1  # In-stream *Escherichia Coli* Modeling Using high-

2  # temporal-resolution data with deep learning and

3  # process-based models

4  *Ather Abbas[1], Sangsoo Baek[1], Norbert Silvera[2], Bounsamay Soulileuth[3], Yakov Pachepsky[4]*

5  *Olivier Ribolzi[5], Laurie Boithias[5†], Kyung Hwa Cho[1†]*

6  [1] School of Urban and Environmental Engineering, Ulsan National Institute of Science and

7  Technology, Ulsan 689-798, Republic of Korea

8  [2] Institute of Ecology and Environmental Sciences of Paris (iEES-Paris), Sorbonne Université,

9  Univ Paris Est Creteil, IRD, CNRS, INRA, Paris, France

10  [3] IRD, IEES-Paris UMR 242, c/o National Agriculture and Forestry Research Institute,

11  Vientiane, Lao PDR

12  [4] Environmental Microbial and Food Safety Laboratory, USDA-ARS, Beltsville, MD, USA

13  [5] Géosciences Environnement Toulouse, Université de Toulouse, CNRS, IRD, UPS, Toulouse,

14  France

15  [†] *Co-corresponding authors: Kyung Hwa Cho (khcho@unist.ac.kr), Laurie Boithias*

16  *(laurie.boithias@get.omp.eu)*

17

## Abstract

20

21       Contamination of surface waters through microbiological pollutants is a major concern to

22 public health. Although long-term and high-frequency *E. coli* monitoring can help prevent diseases

23 from fecal pathogenic microorganisms, this monitoring is time consuming and expensive. Process-

24 driven models are an alternative method for determining fecal pathogenic microorganisms.

25 However, process-based modeling still has limitations in improving the model accuracy because

26 of the complex mechanistic relationships among hydrological and environmental variables. On the

27 other hand, with the rise in data availability and computation power, the use of data-driven models

28 is increasing. Therefore, in this study, we simulated the transport of *Escherichia coli* (*E. coli*) in a

29 0.6 km² tropical headwater catchment located in Lao PDR using a deep learning model and a

30 process-based model. The deep learning model was built using the long short-term memory

31 (LSTM) technique, whereas the process-based model was constructed using the Hydrological

32 Simulation Program–FORTRAN (HSPF). First, we calibrated both models for surface as well as

33 for subsurface flow. Then, we simulated the *E. coli* transport with 6 min time steps with both the

34 HSPF and LSTM models. The LSTM provided accurate results for surface and subsurface flow,

35 by showing 0.51 and 0.64 of Nash–Sutcliffe Efficiency (NSE), respectively, whereas the NSE

36 values yielded by the HSPF were -0.7 and 0.59 for surface and subsurface flow. The simulated *E.

37 coli* concentration from LSTM also improved, yielding an NSE of 0.35, whereas the HSPF showed

38 an unacceptable performance, with an NSE value of -3.01. This is because of the limitation of

39 HSPF in capturing the dynamics of *E. coli* with land-use change. The simulated *E. coli*

40 concentration showed rise and drop patterns corresponding to annual changes in land use. This

41 study shows the application of deep learning-based models as an efficient alternative to process-

42 based models for *E. coil* fate and transport simulation at the catchment scale.

43    **Keywords**: hydrological modeling; neural networks; fecal contamination; tropical rivers; South-

44    East Asia; hydrograph separation

45

46

47    **1 Introduction**

48         Contamination of surface waters through microbiological pollutants is a major public

49    health concern (Bain et al., 2014). Worldwide, pathogens have a propensity to wreak havoc on

50    human health because of the diseases they cause, such as diarrhea, resulting in infant mortality. In

51    particular, developing countries are vulnerable to pathogen-related diseases due to the deficit of

52    sanitation facilities (Boithias et al., 2016). *Escherichia coli* (*E. coli*) has been frequently used as

53    an indicator of fecal bacteria because it is easy to culture and less dangerous than other pathogens

54    (Rochelle-Newall et al., 2015). Higher concentrations of *E. coli* in water tend to be linked to fecal

55    pathogenic microorganisms, which are harmful to human health. Although long-term and high-

56    frequency *E. coli* monitoring can help prevent waterborne diseases from fecal pathogenic

57    microorganisms, the monitoring of *E. coli* concentration is time consuming and expensive (Cho et

58    al., 2016; Frolich et al., 2017; Kim et al., 2017). High-frequency datasets of *E. coli* concentration

59    are scarce, and available long-term datasets are often inadequate to yield a continuous

60    concentration of fecal pathogenic microorganisms (van der Leeuw, 2004). This drawback in

61    monitoring can be overcome by modeling approaches. Thus, they can be an alternative to

62    determinin the fate and transport of fecal pathogenic microorganisms at the catchment scale by

63    simulating *E. coli* in each one of the environmental compartments, for example the soil surface

64    and streams (Ligaray et al., 2016; Perez-Pedini et al., 2005; Pacehpsky et al., 2011).

65       Several process-based models have been developed to model stream water contamination

66    by *E. coli*. Popular models to simulate *E. coli* are the Soil and Water Assessment Tool (SWAT)

67    (Neitsch et al., 2011), Hydrological Simulation Program–FORTRAN (HSPF) (Bicknell et al.,

68    1997), INCA-pathogen (Whitehead et al., 2016), and Pathogen Catchment Budget (PCB)

69    (Ferguson et al., 2007). The fate and transport of *E. coli* is a complex phenomenon that depends

70    on several drivers (Pachepsky et al., 2018), such as the hydrological regime (Boithias et al., 2016;

71    Pachepsky et al., 2017), relative contributions of both surface runoff and subsurface flow to the

72    overall in-stream discharge (Boithias et al., 2021), concentration and sources of suspended

73    sediment (Ribolzi et al., 2016; Nguyen et al., 2016), land use (Causse et al., 2015; Nakhle et al.,

74    2021), intrinsic properties of the bacterium (Pachepsky et al., 2014), and economic conditions

75    (Iqbal et al., 2019). However, the process-based model still has limitations in terms of high

76    accuracy due to complex mechanistic relationships among hydrological and environmental

77    variables (Abimbola et al., 2020). In addition, the simplified equations of these models might

78    increase the inherent uncertainties, resulting in simulation errors. The *E. coli* concentration in

79    surface water varies significantly within a very short span of time (Chen et al., 2014; Boithias et

80    al., 2021). Daily and weekly simulations cannot capture the dynamics of *E. coli* in a short duration.

81    In particular, the simulation with high-resolution frequency is important in small headwater

82    catchments because the duration of flood events might be less than one day (Gassman et al., 2007).

83    Therefore, an *E. coli* concentration simulation with high-frequency resolution should be conducted

84    to determine the temporal distribution of *E. coli*.

85       Recently, deep learning (DL) has become a promising alternative approach for estimating

86    water quality by using features of water constituent dynamics (Pyo et al., 2021). Long short-term

87    memory (LSTM) networks have an advantage over other deep learning-based models in that they

88   can extract complex patterns from sequence data (Schmidthuber and Hochreiter, 1997). Several

89   studies have applied deep learning to water quality modeling and prediction (Peterson et al., 2020;

90   Isikdogan et al., 2017; Solanki et al., 2015). Dong et al. (2019) used LSTM to predict dissolved

91   oxygen and showed that LSTM performs better than machine learning methods, such as

92   autoregressive integrated moving average or artificial neural networks. Although LSTM has been

93   used extensively for building hydrological models (Abbas et al., 2020), its potential has not yet

94   been explored to estimate *E. coli* concentration in stream waters. Deep learning-based models have

95   also not been developed for the simulation of water quality with high-resolution frequency.

96        This study aims to evaluate the applicability of LSTM to simulate in-stream *E. coli*

97   concentration with high temporal resolution. In addition, the process-based model HSPF was used

98   as a benchmark to compare and assess the performance of LSTM. Both models were applied in a

99   0.6 km² tropical headwater catchment from the northern Lao People's Democratic Republic (PDR).

100  The temporal resolution of the simulations was 6 min in both models. Thus, the specific objectives

101  of this study were to compare the performance of a process-based model and a deep learning model

102  1) to simulate both surface and subsurface flow, 2) to simulate *E. coli* concentration, and 3) to

103  analyze the response of *E. coli* by changing land use.

## 2 Materials and Methods

### 2.1 Study site and data acquisition

The study area is the Houay Pano headwater catchment, located 10 km south of the city of Luang Prabang, Lao PDR (Boithias et al., 2021) (Fig. 1). This catchment is representative of a montane agroecosystem in Southeast Asia and is part of the long-term critical zone observatories' network called multiscale TROPIcal CatchmentS (M-TROPICS), which is affiliated with the French research infrastructure OZCAR (Gaillardet et al., 2018). This site had undergone rapid land-use changes from 2011 to 2018 (Fig. S1a). The characteristics of this area, including land use information, are provided in the supplementary information (Text S1). We collected climate, hydrological, *E. coli* concentration, and electrical conductivity data at 6 min time steps from 2011 to 2018. Rainfall, relative humidity, solar radiation, wind speed, and air temperature were measured with an automatic weather station Campbell Scientific BWS200, which was equipped with ARG100 (a 0.2 mm capacity tipping bucket). The potential evapotranspiration was calculated using the Penman–Monteith method. We measured the stream water level at the monitoring station using a V-notch and water-level recorder (OTT Thalimedes). The discharge was estimated based on the rating curve between the discharge and water levels. The surface and subsurface flow were calculated using the electrical conductivity method (Ribolzi et al., 2018). A detailed description of this method is provided in the supplementary information (Text S2). *E. coli* concentration was measured based on the standardized microplate method (ISO 9308–3). A detailed explanation of *the E. coli* experiment can be found in the supplementary information (Text S3). In this study, we carried out biweekly grab sampling of *E. coli* from 2011 to 2018. Over the same period, we also specifically sampled 11 flood events to assess *E. coli* dynamics during flood events by using an

Hydrology and
Earth System
Sciences
Discussions
Open Access
EGU

126    automated sampler (ICRISAT) triggered by the water level recorder to collect water after every 2

127    cm water level change during flood rising and every 5 cm water level change during flood

128    recession. The total number of *E. coli* samples collected over the 2011–2018 period was 255. In

129    addition, we collected the monthly number of poultries, swine, goats, and the number of humans

130    who visited the study area. These data were used to quantify the source of *E. coli* in this catchment

131    (Rochelle-Newall et al., 2016) (Fig. S1b).

