth column is the averaged coherence with $p < 0.05$ values indicated in red.

Since the dynamics of the system are driven by the river stage, we perform magnitude-squared wavelet coherence analysis via the Morlet wavelet to reveal dynamic correlations between the SpC and river stage time series (Grinsted et al., 2004; Vacha and Barunik, 2012). Wavelet coherence in the time-frequency domain is plotted in the third column in Figure 2 and the average coherence is plotted in the fourth column; statistically significant values at the 95th percent confidence interval are indicated with red points. A larger coherence at a given frequency indicates a stronger correlation at that frequency between the SpC at a well and the river stage. We consider these two variables highly correlated when the coherence is larger than 0.7 (shown in green to red colors in Coherence plots). We found that such high correlations exist at multiple frequencies, from subdaily to daily to yearly, at all the wells close to the river (e.g., 1-1, 1-10A, 2-2, and 2-3), while the higher correlation regimes in wells farther from the river (e.g., 1-15 and 2-5) are shifted towards longer periods at semi-annual and annual frequencies and less persistent in time.

As can be seen in Figure 2, many of the wells have long data gaps, which have unknown effects on our ability to estimate dynamics from the wavelet spectra. As such, gap filling is needed to infer observations and guide modeling of the underlying system. Figure 3 provides a summary of gap lengths for the overall network of monitoring wells. The majority of the gap lengths of all the three monitored variables are less than 50 hours. Therefore, in our investigations we explore the ability of the methods in filling gaps of 1-, 6-, 12-, 24-, 48-, and 72-hour lengths using hourly data over an input window to capture the multi-frequency fluctuations.

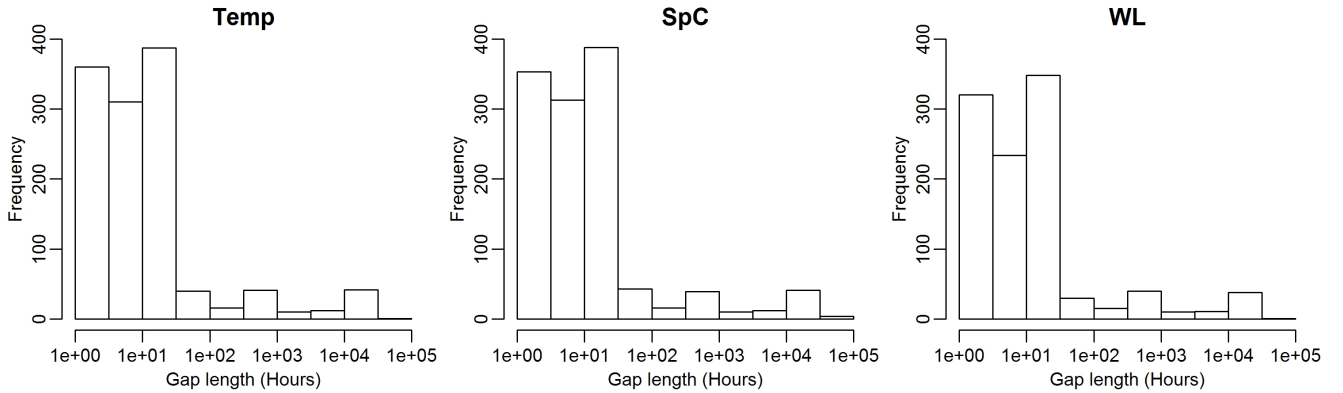


Figure 3. Histograms of gap lengths for each monitored variable, aggregated across all wells in the monitoring network during 2008-2018.

3 Gap-Filling Methods

In this section, we describe two methods we implemented to fill gaps of various lengths in SpC measurements at selected wells: an LSTM model and the traditional ARIMA model. ~~While the LSTM model can be used to fill in gaps in all three variables, we~~ We focused our analyses on filling gaps in SpC because of its importance to reveal river water and groundwater mixing.

- 5 Same set of analyses can be performed on water level and temperature. An input with M time steps (input window length) is provided to both ~~models for estimating the missing SpC measurements. The LSTM model predicts the~~ LSTM and ARIMA models for predicting the next time step that immediately follows the input window. For gaps larger than one hour, ~~this LSTM model is the~~ gap-filling models are applied to fill in one missing value at a time. The entire gap is filled by sliding the input window forward hour by hour and treating the gap-filled values of the previous missing hours as observed values. ~~The ARIMA model, on the other hand, fills the entire gap length using the input of M time steps preceding the gap.~~

- The input window may contain multiple variables from a single or multiple wells that are relevant to the prediction. After experimenting with different sets of input variables (SpC only, SpC and water level, SpC plus water level and temperature), we found that including SpC and water level measurements in the input window yielded the most robust performance. Therefore, we used historic water level (m) and SpC (mS/cm) observations to fill gaps in SpC time series of a single well. Using
- 15 measurements from multiple wells as input allows the models to account for both the temporal and spatial correlations in the data to impact gap-filling performance. Wells were selected based on adequate data availability and their distances from the target well. Assuming the observations from W ($W \geq 1$) wells are used to fill in data gaps, the input size of the model is then $M \times 2W$.

3.1 LSTM Models for Gap Filling

- 20 We designed an LSTM architecture, as shown in Figure 4, to train models of an input size of M time steps and an output size of one time step. The LSTM model contains a single LSTM layer followed by an output dense layer. The detailed structures of the LSTM layer is provided in the supplemental materials (Figures S1 and S2).

- Training data for the LSTM models were created by finding data segments of $M + 1$ hours that have no missing values, i.e., no gaps in the data, for both measurements over a specified monitoring window. The well data were then pre-processed
- 25 by normalizing all measurements ~~via zero-mean~~ to zero mean and unit variance for each variable, as SpC is on a scale of 10^{-1} , and water level is on a scale of 10^2 . Validation datasets were used to select the best model hyperparameters and the optimal input window size M (3.1.1) for gap filling at each well. Another independent testing period was selected at each well, depending on data availability, to compare the gap filling performance using the LSTM and ARIMA methods. The complete set of alternatives we considered for each LSTM model configuration is shown in Table 1. ~~We used an~~ We used the Adam
- 30 optimizer (Kingma and Ba, 2014) for training and the mean-squared error as the loss function. The models were trained for 50 iterations (i.e., epochs) over the training data.

To evaluate the accuracy of the trained LSTM models in filling SpC data gaps during the validation and testing processes, we assumed that synthetic gaps of various lengths (e.g., 1, 6, 12, 24, 48, and 72 hours, referred to as gap scenarios hereafter)

