

Supplementary Materials for paper ‘Exploring the generalisation ability and interpretability of Long Short-Term Memory (LSTM) networks for large-sample groundwater level predictions’

Qidong Fang¹, Mostaqimur Rahman¹, Thorsten Wagener², Francesca Pianosi¹

¹School of Civil, Aerospace and Design Engineering, University of Bristol, Bristol, BS8 1US, UK

²Institute of Environmental Science and Geography, University of Potsdam, Potsdam, 14476, Germany

Correspondence to: Qidong Fang (Qidong.fang@bristol.ac.uk)

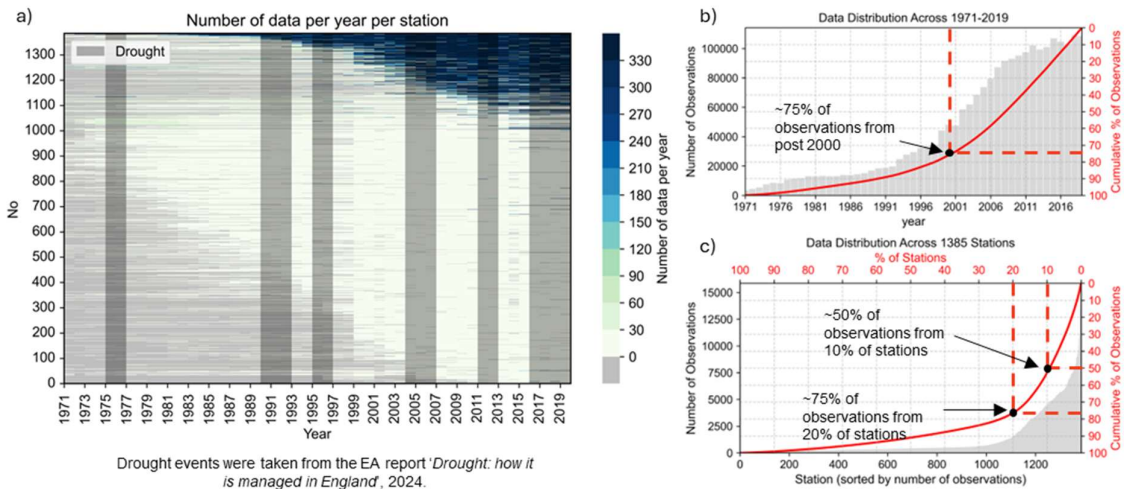


Fig. S1. Overview of the groundwater dataset over time and across stations. (a) Distribution of groundwater observations over time across 1384 stations. Colour indicates the by number of observations per year, while light grey denotes years with no data. Dark grey shading marks the drought events that occurred during the 1975 - 2019, based on records from the Environment Agency (<https://www.gov.uk/government/publications/drought-management-for-england/drought-how-it-is-managed-in-england#drought-in-england-an-overview>). (b) Distribution of observations by year. (c) Distribution of observations by station.

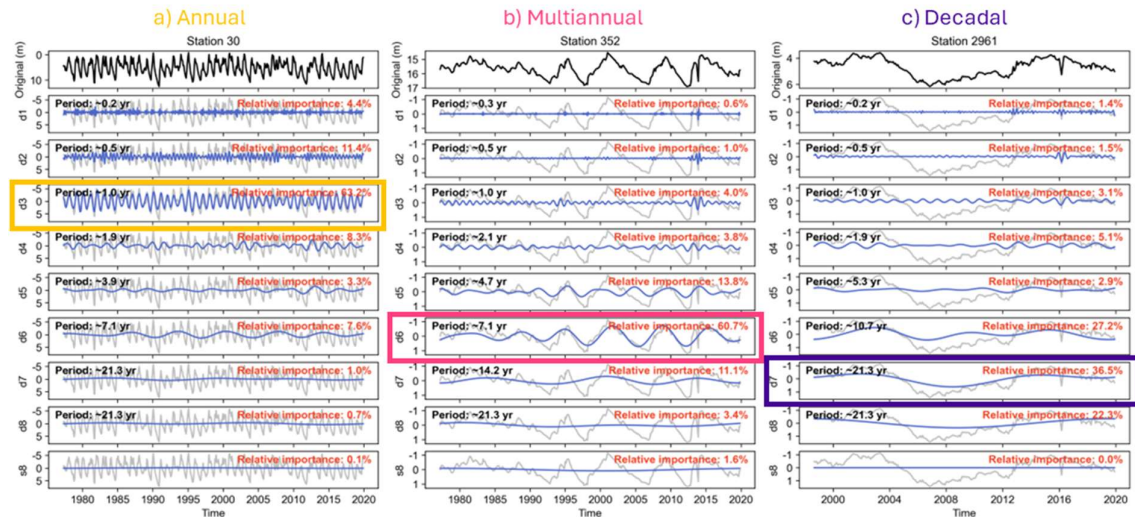


Fig. S2. Examples of wavelet transform decomposition of monthly groundwater level timeseries. (a) Example of a station dominated by an annual cycle. (b) Example of a station dominated by a multi-annual cycle.

(c) Station dominated by a decadal cycle. In each panel, the black line at the top and the grey lines show the original groundwater timeseries. The blue lines represent the wavelet detail components (d1-d8) and the final approximation component (s8). The period (in black) and relative importance (in red) associated with each detail can be calculated. The bold coloured box highlights the dominant fluctuation period for each example.

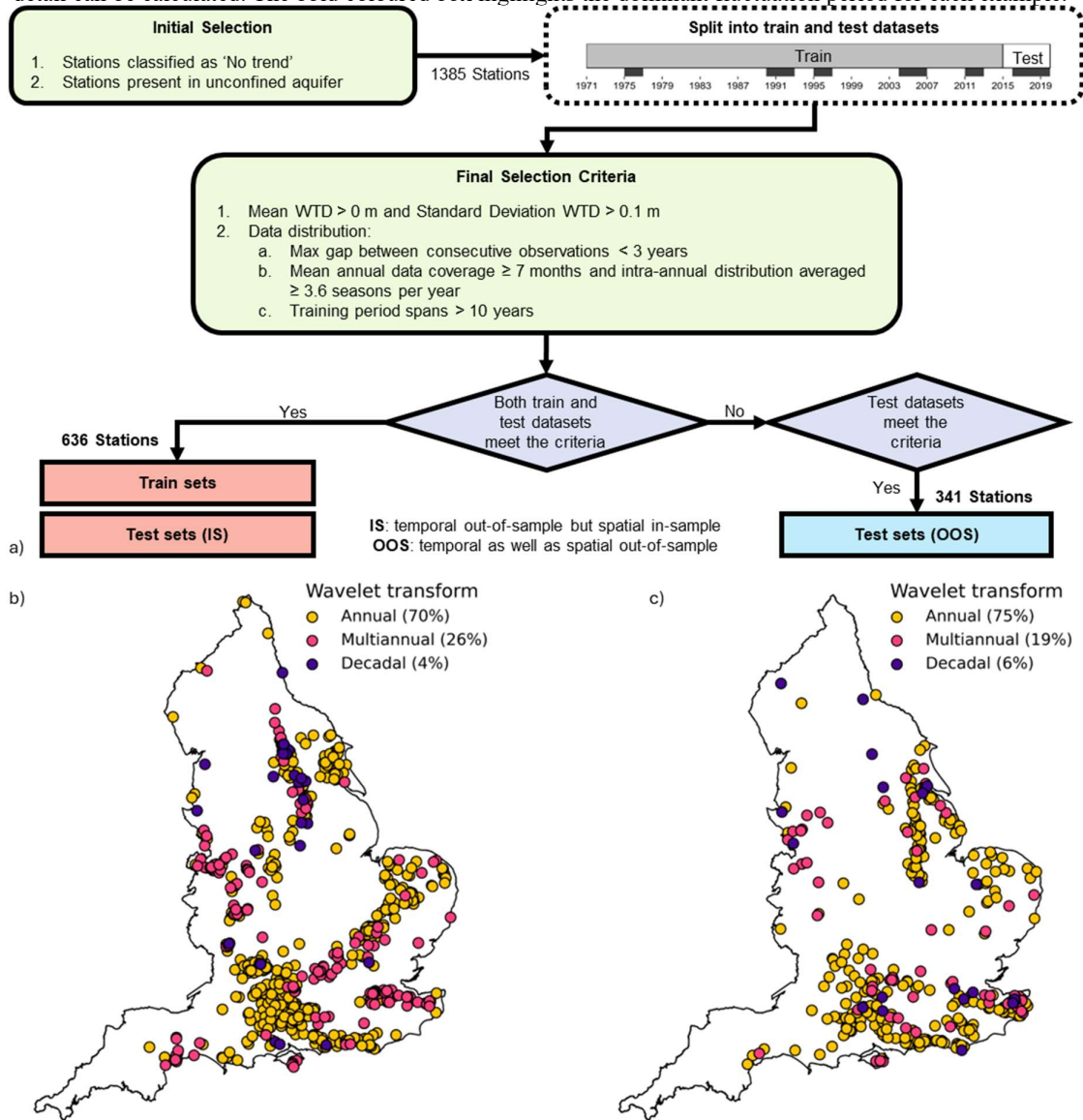


Fig. S3. Schematic of groundwater dataset construction for model training and in-sample / out-of-sample testing, and spatial distribution of stations. (a) Workflow used to construct the groundwater dataset for training, in-sample testing, and out-of-sample testing. (b) Map of groundwater stations included in the train and in-sample test sets. (c) Map of groundwater stations included in the out-of-sample test set. Stations in (b) and (c) are classified based on wavelet transform analysis, indicating whether groundwater level is primarily influenced by annual (yellow), multi-annual (pink) or decadal (purple) cycles.

