1    **A neighborhood statistics model for predicting stream pathogen**

2    **indicator levels**

3    Pramod K Pandey [1*], Gregory B Pasternack[2], Mahbubul Majumder[3], Michelle L Soupir[4],

4    Mark S. Kaiser[3],

5    *[1]Department of Population Health and Reproduction, University of California, Davis,*

6    *California 95616, USA*

7    *[2]Department of Land, Air and Water Resources, University of California, Davis,*
8    *California 95616, USA*

9    *[3]Department of Statistics & Statistical Laboratory, Iowa State University,*

10   *Ames, IA 50014, USA*

11   *[4]Department of Agricultural and Biosystems Engineering, Iowa State University,*

12   *Ames, Iowa 50014, USA*

13

14   *[*]Corresponding author email: pkpandey@ucdavis.edu*

15

16

17

18

19  **ABSTRACT**

20  Because elevated levels of water borne *E. coli* in streams are a leading cause of water quality

21  impairments in the U.S., water quality managers need tools for predicting aqueous *E. coli* levels.

22  Presently *E. coli* levels may be predicted using complex mechanistic models that have a high

23  degree of unchecked uncertainty or simpler statistical models.   To assess spatio-temporal

24  patterns of instream *E. coli* levels, herein we measured *E. coli*, a pathogen indicator, at 16 sites

25  (at four different times) within the Squaw Creek watershed, Iowa, and subsequently the Markov

26  Random Field model was exploited to develop a neighborhood statistics model for predicting

27  instream *E. coli* levels. Two observed covariates, local water temperature ($^{o}$C) and mean cross-

28  sectional depth (m), were used as inputs to the model. Predictions of *E. coli* levels in the water

29  column were compared with independent observational data collected from sixteen in-stream

30  locations.  The results revealed that spatio-temporal averages of predicted and observed *E. coli*

31  levels were extremely close (all within factor of 2), while 66% of individual predicted *E. coli*

32  concentrations were within a factor of 2 of the observed values. In only one event, difference

33  between prediction and observation was beyond 1 order of magnitude.  The mean of all predicted

34  values at sixteen locations was approximately 1% higher than the mean of the observed values.

35  The approach presented here will be useful while assessing instream contaminations such as

36  pathogen/pathogen indicator levels at watershed scale.

37  *Keywords*: stream water; *E. coli*; neighborhood structures; Markov Random Field model

38  **1. INTRODUCTION**

39  Unsafe levels of pathogens in ambient water bodies such as streams, ground water, lakes and

40  reservoirs, estuaries, and coastal waters are a major concern for the environment and pose a

41  serious risk to public health (U.S. EPA 2012a). Predictive models have been developed to

42  simulate watershed-scale hydrological processes and associated bacterial transport and

43  interactions. In this study, we report spatio-temporal patterns of *E.coli* levels in a stream network

44  and then introduce the use of a neighborhood statistics model for predicting stream pathogen

45  indicator levels.

46  **1.1 Research motivation**

47  Water borne pathogens have been linked to various diseases, including diarrhea, malaria, yellow

48  fever, dengue, hepatitis A, Hepatitis E, and typhoid fever. For example, approximately 37.5% of

49  diarrhea cases in developing countries are due to contaminated water. Even in a developed

50  country such as United States, approximately 60% of total diarrhea cases are attributable to

51  unsafe water and poor hygiene. According to the World Health Organization (WHO)

52  approximately 4% of the global disease burden is caused by contaminated water (WHO, 2010);

53  improving water quality is a viable option for mitigating health risk.

54      One of the leading causes of stream water quality impairment in the U.S. is elevated

55  levels of pathogens such as *E. coli*, which is also an indicator of the presence of other pathogens.

56  According to the EPA's national summary of impaired waters, approximately 40,235 water

57  bodies are impaired; approximately 15% of the total 71,363 impairments are due to water borne

58  pathogens (EPA 2013). One major source of *E. coli* in streams is diffuse pollution (i.e., non-point

59  source pollution from agriculture).  For example, in Iowa, where approximately 75% of the

60  watershed is dominated by cropping land and precipitation is a major source of water for

61  agriculture, 69% of assessed streams are impaired, with 27.5% of those due to high levels of

62  pathogens (U.S. EPA 2012a). In California, where approximately 43% of land is dominated by

63  agriculture, but irrigation is the major source of agriculture water, approximately 89% of

64  assessed streams are contaminated (U.S. EPA 2013). Currently, 15.8% of total streams in

65  California are assessed, and 15% of assessed streams are impaired by pathogens.

66  **1.2 Predicting bacterial concentrations**

67  Evaluating public health risks caused by water borne pathogens requires predictions of pathogen

68  levels in ambient water bodies such as streams.  In turn, predicting instream pathogen

69  concentrations requires understanding fate and transport of pathogens at watershed scale.

70  Previously process based modelling approaches have been used extensively for predicting

71  pathogen levels in streams (Hipsey et al. 2008; Rehmann and Soupir, 2009; Pandey et al.,

72  2012a,b; Droppo et al. 2009; Jamieson et al. 2005; Schilling et al. 2009; Wilkes et al. 2011).