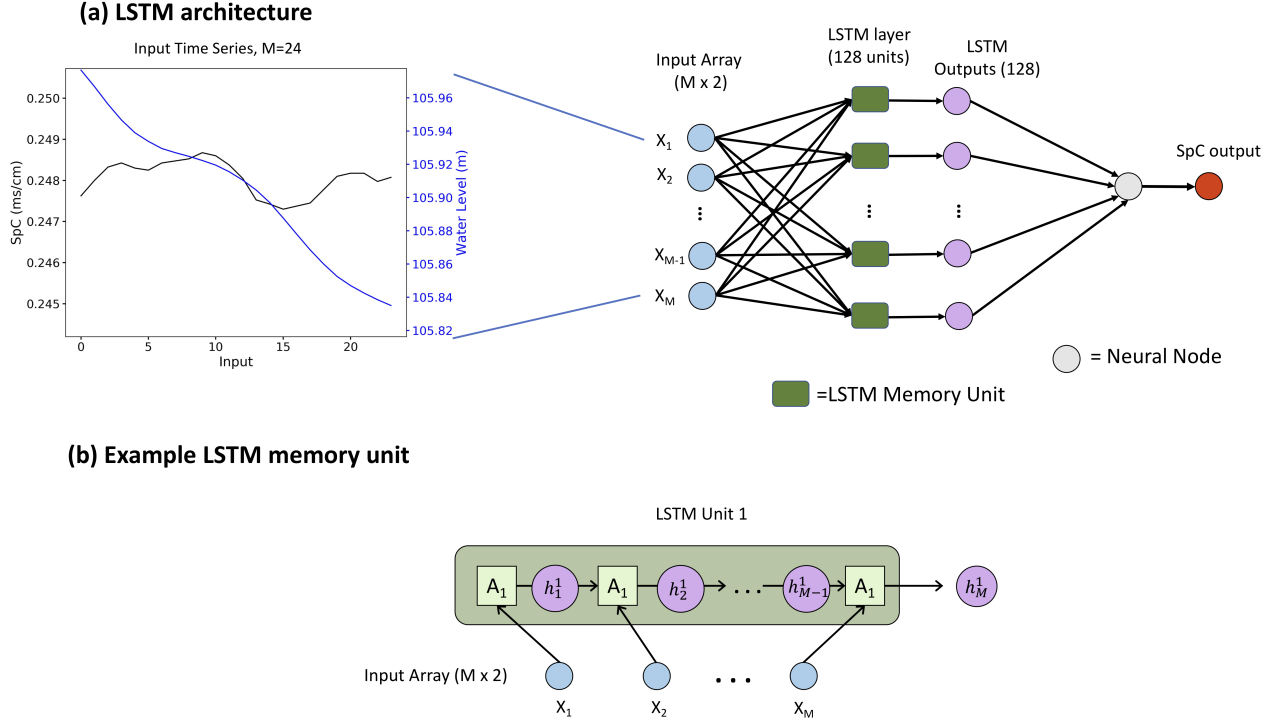


Figure 4. Illustration of LSTM models for gap filling. (a) Architecture of the LSTM models, where M is the input window size. Includes example input with $M = 24$ and example LSTM layer with 128 units; (b) Example of an LSTM unit, where A is the repeating module of the LSTM unit and h is the output.

exist in the validation or testing dataset of a well. We assume only the SpC measurements are missing while the water level measurements are available. Then an LSTM model configured with an input of M is given M hours of data from the time series preceding the occurrence of a gap (assuming no missing values in these M hours) to fill the gap hour by hour. The accuracy of the gap-filling model is evaluated by calculating the mean absolute percentage error (MAPE; %) between the SpC values that

5 are filled in (i.e., predicted) and that were observed:

$$MAPE = 100 \times \frac{1}{n} \sum_{t=1}^n \left| \frac{\text{Prediction} - \text{Observation}}{\text{Observation}} \right|, \quad (1)$$

where n is the total number of synthetic missing data points during the evaluation period.

In addition to the LSTM models trained for the single-well setup, we also trained multi-well [LSTM](#) models that used observations from wells 1-1, 1-10A, and 1-16A to fill in data gaps for well 1-1. We explored the same set of configuration parameters for single-well models as shown in Table 1 in multi-well LSTM models. We then compared the gap filling performance of the multi-well LSTM with the single-well LSTM model for well 1-1. The multi-well models were not explored for the other wells due to lack of neighboring wells in close proximity.

Table 1. Parameters used in training single-well LSTM models.

Parameter	Values
Training wells	1-1, 1-10A, 1-15, 2-2, 2-3, 2-5
Synthetic gap length (hours)	1, 6, 12, 24, 48, 72
Model input window (M hours)	24, 48, 72, 96, 120, 144, 168
LSTM Units (U units)	32, 64, 128
Learning Rate (L)	1e-3, 1e-4, 1e-5
Training period	2012-2015
Validation period ^a	2011
Testing Period ^b	2008 for well 2-5; 2017 for well 1-15; 2016 for all other wells

^a used to select the best LSTM model configurations and hyperparameters.

^b used to evaluate performance of LSTM vs ARIMA.

3.1.1 Optimizing LSTM model configuration

We performed a hyperparameter search using a grid-search approach to explore different LSTM model configurations, including the hyperparameter configurations to find the best model for a given gap length at each well. This involved iterating over all combinations of input time window size (M), the number of units (U) in the LSTM layer, and the learning rate (L) at-listed in Table 1 for each well. We chose the optimal LSTM configuration using model performance on the validation data set (see Table 1) based on the MAPE defined in Eq. (1). The combination that yielded the lowest SpC MAPE were selected as the best configuration for a given gap length at each well. These configurations are shown in Table S1 in the online supplemental material. The best model configurations were then used to evaluate the LSTM-based gap filling method against the ARIMA-based method (3.2) using relative errors (similar to MAPE by setting $n=1$ in Eq. (1) and removing the absolute value operation) calculated for each data point in the testing period (table Table 1), which varied among the wells due to availability of continuous data required for testing.

3.2 ARIMA Models for Gap Filling

ARIMA is one of the most general classes of models for extrapolating time series to produce forecasts and we used it as a baseline to compare and assess the LSTM gap-filling method. ARIMA can be applied to nonstationary time series data using a combination of differencing, autoregressive, and moving average components. A nonseasonal ARIMA(p, d, q)

model is given by:

$$Y_t = c + \phi_1 Y_{t-1}^d + \phi_p Y_{t-p}^d + \dots + \theta_1 e_{t-1} + \theta_q e_{t-q} + e_t, \quad (2)$$

where ϕ s and θ s are polynomials of orders p and q , respectively, each containing no roots inside the unit circle. e s are the error terms, Y_t^d is Y_t differenced d times, and c is a constant. Note that only non-seasonal terms (p, d, q) are included in Eq. (2).

5 Seasonal structure can be added with parameters $(P, D, Q)_m$ to the base ARIMA model to become $\text{ARIMA}(p, d, q)(P, D, Q)_m$, with a periodic component containing m periods. $c \neq 0$ implies a polynomial of order $d + D$ in the forecast function. The detailed mathematical equations for the seasonal ARIMA are provided in the online supplemental materials.

The main task in ARIMA-based forecasting is to select appropriate model orders, i.e., the values of p, q, d, P, Q, D . If d and D are known, we can select the orders p, q, P, Q via an information criterion such as the Akaike Information Criterion (AIC):

$$10 \quad AIC = -2\log(L) + 2(p + q + P + Q + k), \quad (3)$$

where $k = 1$ if $c \neq 0$ and 0 otherwise, and L is the ~~maximized~~ likelihood of the model fitted to the differenced data. The ~~best-fitted parameters of the ARIMA model can be determined by minimizing the AIC. The ARIMA-ARIMA~~ models were built using the `auto.arima` function from the R package `forecast` (Hyndman et al., 2007), which applies the Hyndman-Khandakar algorithm (Hyndman and Khandakar, 2008) that minimizes the AIC to obtain the best-fit parameters of the ARIMA model.

15 Similar to the LSTM-based gap filling, we explored various input window sizes, from 24 to 504 hours in an increment of 24 hours, for the ARIMA model at each well to identify optimal input windows within the search range. An optimal input window size is chosen for each gap length of each well using the same MAPE metric (i.e., Eq. (1)) on the validation dataset.

4 Results and Discussion

20 4.1 Performance of single-well LSTM models

We selected the best combination of LSTM units (U) and learning rate (L) for each input time window (M) under each gap length at each well using the MAPE metric. The validation MAPEs of those selected models were summarized in boxplots under different grouping, as shown in Figure 5. Each MAPE boxplot was drawn from a group of models with one parameter (corresponding to each x-axis) fixed at the given value while all the other parameters, including the training wells and gap scenarios, cycle through their possible combinations.

25 As shown in Figure 5 (a), model performance deteriorates as the gap length increases, indicating that the LSTM-based method tends to lose ground truth information from its input to inform prediction. In comparing MAPEs across various input window sizes shown in Figure 5 (b), we observe that models with all input windows have comparable MAPE summary statistics, with those larger input windows (> 96 hours) leading to slightly smaller MAPE quartiles. The larger input windows also yield fewer outliers on the larger MAPE end, indicating that the memory units in the LSTM layers are capturing important daily to weekly signatures (evident in WPS plots in Figure 2 for all wells except for Well 1-15) for some wells.

The performance of single-well LSTM models varied among the wells as shown in Figure 5 (c). The LSTM models for well 1-15 lead the performance with the smallest MAPEs, while those for well 2-2 yield the worst performance. The LSTM models for wells 1-1, 1-10A, 2-3 and 2-5 performed comparably overall, with slightly more large MAPE outliers for well 2-3.

4.2 Single-well LSTM and ARIMA comparisons

- 5 The single-well LSTM gap filling approach was compared to the ARIMA approach using relative errors calculated for each data point that was assumed to be missing in the testing data by setting $n=1$ in Eq. (1) for MAPE. Relative errors were used to show overestimations or underestimations by both approaches. Their respective best model configurations determined on the validation dataset (i.e., data from year 2011) were used in comparing the two approaches. Figure 6 illustrates the optimal input windows for the LSTM and ARIMA methods. We observe that LSTM models require ~~less or equal~~ much less input information
- 10 than the ARIMA method under all gap lengths for all the wells ~~except well 2-5~~.