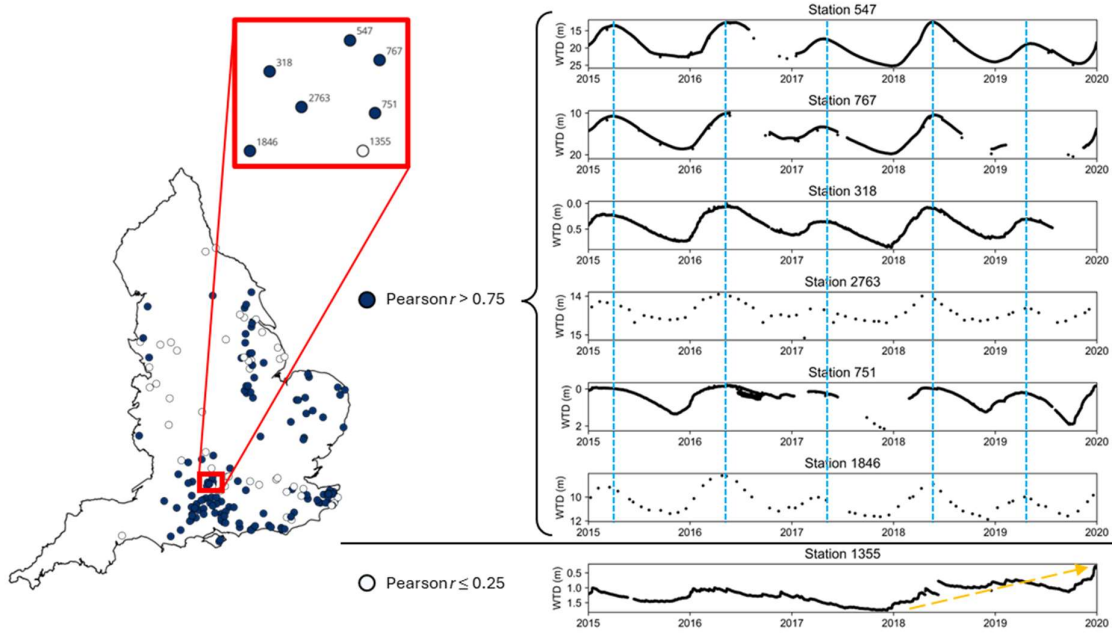


Fig. S4. Representative nearby stations showing opposite performance extremes: Pearson $r \leq 0.25$ (poor performance) and Pearson $r > 0.75$ (good performance, predictable). The map and zoomed-in area (red rectangle) show the locations of all qualified stations and seven example stations. Dark circles represent stations achieving good performance, and white circles represent stations achieving poor performance. In the right panel (top six), we showed six timeseries for stations achieving good performance (stations 547, 767, 318, 2763, 751, and 1846), which exhibit clear seasonal variation (blue dashed lines denote the approximate annual peaks of predictable timeseries). In the right panel bottom, we showed one timeseries for a station achieving poor performance (station 1355), where the station does not exhibit clear seasonal variation and even shows an increase during 2018 – 2020 (yellow arrow).

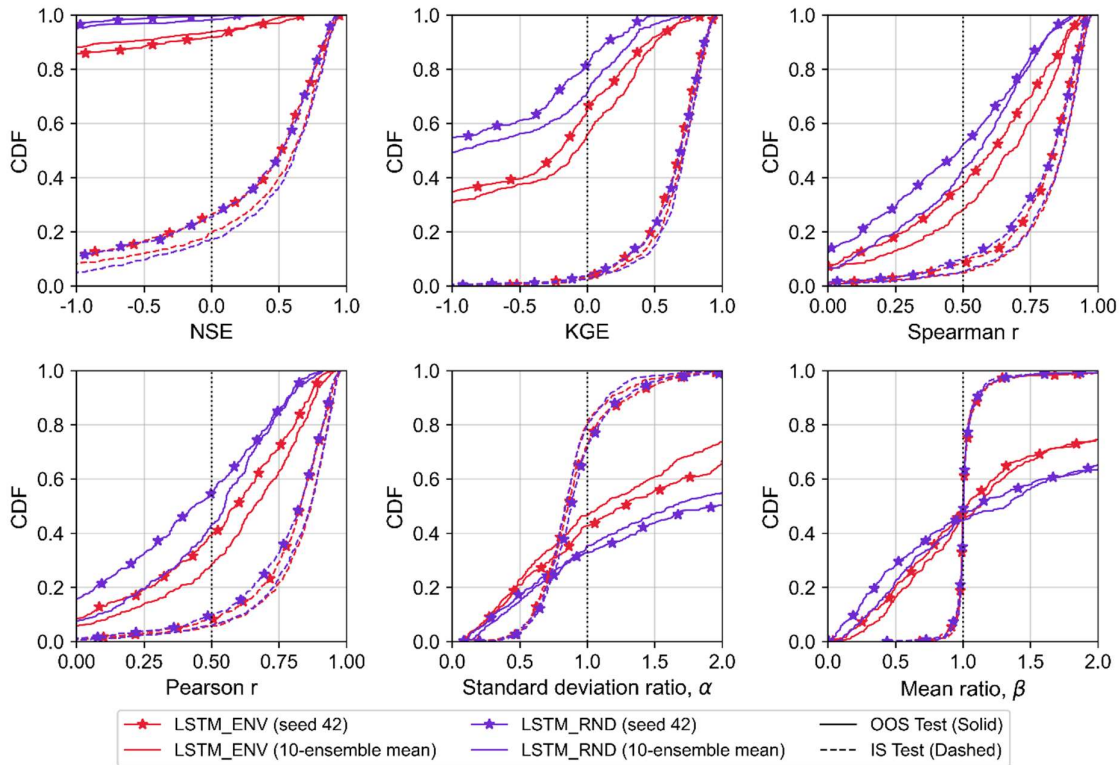


Fig. S5. Cumulative distribution functions (CDFs) of performance metrics for LSTM_ENV (red) and LSTM_RND (purple). Comparisons are shown between the 10-ensemble mean (no markers) and one of the

ensemble member (star markers). Performances is shown for 636 in-sample stations (dashed lines) and 341 out-of-sample stations (i.e., stations not used for the LSTM training, solid lines). For NSE, KGE, Spearman r , and Pearson r , superior performance is indicated by curves shifted toward the bottom-right. For the Standard deviation ratio (α) and Mean ratio (β), proximity to the vertical (ratio=1) indicates higher accuracy.