73  Jamieson et al. (2005) used stream bed stresses and stream flow while computing stream water

74  column *E. coli* levels. A study by Hipsey et al. (2008) emphasized sediment properties that

75  potentially affects stream water column *E. coli* levels. Rehmann and Soupir (2009) used a one-

76  dimensional approach to understand the impacts of interactions between water column and

77  streambed sediment on *E. coli* concentrations in streams. Pandey et al. (2012a) calculated *E. coli*

78  resuspension rate, while predicting water column *E. coli* levels at watershed scale. Similarly,

79  Kim et al. (2010) embedded a resuspension of *E. coli* to the existing Soil Water Assessment Tool

80  (SWAT) for predicting in stream *E. coli* levels, while Parajuli et al. (2009) used the SWAT

81  model for predicting instream *E. coli* levels without adding resuspension process. While previous

82  approaches considerably enhanced the understanding of bacteria fate and transport in streams,

83  the development of relatively simpler approaches, such as the statistical model described herein,

84  can be another option for predicting instream *E. coli* levels.

85    In addition to using process-based models, many previous studies implemented such

86  models of instream *E. coli* levels within geographical information systems (GIS) taking

87  advantage of geospatial data. For example, Pandey et al. (2012b) estimated watershed indexes

88  considering undisturbed land cover (e.g., wetlands, vegetated streams) and disturbed land cover

89  (e.g., crop land, crop land receiving animal manure, urban land) for identifying the relationships

90  between in-stream *E. coli* levels and watershed characteristics. Studies by Rothwell et al.

91  (2010a;b) exploited GIS tools to identify the relationships between water chemistry (e.g., pH,

92  sulphate, cations, and nutrients) and a watershed's land cover, topography, soil, and hydrology.

93  These studies reported that stream water quality is significantly linked to watershed

94  characteristics.

95    Understanding how climate and land surface characteristics (e.g., land cover, soil,

96  topography, and geology) interact at the watershed scale to generate runoff and transport

97  materials is crucial for predicting and ultimately mitigating in-stream pathogen and pathogen-

98  indicator levels. Watershed-scale models that account for these relationships to simulate

99  processes and fluxes can help with development and implementation of a watershed management

100  plan for improving in-stream water quality. For example, SWAT has been extensively used

101  (Parajuli et al. 2009; Cho et al. 2010; Kim et al. 2010) to predict in-stream water *E. coli* levels. In

102  the SWAT model, watershed characteristics such as cropland, grazing land, livestock density,

103  decay of bacteria, and climate of the watershed (rainfall and temperature) are used as inputs for

104  predicting bacteria levels in streams. Previous studies have shown that SWAT can help deriving

105  suitable land management plans and guidelines supportive for mitigating instream pathogen

106  levels.

107         Despite the potential opportunities for a model-based management approach, comparison

108 between model predictions and observations of in-stream bacteria levels clearly indicates that

109 considerable improvements in the existing models are required before their potential is reached

110 (Nagels et al. 2002; Rehmann and Soupir 2009; Hipsey et al. 2008; Droppo et al. 2009; Dorner et

111 al. 2006; Pachepsky and Shelton 2011). For instance, Dorner et al. (2006) developed a

112 hydrological model (WATFLOOD model was augmented with a pathogen transport model) and

113 found that daily predictions of *E. coli* levels varied from 1 to 4 orders of magnitude of observed

114 values (more than 70 observations were compared with predicted values). Similarly Kim et al.

115 (2010) predictions using SWAT model varied from 1 to 3 orders of magnitude of the observed

116 values (more than 150 observations were compared with predicted values).

117         To address the underlying deficiencies, studies have suggested everything from adding

118 more physical processes to improving statistical methods. For example, one idea has been to

119 improve the formulations of in-stream processes such as resuspension of *E. coli* from streambed

120 sediment to the water column in order to improve existing water quality models for bacteria

121 predictions (Muirhead et al. 2004; Bai and Lung 2005; Jamieson et al. 2005). Another approach

122 for improving in-stream *E. coli* predictions could be combining the capability of GIS data and

123 statistics.

124 **1.3 Research objectives**

125 The overall goal of this study was to explore *E. coli* levels in a watershed stream network and

126 test the value of a spatial neighborhood statistics model, the Markov Random Field model, to

127 predict *E. coli* levels in streams. The model was formulated and tested for Iowa's Squaw Creek

128 Watershed, which is an agriculture-dominated watershed. Non-point source pollution is known

129    to be the leading cause of bacterial contamination in the streams. This study builds upon the

130    work of Kaiser (2010), who previously used this approach successfully to predict nitrate

131    concentrations in the Des Moines River, Iowa prior to impoundment in Saylorville Reservoir.

132    The study used stream flow and nitrate data (2954 observations from January, 1982 to

133    December, 1996) from seven gaging stations along the Des Moines River from Boone to Pella

134    (about 116 miles). The specific objectives of this study were to (i) observe and analyze how *E.*

135    *coli* levels vary between four different times in relation to the covariates of water temperature

136    and water depth, (ii) compare *E. coli* levels in tributaries versus the mainstem channel, and (iii)

137    develop and assess the predictive prowess of a neighborhood statistics model.

138    **2. METHODS**

139    **2.1 Field setting and observations**

140    Squaw Creek passes through Story, Webster, Hamilton, and Boone Counties of Iowa (Figure 1).

141    The Squaw Creek watershed, Hydrologic Unit Code (HUC) 10 (ID 0708010503), has a total

142    drainage area of 592.4 sq km and average slope of 2%. The watershed's humid continental

143    climate, Köppen climate classification *Dfa*, receives an average annual precipitation of 910 mm.

