

**1 Leveraging gauge networks and strategic discharge measurements to aid development of
2 continuous streamflow records**

3

4

5 Michael J. Vlah¹, Matthew R. V. Ross², Spencer Rhea¹, Emily S. Bernhardt¹

6

7 ¹Duke University

8 ²Colorado State University

9

10 *Correspondence to:* Michael J Vlah (michael.vlah@duke.edu)

11

12 Grant sponsor: National Science Foundation, MacroSystems, 1926420

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45 Abstract

46

47 Quantifying continuous discharge can be difficult, especially for nascent monitoring efforts, due to the
48 challenges of establishing gauging locations, sensor protocols, and installations. Some continuous
49 discharge series generated by the National Ecological Observatory Network (NEON) during its pre- and
50 early-operational phases (2015-present) are marked by anomalies related to sensor drift, gauge movement,
51 and incomplete rating curves. Here, we investigate the potential to estimate continuous discharge when
52 discrete streamflow measurements are available at the site of interest. Using field-measured discharge as
53 truth, we reconstructed continuous discharge for all 27 NEON stream gauges via linear regression on
54 nearby donor gauges and/or prediction from neural networks trained on a large corpus of established
55 gauge data. Reconstructions achieved median efficiencies of 0.83 (Nash-Sutcliffe, or NSE) and 0.81
56 (Kling-Gupta, or KGE) across all sites, and improved KGE at 11 sites versus published data, with linear
57 regression generally outperforming deep learning approaches due to the use of target site data for model
58 fitting, rather than evaluation only. Estimates from this analysis inform ~199 site-months of missing data
59 in the official record, and can be used jointly with NEON data to enhance the descriptive and predictive
60 value of NEON's stream data products. We provide 5-minute composite discharge series for each site that
61 combine the best estimates across modeling approaches and NEON's published data. The success of this
62 effort demonstrates the potential to establish "virtual gauges," or sites at which continuous streamflow can
63 be accurately estimated from discrete measurements, by transferring information from nearby donor
64 gauges and/or large collections of training data.

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89 1. Introduction

90

91 Discharge, or streamflow, is a fundamental measure in hydrology, biogeochemistry, and river science
92 more broadly. A measure of water volume over time, discharge is used to infer theoretical watershed
93 runoff (depth of water “blanketing” the land surface, or depth over time), which in turn is integral to
94 understanding watershed processes such as chemical weathering (White & Blum 1995). Accurate, and at
95 least daily, discharge estimates are essential components of nearly any quantitative study of physical or
96 chemical watershed or river processes at the ecosystem scale. Determination of solute fluxes (Bukaveckas
97 et al. 1998), gas exchange rates (Hall, 2016), ecosystem metabolism (Odum 1956), and sediment transport
98 (Graf 1984) all require well constrained estimates of discharge.

99

100 Despite its centrality to so many fields of study, discharge is a notoriously difficult metric to capture on a
101 regular basis, especially in free-flowing systems, as it may vary greatly with annual cycles and weather
102 events (Turnipseed & Sauer 2010). Established institutions like the USGS (USA), ECCC (Canada), and
103 ANA (Brazil) have honed their instrumentation, methods, and monitoring locations over decades to
104 generate reasonable discharge estimates even under extreme conditions (Benson & Dalrymple 1967;
105 Costa 2004); however, nascent and/or small-budget monitoring efforts face several challenges. Critically,
106 hundreds of these efforts are constantly occurring within academic research groups, municipalities,
107 counties, and other entities building smaller gauge networks, with much less expertise, support, and
108 budget than gauging programs supported by dedicated national programs.

109

110 Not including purely model-based methods for discharge prediction (Manning 1891; Hsu et al. 1995,
111 Durand et al. 2022), automated discharge estimation requires the careful construction of an empirical
112 “rating curve,” by which discharge can be continuously inferred from water level, or “stage” (but see
113 Shen 1981). To build such a relationship, technicians must sample discharge and stage at points covering
114 the range of observable flow, ideally including flood stage. In dynamic systems, this rating curve must be
115 regularly updated. Point estimates of discharge can be collected using Acoustic Doppler current profiling
116 (Moore et al. 2017), manual flow meter profiling, or light-based methods (Wang 1988) to determine
117 average cross-sectional velocity, or via conservative tracer injections (Tazioli 2011). In many streams, two
118 or more of these methods must be employed, depending on conditions (Turnipseed & Sauer 2010).
119 During 10-year or 100-year floods, no method may be viable or safe. Even under regular storm
120 conditions, a technician may be unable to mount a sampling effort quickly enough to capture peak flow,
121 or may produce an inaccurate measurement. As a result, rating curves may remain in a state of
122 insufficiency for years, during which time high discharge estimates are unreliable, especially where they
123 are made by extrapolating beyond observed maximum flow.

124

125 Gauge placement presents another obstacle to the rapid deployment of discharge monitoring stations
126 (Isaacson & Coonrod 2011). Stage measured via pressure transduction is susceptible to bias and
127 nonlinearity under turbulent flow conditions (Horner et al. 2018). Sensors placed in a depositional area
128 may be buried by sediment, and installations in forested watersheds or debris flow regions may be
129 destroyed during floods. Often, equipment must be relocated at least once before a new gauge site can be
130 properly established. Even an established stage-discharge rating curve must be regularly updated and
131 maintained because the bed of the river can change as sediment is deposited or excavated, altering the
132 relationship between stage and flow.

133

134 For some studies aiming to quantify stream or watershed processes that require continuous discharge time
135 series, establishment of a high-quality monitoring station may be infeasible. Where co-location of the site
136 of interest with an existing stream gauge is also infeasible, record extension (Hirsch 1982; Nalley et al.
137 2020) and gap-filling (Harvey et al. 2012; Arriagada et al. 2021) techniques cannot be employed, as these
138 rely on prior knowledge of the statistical properties of the discharge time series being augmented. In such
139 scenarios, streamflow reconstruction or prediction techniques are suitable, as these may proceed a priori
140 or from minimal observation. Reconstruction typically involves methods that leverage the correlation
141 between a partially measured target site and nearby “donor” (predictor) gauges. Discharge may also be
142 quantified in the absence of direct measurements at the target location via statistical (Chokmani & Ouarda
143 2004), mechanistic (Regan et al. 2019), or machine learning (Kratzert et al. 2022) modeling techniques.

144

145 Here, we use both linear regression (OLS, L2/ridge, segmented) and deep learning (LSTM-RNN)
146 approaches to reconstruct discharge from the early operational phase (2015-2022) of the National
147 Ecological Observatory Network (NEON), a time during which site selection issues and rating curve
148 development rendered potentially unreliable many site-months of discharge estimates (Rhea et al. 2023a).
149 Our goal was to achieve Kling-Gupta Efficiency (KGE) scores greater than those of the official NEON
150 continuous discharge product at as many sites as possible. A secondary goal was to improve temporal
151 coverage of the official record where it contains gaps. For researchers intending to use NEON continuous
152 discharge data between 2015 and 2022, the results of this effort, as well as efforts by Rhea et al. (2023a),
153 can ensure that data gaps and questionable periods in the official record are replaced by high-quality
154 estimates wherever possible. We provide composite discharge series for all 27 NEON stream gauge
155 locations, built from the best NEON-published estimates and the best estimates generated by this study
156 (<https://doi.org/10.6084/m9.figshare.c.6488065>). Composite series can be visualized at
157 https://macrosheds.org/data/vlah_et al_2023_composites/.

158

159 The success of this effort demonstrates the viability of “virtual gauges” (*sensu* Philip & McLaughlin
160 2018; not to be confused with the “virtual staff gauges” of Seibert et al. 2019). In this study, we use the
161 term to describe sites at which discrete discharge observations can be used to fit or evaluate models that
162 generate continuous flow. For accurate results, field measurement campaigns should prioritize
163 characterizing the distribution of possible flow conditions, rather than achieving any particular threshold
164 number of observations. Methods like those presented could be used to reduce the cost and simplify the
165 process of establishing streamflow monitoring sites, especially in river networks that are already partially
166 gauged.

167

168 2. Methods

169

170 2.1 Data selection, acquisition, and processing

171

172 We used the “neonUtilities” package (Lunch et al. 2022) in R to retrieve NEON discharge data. Officially
173 released (NEON 2023c) and provisional (NEON 2023b) field measurements were used to fit linear
174 regression models and evaluate all models, as these data were collected directly by NEON technicians,
175 using a combination of state-of-the-art methods including acoustic Doppler current profiling (ADCP;
176 Moore et al. 2017), conservative salt tracer releases (Tazioli 2011), and flow meter measurements

177 (Pantelakis et al. 2022). We used quality-controlled “finalQ” values where available, or “totalQ” values
178 (taken directly from the flowmeter) in their absence. We refer to NEON’s discharge field measurements
179 hereafter as e.g. “the response variable”, or “response discharge time series,” in the context of linear
180 regression, or as the “target” variable in the context of machine learning. In either context, we refer to the
181 27 NEON sites for which discharge predictions were generated as “target sites” or “target gauges” (Table
182 1).

183

184 Continuous discharge data (NEON 2023a) were also retrieved via neonUtilities. We used
185 RELEASE-2023 and *not* provisional data in this case. These data were used to finetune a subset of
186 site-specific neural network models, and to construct composite discharge series. Provisional continuous
187 discharge data were not used. Evaluation results used to distinguish likely reliable vs. potentially
188 unreliable subsets of NEON’s RELEASE-2023 continuous discharge time series, per site-month, were
189 provided by Rhea et al. (2023a) and accessed through HydroShare (Rhea 2023). Continuous elevation of
190 surface water data are available, but approximately one third of all site-months are marked by
191 disagreement between reported surface elevation and measured stage, or by likely sensor drift (Rhea et al.
192 2023a). We therefore chose not to use surface elevation to inform our models, though it no doubt contains
193 predictive value.

194

195 Donor gauge data for linear regression analysis were acquired primarily from the US Geological Survey’s
196 National Water Information System (NWIS), using the “dataRetrieval” package (DeCicco et al. 2022) in
197 R. NWIS gauge ID numbers are provided in `cfg/donor_gauges.yml` at the GitHub and Zenodo links
198 below. Additional donor gauge data from Niwot Ridge LTER and Andrews Forest LTER were retrieved
199 from the MacroSheds dataset (Vlah et al. 2023) via package “macrosheds” (Rhea et al. 2023b), and from
200 the EDI data portal (Johnson et al. 2020), respectively.

201

202 We used the original CAMELS dataset (Newman et al. 2014; Addor et al. 2017), the USGS National
203 Hydrologic Model with Precipitation-Runoff Modeling System (NHM-PRMS; hereafter NHM; Regan et
204 al. 2019), and the MacroSheds dataset as training data for neural network simulations of discharge data at
205 each target site. CAMELS watershed attributes were generated for MacroSheds and NHM sites using the
206 code provided at <https://github.com/naddor/camels>, except where otherwise indicated in Table 2, and
207 daily Daymet meteorological forcings (Thornton et al. 2022; *sensu* Newman et al. 2015) were retrieved
208 via Google Earth Engine (Gorelick et al. 2017). All code for this project can be found on GitHub, at
209 https://github.com/vlahm/neon_q_sim, or in the Zenodo archive at
210 <https://doi.org/10.5281/zenodo.10067683>. All data sources and links are provided in Table A2.

211

212 2.2 Donor Gauge Selection

213

214 Candidate donor gauges were identified by visually examining an interactive map of NEON gauges,
215 USGS gauges, and MacroSheds gauges
216 (https://macrosheds.org/ms_usgs_etc_reference_map/megamap.html), generated with package “mapview”
217 (Appelhans et al. 2022) in R. We also used the National Water Dashboard of the USGS
218 (<https://dashboard.waterdata.usgs.gov/app/nwd/en/?aoi=default>) to identify active gauges in Alaska, USA.
219 For each target site, up to four donor gauge candidates were selected on the basis of spatial proximity and
220 geographic similarity to the target site (Figure 1). Generally, no greater than this number of gauges were

221 even remotely reasonable candidates (i.e. within 50 km of the target site; not in an urban area; not
 222 downstream of a reservoir), but for one target site (MCRA) we had ten nearby candidate gauges to select
 223 from—all associated with the Andrews Experimental Forest in western Oregon State, USA. In this case we
 224 chose three candidate sites representing a catchment upstream of the target site (GSWS08), downstream
 225 of the target site on the MCRA mainstem (GSLOOK), and downstream on a tributary of MCRA
 226 (GSWS01).

227

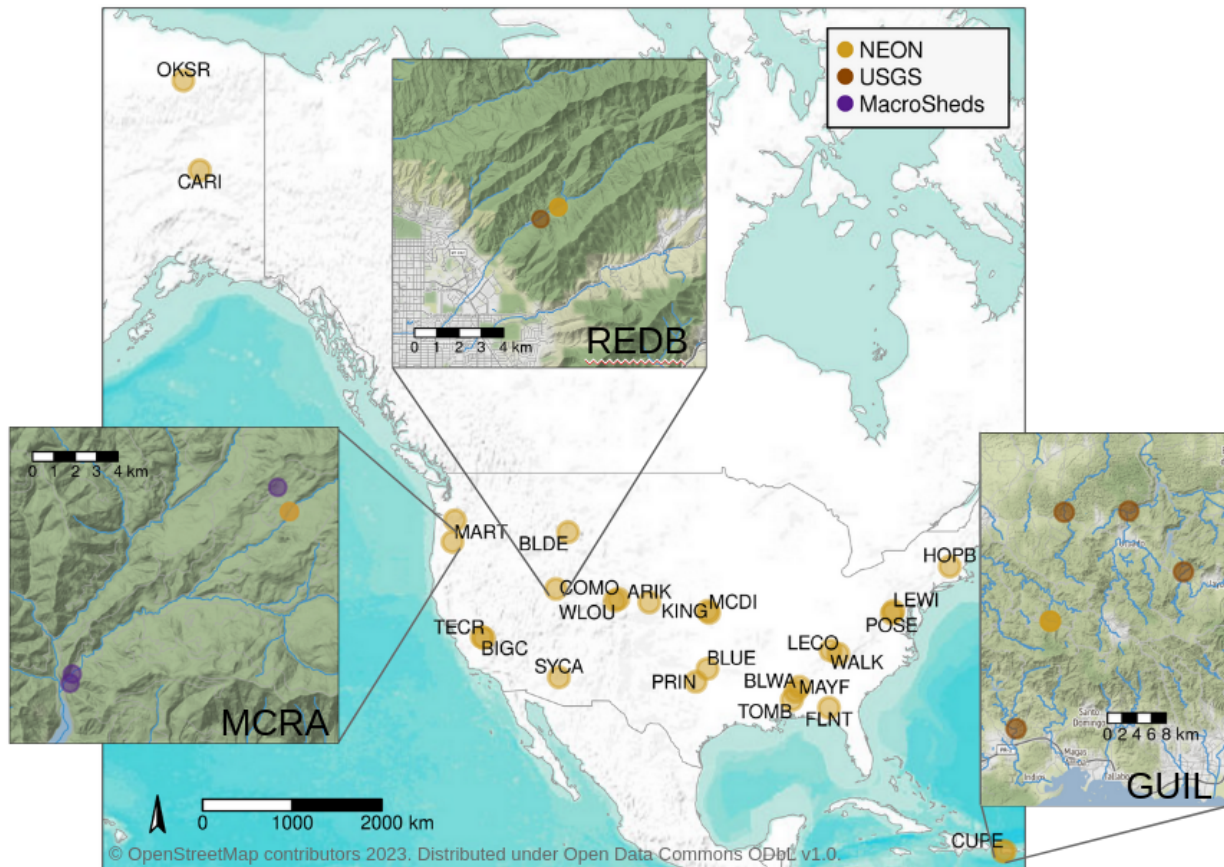
228 Barring gauges on reaches that are subject to overt human influence, the exact methods used to choose
 229 donor gauges are of little consequence, so long as informative donor gauges are not overlooked. In
 230 practice, there will usually be just a few, if any, potential donor gauges available for a given location. If
 231 multiple donor gauges are included in a regression, L2 regularization (ridge regression) should be used to
 232 account for their covariance (see Sect. 2.4)

233

234 2.3 Target sites

235

236



237

238 **Figure 1: Map of target sites (NEON) and donor gauge candidates for three target sites: MCRA =**
 239 **McRae Creek, state of Oregon; REDB = Red Butte Creek, state of Utah; GUIL = Rio Guilarte,**
 240 **Puerto Rico.**