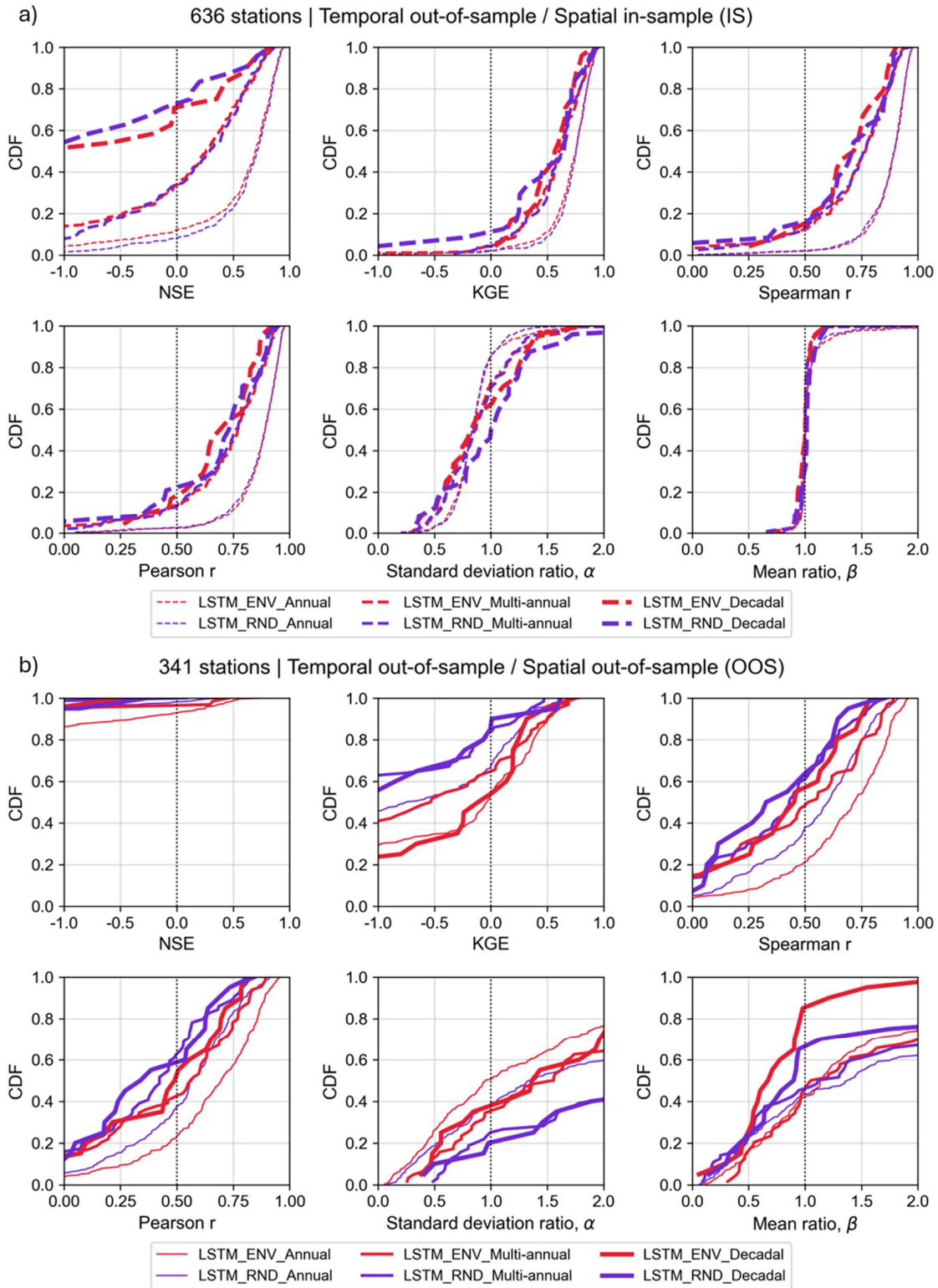


Fig. S6. Cumulative distribution functions (CDFs) of performance metrics for the 10-ensemble mean of LSTM_ENV (red) and LSTM_RND (purple), categorised by dominant variation cycle (i.e. annual, multi-

annual, and decadal) of the stations. Performances is shown for 636 in-sample stations (a) and 341 out-of-sample stations (b). For NSE, KGE, Spearman r , and Pearson r , superior performance is indicated by curves shifted toward the bottom-right. For the Standard deviation ratio (α) and Mean ratio (β), proximity to the vertical (ratio=1) indicates higher accuracy.

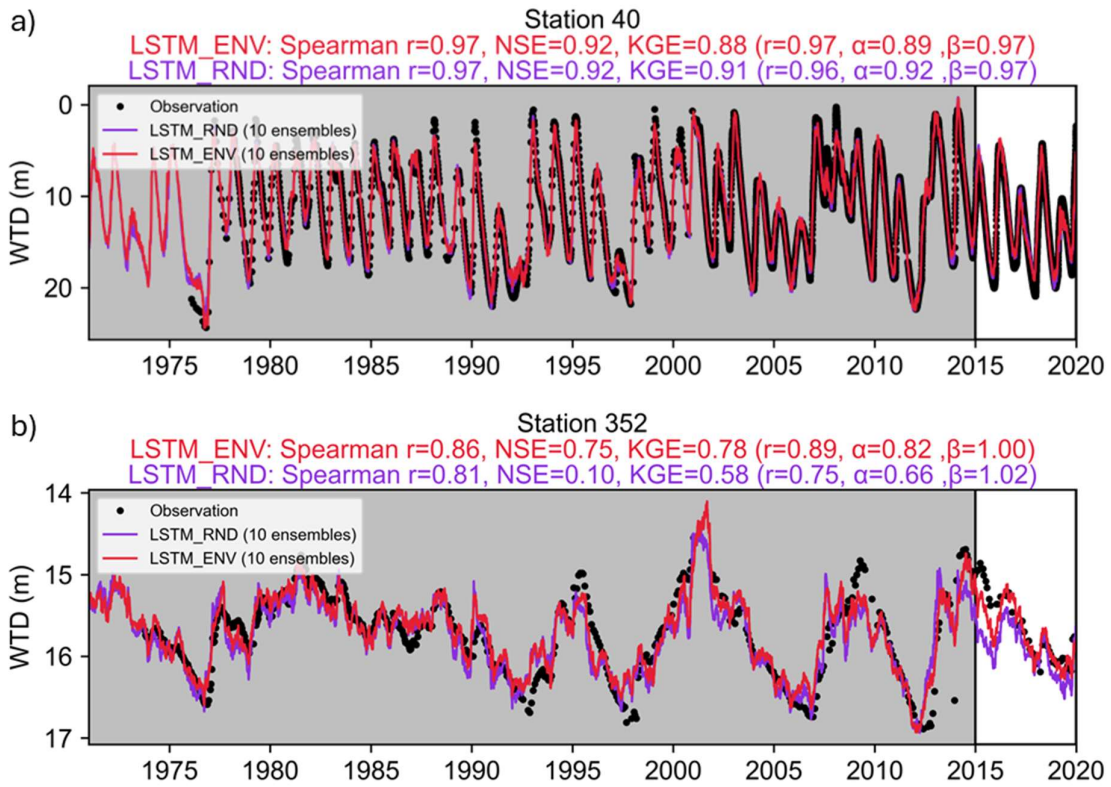


Fig. S7. Representative water table depth (WTD) timeseries for spatially in-sample stations. (a) Station 40, exhibiting a dominant annual cycle, and (b) Station 352, exhibiting a dominant multi-annual cycle. Model simulations and performance metrics for LSTM_ENV and LSTM_RND are shown in red and purple, respectively. Black dots represent observations. The grey background indicates the training period (1971-2014), followed by the tested period (2015-2019).

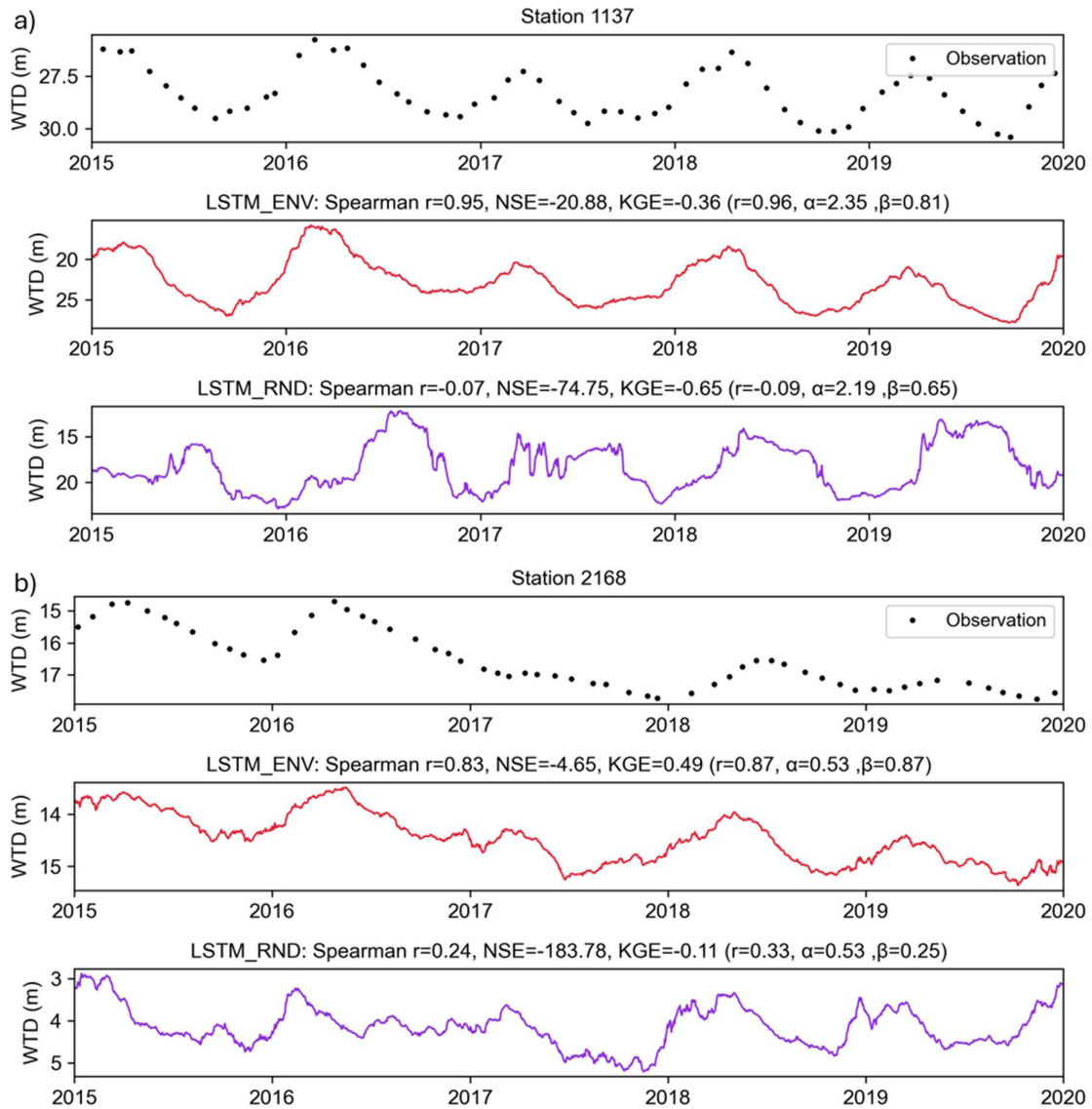


Fig. S8. Representative water table depth (WTD) timeseries for spatially out-of-sample stations. (a) Station 1137, exhibiting a dominant annual cycle, and (b) Station 2168, exhibiting a dominant multi-annual cycle according to wavelet transform analysis. Model simulations and performance metrics for LSTM_ENV and LSTM_RND are shown in red and purple, respectively. Black dots represent observations. Simulations and observations are plotted on separate vertical axes to better compare temporal dynamics, as the means and ranges of observations and simulations differ substantially.

