



Numerical and qualitative contrasts of two statistical models for water quality change in tidal waters

Journal:	<i>Journal of the American Water Resources Association</i>
Manuscript ID	JAWRA-16-0152-P.R2
Manuscript Type:	Technical Paper
Date Submitted by the Author:	17-Oct-2016
Complete List of Authors:	Beck, Marcus; US Environmental Protection Agency, National Health and Environmental Effects Research Laboratory, Gulf Ecology Division Murphy, Rebecca; University of Maryland Center for Environmental Science at Chesapeake Bay Program
Key Terms:	estuaries < GEOGRAPHY, additive models, nutrients < WATER QUALITY, Patuxent River Estuary, statistics < MODELING, time series analysis < MODELING, weighted regression

Numerical and qualitative contrasts of two statistical models for water quality
change in tidal waters

Marcus W. Beck, Rebecca R. Murphy

Research Ecologist (Beck), US Environmental Protection Agency, National Health and
Environmental Effects Research Laboratory, Gulf Ecology Division, 1 Sabine Island Drive, Gulf
Breeze, FL 32561; Estuarine Data Analyst (Murphy), University of Maryland Center for
Environmental Science at Chesapeake Bay Program, 410 Severn Avenue, Suite 112, Annapolis,
MD 21403 (Email/Beck: beck.marcus@epa.gov)

ABSTRACT

Two statistical approaches, weighted regression on time, discharge, and season (WRTDS) and
generalized additive models (GAMs), have recently been used to evaluate water quality trends in
estuaries. Both models have been used in similar contexts despite differences in statistical
foundations and products. This study provided an empirical and qualitative comparison of both
models using 29 years of data for two discrete time series of chlorophyll-*a* (chl-*a*) in the Patuxent
River Estuary. Empirical descriptions of each model were based on predictive performance
against the observed data, ability to reproduce flow-normalized trends with simulated data, and
comparisons of performance with validation datasets. Between-model differences were apparent
but minor and both models had comparable abilities to remove flow effects from simulated time
series. Both models similarly predicted observations for missing data with different
characteristics. Trends from each model revealed distinct mainstem influences of the Chesapeake
Bay with both models predicting a roughly 65% increase in chl-*a* over time in the lower estuary,
whereas flow-normalized predictions for the upper estuary showed a more dynamic pattern, with
a nearly 100% increase in chl-*a* in the last 10 years. Qualitative comparisons highlighted
important differences in the statistical structure, available products, and characteristics of the
data and desired analysis.

KEY TERMS: estuaries, additive models, nutrients, Patuxent River Estuary, statistics, time
series analysis, weighted regression

INTRODUCTION

The interpretation of environmental trends can have widespread implications for the management of natural resources and can facilitate an understanding of ecological factors that mediate system dynamics. An accurate interpretation of trends can depend on the chosen method of analysis, and more importantly, its ability to consider effects of multiple drivers on response endpoints that may be particular to the system of interest. The need to interpret potential impacts of nutrient pollution has been a priority issue for managing aquatic resources (Nixon 1995), particularly for estuaries that serve as focal points of human activities and receiving bodies for upstream hydrologic networks (Paerl *et al.* 2014). Common assessment endpoints for eutrophication in estuaries have included seagrass growth patterns (Steward and Green 2007), frequency and magnitude of oxygen depletion in bottom waters (Paerl 2006), and trophic network connectivity (Powers *et al.* 2005). Additionally, chlorophyll-*a* (chl-*a*) concentration provides a measure of the release of phytoplankton communities from nutrient limitation with increasing eutrophication. Chlorophyll time series have been collected for decades in tidal systems (e.g., Tampa Bay, TBEP (Tampa Bay Estuary Program), (2011); Chesapeake Bay, Harding (1994); datasets cited in Monbet (1992), Cloern and Jassby (2010)), although the interpretation of trends in observed data has been problematic given the inherent variability of time series data. Identifying the response of chl-*a* to different drivers, such as management actions or increased pollutant loads, can be confounded by natural variation from freshwater inflows (Borsuk *et al.* 2004) or tidal exchange with oceanic outflows (Monbet 1992). Seasonal and spatial variability of chl-*a* dynamics (see Cloern (1996)) can further complicate trend evaluation, such that relatively simple analysis methods may insufficiently describe variation in long-term datasets (Hirsch 2014). More rigorous quantitative tools are needed to create an

unambiguous characterization of chl-*a* response independent of variation from confounding variables.

Recent applications of statistical methods to describe water quality dynamics have shown promise in estuaries, specifically weighted regression on time, discharge, and season (WRTDS) and generalized additive models (GAMs). The WRTDS method was initially developed to describe water quality trends in rivers (Hirsch *et al.* 2010, Hirsch and De Cicco 2014) and has recently been adapted to describe chl-*a* trends in tidal waters (Beck and Hagy III 2015). A defining characteristic of WRTDS is a weighting scheme that fits a continuous set of parameters to the time series by considering the influence of location in the record and flow values relative to the period of interest. The WRTDS model has been used to model pollutant delivery from tributary sources to Chesapeake Bay (Hirsch *et al.* 2010, Moyer *et al.* 2012, Zhang *et al.* 2015), Lake Champlain (Medalie *et al.* 2012), the Mississippi River (Sprague *et al.* 2011), and is now being used operationally at the US Geological Survey (USGS) to produce nutrient load and concentration trend results annually for tributaries of the Chesapeake Bay (USGS, Water Quality Loads and Trends at Nontidal Monitoring Stations in the Chesapeake Bay Watershed. Accessed November, 2015, <http://cbrim.er.usgs.gov/>). A comparison to an alternative regression-based model for evaluating nutrient flux, ESTIMATOR, suggested that WRTDS can produce more accurate trend estimates (Moyer *et al.* 2012). As opposed to WRTDS, GAMs were initially developed in a more general context as a modification to generalized linear models to model a response variable as the sum of smoothing functions of different predictors (Hastie and Tibshirani 1990, Wood 2006a). GAMs have recently been used to describe eutrophication endpoints in tidal waters (Haraguchi *et al.* 2015, Harding *et al.* 2016), and exploratory analyses are underway to use GAMs for long-term trend analysis in Chesapeake Bay tidal waters at the

Chesapeake Bay Program. Although the approach was not developed specifically for application to water quality problems, GAMs are particularly appealing because they are less computationally intense and provide more accessible estimates of model uncertainty than WRTDS. Both approaches appear to have similar potential to characterize system dynamics, but the relative merits of each have not been evaluated. Comparisons that describe the accuracy of empirical descriptions and the desired products could inform the use of each model to describe long-term changes in ecosystem characteristics.

The goal of this study is to provide an empirical and qualitative description of the relative abilities of WRTDS and GAMs to describe long-term changes in time series of eutrophication response endpoints in tidal waters. Two discrete time series covering 1986-2014 from two stations in the Patuxent River estuary are used as a common dataset for evaluating each model. The Patuxent Estuary is a well-studied tributary of the Chesapeake Bay system that has been monitored for several decades with fixed stations along the longitudinal axis. Two stations were chosen as representative time series that differed in the relative contributions of watershed inputs and influences from the mainstem of the Chesapeake, in addition to known historical events that have impacted water quality in the estuary. The specific objectives of the analysis were to 1) provide a narrative comparison of the statistical foundation of each model, both as a general description and as a means to evaluate water quality time series, 2) use each model to develop an empirical description of water quality changes at each monitoring station given known historical changes in water quality drivers, 3) evaluate each model's ability to reproduce flow-normalized trends as known components of simulated time series, and 4) compare each technique's ability to describe changes, as well as the differences in the information provided by each. We conclude with recommendations on the most appropriate use of each method, with particular attention

given to the desired products and characteristics of a dataset that could influence interpretation of model results.

METHODS

Study site and water quality data

The Patuxent River estuary, Maryland, is a tributary to Chesapeake Bay on the Atlantic coast of the United States (Figure 1). The longitudinal axis extends 65 km landward from the confluence with the mesohaline portion of Chesapeake Bay. Estimated total volume at mean low water is $577 \times 10^6 \text{ m}^3$ and a surface area of $126 \times 10^6 \text{ m}^2$. The lower estuary (below 45 km from the confluence) has a mean width of 2.2 km and depth of 6 m (Cronin and Pritchard 1975), whereas the upper estuary has a mean width of 0.4 km and mean depth of 2.5 m (Hagy 1996). The lower estuary is seasonally stratified and the upper estuary is vertically mixed. A two-layer circulation pattern occurs in the lower estuary characterized by an upper seaward-flowing layer and a lower landward-flowing layer. A mixed diurnal tide dominates with mean range varying from 0.8 m in the upper estuary to 0.4 m near the mouth (Boicourt and Sanford 1998). The estuary drains a 2300 km^2 watershed that is 51% wooded, 35% developed, 9% cropland, and 5% pasture/hay (USEPA 2010). The USGS stream gage on the Patuxent River at Bowie, Maryland measures discharge from 39% of the watershed. Daily mean discharge from 1985 to 2014 was $11.0 \text{ m}^3 \text{ s}^{-1}$, with abnormally high years occurring in 1996 (annual mean $20.0 \text{ m}^3 \text{ s}^{-1}$) and 2003 (annual mean $22.5 \text{ m}^3 \text{ s}^{-1}$).