132

133    **2.2 Flow and *E. coli* concentration simulation.**

134    In this study, HSPF and LSTM models were used to simulate in-stream surface flow,

135    subsurface flow, and *E. coli* concentration. HSPF and LSTM are popular models among the

136    process-based and DL models (Bicknell et al., 1997; Ahmadisharaf and Benham, 2020; Kratzert

137    et al., 2019). Both models have been adopted for hydrological and water quality simulations

138    (Peterson et al., 2020; Isikdogan et al., 2017; Ahmed et al., 2014). In the HSPF, the simulation of

139    surface and subsurface flow and of *E. coli* concentration was carried out in three steps: (1) building

140    the model, (2) conducting sensitivity analysis based on the Latin-Hypercube–One-factor-At-a-

141    Time (LH-OAT), and 3) calibrating the model using the Newton algorithm (Nash, 1984). A

142    schematic of the LSTM simulation is shown in Fig. 2. The first step in building this model was

143    data preparation (Fig. 2a). LSTM then simulated surface and subsurface flow with climate data

144    (Fig. 2b). Finally, we estimated the *E. coli* concentration at 6 min intervals using rainfall, bacteria

145    source, land-use change, and surface and subsurface flow (Fig. 2c). Both models considered the

146    source of *E. coli* to simulate its concentration at the catchment outlet. The fecal matter from the *E.*

147    *coli* sources was assumed to be evenly distributed in the catchment. The monthly *E. coli* source

148    data is presented in Fig. S1b. The time series data of the *E. coli* source was used as input for the *E.*

149    *coil* simulation.

150

**2.2.1** Hydrological Simulation Program Fortran (HSPF)

151

152      The HPSF model is a process-driven model that simulates processes at the catchment scale

153 (Bicknell et al., 1997). It has been extensively used to model the fate and transport of *E. coli* in

154 catchments (Ahmadisharaf and Benham, 2020; Chin et al., 2009) and to develop total maximum

155 daily loads of *E. coli* at various locations (Mishra et al. 2018; Yagow et al., 1998). The original

156 software was written in the FORTRAN programming language. Recently, the Hydrological

157 Simulation Program Python (HSP2) was developed based on the Python programming language

158 (van Rossum, 2007). HSP2 is a platform-independent software that extends the functionality of

159 HSPF by allowing the use of dynamic variables and easier management of input and output files

160 (Heaphy et al., 2015). The HSPF simulates the hydrological cycle by discretizing the catchment

161 into pervious and impervious hydrological response units (HRUs). Previous HRUs simulate

162 evapotranspiration, surface detention, surface infiltration, interflow, baseflow, and deep

163 percolation, whereas impervious HRUs simulate surface detention and surface flow (Bicknell et

164 al., 1997). The simulation of in-stream *E. coli* concentration in HSPF is based on a first-order

165 kinetics approach, considering the decay rate (Fonseca et al., 2014). Detailed descriptions of

166 hydrological and *E. coli* simulations can be found in Bicknell et al. (1997). For this study, we

167 rewrote the modules of *E. coli* simulation, and the simulation was carried out in the Python

168 programming language. This allowed us to incorporate more dynamic use of input data, such as

169 the annual change in land use and the monthly bacterial source.

170      In our study, HRUs were divided into four units based on land use: Forest, Fallow, Teak,

171 and Annual crop. Among land uses, we did not consider any imperviousness in the Forest and

172 Fallow. We considered 2 % and 1 % imperviousness for the Teak and Annual crop land uses (Patin

173   et al., 2018). We selected 13 and 4 parameters for each land use for the sensitivity analysis of

174   hydrological and *E. coil* simulations, respectively (Table 1 and Table S1). The total number of

175   parameters for hydrological and *E. coil* simulation were 52 and 18, respectively. In model

176   calibration, we selected the 25 most sensitive parameters of the hydrological simulation and all

177   parameters of the *E. coli* simulation. Sensitivity analysis and model calibration were conducted

178   based on the LH-OAT and the Newton algorithm, respectively. A detailed explanation of the LH-

179   OAT and the Newton algorithm can be found in the Supplementary Information (Text S4).

180

Hydrology and
Earth System
Sciences
Discussions

### 2.2.2 Long short-term memory (LSTM)

181

182    In the data preparation step (Fig. 2a), our data were converted to a 6 min frequency. We

183    then built the LSTM model to simulate surface and subsurface flow using the validated model

184    structure (Abbas et al., 2020) (Fig. 2b). It uses historical data of rainfall, solar radiation, air

185    temperature, and potential evapotranspiration to simulate surface and subsurface flow. To simulate

186    the output at a time-step "t," LSTM uses the data of previous "n" time steps as inputs (Chollet,

187    2018). The inputs from previous time steps are used by LSTM to predict the output at the next

188    time step (t+1). The number of these time steps "n" are called lookback steps (Chollet, 2018). The

189    simulated surface and subsurface flow from the LSTM were applied to simulate the *E. coli*

190    concentration (Fig. 2c). We adopted a bacterial source and land-use information as input for the

191    LSTM. To investigate the impact of land-use change on in-stream *E. coli* concentration, we

192    conducted *E. coli* simulations in two scenarios. In scenario 1, we used the land-use change and *E.*

193    *coli* source information separately. In scenario 2, we calculated the *E. coli* source per area for each

194    land use.

195    LSTM is a special type of recurrent neural network designed to extract temporal features

196    from sequence data (Hochreiter and Schmidhuber, 1997). An LSTM cell is the basic building block

197    of the LSTM (Fig. S2). It consists of three "gates" and two "states" The gates are "forget," "update,"

198    and "output," which decide what information to forget, allow in, and allow out from the LSTM

199    "memory," respectively. The states act as a memory or information carrier across time. The

200    equations describing the functions of gates and states are as follows:

$$C_c^{<t>} = tanh(W_c[\,h^{<t-1>}, x^{<t>}] + b_c), \tag{1}$$

$$\Gamma_f = \sigma(W_f[c^{<t-1>}, x^{<t>}] + b_{f)}, \tag{2}$$

$$\Gamma_o = \sigma(W_o[c^{<t-1>}, x^{<t>}] + b_{o)}, \qquad (3)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_{u)}, \qquad (4)$$

$$C^{<t>} = \Gamma_u * C_c^{<t>} + \Gamma_f * C^{<t-1>} \qquad (5)$$

$$h^{<t>} = \Gamma_o * tanh\, C^{<t>}. \qquad (6)$$

201

202   The symbol $*$ in the above equations represents elementwise multiplication. The behavior

203 of each gate is controlled by the weights (W) and biases (b) associated with them. Their output

204 was further modified by a nonlinear function ($\sigma$). At each time step (t), the prospective cell state

205 ($C_c^{<t>}$) is calculated based on the output from the previous time step ($h^{<t-1>}$) and the input from

206 the current time step ($x^{<t>}$) (Eq. 1). The notation $W_c[h^{<t-1>}, x^{<t>}]$ represents pointwise

207 multiplication of new inputs and previous hidden state with the weight matrix $W_c$ separately and

208 then adding their output. This prospective cell state ($C_c^{<t>}$), along with the output from the "forget"

209 and "update" gate decides the current cell state ($c^{<t>}$) (Eq. 5). The current cell state and output

210 gate control the output values from LSTM ($h^{<t>}$), the so-called hidden state (Eq. 6). The

211 hyperbolic tangent ($tanh$) is another nonlinearity used in LSTM for the calculation of the cell state

212 (Eq. 1) and the output state (Eq. 6). Equations 1–6 are used to calculate the LSTM output, which

213 is then compared with observed values to calculate the error. This study used the mean square error

214 (MSE) as the error function.

215   We used the TensorFlow software v1.15 for building the LSTM model (Abadi et al., 2016).

216 We used an Intel® Core™ i7-9700 processor with a graphics card of NVIDIA GeForce RTX 2080

217 having 12 GB of dedicated GPU memory, along with 64 GB of Random-Access Memory for

218 simulating surface, subsurface, and *E. coli*.

219

### 2.2.3 Hyperparameters of LSTM

221        The structure and performance of the LSTM were controlled by hyperparameters,

222   including the dropout rate, LSTM units, learning rate, lookback steps, and activation functions for

223   both LSTM and the fully connected layer (Table 2). Dropout is a regularization technique that

224   switches off a certain number of nodes in the LSTM (Goodfellow et al., 2016). This simple

225   technique helps break the brittle coadaptation of weights, which hinders generalization to unseen

226   data. This way, dropout prevents overfitting (Srivastava et al., 2014). In overfitting, the model

227   performs better on calibration data, but its performance deteriorates on new unseen data. The

228   number of LSTM units directly corresponds to the learning capacity of LSTM, but it also accounts

229   for more memory and computation. This number determines the size of the weight matrix of an

230   LSTM. The learning rate defines the change in the weights of the neural network during calibration

231   (Goodfellow et al., 2016). A higher number of lookback steps allows LSTM to capture long-term

232   patterns at the cost of an increase in memory consumption and computation. The activation

233   function determines the nonlinearity in the model.