241

242 All 27 lotic (flowing) aquatic sites associated with NEON were included as target sites for discharge
 243 prediction in this study (Figure 1). Sites TOMB, BLWA, and FLNT are installed on major rivers,
 244 downstream of hydropower dams. All other sites have been free of dam influence since 2012 at the
 245 latest, and are designated “wadeable streams” by NEON. In addition to the three sites above,
 246 hydrology at BLUE, GUIL, KING, MCDI, and ARIK may be influenced by agricultural activity,
 247 especially in the relatively arid Midwest (i.e. states KS, CO, OK). Continuous discharge data for
 248 TOMB are provided by a nearby gauge of the U.S. Geological Survey’s National Water Information
 249 System, and are given at hourly intervals, rather than NEON’s customary 1-minute intervals.
 250

251 **Table 1: Target sites for discharge prediction. See <https://www.neonscience.org/field-sites> for more**
 252 **information.**

Site code	Full name	State (USA)	Watershed area (km ²)	Mean watershed elevation (m)
TOMB	Lower Tombigbee River	AL	47085.3	20
BLWA	Black Warrior River	AL	16159.4	22
FLNT	Flint River	GA	14999.4	30
ARIK	Arikaree River	CO	2631.8	1179
BLUE	Blue River	OK	322.2	289
SYCA	Sycamore Creek	AZ	280.3	645
OKSR	Oksrukuyik Creek	AK	57.8	766
PRIN	Pringle Creek	TX	48.9	253
BLDE	Blacktail Deer Creek	WY	37.8	2053
CARI	Caribou Creek	AK	31.0	225
MCDI	McDiffett Creek	KS	22.6	396
REDB	Red Butte Creek	UT	16.7	1694
MAYF	Mayfield Creek	AL	14.4	77
KING	Kings Creek	KS	13.0	324
HOPB	Lower Hop Brook	MA	11.9	203
LEWI	Lewis Run	VA	11.9	152
BIGC	Upper Big Creek	CA	10.9	1197

GUIL	Rio Guilarte	PR	9.6	551
LECO	LeConte Creek	TN	9.1	579
MART	Martha Creek	WA	6.3	337
WLOU	West St Louis Creek	CO	4.9	2908
CUPE	Rio Cupeyes	PR	4.3	157
MCRA	McRae Creek	OR	3.9	876
COMO	Como Creek	CO	3.6	3021
TECR	Teakettle Creek - Watershed 2	CA	3.0	2011
POSE	Posey Creek	VA	2.0	276
WALK	Walker Branch	TN	1.1	264

253

254

255 **2.4 Linear regression and model selection**

256

257 All donor and response discharge time series were neglog transformed (Eq. 1; Whittaker et al. 2005)
258 before fitting linear regression models.

259

$$260 x_{\text{neglog}} = \text{sign}(x) \log(|x| + 1)$$

$$261 \quad (1)$$

262

263 Series were scaled by 1000 before transformation, in order to reduce the disproportionate impact of
264 adding one to every value. Response observations were synchronized to the interval of the predictor series
265 by approximate datetime join, allowing forward or backward time-shifts of up to 12 hours if necessary.

266

267 One of three forms of linear regression was employed at each site, depending on the number and location
268 of donor gauges, and the donor-target gauge relationships. For sites with a single donor gauge (REDB,
269 HOPB, BLUE, SYCA, LECO), considered predictors were: discharge from the donor gauge, a 4-season
270 categorical variable, and their interaction. Additionally, an intercept parameter could be estimated, or not,
271 for each specification. Thus, up to six models were fit using Ordinary Least Squares (OLS) regression
272 (Galton 1886), ensuring at least 15 observations per model parameter. At LECO, an additional dummy
273 variable was included to address an intercept change due to a wildfire in November of 2016. The best
274 model was selected via 10-fold cross-validation, minimizing mean squared error (MSE). MSE, being a
275 squared-error term, disproportionately penalizes inaccurate prediction of high discharge values, and helps
276 to balance against the relative rarity of high discharge measurements in the field data. At site SYCA, the
277 log-log relationship between discharge at the target gauge and a single donor gauge exhibited a distinct

278 breakpoint, and segmented least-squares regression was used (R package “segmented”; Muggeo 2008). At
 279 all other sites (19 in total), predictors included discharge series from 2-4 donor gauges, season, and all
 280 interactions. To control overfitting and shrink covarying coefficients toward zero, we used L2
 281 regularization (ridge regression; Gruber 2017) via R package “glmnet” (Friedman et al. 2010). As with
 282 the other regression approaches, 10-fold cross-validation and MSE loss were used for model parameter
 283 selection—in this case for the value of the penalty hyperparameter λ , which was set to the mean across
 284 folds of λ producing minimum cross-validated error. Unlike OLS and segmented regression, ridge
 285 regression uses biased estimators that complicate calculation of prediction intervals. We generated 95%
 286 prediction intervals for ridge regression discharge estimates using the 95th percentiles of 1000 bootstrap
 287 predictions at each prediction point, generated from 1000 resamples of the fitting data, stratified by
 288 season. We emphasize that these prediction intervals should be conservative estimates of the true
 289 uncertainty, as they do not fully account for uncertainty due to bias (Goeman et al. 2012).

290

291 For each site, we fit two sets of models as described above, one with discharge scaled by watershed area
 292 (i.e. “specific discharge” in the surface water hydrology sense) prior to transformation, and one without
 293 areal scaling. Only one model from each set was ultimately selected for each target site, on the basis of
 294 Kling-Gupta efficiency (KGE; Gupta et al. 2009), a composite model efficiency metric that incorporates
 295 measures of correlation, variance, and bias. We also report percent bias and Nash-Sutcliffe efficiency
 296 (NSE; Nash & Sutcliffe 1970), a measure of predictive accuracy that implicitly compares predictions to a
 297 mean-only reference model.

298

299 Predictions were generated for all time points during which data were available at the selected donor
 300 gauges. At target site COMO, a secondary model omitting one donor gauge was able to produce 36%
 301 more predictions than the selected model, so our predicted discharge at COMO is a composite of both
 302 models, preferring the better model’s predictions where available. We were unable to locate sub-daily
 303 donor gauge data near COMO, so regression predictions for this site are at a daily interval. Regression
 304 predictions for all other sites were generated at sub-daily intervals matching the coarsest interval across
 305 predictor gauges—generally 15 minutes, though note that in most cases these predictions were interpolated
 306 to five minutes for our composite discharge product.

307

308 **2.5 Neural network setup and operation**

309

310 Supplementing the linear regression methods described above, we simulated discharge data at all 27 target
 311 sites using long short-term memory recurrent neural networks (LSTM-RNN; hereafter “LSTM”;
 312 Hochreiter & Schmidhuber 1997). Four LSTM strategies were employed, all of which involved training
 313 on a large and diverse corpus of stream discharge data (Table 3). Two of these strategies included further
 314 finetuning to the time-series dynamics of each target site in turn. Due to the relative scarcity of
 315 field-measured discharge observations (between 39 and 213 per site; mean 122), none were used in
 316 LSTM training. Instead, these measurements were used only to evaluate predictions. LSTMs trained in
 317 this study are intended only for discharge prediction within the temporal and spatial bounds of NEON’s
 318 early operational phase, not for forecasting or application to other sites. Therefore, all available, daily
 319 training data were used as such; no validation set was kept for hyperparameter tuning, and no holdout set
 320 of daily estimates was kept for evaluation (note that split-sample designs may be undesirable more
 321 generally: Arsenault et al. 2018; Guo et al. 2018; Shen et al. 2022). See Kratzert et al. (2019b) and Read

et al. (2019) for split-sample considerations in the context of a generalist and process-guided generalist LSTM, respectively.

324

After a hyperparameter search routine, described below, potentially skilled models were identified as those achieving at least 0.5 KGE and 0.4 NSE. The best performing, potentially skilled LSTM for each site (if applicable) was then re-trained 30 times, forming an ensemble. Ensembles were trained for 18 of 27 sites. LSTM predictions included in our composite discharge product are means taken across the distributions of ensemble point predictions. Uncertainty bounds were computed as the 2.5 and 97.5% quantiles of these distributions. LSTM skill was evaluated on the basis of mean ensemble efficiency (KGE) with respect to field-measured discharge (Table A1).

332

Daily discharge time series (training data) and field-measured discharge were scaled by watershed area. For each predicted day, LSTMs received 5 dynamic Daymet meteorological forcing variables and 11 static watershed attribute summary statistics (Table 2). Multitask learning (Caruana 1998; Sadler et al. 2022) was found to improve discharge prediction broadly in a preliminary analysis, so Daymet minimum air temperature was used as a secondary target variable. Kratzert et al. (2019a) found that a maximum of about 150 preceding days were able to influence LSTM output on a similar prediction problem, so we set the input sequence length to 200 days to ensure full utilization of available information. In other words, for each day of prediction, the model was able to leverage information from the preceding 200 days.

341

We employed four different training pipelines described in Table 3. Of the 671 CAMELS watersheds (i.e. basins), we used a subset of 531 with undisputed areas less than 2000 km² (Newman et al. 2017). For finetuning data, we used version 1 of the MacroSheds dataset (Vlah et al. 2023). We excluded MacroSheds sites outside North America, or with coastal or urban hydrological influence, for a total of 133 sites out of the 169 that are currently available. We chose MacroSheds sites for finetuning because the MacroSheds and NEON datasets focus primarily on small watersheds, often smaller than 10 km² in area, while only eight CAMELS watersheds are smaller than 10 km² and most are larger than 100 km² (Vlah et al. 2023). Daily mean discharge computed from NEON's continuous discharge product, only for those site-months deemed Tier 1 or Tier 2 by Rhea et al (2023a), was used alongside MacroSheds data for finetuning.

352

For the process-guided strategies, we used NHM estimates for all reaches coinciding with a CAMELS or MacroSheds gauge, for a total of 551 reaches. Only nine target sites on relatively high-order streams were amenable to the process-guided specialist approach, as these sites are on reaches large enough to be modeled by the NHM. The most recent version of the NHM at the time of this writing provides discharge estimates beginning in 1980, and ending in 2016, just before the installation of most NEON target sites.

358

Table 2: LSTM input data. * = Attribute tested as an afterthought, but not included in this study due to negligible improvement in trial parameter search.

Meteorological forcing data (watershed-average time series)	
Maximum air temp	2-meter daily maximum air temperature (°C)
Precipitation	Mean daily precipitation (mm/day)

Solar radiation	Daily surface-incident solar radiation (W/m ²)
Vapor pressure	Near-surface daily average vapor pressure (Pa)
PET	Potential evapotranspiration (mm); estimated using Priestley-Taylor formulation with gridded alpha product (Aschonitis et al. 2017)
Watershed attributes (statistics computed over full record)	
Precipitation mean	Mean daily precipitation (mm/day)
PET mean	Mean daily potential evapotranspiration (mm/day); estimated using Priestley-Taylor formulation with gridded alpha product (Aschonitis et al. 2017)
Aridity index	Ratio of PET mean to Precipitation mean
Precip seasonality	Seasonality of precipitation; estimated by representing annual precipitation and temperature as sine waves. Positive values indicate summer peaks, while negative values indicate winter peaks. Values near 0 indicate uniform precipitation throughout the year.
Snow fraction	Fraction of precipitation falling on days with temp < 0 °C
High precipitation frequency	Frequency of high precipitation days (days with ≥ 5x mean daily precipitation)
High precip duration	Average duration of high precipitation events (number of consecutive days ≥ 5x mean daily precipitation)
Low precip frequency	Frequency of dry days (days with precipitation < 1 mm/day)
Low precip duration	Average duration of dry periods (number of consecutive days with precipitation < 1 mm/day)
Elevation	Catchment mean elevation (m)
Slope	Catchment mean slope (m/km)
Area	Catchment area (km ²)
Source*	Binary indicator for NHM estimates–process-guided LSTMs only.
Target data (time series)	
Discharge	Specific discharge, or discharge normalized by watershed area. The same

	quantity may be referred to as “runoff” in other studies (mm/day).
Minimum air temp	2-meter daily minimum air temperature (°C)

361

362 **Table 3: LSTM model training pipelines used in the simulation of discharge at target sites. Here,**
 363 **“NEON” refers to NEON’s continuous discharge product, RELEASE-2023, with quality-flagged**
 364 **estimates and < Tier-2 site-months (according to Rhea et al. 2023a) removed.**

Model type	Phase 1	Phase 2	Phase 3
Generalist	Pretrain on CAMELS	Finetune on MacroSheds + NEON	N/A
Specialist	Pretrain on CAMELS	Finetune on MacroSheds + NEON	Finetune on NEON target site
Process-guided generalist	Pretrain on CAMELS + CAMELS-NHM	Finetune on MacroSheds + MacroSheds-NHM + NEON + NEON-NHM	N/A
Process-guided specialist	Pretrain on CAMELS + CAMELS-NHM	Finetune on MacroSheds + MacroSheds-NHM + NEON + NEON-NHM	Finetune on NHM estimates for target site

365

366 LSTMs were configured in R, and trained using v1.3.0 of the NeuralHydrology library in Python
 367 (Kratzert et al. 2022; Van Rossum & Drake 2009) on the Duke Compute Cluster at Duke University,
 368 Durham NC, USA. All trained models used the Adam optimizer (Kingma & Ba 2014) and
 369 NeuralHydrology’s “NSE loss” function, after an initial evaluation in which we compared it to MSE and
 370 root mean squared error (Table 4). Learning was annealed using series of three fixed rates for pretraining
 371 and for round one of finetuning, according to Eq. (2):

372

$$r = \begin{cases} a, & e \in \{0, \dots, \lfloor \frac{E}{3} \rfloor\} \\ \frac{a}{10}, & e \in \{\lfloor \frac{E}{3} \rfloor, \dots, \lfloor \frac{2E}{3} \rfloor\} \\ \frac{a}{100}, & e \in \{\lfloor \frac{2E}{3} \rfloor, \dots, E\} \end{cases} \quad (2)$$

373

374

375

376 Where r is the learning rate, a is any power of 10 between 0.1 and 10^{-7} , and E is the number of training
 377 epochs. Learning rate was annealed using series of two fixed rates for round two of finetuning, according
 378 to Eq. (3):

379

$$r = \begin{cases} \frac{a}{10}, & e \in \{0, \dots, \lfloor \frac{E}{2} \rfloor\} \\ \frac{a}{100}, & e \in \{\lfloor \frac{E}{2} \rfloor, \dots, E\} \end{cases} \quad (3)$$

380

381

382

383 Learning rate and other hyperparameters were selected via an inexhaustive (pseudo) grid search (Table 4),
 384 i.e. we specified a sequence of possible values for each hyperparameter and randomly selected from them
 385 to specify 30 models for each generalist. For each site, one specialist model was then configured to further
 386 finetune each of the 30 generalists, again using partial grid search to define any mutable hyperparameters.
 387 Otherwise, hyperparameters were inherited from the previous training period (Table 4). Due to our
 388 incomplete hyperparameter search procedure, better combinations probably exist. We elected not to
 389 exhaustively pursue optimal hyperparameter combinations due to the computational demand of a full grid
 390 search, and a lack of access via NeuralHydrology to callback methods necessary for implementation of
 391 true random search (Bergstra & Bengio 2012).

392

393 **Table 4: LSTM hyperparameter search space for all model types, and selected values (bold, italic)**
 394 **used for pretraining. These were observed to allow for both malleability and high performance of**
 395 **subsequent finetuning iterations over nearly 2000 exploratory LSTM trials. The ditto mark “``”**
 396 **indicates that a finetuning parameter is inherited from the preceding training iteration. The**
 397 **relationship of a to the learning_rate is defined in Equations 2 and 3. See the NeuralHydrology**
 398 **documentation for parameter definitions:**

399 <https://neuralhydrology.readthedocs.io/en/latest/usage/config.html>.