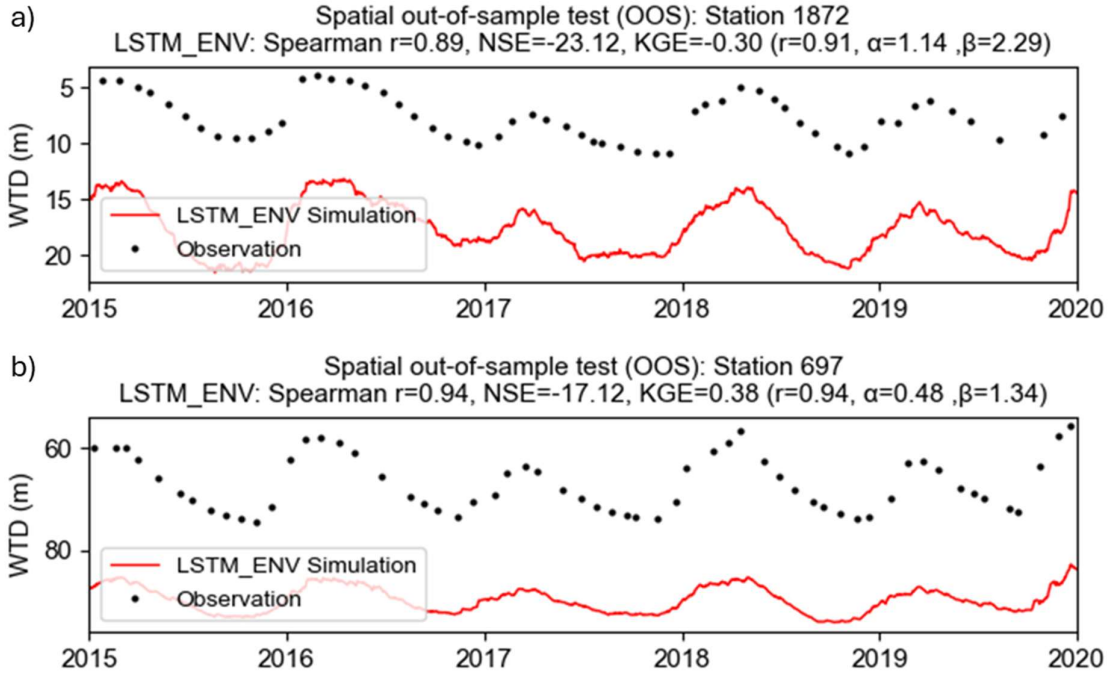


Fig. S9. Representative water table depth (WTD) stations that achieved high correlation but negative NSE. (a) Station 1872, and (b) Station 697.

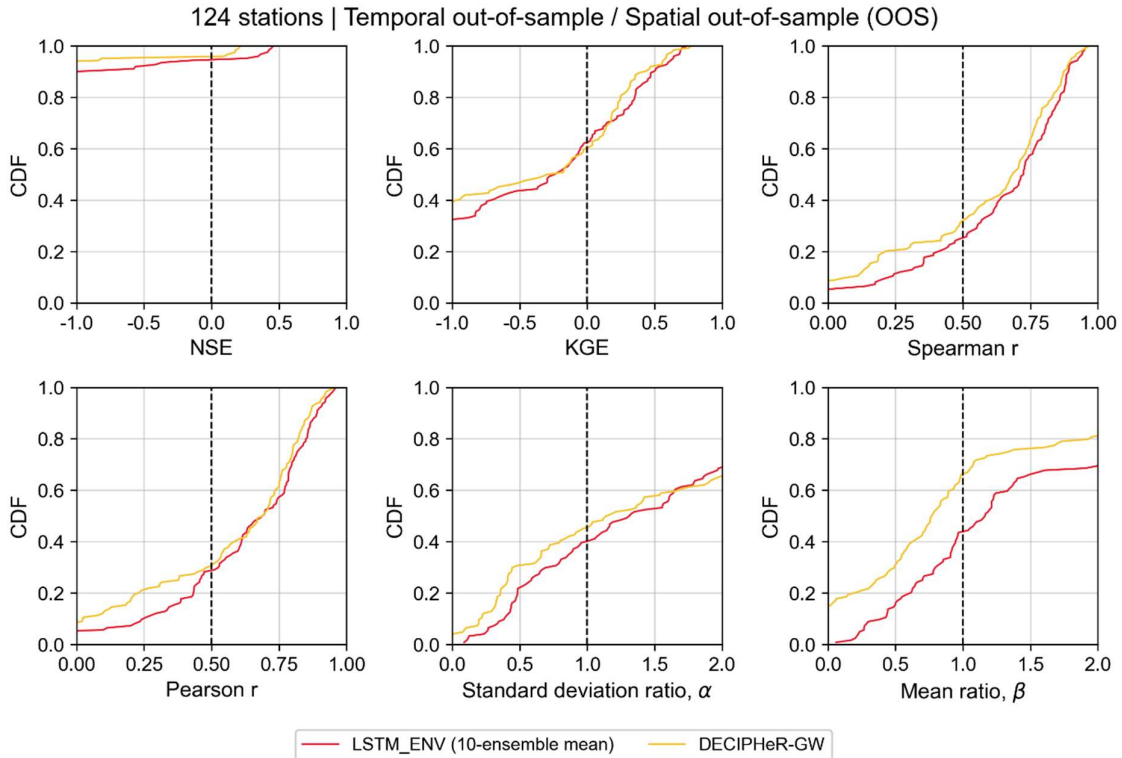


Fig. S10. Cumulative distribution functions (CDFs) of performance metrics for LSTM_ENV (red) and DECIPHER-GW (yellow). CDFs compile the results from 124 out-of-sample stations. For NSE, KGE,

Spearman r , and Pearson r , superior performance is indicated by curves shifted toward the bottom-right. For Standard deviation ratio (α) and mean ratio (β), proximity to the vertical (ratio=1) indicates higher accuracy.