144    In general, December and January are the coldest months (temperature variation from -1 to -10

145    ⁰C), and June and July are the warmest month (temperature variation from 30 to 35 ⁰C). The

146    mainstem length (i.e., Squaw Creek) is 60.5 km and the total stream length (including tributaries)

147    within the watershed is 346.7 km. There are 75 first order streams. Approximately 74% of the

148    watershed is under agriculture: corn 41% and soybeans 33%. Forest cover is about 2.7%, and the

149    total grassland is about 17% of the total watershed.

150

151    Corn and soybean are two major crops grown in the watershed. Planting and harvesting

152    of corn in Iowa are done generally between April and October. Soybeans are usually planted in

153    May after completing corn planning, with soybean harvesting in early-mid October. Corn is the

154    major crop receiving liquid manure (mostly in fall) from confined animal feeding operations.

155    Water samples (total 64 observations) in support of model development were collected at 16

156    locations along the stream on 27[th] June (t = 1), 6[th] July (t = 2), 17[th] July (t = 3), and 17[th] October

157    2009 (t = 4). Eight locations (1 − 8) were located in tributaries and another eight (9 − 16) were

158    located along the mainstem (shown in Figure 1). Samples were collected using a Horizontal

159    Polycarbonate Water Bottle Sampler (2.2 L, Forestry Suppliers Inc., Jackson, Mississippi City,

160    USA) by lowering the instrument from a bridge into the center ($\approx$ 15 cm below surface water) of

161    the stream at the sampling location. After sample collections, samples were stored at 4 $^0$C (in a

162    cooler) immediately and were analyzed (triplicate) within 24 hours. Membrane filtration

163    technique (US EPA method 1603) has been used for *E. coli* enumeration using modified mTEC

164    agar (Difco$^{TM}$, Modified mTEC agar, Becton, Dickinson and Company, Sparks, MD, USA)

165    (APHA 1999). In addition to *E. coli* enumeration, we also measured the stream water column

166    depth (m) and temperature ($^0$C), while collecting water samples. Average stream water column

167    depth along the transact at each sampling location was determined by marking off equal intervals

168    of approximately 60 cm along the measuring string and then the mean of the water depths was

169    used for analysis. Streamflow data was obtained for the U. S. Geological Survey gaging station

170    (ID 05470500) at site 16 (Figure 1). Climate data, precipitation and temperature, were obtained

171    for Ames City (lat 42.02, long − 93.77) using Iowa Mesonet, Iowa State University, Ames, Iowa

172    (IEM 2012).

173        To address the first objective, we performed a comparative analysis of event based

174    observation data of water column *E. coli* levels, stream water depths, and stream water

175    temperatures that were collected at 16 locations along the stream at four different times. The

176    second objective was addressed by exploiting the use of Mann-Whitney U test, a non-parametric

177    test. The test was used to compare *E. coli* levels across all sites in tributaries and main stem at

178    each time. Further, Pearson correlation coefficients were estimated to relate *E. coli* levels among

179    sampling locations. The third objective was resolved by developing a statistics model that uses a

180    neighborhood structure linking *E. coli* levels in downstream locations with upstream sampling

181    locations. Subsequently model predictions were compared with observations to verify the

182    model's predictability. In addition, EPA's water quality criteria of indicator organisms (*E. coli*)

183    of fresh water were used as reference points, while comparing the model predictions and

184    observations.

185    **2.2 Neighborhood statistics model**

186    To develop the model for the study area, we developed the conditionally specified model for

187    Squaw Creek. In equation 1, *Y* is a random variable, and $s_i \equiv (l, t)$ where $l$ is sampling locations

188    (1–16), and t is sampling events (1 – 4).

$$Y \equiv \{Y(s_i) : i = 1,...64\}$$
189    $$= \{Y(l,t) : l = 1,...16; t = 1,...4\} \tag{1}$$

190    We assume that the temporal distributions of *E. coli* at a station, conditional on all stations

191    upstream depends only on closest upstream stations. Based on sampling locations shown in

192    Figure 1, neighborhood structures were developed, which are shown in Table 1. The criteria of

193 neighbor selection were defined based on inflowing tributaries and sampling locations. For each

194 sampling location, upstream tributaries, and immediate upstream and downstream locations were

195 defined as neighbors. For example, location 10 has two tributaries just upstream, therefore these

196 two tributaries (locations 1 and 2) are considered neighbors in addition to location 11

197 (immediately downstream).

198 We also assume that measurements of *E. coli* concentrations are independent in time. This leads

199 us to define neighbors of $Y(s_i)$ as:

200
$$N_i \equiv \left\{ s_j : s_j \in \left\{ (l-1,t), (l+1,t) \right\}; i = 1, ..., n \right\}$$
(2)

201 Then

202
$$\left[ Y(s_i) \middle| \left\{ Y(s_j) : j \neq i \right\} \right] = \left[ Y(s_i) \middle| Y(N_i) \right]$$
(3)

203 For $i = 1, ...., n$ let $Y(s_i)$ have conditional density

204
$$f\left( y(s_i) \middle| y(N_i) \right) = Gau\left( \mu_i, \tau^2 \right)$$
(4)

205 where

206
$$\mu_i = \theta_i + \sum_{j \in N_i} c_{i,j} \left( y(s_j) - \theta_j \right)$$
(5)

207 subject to $c_{i,j} = c_j$ where

208     $\theta \equiv (\theta_1,...,\theta_n)^T$ $\theta = \{\theta\}$ is the parameter vector of marginal mean that incorporate the covariates

209     $X_i$, i = 1, 2, .., p. We used two covariates (p = 2) temperature ($^oC$) and stream water depth (m).