LSTM parameter	Pretrain	Finetune 1	Finetune 2 (specialists only)
hidden_size	20, 30 , 40, 50	``	``
output_dropout	0.1, 0.2, 0.3, 0.4, 0.5 , 0.6	0.2, 0.3, 0.4, 0.5	``
learning_rate a	10^{-2} , 10^{-3} , 10^{-4} , 10^{-5}	10^{-2} , 10^{-3} , 10^{-4} , 10^{-5}	10^{-2} , 10^{-3} , 10^{-4} , 10^{-5}
batch_size	32, 64, 128, 256, 512 , 1024	32, 64, 128, 256, 512	``
epochs	20, 30 , 40, 50, 60	20, 30, 40	10, 20, 30
finetune_modules	N/A	head, lstm, head & lstm	head, lstm
target_variables	discharge, <i>discharge & min air temp</i>	``	``
loss	<i>NSE</i> , MSE, RMSE	``	``

400

401 All LSTM models were outfitted with fully connected, single-layer embedding networks to efficiently
 402 encode inputs as fixed-length numerical vectors (Arsov & Mirceva 2019). Separate embedding networks
 403 were used for static and dynamic inputs, with 20 neurons for static inputs and 200 neurons for dynamic
 404 inputs. All embedding neurons used the hyperbolic tangent activation function. Another advantage of
 405 embedding networks in the context of the NeuralHydrology library is that they provide one of few
 406 opportunities to introduce dropout, which can improve training efficiency and reduce overfitting
 407 (Srivastava et al. 2014).

408

409 2.6 Composite discharge data product

410

411 This study generated time-series predictions of discharge for each lotic NEON site using up to three
412 distinct processes: linear regression on absolute discharge, linear regression on specific discharge, and one
413 of four LSTM strategies. We provide regression predictions wherever applicable (24 of 27 sites). LSTM
414 predictions are provided only for sites that had promising model performance after a hyperparameter
415 search, and for which ensemble models were therefore trained (18 of 27). All model outputs and results
416 from this study are archived at <https://dx.doi.org/10.6084/m9.figshare.22344589>.

417

418 In addition to predictions from individual modeling strategies, we provide an analysis-ready discharge
419 dataset for all 27 sites that splices the best available predictions across methods, including published
420 NEON estimates (NEON 2023a), into composite series
421 (<https://dx.doi.org/10.6084/m9.figshare.23206592>), which can be visualized interactively at
422 https://macrosheds.org/data/vlah_et_al_2023_composites/. Composite series for each NEON site begin at
423 the start of site operation and extend to at most September 30, 2021, the last date included in the 2023
424 release of NEON's continuous discharge product. We also provide individual model predictions extending
425 through 2022. A complete list of products from this study, and their links, can be found in Table A3.

426

427 To construct composite series, we first distinguished as “good” site-months of NEON discharge estimates
428 categorized as Tier 1 or Tier 2 by Rhea et al. (2023a). For a NEON site-month to meet the requirements
429 for at least Tier 2, four requirements must be met. The linear relationship between stage, determined from
430 pressure transducer readings, and field-measured gauge height must score at least 0.9 NSE. The
431 transducer-derived stage series must also pass a drift test, relative to gauge height, but only if sufficient
432 data exist to perform such a test. The rating curve used to relate stage to discharge must score at least 0.75
433 NSE, and fewer than 30% of predicted discharge values may exceed the range of measured discharge
434 used to build the curve. See Rhea et al. (2023a) for further details.

435

436 Although only 50% of NEON's RELEASE-2023 estimates are classified as Tier 1 or Tier 2, the remainder
437 may still be of high analytical value if NEON's quality control indicators and uncertainty bounds are
438 observed. We also stress that NEON rating curves and protocols have improved over the course of its
439 early operational phase, and continue to do so.

440

441 We then ranked the available predictions for each site, assigning rank 1 either to predictions from linear
442 regression, or to NEON's continuous data product, depending on overall KGE and NSE against field
443 measured discharge. KGE was considered first, and used to determine preference except in cases where
444 the difference between NSE scores was greater than that between KGE scores, and opposite in sign. Rank
445 2 predictions were then used to fill gaps of 12 or more hours in the rank 1 series, but only “good” NEON
446 site-months were included. Only after this first round of gap-filling were the remaining NEON data
447 incorporated, with site-years achieving at least 0.5 KGE and 0.5 NSE against field-measured discharge
448 being used to fill still-remaining gaps. Finally, daily LSTM predictions (placed at 12:00:00 UTC on the
449 day of prediction) were used to fill any recalcitrant gaps, but only if produced by an ensemble model
450 achieving at least 0.5 KGE and 0.5 NSE across all field discharge observations. Note that while such
451 benchmarks are in common use (Moriassi et al. 2015), the efficiency that any model can or should achieve
452 varies substantially with the hydroclimate and watershed characteristics of a given site (Seibert et al.
453 2018). We provide all data and code for modifying the composite discharge product in accordance with

454 alternative benchmarks as users see fit. After visual examination of composite series plots, we chose to
 455 prefer NEON predictions to linear regression predictions at site ARIK, “good” or not, due to frequent
 456 sharp disjoints between the two predicted series. See Table A1 for an account of linear regression and
 457 LSTM methods used in the construction of ensemble series.

458

459 The prevailing interval varies across data sources used to assemble our composite discharge product, from
 460 one minute (NEON) to one day (LSTM predictions; regression predictions at site COMO). Regression
 461 predictions were primarily generated at 15-minute intervals, and their timestamps are always divisible by
 462 15 minutes. Around the prevailing NEON interval there is considerable variation due to data gaps and
 463 sensor reconfigurations, both across sites and across the temporal ranges of each site’s record. To reduce
 464 the complexity associated with irregular time-series analysis, we synchronized the interval across data
 465 sources to five minutes. Regression estimates were linearly interpolated to five minutes, though gaps
 466 larger than 15 minutes were not interpolated. NEON estimates were first smoothed with a triangular
 467 moving average window of 15 minutes to remove unrealistic minute-to-minute noise associated with
 468 Bayesian error propagation. They were then interpolated the same way as the regression estimates, and
 469 finally downsampled to five minutes, with some timestamps being shifted by up to two minutes. For
 470 example, a duration of 30-minute sampling, with a sample taken at 00:03:00, would be shifted by two
 471 minutes, by rounding each timestamp up to the nearest minute divisible by five.

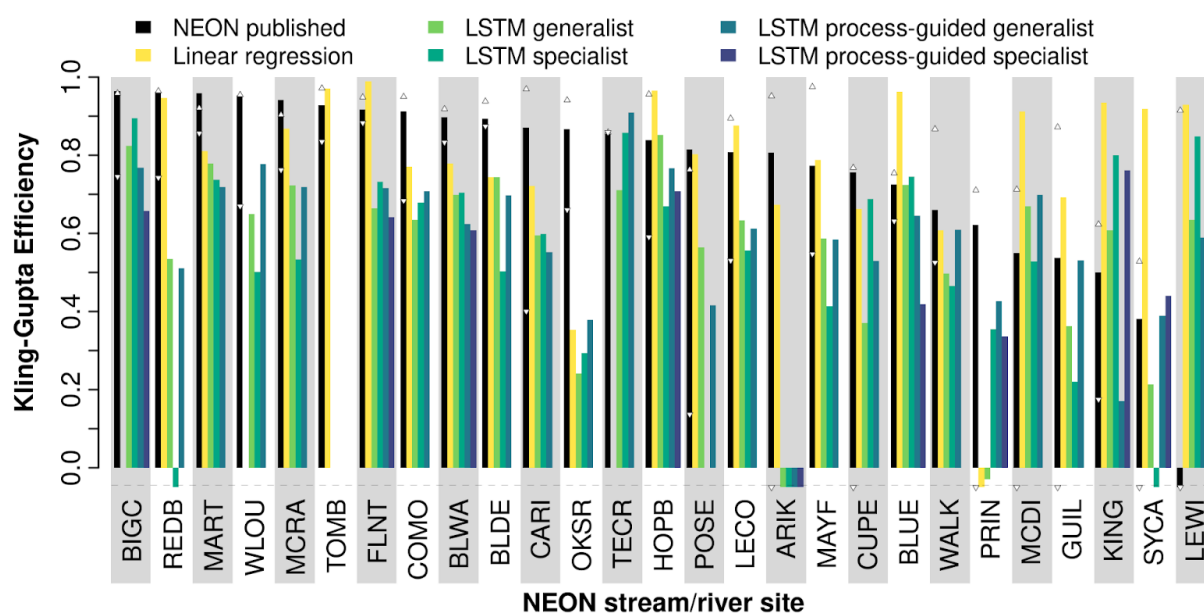
472

473 3. Results

474

475 A performance comparison of linear regression on discharge from donor gauges, and four LSTM
 476 strategies, is shown in Figure 2 and Figure A1, and detailed in Table A1. Via linear regression, we were
 477 able to produce 15-minute discharge estimates at 11 sites with overall KGE scores higher than those of
 478 published series (Figure 2). At four of the same sites, we achieved higher KGE via LSTM methods, which
 479 generated daily discharge series. Of the ten sites at which published discharge KGE was less than 0.8, we
 480 improved five to above that mark (mean 0.932, $n = 5$).

481



483 **Figure 2: Efficiency of five stream discharge prediction methods and NEON’s published continuous**
 484 **discharge product at 27 NEON gauge locations, versus field-measured discharge. Small, white**
 485 **triangles represent max/min KGE of published discharge by water year (Oct 1 through Sept 30)**
 486 **with at least 5 field measurements (or 2 for site OKSR). KGE was computed on all available**
 487 **observation-estimate pairs except those with quality flags (dischargeFinalQF or**
 488 **dischargeFinalQFSciRvw of 1). For the best performing LSTM method, at all sites except TECR,**
 489 **FLNT, REDB, WALK, POSE, and KING, displayed KGE is averaged over 30 ensemble runs with**
 490 **identical hyperparameters. For the sites just named, performance of a chosen method, after**
 491 **ensembling, dropped below that of at least one other method’s optimal KGE from parameter**
 492 **search. For all other LSTM site-method pairs, which were not ensembled, displayed performance is**
 493 **that of the best model trained during the parameter search phase. Sites are ordered by the KGE of**
 494 **NEON continuous discharge. See Table 3 for LSTM model definitions. KGE of 1 is a perfect**
 495 **prediction, while KGE of -0.41 is similar in skill to prediction from the mean. Negative values are**
 496 **truncated at -0.05 in this plot to improve visualization.**

497

498

499 For 12 of 27 sites, linear regression on specific discharge (i.e. scaled by watershed area) provided the
 500 most accurate discharge predictions, while linear regression on absolute discharge performed better at the
 501 other 12 sites with donor gauges. LSTM models (as proper ensembles) outperformed linear regression at
 502 only 2 sites. In general, linear regression provided more accurate predictions than all LSTM methods.
 503 Linear regression on absolute discharge produced estimates with median NSE of 0.848 and median KGE
 504 of 0.806, across sites ($n = 24$; Table 5). Linear regression on specific discharge produced similar median
 505 scores (Table 5), but with deviations of up to 0.05 NSE and 0.08 KGE at individual sites.

506

507

508 **Table 5: Performance of five stream discharge prediction methods, and official continuous**
 509 **discharge time-series data, across n of 27 NEON gauge locations (final column). For both the**
 510 **Nash-Sutcliffe and Kling-Gupta Efficiency coefficients, a value of 1 indicates perfect prediction. A**
 511 **value of 0 NSE indicates that predictive skill is equivalent to prediction from the mean, while**
 512 **negative NSE is worse than mean prediction. This threshold lies at approximately -0.41 for KGE**
 513 **(Knoben et al. 2019). “Linreg” = linear regression on donor gauge discharge series, and “scaled”**
 514 **means predictor and response discharge were scaled by their respective watershed areas.**

Model/Data	NSE				KGE				n
	Median	Mean	Min	Max	Median	Mean	Min	Max	
Official record	0.880	0.417	-9.95	0.989	0.839	0.711	-1.50	0.964	27
Linreg	0.848	0.760	-0.038	0.993	0.806	0.746	-0.697	0.988	24
Linreg scaled	0.847	0.757	-0.037	0.993	0.807	0.743	-0.695	0.989	24
Generalist LSTM	0.473	-18.8	-498	0.904	0.634	-0.220	-20.2	0.852	26
Specialist LSTM	0.477	-12.6	-307	0.920	0.556	-0.256	-15.7	0.895	25

Model/Data	NSE				KGE				n
	Median	Mean	Min	Max	Median	Mean	Min	Max	
Process-guided generalist LSTM	0.434	-31.3	-824	0.848	0.618	-0.453	-26.4	0.869	26
Process-guided specialist LSTM	0.329	-92.0	-831	0.749	0.652	-2.40	-26.5	0.866	9

515

516

517 Linear regression was not applicable at sites TECR, BIGC, or WLOU due to the lack of donor gauges
518 contemporary with target gauge data. Donor gauges associated with Kings River Experimental
519 Watersheds exist within close proximity to TECR and BIGC, but we were unable to access up-to-date
520 discharge records for these gauges.

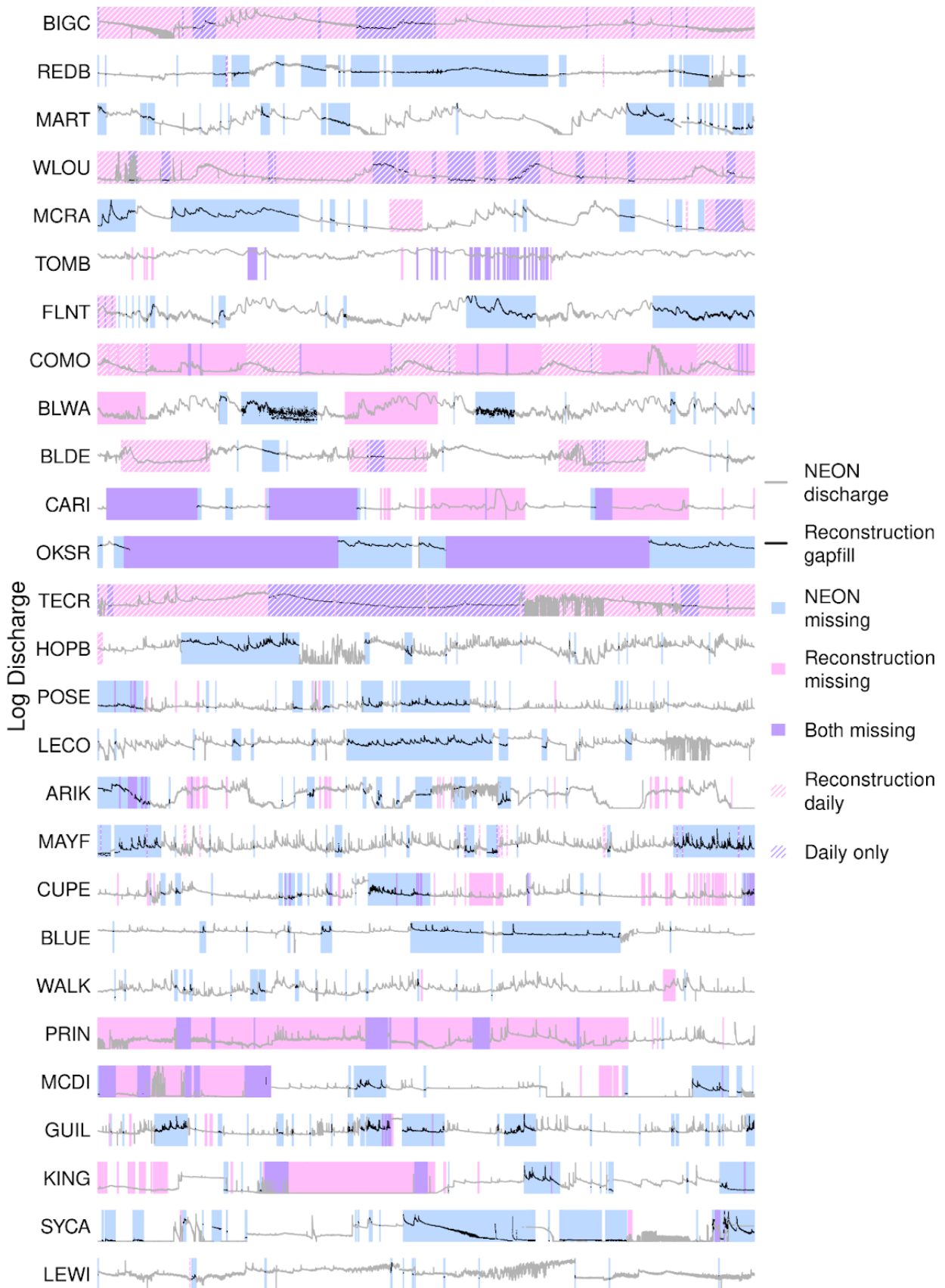
521

522 The process-guided specialist LSTM yielded predictions on par with those of the other LSTM strategies
523 in terms of KGE, (median 0.652; $n = 9$), but performed worst of the four in terms of NSE (median 0.329;
524 $n = 9$). Conversely, the specialist performed better than the generalist in terms of NSE, but not KGE. The
525 process-guided specialist LSTM strategy was viable at nine sites for which discharge estimates were
526 available from the National Hydrologic Model.