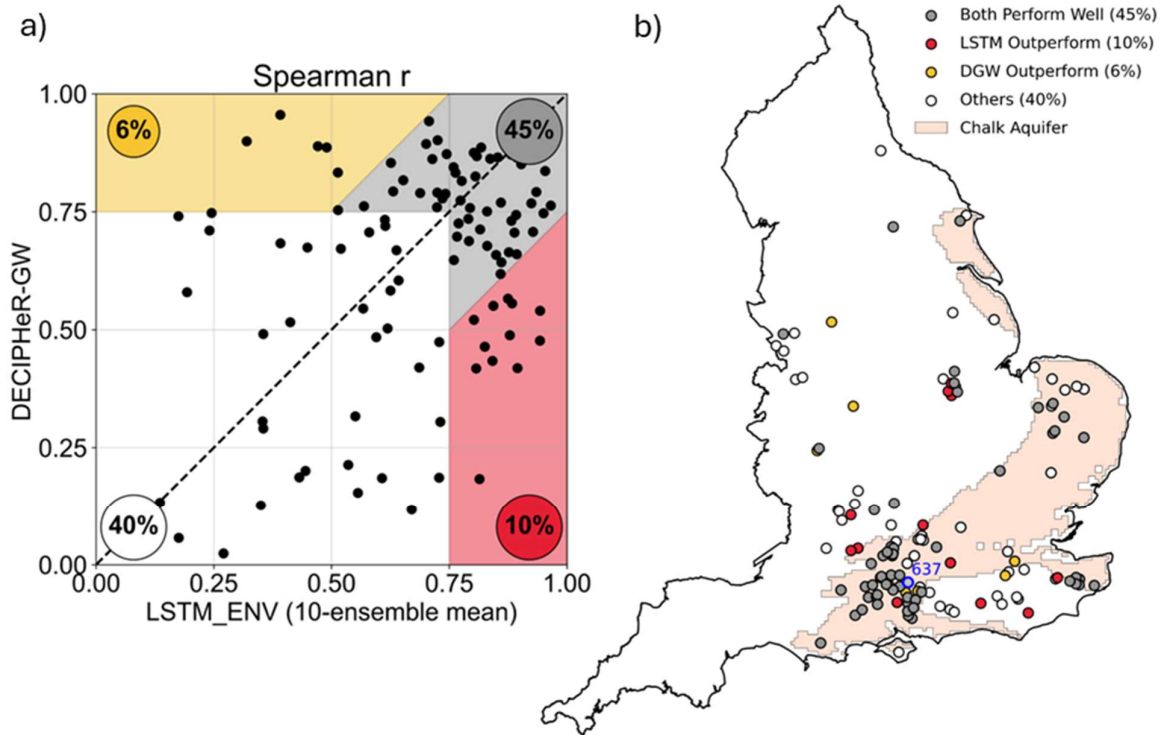


Fig. S11. Performance comparison (Spearman correlation coefficient) between LSTM_ENV (red) and DECIPHeR-GW (yellow) across 124 common stations. (a) Comparison of Spearman correlation coefficient. (b) Spatial distribution of performance classes. Stations (circles) are categorised by their relative Spearman r : grey indicates high performance for both models (Spearman r of one model > 0.75 and the difference between models $\Delta r < 0.25$); yellow indicates that DECIPHeR-GW outperforms LSTM_ENV, red indicates that LSTM_ENV outperforms DECIPHeR-GW, and white represents all remaining cases. The Chalk aquifer is shown in the background.

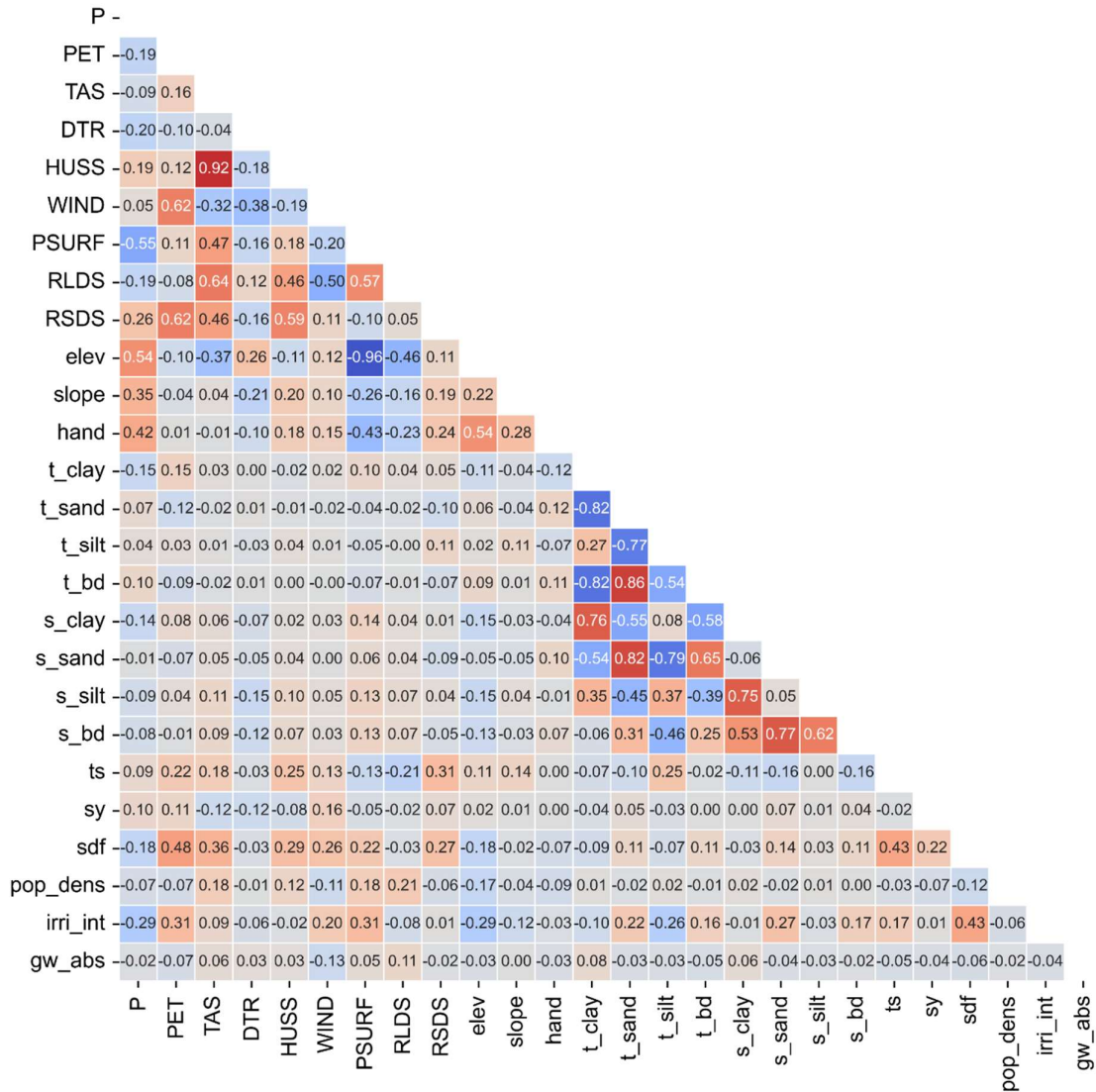
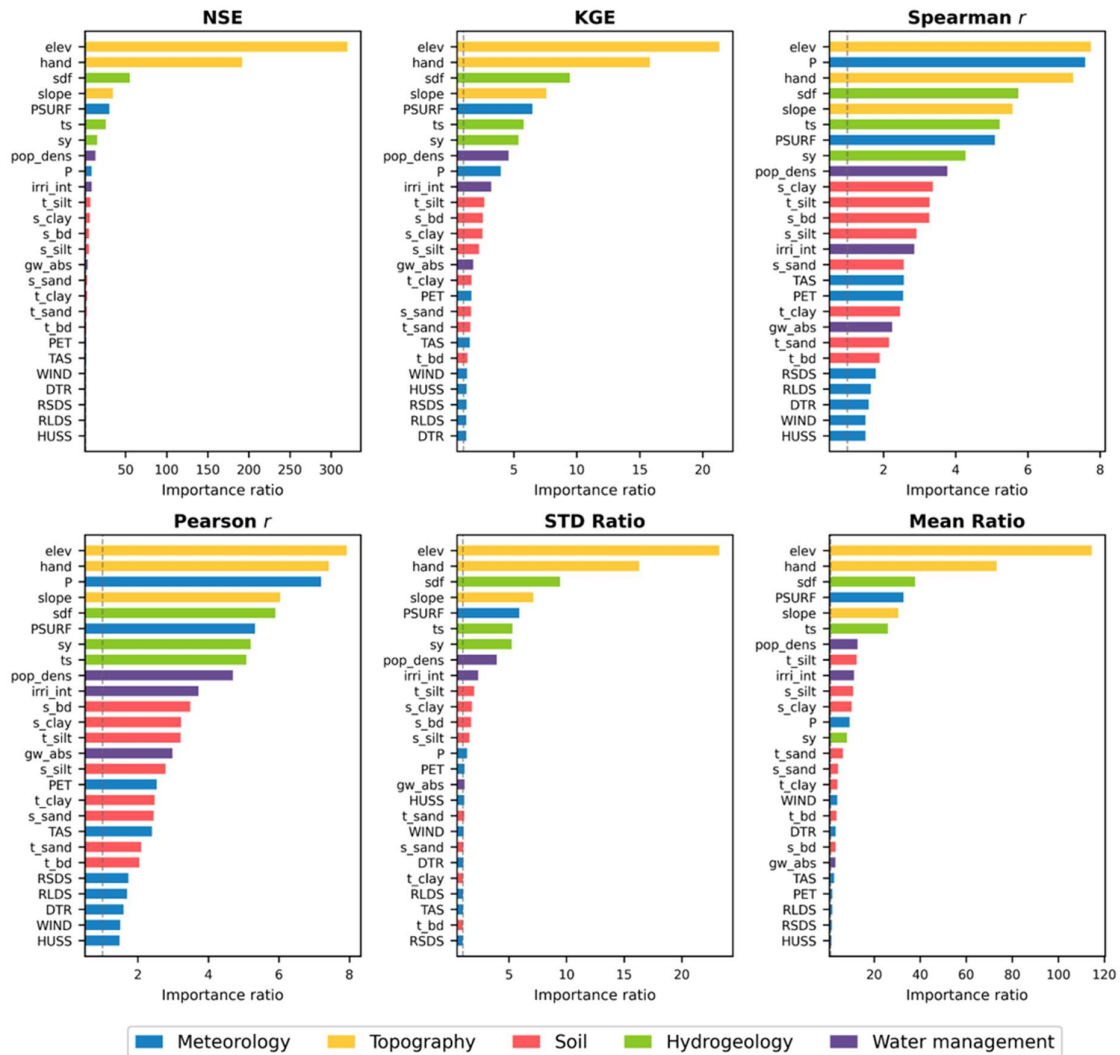


