

Reviewer Comment

Author Response

Review of:

Spatially Referenced Bayesian State-Space Model of Total Phosphorus in western Lake Erie

July 16, 2021

Referee: Ken Newman, Biomathematics & Statistics Scotland, and School of Mathematics, University of Edinburgh

1 Overall comments

Developing a quantitative, explanatory model that links total phosphorus (TP) loads in rivers that discharge into Lake Erie to the TP loads in the lake is useful for both increasing scientific understanding of the processes and for assessing potential management actions. The compilation and "wrangling" of both the lake and river TP measurements and the surface current data was no small task. The state-space model (SSM) framework seems quite appropriate given the time series nature of the data, and the spatial structure of the process model for these data is a crucial feature. The subsequent use of the fitted SSM to conduct "what if" exercises on changes in riverine TP loads shows the work's potential utility as a decision support tool. I have major concerns about the formulation and fitting of the SSM and the explanation of the results.

R2C1 Adjacency matrix: instead of calculating mean daily eastward and northward water velocity (m/sec) for a given node, using the hourly measurements at that node, why not follow the hour by hour trajectory of points across a day? On a per cell basis at hour 1 start with a point at the cell center (the node), say (Lonc;1; Latc;1) and use eq'ns 3 and 4 to advance that point to its longitude and latitude at hour 2, (Lonc;2; Latc;2). Then for whatever cell that point is in, apply eqn's 3 and 4 again to move it to (Lonc;3; Latc;3), and so on until reaching (Lonc;24; Latc;24).

We changed our approach and remade the adjacency matrix as the reviewer suggested. The difference was very slight, e.g., in 2018 of the ~35,000 adjacencies only 5 were different using the suggested methodology. This suggests that the computationally intense hourly approach does not finely resolve information on water movement, at least in this location. The "hour-step" method, while not a deterministic particle-tracking model, did not improve performance in our system potentially because of our dominated west to east movement patterns. We used the reviewer's suggested method regardless because while our surface current data was insensitive to which method was used, differences between the two methods may be more extreme in other environments and thus the reviewers "hour-step" method would be more widely applicable.

We propose this text added to **Section 2.2** amending the adjacency matrix description.

"Hourly northward and eastward velocity (m day^{-1}) for each node for years 2008 to 2018 defined surface current direction in radians (dLat and dLon) using the node latitude (Lat0) and longitude (Lon0), the Earth's radius (R, 6378137 m), the northward velocity offset in meters (dN), and eastward offset in meters (dE) (Eqs 1 and 2). The direction the surface water travelled in radians was used to determine the latitude (Lat1) and longitude (Lon1) which represented by each hourly movement (Eqs 3 and 4), and was repeated for

24-hours until the final position of the surface water movement from each node was determined.”

R2C2 A truncated normal distribution for the process model does not make sense. The state component is the “true” log TP concentration and that is not a truncated value. A truncated normal could be used for an observation model, however.

We argue that using the truncation on the “true” log TP value is appropriate. This truncation forces the model to stay within concentration bounds that must exist in the lake. At no time step or location in the lake will “true” concentrations exceed those observed at the maximum Maumee River concentration. Additionally, the “true” value lower bound within our lake is not going to be zero or negative and a boundary condition at some logical value is needed. Additionally, most of the water entering western Lake Erie is from the Detroit River and its concentration typically ranges between 10-25 $\mu\text{g l}^{-1}$, amongst our 11-years of observations 0.7% of samples were under 10 $\mu\text{g l}^{-1}$. Thus, our 5 $\mu\text{g l}^{-1}$ lower bound is reasonable. Allowing the model to predict concentrations higher or lower than what we know to exist would decrease the applicability of SSMs in water quality modeling.

R2C3 Given that year-specific SSMs were fit, it would be useful to compare the posterior distributions for the four process model parameters, β_{mau} , β_{rai} , β_{self} , and β_{lake} . At a minimum there needs to be summaries about the posteriors of these parameters, and some discussion and interpretation of the values are needed.