527

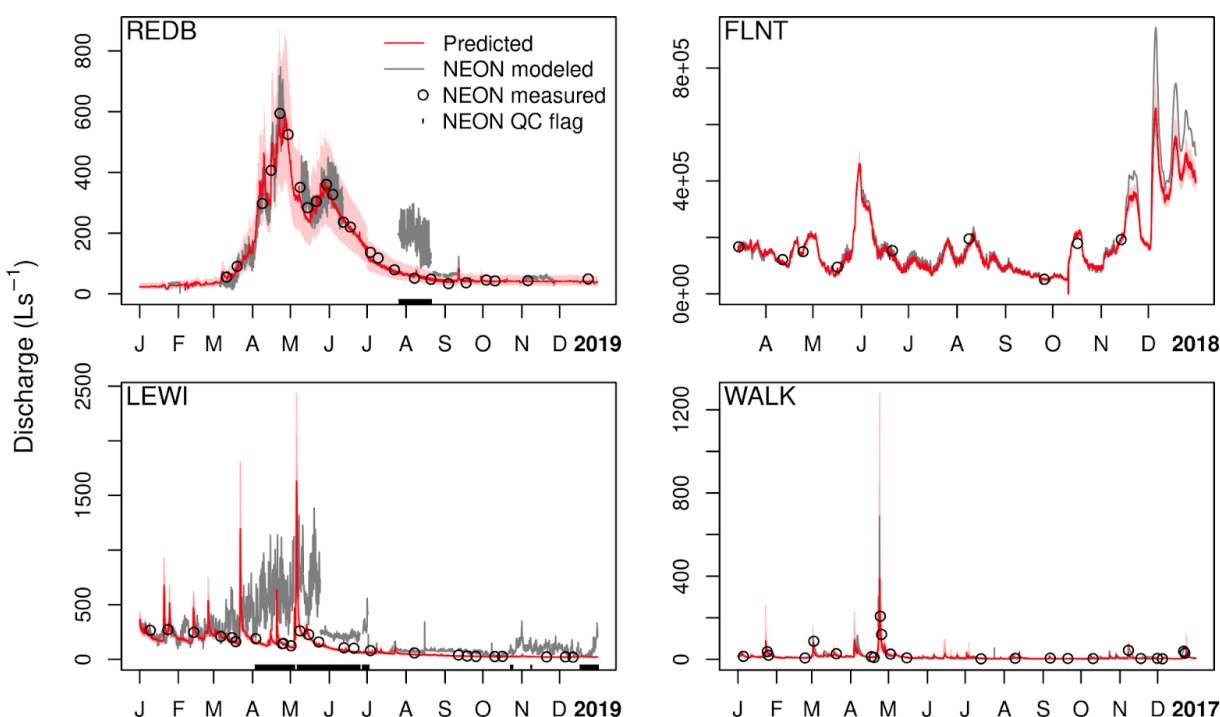
528 In addition to improvements in accuracy, estimates from this study inform ~5,981 site-days (75%) of
529 missing data in the official discharge record (Figure 3), though note that they also omit ~4,486 site-days
530 otherwise present in NEON's official record. Omissions occur wherever observations are missing from
531 the records of one or more donor gauges, and LSTM methods did not achieve desired efficiencies.
532 Approximately 1,221 site-days are missing from the official record and from our reconstructions.

533



535 **Figure 3: Durations of missing values (gaps) in NEON’s 2023 release of continuous discharge time**
 536 **series, illustrating gaps filled or informed by estimates from this analysis. All officially published**
 537 **values are shown, including those with quality control flags. Sites are ordered as in Figure 2. Gaps**
 538 **smaller than six hours are not indicated. Figure A10 is the same, but with a fixed and labeled x-axis.**

539 Estimated discharge time series from this study are of practical value for any researcher using NEON
 540 continuous discharge data, especially for those sites and site-months at which published data from
 541 NEON’s early operational phase may be unreliable (Rhea et al. 2023a). Figure 4 shows that official
 542 records at sites REDB and LEWI are compromised by disagreement (erratic sections of gray lines)
 543 between pressure transducer stage readings and manual gauge height recordings, discussed in Rhea et al
 544 (2023a). Red lines show improved estimates via linear regression on discharge from donor gauges. Sites
 545 FLNT and WALK show generally close agreement between NEON discharge and our regression
 546 estimates, but note uncertainty associated with high discharge values.



547

548 **Figure 4: Best linear regression predictions of continuous discharge for four NEON gauge-years,**
 549 **compared with official NEON discharge data. All officially published values are shown, including**
 550 **those with quality control flags, indicated by black marks on lower border. Light red polygons**
 551 **represent 95% prediction intervals. NEON uncertainty is not shown.**

552

553

554 4. Discussion

555

556 This study was designed to produce high-quality estimates of continuous discharge for NEON stream
 557 gauges, especially at ten gauges for which the KGE of published continuous discharge was lower than

558 0.8, over the full record, when compared to field-measured discharge. A secondary goal was to improve
559 temporal coverage of the official discharge record where possible.

560

561 We treat NEON field-measured discharge as truth, which means there are 39-213 observations for each
562 target site. Although these numbers represent a tremendous investment of time and technical effort, they
563 do not meet the high data volume requirements for most machine learning approaches, so we used field
564 discharge only to evaluate, rather than train, LSTM models. By contrast, in linear regression, regardless of
565 the details of any particular method, we ultimately fit a line to the relationship between donor gauge data
566 and field measurements at each target site. Because the linear regression models are allowed to “see” all
567 of the target site data (after a model is selected via cross-validation), they have a powerful advantage over
568 the LSTM approaches, which in this context must essentially treat target watersheds as if they are
569 ungauged. Furthermore, whereas the LSTM models must parameterize each day of prediction
570 individually, the regression models need only parameterize relationships between flow regimes. Still, if
571 given enough training data, including examples of watersheds and streams similar to each of those
572 modeled in this study, the LSTM approaches would eventually close the performance gap. See Figures
573 A2, A3, A4, A5, A7, and A8 for linear regression diagnostics.

574

575 In this study, discharge estimates produced by linear regression were more accurate than those generated
576 by LSTM models in 21 of 23 comparisons (Figure 2). This demonstrates the value of existing gauge
577 networks in advancing discharge estimation at newly or partially gauged locations; however, there is a
578 limit to the predictive potential of linear regression methods, as they depend on strong correlation
579 between streamflow at target and donor gauges. In principle, there is no such limit for machine learning
580 approaches, which are instead limited by the quality and quantity of training data.

581

582 The process-guided specialist LSTM yielded predictions on par with those of the other LSTM strategies
583 in terms of KGE, but performed worst of the four in terms of NSE, possibly indicating that information
584 gleaned from NHM estimates helped this strategy to accurately capture discharge variance and reduce
585 prediction bias, without ultimately improving the correlation between predictions and observations.
586 Unlike KGE, NSE only explicitly captures this latter metric (Nash & Sutcliffe 1970; Gupta et al. 2009).
587 Conversely, the specialist performed better than the generalist in terms of NSE, but not KGE, suggesting
588 information contained in NEON’s continuous discharge product was of disproportionate predictive value
589 relative to each of correlation, variance, and bias, favoring correlation.

590

591 The specialist may have been affected by data filtering choices. After filtering NEON continuous
592 discharge for rating curve issues, drift, and quality flags, relatively few daily estimates were available for
593 some sites (47-1642). Annual and seasonal variation in meteorological forcings and discharge in NEON
594 sites’ generally small, often mountainous watersheds may be large enough that finetuning a pretrained
595 LSTM on a few hundred days of site-specific data reduces its ability to generalize at that site. Our
596 specialist LSTM strategy in particular might be improved with a broader hyperparameter search,
597 especially one that explores smaller learning rates. Ideally, site-specific finetuning should enable better
598 prediction by allowing the network to assimilate information unique to the target site without corrupting
599 previously learned generalities. For validation plots of all ensembled LSTMs, see Figure A6.

600

601 The process-guided specialist LSTM strategy was viable at nine sites for which discharge estimates were
602 available from the National Hydrologic Model. By using a mechanistic (i.e. process-based) model with
603 higher spatial resolution than the NHM, it should be possible to apply this process-guided approach at
604 more of the NEON sites. A potentially stronger process-guided approach would use mechanistic model
605 predictions as features (predictors), rather than training targets, but that would require mechanistic model
606 predictions concurrent with discharge series at target sites, whereas NHM predictions at the time of this
607 writing are available only through the year 2016. For a summary of process-guided deep learning
608 strategies, see the “Integrating Design” subsection of Appling et al. (2022).

609

610 We caution that evaluation scores for both NEON’s published estimates and ours are computed on a small
611 fraction of each series for which both an estimate and a direct field measurement are available (39-213 per
612 site), and that measurements tend to be collected disproportionately at low flow. This often occurs for
613 practical reasons such as site access and technician safety, but may also reflect a need to characterize the
614 low-flow variability of the stage-discharge relationship in streams with unstable low-flow hydrologic
615 controls, such as unconsolidated bed material.

616

617 Whatever the reason for less sampling at high flow, any model attempting to use field measurements to
618 reconstruct continuous discharge will estimate with greater uncertainty at high flow than at low, and users
619 of our composite discharge product should observe uncertainties associated with estimates from all
620 methods. Mechanistic models that proceed from physical principles, or data-driven approaches that can
621 generalize from prior observations, do not in principle suffer this disadvantage, as they do not depend on
622 observations from a target site. However, these approaches may not reliably generate strong predictions at
623 all sites or under all conditions (Razavi & Coulibaly 2013; Kratzert et al. 2019b), and may produce erratic
624 point estimates where conditions diverge from past observations. Hybrid approaches that successfully
625 leverage field measurements, as well as physical principles or learned relationships, are likely to yield
626 well-constrained predictions where our efforts did not.

627 This study demonstrates that, in proximity to established streamflow gauges, even simple statistical
628 methods can be used to generate accurate, continuous discharge at “virtual gauges,” where discrete
629 discharge has been measured. The number of field measurements across sites in this study varies from 39
630 to 213, but the number required for virtual gauging may be substantially smaller even than the minimum
631 of this range. If the discharge relationships between a target site and all donor gauges were perfectly linear
632 or log-linear, they could in principle be established with only two precise measurements at the target site.
633 More important than the quantity is the distribution of measurements across flow conditions, which
634 should be sufficient to fully characterize all modeled discharge relationships and their linearity or lack
635 thereof (Sauer 2002; Zakwan et al. 2017). Concretely, we advocate for “storm chasing,” or
636 disproportionately seeking to sample discharge under high-flow conditions, and during both rising and
637 falling limbs of storm events, rather than routine sampling. Observed NEON flow conditions relative to
638 predicted discharge can be seen in Figure A9. See Philip & McLaughlin (2018) for further commentary
639 on establishing a virtual gauge network, and Seibert & Beven (2009) and Pool & Seibert (2021) for
640 information on the number and statistical properties of discharge samples required to establish strong
641 stage-discharge or discharge-discharge relationships.

642

643 **5. Conclusions**

644

645 Using linear regression on donor gauge data and LSTM-RNNs, we reconstructed continuous discharge at
646 5-minute and/or daily frequency for the 27 stream and river monitoring locations of the National
647 Ecological Observatory Network (NEON) over the water years 2015-2022. Relative to field-measured
648 discharge as ground truth, our estimates achieve higher Kling-Gupta efficiency than NEON's official
649 continuous discharge at 11 sites. We also provide continuous discharge estimates for ~199 site-months for
650 which no official values have been published. Estimates from this study can be used in conjunction with
651 officially released NEON continuous discharge data to enhance the analytical potential of NEON's river
652 and stream data products during its early operational phase. Toward that end, we provide composite
653 discharge series for each site, incorporating the best available estimates across all methods used in this
654 study and NEON's published estimates. Considering the lag of up to 2.5 years before provisional
655 discharge data become fully quality controlled and officially released by NEON, our methods may also be
656 used to increase the rate at which discharge-associated stream chemistry, dissolved gas, and water quality
657 products become fully usable by the community. All data and results from this study can be downloaded
658 from the Figshare collection at <https://doi.org/10.6084/m9.figshare.c.6488065>. Composite series can be
659 visualized interactively at https://macrosheds.org/data/vlah_etal_2023_composites/. All code necessary to
660 reproduce this analysis is archived at <https://doi.org/10.5281/zenodo.10067683>. A complete list of
661 products and URLs can be found in Table A3.

662 In general, linear regression methods produced more accurate discharge estimates (median KGE: 0.79;
663 median NSE: 0.81; $n = 24$ sites) than LSTM approaches due to the fact that regression models were able
664 to fully leverage available field measurements as well as highly informative donor gauge data.
665 Nonetheless, LSTM methods achieved median ensemble KGE of 0.71 and NSE of 0.56 across 18 sites,
666 making their estimates a valuable supplement. Although LSTM-generated discharge series are of daily
667 frequency, some users will prefer them to higher resolution regression estimates, as the latter may be
668 subject to error in the event of highly localized precipitation events affecting either donor or target
669 gauges, but not both.

670 Improvements to our design could be made in several ways. LSTM models could be exposed to additional
671 training data, such as the recently published Caravan compendium of CAMELS offshoots (Kratzert et al.
672 2023) or future expansions of the MacroSheds dataset (Vlah et al. 2023). Neural networks trained on
673 sub-daily inputs might be better equipped to exploit atmospheric-hydrological dynamics that respond to
674 both daily and annual cycles. Linear regression methods too might be improved with the use of additional
675 predictors, such as continuous water level or precipitation.

676 The success of simple statistical methods in generating high-quality continuous discharge time series
677 demonstrates the viability of "virtual gauges," or locations at which a small number of field discharge
678 measurements, in proximity to one or more established gauges, provide a basis for continuous discharge
679 estimation in lieu of a gauging station. Virtual gauges have the potential to greatly expand the spatial
680 coverage of continuous discharge data throughout the USA and any richly gauged region of the world.

681 **Author contribution**

682 MRVR, ESB, and MJV originated the project and identified its goals and methods. MJV carried out all
 683 analyses and drafted the manuscript. SR assisted in data collection. All authors took part in steering the
 684 project and editing the manuscript.

685

686 **Acknowledgements**

687

688 The authors are grateful to the NeuralHydrology team for their efforts in democratizing deep learning for
 689 the hydrology community. We thank NEON, NCAR, NWIS, Niwot Ridge LTER, Andrews Forest LTER,
 690 and the USGS for generating the data, and the National Science Foundation for providing the funding that
 691 made this analysis possible. Special thank-yous to Dr. Parker Norton of the USGS for extracting all
 692 NHM-PRMS outputs used in this study.

693

694 The National Ecological Observatory Network is a program sponsored by the National Science
 695 Foundation and operated under cooperative agreement by Battelle. This material is based in part upon
 696 work supported by the National Science Foundation through the NEON Program.

697

698 **Code availability**

699

700 All project code is on GitHub at https://github.com/vlahm/neon_q_sim.

701 The code repository is archived on Zenodo: <https://doi.org/10.5281/zenodo.10067683>

702

703 **Data availability**

704

705 All model input, output, and diagnostics are archived on Figshare:

706 <https://doi.org/10.6084/m9.figshare.c.6488065.v1>. See Tables A2 and A3 for details.

707

708 **Competing interests**

709

710 The authors declare that they have no conflict of interest.

711

712 **References**

713

714 Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment
 715 attributes and meteorology for large-sample studies, *Hydrology and Earth System Sciences*, 21,
 716 5293–5313, <http://dx.doi.org/10.5194/hess-21-5293-2017>, 2017.