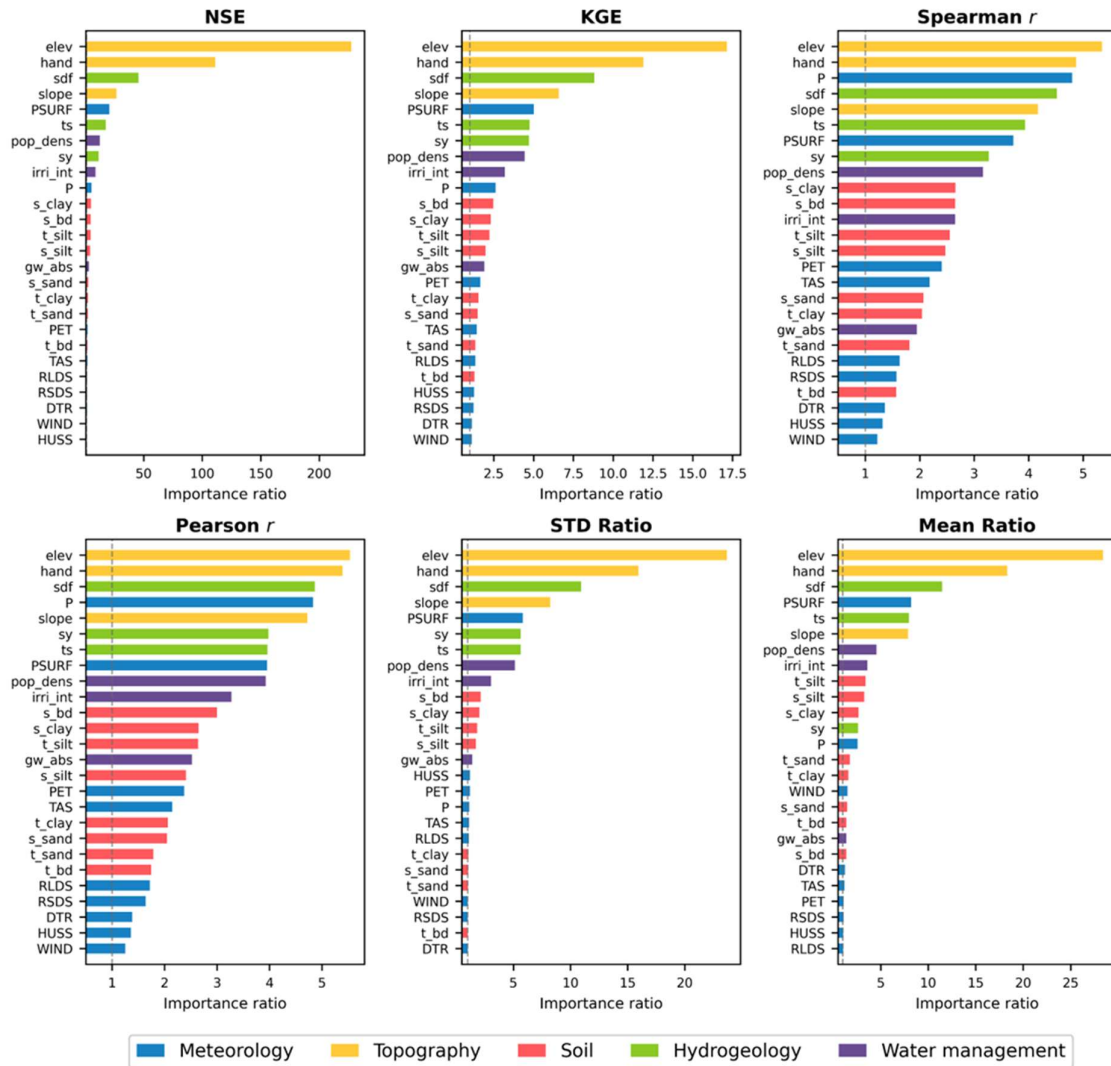
Fig. S12. Correlation between features (i.e. dynamic inputs averaged over 1961–2019, UPPERCASE, and static inputs, lowercase). Pearson correlations are calculated based on 636 stations used for training and 341 unseen stations during training (977 stations in total). Value close to 1 and -1 represent stronger positive (red) and negative (blue) relationships, respectively.

Fig. S13. (Figure across three pages) Individual feature importance for LSTM_ENV on (a) the training set, (b) the in-sample test set, and (c) the out-of-sample test set. Dynamic inputs are timeseries of meteorological variables (blue, UPPERCASE). Static inputs (lowercase) are classified into four categories: topography (yellow), soil (red), hydrogeology (green), and water management (purple).