We will add both a visual representation of the β s, the uncertainty of the data model, the uncertainty of the process model, and a discussion of their values. The visual will be added as **Appendix C** and the discussion will be added to **Section 4.1**.

We propose the following text:

“TP is a conservative water quality constituent. TP observations are insensitive to biogeochemical transformations of phosphorus form because these data represent both the organic and inorganic forms of phosphorus occurring in the water column. β s near 1 would then be expected in the absence of dilution. Dilution of TP would happen west to east across our spatial model window, however the depth gradient within western Lake Erie is muted. β s larger than 1 would indicate in-lake sources of TP. Every β_{Mau} , β_{Ras} , β_{Lake} , and β_{Self} fit in our models had 95% predictive intervals encompassing a value of 1. No identifiability issues identified by priors dominating the fit of coefficients were observed (Auger-Méthé et al., 2020; Appendix C).

Fit process model or data model uncertainty were well identified (Appendix C). The apportionment of uncertainty between the process model and the data model varied from year to year. This was driven in annual variation in the data model uncertainty. We propose these annual differences were due to the combination of the number of samples collected and they relative position to the surface currents. However, the uncertainty within our models did not prevent accurate outputs estimating TP concentrations at observed and unobserved nodes.”

Appendix C

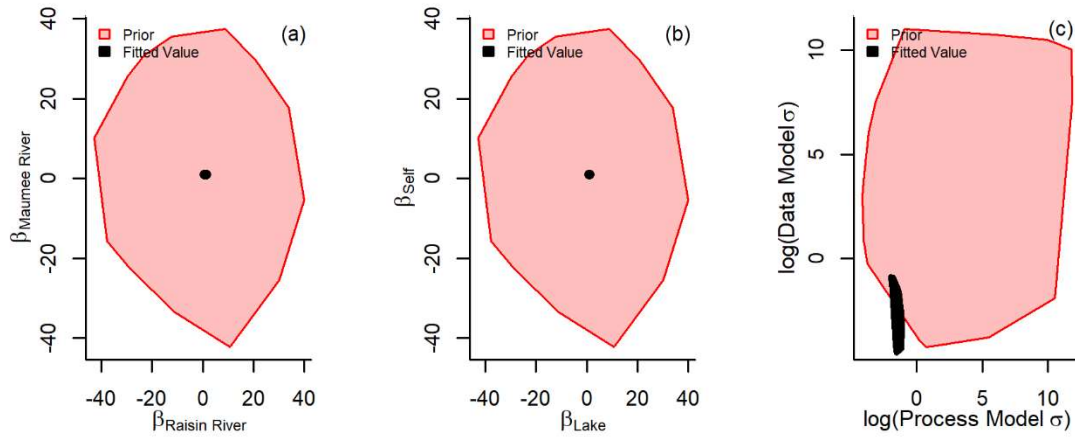


Figure C1. Uninformed priors were used to fit the state space coefficients of the Maumee River (β_{Mau} , a. all values represented as red polygon), River Raisin (β_{Ras}), the western Lake Erie nodes subject to movement (β_{Lake} , b. red polygon), and those Lake Erie nodes which did not encounter sufficient water movement to associate to an “upstream” node (β_{Self}). The fitted values for every year (a. and b., all values represented as black polygon) did not significantly overlap. The log data model and process model uncertainty (c) for every year also were well identified (the uncertainty σ was logged to aid in visually comparing prior and fitted values).

Table C1. State space models in western Lake Eire fit coefficients predicting TP concentrations from the previous time-step within the Maumee River, the River Raisin, Lake Erie, and lake locations where the current time-step is informed by the same location, β_{Maumee} , β_{Raisin} , β_{Lake} , and β_{Self} , respectively. The 95% predictive interval (PI) for each year and coefficient was examined. The same models fit annual process model and data model precision.