210     Thus we have

211     $\theta = \beta_0 + \beta_1 X_1 + \beta_2 X_2 = X\beta$.                                           (6)

212     The joint distribution (Besag 1974; Cressie 1993) is:

213     $Y \approx Gau\left(\theta;(I-C)^{-1} I\tau^2\right)$                                         (7)

214     where

215     $C \equiv \left[c_{i,j}\right]_{N \times N}$ and $c_{i,j} = 0$ if $j \notin N_i$                           (8)

216     For this model, C has the form

217     $C = \eta I_n H$                                                 (9)

218     where H is a block diagonal matrix of size 64 × 64; and each block (size = 16 × 16) consists of

219     the neighborhood structures based on inflowing tributaries and sampling locations. The

220     neighborhood was defined as 1 if two locations are neighbor; and 0 otherwise.

221     To obtain the estimates of these parameters we apply the maximum likelihood approach (Kaiser

222     and Nordman 2012; Kaiser 2010). The Log likelihood function for the above model is:

223     $L(\beta,\tau^2,\eta) = (1/2)Log(|I-C|) - (N/2)Log(2\pi\tau^2)$                 (10)
$\qquad\qquad -(1/(2\tau^2))(y-X\beta)^T(1-C)(y-X\beta).$

224    An advantage of this model specification is that for any given η the maximum likelihood

225    estimate (MLE) β and $\tau^2$ are the closed form solutions, which are given by:

226 $$\hat{\beta} = \left[ X^T (I - C) X \right]^{-1} \left[ X^T (1 - C) y \right] \tag{11}$$

227

228 $$\hat{\tau}^2 = \frac{1}{n} \left( y - X \hat{\beta} \right)^T (I - C) \left( y - X \hat{\beta} \right) \tag{12}$$

229    Once we have MLE of β and $\tau^2$ we can plug these values into Equation 9 to get the likelihood

230    for η that gives us MLE for η as shown in Figure 2. Again plugging in the MLE of η in above

231    two equations, we obtain the estimates of β and $\tau^2$. We obtain the confidence intervals of the

232    parameters using the maximum likelihood approach described by Kaiser (2010), Cressie (1993)

233    and Besag (1974). The values of estimated parameters $\tau^2$, η, $\beta_0$, $\beta_1$, $\beta_2$ are shown in Table 2.

234    These values were used for predicting the *E. coli* concentrations at each sampling location.

235    **3. RESULTS AND DISCUSSION**

236    **3.1 Event-based observations**

237    Comparing the four sampling events, each one showed a different range of *E. coli* concentrations

238    and there were identifiable factors explaining the observed differences, which are shown in

239    Figure 3. As an example, *E. coli* levels varied from 144 to 944 CFU/100 mL in the first spatial

240    sampling event (t = 1). During the second sampling event (t = 2), *E. coli* levels varied from 336

241    to 633 CFU/100 mL, which was a narrower range than for event 1. Prior to this sampling event,

242    the watershed witnessed around 50 mm of cumulative rainfall in the first two weeks of June.

243    Between the first two sampling events cumulative rainfall was less than 1 mm. The average of *E.*

244    *coli* levels at 16 locations (shown in Figure 3) at t = 1 was 30% higher than that of at t = 2.

245        During t = 3, *E. coli* levels varied from 225 to 5467 CFU/100 mL. One location

246    (sampling point 7 of Figure 1) showed the maximum large *E. coli* level. Though between t = 2

247    and t = 3, there was no additional rainfall, and streamflow was also identical to preceding

248    sampling events, at t = 3, the average of *E. coli* levels at 16 locations was 78% and 132% higher

249    than that of during t = 1 and t = 2, respectively.

250        During t = 4, *E. coli* levels in tributaries varied from 53 − 333 CFU/100 mL, while in

251    mainstem variation was from 17 − 120 CFU/100 mL. *E. coli* levels during t = 4 were

252    considerably low compared to t = 1, 2, 3. The average *E. coli* level at t = 4 was only 14% of the

253    *E. coli* levels at t = 3. Between t = 3 and t = 4, the cumulative rainfall was only 6.5 mm;

254    however, temperature was considerably lower. For instance, the minimum and maximum daily

255    air temperatures at t = 3 were 11.3 and 19.2 °C, respectively, while at t = 4, these values were −

256    0.5 and 7.4 °C, respectively. During this sampling event, stream flow was 0.13 m³/s, which is

257    about 90% lower than that during t = 3. Overall, event-scale results indicated that winter season

258    (i.e., low temperature) could be the potential reason for low *E. coli* levels.