717 Appelhans, T., Detsch, F., Reudenbach, C., and Woellauer, S.: mapview: Interactive Viewing of
 718 Spatial Data in R, <https://CRAN.R-project.org/package=mapview>, 2022.

719 Appling, A. P., Oliver, S. K., Read, J. S., Sadler, J. M., and Zwart, J.: Machine learning for
 720 understanding inland water quantity, quality, and ecology,
 721 <http://dx.doi.org/10.1016/B978-0-12-819166-8.00121-3>, 2022.

722 Arriagada, P., Karelavic, B., and Link, O.: Automatic gap-filling of daily streamflow time series in
 723 data-scarce regions using a machine learning algorithm, *Journal of Hydrology*, 598, 126454,
 724 <http://dx.doi.org/10.1016/j.jhydrol.2021.126454>, 2021.

725 Arsenaault, R., Brissette, F., and Martel, J.-L.: The hazards of split-sample validation in hydrological
 726 model calibration, *Journal of hydrology*, 566, 346–362,
 727 <http://dx.doi.org/10.1016/j.jhydrol.2018.09.027>, 2018.

- 728 Arsov, N. and Mirceva, G.: Network Embedding: An Overview,
729 <https://doi.org/10.48550/ARXIV.1911.11726>, 2019.
- 730 Aschonitis, V. G., Papamichail, D., Demertzi, K., Colombani, N., Mastrocicco, M., Ghirardini, A.,
731 Castaldelli, G., and Fano, E.-A.: High resolution global grids of revised Priestley-Taylor and
732 Hargreaves-Samani coefficients for assessing ASCE-standardized reference crop
733 evapotranspiration and solar radiation, links to ESRI-grid files,
734 <https://doi.org/10.1594/PANGAEA.868808>, 2017.
- 735 Benson, M. A. and Dalrymple, T.: General field and office procedures for indirect discharge
736 measurements, US Govt. Print. Off., <https://doi.org/10.3133/twri03A1>, 1967.
- 737 Bergstra, J. and Bengio, Y.: Random search for hyper-parameter optimization., Journal of machine
738 learning research, 13, <http://dblp.uni-trier.de/db/journals/jmlr/jmlr13.html#BergstraB12>, 2012.
- 739 Bukaveckas, P., Likens, G., Winter, T., and Buso, D.: A comparison of methods for deriving solute
740 flux rates using long-term data from streams in the Mirror Lake watershed, Water, Air, and Soil
741 Pollution, 105, 277–293, http://dx.doi.org/10.1007/978-94-017-0906-4_26, 1998.
- 742 Caruana, R.: Multitask learning, Springer, http://dx.doi.org/10.1007/978-1-4615-5529-2_5, 1998.
- 743 Chokmani, K. and Ouarda, T. B.: Physiographical space-based kriging for regional flood frequency
744 estimation at ungauged sites, Water Resources Research, 40,
745 <http://dx.doi.org/10.1029/2003WR002983>, 2004.
- 746 DeCicco, L. A., Lorenz, D., Hirsch, R. M., Watkins, W., and Johnson, M.: dataRetrieval: R packages
747 for discovering and retrieving water data available from U.S. federal hydrologic web services,
748 <https://doi.org/10.5066/P9X4L3GE>, 2022.
- 749 Durand, M., Gleason, C. J., Pavelsky, T. M., de Moraes Frasson, R. P., Turmon, M. J., David, C. H.,
750 Altenau, E. H., Tebaldi, N., Larnier, K., Monnier, J., and others: A framework for estimating
751 global river discharge from the Surface Water and Ocean Topography satellite mission,
752 Authorea Preprints, <http://dx.doi.org/10.1029/2021WR031614>, 2022.
- 753 Friedman, J., Tibshirani, R., and Hastie, T.: Regularization Paths for Generalized Linear Models via
754 Coordinate Descent, Journal of Statistical Software, 33, 1–22,
755 <https://doi.org/10.18637/jss.v033.i01>, 2010.
- 756 Galton, F.: Regression towards mediocrity in hereditary stature., The Journal of the Anthropological
757 Institute of Great Britain and Ireland, 15, 246–263, <http://dx.doi.org/10.2307/2841583>, 1886.
- 758 Goeman, J., Meijer, R., and Chaturvedi, N.: L1 and L2 penalized regression models, cran. r-project.
759 or, <https://cran.r-project.org/web/packages/penalized/vignettes/penalized.pdf>, 2012.
- 760 Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., and Moore, R.: Google Earth
761 Engine: Planetary-scale geospatial analysis for everyone, Remote Sensing of Environment,
762 <https://doi.org/10.1016/j.rse.2017.06.031>, 2017.
- 763 Graf, W. H.: Hydraulics of sediment transport, Water Resources Publication, ISBN 13:
764 978-1-887201-57-5, 1984.
- 765 Gruber, M.: Improving efficiency by shrinkage: The James–Stein and Ridge regression estimators,
766 Routledge, <http://dx.doi.org/10.1201/9780203751220>, 2017.
- 767 Guo, D., Johnson, F., and Marshall, L.: Assessing the potential robustness of conceptual
768 rainfall-runoff models under a changing climate, Water Resources Research, 54, 5030–5049,
769 <http://dx.doi.org/10.1029/2018WR022636>, 2018.
- 770 Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared
771 error and NSE performance criteria: Implications for improving hydrological modelling,
772 Journal of hydrology, 377, 80–91, <http://dx.doi.org/10.1016/j.jhydrol.2009.08.003>, 2009.
- 773 Hall Jr, R. O.: Metabolism of streams and rivers: Estimation, controls, and application, in: Stream
774 ecosystems in a changing environment, Elsevier, 151–180,
775 <https://doi.org/10.1016/B978-0-12-405890-3.00004-X>, 2016.
- 776 Harvey, C. L., Dixon, H., and Hannaford, J.: An appraisal of the performance of data-infilling
777 methods for application to daily mean river flow records in the UK, Hydrology Research, 43,
778 618–636, <http://dx.doi.org/10.2166/nh.2012.110>, 2012.

- 779 Hirsch, R. M.: A comparison of four streamflow record extension techniques, *Water Resources*
780 *Research*, 18, 1081–1088, <http://dx.doi.org/10.1029/WR018i004p01081>, 1982.
- 781 Hochreiter, S. and Schmidhuber, J.: Long short-term memory, *Neural computation*, 9, 1735–1780,
782 <http://dx.doi.org/10.1162/neco.1997.9.8.1735>, 1997.
- 783 Horner, I., Renard, B., Le Coz, J., Branger, F., McMillan, H., and Pierrefeu, G.: Impact of stage
784 measurement errors on streamflow uncertainty, *Water Resources Research*, 54, 1952–1976,
785 <http://dx.doi.org/10.1002/2017WR022039>, 2018.
- 786 Hsu, K., Gupta, H. V., and Sorooshian, S.: Artificial neural network modeling of the rainfall-runoff
787 process, *Water resources research*, 31, 2517–2530, <http://dx.doi.org/10.1029/95WR01955>,
788 1995.
- 789 Isaacson, K. and Coonrod, J.: USGS streamflow data and modeling sand-bed rivers, *Journal of*
790 *Hydraulic Engineering*, 137, 847–851,
791 [http://dx.doi.org/10.1061/\(ASCE\)HY.1943-7900.0000362](http://dx.doi.org/10.1061/(ASCE)HY.1943-7900.0000362), 2011.
- 792 Johnson, S. L., Rothacher, J. S., and Wondzell, S. M.: Stream discharge in gaged watersheds at the
793 HJ Andrews Experimental Forest, 1949 to present,
794 <https://doi.org/10.6073/PASTA/0066D6B04E736AF5F234D95D97EE84F3>, 2020.
- 795 Kingma, D. P. and Ba, J.: Adam: A method for stochastic optimization, arXiv preprint
796 arXiv:1412.6980, <https://doi.org/10.48550/arXiv.1412.6980>, 2014.
- 797 Knoben, W. J., Freer, J. E., and Woods, R. A.: Inherent benchmark or not? Comparing
798 Nash–Sutcliffe and Kling–Gupta efficiency scores, *Hydrology and Earth System Sciences*, 23,
799 4323–4331, <http://dx.doi.org/10.5194/hess-23-4323-2019>, 2019.
- 800 Kratzert, F., Herrnegger, M., Klotz, D., Hochreiter, S., and Klambauer, G.:
801 *NeuralHydrology*—interpreting LSTMs in hydrology, *Explainable AI: Interpreting, explaining*
802 *and visualizing deep learning*, 347–362, http://dx.doi.org/10.1007/978-3-030-28954-6_19,
803 2019a.
- 804 Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward
805 improved predictions in ungauged basins: Exploiting the power of machine learning, *Water*
806 *Resources Research*, 55, 11344–11354, <http://dx.doi.org/10.1029/2019WR026065>, 2019b.
- 807 Kratzert, F., Gauch, M., Nearing, G., and Klotz, D.: *NeuralHydrology* — A Python library for Deep
808 Learning research in hydrology, *Journal of Open Source Software*, 7, 4050,
809 <https://doi.org/10.21105/joss.04050>, 2022.
- 810 Kratzert, F., Nearing, G., Addor, N., Erickson, T., Gauch, M., Gilon, O., Gudmundsson, L.,
811 Hassidim, A., Klotz, D., Nevo, S., and others: Caravan-A global community dataset for
812 large-sample hydrology, *Scientific Data*, 10, 61,
813 <http://dx.doi.org/10.1038/s41597-023-01975-w>, 2023.
- 814 Lunch, C., Laney, C., Mietkiewicz, N., Sokol, E., Cawley, K., and NEON (National Ecological
815 Observatory Network): *neonUtilities: Utilities for Working with NEON Data*,
816 <https://CRAN.R-project.org/package=neonUtilities>, 2022.
- 817 Manning, R.: On the flow of water in open channels and pipes, 20, 161–207, 1891.
- 818 Moore, S. A., Jamieson, E. C., Rainville, F., Rennie, C. D., and Mueller, D. S.: Monte Carlo
819 approach for uncertainty analysis of acoustic Doppler current profiler discharge measurement
820 by moving boat, *Journal of Hydraulic Engineering*, 143, 04016088,
821 [http://dx.doi.org/10.1061/\(ASCE\)HY.1943-7900.0001249](http://dx.doi.org/10.1061/(ASCE)HY.1943-7900.0001249), 2017.
- 822 Moriasi, D., Gitau, M., Pai, N., and Daggupati, P.: Hydrologic and Water Quality Models:
823 Performance Measures and Evaluation Criteria, *Transactions of the ASABE (American Society*
824 *of Agricultural and Biological Engineers)*, 58, 1763–1785,
825 <https://doi.org/10.13031/trans.58.10715>, 2015.
- 826 Muggeo, V. M. R.: segmented: an R Package to Fit Regression Models with Broken-Line
827 Relationships., *R News*, 8, 20–25, <https://cran.r-project.org/doc/Rnews/>, 2008.

- 828 Nalley, D., Adamowski, J., Khalil, B., and Biswas, A.: A comparison of conventional and wavelet
829 transform based methods for streamflow record extension, *Journal of Hydrology*, 582, 124503,
830 <http://dx.doi.org/10.1016/j.jhydrol.2019.124503>, 2020.
- 831 Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models part I—A
832 discussion of principles, *Journal of hydrology*, 10, 282–290,
833 [http://dx.doi.org/10.1016/0022-1694\(70\)90255-6](http://dx.doi.org/10.1016/0022-1694(70)90255-6), 1970.
- 834 NEON (National Ecological Observatory Network): Continuous discharge (DP4.00130.001),
835 RELEASE-2023, <https://doi.org/10.48443/H2ZE-2F12>. Data accessed from
836 <https://data.neonscience.org/data-products/DP1.00130.001/RELEASE-2023> on May 5, 2023.,
837 2023a.
- 838 NEON (National Ecological Observatory Network): Discharge field collection (DP1.20048.001),
839 PROVISIONAL, <https://dx.doi.org/10.6084/m9.figshare.22344589>. Data accessed from
840 <https://data.neonscience.org/data-products/DP1.20048.001> on January 1, 2023., 2023b.
- 841 NEON (National Ecological Observatory Network): Discharge field collection (DP1.20048.001),
842 RELEASE-2023, <https://doi.org/10.48443/TYS0-ZE83>. Data accessed from
843 <https://data.neonscience.org/data-products/DP1.20048.001/RELEASE-2023> on January 1,
844 2023., 2023c.
- 845 Newman, A., Sampson, K., Clark, M., Bock, A., Viger, R., and Blodgett, D.: A large-sample
846 watershed-scale hydrometeorological dataset for the contiguous USA, UCAR/NCAR: Boulder,
847 CO, USA, <https://dx.doi.org/10.5065/D6MW2F4D>, 2014.
- 848 Newman, A., Clark, M., Sampson, K., Wood, A., Hay, L., Bock, A., Viger, R., Blodgett, D., Brekke,
849 L., Arnold, J., and others: Development of a large-sample watershed-scale hydrometeorological
850 data set for the contiguous USA: data set characteristics and assessment of regional variability
851 in hydrologic model performance, *Hydrology and Earth System Sciences*, 19, 209–223,
852 <https://doi.org/doi:10.5194/hess-19-209-2015>, 2015.
- 853 Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., and Nearing, G.:
854 Benchmarking of a physically based hydrologic model, *Journal of Hydrometeorology*, 18,
855 2215–2225, <http://dx.doi.org/10.1175/JHM-D-16-0284.1>, 2017.
- 856 Odum, H. T.: Primary production in flowing waters 1, *Limnology and oceanography*, 1, 102–117,
857 <https://doi.org/10.4319/lo.1956.1.2.0102>, 1956.
- 858 Pantelakis, D., Doulgeris, C., Hatzigiannakis, E., and Arampatzis, G.: Evaluation of discharge
859 measurements methods in a natural river of low or middle flow using an electromagnetic flow
860 meter, *River Research and Applications*, 38, 1003–1013, <http://dx.doi.org/10.1002/rra.3966>,
861 2022.
- 862 Philip, E. and McLaughlin, J.: Evaluation of stream gauge density and implementing the concept of
863 virtual gauges in Northern Ontario for watershed modeling, *Journal of Water Management
864 Modeling*, <http://dx.doi.org/10.14796/JWMM.C438>, 2018.
- 865 Pool, S. and Seibert, J.: Gauging ungauged catchments—Active learning for the timing of point
866 discharge observations in combination with continuous water level measurements, *Journal of
867 Hydrology*, 598, 126448, <http://dx.doi.org/10.1016/j.jhydrol.2021.126448>, 2021.
- 868 Razavi, T. and Coulibaly, P.: Streamflow prediction in ungauged basins: review of regionalization
869 methods, *Journal of hydrologic engineering*, 18, 958–975,
870 [http://dx.doi.org/10.1061/\(ASCE\)HE.1943-5584.0000690](http://dx.doi.org/10.1061/(ASCE)HE.1943-5584.0000690), 2013.
- 871 Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., Karpatne, A., Hansen, G. J.
872 A., Hanson, P. C., Watkins, W., Steinbach, M., and Kumar, V.: Process-Guided Deep Learning
873 Predictions of Lake Water Temperature, *Water Resources Research*, 55, 9173–9190,
874 <https://doi.org/10.1029/2019WR024922>, 2019.
- 875 Regan, R. S., Juracek, K. E., Hay, L. E., Markstrom, S., Viger, R. J., Driscoll, J. M., LaFontaine, J.,
876 and Norton, P. A.: The US Geological Survey National Hydrologic Model infrastructure:
877 Rationale, description, and application of a watershed-scale model for the conterminous United