- Figure 7 shows the boxplots of relative errors under different gap lengths for all testing wells. For both approaches, the relative errors increase as the gap length increases as expected. The ARIMA models tend to perform better than the LSTM models in terms of interquartile range. However, the ARIMA models, in general, produce more outliers of large positive or negative relative errors than the LSTM models, especially for larger gap lengths (48 and 72 hours). For well 1-15, the relative errors of
- 15 both approaches are small for all gap lengths. Both approaches appear to have larger error outliers at well 2-3.

- For each well, we performed a T-Test to calculate the T-Score and P-Value between the relative errors of the two models to determine the significance in the difference of the performance of the models. As seen in Table 2, each well has a high T-Score and P-Value significantly less than 0.05. Thus, the differences between their relative errors are significant and meaningful.
- In addition to the relative errors, we calculated the MAPE, Root Mean Squared Error (RMSE), Nash–Sutcliffe model efficiency coefficient (NSE) (Nash and Sutcliffe, 1970), and Kling-Gupta Efficiency (KGE) (Gupta et al., 2009) for the best LSTM and
- 20 ARIMA model per gap length to compare the two approaches. Table 2 compares the performance of the LSTM and ARIMA models filling in gap lengths of 24 hours. The table for all gap lengths is in the online supplemental material (Table S2).
- NSE is a metric used to assess the predictive skills and accuracy of hydrological models. Values range from $-\infty$ to 1, where 1 indicates a perfect model fit, 0 indicates that the model has the same predictive power as the mean of the observations, and less
- 25 than 0 indicates that the model is a worse predictor than the mean of the observations. NSE is calculated on the SpC predictions by the following equation:

$$NSE = 1 - \frac{\sum_{t=1}^n (P_t - O_t)^2}{\sum_{t=1}^n (O_t - \mu(O))^2}, \quad (4)$$

- where n is the total number of synthetic missing data points during the evaluation period, P_t and O_t are the predicted and observed SpC values at time t , and $\mu(O)$ is the mean observed SpC value.
- 30 KGE is another goodness-of-fit metric used to evaluate hydrological models by combining the three components of NSE of model errors (i.e. correlation, bias, ratio of variances or coefficients of variation) in a more balanced way. It has the same range

of values as NSE, where 1 indicates a perfect model fit. KGE is calculated on a model's SpC predictions by the following equations:

$$KGE = 1 - \sqrt{(r-1)^2 + (\alpha-1)^2 + (\beta-1)^2} \quad (5)$$

$$r = \frac{\text{cov}(O, P)}{\sigma(O) * \sigma(P)} \quad (6)$$

$$\alpha = \frac{\sigma(P)}{\sigma(O)} \quad (7)$$

$$\beta = \frac{\mu(P)}{\mu(O)}, \quad (8)$$

where cov is the covariance, σ is the standard deviation, and μ is the arithmetic mean.

The LSTM and ARIMA models yielded comparable average metrics at all the wells for the gap length of 24 hours, as can be seen in Table 2. The NSE and KGE resulted from both models are close to 1 for all the wells with negligible differences between the two models. The difference in MAPE and RMSE is also small, with relatively more notable differences for wells 2-2 and 2-3, where the ARIMA models resulted in lower MAPE and RMSE.

Table 2. Comparison of single-well LSTM and ARIMA models on 24-hour synthetic gap in the SpC data on the test set for each well. The models are the same ones used in Figure 8. The calculated statistics are: MAPE, Root Mean Squared Error (RMSE), Nash–Sutcliffe model efficiency coefficient (NSE), and Kling-Gupta Efficiency (KGE). The T-Score and P-Value are calculated on the relative errors of the two models per well.

Well	Model Type	MAPE	RMSE	NSE	KGE	T-Score	P-Value
1-1	LSTM	1.38	8.33×10^{-3}	0.991	0.988	19.1	1.00×10^{-80}
	ARIMA	1.36	8.98×10^{-3}	0.989	0.994		
1-10A	LSTM	1.37	8.07×10^{-3}	0.986	0.968	-24.6	1.48×10^{-131}
	ARIMA	1.5	9.60×10^{-3}	0.98	0.987		
1-15	LSTM	0.259	1.88×10^{-3}	0.989	0.982	-48.9	0.00
	ARIMA	0.119	1.18×10^{-3}	0.996	0.997		
2-2	LSTM	2.97	1.87×10^{-2}	0.922	0.962	48.1	0.00
	ARIMA	2.23	1.64×10^{-2}	0.939	0.967		
2-3	LSTM	2.15	1.63×10^{-2}	0.945	0.965	21.6	4.69×10^{-102}
	ARIMA	1.72	1.48×10^{-2}	0.954	0.971		
2-5	LSTM	0.929	6.86×10^{-3}	0.976	0.988	-9.6	9.22×10^{-22}
	ARIMA	0.866	7.45×10^{-3}	0.971	0.977		

In addition to the error statistics, it is also important to examine how well a gap-filling method captures the desired dynamics patterns in the gap-filled time series. Therefore, the SpC time series reproduced by the gap-filling methods during the testing period (2016 for wells 1-1, 1-10A, 2-2, 2-3; 2017 for well 1-15; 2008 for well 2-5) with 24-hour synthetic gaps are evaluated against the real time series. Model configurations are the same as those used in error statistics comparison (Figure 6). A gap length of 24 hours is selected as an example because we consider it as a reasonably challenging case to fill gaps in time series data exhibiting significant nonstationarity, such as the SpC data at well 2-3. Moreover, the relative performance between the two approaches are similar at other gap lengths with varying error magnitudes.

As shown in the first column of Figure 8, both approaches capture the general dynamic patterns in the data fairly well. The time series of relative errors for both methods are provided in Figure S3 in the online supplemental materials for more details.

The ARIMA approach (blue lines in column 1) missed some abrupt changes that occur over a short time window (i.e., at higher frequency), leading to more error spikes in all wells, consistent with relative error outliers in Figure 7. This is an indication that the ARIMA models lack mechanisms to represent such high-frequency changes. The LSTM approach (red lines in column 1 plots), on the other hand, is able to better capture such dynamics in all the wells. However, the inclusion of such high-frequency fluctuations may also lead to overly dynamic predictions in time windows dominated by lower-frequency fluctuations, which contributed to less desirable relative errors distributed between the first and third quartiles in some wells (i.e., wells 1-1, 2-3, 2-2, and 1-15), as shown in Figure S4 in the supplemental materials. This is likely caused by the variability in dynamics signatures among the training, validation and testing periods, as well as the selection of training loss functions and validation metric that balance between the occurrence of small vs large errors to achieve optimal solutions.

To further investigate how the relative performance of the two gap-filling methods depends on the inherent dynamics in each time series, spectral analyses for the testing SpC datasets were performed using the same wavelet decomposition method for the multi-year analyses (shown earlier in Figure 2). As shown in Figure 8, the time windows of high relative errors are found to approximately co-locate with the time when the high-frequency (daily and subdaily) signals are gaining more power. The difference between the LSTM and ARIMA models tend to be amplified during those time windows. Wells 1-1, 1-10A, and 2-2 share similar seasonal patterns in WPS, with the highest intensity bin above 1024 hours across February to July. Their average WPSs all show peaks around daily and subdaily frequencies. Well 2-3 has its greatest energy between 16 to 256 hours from January to July. Well 2-5 has low intensities of variability at daily and subdaily frequencies with the low-frequency variations (monthly and seasonal) dominating the Jan to March time frame. For well 1-15, one of its strongest intensities is above 2048 hours across the entire year, and the other strong intensities are narrow bands between 16 to 256 hours. In general, both LSTM and ARIMA are effective at capturing low-frequency variability (monthly and seasonal). Although LSTM is more effective at capturing high-frequency (daily and subdaily) fluctuations and nonlinearities in the datasets, it may also lead to overly dynamic predictions in time windows with no considerable high-frequency fluctuations. However, the errors during these time windows are small and can be improved by smoothing if such fluctuations are not desirable.