259        Stream water temperatures and stream water depths are shown in Figure 4. The average

260    daily temperatures during t = 1, 2, 3, and 4 were 24.5, 22.9, 15.3, and 3.5 °C, while average

261    stream water temperatures were 21.3 ± 2.5 °C, 24.2 ± 1.2 °C, 19.9 ±1.9 °C, and 12.2 ± 3 °C,

262    respectively. The stream water depths during these sampling events (0.5, 0.5, 0.4, 0.3 m,

263    respectively) were generally similar, but declined for the last two events, with the value at t = 4

264    being the lowest observed. A total rainfall of 146 mm occurred May 1 to June 27, 2009. The

265    average streamflow for the same period was 9.6 (± 6.4) m³/s with a range from 3 to 30.4 m³/s.

266    **3.2 Tributaries vs. mainstem *E. coli* level analyses**

267 Results of Mann-Whitney U test indicated that there was no significant difference (significant

268 level of 0.05) in *E. coli* levels among t = 1 and t = 2 (all 16 sampling locations). *E. coli* levels

269 among t = 2 and t = 3 were significantly different. There was also significant difference among *E.*

270 *coli* levels of t = 3 and t = 4, and t = 2 and t = 4. Further, there was significant difference in *E.*

271 *coli* levels among t = 1 and t = 4.

272       A Pearson correlation matrix relating *E. coli* observations among 16 sites are shown in

273 Table 3. Analysis showed significant correlations among the sampling locations (Table 3). Out

274 of 120 correlations of sixteen sampling locations, 65% have shown high correlation (r > 0.70; p

275 = 0.05). Relatively greater levels of correlation existed among proximal locations, particularly

276 along the Squaw Creek. For example, locations: 1, 2 and 10; and 12, 13, and 14; and 14, 15 and

277 16.

278 **3.3 Model results**

279 The neighborhood statistics model implemented to predict in-stream *E. coli* levels yielded values

280 within the range observed. Whereas the model produced similar *E. coli* concentrations within a

281 relatively narrow range or both tributaries (Figure 5A) and the mainstem stream (Figure 5B),

282 sampling observations showed a wider range in both settings. As shown in the figure, the model

283 was not able to predict very high and low values. Compared to low values, the model predictions

284 were reasonable well for higher *E. coli* levels.

285       Time averages of observed and predicted *E. coli* levels for tributaries and mainstem were

286 very similar, but those for predictions showed less spatial variability (Figure 5C,D). The spatio-

287 temporal average of tributary observations was 341 CFU/100 ml, while that of tributary

288    predictions was 343 CFU/100 ml. Similarly, that of mainstem observations and predictions was

289    337 and 339 CFU/100 ml, respectively.  These averages are all extremely close.

290         Besides local water temperature and depth, many other local parameters of natural

291    streams, such as channel geometry, nutrient concentrations, solar radiation, and dissolved oxygen

292    also impact *E. coli* levels (Hipsey et al. 2008). In this model we use only two covariates stream

293    water depth and temperature, which might be the reason for the relatively large difference (in

294    few predictions) between measured and predicted values (Fig 5A & B); however, considering the

295    uncertainties involved in predicting *E. coli* levels in natural streams, which is influenced by

296    many factors such as grazing operations, livestock density, cropping land, and land management

297    practices, this parsimonious model can be considered reasonably good for predicting instream *E.*

298    *coli* concentrations. We anticipate availability of a larger dataset could improve the model

299    results.

300         Comparing the predictions of this study with previous ones (Kim et al. 2010; Dorner et

301    al. 2006), the model predictions fit reasonably well. For instance, in the referenced studies

302    predictions were only within $1 - 4$ orders of magnitude of the observations, while in this study

303    the average of predicted values were within a factor of 2 of the observed values (Figs. 5C,D),

304    which is substantially better. Figure 6 compares average observations with predictions of 16

305    sampling locations in reference to EPA guidelines (based on the 1986 RWQC) (U.S. EPA

306    2012b) that say geometric mean (GM) coliform density and statistical threshold value (STV) of

307    indicator organisms for waters designated for primary contact recreation should be less than 126

308    CFU/100 mL and 410 CFU/100 mL, respectively. The figure showed that both average

309    predictions and observations exceeded EPA's GM criteria. About 18% observations exceeded

310    EPA's STV criteria, while all predictions were lower than the STV value indicating model's

311    under predictions for few locations. Nevertheless, 81% of both predictions and observations were

312    lower than the STV value indicating the model's suitability for assessing instream water quality.

313        In addition, the neighborhood statistics model proposed here does not requires intensive

314    calibration, which is necessary in hydrological models while implementing for predicting in-

315    stream *E. coli* levels at watershed scale. Even though it is not expected to fit the observations

316    with predictions very well, while predicting in-stream bacteria levels (Dorner et al. 2006),

317    advancing existing modelling approaches are necessary in order to derive/identify efficient

318    watershed management plans for improving stream water quality. The approach we presented

319    here requires further improvement, and we anticipate that using a larger observed dataset will

320    potentially enhance the predictions. One major challenge in stream bacteria modeling is the

321    availability of limited observed data. Therefore, future studies carrying out extensive monitoring

322    as well as modeling based on the field observations will certainly improve the existing models.