- 878 States, *Environmental Modelling & Software*, 111, 192–203,
879 <https://doi.org/10.1016/j.envsoft.2018.09.023>, 2019.
- 880 Rhea, S.: NEON Continuous Discharge Evaluation,
881 <https://doi.org/10.4211/hs.03c52d47d66e40f4854da8397c7d9668>, 2023.
- 882 Rhea, S., Vlah, M., Slaughter, W., and Gubbins, N.: macrosheds: Tools for interfacing with the
883 MacroSheds dataset, <https://github.com/MacroSHEDS/macrosheds>, 2023a.
- 884 Rhea, S., Gubbins, N., DelVecchia, A. G., Ross, M. R., and Bernhardt, E. S.: User-focused
885 evaluation of National Ecological Observatory Network streamflow estimates, *Scientific Data*,
886 10, 89, <http://dx.doi.org/10.1038/s41597-023-02026-0>, 2023b.
- 887 Sadler, J. M., Appling, A. P., Read, J. S., Oliver, S. K., Jia, X., Zwart, J. A., and Kumar, V.:
888 Multi-task deep learning of daily streamflow and water temperature, *Water Resources Research*,
889 58, e2021WR030138, <http://dx.doi.org/10.1029/2021WR030138>, 2022.
- 890 Sauer, V. B.: Standards for the analysis and processing of surface-water data and information using
891 electronic methods, US Geological Survey, 2002.
- 892 Seibert, J. and Beven, K. J.: Gauging the ungauged basin: how many discharge measurements are
893 needed?, *Hydrology and Earth System Sciences*, 13, 883–892,
894 <http://dx.doi.org/10.5194/hess-13-883-2009>, 2009.
- 895 Seibert, J., Vis, M. J. P., Lewis, E., and van Meerveld, H. J.: Upper and lower benchmarks in
896 hydrological modelling, *Hydrological Processes*, 32, 1120–1125,
897 <https://doi.org/10.1002/hyp.11476>, 2018.
- 898 Seibert, J., Strobl, B., Etter, S., Hummer, P., and van Meerveld, H. J. (Ilja): Virtual Staff Gauges for
899 Crowd-Based Stream Level Observations, *Frontiers in Earth Science*, 7,
900 <https://doi.org/10.3389/feart.2019.00070>, 2019.
- 901 Shen, H., Tolson, B. A., and Mai, J.: Time to update the split-sample approach in hydrological model
902 calibration, *Water Resources Research*, 58, e2021WR031523,
903 <http://dx.doi.org/10.1029/2021WR031523>, 2022.
- 904 Shen, J.: Discharge characteristics of triangular-notch thin-plate weirs, United States Department of
905 the Interior, Geological Survey, 1981.
- 906 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.: Dropout: a simple
907 way to prevent neural networks from overfitting, *The journal of machine learning research*, 15,
908 1929–1958, <https://doi.org/10.5555/2627435.2670313>, 2014.
- 909 Tazioli, A.: Experimental methods for river discharge measurements: comparison among tracers and
910 current meter, *Hydrological Sciences Journal*, 56, 1314–1324,
911 <http://dx.doi.org/10.1080/02626667.2011.607822>, 2011.
- 912 Thornton, M. M., Shrestha, R., Wei, Y., Thornton, P. E., Kao, S.-C., and Wilson, B. E.: Daymet:
913 Daily Surface Weather Data on a 1-km Grid for North America, Version 4 R1,
914 https://daac.ornl.gov/cgi-bin/dsvviewer.pl?ds_id=2129,
915 <https://doi.org/10.3334/ORNLDAAAC/2129>, 2022.
- 916 Turnipseed, D. P. and Sauer, V. B.: Discharge measurements at gaging stations, US Geological
917 Survey, 2010.
- 918 Van Rossum, G. and Drake, F. L.: Python 3 Reference Manual, CreateSpace, Scotts Valley, CA,
919 ISBN 1-4414-1269-7, 2009.
- 920 Vlah, M. J., Rhea, S., Bernhardt, E. S., Slaughter, W., Gubbins, N., DelVecchia, A. G., Thellman, A.,
921 and Ross, M. R.: MacroSheds: A synthesis of long-term biogeochemical, hydroclimatic, and
922 geospatial data from small watershed ecosystem studies, *Limnology and Oceanography Letters*,
923 <http://dx.doi.org/10.1002/lo12.10325>, 2023.
- 924 Wang, C. P.: Laser doppler velocimetry, *Journal of Quantitative Spectroscopy and Radiative*
925 *Transfer*, 40, 309–319, [http://dx.doi.org/10.1016/0022-4073\(88\)90122-7](http://dx.doi.org/10.1016/0022-4073(88)90122-7), 1988.
- 926 White, A. F. and Blum, A. E.: Effects of climate on chemical_ weathering in watersheds,
927 *Geochimica et Cosmochimica Acta*, 59, 1729–1747,
928 [http://dx.doi.org/10.1016/0016-7037\(95\)00078-E](http://dx.doi.org/10.1016/0016-7037(95)00078-E), 1995.

- 929 Whittaker, J., Whitehead, C., and Somers, M.: The neglog transformation and quantile regression for
 930 the analysis of a large credit scoring database, *Journal of the Royal Statistical Society: Series C*
 931 (Applied Statistics), 54, 863–878, <http://dx.doi.org/10.1111/j.1467-9876.2005.00520.x>, 2005.
 932 Zakwan, M., Muzzammil, M., and Alam, J.: Developing stage-discharge relations using optimization
 933 techniques, *Aquademia: Water, Environment and Technology*, 1, 05,
 934 <http://dx.doi.org/10.20897/awet/81286>, 2017.

935

936

937 **Appendix A**

938

939 **Tables**

940

941 Table A1: Methods from this study used in the construction of composite discharge series. Composite
 942 series also incorporate NEON continuous discharge product DP4.00130.001 (NEON 2023a). “Linreg” =
 943 linear regression; “glmnet” = ridge regression; “lm” = OLS regression; “segmented” = segmented
 944 regression; “abs” = absolute discharge; “spec” = specific discharge; “pgdl” = process-guided deep
 945 learning.

946

Site	KGE linreg	NSE linreg	Method linreg	KGE LSTM	NSE LSTM	Method LSTM
FLNT	0.989	0.980	glmnet_spec	0.664	0.507	generalist
TOMB	0.970	0.993	glmnet_abs			
HOPB	0.966	0.937	lm_abs	0.852	0.704	generalist
BLUE	0.962	0.932	lm_spec	0.746	0.567	specialist
REDB	0.946	0.973	lm_abs	0.511	0.551	generalist_pgdl
KING	0.935	0.888	glmnet_abs			
LEWI	0.929	0.875	glmnet_abs	0.848	0.724	specialist
SYCA	0.919	0.938	segmented_spec			
MCDI	0.912	0.897	glmnet_spec			
LECO	0.877	0.833	lm_spec			
MCRA	0.868	0.866	glmnet_spec	0.723	0.531	generalist
MART	0.811	0.706	glmnet_spec	0.779	0.566	generalist
POSE	0.803	0.648	glmnet_spec			
MAYF	0.787	0.806	glmnet_abs	0.586	0.666	generalist
BLWA	0.779	0.892	glmnet_abs			
COMO	0.771	0.806	glmnet_composite_spec			
BLDE	0.744	0.863	glmnet_abs	0.744	0.687	generalist
CARI	0.721	0.637	glmnet_abs			
GUIL	0.692	0.653	glmnet_abs			
ARIK	0.674	0.596	glmnet_abs			
CUPE	0.663	0.728	glmnet_spec			

WALK	0.607	0.532	glmnet_spec			
BIGC				0.895	0.827	specialist
WLOU				0.778	0.596	generalist_pgdl
TECR				0.711	0.904	generalist
PRIN						
OKSR						

947

948

949 Table A2: Model input data used in this study.

Resource	Description	Source/Link
NEON discharge field collection	Discharge measurements from field-based surveys	NEON 2023b, NEON 2023c
NEON continuous discharge	Discharge calculated from a rating curve and sensor measurements of water level	NEON 2023a
User-focused evaluation of NEON streamflow estimates	3-tier classification of the reliability of NEON continuous discharge by site-month	https://www.nature.com/articles/s41597-023-01983-w
CAMELS dataset	Catchment Attributes, Meteorology, (and streamflow) for Large-sample Studies	https://ral.ucar.edu/solutions/products/camels
National Hydrologic Model (NHM)	USGS infrastructure that, when coupled with the Precipitation-Runoff Modeling System, can produce streamflow simulations at local to national scale	https://www.usgs.gov/mission-areas/water-resources/science/national-hydrologic-model-infrastructure
MacroSheds	A synthesis of long-term biogeochemical, hydroclimatic, and geospatial data from small watershed ecosystem studies	https://portal.edirepository.org/nis/mapbrowse?scope=edi&identifier=1262

Daymet	Gridded estimates of daily weather parameters	https://developers.google.com/earth-engine/datasets/catalog/NASA_ORNL_DAYMET_V4
HJ Andrews Experimental Forest stream discharge	Stream discharge in gaged watersheds, 1949 to present	https://portal.edirepository.org/nis/mapbrowse?packageid=knb-lter-and.4341.33
USGS National Water Information System	Streamflow and associated data for thousands of gauged streams and rivers within the USA	https://waterdata.usgs.gov/nwis , e.g. https://waterdata.usgs.gov/monitoring-location/06879100/

950

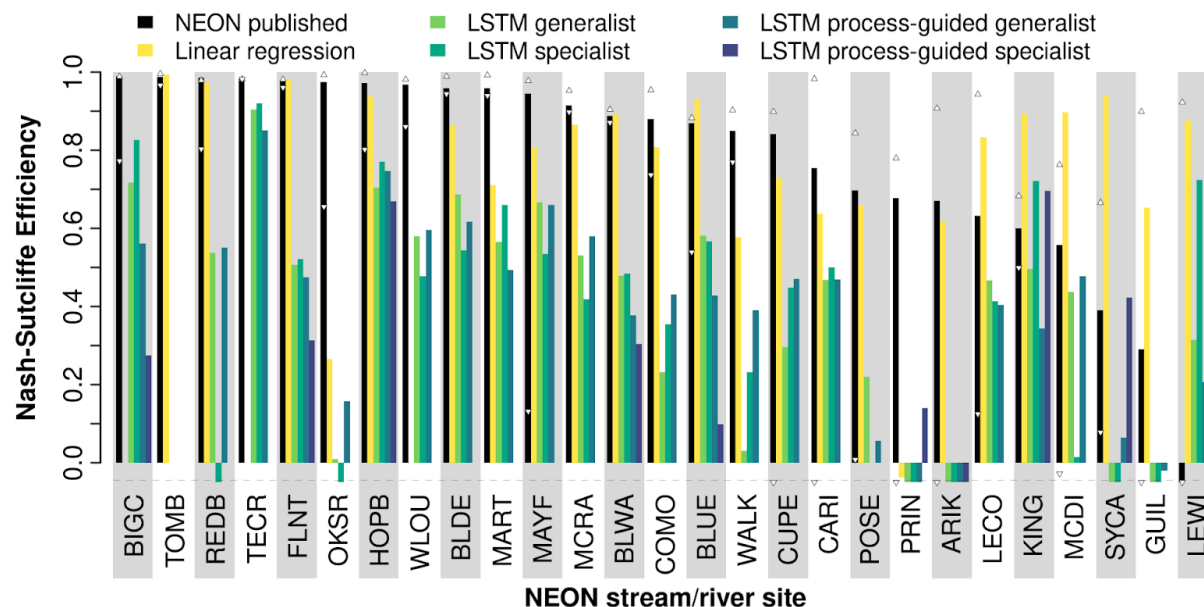
951 Table A3: Products of this study.

Product	Description	Link
Data archive landing page	Figshare page linking to each of four archives described below	https://doi.org/10.6084/m9.figshare.c.6488065
Composite discharge timeseries	Analysis-ready CSVs combining the best available discharge estimates across linear regression and LSTM approaches from this study, and NEON's published data	https://doi.org/10.6084/m9.figshare.23206592.v1
Composite discharge plots	Interactive plots of our composite discharge product	https://macrosheds.org/data/vlah_etal_2023_composites
All model outputs and results	Complete predictions from all linear regression and LSTM models, run results, and diagnostics	https://doi.org/10.6084/m9.figshare.22344589.v1
All model input data	Donor gauge streamflow, training data for LSTMs, model configurations, etc.	https://doi.org/10.6084/m9.figshare.22349377.v1
All code associated with this paper	Zenodo archive of GitHub repository	https://doi.org/10.5281/zenodo.10067683

All figures associated with this paper	High-resolution images of all figures from the main body and appendix	https://doi.org/10.6084/m9.figshare.23169362.v1
--	---	---