234

### 2.3 Performance statistics

236        Evaluations to assess the performance of the HSPF and LSTM were conducted

237   using MSE, Nash–Sutcliffe efficiency (NSE), and percent bias (PBIAS) (Nash and Sutcliffe,

238   1970; Gupta et al., 1999). NSE is useful for interpreting the model performance by generating a

239   dimensionless value as the performance index (Lin et al., 2017). The PBIAS measures the

240   average tendency of the simulated data to be overestimated or underestimated than observed

OK, producing final.

241 values (Moriasi 2007). The MSE, NSE, and PBIAS were calculated using the following

242 equations:

$$MSE = \frac{[\sum_{i=1}^{n}(o_i - p_i)^2]}{n} \tag{7}$$

$$NSE = 1 - \frac{\sum(o_i - p_i)^2}{\sum(o_i - \bar{o})^2} \tag{8}$$

$$PBIAS = 1 - \frac{\sum_{i=1}^{n} o_i - p_i}{\sum_{i=1}^{n} o_i} \tag{9}$$

246 where $p_i$ is the simulated data, $o_i$ is the observed data, and $n$ is the number of points in the data.

## 3 Results and discussion

### 3.1 Land use change and *E. coli* source

The land-use change from 2011 to 2018 is shown in Fig. S1a. The area of Fallow land-use increased from 2011 to 2016, whereas Annual crop area decreased. Teak tree plantations were expanded until 2013 and were retained. Forest land use accounted for about 10 % of the study area from 2011 to 2018. In general, the land-use change has been dynamic from 2011 to 2013, whereas its variation diminished from 2016 to 2018. Previous studies have demonstrated that the expansion of Teak trees might increase the surface flow (Ribolzi et al., 2017; Song et al., 2020). Higher runoff at the soil surface may cause a higher inflow of *E. coli* with surface flow. The monthly *E. coli* source in the catchment decreased from $2 \times 10^{15}$ in 2011 to $3 \times 10^{14}$ in 2018 (Fig. S1b). This decrease in *E. coli* source is caused by the decrease in manpower needed in Teak tree plantations and in Fallow plots, compared to the Annual crop (Fig. S1a) (Boithias et al., 2021).

## 3.2 Sensitivity analysis and optimization result

261

262      The sensitivity results for the flow simulation are shown in Fig. S3, and the most sensitive

263 parameters are listed in Table S2. The interflow and infiltration-related parameters were the most

264 sensitive parameters for surface and subsurface flows. The Manning's "n" value (NSUR) for Teak

265 and Fallow land uses was among the 10 most sensitive parameters. Kim et al. (2017) suggested

266 that Manning's value is the most sensitive parameter in the hydrological simulation of tropical

267 headwater catchments, such as the Houay Pano catchment in northern Lao PDR. The groundwater

268 recession rate (AGWRC) and soil infiltration capacity (INFILD) were sensitive to subsurface flow.

269 In Annual crop land use, infiltration capacity (INFILT) and upper zone storage (UZSN) were the

270 most sensitive parameters. Abbas et al. (2020) demonstrated that INFILT is the most sensitive

271 parameter for subsurface flow in tropical subcatchments.

272      The sensitivity analysis results for *E. coli* are shown in Fig. S4 and Table S3. The

273 parameters related to the transport of *E. coli* on the land surface (e.g., WSQOP, SQOLIM_MF)

274 were more sensitive than other parameters. IOQC and AOQC were the least sensitive parameters.

275 These parameters are related to *E. coli* transport in interflow and baseflow (Bicknell et al., 2011).

276 This implies that the in-stream *E. coli* concentration at the study site is mainly driven by surface

277 flow (Boithias et al., 2021). A previous study also demonstrated that 89 % of in-stream *E. coli*

278 concentrations were driven by surface flow (Boithias et al., 2021). Figure 3 shows the model

279 performance dependent on different objective functions. We found that the model performance

280 was better when the NSE was selected as the objective function. The NSE of the surface and

281 subsurface flow was positive by optimizing with NSE. However, the NSE value for surface flow

282     was negative when the objective function was MSE during the optimization. Negative NSE

283     indicated an "unsatisfactory" performance range (Moriasi et al., 2015).

284

285     **3.3 Flow simulation**

286          The simulated surface and subsurface flow using the HSPF are plotted in Fig. 4. We found

287     that the simulated subsurface flow was underestimated compared to the observations. Although

288     surface flow from the HSPF followed the trend and peaks of observations, this model yielded a

289     negative NSE value, indicating that the model simulation was unacceptable (Moriasi et al., 2015)

290     (Table 3). The NSE values for subsurface flow from HSPF were 0.49 and 0.59 for calibration and

291     validation, respectively. Hence, the HSPF model is better at simulating subsurface flow than

292     surface flow. In particular, the simulated surface flow was underestimated compared to the

293     observations. The average values of INFILT and UZSN were 0.36 and 1.22, respectively, which

294     were larger than those reported in previous studies (Lee et al., 2020). INIFILT controls the overall

295     division of available moisture into the surface and subsurface (Bicknell et al., 2001). The parameter

296     UZSN influences the evapotranspiration process (Bicknell et al., 2001). This underestimation of

297     surface flow using HSPF is consistent with a previous study (Kim et al., 2017). We also

298     investigated the impact of underestimation and overestimation of the flow by plotting flow

299     duration curves (Fig. S5). Although both flows can capture the peak flow, the simulated subsurface

300     flow was still underestimated compared to the observed subsurface flow.

301          The simulated surface and subsurface flows using the LSTM model are plotted in Fig. 5.

302     The NSE values for the calibration period were 0.56 and 0.69 for surface and subsurface flow,

303     respectively. The corresponding validation NSE of the surface and subsurface flow were 0.51 and

304    0.64, respectively. These results indicate that the LSTM had a satisfactory performance for both

305    the calibration and validation periods according to the criteria of Moriasi et al. (2015). LSTM

306    overcame the problem of the HSPF model underestimating subsurface flow. In addition, the peak

307    surface flows from the LSTM were similar to observations. The observed and simulated flows in

308    storm events are presented in Figs. S6–S11. LSTM can follow the observed trends in surface and

309    subsurface flow more closely than the HSPF. This leads to increased NSE values for both surface

310    flow as well as for subsurface flow. The hyperparameters of the LSTM are described in Table 2.

311    The rectified linear unit (ReLU) was chosen as the activation function for the LSTM output.

312    Because the simulated *E. coli* should be positive, we chose ReLU, which cannot produce negative

313    values from the model (Nair and Hinton, 2010). The optimal batch size and LSTM units were 16

314    and 100, respectively. The optimal value of the lookback steps was 50, which is equal to 5 h of

315    input data.

316        We analyzed the model performance for surface and subsurface flows during storm events

317    (Fig. 6). These events were selected where the peak flow exceeded 0.2 m per s. The performance

318    of LSTM is considerably better than that of HSPF for most storm events. In surface flow, the

319    average MSE of LSTM and HSPF was 1.1e-4 and 6.1e-4 ($m^3s^{-1}$), respectively. The NSE values

320    from LSTM varied from 0.2 to 0.6, whereas that of HSPF ranged from -1.0 to 0.4. We found that

321    the NSE values from the HSPF vary considerably depending on storm events. On June 11, 2015,

322    the NSE value of HSPF was as high as 0.4, whereas for some others it was below 0. Although the

323    subsurface flow of the HSPF provided better model performance than surface flow simulation, this

324    model still presented an unacceptable result with a negative NSE value.

325

Hydrology and
Earth System
Sciences
Discussions

### 3.4 *E. coli* simulation

326

327     Figure 7 shows the temporal distribution of *E. coli* concentration using HSPF and LSTM.

328 The *E. coli* concentration from HSPF was overestimated compared to the observed *E. coli*

329 concentration. The performance matrices of the HSPF were also worse than those of the LSTM

330 (Table 4). In particular, the HSPF simulation presented a PBIAS value of 73, indicating an

331 overestimation of *E. coli* concentration (Moriasi et al., 2015). Ackerman and Weisman (2014)

332 reported that the *E. coli* simulation from HPSF was overestimated compared to observation. The

333 overestimation of simulated *E. coli* at tropical sites has also been observed by Kim et al. (2017).

334 *E. coli* simulation from LSTM is satisfactory in both calibration and validation periods according

335 to the criteria set by Moriasi et al. (2015). In contrast, the HSPF result can be regarded as

336 "unsatisfactory" in both the calibration and validation periods. These results implied that LSTM

337 could generate acceptable performances and had good agreement between the observed and

338 simulated *E. coli*.

339     The simulation during the storm events using both the HSPF and LSTM models are

340 shown in Fig. 8 and Figs. S6–S11. Figure 8 shows the storm events from the validation data,

341 whereas the other figures show the storm events from the calibration data. In general, the simulated

342 *E. coli* by HPSF and LSTM were overestimated and underestimated, respectively. This difference

343 might be caused by the fact that *E. coli* from HSPF is more responsive to surface flow, wheras *E.*

344 *coli* from LSTM is more influenced by subsurface flow (Ackerman and Weisman, 2014). The

345 sensitivity analysis of HSPF also demonstrated that the influence of interflow and baseflow on *E.*

346 *coli* is weaker than surface flow because the parameters IOQC and AOQC are the least sensitive

347 parameters for *E. coli* simulation. Both parameters affect the *E. coli* concentration in interflow and

Hydrology and
Earth System
Sciences
Discussions

Open Access
EGU

348    baseflow (Bicknell et al., 2001). The simulated *E. coli* of LSTM rose sharply and dropped slowly,

349    similar to the observations, whereas that of the HSPF decreased steeply (Figures S6–S11).