323    **4. CONCLUSIONS**

324    To predict in-stream *E. coli* levels, we have developed a neighborhood statistics model, Markov

325    Random Field model, which was implemented in the Squaw Creek watershed, Iowa. The model

326    predictions were compared with the observed *E. coli* levels at 16 different locations. The two

327    independent parameters water temperature ($^{o}$C) and stream water depth (m) were used for

328    predicting the *E. coli* levels. Results indicated that the method used here is a potentially useful

329    approach to predict instream *E. coli* levels at watershed scale with certain degree of

330    predictability. The approach can be useful in understanding of the spatial variability of *E. coli*

331    levels at watershed scale.

336 **References**

337 American Public Health Association (APHA). (1999). Standard methods for the examination of

338       water and wastewater, AWWA, Water Environment Federation.

339 Bai, S., & Lung, W. S. (2005). Modeling sediment impact on the transport of fecal bacteria.

340       *Water Research, 39 (20),* 5232–5240.

341 Besag, J. (1974). Spatial integration and the statistical analysis of lattice systems. *Journal of*

342       *Royal Statistical Society, 36 (2)*, 192–236.

343 Cho, K. H., Pachepsky, Y. A. Kim, J. H., Guber, A. K., Shelton, D. R. & Rowland, R. (2010).

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

345       and reach-specific modeling. *Journal of Hydrology, 391*, 322–332.

346 Cressie, N. (1993). Statistics for spatial data. John Wiley & Sons, New York.

347 Dorner, S.M., Anderson, W. B., Slawson, R. M., Kouwen, N., & Huck, P. M. (2006). Hydrologic

348       modeling of pathogen fate and transport. *Environmental Science & Technology 40 (15),*4746–

349       4753.

350 Droppo, I. G., Liss, S. N., Williams, D., Nelson, T., Jaskot, C., & Trapp, B. (2009). Dynamic

351       existence of waterborne pathogens within river sediment compartments: Implications for

352       water quality regulatory affairs. *Environmental Science & Technology, 43 (6)*, 1737–1743.

353

354     Hipsey, M. R., Antenucci, J. P., & Brookes, J. D. ( 2008). A generic, process-based model of

355           microbial pollution in aquatic systems. *Water Resources Research, 44 (7),* W07408.

356     Iowa Environment Mesonet (IEM). (2012). Iowa AG Climate Network, Iowa State University.

357           http://mesonet.agron.iastate.edu/agclimate/hist/hourlyRequest.php.

358     Jamieson, R.C., Joy, D. M., Lee, H., Kostaschuk, R., & Gordon, R. J. ( 2005). Resuspension of

359           sediment-associated *Escherichia coli* in a natural stream. *Journal of Environmental*

360           *Quality, 34 (2)*, 581–589.

361     Kaiser, M.S. (2010). Statistical methods for spatial data. Course lecture, Department of Statistics,

362           Iowa  State University.

363     Kaiser, M. S., & Nordman, D. J. (2012). Blockwise empirical likelihood for spatial Markov

364           Model assessment. *Statistics and Its Interface, 0,*1–8.

365     Kim, J.W., Pachepsky, Y.A., Shelton, D. R., & Coppock, C. (2010). Effect of streambed bacteria

366           release on *E. coli* concentrations: Monitoring and modeling with the modified SWAT.

367           *Ecological Modelling, 221 (12)*,1592–1604.

368     Muirhead, R.W., Davies-Colley, R. J., Donnison, A. M. & Nagels, J. W. (2004). Fecal bacteria

369           yields in artificial flood events: quantifying in-stream stores. *Water Research,* 38,1215–

370           1224.

371     Nagels, J.W., Davies-Colley, R. J., Donnison, A. M., & Muirhead, R.W. (2002). Faecal

372           contamination over flood events in a pastoral agricultural stream in New Zealand. *Water*

373           *Science and Technology, 45 (12)*, 45–52.

374     Pachepsky, Y.A., & Shelton, D. R.  (2011). *Escherichia coli* and fecal coliforms in freshwater

375           and estuariene sedients. *Critical reviews in Environmental Science and technology,*

376           *41(12)*,1067–111.

377  Pandey, P. K., Soupir, M. L., & Rehmann, C. R. (2012a). Predicting resuspension of *Escherichia*

378      *coli* from streambed sediments. *Water Research, 46*, 115–126.

379  Pandey, P. K., Soupir, M. L., Haddad, M., & Rothwell, J. J. (2012b). Assessing the impacts of

380      watershed indexes and precipitation on spatial in-stream *E. coli* concentrations.

381      *Ecological Indicator, 23*, 641–652.

382  Parajuli, P. B., Douglas-Mankin, K. R., Barnes, P. L., & Rossi, C.G. (2009). Fecal bacteria

383      source characterization and sensitivity analysis of SWAT 2005. *Transaction of ASABE,*

384      *52 (6)*,1847–1858.

385  Rehmann, C.R., & Soupir, M.L. 2009. Importance of interactions between the water column and

386      the sediment for microbial concentrations in streams. *Water Research,  43(18)*, 4579–

387      4589.

388  Rothwell, J.J., Dise, N.B., Taylor, K.G., Allott, T.E.H., Scholefield, P., Davies, H., & Neal, C.

389      (2010a). A spatial and seasonal assessment of river water chemistry across North West

390      England. *Science of Total Environment, 408*, 841–855.

391  Rothwell, J.J., Dise, N.B., Taylor, K.G., Allott, T.E.H., Scholefield, P., Davies, H., & Neal, C.

392      (2010b). Predicting river water quality across North West England using catchment

393      characteristics. *Journal of Hydrology, 395(3-4)*, 153–162.