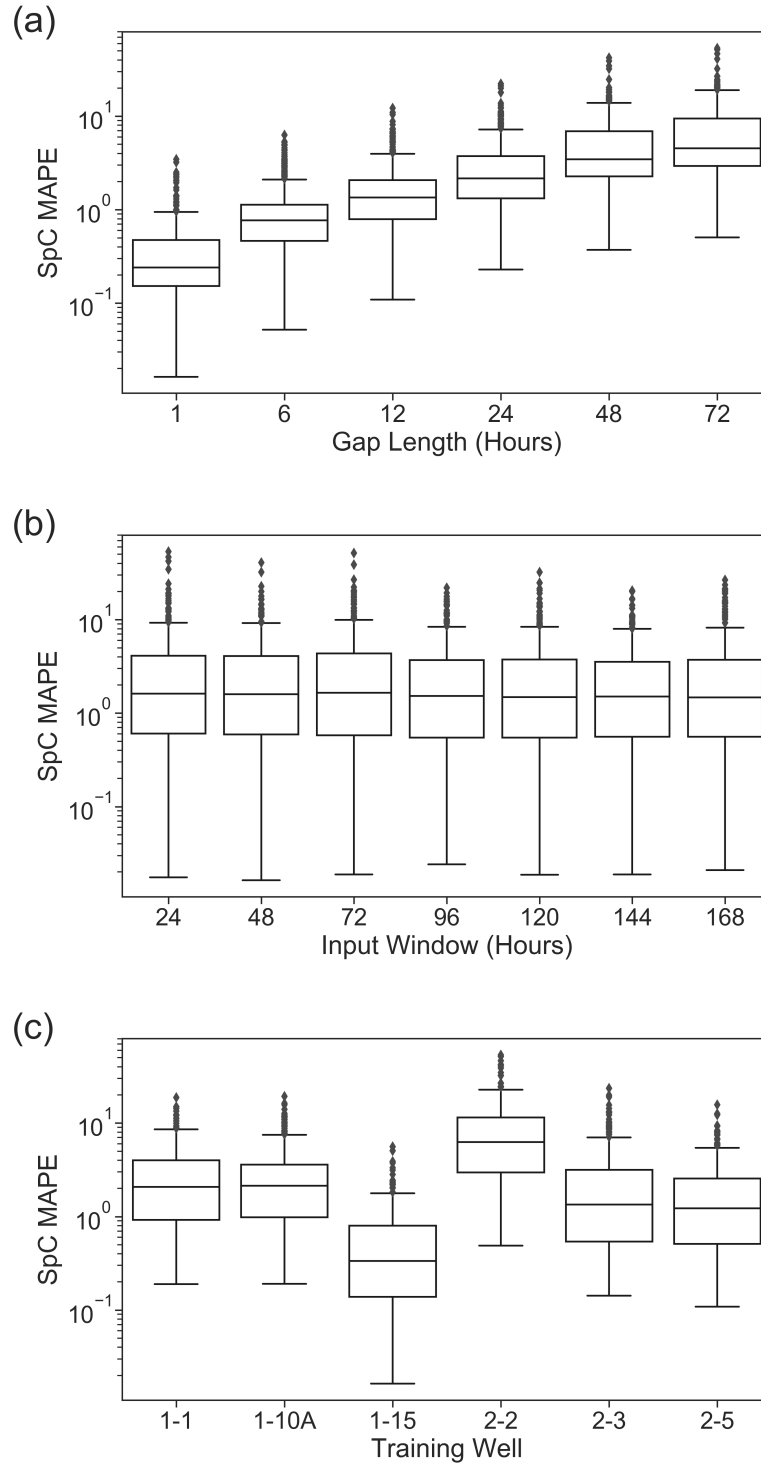


Figure 5. Gap filling performance for SpC evaluated against the validation datasets, grouped by gap lengths (a), model input window size M (b), and training wells (c).

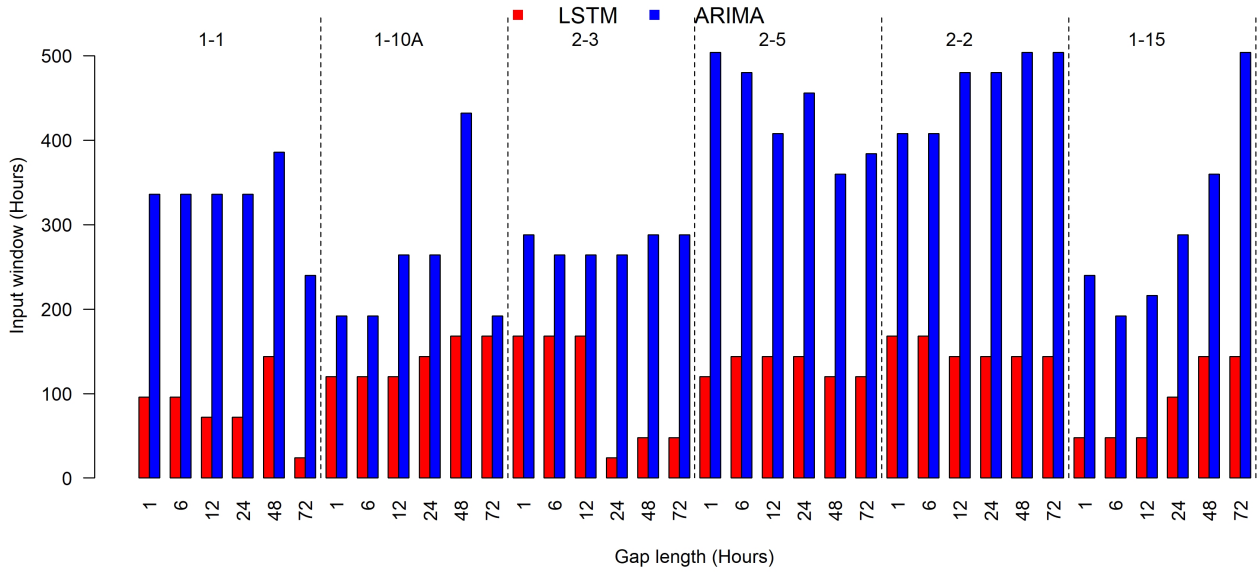


Figure 6. Optimal input windows for the LSTM and ARIMA models to fill gaps of various lengths at each well.

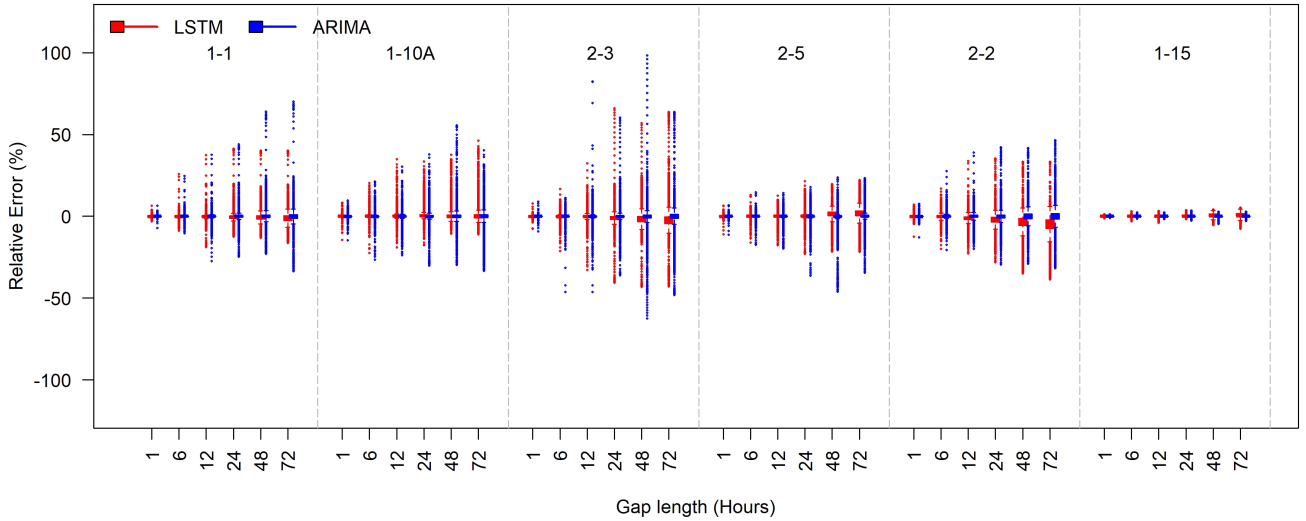


Figure 7. Boxplots of relative errors for filling SpC gaps of various lengths (1, 6, 12, 24, 48, and 72 hours) at each well during the test periods. The best LSTM and ARIMA models were used for evaluation. The LSTM and ARIMA models are represented by red bars and blue bars, respectively.