Year	β_{Maumee} 95% PI		β_{Raisin} 95% PI		β_{Lake} 95% PI		β_{Self} 95% PI		Process σ 95% PI		Data σ 95% PI	
2008	0.989	1.008	0.985	1.013	0.986	0.999	0.963	1.012	0.263	0.331	0.02	0.098
2009	0.993	1.005	0.992	1.007	0.991	1	0.984	1.002	0.168	0.193	0.125	0.165
2010	0.994	1.003	0.993	1.006	0.988	0.994	0.992	1.001	0.13	0.183	0.326	0.404
2011	0.993	1.006	0.99	1.008	0.991	1.006	0.985	1.007	0.174	0.209	0.017	0.061
2012	0.993	1.007	0.804	1.186	0.997	1.008	0.983	1.005	0.177	0.218	0.017	0.062
2013	0.993	1.007	0.992	1.008	0.992	1.003	0.978	1.008	0.181	0.277	0.135	0.25
2014	0.992	1.005	0.99	1.006	0.991	1.001	0.991	1.018	0.173	0.235	0.145	0.232
2015	0.992	1.006	0.989	1.007	0.992	1.003	0.982	1.003	0.212	0.264	0.156	0.238
2016	0.993	1.005	0.991	1.005	0.996	1.004	0.983	1.001	0.161	0.2	0.146	0.214
2017	0.993	1.006	0.991	1.008	0.995	1.003	0.978	0.998	0.197	0.24	0.019	0.079
2018	0.992	1.007	0.99	1.008	0.983	0.991	0.966	0.997	0.221	0.26	0.018	0.065

R2C4 Related to item 3, a more comprehensive approach to fitting these data would be a hierarchical SSM. Model the slope coefficients (the β 's) as random variables coming from a generating distribution that reflects between year (or environmental) variation; eg., $\beta_{\text{mau};y} \sim \text{Normal}(\mu_{\beta_{\text{mau}}}; \sigma_{\beta_{\text{mau}}})$.

We agree that our SSM could be represented by a hierarchical form where each year is fitted together. The limitation is computation. As the reviewer notes in comment **R2C7**, the run time for each model year is very long (+/- 24 hours), so fitting all the years together would at a minimum require 18-days to run, and the outputs from such a hierarchical model would be matrices so big that exponentiating the outputs would surpass the memory available to us on our computer cluster. We will add text to **Section 4.1** suggesting a hierarchical SSM for all years as an approach readers with smaller datasets could employ.

We propose the following text:

“This framework could be implemented with the coefficients (β_{Mau} , β_{Ras} , β_{Lake} , and β_{Self}) fit hierarchically by year, current restrictions on computer memory prevented that use here. However for smaller spatial and temporal models it could be effective.”

R2C5 Nothing has been said about the estimates of the process model precision (Q) and the obs'n model precision (R). How do they vary between years? There are often weak identifiability problems with these kinds of linear Gaussian SSMs (see Auger-Methe et al. (2016)).

We propose added a figure as **Appendix C** with the estimates of Q and R, visually presented in our response to **R2C3**. Non-identifiability represented by overlap of the fitted values and the prior is absent across all models. The uncertainty attribution between process model and observation model differed from year to year, some years uncertainty is dominated by process error while other years it is equally attributed to process and observation. Process model precision is consistent while observation error changes year-to-year. These year-to-year variations in the observation model are potentially due to the number of samples collected or the location of the samples relative to the surface currents. We propose new text for **Section 4.1** in the above response to **R2C3**.

R2C6 Regarding application of the SSM to assess the effect of reducing TP loads in the Maumee River, assuming year specific SSMs were fit, then the particular SSM that was used needs to be stated. Also related to 3, it would be good to report results for each of the year specific SSMs to show (some of) the uncertainty in the assessments.