394  Schilling, K. E., Zhang, Y., Hill, D. R., Jones, C. S., & Wolter, C. F. (2009). Temporal variations

395      of *E. coli* concentrations in a large Midwestern river.  *Journal of hydrology, 365*,79–85.

396  U.S. Environmental Protection Agency (USEPA). (2012a). WATERS (Watershed Assessment,

397      Tracking & Environmental ResultS). Washington, D.C. (accessed on 10/22/2012).

398    U.S. Environmental Protection Agency (U.S. EPA). (2012b). Recreational water quality criteria.

399    http://water.epa.gov/scitech/swguidance/standards/criteria/health/recreation/index.cfm

400    (accessed on 2/22/2013).

401    U.S. Environmental Protection Agency (USEPA). (2013). WATERS (Watershed Assessment,

402    Tracking & Environmental ResultS). Washington, D.C. (accessed on 2/22/2013).

403    Wilkes, G., Edge, T. A., Gannon, V. P. J., Jokinen, C., Lyautey, E., Neumann, N. F., Ruecker,

404    N., Scott, A., Sunohara, M., Topp, E., & Lapen, D. R. (2011). Associations among

405    pathogenic bateria, parasites, and environmental and land use factors in multiple mixed-

406    use watersheds. *Water Research, 45(18)*, 5807–5825.

407    World Health Organization. (2010). Water Sanitation and Health.

408    http://www.who.int/water_sanitation_health/diseases/en/

409
410

411

412

413

414

415

416

417

418

419

420

421

422

423 **Table 1** Neighborhood structures of main channel

| Sampling locations | Neighbors |
|---|---|
| 1 | 10 |
| 2 | 10 |
| 3 | 11 |
| 4 | 12 |
| 5 | 12 |
| 6 | 13 |
| 7 | 13 |
| 8 | 16 |
| 9 | 16 |
| 10 | 1, 2,11 |
| 11 | 10, 3,12 |
| 12 | 11, 5, 4,13 |
| 13 | 12, 7, 6,14 |
| 14 | 13,15 |
| 15 | 14,16 |
| 16 | 15, 8, 9 |

424

425

426 **Table 2** Parameter values of neighborhood structures

| | $\tau^2$ | $\eta$ | $\beta_0$ | $\beta_1$ | $\beta_2$ |
|---|---|---|---|---|---|
| Estimate | 4.0E+04 | -0.02 | -18.2 | 15.9 | 109.9 |
| Lower limit | 2.6E+04 | -0.27 | -219 | 5.3 | -153 |
| Upper limit | 5.4E+04 | 0.23 | 183.5 | 26.7 | 373 |
| p-value | 7.8 E-09 | 0.56 | 0.42 | 0.001 | 0.20 |

427

428
429
430
431
432
433
434

435 **Table 3** Correlation coefficients of *E. coli* levels at different locations
436

| Locations | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | | 0.97 | | 0.98 | 0.93 | | 0.83 | 0.92 | 0.82 | 0.96 | 1.0 | 0.75 | 0.67 | 0.78 | | |
| **2** | | | | 0.99 | 0.98 | | 0.76 | 0.98 | 0.75 | 0.96 | 0.97 | 0.70 | 0.68 | 0.66 | | |
| **3** | | | | | | | | | | | | | | | | |
| **4** | | | | | 0.99 | | 0.71 | 0.97 | 0.70 | 0.93 | 0.98 | | | | | |
| **5** | | | | | | | | 0.99 | | 0.90 | 0.93 | | | | | |
| **6** | | | | | | | 0.94 | | 0.95 | 0.70 | | 0.95 | 0.85 | 0.96 | 0.94 | 0.87 |
| **7** | | | | | | | | | 1.0 | 0.9 | 0.82 | 0.99 | 0.91 | 0.97 | 0.86 | 0.86 |
| **8** | | | | | | | | | | 0.91 | 0.91 | | | | | |
| **9** | | | | | | | | | | 0.89 | 0.82 | 0.99 | 0.90 | 0.98 | 0.87 | 0.86 |
| **10** | | | | | | | | | | | 0.96 | 0.87 | 0.84 | 0.81 | | 0.70 |
| **11** | | | | | | | | | | | 0.75 | 0.66 | 0.78 | | | |
| **12** | | | | | | | | | | | | | 0.95 | 0.95 | 0.92 | 0.92 |
| **13** | | | | | | | | | | | | | | 0.81 | 0.91 | 0.97 |
| **14** | | | | | | | | | | | | | | | 0.84 | 0.78 |
| **15** | | | | | | | | | | | | | | | | 0.97 |
| **16** | | | | | | | | | | | | | | | | |

437
438 NOTE: Only statistically significant numbers are shown in the table ($p < 0.05$).
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468

**Figure captions:**

469

470

471 **Fig. 1** Study Area (Squaw Creek Watershed, Iowa, U.S.A). Corn and soybean crops dominate

472 the watershed

473 **Fig. 2** Maximum likelihood of η

474 **Fig. 3** Spatial observations of *E. coli* levels and climate of watershed. Top four figures shows *E.*

475 *coli* levels along the stream, and bottom figure shows temperature, precipitation, and stream flow

476 (stream flow was observed at the lowest end of the watershed i.e., location 16 of Figure 1)

477 **Fig. 4** Observed stream water depth and water temperature at 16 locations of the watershed

478 **Fig. 5** Comparison between observations and predictions of in-stream *E. coli* levels in Squaw

479 Creek watershed

480 **Fig. 6** Comparison between observed and predicted *E. coli* levels, and EPA's Geometric Mean

481 (GM) and Statistical Threshold Value (STV) criteria of indicator organisms for fresh water.