The Chesapeake Bay Program and Maryland Department of Natural Resources (MDDNR) maintain a continuous monitoring network for the Patuxent at multiple fixed stations that cover the salinity gradient from estuarine to tidal fresh (<http://www.chesapeakebay.net/>, Figure 1 and Table 1). Water quality samples have been collected by MDDNR since 1985 at

1
2
3 monthly or bimonthly intervals and include salinity, temperature, chl-*a*, dissolved oxygen, and
4
5 additional dissolved or particulate nutrients and organic carbon. Seasonal variation in chl-*a* is
6
7 observed across the stations with spring and summer blooms occurring in the upper, oligohaline
8
9 section, whereas chl-*a* is generally higher in the lower estuary during winter months (Figure 2).
10
11 Chl-*a* concentrations are generally lowest for all stations in late fall and early winter. Periods of
12
13 low flow are associated with higher chl-*a* concentrations in the upper estuary, whereas the
14
15 opposite is observed for high flow. Stations TF1.6 and LE1.2 were chosen as representative time
16
17 series from different salinity regions to evaluate the water quality models. Observations at each
18
19 station capture a longitudinal gradient of watershed influences at TF1.6 to main stem influences
20
21 from the Chesapeake Bay at LE1.2. Long-term changes in chl-*a* have also been related to
22
23 historical reductions in nutrient inputs following a statewide ban on phosphorus-based detergents
24
25 in 1984 and wastewater treatment improvements in the early 1990s that reduced point sources of
26
27 nitrogen (Lung and Bai 2003, Testa *et al.* 2008). Therefore, the chosen stations provide unique
28
29 datasets to evaluate the predictive and flow-normalization abilities of each model given the
30
31 differing contributions of landward and seaward influences on water quality.
32
33
34
35
36
37

38
39 Thirty years of chl-*a* and salinity data from 1986 to 2014 were obtained for stations
40
41 TF1.6 (n = 522) and LE1.2 (n = 530, Chesapeake Bay Program data hub, Accessed March 23,
42
43 2015, <http://www.chesapeakebay.net/data>). All data were vertically integrated throughout the
44
45 water column for each date to create a representative sample of water quality. The integration
46
47 averaged all values after interpolating from the surface to the maximum depth. Observations at
48
49 the most shallow and deepest sampling depth were repeated for zero depth and maximum depths,
50
51 respectively, to bound the interpolations within the range of the data. Daily flow data were also
52
53 obtained from the USGS stream gage station at Bowie, Maryland and merged with the nearest
54
55
56
57
58
59
60

date in the chl-*a* and salinity time series. Initial analyses suggested that a moving-window average of discharge for the preceding five days provided a better fit to the chl-*a* data at TF1.6, whereas the salinity record was used as a tracer of discharge at LE1.2. Both chl-*a* and discharge data were log-transformed. Censored data were not present in any of the datasets. Initial quality assurance checks for all monitoring data were conducted following standard protocols adopted by the Chesapeake Bay Program.

Model descriptions

Weighted Regression on Time, Discharge, and Season.

The WRTDS method relates a response variable, typically a nutrient concentration, to discharge and time to evaluate long-term trends (Hirsch *et al.* 2010, Hirsch and De Cicco 2014). Recent adaptation of WRTDS to tidal waters relates chl-*a* concentration to salinity and time (Beck and Hagy III 2015), where salinity is a tracer of freshwater inputs or tidal changes (R package link in Appendix A). The functional form of the model is a simple regression that relates the natural log of chl-*a* (*Chl*) to decimal time (*T*) and salinity (*Sal*) on a sinusoidal annual time scale (i.e., cyclical variation by year).

$$\ln(\text{Chl}) = \beta_0 + \beta_1 T + \beta_2 \text{Sal} + \beta_3 \sin(2\pi T) + \beta_4 \cos(2\pi T) + \varepsilon \quad (1)$$

The tidal adaptation of WRTDS uses quantile regression models (Cade and Noon 2003) to characterize trends in different conditional distributions of chl-*a*, e.g., the median or 90th percentile. For comparison to GAMs, a version of WRTDS created by the authors similar the original model in Hirsch *et al.* (2010) was used to characterize the conditional mean of the response (see Appendix A). Mean models require an estimation of the back-transformation bias parameter for response variables in log-space (Hirsch *et al.* 2010). Although back-transformation is developed for WRTDS, a similar approach has not yet been implemented for GAMs. For

1
2
3 simplicity and ease of comparison, all units for chl-*a* are reported in log-space unless otherwise
4
5 noted.
6
7

8 The WRTDS approach obtains fitted values of the response variable by estimating a
9
10 unique regression model at each point in the time series. Each model is weighted with a three-
11
12 dimensional window by month, year, and salinity (or flow) such that a unique set of regression
13
14 parameters for each observation is obtained. For example, a weighted regression centered on a
15
16 single observation weights other observations in the same year, month, and similar salinity with
17
18 higher importance, whereas observations for different months, years, or salinities receive lower
19
20 importance. This weighting approach allows estimation of regression parameters that vary in
21
22 relation to observed conditions throughout the period of record (Hirsch *et al.* 2010). Optimal
23
24 window widths can be identified using cross-validation, described below, that evaluates the
25
26 ability of the model to generalize results with novel datasets.
27
28
29
30
31

32 Predicted values are based on an interpolation matrix from the unique regressions at each
33
34 time step. A sequence of salinity or flow values based on the minimum and maximum values for
35
36 the data are used to predict chl-*a* using the observed month and year based on the parameters fit
37
38 to the observation. Model predictions are based on a bilinear interpolation from the grid using the
39
40 salinity (flow) and date values closest to observed. Salinity- or flow-normalized values are also
41
42 obtained from the prediction grid that allow an interpretation of chl-*a* trend that is independent of
43
44 variation related to freshwater inputs. Normalized predictions are obtained for each observation
45
46 by collecting the sample of observed salinity or flow values that occur for the same month
47
48 throughout all years in the dataset. These values are assumed to be equally likely to occur across
49
50 the time series at that particular month. A normalized value for each point in the time series is
51
52 the average of the predicted values from each specific model based on the salinity or flow values
53
54
55
56
57
58
59
60

that are expected to occur for each month.

Generalized Additive Models.

A GAM is a statistical model that allows for a linear predictor to be represented as the sum of multiple smooth functions of covariates (Hastie and Tibshirani 1990). In this application, GAMs were constructed with the same explanatory variables as the WRTDS approach: log of chl-*a* was modeled as a function of decimal time, salinity or flow, and day of year (i.e., to capture the annual cycle). Multiple types of smooth functions could be used in a GAM (Hastie and Tibshirani 1990), and our implementation relies on thin plate regression splines (Wood 2006a). A spline is a piece-wise function (e.g., a polynomial) whose pieces are connected at knots, or breakpoints, where the functions are joined smoothly (Hastie and Tibshirani 1990). The thin plate regression spline has the benefit that a user is not required to select knot locations for a spline explicitly, but only selects a reasonable upper limit on the flexibility of the function. Within that limit, the balance between model fit and smoothness is achieved by minimizing both error and “wiggleness” of the function with a smoothness parameter, which minimizes the generalized cross-validation score, controlling the tradeoff (Wood 2006a). To allow for interaction between the model covariates (e.g., seasonal differences in the long-term chl-*a* pattern), a tensor product basis (Wood 2006b) between all three covariates was constructed.

Predictions with GAMs are straightforward after the model is fit and can be estimated along with standard errors based on the Bayesian posterior covariance matrix (Wood 2006a). For this comparison, salinity- or flow-normalized GAM predictions were obtained in a manner consistent with that used for WRTDS. The observed salinity or flow values were compiled that occurred in the same month throughout all years in the dataset. These values were assumed to be equally likely to occur at that particular month. A normalized GAM estimate at each date in the

record was computed as the average of the predictions obtained using all of the flow or salinity values for that month of the year throughout the record.

Methodological contrasts of WRTDS and GAMs.

WRTDS and GAMs are statistical models that have very similar functional forms. Both use core models that empirically describe a response variable as numerical combinations of one or more explanatory variables. As noted above, the core functional model of WRTDS is a simple linear regression that relates pollutant concentration to fixed effects of time, discharge, and season. In a simple regression the fixed effects are parameterized by a single set of model coefficients that describe the linear relationship with the response variable. For WRTDS, this simple regression structure is used, but at each time step a different coefficient set is created based on the relative weighting of the data. By comparison, GAMs link individual explanatory variables with the response using smoothing functions for each variable instead of fixed parameters. As such, the functional forms of both models are conceptually identical where the only difference is the type of parameter used in each (fixed or smoothing function). These functional similarities can potentially explain why both models have been used for the similar purpose of describing pollutant trends over time (e.g., USGS 2015, Harding *et al.* 2016).

Although both models use functional forms with similar structures, the statistical similarities of WRTDS and GAMs depart during model estimation when parameters are fit to the observed data. This difference is critical for understanding the need to describe potential differences between model results and guidance for appropriate use of each. As previously described, WRTDS results are based on repeated multiple linear regressions that are each weighted separately depending on location of an observation in the time, discharge, and season domain. This results in a multi-dimensional parameter set that varies smoothly across the time

series and that has no resemblance to results from a single parameter set that is fit to the entire time series. The final parameter set produces results that are more similar to a locally-estimated (LOESS) polynomial curve (i.e., Cleveland 1979) than a simple regression. By contrast, this implementation of GAMs estimates the smoothing functions for the explanatory variables using a spline-fitting process that results in individual (although quite complicated) spline functions fit across the entire data set for each explanatory variable. Although parallels between GAM fits can be made with both LOESS and WRTDS, the relationship between response and explanatory variables described by the hyper-dimensional smoothing surface from WRTDS is a different theoretical approach than a set of spline functions fit across all the data with GAMs. Therefore, a reasonable expectation is that different estimation techniques used by WRTDS and GAMs can lead to different descriptions of relationships between variables.

Selection of model parameters.