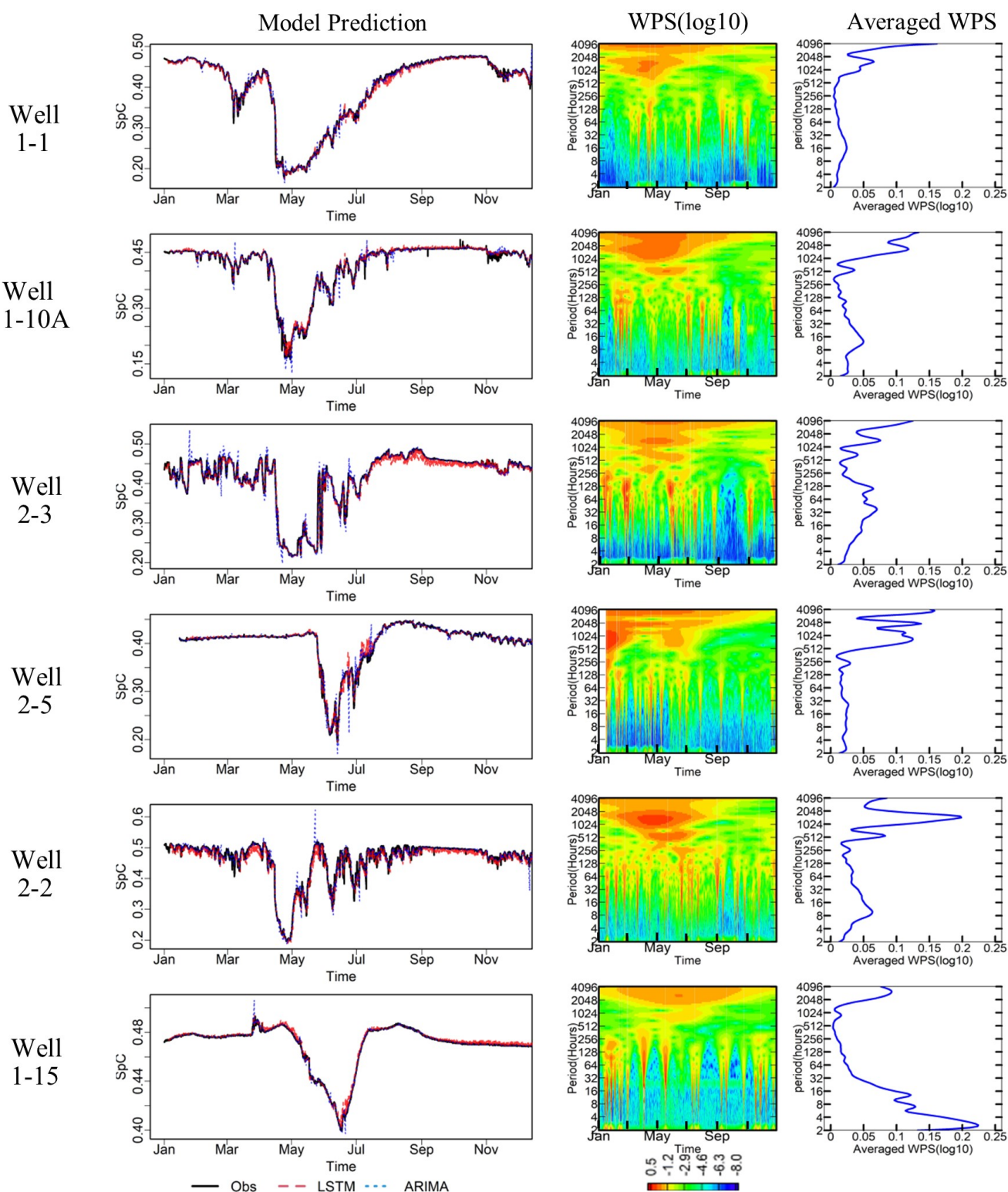


Figure 8. Columns 1 shows time series of model predictions from the LSTM (in red) and ARIMA (in blue) methods, respectively, assuming 24-hour synthetic gap in the SpC data, compared with observations (in black). The best model configurations were used for all models. The testing data come from year 2016 for wells 1-1, 1-10A, 2-2, and 2-3, from year 2017 for well 1-15 and from 2008 for well 2-5. Column 2 is the spectrogram of each well and column 3 is the WPS averaged over for the corresponding year.

5

There is also significant difference in computational cost between the LSTM and ARIMA methods for gap filling. ARIMA requires very little computational resources: the `auto.arima` function in R requires approximately 40 seconds for fit and validate a model for each prediction segments on a personal computer with a 3.00 GHz CPU. Conversely, training and validating a single LSTM model takes approximately 20-30 minutes on dual NVIDIA P100 12GB PCI-e based GPUs.

10 4.3 Performance of multi-well models

We evaluated the predictive ability of the multi-well models using both approaches in filling gaps of various lengths in the SpC data at well 1-1 by comparing the performance against their single-well counterparts. Well 1-1 was chosen because of data availability in nearby wells (wells 1-10A and 1-16A). Similar to the single-well ARIMA and LSTM model for well 1-1, the multi-well models also predict the SpC measurement using water level and SpC from three wells. We adopted the same LSTM architecture from the single-well LSTM model and trained the same set of alternatives considering input window sizes, LSTM units, and learning rates for various gap lengths as listed in Table 1. The same training and validation periods were adopted to select the optimal combination of M , U , and L . The optimal combinations are shown in Table S3 in the online supplemental material. For the multi-well ARIMA models, we included additional variables as regression terms when building and fitting models using the `same` `auto.arima` function. The optimal input window sizes of ARIMA are 216, 240, 288, 288, 288, and 192 hours for gap length 1, 6, 12, 24, 48, and 72 hours, respectively. They are smaller than the optimal window sizes of the single-well models.

The boxplots of relative errors yielded from the single-well and multi-well models using both approaches are provided in Figure 9 for comparison. Additionally, we include performance metrics comparing the single and multi-well models in Table 3. Additional spatial information ~~helps the ARIMA models to reduce the large errors when the gaps are small~~ seems to exacerbate the relative errors by the ARIMA models, except for large gaps (e.g., ~~1 and 6 hours~~), ~~while it exacerbates the errors for gaps larger than 6 hours~~ 72 hours). The LSTM approach, on the contrary, benefits from the information carried by the neighboring wells to fill in those larger gaps, while the performance for small gaps stay unchanged, ~~indicating that~~. The aggregated performance metrics in Table 3 show slightly improved metrics for multi-well ARIMA models for gaps smaller than 24 hours compared to the single-well models, while the turning point in relative performance is at 12 hours for the LSTM models. The deteriorated performance metrics of the multi-well LSTM models at the larger gap lengths are consistent with their larger inter-quartile ranges as revealed by the boxplots of relative errors in Figure 9. However, the multi-well LSTM and ARIMA models can reduce the occurrence of large relative errors for larger gaps, providing more robust gap-filling under those circumstances. This comparisons show that, although the information from a single well is may be sufficient to fill in those small gaps, Therefore gaps smaller than a day, including spatial information from neighbouring wells in the LSTM and ARIMA models

could potentially increase the chance of successes in filling data gaps under more challenging circumstances, such as capturing more complex dynamic patterns with longer data gaps. While the aggregated metrics provide an overall assessment of model performance, examining the distribution of relative errors could provide complementary information on the large error spikes when selecting optimal model configurations.