350    Although both models simulated the peak time of the *E. coil* correctly, the HSPF was limited to

351    simulate a slope in its falling limb. This performance difference between both models was caused

352    by the extent of influence from hydrological variables (e.g., rainfall, surface flow, and subsurface

353    flow) to model output. LSTM was effective in reflecting the response of variables to output

354    (Kratzert et al., 2019).

355         The performance matrices for the LSTM and HSPF models during storm events are shown

356    in Fig. 9. In general, we observed better LSTM performance than HSPF for both NSE and MSE

357    values. The HSPF model performed better than the LSTM for only two storm events on June 15,

358    2014, and June 11, 2015. For the remaining storm events, the NSE values from LSTM are higher

359    than those of the HSPF—an NSE range from 0.20 to 0.65. Similarly, for MSE values, the LSTM

360    was superior to the HSPF for all storm events except for the storm events of June 15, 2014, and

361    June 11, 2015.

362         We observed the impact of logarithmic and minmax transformations on model performance

363    (Fig. 10). The result of the logarithmic transformation was closer to the observation than the

364    minmax transformation by showing an NSE of 0.57. A negative PBIAS value was obtained in

365    logarithmic transformation. This indicated that the simulated *E. coli* from logarithmic

366    transformation was underestimated, whereas the result of the minmax transformation was

367    overestimated. The reason for this behavior can be attributed to the ability of minmax scaler to be

368    more sensitive to outliers (Chuang et al., 2010). As a result, if a better accuracy during storm events

369    is required, the target variable can be transformed on a logarithmic scale prior to calibration. This

370    is because log transformation can reduce the effect of outliers from data (Singh and Kingsbury,

Hydrology and
Earth System
Sciences
Discussions
Open Access

EGU

371  2017). It has been reported that log transformation can improve the performance of data-driven

372  models when the data contain outliers (Zheng and Casari, 2018).

373

### 374  3.5 *E. coli* response to land-use change

375       We investigated the impact of land-use change and bacterial sources on the in-stream *E.*

376  *coli* concentration simulation (Fig. 11). In scenario 1, we used land-use change time-series

377  information (Fig. S2a) and bacterial source information (Fig. S2b). In scenario 2, we divided the

378  bacterial source by the fraction of each land use (Fig. S2c). In scenario 1, we observed a larger

379  variation in *E. coli* concentration from 2014 to 2018 (Fig. 11a), whereas in scenario 2, the variation

380  in *E. coli* was smaller than that in scenario 1 (Fig. 11b). This variation in *E. coli* was due to land-

381  use change in scenario 1. In particular, *E. coli* in 2016 was less than in other years because Annual

382  crop land use decreased. On the other hand, the variation in *E. coli* was not observed in scenario 2

383  from 2015 to 2017. Neither scenario showed a significant response from 2011 to 2014. During

384  these years, the rise in Fallow land use was complemented by a decrease in Annual crop land use.

385

Hydrology and
Earth System
Sciences
Discussions

### 3.6 Limitations and future research

386

387      Transport of soil particles by surface flow and suspended sediments within the stream play

388 a crucial role in the fate and transport of *E. coli* (Thupaki et al., 2013). Several studies have

389 emphasized the importance of particle size (Cho et al., 2010), adsorption to soil and sediment

390 particles (Palmateer et al., 1993), and resuspension of *E. coli* (Kim et al., 2017) with streambed

391 sediments for modeling the fate and transport of *E. coli* at the catchment scale. In this study, we

392 did not consider sediment transport, nor the attachment/detachment of *E. coli* on/from soil particles

393 and suspended sediments. Several studies have been conducted on the monitoring and modeling

394 of *E. coli* without considering sediment transport (Ahmadisharaf and Benham, 2020; Mishra et al.,

395 2018). However, the need for its inclusion has been indicated elsewhere (Pandey and Soupir, 2013).

396 To model sediment transport, additional data on suspended sediment concentration are required to

397 build both the HSPF and deep learning-based models. Therefore, this modeling exercise can be

398 further improved by collecting sediment-related data and modeling sediment transport along with

399 *E. coli* concentration.

400

## 4 Conclusions

401

402    In this study, we simulated the transport of bacteria in a headwater catchment of the northern

403    Lao PDR at 6 min time steps. The main findings of this study are summarized as follows:

404    • Both the LSTM and HSPF models can accommodate land-use change and bacteria-

405      source variation with time.

406    • The performance of the surface and subsurface flow simulation of LSTM was superior for

407      both the calibration and validation steps when compared with the HSPF. The LSTM

408      provided accurate results for surface and subsurface flow by showing NSE values of 0.51

409      and 0.59, respectively, whereas the HSPF showed -0.7 and 0.55 of NSE.

410    • Our LSTM model showed better performance compared to HSPF for *E. coli* simulation.

411      The NSE of the HSPF and LSTM were -3.01 and 0.35, respectively. We found that the

412      LSTM model can respond to changes in land use.

413    This study shows that deep learning-based models are an efficient alternative to process-based

414    models to simulate *E. coli* in a given catchment. Because LSTM can generate reasonable *E. coli*

415    simulations, it could be applied to provide effective strategies for diseases that wreak havoc on

416    human health. Therefore, a deep learning approach can be useful in developing better water

417    sustainability and management.

**Code Availability**

419     Programming Language: Python

420     Software development: PyCharm

421     Year first available: 2021

422     Software Availability: contact the authors

423     Contact Address: School of Urban and Environmental Engineering, Ulsan National Institute of

424     Science and Technology, UNIST-gil 50, Ulsan, 689–798, Republic of Korea. E-mail:

425     khcho@unist.ac.kr

**Author contributions**

427     **Ather Abbas:** Conceptualization, Data curation, Methodology, Visualization, Writing - original

428     draft, Writing - review & editing. **Sangsoo Baek:** Visualization, Writing - review & editing.

429     **Olivier Ribolzi:** review & editing. **Norbert Silvera:** Data curation. **Bounsamay Soulileuth**:

430     Sampling and data preparation. **Yakov Pachepsky:** review & editing. **Laurie Boithias:** Funding

431     acquisition, Supervision, Validation, Writing - review & editing. **Kyung Hwa Cho:**

432     Conceptualization, Funding acquisition, Supervision, Validation, Writing - review & editing.

**Competing Interests**

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

435

436

Hydrology and
Earth System
Sciences

Open Access

EGU

Discussions

## Acknowledgments

# References

Abbas, A., Baek, S., Kim, M., Ligaray, M., Ribolzi, O., Silvera, N., ... and Cho, K. H.: Surface

and sub-surface flow estimation at high temporal resolution using deep neural networks. *Journal*

*of Hydrology*, 125370. https://doi.org/10.1016/j.jhydrol.2020.125370, 2020.

Abimbola, O. P., Mittelstet, A. R., Messer, T. L., Berry, E. D., Bartelt-Hunt, S. L., and Hansen,

S. P.: Predicting Escherichia coli loads in cascading dams with machine learning: An integration

of hydrometeorology, animal density and grazing pattern. *Science of The Total Environment*,

137894. https://doi.org/10.1016/j.scitotenv.2020.137894, 2020.


Abimbola, O., Mittelstet, A., Messer, T., Berry, E., and van Griensven, A.: Modeling and

Prioritizing Interventions Using Pollution Hotspots for Reducing Nutrients, Atrazine and E. coli

Concentrations in a Watershed. *Sustainability*, *13*(1), 103. https://doi.org/10.3390/su13010103,

2021.


Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., ... and Kudlur, M.: Tensorflow: A

system for large-scale machine learning. In *12th {USENIX} symposium on operating systems*

*design and implementation ({OSDI} 16)* (pp. 265-283), 2016.


Ackerman, D., and Weisberg, S. B.: Evaluating HSPF runoff and water quality predictions at

multiple time and spatial scales. *SBW a. K. Miller (Ed.), Southern California coastal water*

*research project biennial report*, *2006*, 293-303, 2005.

465

466 Adomat, Y., Orzechowski, G. H., Pelger, M., Haas, R., Bartak, R., Nagy-Kovács, Z. Á., ... and

467 Grischek, T.: New Methods for Microbiological Monitoring at Riverbank Filtration Sites. *Water*,

468 *12*(2), 584. https://doi.org/10.3390/w12020584, 2020.

469

470 Ahmadisharaf, E., and Benham, B. L.: Risk-based decision making to evaluate pollutant

471 reduction scenarios. *Science of The Total Environment*, *702*, 135022.

472 https://doi.org/10.1016/j.scitotenv.2019.135022, 2020.

473

474 Ahmed, S. I., Singh, A., Rudra, R., and Gharabaghi, B.: Comparison of CANWET and HSPF for

475 water budget and water quality modeling in rural Ontario. *Water Quality Research Journal of*

476 *Canada*, *49*(1), 53-71, 2014.

477

478 Benham, B., Yagow, G., Barham, B., Zeckoski, R., and Dillaha, T.: Total Maximum Daily Load

479 Development: Mill Creek bacteria (E. coli) impairment, Page County, Virginia. Richmond, Va.:

480 Virginia Department of Environmental Quality, 2005.

481

482 Bicknell, B. R., Imhoff, J. C., Kittle Jr, J. L., Donigian Jr, A. S., and Johanson, R. C.:

483 Hydrological simulation program—FORTRAN user's manual for version 11. *Environmental*

484 *Protection Agency Report No. EPA/600/R-97/080. US Environmental Protection Agency, Athens,*

485 *Ga.*, 1997.