The selection of optimal model parameters is a challenge that represents a tradeoff between model precision and ability to generalize to novel datasets. Weighted regression requires identifying optimal half-window widths, whereas the GAM approach used here requires identifying an optimal value for a smoothing parameter that weights the wiggleness of the function against model fit (Wood 2006a). Both models used a form of cross-validation to identify model parameters that maximize the precision of model predictions with novel data. For the GAM approach, generalized cross-validation was used to obtain the optimal smoothing parameter in an iterative process with penalized likelihood maximization to solve for model coefficients. The effective degrees of freedom of the resulting model varied with the smoothing parameter (Wood 2006a). Similarly, the tidal adaptation of WRTDS used k-fold cross-validation ($k = 10$) to identify the optimal half-window widths (Efron and Tibshirani 1993, Arlot and

1
2
3 Celisse 2010). Evaluating multiple combinations of window-widths can be computationally
4
5 intensive. An optimization function was implemented in R (RDCT 2015) to more efficiently
6
7 evaluate model parameters with cross-validation. Window widths were searched using the
8
9 limited-memory modification of the BFGS quasi-Newton method that imposes upper and lower
10
11 bounds for each parameter (Byrd *et al.* 1995, Nocedal and Wright 2006). The chosen parameters
12
13 were based on a selected convergence tolerance for the error minimization of the search
14
15 algorithm. Specifically, the algorithm converged when the reduction in the minimization
16
17 function for a given change in parameters was within an acceptable tolerance without excessive
18
19 search time.
20
21
22
23

24 *Comparison of modelled trends*

25
26
27 Separate WRTDS and GAMs were created using the above methods for the chl-*a* time
28
29 series at TF1.6 and LE1.2. For each model and station, a predicted and flow-normalized
30
31 (hereafter flow-normalized refers to flow at TF1.6 and salinity at LE1.2) time series was
32
33 obtained for comparison. The results were compared with summary statistics that evaluated both
34
35 the predictive performance to describe observed chl-*a* and direct comparisons between the
36
37 models. Emphasis was on agreement between observed and predicted values, rather than
38
39 uncertainty associated with parameter estimates or model results. As of writing, methods for
40
41 estimating confidence intervals of WRTDS have been developed for the original model (Hirsch
42
43 *et al.* 2015), but have not been fully developed for application to WRTDS in tidal waters. In
44
45 addition to visual evaluation of trends over time, summary statistics used to compare model
46
47 predictions to observed chl-*a* included root mean square error (RMSE) and average differences.
48
49 For all comparisons, RMSE comparing each model's predictions to observed chl-*a* (fit) was
50
51 defined as:
52
53
54
55
56
57
58
59
60

$$RMSE = \sqrt{\frac{\sum_{i=1}^n \left(Chl_i - \hat{Chl}_i \right)^2}{n}} \quad (2)$$

where n is the number of observations for a given evaluation, Chl_i is the observed value of chl- a for observation i , and \hat{Chl}_i is the predicted value of chl- a for observation i . RMSE values closer to zero represent model predictions closer to observed. Comparisons between models were performed similarly using the root mean square difference (RMSD):

$$RMSD = \sqrt{\frac{\sum_{i=1}^n \left(\hat{Chl}_{WRTDS,i} - \hat{Chl}_{GAM,i} \right)^2}{n}} \quad (3)$$

where the estimated chl- a values for each model, $\hat{Chl}_{WRTDS,i}$, and $\hat{Chl}_{GAM,i}$, are compared directly.

Similarly, average difference (or bias) of predictions between models as a percentage was defined as:

$$\text{Average difference} = \frac{\sum_{i=1}^n \hat{Chl}_{WRTDS,i} - \sum_{i=1}^n \hat{Chl}_{GAM,i}}{\left(\sum_{i=1}^n \hat{Chl}_{WRTDS,i} + \sum_{i=1}^n \hat{Chl}_{GAM,i} \right) / 2} * 100 \quad (4)$$

Positive values indicate that WRTDS provided higher predictions than GAMs on average, whereas the opposite is true for negative values (Moyer *et al.* 2012). Results between models were also evaluated using regressions comparing WRTDS (as the response) and GAM (as the predictor). The regressions were compared to a null model having an intercept of zero and slope of one. Deviation of either the intercept or slope of the regressions from the null model provided evidence of systematic differences between the models. An intercept significantly different from zero was interpreted as an overall difference between the predictions, whereas a slope different from one was interpreted as a difference that varies with relative magnitude of the predictions.

Although the signs of the slope and intercept estimates for the comparisons depended on which model was used as the predictor, we were primarily concerned with magnitude of the parameter estimates in the regression comparisons as evidence of systematic differences between each model. The statistical comparisons described above were conducted for the entire time series at each station to evaluate overall performance. Different time periods were also evaluated to identify potential temporal variation in results, which included a comparison of results by annual and seasonal aggregations and periods with different levels of flow using the discharge record at Bowie, Maryland. Annual and seasonal aggregations shown in Figure 2 were evaluated between the models, in addition to evaluating the models at different levels of flow defined by the quartile distributions (min–25%, 25%–median, median–75%, and 75%–max).

Flow-normalized time series were compared similarly but only between models because the true flow-independent component of the observed data is not known and can only be empirically estimated. As described below, an evaluation of flow-normalized data for each model was accomplished using simulated datasets with known components that were independent of discharge. However, a simple comparison of flow-normalized trends by different time periods summarized long-term patterns in the Patuxent River estuary. These comparisons evaluated percent changes of flow-normalized estimates at the beginning and end of each time period. Percent changes within each period were based on annual mean estimates for the first and last three years of flow-normalized chl-*a* estimates, excluding the annual aggregations that had limited annual mean data (i.e., seven years per period). For example, percent change for the January-February-March (JFM) seasonal period compared an average of JFM annual means for 1986 through 1988 to an average of JFM annual means for 2012 through 2014. This approach was used to reduce the influence of abnormal years or missing data on trend estimates.

Comparison of flow-normalized trends

The relative abilities of each model to characterize flow-normalized trends in chl-*a* were evaluated using simulated datasets with known components. This approach was used because the flow-independent component of chl-*a* can only be empirically estimated from raw data. Accordingly, the ability of each model to isolate the flow-normalized trend cannot be evaluated with reasonable certainty unless the true signal is known. Following similar concepts in Beck *et al.* (2015), simulated time series of observed chl-*a* (Chl_{obs}) were created as additive components related to flow (Chl_{flo}) and a flow-independent biological component of chl-*a* (Chl_{bio}):

$$Chl_{obs} = Chl_{flo} + Chl_{bio} \quad (5)$$

A distinction between Chl_{flo} and Chl_{bio} is that the former describes variation in the observed time series with changes in discharge (e.g., concentration dilution with increased flow) and the latter describes a true, desired measure of chl-*a* in the water column that is directly linked to productivity. The biological component of chl-*a* is comparable to an observation in a system that is not affected by flow and is the time series that is estimated by flow-normalization with WRTDS and GAMs.

The simulated time series was created using methods similar to those in Hirsch *et al.* (2015) and was based on a stochastic model derived from actual flow and water quality measurements to ensure the statistical properties were comparable to existing datasets. This approach allowed us to evaluate GAMs and WRTDS under different sampling regimes (e.g., monthly rather than daily), while ensuring the simulated datasets had statistical properties that were consistent with known time series. Daily flow observations from the USGS stream gage station 01594440 near Bowie, Maryland (38° 57' 21.3'' N, 76° 41' 37.3'' W) were obtained from 1985 to 2014. Daily chl-*a* records were estimated from fluorescence values from the Jug Bay

station (38° 46' 50.6'' N, 76° 42' 29.1'' W) of the Chesapeake Bay Maryland National Estuarine Research Reserve in the upper Patuxent.

Four time series were estimated or simulated from the actual datasets to create the complete, simulated time series: 1) estimated discharge as a stationary seasonal component (\hat{Q}_{seas}), 2) simulated error structure from the residuals of the seasonal discharge model ($\varepsilon_{Q,sim}$), 3) estimated chl-*a* independent of discharge as a stationary seasonal component (\hat{chl}_{seas}), and 4) simulated error structure from the residuals of the seasonal chl-*a* model ($\varepsilon_{chl,sim}$). Unless otherwise noted, chl-*a* and discharge are in natural log-transformed units. Each of the four components was used to simulate the components in eq. (5):

$$Chl_{flo} = I \left(\hat{Q}_{seas} + \sigma_{\varepsilon} \cdot \varepsilon_{Q,sim} \right) \quad (6)$$

$$Chl_{bio} = \hat{chl}_{seas} + \sigma_{\hat{chl}_{seas}} \cdot \varepsilon_{chl,sim} \quad (7)$$

Standard deviation for the residuals of the seasonal flow component, σ_{ε} , and the random errors,

$\varepsilon_{Q,sim}$, are derived from the observed data (see Appendix B). The estimated flow time series

within the parentheses, $\hat{Q}_{seas} + \sigma_{\varepsilon} \cdot \varepsilon_{Q,sim}$, is floored at zero to simulate an additive effect of

increasing flow on Chl_{obs} . Although the actual relationship of water quality measurements with

flow is more complex, we assumed that a simple addition was sufficient for the simulations

where the primary objective was to create an empirical and linear link between flow and chl-*a*.

The vector I (where $0 < I < 1$) is a weighting and unit-conversion vector that translates the terms

enclosed in parentheses from flow to chl-*a* concentration units and allows for the effect of flow

to be defined as time-varying. For example, a flow effect that changes from non-existent to

positive throughout the period of observation can be simulated by creating a vector ranging from zero to one. For the simulated Chl_{bio} time series, the seasonal and error components were characterized using the daily time series at Jug Bay that likely included an effect of flow in the observed data. For the simulated models, we assumed that the actual flow effect was part of the seasonal component to obtain an accurate estimate of the error component that was independent of both flow and season. In other words, the seasonal component of chl-*a* was modelled with a discharge component to remove any variability related to flow in the residuals. Methods for estimating each of the components in eqs. (6) and (7) are described in detail in Appendix B and Figure B1.

Three time series with monthly sampling frequencies and varying contributions of the flow component (Chl_{flo} in eqs. (5) and (6)) were created from daily simulated time series of Chl_{obs} (Appendix B). One day in each month for each year (e.g., January 5, 2010, February 19, 2010,..., January 28, 2011, February 1, 2011, etc.) was randomly sampled and used as an approximate monthly time step for each time series. Varying effects of the flow component on observed chl-*a* were created by multiplying Chl_{flo} by different indicator vectors (I in eq. (6)). The contribution of the flow component varied from non-existent (vector of zeroes), constant (vector of ones), and steadily increasing (continuous vector from zero to one). This created three monthly time series that were used to evaluate each model that were analogous to no influence, constant, and changing influence of the flow component over time (Figure B2). Results were evaluated by first comparing the predicted (\hat{Chl}_{obs}) and observed (Chl_{obs}) chl-*a* values for each simulation, followed by comparing the flow-normalized results (\hat{Chl}_{bio}) from each model to the original biological chl-*a* (Chl_{bio}) component of each simulated time series (eqs. (5) and (7)). The former comparison provided information on relative fit to validate the simulated data, whereas

the latter comparison to evaluate flow-normalization was the primary focus of the analysis.