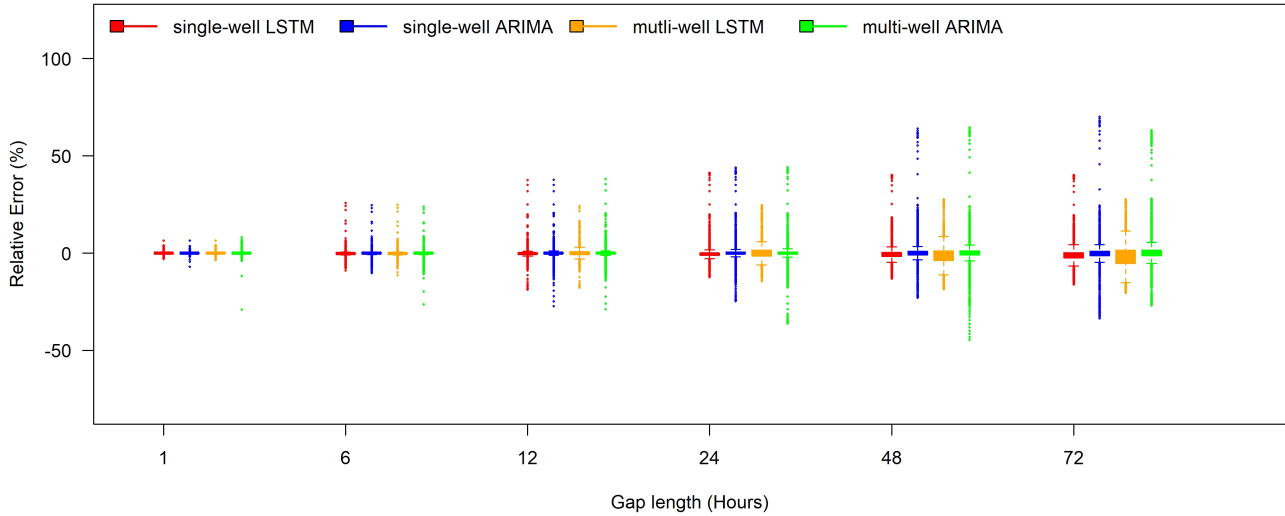


Figure 9. Comparing relative error performance between the best single-well LSTM models (well 1-1, red), multi-well LSTM models (wells 1-1, 1-10A and 1-16A, yellow), single-well ARIMA model (blue), and multi-well ARIMA model (green) for filling in various SpC gap lengths for well 1-1 during the test period (year 2016).

5 Conclusion

In this study, we implemented an LSTM-based gap filling method to account for spatio-temporal correlations in monitoring data. We extensively evaluated the method on filling data gaps in groundwater SpC measurements that are often used to indicate groundwater and river water interactions along river corridors. We optimized an LSTM architecture to take advantage of a 10-year spatially distributed multi-variable time series dataset collected by a groundwater monitoring well network for filling SpC data gaps. A primary advantage of using LSTM is the ability to incorporate spatio-temporal correlations and nonlinearity in model states without assuming a priori an explicit form of correlations or nonlinear functions in advancing system states. We compared the performance of single-well LSTM-based gap-filling method with a traditional gap-filling method, ARIMA, to evaluate how well an LSTM model can capture multi-frequency dynamics. We also trained LSTM and ARIMA models that take input from multiple wells to predict responses at one well. The multi-well models were compared with single-well models to assess the improvement in gap filling performance by including additional spatial correlation from neighboring wells.

Table 3. Comparison of single-well and multi-well LSTM and ARIMA models for all synthetic gap lengths in the SpC data. The models are the same ones used in Figure 9. Calculations are performed on the test data set for well 1-1 (year 2016). The calculated statistics are: MAPE, Root Mean Squared Error (RMSE), Nash–Sutcliffe model efficiency coefficient (NSE), and Kling-Gupta Efficiency (KGE). T-Score and P-Value are calculated on the relative errors of the two single-well models for each gap length and calculated on the relative errors of the two multi-well models.

Gap Length	Model Type	MAPE	RMSE	NSE	KGE	T-Score	P-Value
1	Single-Well LSTM	0.117	7.76×10^{-4}	1.0	0.999	11.6	6.14×10^{-31}
	Single-Well ARIMA	0.183	1.23×10^{-3}	1.0	1.0		
	Multi-Well LSTM	0.117	7.94×10^{-4}	1.0	1.0	9.37	8.54×10^{-21}
	Multi-Well ARIMA	0.134	1.22×10^{-3}	1.0	1.0		
6	Single-Well LSTM	0.435	2.92×10^{-3}	0.999	0.995	13.8	4.46×10^{-43}
	Single-Well ARIMA	0.461	3.29×10^{-3}	0.999	0.999		
	Multi-Well LSTM	0.435	2.98×10^{-3}	0.999	0.998	12.2	5.82×10^{-34}
	Multi-Well ARIMA	0.405	3.16×10^{-3}	0.999	0.999		
12	Single-Well LSTM	0.75	5.07×10^{-3}	0.997	0.996	16.0	2.05×10^{-57}
	Single-Well ARIMA	0.781	5.47×10^{-3}	0.996	0.996		
	Multi-Well LSTM	1.19	6.48×10^{-3}	0.994	0.985	5.25	1.55×10^{-7}
	Multi-Well ARIMA	0.77	5.40×10^{-3}	0.996	0.997		
24	Single-Well LSTM	1.38	8.33×10^{-3}	0.991	0.988	19.1	1.00×10^{-80}
	Single-Well ARIMA	1.36	8.98×10^{-3}	0.989	0.994		
	Multi-Well LSTM	2.26	1.17×10^{-2}	0.982	0.968	7.77	8.48×10^{-15}
	Multi-Well ARIMA	1.47	9.55×10^{-3}	0.988	0.99		
48	Single-Well LSTM	2.13	1.21×10^{-2}	0.98	0.988	17.8	3.24×10^{-70}
	Single-Well ARIMA	2.15	1.34×10^{-2}	0.976	0.988		
	Multi-Well LSTM	3.49	1.76×10^{-2}	0.958	0.969	28.4	7.83×10^{-174}
	Multi-Well ARIMA	2.35	1.40×10^{-2}	0.974	0.981		
72	Single-Well LSTM	2.56	1.40×10^{-2}	0.974	0.983	26.7	4.91×10^{-154}
	Single-Well ARIMA	2.57	1.46×10^{-2}	0.971	0.985		
	Multi-Well LSTM	4.41	2.19×10^{-2}	0.936	0.955	31.7	4.78×10^{-214}
	Multi-Well ARIMA	3.02	1.78×10^{-2}	0.958	0.972		

In general, both LSTM and ARIMA methods were highly accurate in filling smaller data gaps (i.e., 1 and 6 hours). They were reasonably effective at filling in medium gaps between 12 and 48 hours, while more work is needed for gaps larger than 48 hours. Both models captured the long-term trends in data (i.e., low-frequency variations at the monthly or seasonal time scales). The ARIMA method was found to have difficulty in capturing abrupt changes. Thus, it is more suitable for time series with

less dynamic behavior. Compared with the ARIMA models, the LSTM models excel in dealing with high-frequency dynamics (daily and subdaily) and nonlinearities, although they require more training data and computational resources. As a side effect of including high-frequency (daily and subdaily) fluctuations in the model, the LSTM approach may produce overly dynamic predictions in time windows that lacks such dynamics. Thus, availability of sufficient training data that cover a wide range of conditions is critical for the success of LSTM methods, as is with any deep learning method. Extrapolating the LSTM models to conditions beyond those in the training data remains as a major challenge.

Wavelet analysis could provide useful insights to the dynamic signatures of the data and the change in composition of their important frequencies over time, which can serve as a prior basis for selecting an appropriate gap-filling method. For example, the ARIMA method would work well if the dynamics are dominated by seasonal cycles, while more sophisticated approaches like LSTM-based methods could work better if there is evidence of weekly, daily and subdaily fluctuations. Depending on the mixture of high- and low-frequency variability inherent in the time series, different LSTM architecture and configurations can be explored and evaluated through hyperparameter search with respect to LSTM layers, dense layers and activation functions to achieve better performance in capturing more complex dynamics. The optimal LSTM model configuration and performance that could be achieved would vary case by case.

We also demonstrated that incorporating spatial information from neighboring stations in LSTM models could contribute to performance improvement under challenging scenarios with dynamic system behaviours with reasonably longer data gaps up to 2 days. However, other alternatives need to be explored for gaps beyond 2 days. The bidirectional and convolutional LSTMs are two promising methods to leverage information from the future time window and from spatially distributed networks. While we introduced a new method that can be broadly applied to fill in gaps in irregularly spaced network for monitoring groundwater and surface water interactions, the transferrability of this method to other monitoring systems could be evaluated more extensively by community participation. Capturing spatio-temporal dynamics in system states is essential for generating the most valuable insights to advance our understanding of dynamic complex systems.

Code and data availability. The well observations have been made accessible at <https://sbrsfa.velo.pnnl.gov/datasets/?UUID=14febd81-05b6-47fb-be52-439c4382decd>

Author contributions. HR and EC developed scripts and performed the analyses and they have equal contribution to the paper. BK contributed on interpretation of the results. XC conceived and designed the study. All authors contributed to writing the manuscript.

Competing interests. The authors declare that they have no conflicts of interest.