486

487 Boithias, L., Choisy, M., Souliyaseng, N., Jourdren, M., Quet, F., Buisson, Y., ... and Pierret, A.:

488 Hydrological regime and water shortage as drivers of the seasonal incidence of diarrheal diseases

489 in a tropical montane environment. *PLoS neglected tropical diseases*, *10*(12), e0005195.

490 https://doi.org/10.1371/journal.pntd.0005195, 2016.

491

492 Causse, J., Billen, G., Garnier, J., Henri-des-Tureaux, T., Olasa, X., Thammahacksa, C., ... and

493 Ribolzi, O.: Field and modelling studies of Escherichia coli loads in tropical streams of montane

494 agro-ecosystems. *Journal of Hydro-Environment Research*, *9*(4), 496-507.

495 https://doi.org/10.1016/j.jher.2015.03.003, 2015.

496

497 Chen, H. J., and Chang, H.: Response of discharge, TSS, and E. coli to rainfall events in urban,

498 suburban, and rural watersheds. *Environmental Science: Processes & Impacts*, *16*(10), 2313-

499 2324, 2014.

500

501 Chen, K., Chen, H., Zhou, C., Huang, Y., Qi, X., Shen, R., ... and Zhang, Y.: Comparative

502 analysis of surface water quality prediction performance and identification of key water

503 parameters using different machine learning models based on big data. *Water Research*, *171*,

504 115454. https://doi.org/10.1016/j.watres.2019.115454, 2020.

505

506    Chin, D. A., Sakura-Lemessy, D., Bosch, D. D., and Gay, P. A.: Watershed-scale fate and

507    transport of bacteria. *Transactions of the ASABE*, *52*(1), 145-154.

508    https://doi.org/10.13031/2013.25955, 2009.

509

510    Cho, K. H., Pachepsky, Y. A., Kim, J. H., Guber, A. K., Shelton, D. R., and Rowland, R.:

511    Release of Escherichia coli from the bottom sediment in a first-order creek: Experiment and

512    reach-specific modeling. *Journal of Hydrology*, *391*(3-4), 322-332.

513    https://doi.org/10.1016/j.jhydrol.2010.07.033, 2010.

514

515    Cho, K. H., Pachepsky, Y. A., Oliver, D. M., Muirhead, R. W., Park, Y., Quilliam, R. S., and

516    Shelton, D. R.:. Modeling fate and transport of fecally-derived microorganisms at the watershed

517    scale: state of the science and future opportunities. *Water research*, *100*, 38-56.

518    https://doi.org/10.1016/j.watres.2016.04.064, 2016.

519

520    Chuang, C. C., Wang, C. M., and Li, C. W.: Weighted linear regression for symbolic interval-

521    values data with outliers. In *2010 5th IEEE Conference on Industrial Electronics and

522    Applications* (pp. 2238-2242). IEEE, 2010.

523

524    Clevert, D. A., Unterthiner, T., and Hochreiter, S.: Fast and accurate deep network learning by

525    exponential linear units (elus). *arXiv preprint arXiv:1511.07289*. 2015.

526

527    Chollet, F.: *Deep Learning mit Python und Keras: Das Praxis-Handbuch vom Entwickler der*

528    *Keras-Bibliothek*. MITP-Verlags GmbH & Co. KG. 2018.

529

530    Dosovitskiy, A., and Djolonga, J.: You Only Train Once: Loss-Conditional Training of Deep

531    Networks. In *International Conference on Learning Representations*.

532    https://openreview.net/pdf?id=HyxY6JHKwr, September 2019.

533

534    Dong, Q., Lin, Y., Bi, J., and Yuan, H.: An Integrated Deep Neural Network Approach for

535    Large-Scale Water Quality Time Series Prediction. In *2019 IEEE International Conference on*

536    *Systems, Man and Cybernetics (SMC)* (pp. 3537-3542). IEEE, October 2019.

537

538    Frolich, L., Vaizel-Ohayon, D., and Fishbain, B.: Prediction of Bacterial Contamination

539    Outbursts in Water Wells through Sparse Coding. *Scientific Reports*, *7*(1), 1-11.

540    https://doi.org/10.1038/s41598-017-00830-4, 2017.

541

542    Fujioka, R. S., Solo-Gabriele, H. M., Byappanahalli, M. N., and Kirs, M. US recreational water

543    quality criteria: a vision for the future. *International journal of environmental research and*

544    *public health*, *12*(7), 7752-7776. 10.3390/ijerph120707752, 2015.

545

546     Gaillardet, J., Braud, I., Hankard, F., Anquetin, S., Bour, O., Dorfliger, N., ... and Zitouna, R.:

547     OZCAR: The French network of critical zone observatories. *Vadose Zone Journal*, *17*(1), 1-24.

548     https://doi.org/10.2136/vzj2018.04.0067, 2018.

549

550     Gassman, P. W., M.R. Reyes, C.H. Green, J.G. Arnold.: The soil and water assessment tool:

551     historical development, applications, and future research directions Trans. ASABE, 50 (4), pp.

552     1211-1250, 10.13031/2013.23637, 2007.

553

554     Goodfellow, I., Bengio, Y., and Courville, A.: *Deep learning*: MIT press, 2016.

555

556     Gupta, H. V., Sorooshian, S., and Yapo, P. O.: Status of automatic calibration for hydrologic

557     models: Comparison with multilevel expert calibration. *Journal of hydrologic engineering*, *4*(2),

558     135-143. https://doi.org/10.1061/(ASCE)1084-0699(1999)4:2(135, 1999.

559

560     Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared

561     error and NSE performance criteria: Implications for improving hydrological modelling. *Journal

562     of hydrology*, *377*(1-2), 80-91. https://doi.org/10.1016/j.jhydrol.2009.08.003, 2009.

563

564     Heaphy, R. T., Burke, M. P., and Love, J. T.: Conversion of HSPF Legacy Model to a Platform-

565     Independent, Open-Source Language. *AGUFM*, *2015*, H13C-1529, 2015.

566

567  Hinton, G. E., Osindero, S., and Teh, Y. W.: A fast learning algorithm for deep belief nets.

568  *Neural computation*, *18*(7), 1527-1554. https://doi.org/10.1162/neco.2006.18.7.1527, 2006.

569

570  Iqbal, M. S., Islam, M. M., and Hofstra, N.: The impact of socio-economic development and

571  climate change on E. coli loads and concentrations in Kabul River, Pakistan. *Science of the Total*

572  *Environment*, *650*, 1935-1943. https://doi.org/10.1016/j.scitotenv.2018.09.347, 2019.

573

574  Kim, M., Boithias, L., Cho, K. H., Silvera, N., Thammahacksa, C., Latsachack, K., ... and

575  Ribolzi, O.: Hydrological modeling of fecal indicator bacteria in a tropical mountain catchment.

576  *Water research*, *119*, 102-113. https://doi.org/10.1016/j.watres.2017.04.038, 2017.

577

578  Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., and Nearing, G.: Towards

579  learning universal, regional, and local hydrological behaviors via machine learning applied to

580  large-sample datasets. *Hydrology & Earth System Sciences*, *23*(12). https://doi.org/10.5194/hess-

581  23-5089-2019, 2019.

582

583  Lerios, J. L., and Villarica, M. V.: Pattern Extraction of Water Quality Prediction Using Machine

584  Learning Algorithms of Water Reservoir. *International Journal of Mechanical Engineering and*

585  *Robotics Research*, *8*(6), 2019.

586

587    Liu, P., Wang, J., Sangaiah, A. K., Xie, Y., and Yin, X.: Analysis and prediction of water quality

588    using LSTM deep neural networks in IoT environment. *Sustainability*, *11*(7), 2058, 2019.

589

590    Lin, F., Chen, X., and Yao, H.: Evaluating the use of Nash-Sutcliffe efficiency coefficient in

591    goodness-of-fit measures for daily runoff simulation with SWAT. *Journal of Hydrologic*

592    *Engineering*, *22*(11), 05017023. https://doi.org/10.1061/(ASCE)HE.1943-5584.0001580, 2017.

593

594    Mazzocchi, F.: Could Big Data be the end of theory in science? A few remarks on the

595    epistemology of data-driven science. *EMBO reports*, *16*(10), 1250-1255.

596    https://doi.org/10.15252/embr.201541001, 2015.

597

598    Mishra, A., Ahmadisharaf, E., Benham, B. L., Wolfe, M. L., Leman, S. C., Gallagher, D. L., ...

599    and Smith, E. P.: Generalized likelihood uncertainty estimation and Markov chain Monte Carlo

600    simulation to prioritize TMDL pollutant allocations. *Journal of Hydrologic Engineering*, *23*(12),

601    05018025. https://doi.org/10919/93506, 2018.

602

603    Morris, M. D.: Factorial sampling plans for preliminary computational experiments.

604    *Technometrics*, *33*(2), 161-174, 1991.

605

606 Nakhle, P., Ribolzi, O., Boithias, L., Rattanavong, S., Auda, Y., Sayavong, S., Zimmermann, R.,

607 Soulileuth, B., Pando, A., Thammahacksa, C., Rochelle-Newall, E., Santini, W., Martinez, J.M.,

608 Gratiot, N., Pierret, A.: Effects of hydrological regime and land use on in-stream Escherichia coli

609 concentration in the Mekong basin, Lao PDR. Sci. Rep. In press, 2021.