952

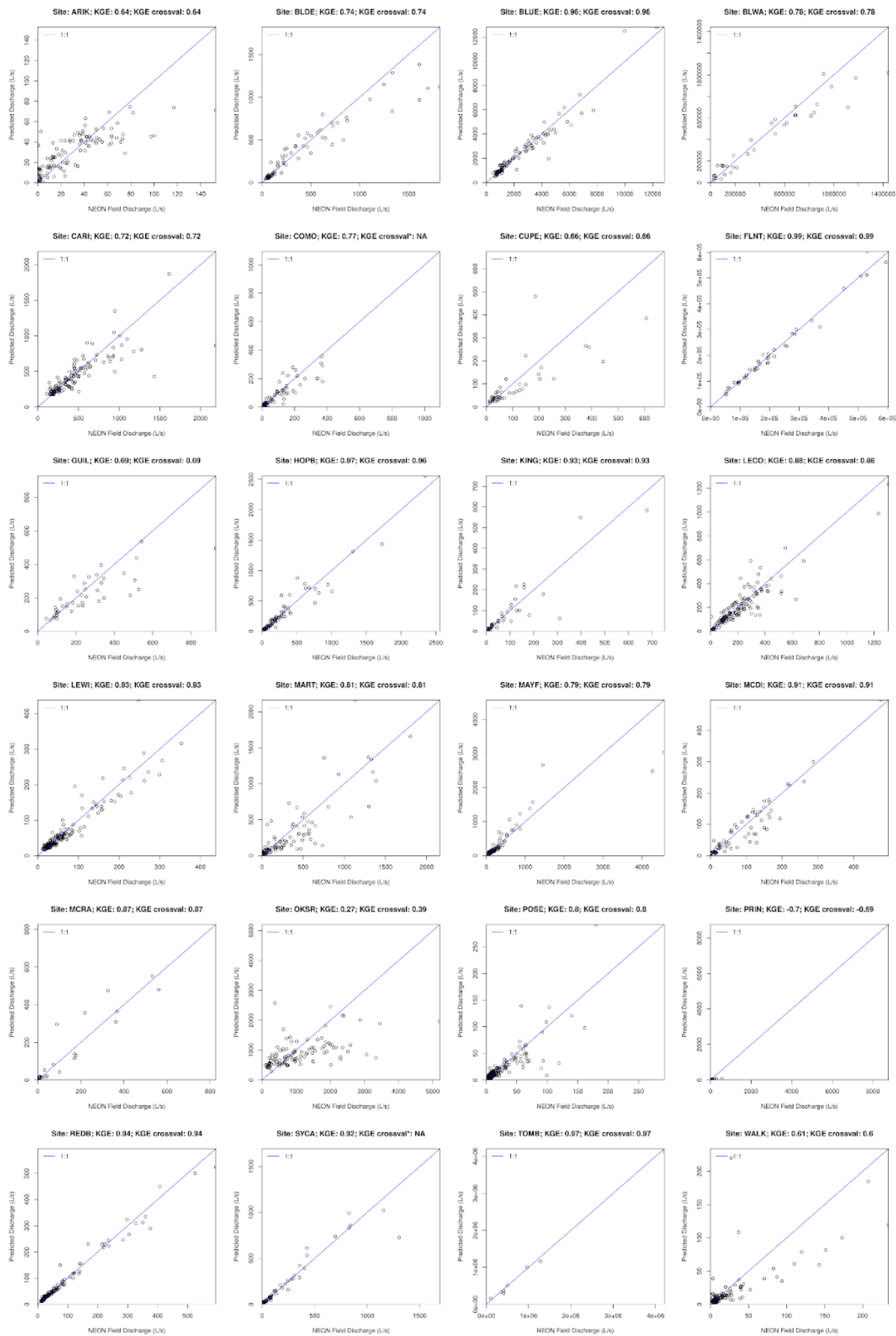
953

954 **Figures**

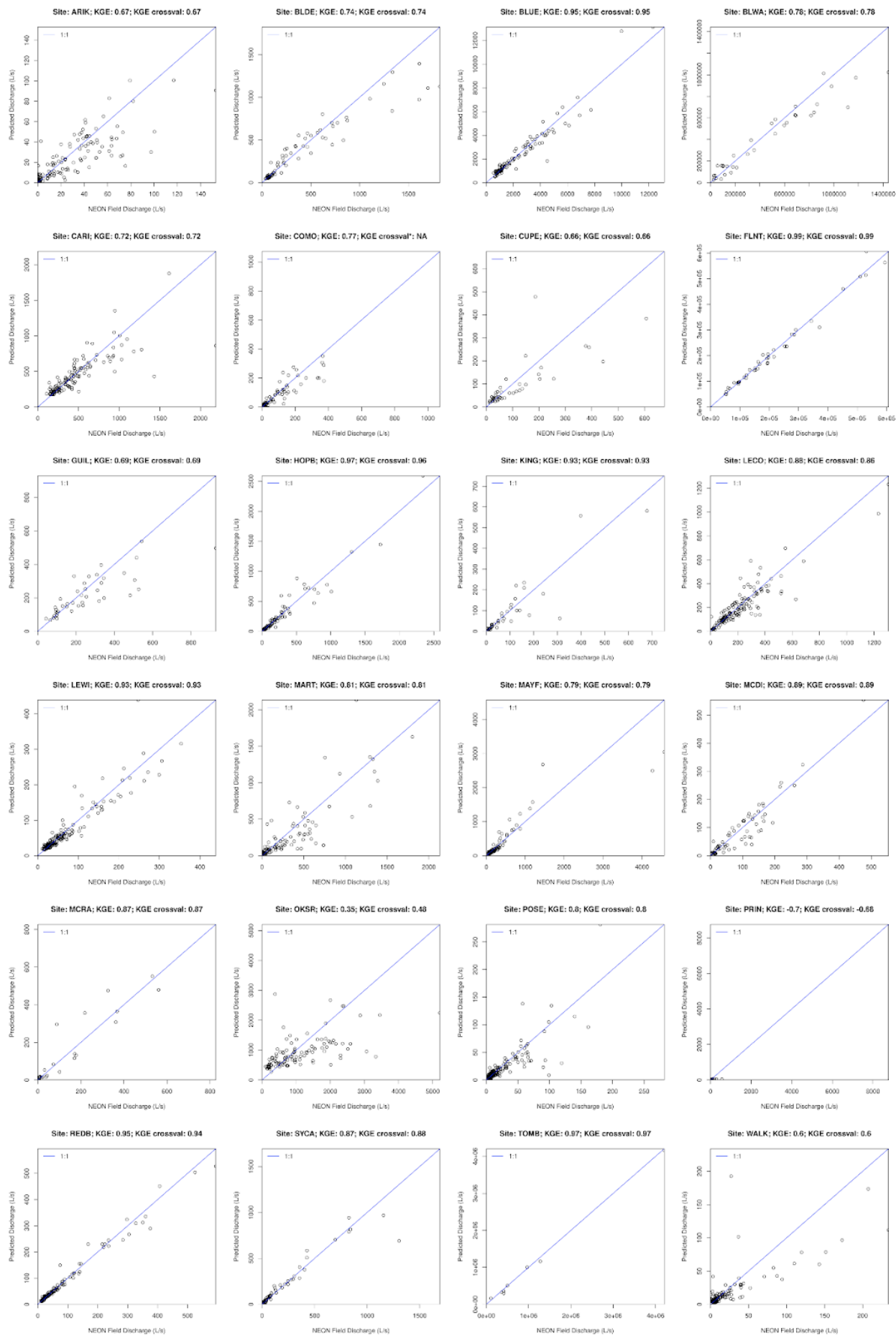
955

956 Figure A1: Efficiency of five stream discharge prediction methods and NEON's published continuous
 957 discharge product at 27 NEON gauge locations, versus field-measured discharge. Small, white triangles
 958 represent max/min NSE of published discharge by water year (Oct 1 through Sept 30) with at least 5 field
 959 measurements (or 2 for site OKSR). NSE was computed on all available observation-estimate pairs except
 960 those with quality flags (dischargeFinalQF or dischargeFinalQFSciRvw of 1).. For the best performing
 961 LSTM method, at all sites except TECR, FLNT, REDB, WALK, POSE, and KING, displayed NSE is
 962 averaged over 30 ensemble runs with identical hyperparameters. For the sites just named, performance of
 963 a chosen method, after ensembling, dropped below that of at least one other method's optimal NSE from
 964 parameter search. For all other LSTM site-method pairs, which were not ensembled, displayed
 965 performance is that of the best model trained during the parameter search phase. Sites are ordered by the
 966 NSE of NEON continuous discharge. See Table 3 for LSTM model definitions. NSE of 1 is a perfect
 967 prediction, while NSE of 0 is equivalent in skill to prediction from the mean. Negative values are
 968 truncated at -0.05 in this plot to improve visualization.

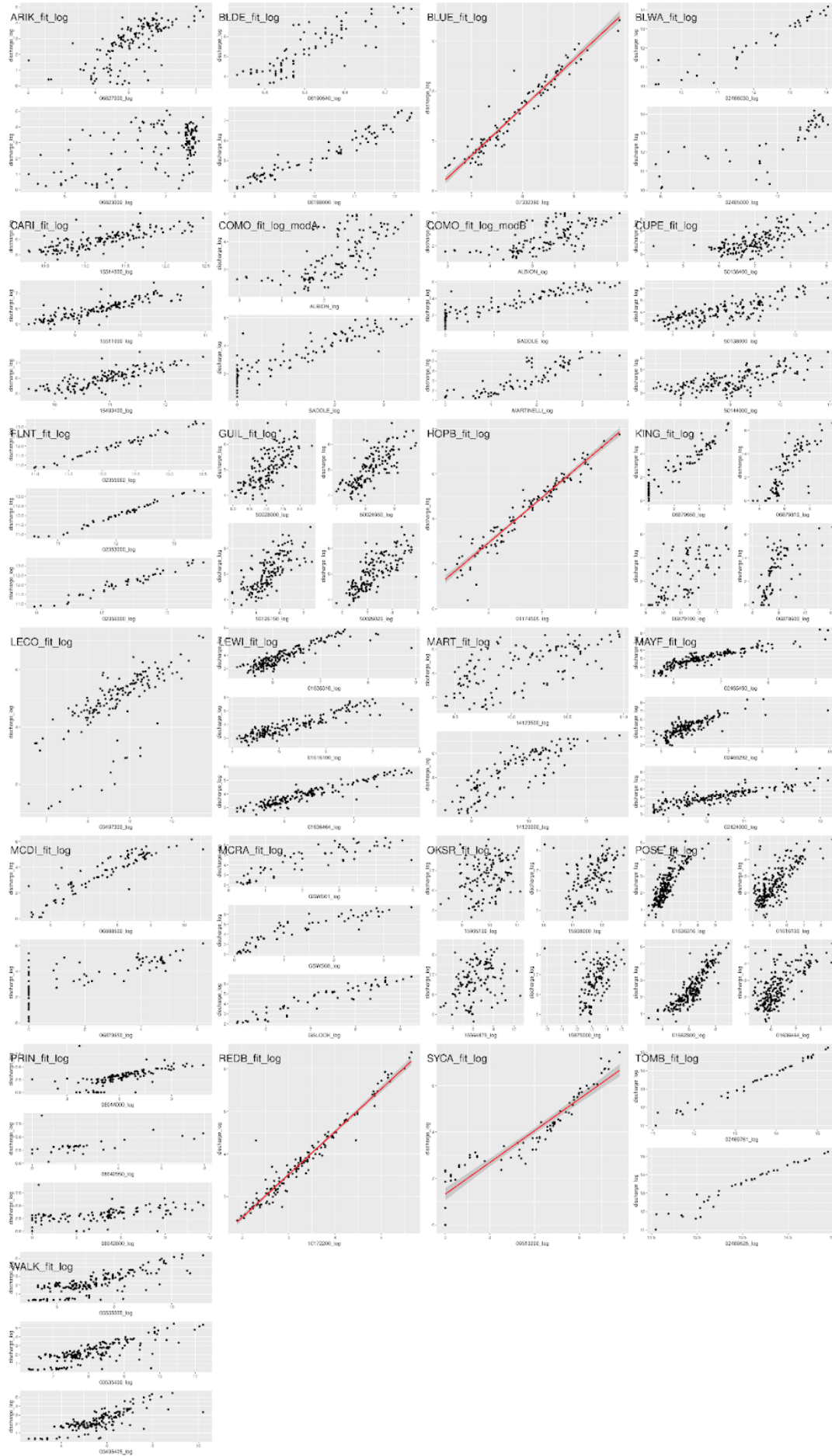
969



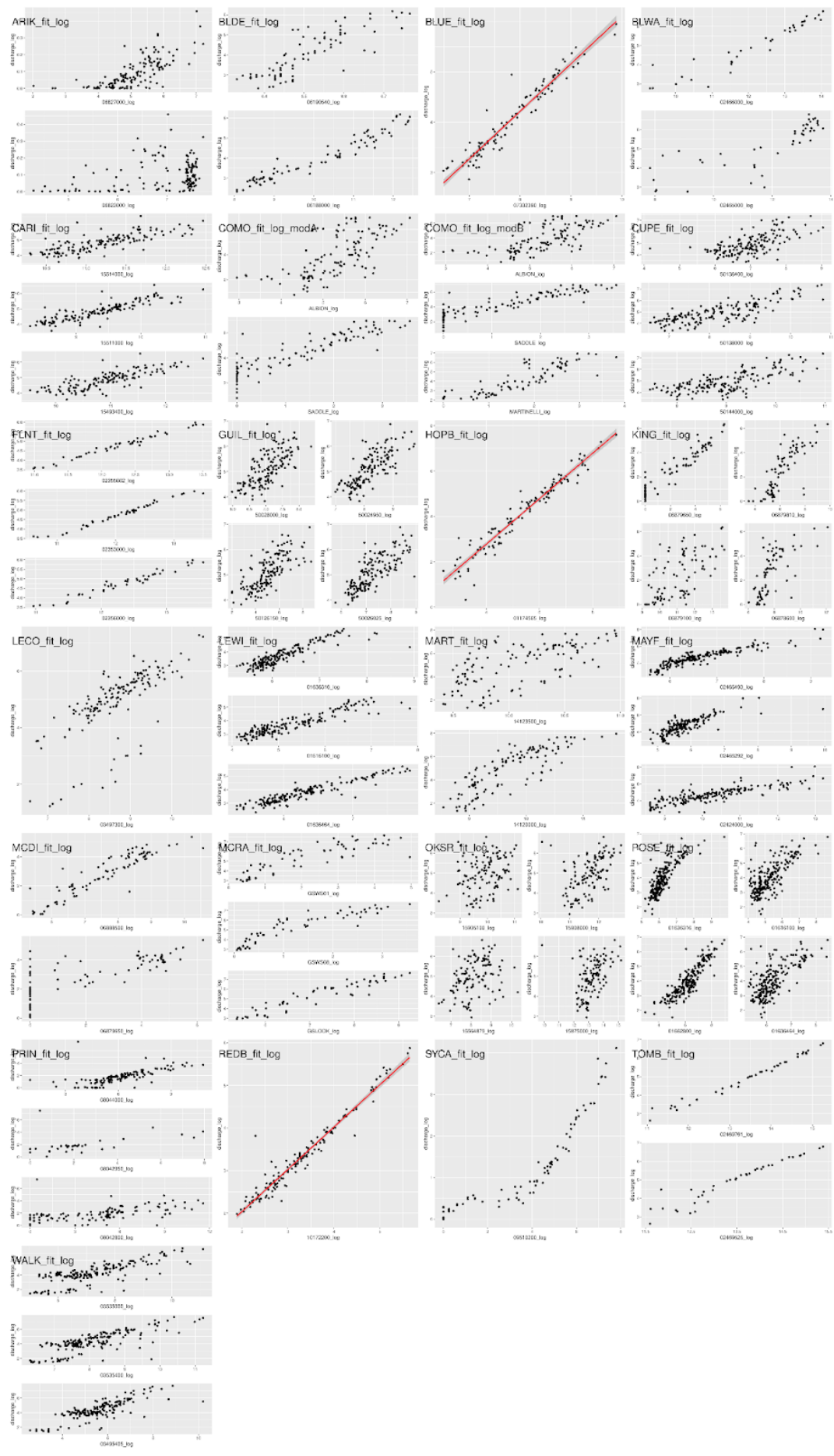
971 Figure A2: Observed (field) discharge vs. predictions from linear regression on specific discharge (i.e.
972 scaled by watershed area).
973



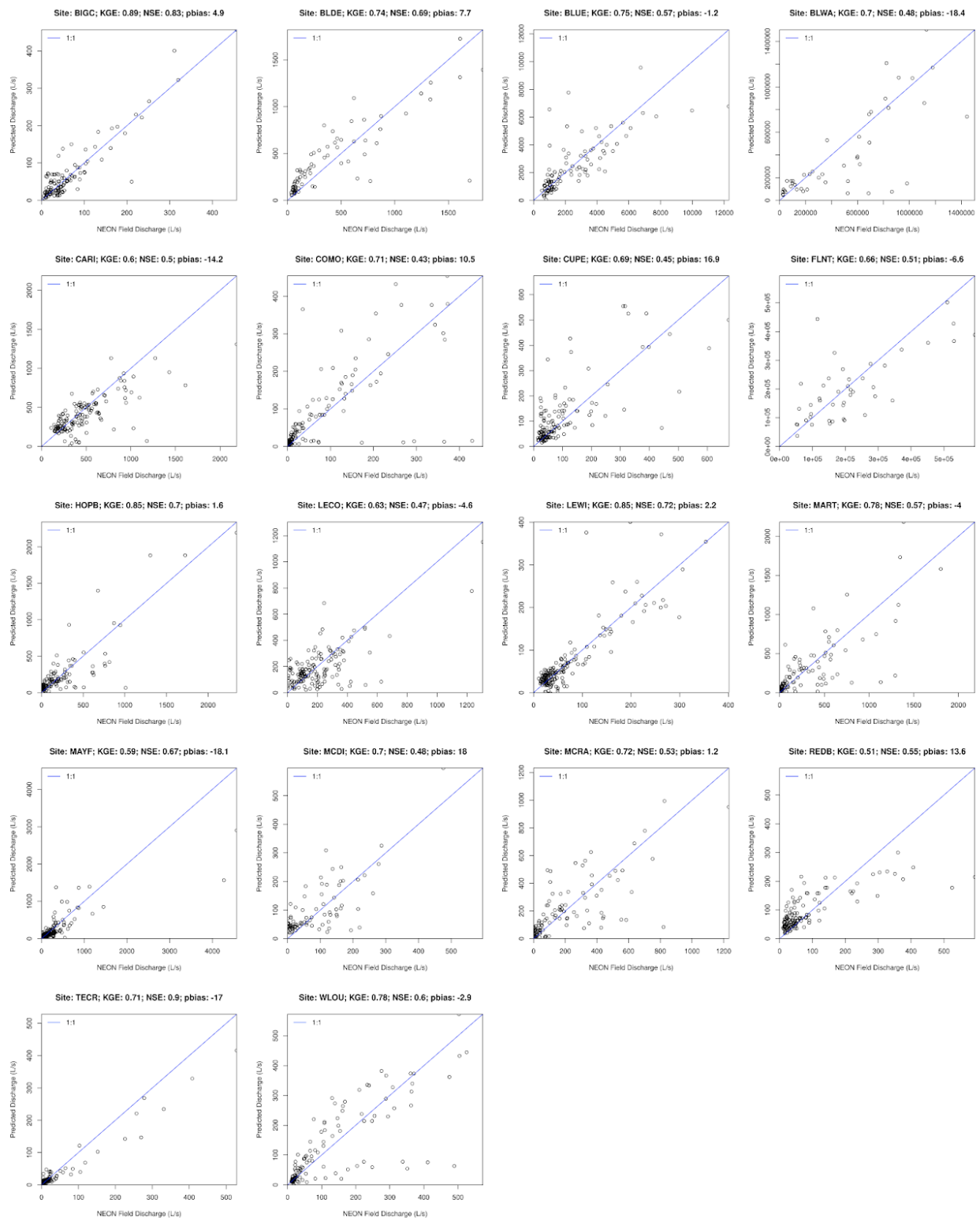
975 Figure A3: Observed (field) discharge vs. predictions from linear regression on absolute discharge (i.e.
976 not scaled by watershed area).
977



979 Figure A4: Marginal relationships between donor and target gauges for regression on specific discharge.
980 Regression lines are shown only for single-donor regressions, fitted via OLS. Site SYCA, here exhibiting
981 a breakpoint, was modeled with segmented regression, and thus the regression line shown has no
982 relevance.
983