The SSMs per year were run twice, first with all the unaltered data, and second with Maumee River data halved. The values for process and observation model precision, as well as the value of every year's β s, will be added in an Appendix C as a table following the visual comparing the priors to the fitted values.

Proposed text in **Section 2.3**:

The Maumee River impact plume was estimated by artificially reducing the Maumee River TP concentrations by 50% ($\hat{y}_{\text{Maumee},t,y}$), each years model was then again (Eq 5-7).

Proposed text in **Section 4.1**:

The coefficients (β s) and uncertainties (process and data model σ) varied only slightly from year to year (Table C1). The 2012 Raisin River coefficient (β_{Ras}) predictive interval was larger than other years because of a lack of data in that year. The proportion of uncertainty between process and data model also varied only slightly (Table C1), possibly because of the number of or spatial position of observations.

R2C7 While it is good to have included the R code and data for fitting the SSM to the 2018 data, something needs to be said about computational time. I ran the code, but it had not completed after 14 hours, and I could not therefore verify any results.

We noted this issue in our response to **R2C4**, the model takes our system +/- 24-hours to run. We will note that within the publicly available code.

Detailed technical comments

p=page, L=line.

R2C8 Introduction, p3, L70: Why test a hypothesis of linearity? A more general aim would be to quantify the nature of a relationship, be it linear or nonlinear.

This hypothesis is based on a conservative water quality constituent (TP) and a single source (the Maumee River). If flow is constant, the more TP imported from the Maumee River the higher concentrations should be in the lake. While this is simplistic, it is widely examined by previous research, and supporting these earlier works we found a linear relationship.

R2C9 Methods (a) p3, L73-79. Is there overlap in the sampling on Lake Erie? For example, do two agencies collect data on the same cells?

There are instances where the same agency collects multiple samples on a single day at a single node, however no instances where differing agencies collect samples on the same day and node.

R2C10(b) p3, L81. By station is that just referring to the rivers, or does that include the lake? When multiple samples were collected from a station on a single day, how much variation was there in the measure value? Such information could be used in the observation model of the SSM.

“Station” in that sentence should be replaced with “node”, the multiple sampling days are rare, and the values are always very close to each other.

R2C11(c) pp 3-4, L84-100. It would help to create an example figure showing the geometry (or trigonometry) underlying the calculation of a particle's change in position at time t to time $t + 1$. For example, draw an x-y plot with points p_t and p_{t+1} with coordinates x =longitude and y =latitude. Draw a right triangle, where the vertical and horizontal sides are parallel to the y-axis and x-axis respectively, and the hypotenuse connects p_t and p_{t+1} . So the vertical side indicates the “northerly” movement, the horizontal side indicates “easterly” movement. (For what it's worth, I find thinking of this as a step direction and step size process as is done in animal movement modeling.)

A figure of an example track of surface current will be added as **Appendix D**.

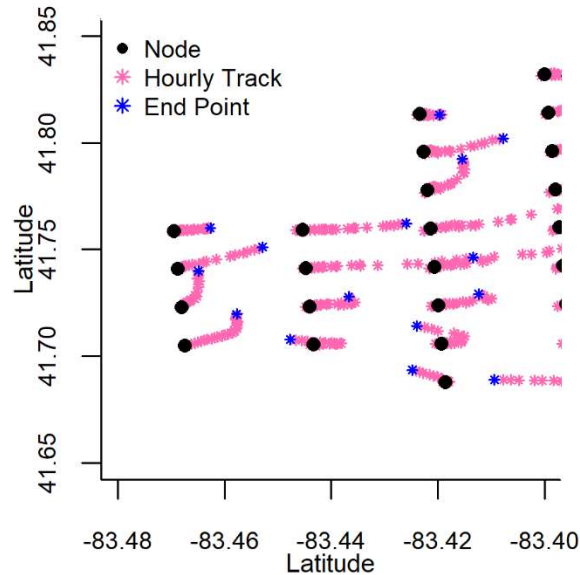


Figure D1. Surface current data was available hourly within western Lake Erie. These 24-hour data were used to track the daily movement of water from each node.