610

611 Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models part I—A

612 discussion of principles. *Journal of hydrology*, *10*(3), 282-290, 1970.

613

614 Nash, S. G. (1984). Newton-type minimization via the Lanczos method. *SIAM Journal on

615 Numerical Analysis*, *21*(4), 770-788.

616

617 Nguyen, H. T. M., Le, Q. T. P., Garnier, J., Janeau, J. L., and Rochelle-Newall, E.: Seasonal

618 variability of faecal indicator bacteria numbers and die-off rates in the Red River basin, North

619 Viet Nam. *Scientific Reports*, *6*(1), 1-12. https://doi.org/10.1038/srep21644, 2016.

620

621 Odonkor, S. T., and Ampofo, J. K.: Escherichia coli as an indicator of bacteriological quality of

622 water: an overview. *Microbiology research*, *4*(1), e2-e2. https://doi.org/10.4081/mr.2013.e2,

623 2013.

624

34

625   Muirhead, R. W., and Meenken, E. D.: Variability of Escherichia coli Concentrations in Rivers

626   during Base-Flow Conditions in New Zealand. *Journal of environmental quality*, *47*(5), 967-973.

627   https://doi.org/10.2134/jeq2017.11.0458, 2018.

628

629   Nair, V., and Hinton, G. E.: Rectified linear units improve restricted boltzmann machines. In

630   *ICML*, January 2010.

631

632   Palmateer, G., McLean, D., Kutas, W. L., and Meissner, S. M.: Suspended particulate/bacterial

633   interaction in agricultural drains. *SS RAO*, 1-40, 1993.

634

635   Pachepsky, Y. A., Blaustein, R. A., Whelan, G., and Shelton, D. R.: Comparing temperature

636   effects on Escherichia coli, Salmonella, and Enterococcus survival in surface waters. *Letters in*

637   *applied microbiology*, *59*(3), 278-283. https://doi.org/10.1111/lam.12272, 2014.

638

639   Pachepsky, Y., Stocker, M., Saldaña, M. O., and Shelton, D.: Enrichment of stream water with

640   fecal indicator organisms during baseflow periods. *Environmental monitoring and assessment*,

641   *189*(2), 51. https://doi.org/10.1007/s10661-016-5763-8, 2017.

642

643  Pachepsky, Y. A., Allende, A., Boithias, L., Cho, K., Jamieson, R., Hofstra, N., and Molina, M.:

644  Microbial water quality: monitoring and modeling. *Journal of environmental quality*, *47*(5), 931-

645  938. https://doi.org/10.2134/jeq2018.07.0277, 2018.

646

647  Pandey, P. K., and Soupir, M. L.: Assessing the impacts of E. coli laden streambed sediment on

648  E. coli loads over a range of flows and sediment characteristics. *JAWRA Journal of the American*

649  *Water Resources Association*, *49*(6), 1261-1269. https://doi.org/10.1111/jawr.12079, 2013.

650

651  Park, Y., Kim, M., Pachepsky, Y., Choi, S. H., Cho, J. G., Jeon, J., and Cho, K. H.: Development

652  of a nowcasting system using machine learning approaches to predict fecal contamination levels

653  at recreational beaches in Korea. *Journal of environmental quality*, *47*(5), 1094-1102.

654  https://www.doi.org/10.2134/jeq2017.11.0425, 2018.

655

656  Pool, S., Vis, M., and Seibert, J.: Evaluating model performance: towards a non-parametric

657  variant of the Kling-Gupta efficiency. *Hydrological Sciences Journal*, *63*(13-14), 1941-1953.

658  2018.

659

660  Pyo, J., Park, L. J., Pachepsky, Y., Baek, S. S., Kim, K., and Cho, K. H.: Using convolutional

661  neural network for predicting cyanobacteria concentrations in river water. *Water Research*, *186*,

662  116349. https://doi.org/10.1016/j.watres.2020.116349, 2020.

663

664    Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., ... and Steinbach, M.:

665    Process-guided deep learning predictions of lake water temperature. *Water Resources Research*,

666    *55*(11), 9173-9190. https://doi.org/10.1029/2019WR024922, 2019.

667

668    Rochelle-Newall, E., Nguyen, T. M. H., Le, T. P. Q., Sengtaheuanghoung, O., and Ribolzi, O.: A

669    short review of fecal indicator bacteria in tropical aquatic ecosystems: knowledge gaps and

670    future directions. *Frontiers in microbiology*, *6*, 308. https://doi.org/10.3389/fmicb.2015.00308,

671    2015.

672

673    Ribolzi, O., Evrard, O., Huon, S., Rochelle-Newall, E., Henri-des-Tureaux, T., Silvera, N., ... and

674    Sengtaheuanghoung, O.: Use of fallout radionuclides (7 Be, 210 Pb) to estimate resuspension of

675    Escherichia coli from streambed sediments during floods in a tropical montane catchment.

676    *Environmental Science and Pollution Research*, *23*(4), 3427-3435, 2016.

677

678    Rumelhart, D. E., Hinton, G. E., and Williams, R. J.: Learning representations by back-

679    propagating errors. *nature*, *323*(6088), 533-536. https://doi.org/10.1038/323533a0, 1986.

680

681    Seong, C. H., Benham, B. L., Hall, K. M., and Kline, K.: Comparison of alternative methods to

682    simulate bacteria concentrations with HSPF under low-flow conditions. *Applied Engineering in*

683    *Agriculture*, *29*(6), 917-931. https://doi.org/10.13031/aea.29.10203, 2013.

684

Hydrology and
Earth System
Sciences
Discussions

Open Access
EGU

685    Singh, A., and Kingsbury, N.: Dual-tree wavelet scattering network with parametric log

686    transformation for object classification. In *2017 IEEE International Conference on Acoustics,*

687    *Speech and Signal Processing (ICASSP)* (pp. 2622-2626). IEEE. 10.1109/ICASSP.2017.7952631,

688    March 2017.

689

690    Song, L., Boithias, L., Sengtaheuanghoung, O., Oeurng, C., Valentin, C., Souksavath, B., ... and

691    Ribolzi, O.: Understory Limits Surface Runoff and Soil Loss in Teak Tree Plantations of

692    Northern Lao PDR. *Water*, *12*(9), 2327. https://doi.org/10.3390/w12092327, 2020.

693

694    Sowah, R. A., Bradshaw, K., Snyder, B., Spidle, D., and Molina, M.: Evaluation of the soil and

695    water assessment tool (SWAT) for simulating E. coli concentrations at the watershed-scale.

696    *Science of the Total Environment*, *746*, 140669. https://doi.org/10.1016/j.scitotenv.2020.140669,

697    2020.

698

699    Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.: Dropout: a

700    simple way to prevent neural networks from overfitting. *The journal of machine learning*

701    *research*, *15*(1), 1929-1958, 2014.

702

703    Thupaki, P., Phanikumar, M. S., Schwab, D. J., Nevers, M. B., and Whitman, R. L.: Evaluating

704    the role of sediment-bacteria interactions on Escherichia coli concentrations at beaches in

705    southern Lake Michigan. *Journal of geophysical research: oceans*, *118*(12), 7049-7065.

706    https://doi.org/10.1002/2013JC008919, 2013.

Hydrology and
Earth System
Sciences
Open Access
Discussions
EGU

707

708    Troeger, C., Forouzanfar, M., Rao, P. C., Khalil, I., Brown, A., Reiner Jr, R. C., ... and

709    Alemayohu, M. A.: Estimates of global, regional, and national morbidity, mortality, and

710    aetiologies of diarrhoeal diseases: a systematic analysis for the Global Burden of Disease Study

711    2015. *The Lancet Infectious Diseases*, *17*(9), 909-948.

712    doi:https://doi.org/10.1073/pnas.1109326109, 2017.

713

714    Van Rossum, G.: Python programming language. In *USENIX annual technical conference* (Vol.

715    41, p. 36), June 2007.

716

717    Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., ... and

718    van der Walt, S. J.: SciPy 1.0: fundamental algorithms for scientific computing in Python.

719    *Nature methods*, *17*(3), 261-272. https://doi.org/10.1038/s41592-019-0686-2, 2020

720

721    Wang, X., Zhang, F., and Ding, J.: Evaluation of water quality based on a machine learning

722    algorithm and water quality index for the Ebinur Lake Watershed, China. *Scientific reports*, *7*(1),

723    1-18. https://doi.org/10.1038/s41598-017-12853-y, 2017

724

725  Xiang, Z., Yan, J., and Demir, I.: A rainfall-runoff model with LSTM-based sequence-to-

726  sequence learning. *Water resources research*, *56*(1), e2019WR025326.

727  https://doi.org/10.1029/2019WR025326, 2020.

728

729  Van der Leeuw, S. E.: Why model? *Cybernetics and Systems*, *35*(2-3), 117-128.

730  https://doi.org/10.1080/01969720490426803, 2004.

731

732  Yagow, G., Dillaha, T., Mostaghimi, S., Brannan, K., Heatwole, C., and Wolfe, M. L.: TMDL

733  modeling of fecal coliform bacteria with HSPF. In *2001 ASAE Annual Meeting* (p. 1). American

734  Society of Agricultural and Biological Engineers, 1998.

735

736  Zheng, A., and Casari, A.: *Feature engineering for machine learning: principles and techniques*

737  *for data scientists*. " O'Reilly Media, Inc.", 2018.