Model comparisons with independent data

The final analysis provided a complementary comparison to those described above for model performance by evaluating the ability of both models to predict missing or novel data. Prediction performance was evaluated for validation datasets to provide a measure that was completely independent of the data used to train the models. This analysis used the simulated time series with a constant flow effect that was described in the previous section. Weekly samples at a fixed interval were taken from the daily time series to ensure sufficient data (compared to a longer time step), minimal processing time (compared to using daily data), and sampling structure similar to common monitoring datasets in tidal waters. The weekly time series was split multiple times into training and validation datasets to evaluate effects of different ratios of training-to-validation (1:1, 2:1, etc.), and characteristics of the missing data such as random missing data or data missing in blocks. Block sampling, in addition to completely random sampling, was used to account for temporal correlation, i.e., missing data frequently occur in blocks of time due to equipment failure or funding changes. RMSE of model predictions for GAMs and WRTDS were evaluated for validation-to-training split ratios ranging from 5-50% (e.g., validation was 5% and training was 95% of total, validation 10% and training 90%, etc.) and sampling from completely random to blocks of increasing size. Because the data splits and blocks were stochastic, 100 replicates were created for each split ratio and block sampling level to place a range on model performance.

RESULTS

Observed trends and relative fit

The optimal half-window widths and degrees of freedom for smoothing varied at each

station for WRTDS and GAMs, respectively. For WRTDS, optimal half-window widths identified by generalized cross-validation were 0.25 as a proportion of each year (seasonal component, sinusoidal terms in eq. (1)), 13.59 years (T in eq. (1)), and 0.25 as a proportion of the total range of salinity (Sal in eq. (1)) for LE1.2, and 0.25 of each year, 6.28 years, and 0.50 of flow at TF1.6. For both stations, the optimization method selected relatively wide windows for the year weights while minimizing the seasonal and flow components. For GAMs, the optimal smoothing procedure resulted in a smoother model at LE1.2 than TF1.6 with effective degrees of freedom of 35.5 and 71.4, respectively. The smoothing method used for the GAMs does not split the degrees of freedom among the three interacting variables.

The predicted chl- a from each model generally followed patterns in observed chl- a from 1986 to 2014 (Figure 3). At LE1.2, each model showed seasonal minima typically in November, whereas maximum chl- a was observed in a spring bloom, typically March or April (Figure 4). A secondary, smaller seasonal peak was also observed in late summer from bottom-layer regeneration and upward nutrient transport (Testa *et al.* 2008). Seasonal variation at TF1.6 was noticeably different with an initial peak typically observed in May and a larger dominant bloom occurring in September or October (Figure 4). Elevated chl- a concentrations were also more prolonged than those at LE1.2 with only a slight decrease between the two seasonal blooms. A seasonal minimum was typically observed in December or January, followed by a rapid increase in the following months. Differences in magnitude of the seasonal range were also less pronounced at LE1.2 compared to TF1.6, with differences throughout the year approximately 3 $\mu\text{g/L}$ of chl- a (arithmetic) at LE1.2 and 7 $\mu\text{g/L}$ of chl- a at TF1.6. Visual evaluation of seasonal trends suggested each model provided similar results, although WRTDS predictions had slightly better fits at the extreme ends of the distribution of chl- a (Figure 3a). Normalized predictions for

both models were visually distinct from the non-flow-normalized predictions such that seasonal minima and maxima and extreme predictions were not observed with the normalized values. Overall, both models had predictions that provided a more adequate visual description of the range of chl-*a* at TF1.6 as compared to LE1.2 where observed values lower or higher than the predictions were more common.

Quantitative summaries of model fit by site indicated that performance between sites and models was similar, with RMSE ranging from a minimum of 0.50 at TF1.6 for GAM predictions and a maximum of 0.52 at TF1.6 for WRTDS predictions (Table 2). Overall, both models performed similarly, although WRTDS had slightly better performance at LE1.2 and GAMs had slightly better performance at TF1.6 (Table 2). Fit by different time periods generally showed agreement between methods during periods when performance was relatively high or low. For LE1.2, both models had the worst fit during the 2001-2007 annual period (RMSE 0.61 for GAMs, RMSE 0.60 for WRTDS), the April-May-June (AMJ) seasonal periods (0.64 for GAMs, 0.64 for WRTDS), and periods of high flow (0.64 for GAMs, 0.63 for WRTDS). For TF1.6, models had the worst fit during the 1994-2000 annual period (0.55 for GAMs, 0.58 for WRTDS) and the AMJ seasonal period (0.54 for GAMs, 0.58 for WRTDS). Errors between models were comparable for all flow periods at TF1.6, with the exception of lower errors during low flow (0.45 for GAMs, 0.46 for WRTDS). In general, model performance was partially linked to flow such that fit was improved during periods of low flow, including seasonal or annual periods of low flow. For example, both models at both sites had the best fit during the July-August-September (JAS) period when seasonal flow was minimized (Table 2 and Figure 2).

Results as annual aggregations suggested that chl-*a* patterns between years have not been constant and are considerably different between sites (Figure 3b). Both models showed a gradual

and consistent increase in chl-*a* at LE1.2, with values increasing by approximately 1.5 $\mu\text{g/L}$ (arithmetic) from 1986 to 2014. Predictions at TF1.6 did not show a similar increase from the beginning to the end of the time series, although a dramatic decrease from approximately 12 $\mu\text{g/L}$ to 6 $\mu\text{g/L}$ from 2000 to 2006 was observed. By 2014, chl-*a* returned to values similar to those prior to the initial decrease. Flow-normalized predictions that were annually averaged at each site allowed an interpretation of trends that were independent of variation in discharge or salinity (Tables 3 and 4). Overall percent change of ln-transformed chl-*a* concentration from the beginning to the end of the time series at LE1.2 was approximately 20% (Table 3). A slight decrease in chl-*a* at TF1.6 was observed from 1986 to 2014 (Table 4). Changes by annual, seasonal, and flow time periods at LE1.2 were comparable for each time period and model type, although some differences were observed. For example, both models had maximum increases in chl-*a* for the different flow periods for high levels of flow at LE1.2 (25.1% for GAMs, 22.3% for WRTDS). Trends by different time periods were more apparent for TF1.6, particularly as an overall decrease in chl-*a* for both models during the 2001–2007 period and an overall increase during the 2008–2014 period (Table 4). Seasonal changes were especially pronounced during the January-February-March (JFM) and October-November-December (OND) periods where both models showed an increase and decrease, respectively, with differences between the two (JFM period, 9% for GAMs, 32.7% for WRTDS; OND period, –18.2% for GAMs, –17.5% for WRTDS). Percent changes by flow quantile were also observed at TF1.6, with the most noticeable difference from LE1.2 being a decrease in chl-*a* during both high and low flow (both models) and relatively larger increases in chl-*a* during moderate flow.

Comparison of model predictions

The following describes direct comparisons of model results, whereas the previous

section emphasized results relative to trends over time and fit to the observed data. Accordingly, direct comparisons were meant to identify instances when models results were systematically different from each other. Table 5 compares average differences and RMSD of results between each model for the complete time series and different subsets by annual, seasonal, and flow periods. Overall, differences between the models were minor with most percent differences not exceeding 1% and no RMSD values exceeding 0.15. Model differences between different time periods were not apparent for either station, although the largest average difference was observed at TF1.6 for the 2008–2014 time period (3.1%, WRTDS greater than GAMs).

Regressions comparing model results (WRTDS as response, GAMs as predictor) provided additional information about overall differences (significantly different intercept) and differences between the models that varied for different values (significantly different slope) (Table 6, Appendix C). Significant differences were observed for the entire time series such that estimated intercepts and slopes were different from zero and one, respectively, for both stations and model predictions (observed and flow-normalized), excluding intercepts and slopes for the flow-normalized predictions at TF1.6 ($\beta_{0, norm}$ and $\beta_{1, norm}$). Differences were also observed for the time period subsets, with the most obvious differences occurring for the seasonal aggregations. For example, all comparisons between the models for both sites and model predictions had intercept estimates significantly greater than zero and slope estimates significantly less than one for the AMJ period (Table 6). For almost all significant differences, intercept estimates were greater than zero and slope estimates were less than one. Visual comparisons of results in Appendix C confirm those in Table 6, particularly differences in the seasonal aggregations.

Changes in chl-a response to flow over time

Both models described chl-a response with sufficient parameterization of input variables

to evaluate variation with flow changes over time. As in Beck and Hagy III (2015), changes in the relationship of chl-*a* to flow can be evaluated by predicting observed chl-*a* across the range of observed flow (or salinity) values for each year in the time series. Visual information obtained from these plots are useful to identify periods of time when chl-*a* was or was not related to changes in flow and may also lead to the development of hypotheses regarding changes in drivers of water quality, e.g., temporal shifts in point-sources to non-point sources of pollution (Hirsch *et al.* 2010, Beck and Hagy III 2015). The only difference between the models in creating the plots is that the three-dimensional prediction grid of chl-*a*, flow, and time created during model fitting is used for WRTDS, whereas the plots for GAMs are based on post-hoc model predictions with covariates defined to vary over a regular grid.

Figure 5 shows the estimated changes from each model in predicted chl-*a* for salinity (LE1.2) or flow (TF1.6) across all years in the study period. The plots are also separated by months of interest to isolate effects of seasonal variation. Visual assessment of the plots suggests that the relationships were dynamic across the study years and varied considerably between LE1.2 and TF1.6. For example, the October plots show decreasing sensitivity of chl-*a* with increasing flow (decreasing salinity) at LE1.2 from early to late in the time series (i.e., a strong, positive relationship changing to a weak relationship over time). Conversely, the opposite trend is observed at TF1.6 in October such that a weak relationship with flow is observed early in the time series and a strong, negative relationship is observed later in the time series, although overall chl-*a* has decreased over time. Additionally, both models provided similar indications of the changes over time, regardless of site or time of year. However, some differences between the models were observed. For example, WRTDS results for January at LE1.2 provided a wider range, or potentially less stable fit of chl-*a* to salinity changes in the earlier years.