a) LSTM_ENV (10-ensemble mean) – Training set



b) LSTM_ENV (10-ensemble mean) – In-sample test set



c) LSTM_ENV(10-ensemble mean) – Out-of-sample test set

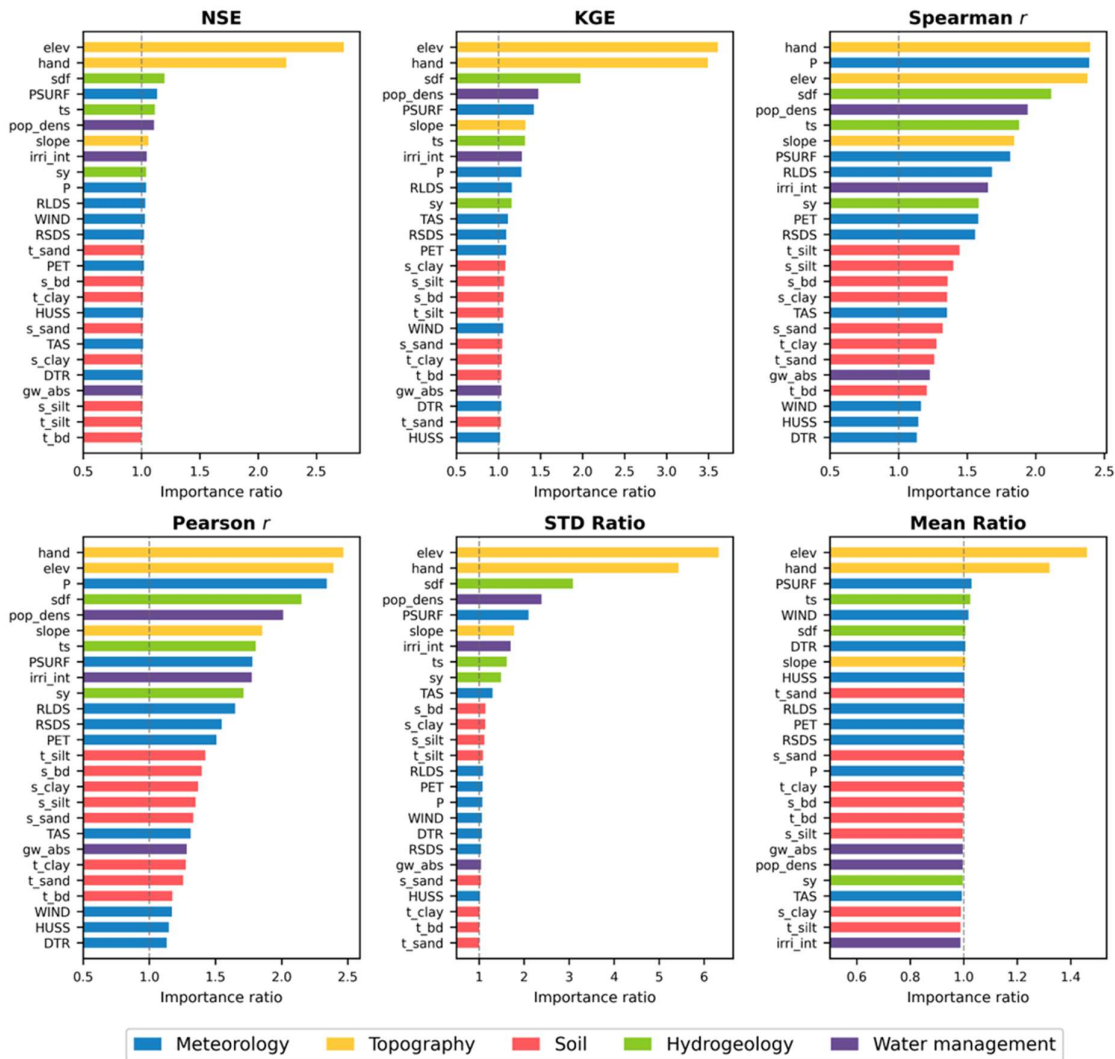
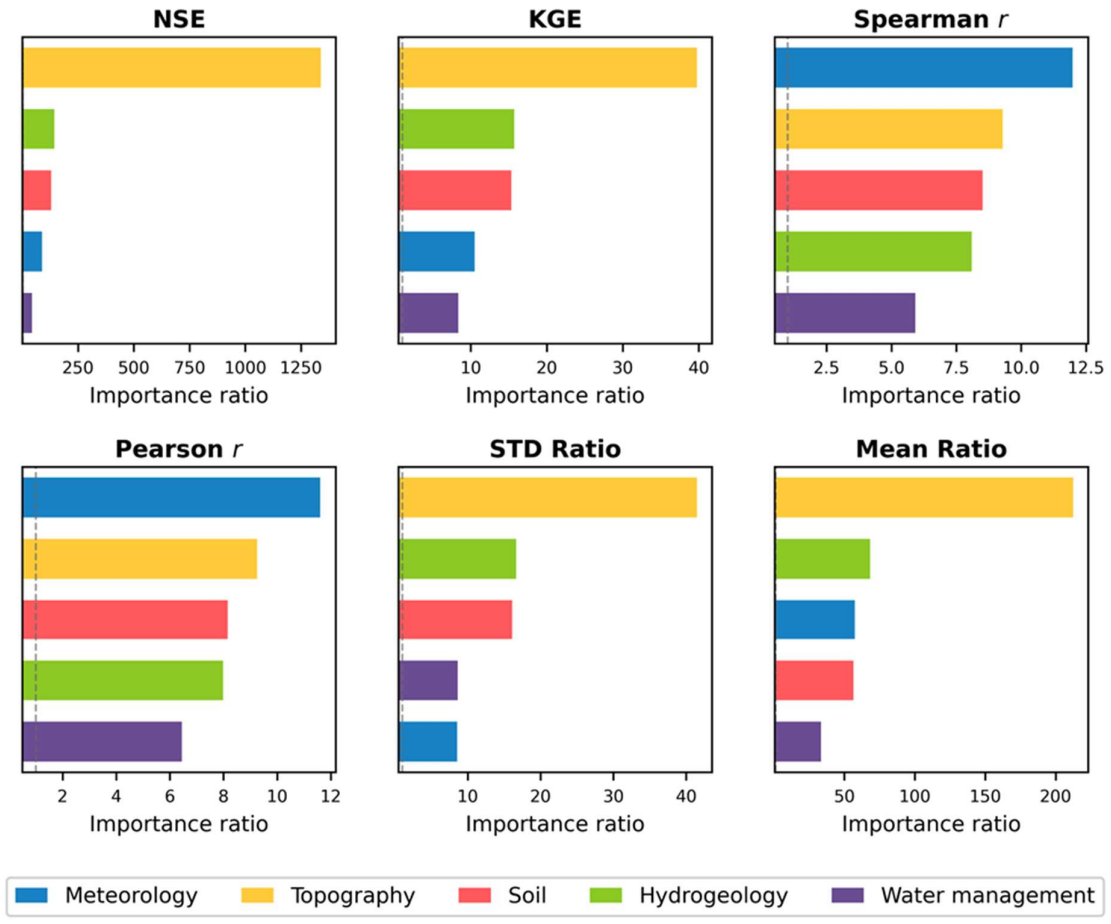
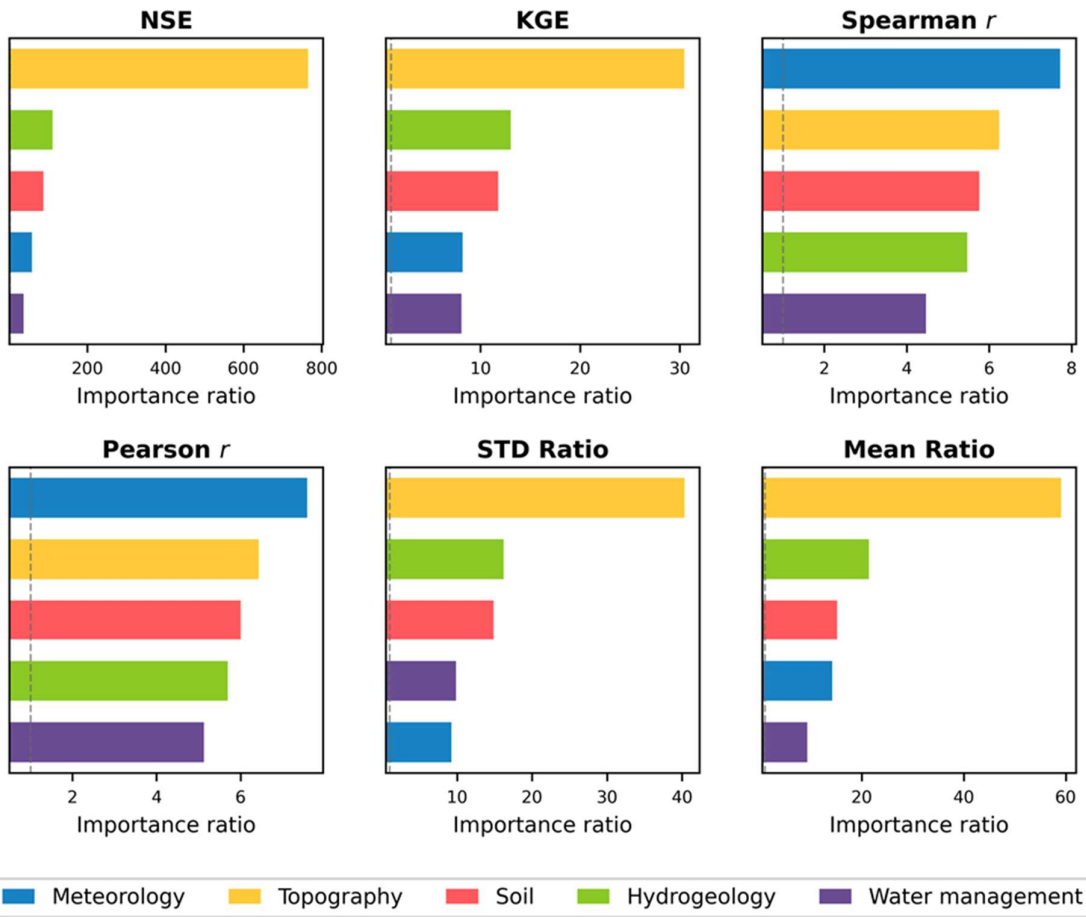


Fig. S14. (Figure across three pages) Grouped feature importance for LSTM_ENV on (a) the training set, (b) the in-sample test set, and (c) the out-of-sample test set. Dynamic inputs are timeseries of meteorological variables (blue, UPPERCASE). Static inputs (lowercase) are classified into four categories: topography (yellow), soil (red), hydrogeology (green), and water management (purple).