738 **Table 1:** Optimal values and range of HSPF parameters for surface and sub-surface flows and *E.*
739 *coli* concentration. Bold parameters were optimized during flow calibration process. All
740 parameters related to *E. coli* were optimized during model calibration.

| | | Land use | | | | Lower Limit | Upper Limit |
|---|---|---|---|---|---|---|---|
| | Parameters | Forest | Teak | Fallow | Annual Crop | | |
| Surface and sub-surface flow | INFILT | **0.31** | **0.39** | **0.39** | **0.36** | 0.001 | 0.5 |
| | INFILD | 2.0 | **1.94** | **1.95** | **1.55** | 1.0 | 3.0 |
| | INTFW | 2.60 | **7.01** | **7.01** | 5.64 | 1.0 | 10.0 |
| | UZSN | 1.36 | 1.47 | **0.84** | **1.24** | 0.05 | 2.0 |
| | LZSN | **8.88** | **9.43** | **4.18** | **8.66** | 2.0 | 10.0 |
| | AGWETP | 0.02 | **0.007** | 0.02 | **0.06** | 0.0 | 0.2 |
| | NSUR | **0.18** | **0.39** | **0.15** | **0.43** | 0.05 | 0.5 |
| | BASETP | **0.05** | **0.09** | 0.095 | **0.003** | 0.0 | 0.2 |
| | DEEPFR | 0.28 | 0.16 | 0.21 | **0.20** | 0.0 | 0.5 |
| *E. coli* concentration | SQOLIM MF | 4.99 | 1.35 | 2.04 | 0.53 | 0.5 | 10 |
| | WSQOP | 9.12 | 9.31 | 8.87 | 9.38 | 0.1 | 10.0 |
| | IOQC | 5367 | 8337 | 8380 | 8756 | 1000 | 10000 |
| | AOQC | 8672 | 7474 | 5465 | 8776 | 1000 | 10000 |
| | FSTDEC | | | 3.04 | | 0.1 | 10.0 |
| | THFST | | | 1.92 | | 1.01 | 2.0 |

741

742

**Table 2:** Hyper-parameters of LSTM for surface flow, sub-surface flow and *E. coli*

concentration simulation.

| Parameter | Surface and sub-surface flow | *E. coli* |
|---|---|---|
| **Activation function (LSTM Layer)** | Rectified Linear Unit (ReLU) | Rectified Linear Unit (ReLU) |
| **Activation function (Dense Layer)** | Rectified Linear Unit (ReLU) | Rectified Linear Unit (ReLU) |
| **Batch size** | 128 | 16 |
| **Learning rate** | 1e-5 | 1e-6 |
| **lookback steps** | 5 hours | 5 hours |
| **Dropout** | 0.3 | 0.3 |
| **Hidden units** | 64 | 100 |
| **Input data** | Rainfall, Solar Radiation, Air Temperature, Potential Evapotranspiration | Rainfall, Surface flow, Sub-surface flow, Land use, Bacteria source |
| **Calibration epochs** | 500 | 7000 |
| **Training samples** | 490000 | 182 |
| **Test samples** | 210000 | 73 |

745

746 **Table 3:** Performance metrics of HSPF and LSTM model for surface and sub-surface flow.

| Model | Flow Type | Scenario | MSE $(m^3s^{-1})$ | NSE | PBIAS |
|-------|-----------|----------|-----|-----|-------|
| **HSPF** | **Surface Flow** | Calibration | 6.4e-4 | -0.02 | -59 |
| | | Validation | 4.7e-5 | -0.7 | -28 |
| | **Sub-surface Flow** | Calibration | 2.7e-4 | 0.49 | -51 |
| | | Validation | 5.e-4 | 0.59 | -22 |
| **LSTM** | **Surface Flow** | Calibration | 1.4e-4 | 0.56 | -48 |
| | | Validation | 1.9e-4 | 0.51 | -63 |
| | **Sub-surface Flow** | Calibration | 5.4e-3 | 0.69 | -42 |
| | | Validation | 5.9e-3 | 0.64 | -46 |

747

748    **Table 4:** Performance metrics of HSPF and LSTM for *E. coli* concentration simulation.

| Model | Scenario | MSE (MPN 100 mL$^{-1}$) | NSE | PBIAS |
|-------|----------|-------------------------|-----|-------|
| **HSPF** | Calibration | $1.4e^8$ | -0.29 | -58 |
|  | Validation | $1.9e^8$ | -3.01 | 73.01 |
| **LSTM** | Calibration | $7.1e^6$ | 0.39 | -1.49 |
|  | Validation | $3.0e^7$ | 0.35 | 62.72 |

749

750

**Figure 1:** Location of the study area. The study area is located near Luang Prabang in northern

Lao PDR. The monitoring station is located at the outlet of the catchment, where water level is

recorded, and where water samples are collected for *E. coli* concentration measurement. Climate

data was measured at the meteorological station.

755

756

757 **Figure 2:** Structure of the LSTM Model. Environmental data is used to predict surface flow and

758 sub-surface flow. Simulated flows along with bacteria source, land-use information and rainfall

759 is used to simulate *E. coli* concentration. The 'n' represents the length of input data used by

760 LSTM.

**Figure 3:** Performance of HSPF model with different objective functions (e.g., M-surface, N-Surface, M-subsurface and N-subsurface). The color indicates the value of MSE and NSE. M-surface is the objective function based on MSE and surface flow, N-surface is the objective function based on NSE and surface flow, M-subsurface is the objective function based on MSE and sub-surface flow, and N-subsurface is the objective function based on NSE and sub-surface flow.

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

770

**Figure 4:** Hydrological simulation from HSPF: (a) Measured rainfall, (b) Simulated and

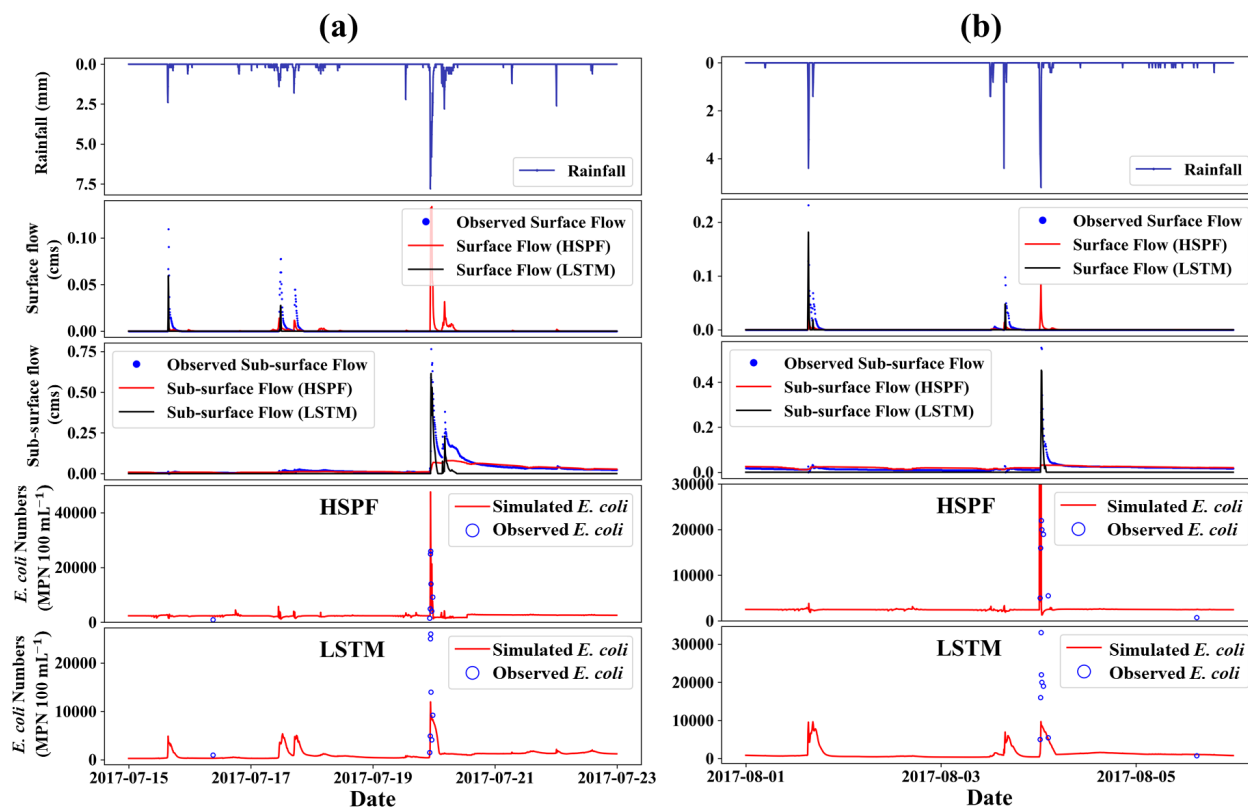observed surface flow, and (c) Simulated and observed sub-surface flow.

Hydrology and
Earth System
Sciences

Open Access

Discussions

EGU

773

**Figure 5:** Hydrological simulation from LSTM: (a) Measured rainfall, (b) Simulated and

775    observed surface flow, and (c) Simulated and observed sub-surface flow.