Flow-normalization with simulated data

WRTDS and GAMs were fit to each of the three simulated datasets, creating six models to evaluate the general fit of observed to predicted ($Chl_{obs} \sim \hat{Chl}_{obs}$) and biological to flow-normalized chl-*a* ($Chl_{bio} \sim \hat{Chl}_{bio}$). Models were fit using identical methods as those for the Patuxent time series such that an optimal window width combination for WRTDS and optimal degrees of freedom for smoothing parameters with GAMs were identified. Figure 6 is similar to Figure 5 and shows an example of the changing relationships between chl-*a* and flow across the simulated time series using the results from three optimal WRTDS models. The plots show the varying effects of flow in each simulated dataset over time (no effect, constant, increasing) and that the models appropriately characterized the relationships. For example, a changing response of chl-*a* to flow is apparent in the third column of Figure 6 such that no response is observed early in the time series followed by an increase in the response of chl-*a* to flow later in the time series. Similar patterns were observed between models, although there is some suggestion that GAMs are not separating the effect of flow and time as completely as WRTDS. Specifically, results for WRTDS with no influence and a constant influence of flow showed less variation than GAMs in the relationship between chlorophyll and flow over time, consistent with the empirical relationships used to create the simulated time series.

Comparisons of fit to the simulated time series showed no systematic differences between the models. Overall, WRTDS results had lower RMSE than GAMs for all comparisons except one ($Chl_{obs} \sim \hat{Chl}_{obs}$, constant flow simulation), although differences in performance were minor (Table 7). Although both models provided similar performance for individual simulations, differences between the simulations were observed. The different effects of flow had a negative effect on the ability of each model to remove the flow component. Comparisons of Chl_{bio} with

\hat{Chl}_{bio} showed the lowest RMSE with no flow effect and the highest with a constant flow effect (Table 7). Different flow effects did not have an influence on the relationship between predicted (\hat{Chl}_{obs}) and observed (Chl_{obs}) chl-*a* such that RMSE for all models and simulations were similar and lower than those comparing the flow-normalized results. Overall, changing the flow component primarily affected the ability of each model to reproduce the flow-normalized component (\hat{Chl}_{bio}) with relatively minor differences between the models.

Model comparisons with independent data

Both models performed similarly for the training datasets based on different splits of the weekly simulated data (median RMSE ~0.52 for both, Figure 7). Overall, median RMSE values decreased slightly as the ratio of validation-to-training data sets size increased (5% to 50% validation), although the range of RMSE values increased. Similar patterns were observed for the validation datasets (median RMSE ~0.54 for both models), although the ranges decreased as more data were included in the validation datasets. Sampling characteristics for the validation datasets (random and block samples) did not have a noticeable effect on training RMSE for either model, although slightly greater variation in the median RMSE was observed for the largest block size (100% of validation data in a single block).

DISCUSSION

Numerical comparisons

A general conclusion from our quantitative comparisons is that both models provided similar information, both in predictive performance and trends over time. Although some instances were observed where one model had lower errors, large differences were not observed and we emphasize that any potential improvement in performance at the scale shown in Table 2 is trivial. Prediction errors for either model could easily be improved by slight adjustments of the

1
2
3 model parameters. This highlights a potential risk of using prediction error as a performance
4
5 metric because the values are sensitive to tuning parameters and the statistical characteristics of
6
7 training datasets. To address this issue, comparable methods for model development were
8
9 implemented to ensure valid comparisons. Both WRTDS and GAMs used a form of cross-
10
11 validation that was meant to identify the most parsimonious parameter space. A more generic
12
13 benefit of cross-validation is that model development was not biased by analyst intervention as
14
15 the parameters were chosen with predefined heuristics. This paper presents the first application
16
17 of a statistical method of selecting optimal window widths for WRTDS. Further work should
18
19 refine use of these methods to develop robust and unbiased parameters for WRTDS.
20
21
22
23

24
25 The comparisons of predictive performance should also be interpreted relative to the
26
27 statistical foundations of each model. The smoothing process in GAMs, although mathematically
28
29 involved, rapidly converges to a solution, whereas the fitting process for WRTDS is much longer
30
31 because a unique regression is estimated for every point in the time series. From a practical
32
33 perspective, the comparable error estimates for each model's predictions suggests that GAMs are
34
35 advantageous because there is no apparent benefit of the added computational time of WRTDS.
36
37 Temporal changes in the relationship between chl-*a* and flow were also comparable. For
38
39 example, Figure 5 shows similar information for each model, although different methods were
40
41 used to characterize chl-*a* variation from salinity or flow over time. Additional insight into trends
42
43 might be a logical expectation with the added computational time required to estimate WRTDS
44
45 interpolation grids. Conventional modelling techniques have been described as 'statistical
46
47 straightjackets' that can inadequately characterize variation in the data with a limited parameter
48
49 space and structural constraints (Hirsch 2014). WRTDS is meant to provide a contrasting
50
51 approach where the data mold the results using multiple parameter sets. In contrast, one might
52
53
54
55
56
57
58
59
60

1
2
3 expect GAMs to be over-constrained by following a potentially less flexible model composed of
4
5 one smoothing function per explanatory variable. However, the results do not provide a
6
7 compelling numeric contrast between GAMs and WRTDS, despite the alternative statistical
8
9 foundations. Both models are extremely flexible through fine control of window widths for
10
11 WRTDS and degrees of smoothing in GAMS, although at the cost of losing generality with
12
13 increased precision.
14
15
16

17
18 Similarity in results for WRTDS and GAMs may suggest that relationships between time,
19
20 season, and flow in the Patuxent were adequately described by each approach, but
21
22 generalizations of the merits of each model should be made sparingly until additional
23
24 assessments with alternative datasets. Site selection of TF1.6 and LE1.2 was meant to capture a
25
26 gradient of watershed to main stem influences at each location. The known historical changes
27
28 from management practices (e.g, wastewater treatment, banning of phosphorus-based detergents)
29
30 and natural events (e.g., storm events, seagrass recovery) that have affected the Patuxent have
31
32 also provided a unique context for the time series. Additionally, a natural conclusion from this
33
34 study is that both models were equally ‘good’ at trend evaluation, although the possibility that
35
36 both were equally inadequate should also be considered as a potential explanation. Alternative
37
38 drivers of chl-*a* response that were not explicitly included in each model could limit explanatory
39
40 power if time, season, and discharge were not the dominant predictors of production. The
41
42 observation that models capture more of the extreme values at TF1.6 than at LE1.2 (Figure 3a)
43
44 suggests this may be the case at LE1.2. For example, Beck and Hagy III (2015) evaluated
45
46 residual variation of WRTDS models in Tampa Bay, Florida in relation to seagrass growth, El
47
48 Niño effects, and nitrogen inputs. A similar analysis of additional variables at LE1.2 could reveal
49
50 insight into factors other than time, season, or flow that influence chl-*a* in the lower estuary.
51
52
53
54
55
56
57
58
59
60

1
2
3 Evaluation of alternative sites with different historical contexts could provide further information
4
5
6 to support our general conclusion of comparability between methods.
7

8 Although our results generally indicated that comparable information was provided by
9
10 both models, some instances were observed when different information was provided. For
11
12 example, significant differences in the regression comparisons between the models (Table 6 and
13
14 Appendix C) typically had intercept estimates greater than zero and slope estimates less than
15
16 one. This suggests that WRTDS estimates were, on average, larger than GAMs (intercept > 0),
17
18 whereas GAMs fit a wider range of values compared to WRTDS (slope < 1). However, these
19
20 conclusions should be interpreted with caution given the certainty of the results in the context of
21
22 the analysis method. More robust approaches to evaluate systematic biases, in addition to
23
24 alternative datasets, should be used to validate these general conclusions. Generally, important
25
26 differences between the models would be those that would result in a different conclusion if one
27
28 model was used instead of the other. Although none of the model differences were large, several
29
30 differences were observed in the patterns of the flow normalized results (Tables 3 and 4). Most
31
32 notably, the LE1.2 annual percent change results from GAMs suggested that the increase in *ln-chl-a*
33
34 has become less steep over time (9.6 to 3.2%), whereas the WRTDS results suggested the
35
36 increase has become steeper over time (1.75 to 6.07%) (Table 3). The seasonal slopes in Table 3
37
38 for LE1.2 also suggested different patterns from the two models. The increase in *chl-a* was the
39
40 smallest in the summer (JAS) from the GAM results, whereas the WRTDS results suggested that
41
42 the smallest increase over time was in the winter months (JFM). For TF1.6 (Table 4), differences
43
44 in the percent changes were also observed, with the JFM change from WRTDS more than three
45
46 times that suggested by GAMs. These slight differences in patterns showed that the models were
47
48 not identical on the fine-scale. Although we cannot know which model was more accurate in
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 depicting flow-normalized trends in Patuxent chl-*a*, these differences reveal that, in fact, a
4
5 multiple models approach could be beneficial when making conclusions on a fine temporal scale.
6
7

8 Finally, initial assessment of Figure 5 suggested that WRTDS provided a more dynamic
9
10 description of chl-*a* response to changes in flow or salinity for specific locations in the record.
11
12 For example, chl-*a* response over time to salinity changes during January at LE1.2 shows
13
14 WRTDS describing greater variation than GAMs, particularly for lower salinity values.
15
16 Additional investigation suggested that these ‘novel’ descriptions were related to low sample size
17
18 for the specific location in the record causing instability in the model predictions. WRTDS
19
20 descriptions may be unstable at extreme or uncommon locations in the data domain where the
21
22 number of observations with non-zero weights may be limited. Methods for WRTDS have been
23
24 developed to address this issue (i.e., automated window width increases with low sample sizes,
25
26 Hirsch *et al.* 2010), although they were not implemented for the current analysis to simplify
27
28 direct comparisons between models. Practical application of WRTDS for trend analysis should
29
30 use the adaptive window-widening scheme provided by the software (i.e., Hirsch and De Cicco
31
32 2014) to ensure enough observations are available for fitting models at extreme locations in the
33
34 data domain. Similar problems may also be avoided with datasets at smaller time steps (e.g.,
35
36 daily), whereas the nutrient time series represent a more coarse resolution at the bimonthly scale.
37
38
39
40
41
42

43 *Qualitative comparisons*
44

45
46 A quantitative description of the predictive performance and apparent trends described by
47
48 the different methods is an incomplete comparison of the relative abilities of each model. For the
49
50 empirical analysis, both models were compared similarly within the constraints of each method
51
52 to provide a more balanced evaluation (e.g., same datasets, similar optimization methods).
53
54
55 However, the decision for using a specific method may be better informed by considering the
56
57
58
59
60

abilities to accommodate characteristics of a dataset or the type of information that is desired, regardless of performance characteristics. Table 8 provides a qualitative comparison of each method to emphasize differences independent of the empirical measures above.