Acknowledgements. This research was supported by the U.S. Department of Energy (DOE), Office of Biological and Environmental Research (BER), as part of BER's ~~Subsurface Biogeochemical Research Program (SBR)~~ Environmental System Science (ESS) program. A portion of methodology development was supported by the Laboratory Directed Research and Development Program at Pacific Northwest National Laboratory, a multiprogram national laboratory operated by Battelle for DOE under contract DE-AC05-76RL01830. This research

5 was performed using PNNL Institutional Computing at Pacific Northwest National Laboratory. This research was also supported in part by the Indiana University Environmental Resilience Institute and the *Prepared for Environmental Change* grand challenge initiative. This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

References

- Alvera-Azcárate, A., Barth, A., Parard, G., and Beckers, J.-M.: Analysis of SMOS sea surface salinity data using DINEOF, Remote sensing of environment, 180, 137–145, 2016.
- Amaranto, A., Munoz-Arriola, F., Corzo, G., Solomatine, D. P., and Meyer, G.: Semi-seasonal groundwater forecast using multiple data-driven models in an irrigated cropland, Journal of Hydroinformatics, 20, 1227–1246, 2018.
- Amaranto, A., Munoz-Arriola, F., Solomatine, D., and Corzo, G.: A spatially enhanced data-driven multimodel to improve semiseasonal groundwater forecasts in the High Plains aquifer, USA, Water Resources Research, 55, 5941–5961, 2019.
- Banerjee, S., Carlin, B. P., and Gelfand, A. E.: Hierarchical modeling and analysis for spatial data, CRC press, 2014.
- Beckers, J.-M. and Rixen, M.: EOF calculations and data filling from incomplete oceanographic datasets, Journal of Atmospheric and oceanic technology, 20, 1839–1856, 2003.
- Beckers, J.-M., Barth, A., and Alvera-Azcárate, A.: DINEOF reconstruction of clouded images including error maps? application to the Sea-Surface Temperature around Corsican Island, 2006.
- Calculli, C., Fassò, A., Finazzi, F., Pollice, A., and Turnone, A.: Maximum likelihood estimation of the multivariate hidden dynamic geostatistical model with application to air quality in Apulia, Italy, Environmetrics, 26, 406–417, 2015.
- Chen, X., Murakami, H., Hahn, M. S., Hammond, G. E., Rockhold, M. L., Zachara, J. M., and Rubin, Y.: Three-dimensional Bayesian geostatistical aquifer characterization at the Hanford 300 Area using tracer test data, Water Resources Research, 48, <https://doi.org/10.1029/2011WR010675>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2011WR010675>, 2012.
- Chen, X., Hammond, G. E., Murray, C. J., Rockhold, M. L., Vermeul, V. R., and Zachara, J. M.: Application of ensemble-based data assimilation techniques for aquifer characterization using tracer data at Hanford 300 area, Water Resources Research, 49, 7064–7076, <https://doi.org/10.1002/2012WR013285>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2012WR013285>, 2013.
- Cheng, T., Haworth, J., and Wang, J.: Spatio-temporal autocorrelation of road network data, Journal of Geographical Systems, 14, 389–413, <https://doi.org/10.1007/s10109-011-0149-5>, <https://doi.org/10.1007/s10109-011-0149-5>, 2012.
- Cheng, T., Haworth, J., Anbaroglu, B., Tanaksaranond, G., and Wang, J.: Spatiotemporal data mining, in: Handbook of regional science, pp. 1173–1193, Springer, 2014.
- Datta, A., Banerjee, S., Finley, A. O., and Gelfand, A. E.: Hierarchical nearest-neighbor Gaussian process models for large geostatistical datasets, Journal of the American Statistical Association, 111, 800–812, 2016.
- Eidsvik, J., Shaby, B. A., Reich, B. J., Wheeler, M., and Niemi, J.: Estimation and prediction in spatial models with block composite likelihoods, Journal of Computational and Graphical Statistics, 23, 295–315, 2014.
- Fang, K., Shen, C., Kifer, D., and Yang, X.: Prolongation of SMAP to Spatiotemporally Seamless Coverage of Continental U.S. Using a Deep Learning Neural Network, Geophysical Research Letters, 44, 11,030–11,039, <https://doi.org/10.1002/2017GL075619>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2017GL075619>, 2017.
- Faruk, D. Ö.: A hybrid neural network and ARIMA model for water quality time series prediction, Engineering Applications of Artificial Intelligence, 23, 586–594, 2010.
- Finley, A. O., Banerjee, S., and Gelfand, A. E.: spBayes for large univariate and multivariate point-referenced spatio-temporal data models, arXiv preprint arXiv:1310.8192, 2013.

- Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., and Yacalis, G.: Could Machine Learning Break the Convection Parameterization Dead-lock?, *Geophysical Research Letters*, 45, 5742–5751, <https://doi.org/10.1029/2018GL078202>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018GL078202>, 2018.
- Ghil, M., Allen, M., Dettinger, M., Ide, K., Kondrashov, D., Mann, M., Robertson, A. W., Saunders, A., Tian, Y., Varadi, F., et al.: Advanced spectral methods for climatic time series, *Reviews of geophysics*, 40, 3–1, 2002.
- Grant, G. E. and Dietrich, W. E.: The frontier beneath our feet, *Water Resources Research*, 53, 2605–2609, 2017.
- Graves, A.: Generating sequences with recurrent neural networks, arXiv preprint arXiv:1308.0850, 2013.
- Graves, A., Mohamed, A.-r., and Hinton, G.: Speech recognition with deep recurrent neural networks, in: *Acoustics, speech and signal processing (icassp)*, 2013 IEEE international conference on, pp. 6645–6649, IEEE, 2013.
- 10 Griffith, D. A.: Modeling spatio-temporal relationships: retrospect and prospect, *Journal of Geographical Systems*, 12, 111–123, 2010.
- Grinsted, A., Moore, J. C., and Jevrejeva, S.: Application of the cross wavelet transform and wavelet coherence to geophysical time series, *Nonlinear processes in geophysics*, 11, 561–566, 2004.
- Grossmann, A. and Morlet, J.: Decomposition of Hardy functions into square integrable wavelets of constant shape, *SIAM journal on mathematical analysis*, 15, 723–736, 1984.
- 15 Güler, C. and Thyne, G. D.: Hydrologic and geologic factors controlling surface and groundwater chemistry in Indian Wells-Owens Valley area, southeastern California, USA, *Journal of Hydrology*, 285, 177–198, 2004.
- Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling, *Journal of Hydrology*, 377, 80–91, <https://doi.org/https://doi.org/10.1016/j.jhydrol.2009.08.003>, <https://www.sciencedirect.com/science/article/pii/S0022169409004843>, 2009.
- 20 Han, P., Wang, P. X., Zhang, S. Y., and Zhu, D. H.: Drought forecasting based on the remote sensing data using ARIMA models, *Mathematical and Computer Modelling*, 51, 1398–1403, 2010.
- Ho, S., Xie, M., and Goh, T.: A comparative study of neural network and Box-Jenkins ARIMA modeling in time series prediction, *Computers & Industrial Engineering*, 42, 371–375, 2002.
- 25 Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, *Neural Comput.*, 9, 1735–1780, <https://doi.org/10.1162/neco.1997.9.8.1735>, <https://doi.org/10.1162/neco.1997.9.8.1735>, 1997.
- Hochreiter, S., Bengio, Y., Frasconi, P., Schmidhuber, J., et al.: Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- Hocke, K. and Kämpfer, N.: Gap filling and noise reduction of unevenly sampled data by means of the Lomb-Scargle periodogram, *Atmospheric Chemistry and Physics Discussions*, 8, 4603–4623, 2008.
- 30 Hyndman, R. J. and Khandakar, Y.: Automatic time series forecasting: the forecast package for R, *Journal of statistical software*, 27, 1–22, 2008.
- Hyndman, R. J., Khandakar, Y., et al.: Automatic time series for forecasting: the forecast package for R, 6/07, Monash University, Department of Econometrics and Business Statistics, 2007.
- 35 JORDAN, M.: Attractor dynamics and parallelism in a connectionist sequential machine, *Proc. of the Eighth Annual Conference of the Cognitive Science Society* (Erlbaum, Hillsdale, NJ), 1986, <https://ci.nii.ac.jp/naid/10018634949/en/>, 1986.
- Kamarianakis, Y. and Prastacos, P.: Forecasting traffic flow conditions in an urban network: Comparison of multivariate and univariate approaches, *Transportation Research Record: Journal of the Transportation Research Board*, pp. 74–84, 2003.