**Figure 6:** Comparison of the hydrological simulation during storm events: (a) MSE value of the surface flow, (b) MSE value of the sub-surface flow, (c) NSE value of the surface flow and, (d) NSE value of the sub-surface flow.
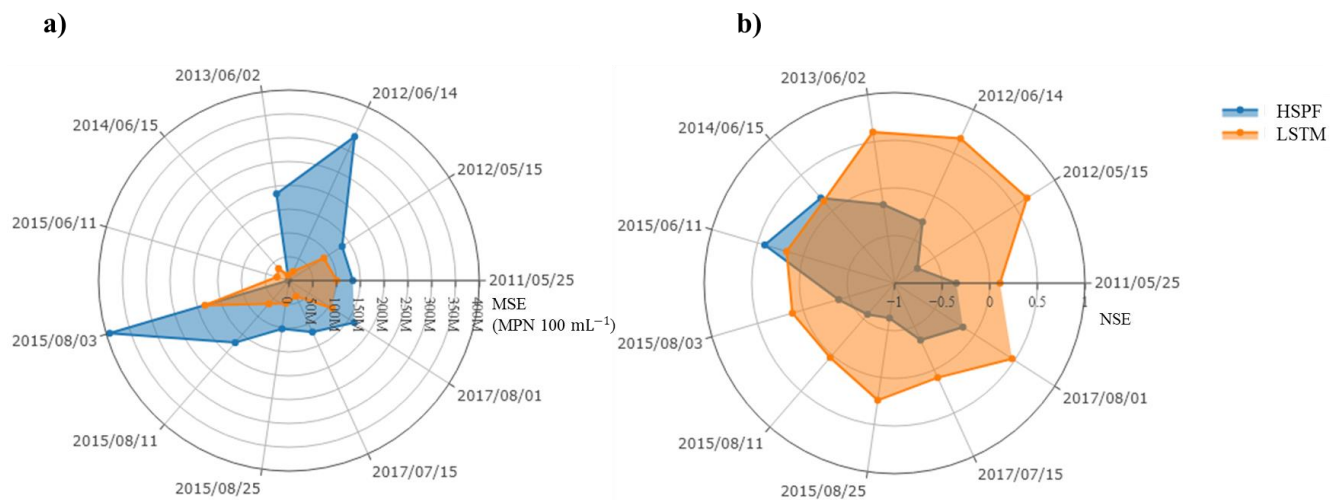
Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

**Figure 7:** *E. coli* simulation from LSTM and HSPF: (a) Measured rainfall, (b) Observed surface and sub-surface flow, (c) Simulated and observed *E. coli* using HSPF, and (d) Simulated and observed *E. coli* using LSTM.
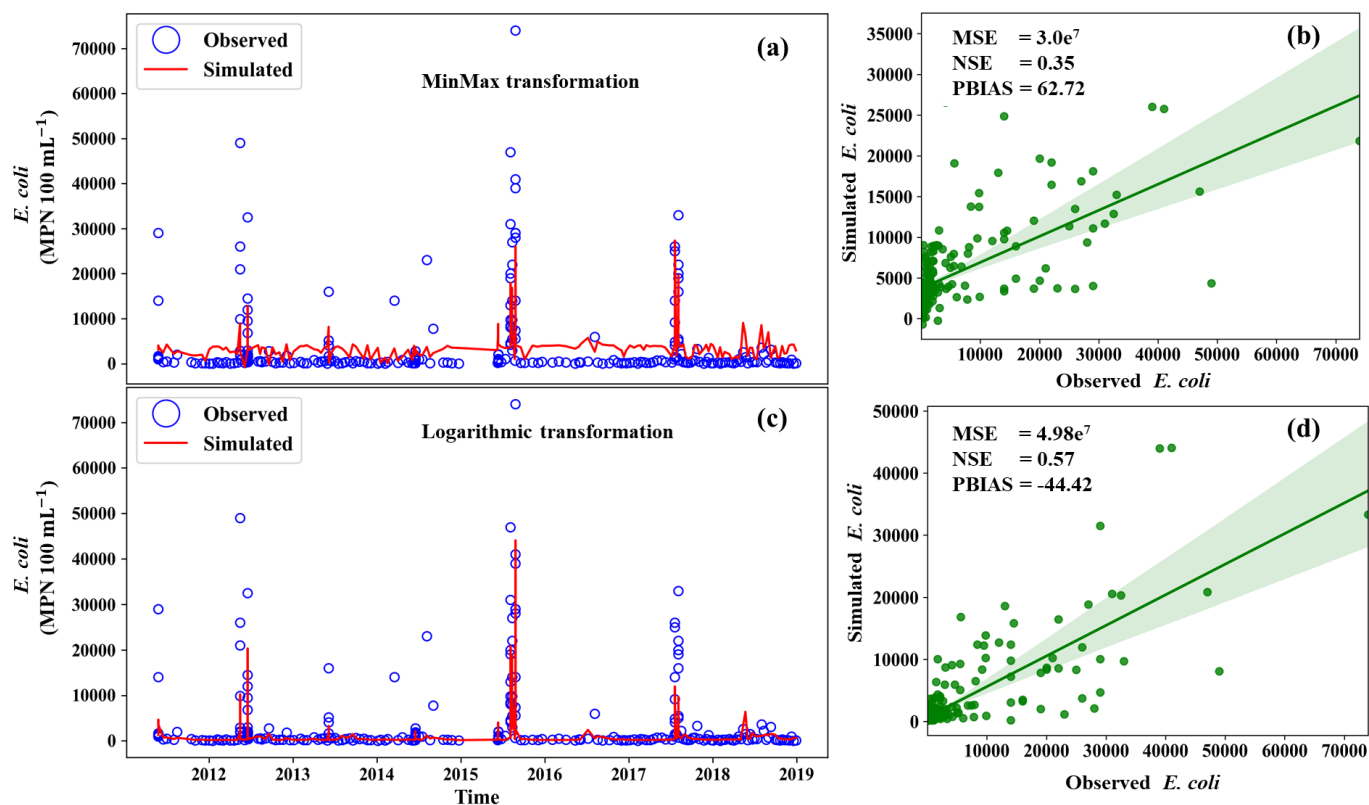
**(a)** **(b)**

784   **Figure 8:** *E. coli* concentration of HSPF and LSTM on July 15 and August 1, 2017. Both storm

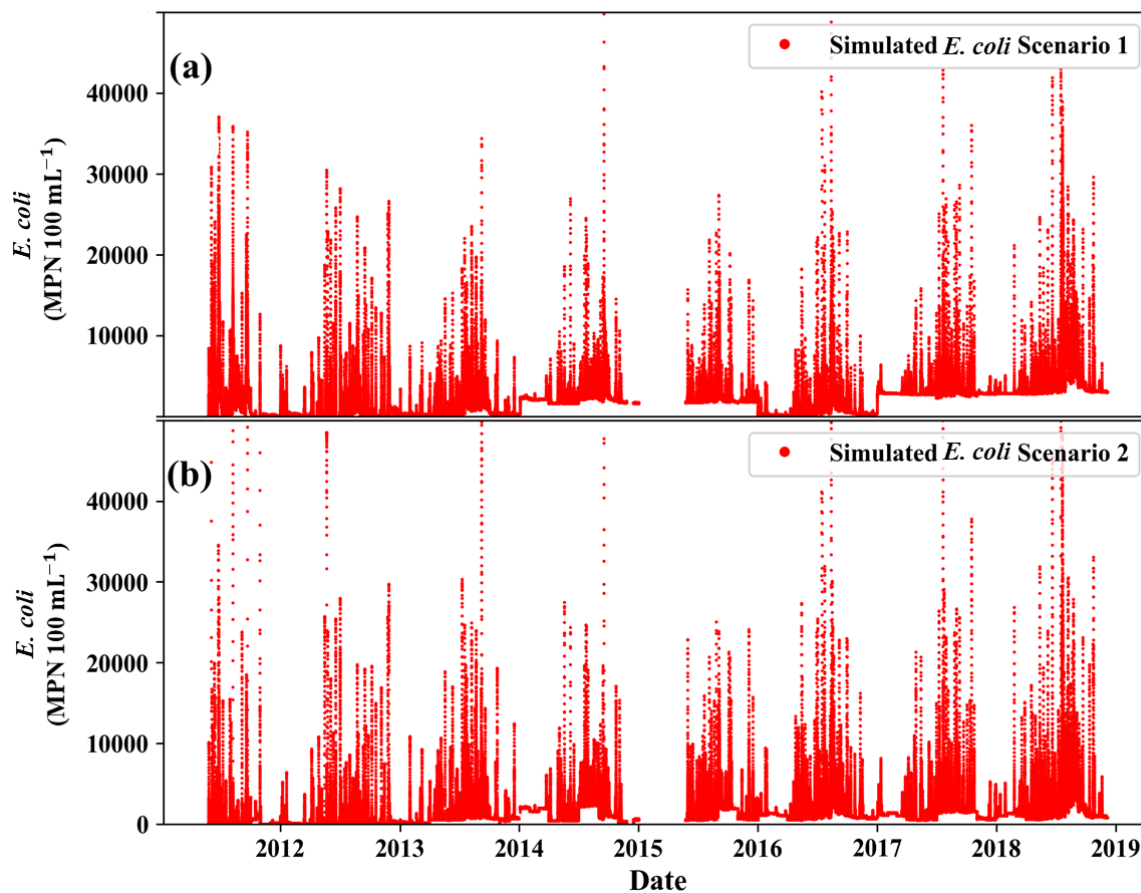785   events were affiliated in validation period.

a)

b)

786



**Figure 9:** Comparison of the *E. coli* simulation during storm events: (a) MSE values and (b)

788     NSE values.

**Figure 10:** Comparison of *E. coil* concentration simulation with the transformation method: (a)

and (c) indicate the *E. coli* simulation using minmax transformation and logarithmic

transformation, respectively. (b) and (d) indicate the scatter plot of *E. coli* using minmax

transformation and logarithmic transformation, respectively.

793

**Figure 11:** Changes in *E. coli* sources with land use change scenarios. Scenario 1 used land use

change and bacterial source information. Scenarios 2 used the bacterial source by the fraction of

each land use.