Ease-of-use for a specific method has importance from an analyst's perspective given constraints on resources or relative skillsets of an individual. Table 8 considers ease-of-use based on computational requirements of each method, interpretation of the statistical basis for a model, availability of software, and tools for visualizing model output. As previously described, GAMs and WRTDS vary significantly in the computational requirements to fit a model. These differences are non-trivial and have direct impacts on how the methods are applied with additional evaluations, e.g., stochastic assessments using bootstrap or Monte Carlo resampling (Efron and Tibshirani 1993, Hirsch *et al.* 2015). Differences in computation time are directly related to statistical differences in parameter estimation for each model. Although the estimation of WRTDS parameters requires more time than a comparable GAM, the underlying math and optimization procedure is simple in comparison. WRTDS is nothing more than a moving window linear regression whereas GAMs use more complex spline-fitting methods. Although there is no objective means to determine which method is 'better' based on complexity alone, the ability to understand the theory of a method is a benefit that will likely have lasting impacts on how results are perceived and applied in decision-making (Carpenter 1995). Moreover, the WRTDS method was developed specifically for trend analysis of water quality and the availability of software, including supporting documentation, far exceeds current resources for GAMs in environmental planning. The WRTDS method for rivers and streams is implemented in the well-documented Exploration and Graphics for RivEr Trends ('EGRET') package (Hirsch and De Cicco 2014) developed by the US Geological Survey using open-source software. A

similar package for tidal waters, including several visualization methods, has also been developed by the authors (see supporting information).

Statistical considerations for each model relate to the products that are provided and the ability to accommodate characteristics of a dataset. As noted above, additional features provided by each model were not directly compared either because such comparisons were impossible (i.e., a feature was unavailable for a method) or they were beyond the analysis scope. For example, WRTDS has been applied using a quantile regression approach to characterize trends at the extreme concentration distributions of the data (Beck and Hagy III 2015). This feature is important for estuaries where the occurrence and magnitude of harmful algal blooms, for example, are often characterized by extreme events as a basis for developing standards (e.g., Schaeffer *et al.* 2013). Although the extension of GAMs to characterize conditional quantiles may be possible (e.g., additive quantile regression, Koenker 2013), comparable applications for water quality analysis have not been developed. An additional concern is the availability of confidence intervals for model estimates that provide direct measures of uncertainty and can facilitate hypothesis-testing. Confidence intervals are readily available from GAMs as standard model output, whereas similar estimates for WRTDS require comprehensive resampling of results with bootstrapping (available as the ‘EGRETci’ package, Hirsch *et al.* 2015). Similar products are not yet available for the tidal adaptation of WRTDS. As such, both methods provide an approach for estimating uncertainty but they differ in implementation that may affect ease of use.

Characteristics of a dataset, questions of interest, and how both can be addressed with WRTDS or GAMs are also important considerations for choosing a method. Water quality data are often characterized by censored observations that are beyond the detection limit of a

1
2
3 monitoring device. With this in mind, WRTDS models were developed using ‘survival analysis’
4
5 as an adaptation of the weighted Tobit model for the original method (Moyer *et al.* 2012, Hirsch
6
7 and De Cicco 2014) and using the Kaplan-Meier approximation for a single-sample survey
8
9 function for conditional regression quantiles in the tidal adaptation (Portnoy 2003, Koenker
10
11 2008). An approach to account for censored data in GAMs is not yet available for water quality
12
13 modeling, although similar methods are feasible and development is anticipated in this area. The
14
15 inclusion of additional variables in a model to describe a response measure may also be a
16
17 concern given the research question. Although both models can theoretically include variables
18
19 other than time, flow, and season, application in GAMs may be much simpler. The ‘mgcv’
20
21 package for GAMs (Wood 2006a) is sufficiently generalizable such that including additional
22
23 variables is a slight modification to the initial function call. Conversely, the available WRTDS
24
25 packages are more specialized and including additional variables would require substantial
26
27 modification. Lastly, sparsity of data including the time step (e.g., continuous, monthly), length
28
29 of the record, or gaps (random or systematic) can affect model performance. For example, Hirsch
30
31 *et al.* (2015) evaluated the effects of sampling intervals and record length on trend comparisons
32
33 between time periods for WRTDS to describe the probability of false positive error rates (Type
34
35 I). Systematic comparisons between WRTDS and GAMs to evaluate effects of data sparsity have
36
37 yet to be done but researchers should be aware of the potential and relatively unknown effects on
38
39 model outcomes.
40
41
42
43
44
45
46
47

48 *Patuxent trends*

49
50 Both models provided a detailed description of water quality changes in the Patuxent
51
52 River estuary. Several trends were described that deserve additional discussion independent of
53
54 the model comparisons. Annual trends at TF1.6 showed a substantial decrease in chl-*a* that lasted
55
56
57
58
59
60

several years, followed by a gradual increase to concentrations similar to those earlier in the time series. By comparison, annual trends in the lower estuary at LE1.2 showed a consistent, linear increase over time. Seasonal patterns and trends related to different flow periods were also described by the models. Spring blooms were commonly observed in the lower estuary, whereas late summer blooms were observed in the upper estuary. Trends related to different flow periods were less obvious, although large increases in chl-*a* were observed for moderate flow levels. Trends in Figure 5 can facilitate an interpretation of changes at each station related to flow effects over time. For example, annual trends in October suggested that the association between flow (decreasing salinity) and chl-*a* have weakened over time at LE1.2. By contrast, trends at TF1.6 showed an increasingly negative relationship between flow and chl-*a* over time. Both models also showed changes in the shape of the relationship between chl-*a* and discharge. For example, a distinct non-linear relationship between chl-*a* and increasing discharge (decreasing salinity) was observed for January predictions at LE1.2 earlier in the record, whereas the trend became more linear near the end of the record. Identifying differences in chl-*a* response at both different flow levels and different seasons could be a first step to identifying influencing factors. The increase over time at LE1.2 is fairly consistent, except for patterns in October at high salinities. Further investigation to reveal what sources are actually being reduced during that period would be insightful.

The results from either model can be used to hypothesize causal links between water quality changes, flow variation, or additional ecosystem characteristics. Previous studies have linked chl-*a* changes and flow relationships to shifts in sources of nutrient pollution (Hirsch *et al.* 2010, Beck and Hagy III 2015). Similarly, historical changes in the Patuxent are likely related to the banning of phosphorus-based detergents in the mid-1980s and wastewater treatment plant

upgrades in the early 1990s (Lung and Bai 2003, Testa *et al.* 2008). An investigation of chl-*a* response to both flow changes and ratios of point-source to non-point sources of nutrients could provide valuable information on system dynamics. Historical changes in flow have also affected water quality in the Patuxent. Flow records for the Patuxent show a drought period from 1999 to 2002 that likely contributed to increases in chl-*a* in the upper estuary and decreases in the lower estuary. By contrast, storm events could be linked to lower chl-*a* from estuarine flushing or shifts in concentration along the longitudinal axis (Hagy *et al.* 2006, Murrell *et al.* 2007). The substantial decline in chl-*a* in the upper estuary in the early 2000s coincides with storm events, including the passage of Hurricane Isabel in 2003. However, low concentrations persisted for several years suggesting additional factors may have had separate or additive effects on chl-*a* response. For example, seagrass growth patterns in the upper estuary have followed a similar but inverse pattern as chl-*a*, with an increase in growth in the late 1990s and early 2000s, followed by a decline in recent years after a peak in coverage in 2005 (J. M. Testa, personal communication). This correlation suggests nutrient sequestration by seagrasses, although definitive links have yet to be shown. Comparison to bay-wide changes for the larger Chesapeake Bay could provide additional explanations, such as the relationship to long-term trends in seagrass growth patterns, additional nutrients, or phytoplankton (Orth *et al.* 2010, Harding *et al.* 2016).

CONCLUSIONS

The use of data-driven statistical techniques that leverage the descriptive potential of long-term monitoring datasets continues to be a relevant research focus in aquatic systems. Both WRTDS and GAMs are actively being developed for application to monitoring time series and our analysis represents the first comparison of WRTDS and GAMs to evaluate trends in tidal

1
2
3 waters. For the Patuxent River estuary, both models had surprisingly similar abilities to describe
4
5 observed and flow-normalized trends in chl-*a*. Some differences in the descriptive capabilities
6
7 were observed, such as specific periods of the time series where data limitations may have
8
9 caused instability in model predictions for WRTDS. Our application to simulated datasets with
10
11 known flow-independent components of chl-*a* provided further indications of similarities
12
13 between the two approaches. Finally, both models had similar abilities to predict observations
14
15 with validation datasets, having implications for the use of either model with missing data.
16
17
18
19

20 We emphasize that simple comparisons of predictive performance with error measures
21
22 provide relatively narrow descriptions of the quantitative abilities of each model. These
23
24 comparison methods were chosen based on the exploratory needs of the analysis and by
25
26 considering that each technique provides a potentially novel approach to trend assessment.
27
28 Inferior performance for one metric does not invalidate an analysis method for all applications
29
30 and alternative comparisons are needed for more specific uses of each method. This analysis was
31
32 the first to rigorously compare both WRTDS and GAMs and further evaluations with alternative
33
34 datasets should be made to compare with our results. Although both models provided similar
35
36 information, the results from either reveal interesting relationships (e.g., flow, nutrient response
37
38 over time, Figure 5) that can lead to additional hypotheses or analysis to investigate ecosystem
39
40 dynamics.
41
42
43
44
45

46 Practical applications of each model should consider alternative characteristics of each
47
48 technique, in addition to the simple quantitative comparisons described above. The use of
49
50 WRTDS to describe water quality trends in tidal waters, particularly with monthly or bimonthly
51
52 time series, is a novel application for which the model was never intended. Hirsch *et al.* (2010)
53
54 developed the original model for streams and rivers using high-resolution, daily time series
55
56
57
58
59
60

where time, discharge, and season are dominant characteristics that influence water quality. Although seasonal and flow effects are important drivers of change in estuaries, other physical or biological characteristics may be equally or more important. For example, the extreme ends of the chl-*a* distribution at LE1.2 were not fit well by either model as compared to TF1.6, which suggests additional predictors besides time, discharge, and season may better describe variation in the lower estuary. As such, recent use of GAMs in tidal waters has followed an alternative paradigm where drivers of change are not necessarily known and the time series may have a larger time step with occasional discontinuous intervals (E. S. Perry, personal communication, Harding *et al.* 2016). Although we have quantitatively compared each method to inform decision-making, choosing a technique should also consider alternative products, characteristics of the dataset, questions of interest, and specifics of the study system.