- Kamarianakis, Y. and Prastacos, P.: Space–time modeling of traffic flow, *Computers & Geosciences*, 31, 119–133, 2005.
- Katzfuss, M. and Cressie, N.: Spatio-temporal smoothing and EM estimation for massive remote-sensing data sets, *Journal of Time Series Analysis*, 32, 430–446, 2011.
- Katzfuss, M. and Cressie, N.: Bayesian hierarchical spatio-temporal smoothing for very large datasets, *Environmetrics*, 23, 94–107, 2012.
- 5 Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, *CoRR*, abs/1412.6980, <http://arxiv.org/abs/1412.6980>, 2014.
- Kondrashov, D. and Ghil, M.: Spatio-temporal filling of missing points in geophysical data sets, *Nonlinear Processes in Geophysics*, 13, 151–159, 2006.
- Kondrashov, D., Shprits, Y., and Ghil, M.: Gap filling of solar wind data by singular spectrum analysis, *Geophysical research letters*, 37, 2010.
- 10 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., and Herrnegger, M.: Rainfall–runoff modelling using Long Short-Term Memory (LSTM) networks, *Hydrology and Earth System Sciences*, 22, 6005–6022, <https://doi.org/10.5194/hess-22-6005-2018>, <https://www.hydrol-earth-syst-sci.net/22/6005/2018/>, 2018.
- Lin, C. Y., Abdullah, M. H., Praveena, S. M., Yahaya, A. H. B., and Musta, B.: Delineation of temporal variability and governing factors influencing the spatial variability of shallow groundwater chemistry in a tropical sedimentary island, *Journal of hydrology*, 432, 26–42, 2012.
- 15 Långkvist, M., Karlsson, L., and Loutfi, A.: A review of unsupervised feature learning and deep learning for time-series modeling, *Pattern Recognition Letters*, 42, 11 – 24, <https://doi.org/https://doi.org/10.1016/j.patrec.2014.01.008>, <http://www.sciencedirect.com/science/article/pii/S0167865514000221>, 2014.
- Nash, J. and Sutcliffe, J.: River flow forecasting through conceptual models part I — A discussion of principles, *Journal of Hydrology*, 10, 282–290, [https://doi.org/https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/https://doi.org/10.1016/0022-1694(70)90255-6), <https://www.sciencedirect.com/science/article/pii/0022169470902556>, 1970.
- 20 Pfeifer, P. E. and Deutch, S. J.: A three-stage iterative procedure for space-time modeling phillip, *Technometrics*, 22, 35–47, 1980.
- Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and Prabhat: Deep learning and process understanding for data-driven Earth system science, *Nature*, 566, 195–204, <https://doi.org/10.1038/s41586-019-0912-1>, <http://www.nature.com/articles/s41586-019-0912-1>, 2019.
- 25 Schmidhuber, J.: Deep learning in neural networks: An overview, *Neural Networks*, 61, 85 – 117, <https://doi.org/https://doi.org/10.1016/j.neunet.2014.09.003>, <http://www.sciencedirect.com/science/article/pii/S0893608014002135>, 2015.
- Shen, C.: A Transdisciplinary Review of Deep Learning Research and Its Relevance for Water Resources Scientists, *Water Resources Research*, 54, 8558–8593, <https://doi.org/10.1029/2018WR022643>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018WR022643>, 2018.
- 30 Shuai, P., Chen, X., Song, X., Hammond, G. E., Zachara, J., Royer, P., Ren, H., Perkins, W. A., Richmond, M. C., and Huang, M.: Dam Operations and Subsurface Hydrogeology Control Dynamics of Hydrologic Exchange Flows in a Regulated River Reach, *Water Resources Research*, 55, 2593–2612, <https://doi.org/10.1029/2018WR024193>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018WR024193>, 2019.
- 35 Song, X., Chen, X., Stegen, J., Hammond, G., Song, H.-S., Dai, H., Graham, E., and Zachara, J. M.: Drought Conditions Maximize the Impact of High-Frequency Flow Variations on Thermal Regimes and Biogeochemical Function in the Hyporheic Zone, *Water Resources Research*, 2018.

- Stockwell, R. G., Mansinha, L., and Lowe, R.: Localization of the complex spectrum: the S transform, *IEEE transactions on signal processing*, 44, 998–1001, 1996.
- Strobl, R. O. and Robillard, P. D.: Network design for water quality monitoring of surface freshwaters: A review, *Journal of environmental management*, 87, 639–648, 2008.
- 5 Stroud, J. R., Stein, M. L., and Lysen, S.: Bayesian and maximum likelihood estimation for Gaussian processes on an incomplete lattice, *Journal of computational and Graphical Statistics*, 26, 108–120, 2017.
- Sun, A. Y.: Discovering State-Parameter Mappings in Subsurface Models Using Generative Adversarial Networks, *Geophysical Research Letters*, 45, 11,137–11,146, <https://doi.org/10.1029/2018GL080404>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018GL080404>, 2018.
- 10 Sun, A. Y., Scanlon, B. R., Zhang, Z., Walling, D., Bhanja, S. N., Mukherjee, A., and Zhong, Z.: Combining Physically Based Modeling and Deep Learning for Fusing GRACE Satellite Data: Can We Learn From Mismatch?, *Water Resources Research*, 55, 1179–1195, <https://doi.org/10.1029/2018WR023333>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018WR023333>, 2019.
- Taylor, C. J. and Alley, W. M.: Ground-water-level monitoring and the importance of long-term water-level data, 1217-2002, *US Geological Survey*, 2002.
- 15 Vacha, L. and Barunik, J.: Co-movement of energy commodities revisited: Evidence from wavelet coherence analysis, *Energy Economics*, 34, 241–247, 2012.
- Valenzuela, O., Rojas, I., Rojas, F., Pomares, H., Herrera, L. J., Guillén, A., Marquez, L., and Pasadas, M.: Hybridization of intelligent techniques and ARIMA models for time series prediction, *Fuzzy sets and systems*, 159, 821–845, 2008.
- Wang, G., Garcia, D., Liu, Y., De Jeu, R., and Dolman, A. J.: A three-dimensional gap filling method for large geophysical datasets: Application to global satellite soil moisture observations, *Environmental Modelling & Software*, 30, 139–142, 2012.
- 20 Wett, B., Jarosch, H., and Ingerle, K.: Flood induced infiltration affecting a bank filtrate well at the River Enns, Austria, *Journal of Hydrology*, 266, 222–234, 2002.
- Wikle, C. K., Berliner, L. M., and Cressie, N.: Hierarchical Bayesian space-time models, *Environmental and Ecological Statistics*, 5, 117–154, 1998.
- 25 Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Łukasz Kaiser, Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., and Dean, J.: Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation, *CoRR*, abs/1609.08144, <http://arxiv.org/abs/1609.08144>, 2016.
- You, Q., Jin, H., Wang, Z., Fang, C., and Luo, J.: Image captioning with semantic attention, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4651–4659, 2016.
- 30 Zachara, J. M., Long, P. E., Bargar, J., Davis, J. A., Fox, P., Fredrickson, J. K., Freshley, M. D., Konopka, A. E., Liu, C., McKinley, J. P., et al.: Persistence of uranium groundwater plumes: contrasting mechanisms at two DOE sites in the groundwater–river interaction zone, *Journal of contaminant hydrology*, 147, 45–72, 2013.
- Zachara, J. M., Chen, X., Song, X., Shuai, P., Murray, C., and Resch, C. T.: Kilometer-scale hydrologic exchange flows in a gravel-bed river corridor and their implications to solute migration, *Water Resources Research*, n/a, e2019WR025 258, <https://doi.org/10.1029/2019WR025258>, <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019WR025258>, e2019WR025258, 2019WR025258, 2020.

Zhang, D., Lindholm, G., and Ratnaweera, H.: Use long short-term memory to enhance Internet of Things for combined sewer overflow monitoring, *Journal of Hydrology*, 556, 409 – 418, <https://doi.org/https://doi.org/10.1016/j.jhydrol.2017.11.018>, <http://www.sciencedirect.com/science/article/pii/S0022169417307722>, 2018.

Zhang, G. P.: Time series forecasting using a hybrid ARIMA and neural network model, *Neurocomputing*, 50, 159–175, 2003.