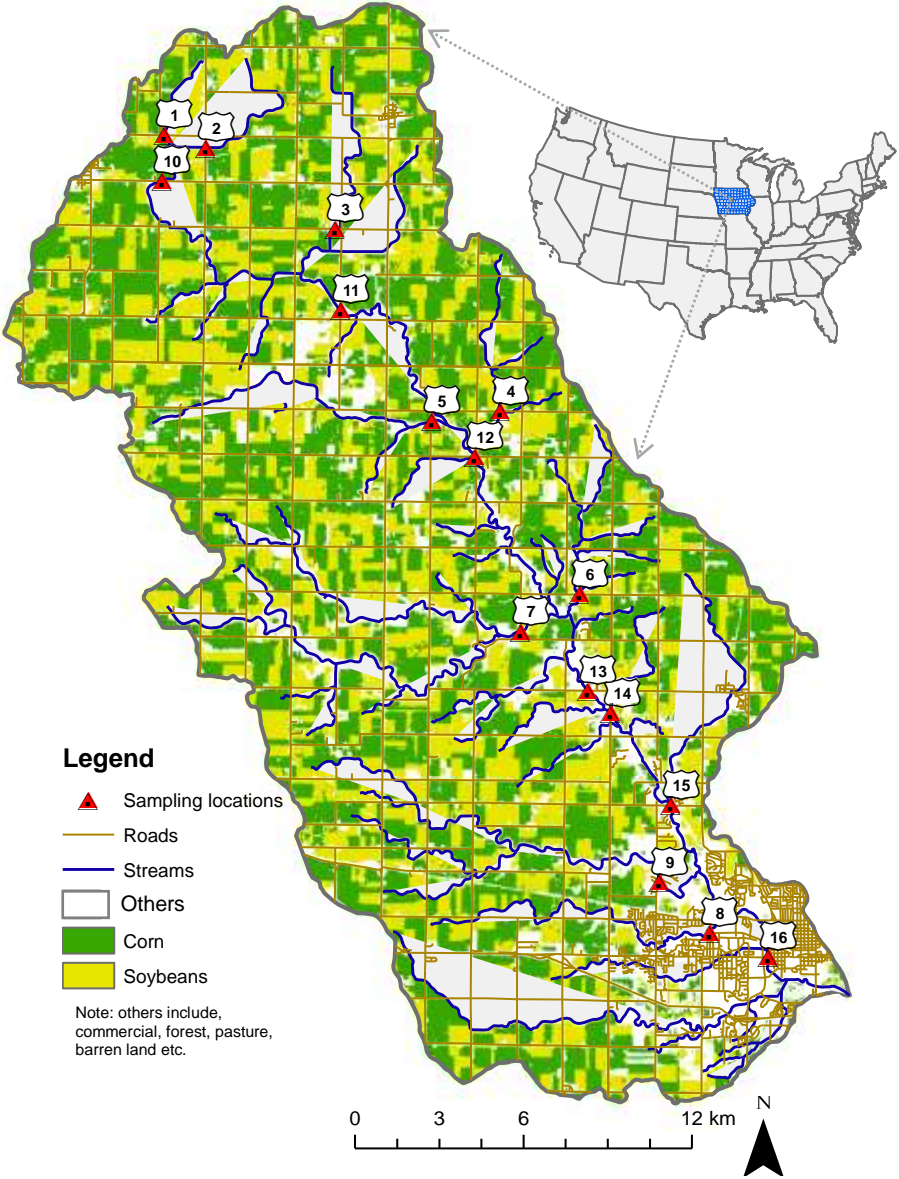482 Average of observed values and predicted values of four sampling events are shown in the figure

483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502

503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548

**Figure** 1.



Legend

- ▲ Sampling locations
- —— Roads
- —— Streams
- ☐ Others
- ■ Corn
- ■ Soybeans

Note: others include, commercial, forest, pasture, barren land etc.

549 **Figure 2**
550
551
552



553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572

573     **Figure 3**



27 June, 2009 (t = 1); 7 July, 2009 (t = 2); 17 July, 2009 (t = 3); 17 October, 2009 (t = 4)

574
575
576
577
578
579

580
581 **Figure 4**
582
583



Legend:
- Depth (June 27, 2009)
- Depth (July 6, 2009)
- Depth (July 17, 2009)
- Depth (October 17, 2009)
- Temp. (June 27, 2009)
- Temp. (July 6, 2009)
- Temp. (July 17, 2009)
- Temp. (October 17, 2009)

Y-axis (left): Stream water depth (m)
Y-axis (right): Stream water temperature (deg C)
X-axis: Sampling locations

584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605

606     **Figure 5**
607



608
609
610
611
612
613
614
615
616
617
618

619
620 **Figure 6.**
621



622
623
624
625
626
627
628
629
630
631
632
633
634
635