Appendix A: The WRTDStidal R package for implementing the tidal adaptation of WRTDS is available for download at <https://github.com/fawda123/WRTDStidal>

SUPPORTING INFORMATION

Additional supporting information may be found online under the Supporting Information tab for this article: **Appendix B:** Additional material describing the simulation of daily discharge and chl-*a* time series. **Appendix C:** Supplementary figure of regression comparisons between WRTDS and GAMs.

ACKNOWLEDGMENTS

We thank Jim Hagy, Bob Hirsch, and Jennifer Keisman for valuable discussions that improved the analysis. Thanks to Jeremy Testa, Jeffrey Chanut, and two anonymous reviewers for providing valuable comments on an earlier draft, and Elgin Perry for aid in implementing GAMs. We thank the Chesapeake Bay Program and Maryland Department of Natural Resources

for providing data.

LITERATURE CITED

Arlot, S. and A. Celisse, 2010. A survey of cross-validation procedures for model selection. *Statistics Surveys* 4:40–79.

Beck, M.W. and J.D. Hagy III, 2015. Adaptation of a weighted regression approach to evaluate water quality trends in an estuary. *Environmental Modelling and Assessment* 20(6):637–655.

Beck, M.W., J.D. Hagy III, and M.C. Murrell, 2015. Improving estimates of ecosystem metabolism by reducing effects of tidal advection on dissolved oxygen time series. *Limnology and Oceanography: Methods* 13(12):731–745.

Boicourt, W.C. and L.P. Sanford, 1998. A hydrodynamic study of the Patuxent River estuary, Technical report, Maryland Department of the Environment, Baltimore, Maryland.

Borsuk, M.E., C.A. Stow, and K.H. Reckhow, 2004. Confounding effect of flow on estuarine response to nitrogen loading. *Journal of Environmental Engineering-ASCE* 130(6):605–614.

Byrd, R.H., P. Lu, J. Nocedal, and C. Zhu, 1995. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing* 16(5):1190–1208.

Cade, B.S. and B.R. Noon, 2003. A gentle introduction to quantile regression for ecologists. *Frontiers in Ecology and the Environment* 1(8):412–420.

Carpenter, D. O., 1995. Communicating with the public on issues of science and public health. *Environmental Health Perspectives* 103(S6):127–130.

Cleveland, W. S., 1979. Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association* 74(368):829–836.

Cloern, J.E., 1996. Phytoplankton bloom dynamics in coastal ecosystems: A review with some general lessons from sustained investigation of San Francisco Bay, California. *Review of Geophysics* 34(2):127–168.

Cloern, J.E. and A.D. Jassby, 2010. Patterns and scales of phytoplankton variability in estuarine-coastal ecosystems. *Estuaries and Coasts* 33(2):230–241.

Cronin, W.B. and D.W. Pritchard, 1975. Additional statistics on the dimensions of the Chesapeake Bay and its tributaries: Cross-section widths and segment volumes per meter depth, Technical Report 42, Reference 75-3, Chesapeake Bay Institute, The Johns Hopkins University, Baltimore, Maryland.

- Efron, B. and R. Tibshirani, 1993. An Introduction to the Bootstrap, Chapman and Hall, New York, first edition.
- Hagy, J.D., 1996. Residence times and net ecosystem processes in Patuxent River estuary, Master's thesis, University of Maryland, College Park, Maryland.
- Hagy, J.D., J.C. Lehrter, and M.C. Murrell, 2006. Effects of hurricane Ivan on water quality in Pensacola Bay, Florida. *Estuaries and Coasts* 29(6A):919–925.
- Haraguchi, L., J. Carstensen, P.C. Abreu, and C. Odebrecht, 2015. Long-term changes of the phytoplankton community and biomass in the subtropical shallow Patos Lagoon Estuary, Brazil. *Estuarine, Coastal and Shelf Science* 162(SI):76–87.
- Harding, L.W., 1994. Long-term trends in the distribution of phytoplankton in Chesapeake Bay - roles of light, nutrients, and streamflow. *Marine Ecology Progress Series* 104:267–291.
- Harding, L.W., C.L. Gallegos, E.S. Perry, W.D. Miller, J.E. Adolf, M.E. Mallonee, and H.W. Paerl, 2016. Long-term trends of nutrients and phytoplankton in Chesapeake Bay. *Estuaries and Coasts* 39:664–681.
- Hastie, T. and R. Tibshirani, 1990. Generalized Additive Models, Chapman and Hall, London, New York.
- Hirsch, R.M., 2014. Large biases in regression-based constituent flux estimates: causes and diagnostic tools. *Journal of the American Water Resources Association* 50(6):1401–1424.
- Hirsch, R.M., S.A. Archfield, and L.A. De Cicco, 2015. A bootstrap method for estimating uncertainty of water quality trends. *Environmental Modelling and Software* 73:148–166.
- Hirsch, R.M. and L. De Cicco, 2014. User guide to Exploration and Graphics for RivEr Trends (EGRET) and dataRetrieval: R packages for hydrologic data, Techniques and Methods book 4, ch. A10, US Geological Survey, Reston, Virginia.
<http://pubs.usgs.gov/tm/04/a10/>.
- Hirsch, R.M., D.L. Moyer, and S.A. Archfield, 2010. Weighted regressions on time, discharge, and season (WRTDS), with an application to Chesapeake Bay river inputs. *Journal of the American Water Resources Association* 46(5):857–880.
- Koenker, R., 2008. Censored quantile regression redux. *Journal of Statistical Software* 27(6):1–25.
- , 2013. quantreg: Quantile Regression R package 5.05. <http://CRAN.R-project.org/package=quantreg>.
- Lung, W. and S. Bai, 2003. A water quality model for the Patuxent estuary: Current conditions and predictions under changing land-use scenarios. *Estuaries* 26(2A):267–279.

Medalie, L., R.M. Hirsch, and S.A. Archfield, 2012. Use of flow-normalization to evaluate nutrient concentration and flux changes in Lake Champlain tributaries, 1990-2009. *Journal of Great Lakes Research* 38(SI):58–67.

Monbet, Y., 1992. Control of phytoplankton biomass in estuaries: A comparative analysis of microtidal and macrotidal estuaries. *Estuaries* 15(4):563–571.

Moyer, D.L., R.M. Hirsch, and K.E. Hyer, 2012. Comparison of two regression-based approaches for determining nutrient and sediment fluxes and trends in the Chesapeake Bay Watershed, Scientific Investigations Report 2012-544, US Geological Survey, US Department of the Interior, Reston, Virginia.

Murrell, M.C., J.D. Hagy, E.M. Lores, and R.M. Greene, 2007. Phytoplankton production and nutrient distributions in a subtropical esuary: Importance of freshwater flow. *Estuaries and Coasts* 30(3):390–402.

Nixon, S.W., 1995. Coastal marine eutrophication: A definition, social causes, and future concerns. *Ophelia* 41:199–219.

Nocedal, J. and S.J. Wright, 2006. Numerical Optimization, Springer-Verlag, New York, New York, 2nd edition.

Orth, R.J., M.R. Williams, S.R. Marion, D.J. Wilcox, T.J.B. Carruthers, K.A. Moore, W.M. Kemp, W.C. Dennison, N. Rybicki, P. Bergstrom, and R.A. Batiuk, 2010. Long-term trends in submersed aquatic vegetation (SAV) in Chesapeake Bay, USA, related to water quality. *Estuaries and Coasts* 33(5):1144–1163.

Paerl, H.W., 2006. Assessing and managing nutrient-enhanced eutrophication in estuarine and coastal waters: Interactive effects of human and climatic perturbations. *Ecological Engineering* 26(1):40–54.

Paerl, H.W., N.S. Hall, B.L. Peierls, and K.L. Rossignol, 2014. Evolving paradigms and challenges in estuarine and coastal eutrophication dynamics in a culturally and climatically stressed world. *Estuaries and Coasts* 37(2):243–258.

Portnoy, S., 2003. Censored regression quantiles. *Journal of the American Statistical Association* 98(464):1001–1012.

Powers, S.P., C.H. Peterson, R.R. Christian, E. Sullivan, M.J. Powers, M.J. Bishop, and C.P. Buzzelli, 2005. Effects of eutrophication on bottom habitat and prey resources of demersal fishes. *Marine Ecology Progress Series* 302:233–243.

RDCT (R Development Core Team), 2015. R: A language and environment for statistical computing, v3.2.0. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org>.