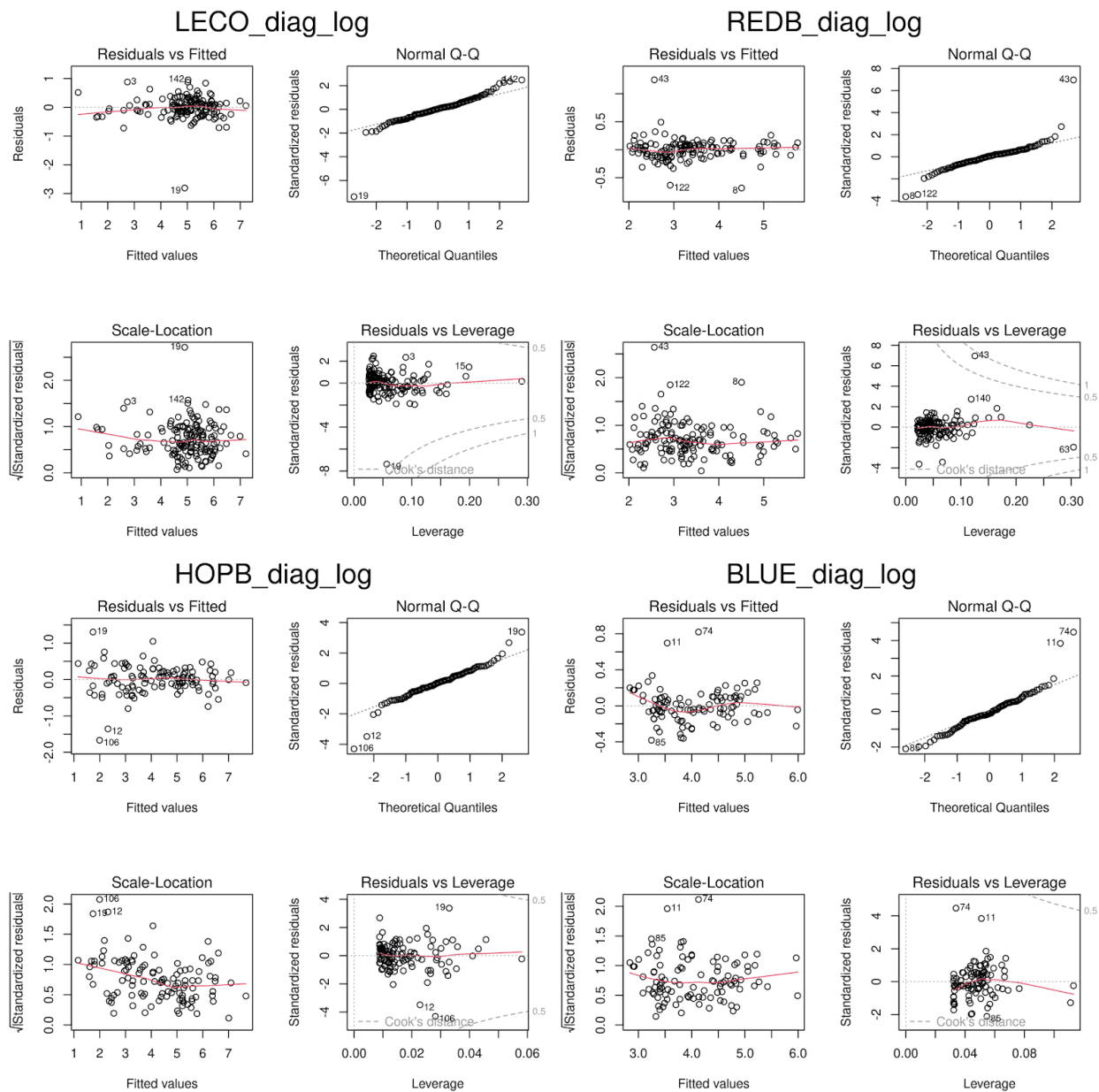


985 Figure A5: Marginal relationships between donor and target gauges for regression on absolute discharge.
 986 Regression lines are shown only for single-donor regressions, fitted via OLS. Site SYCA, here exhibiting
 987 a breakpoint, could not be fitted via segmented regression in the context of absolute discharge.
 988



990 Figure A6: Observed (field) discharge vs. ensemble LSTM predictions.

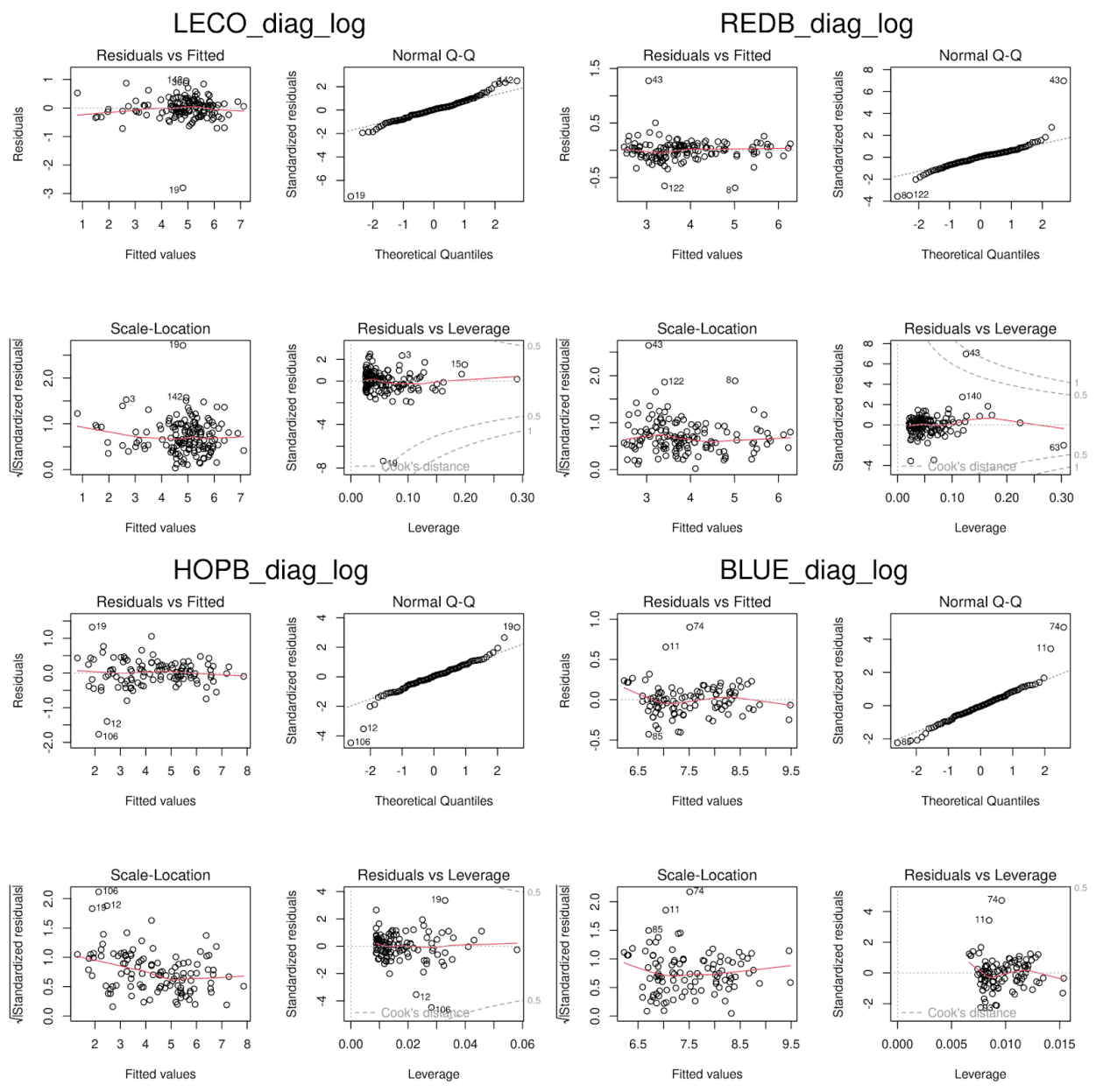
991



992

993 Figure A7: Diagnostic plots for the four sites modeled by OLS regression on specific discharge.

994

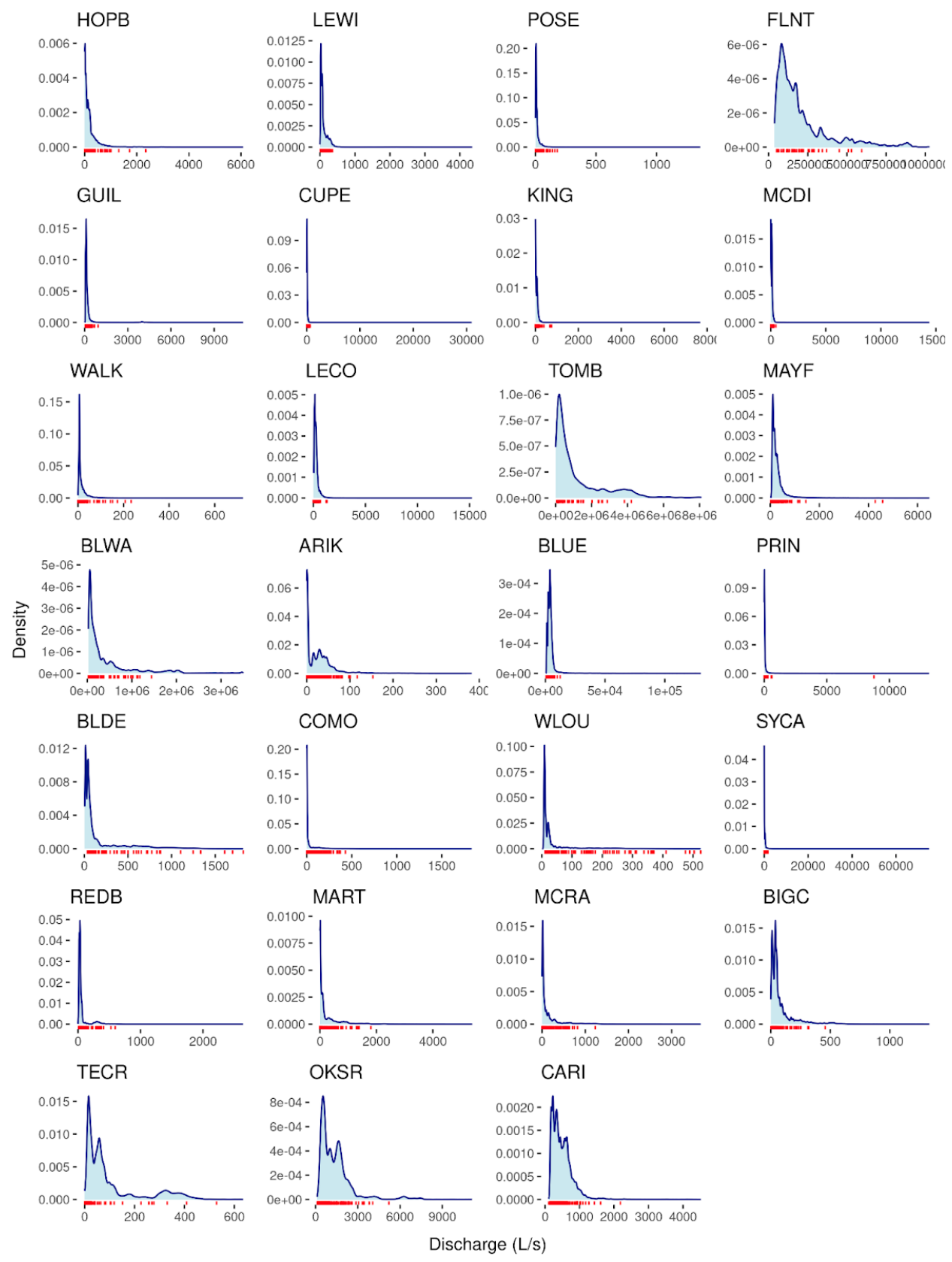


995

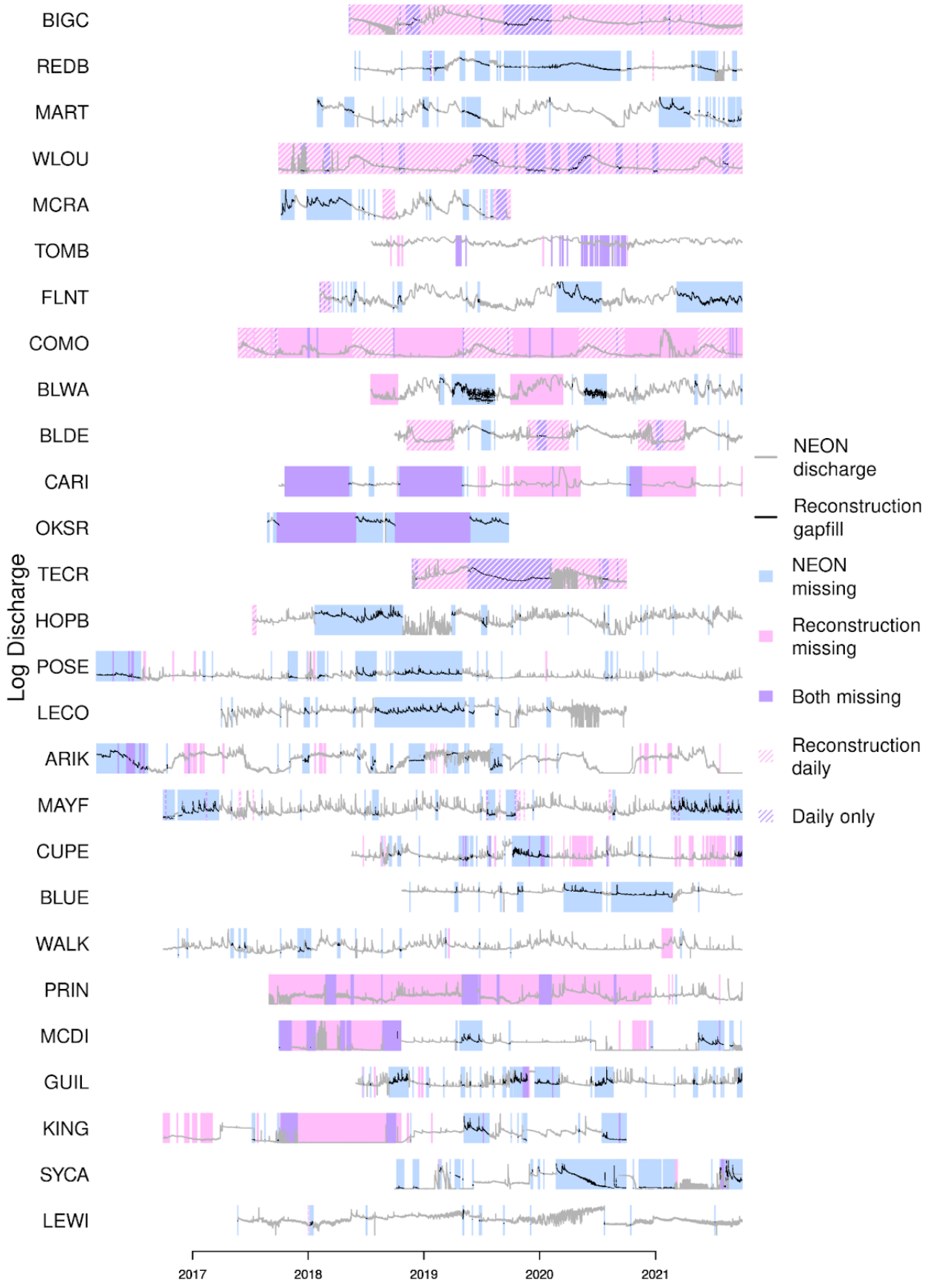
996

997 Figure A8: Diagnostic plots for the four sites modeled by OLS regression on absolute discharge.

998



1000 Figure A9: Density of NEON-estimated discharge (blue polygon) relative to field-measured discharge
1001 observations (red marks).
1002



1004

1005 Figure A10: Durations of missing values (gaps) in NEON's 2023 release of continuous discharge time
1006 series, illustrating gaps filled or informed by estimates from this analysis. All officially published values
1007 are shown, including those with quality control flags. Sites are ordered as in Figure 2. Gaps smaller than
1008 six hours are not indicated.

1009