R2C12(d) p3, L89. A refinement would be fit a spatially smooth velocity map rather than use values at “nodes”, though I’m not sure how much practical effect that would have.

While this would be an improvement, our goal is to present a model that can be constructed easily by others by our example code.

R2C13(e) p4, L101. Does anything need to be handled differently at cells on the perimeter? Such cells could be transferring TP to cells outside the spatial window, presumably.

The nodes along the eastern perimeter do have the additional uncertainty of missing surface water adjacencies due to their location. However, with the sources of TP on the western side of the spatial window there is little practical effect in our system. We will add text in Section 4.2 to bring this issue to the attention of readers if they have systems where perimeter issues are a larger problem.

We propose the following text:

“Nodes along the eastern perimeter of the spatial window of our model have additional uncertainty inherent in their position. Occasionally, they will not be associated with the proper “down gradient” node because the extent removes those nodes. Within our system there is little practical effect as these nodes are far from the Maumee River and are dominated by low concentrations. This is a potential problem in other systems and may necessitate wider spatial windows to eliminate.”

R2C14(f) p4, L109. References to SSM literature are strongly recommended; e.g., Durbin and Koopman (2012); Shumway and Stoffer (2019). Discussion of the distinction between inference about latent states and fixed parameters would be good.

The authors will insert the references suggested at their appropriate location as suggested by the reviewer. Rather than distinguishing between inference and fixed parameters here we want the reader to stay focused on our goal of defining daily TP concentrations at each node on each day irrespective of the availability of observations.

R2C15(g) p4, L119. Say something about the β 's and what they mean. For example, regarding β_{self} and β_{lake} , are they likely less than 1? Is there some loss of TP between a source point k and an end point n from one day to the next? Is any sort of seasonality expected in the river values? This would make the assumption of constant β_{mau} and β_{rai} suspect.

We agree with the reviewer that this should be discussed further, please see our response to **R2C3**.

R2C16(h) p5 L126-128: The joint prior described for β_{mau} and β_{rai} does not exactly match what the R code indicates where the 2nd parameter in the normal is the standard deviation.

The priors described in the text will be corrected as indicated by the reviewer.

R2C17(i) p5, L137-138: Why is cross-validation needed to make comparisons of goodness of fit across years? Doesn't the Bayesian R2 do that? Cross-validation is more often used for model selection, which is the focus of both the referenced Vehtari et al and the Piironen and Vehtari papers.

The cross-validation via leave one-node-out compared how well the model predicted the values not available during model fitting. Estimates of R^2 report how well the model does while given all the available data. Cross-validation's strength is that it shows how the model output accurately gives estimates at unobserved locations, which dominate the dataset. Additionally, the cross-validation by node shows that the model does well estimating TP values irrespective of the nodes proximity to the TP river sources and that the model does well estimating TP values every year irrespective of the number or location of in-lake observations.

R2C18(j) p6,L155-164: I cannot tell what is being done here. Is this necessary?

These steps describe the performance of the cross-validation as a fitted distribution. We did this so that predictive intervals about the difference between cross-validated nodes and years could be used.

R2C19(k) p6, L171-174: Write down an equation for deflection, $d_{n,y}$, and for the normalized estimate.

The authors will change the manuscript as suggested by the reviewer and insert the equation.

2

R2C20(l) p6, L179-185: Are 252 regressions being fit each year? Write down what effective load, $\sim \ln; y$ means- is it an average? Why is there a subscript n if it is measured on the Maumee? I'd be curious about identifiability issues/posterior correlations with the parameters in eq'n 11,

too| seems overly complicated. Would it make sense to take average \raw" y and take the log of that?