a) LSTM_ENV (10-ensemble mean) – Training set



b) LSTM_ENV(10-ensemble mean) – In-sample test set



c) LSTM_ENV (10-ensemble mean) – Out-of-sample test set

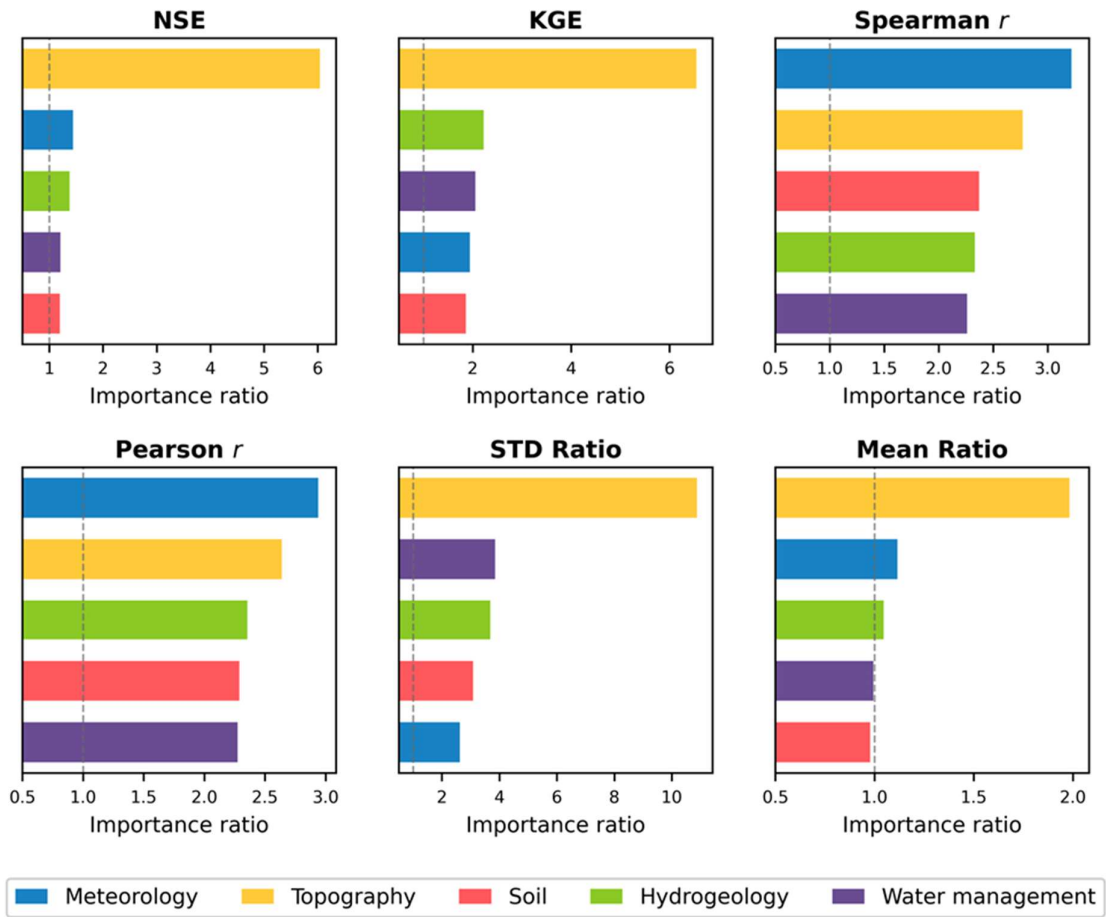
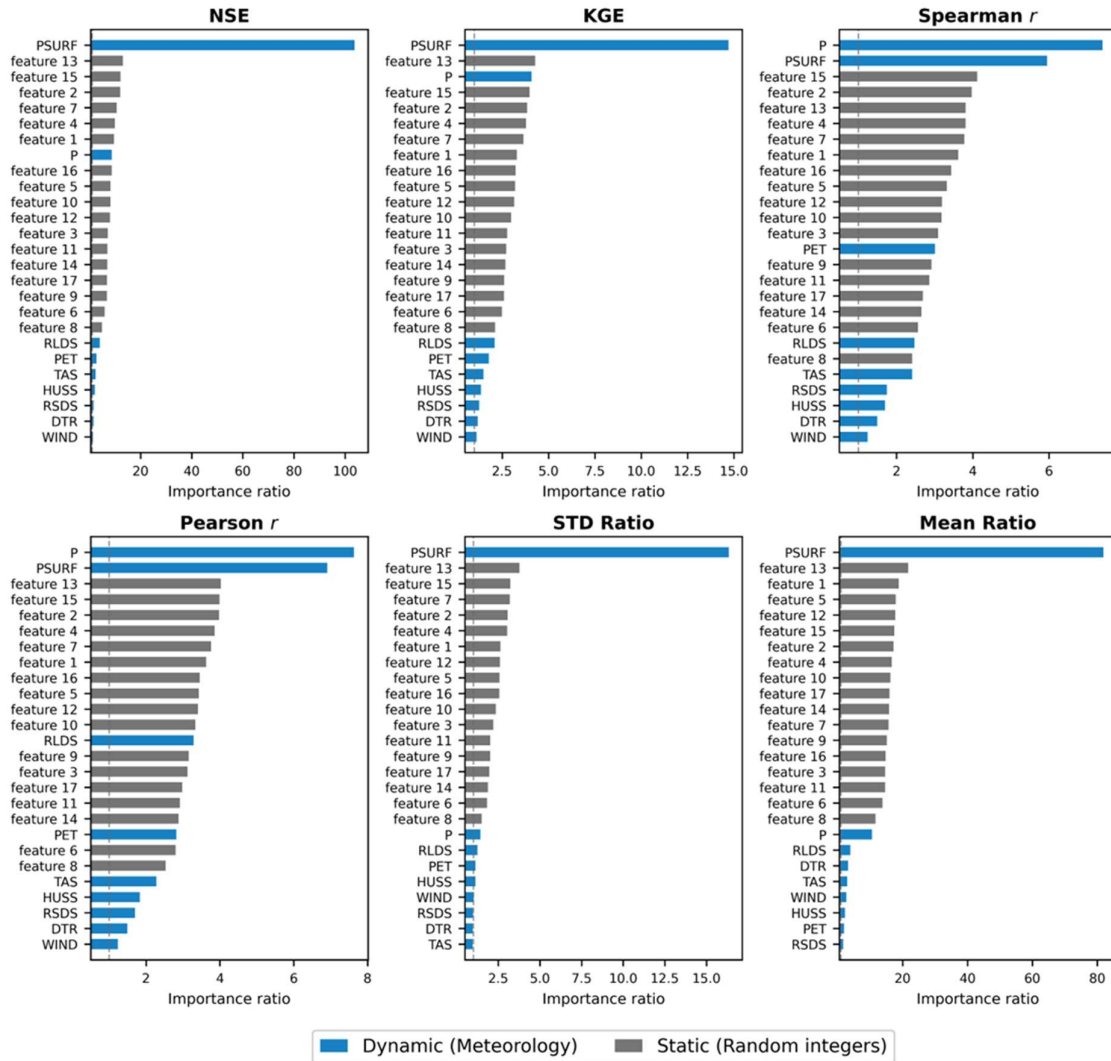
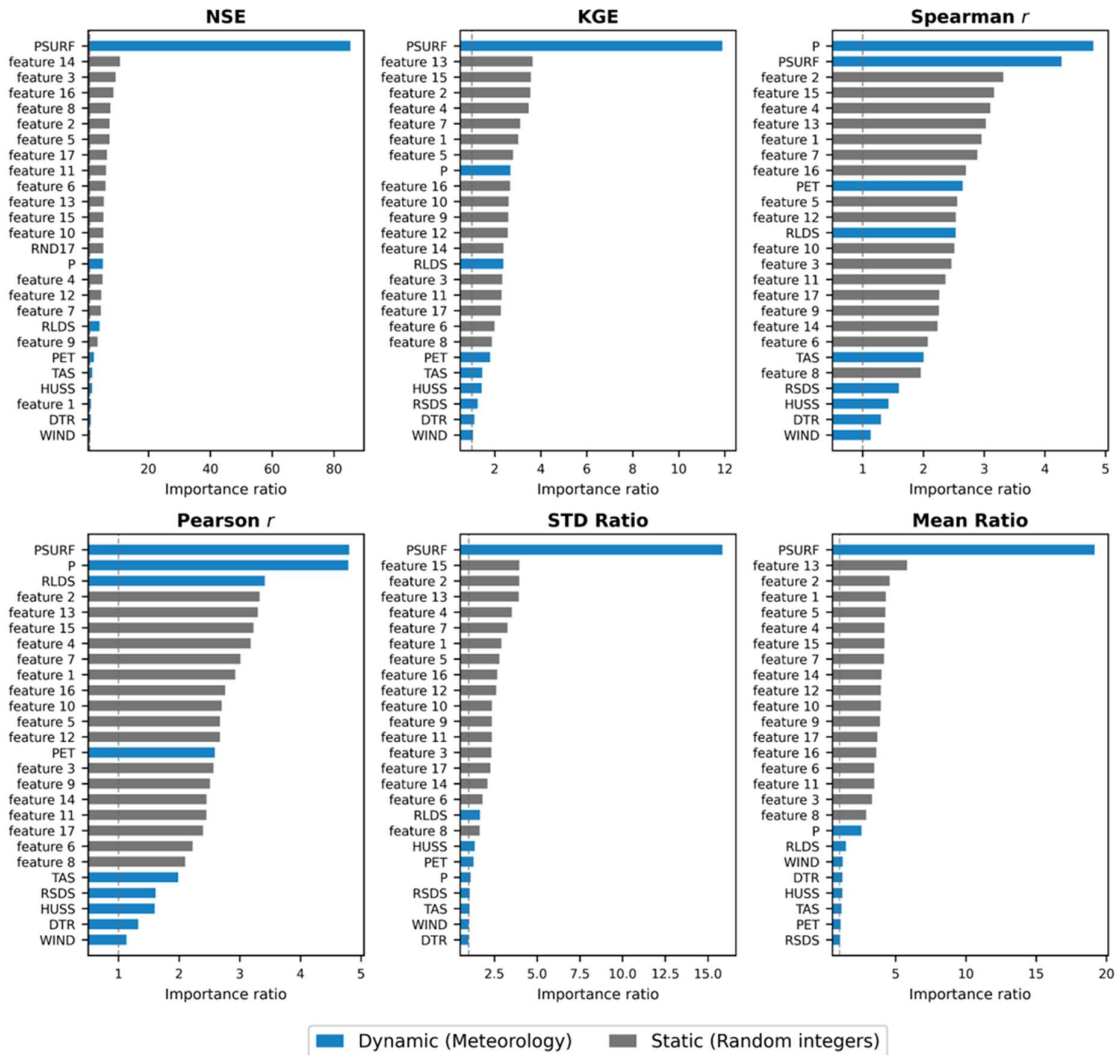


Fig. S15. (Figure across three pages) Individual feature importance for LSTM_RND on (a) the training set, (b) the in-sample test set, and (c) the out-of-sample test set. Dynamic inputs are timeseries of meteorological variables (blue, UPPERCASE). Static inputs (grey) are random integers.

a) LSTM_RND (10-ensemble mean) – Training set



b) LSTM_RND (10-ensemble mean) – In-sample test set



c) LSTM_RND(10-ensemble mean) – Out-of-sample test set