One regression is fit through the data. The effective load of the Maumee River (l) has a subscripts n and y because each node (n) will have its own effective load each year (y). A node close to the river would have a high effective load compared to a node far away. Similarly, a node may have a high effective load in 2008 but a low effective load in 2009 depending on direction and strength of the surface currents.

R2C21(m) R code: The priors for precision to the Lake, Maumee River, and River Raisin are calculated as the inverse of standard deviation; shouldn't that be the variance? Also using the data to set priors is questionable.

We have changed the way priors are set for the initial states of the Lake, Maumee, and Raisin River nodes. The first year we have available (year = 2008) is has a prior for all three node types given from values reported in Rockwell et al (2005). For all the other years the priors are defined as the mean and precision of the first 20-days of the previous year.

Rockwell, David C., et al. "The US EPA Lake Erie indicators monitoring program 1983–2002: trends in phosphorus, silica, and chlorophyll a in the central basin." *Journal of Great Lakes Research* 31 (2005): 23-34.

We propose the following text:

“Initial conditions for the latent state $x_{n,t=1,y}$ were defined as the mean and variance of the previous year first 20 days. The first year (year = 2008) initial conditions were estimated as $N(12, 5)$ (Rockwell et al., 2005).”

R2C22 (n) Data for R code: why does the adjacency object use have 4 values for the \next" cell, at time $t + 1$? Only the first is used in the code.

These columns are hold overs from our original data curation, the new adjacency matrix defined in response to **R2C1** has only one column value.

R2C23 (o) Sensitivity analysis for the priors needs to be conducted.

The sensitivity of fitted values to the priors is now incorporated in a figure generated in response to **R2C3**.

R2C243. Results (a) p7, around L185. Say something about the ranges and averages of northerly and easterly velocities (m/sec), of TP concentrations (distinguishing between Maumee River, River Raisin, and western Lake Erie), and of calculated distances moved in a single day.

In our response to **R2C11** we display an example of the water movement, we believe that further tables of northing and easting data were not aid readers in understanding our SSM approach. We will add text to report the range of TP values within each river and western Lake Erie.

Section 3 Proposed text:

The mean values of the observed TP concentrations within the Maumee River, River Raisin, and western Lake Erie were 170 ug l-1 (95% interval, 3.5 to 438 ug l-1), 80 ug l-1 (95%, 40 to 215 ug l-1), and 38 ug l-1 (95%, 10 to 203 ug l-1), respectively.

R2C25(b) p7, L195-199. In addition to the summaries in Table 1, show some plots of posterior mean values for some of the cells across the 136 day period against corresponding observations. Show an example histogram (or two) of the distribution of predicted log concentrations with the observation. Presumably these R2 are only calculated on the cells with observations. I don't understand what the cross-validated measure is doing (relates to earlier comment in methods).

Posterior mean values for a select number of nodes each year are represented in **Appendix B**. Please refer to our previous comments regarding the cross-validation in **R2C17**.

R2C26(c) p7, L200-205. As mentioned above need to report out results on parameter estimates for the β (each of the 4 for all 11 years) and report out the process and observation model standard deviations ($1=pQ$ and $1=pR$, based on R code).

We will generate figures depicting the values of the four β s and Q & R, for all years in an **Appendix C** as noted in our response to **R2C3**.

Minor editorial remarks

R2C271. Throughout, consider using the word "cell" instead of "node" as a node is usually interpreted as a point. Node could then refer to cell center.

We prefer the term "node" but will consider using "cell" during the further editing steps in submitting the manuscript.

R2C282. Section 2 Methods.

(a) pp 3-4. Consider creating a new subsection for the material in the first paragraph (L84-107) of the Model Description subsection, maybe naming it Construction of an Adjacency matrix.

We prefer the current organization as it aids the organization of results and discussion.

R2C29(b) p4, L109. Emphasize that 11 different SSMs will be fit.

Line 109 states "...state-space models for each year...", we will emphasize the separate years in our new table (Table C1) of Q, R, and β s.

R2C30(c) p4, L119. The observations are modeled with a normal distribution; they are not estimated with a normal. As mentioned previously, a truncated normal (perhaps just on the left) could be used to account for measurement limitations.

Please see our response to **R2C2**.

R2C31(d) p5, L139. Consider renaming Section 2.2.2 "Fitting the SSM" as "SSM Fit" could be interpreted as a result not a method.

The authors will change the manuscript as suggested by the reviewer.

R2C32 (e) p5,L133: Perhaps move text beginning with "The model was run.." into Section 2.2.2. Could delete/move material in the sentence beginning "The efficacy of.." as it is redundant with material on L140-151.

We prefer the current organization.

R2C33 (f) p5, L136. "efficacy" seems an odd choice, why not goodness of fit?

The authors will change the manuscript as suggested by the reviewer.

R2C34 (g) p5,L137: The Vehtari, et al., 2017 paper does not refer to Bayesian R2. Give a mathematical definition of the R2 here: not clear to me what resolved and residual variances mean.

R^2 is defined on line 142. The authors will change the reference to (Gelman et al., 2019). "Resolved" is a typo, it should read "fitted".

R2C35 (h) p5,L138 What does "utility" mean? And what value of K was used?

Please refer to our response to **R2C17**.

R2C36 (i) p5,L145: Cross-validation is meant for CV not coefficient of variation?

The authors will change the manuscript as suggested by the reviewer.

R2C37(j) p5, L154: What does preferentially mean? 3

Please refer to our response to **R2C17**.

R2C38 (k) p6, L165: Instead of Model Experimentation, wouldn't Model Usage or Application make more sense here (and elsewhere)?

We prefer the current organization.

R2C39 (l) 6, L166: As said previously, the SSM for which year was used?

Please refer to our response to **R2C20**.

3. Section 3 Results.

R2C40 (a) p7, L187-193. Re: the degree of missing data, I think it would be easier to follow by first saying how many space-time cells there are and then give the number with data: "For the Lake, there $252 \times 11 \times 136 = 376,992$ cells of which 1218 had data, and for the two rivers, there were $2 \times 11 \times 136 = 2992$ cells, of which 2258 had observations". Save the discussion of the inference for cells without data till later.

We prefer the current organization.

R2C41(b) p7, around L188. It would help to see an example plot that shows the spatial dist'n of cells, in a given year, that had at least one observation (see Figure 1).

Appendix B visually shows all the nodes which contain data through the years.

R2C42(c) Also a plot showing "source" cells and "end" cells would be good to show the adjacency. (Note: in the R code, use object has 4 values in the second dimension, but it appears that only the 1st value is used as the adjacency matrix|what are the other 3 values for?)

Please refer to our response to **R2C22**.

R2C43(d) p7, around L188. Add a plot or two of the log(TP) concentrations. See Figure 2.

We prefer the current organization.

R2C44 (e) p7, L195: \efficiency" here, but \efficacy" in Methods, but would model quality or goodness-of-fit be more appropriate?

The authors will change the manuscript as suggested by the reviewer.

R2C45 (f) p7, L197: typo: Table 1 not Table 2.

The authors will change the manuscript as suggested by the reviewer.

R2C46 (g) p7, L202: Say that Figure 2 shows 2018 and Figure B.1 in Appendix B shows 2006. The captions in those figures need to indicate that the black dots are observed values.

The authors will change the manuscript as suggested by the reviewer. See new figures below.

R2C474. p8,L222 (and p10, L294): \amending" seems an odd choice: to amend would mean to modify data in such a way that the modified data are an improvement.

We propose changing "amending" to "combining".

R2C48p8, L231: Was the notation k for the derived adjacency matrix used before? k was just the \source" cell from time t which feed a \sink" cell at time $t + 1$?

We propose deleting "(k)".

Appendix B Figures