- Schaeffer, B. A., J. D. Hagy, III, and R. P. Stumpf, 2013. Approach to developing numeric water quality criteria for coastal waters: transition from SeaWiFS to MODIS and MERIS satellites. *Journal of Applied Remote Sensing* 7(1):073544.
- Sprague, L.A., R.M. Hirsch, and B.T. Aulenbach, 2011. Nitrate in the Mississippi River and its tributaries, 1980 to 2008: Are we making progress? *Environmental Science and Technology* 45(17):7209–7216.
- Steward, J.S. and W.C. Green, 2007. Setting load limits for nutrients and suspended solids based upon seagrass depth-limit targets. *Estuaries and Coasts* 30(4):657–670.
- TBEP (Tampa Bay Estuary Program), 2011. Tampa Bay Water Atlas. <http://www.tampabay.wateratlas.usf.edu/>. (Accessed October, 2013).
- Testa, J.M., W.M. Kemp, W.R. Boynton, and J.D. Hagy, 2008. Long-term changes in water quality and productivity in the Patuxent River Estuary: 1985 to 2003. *Estuaries and Coasts* 31(6):1021–1037.
- USEPA (U.S. Environmental Protection Agency), 2010. Chesapeake Bay Phase 5.3 Community Watershed Model. EPA 903S10002 - CBP/TRS-303-10. U.S. Environmental Protection Agency, Chesapeake Bay Program Office, Annapolis MD. December 2010.
- USGS (US Geological Survey), 2015. Water Quality Loads and Trends at Nontidal Monitoring Stations in the Chesapeake Bay Watershed. <http://cbrim.er.usgs.gov/>. (Accessed November, 2015)
- Wood, S.N., 2006a. Generalized Additive Models: An Introduction with R, Chapman and Hall, CRC Press, London, United Kingdom.
- , 2006b. Low-rank scale-invariant tensor product smooths for generalized additive mixed models. *Biometrics* 62(4):1025–1036.
- Zhang, Q., D.C. Brady, W.R. Boynton, and W.P. Ball, 2015. Long-Term Trends of Nutrients and Sediment from the Nontidal Chesapeake Watershed: An Assessment of Progress by River and Season. *Journal of the American Water Resources Association* 51(6):1534-1555.

TABLE 1: Summary characteristics of monitoring stations on the Patuxent River estuary. Chlorophyll-*a* and salinity values are based on averages from 1986 to 2014. Stations used for the analysis are in bold. Segments are salinity regions in the Patuxent for the larger Chesapeake Bay area (TF = tidal fresh, OH = oligohaline, MH = mesohaline).

Station	Lat	Long	Segment	Distance (km)	Depth (m)	ln-Chl ($\mu\text{g/L}$)	Sal (ppt)
TF1.3	38.81	-76.71	TF	74.90	2.90	1.52	0.00
TF1.4	38.77	-76.71	TF	69.50	2.00	2.31	0.02
TF1.5	38.71	-76.70	TF	60.30	10.60	2.88	0.27
TF1.6	38.66	-76.68	OH	52.20	6.20	2.44	0.90
TF1.7	38.58	-76.68	OH	42.50	3.00	2.09	4.09
RET1.1	38.49	-76.66	MH	32.20	11.20	2.47	10.25
LE1.1	38.43	-76.60	MH	22.90	12.10	2.31	12.04
LE1.2	38.38	-76.51	MH	13.40	17.10	2.16	12.73
LE1.3	38.34	-76.48	MH	8.30	23.40	2.12	12.89
LE1.4	38.31	-76.42	MH	0.00	15.40	2.21	13.46

TABLE 2: Summaries of model performance using RMSE (deviance in parentheses) of observed to predicted ln-chl-*a* for each station (LE1.2 and TF1.6). Overall performance for the entire time series is shown at the top with groupings by different time periods below. Time periods are annual groupings every seven years (top), seasonal groupings (middle), and flow periods based on quantile distributions of discharge (bottom).

Period	LE1.2		TF1.6	
	GAM	WRTDS	GAM	WRTDS
All	0.51 (139.5)	0.51 (135.1)	0.50 (128.4)	0.52 (138.6)
Annual				
1986-1993	0.50 (41.1)	0.50 (40.9)	0.48 (37.2)	0.49 (39.1)
1994-2000	0.51 (34.7)	0.50 (33.2)	0.55 (39.3)	0.58 (44.9)
2001-2007	0.61 (51.5)	0.60 (49.6)	0.50 (33.7)	0.53 (37.5)
2008-2014	0.37 (12.1)	0.36 (11.4)	0.45 (18.2)	0.44 (17.1)
Seasonal				
JFM	0.60 (38.1)	0.58 (35.3)	0.49 (24.4)	0.49 (23.8)
AMJ	0.64 (65.2)	0.64 (65.3)	0.54 (45.7)	0.58 (51.9)
JAS	0.35 (19.3)	0.35 (18.6)	0.45 (30.4)	0.46 (32.2)
OND	0.39 (16.8)	0.38 (15.9)	0.52 (27.9)	0.54 (30.7)
Flow				
1 (Low)	0.36 (17.4)	0.36 (16.7)	0.45 (26.5)	0.46 (27.7)
2	0.43 (24.4)	0.42 (23.5)	0.53 (36.6)	0.54 (37.8)
3	0.58 (43.8)	0.57 (42.9)	0.49 (31.3)	0.52 (35.4)
4 (High)	0.64 (53.9)	0.63 (52.0)	0.51 (34.0)	0.54 (37.7)

TABLE 3: Summaries of flow-normalized trends from each model at LE1.2 for different time periods. Summaries are averages and percentage changes of $\ln\text{-chl-}a$ ($\mu\text{g/L}$) based on annual means within each category. Percentage changes are the differences between the last and first years in the periods. Time periods are annual groupings every seven years (top), seasonal groupings (middle), and flow periods based on quantile distributions of discharge (bottom).

Period	GAM		WRTDS	
	Ave.	% Change	Ave.	% Change
All	2.17	24.28	2.18	18.85
Annual				
1986-1993	1.99	9.60	2.03	1.75
1994-2000	2.12	5.49	2.12	5.50
2001-2007	2.24	5.50	2.24	5.35
2008-2014	2.37	3.20	2.37	6.07
Seasonal				
JFM	2.57	20.06	2.58	14.04
AMJ	2.32	31.20	2.33	22.47
JAS	2.01	18.48	2.01	19.91
OND	1.82	25.29	1.83	15.14
Flow				
1 (Low)	1.90	20.86	1.93	16.77
2	2.10	13.71	2.11	7.73
3	2.28	15.66	2.29	9.24
4 (High)	2.34	25.09	2.33	22.29

TABLE 4: Summaries of flow-normalized trends from each model at TF1.6 for different time periods. Summaries are averages and percentage changes of ln-chl-*a* ($\mu\text{g/L}$) based on annual means within each category. Percentage changes are the differences between the last and first years in the periods. Time periods are annual groupings every seven years (top), seasonal groupings (middle), and flow periods based on quantile distributions of discharge (bottom).

Period	GAM		WRTDS	
	Ave.	% Change	Ave.	% Change
All	2.43	-4.81	2.44	-2.28
Annual				
1986-1993	2.62	-4.93	2.60	-3.06
1994-2000	2.69	-5.05	2.65	-3.55
2001-2007	2.15	-22.42	2.19	-21.51
2008-2014	2.24	47.10	2.30	38.35
Seasonal				
JFM	1.52	9.03	1.48	32.72
AMJ	2.63	5.47	2.62	5.14
JAS	3.06	0.04	3.08	0.79
OND	2.17	-18.16	2.20	-17.55
Flow				
1 (Low)	2.89	-4.78	2.93	-0.42
2	2.41	16.71	2.43	20.31
3	2.28	6.53	2.27	15.20
4 (High)	2.22	-11.58	2.21	-11.27

TABLE 5: Comparison of predicted results between WRTDS and GAMs using average differences (%) and RMSD. Overall comparisons for the entire time series are shown at the top with groupings by different time periods below. Time periods are annual groupings every seven years (top), seasonal groupings (middle), and flow periods based on quantile distributions of discharge (bottom). Negative percentages indicate WRTDS predictions were lower than GAM predictions (eq. (4)).

Period	LE1.2		TF1.6	
	Ave. diff.	RMSD	Ave. diff.	RMSD
All	-0.11	0.09	0.01	0.13
Annual				
1986-1993	0.20	0.10	-0.75	0.11
1994-2000	0.34	0.09	-1.29	0.15
2001-2007	-0.55	0.07	0.68	0.13
2008-2014	-0.53	0.08	3.06	0.14
Seasonal				
JFM	0.39	0.12	-2.02	0.14
AMJ	0.22	0.10	-0.66	0.14
JAS	-0.71	0.06	0.76	0.10
OND	-0.46	0.05	1.03	0.15
Flow				
1 (Low)	-0.27	0.07	-0.15	0.10
2	-0.14	0.09	0.70	0.13
3	0.48	0.11	1.06	0.14
4 (High)	-0.53	0.09	-1.77	0.15

TABLE 6: Regression fits comparing predicted (*pred*) and flow-normalized (*norm*) results for WRTDS and GAMs Values in bold-italic are those where the intercept (β_0) estimate was significantly different from zero or the slope (β_1) estimate was significantly different from one. Fits for the entire time series are shown at the top. Time periods are annual groupings every seven years (top), seasonal groupings (middle), and flow periods based on quantile distributions of discharge (bottom).

Period	LE1.2		TF1.6		LE1.2		TF1.6	
	$\beta_{0, pred}$	$\beta_{1, pred}$	$\beta_{0, pred}$	$\beta_{1, pred}$	$\beta_{0, norm}$	$\beta_{1, norm}$	$\beta_{0, norm}$	$\beta_{1, norm}$
All	<i>0.05</i>	<i>0.97</i>	<i>0.08</i>	<i>0.97</i>	<i>0.15</i>	<i>0.94</i>	0.02	0.99
Annual								
1986-1993	0.02	0.99	-0.02	1.00	<i>0.20</i>	<i>0.92</i>	<i>-0.12</i>	<i>1.03</i>
1994-2000	<i>0.16</i>	<i>0.93</i>	-0.03	0.99	<i>0.17</i>	<i>0.92</i>	<i>-0.12</i>	<i>1.02</i>
2001-2007	0.02	0.99	<i>0.13</i>	<i>0.95</i>	<i>0.06</i>	<i>0.98</i>	<i>0.11</i>	<i>0.97</i>
2008-2014	0.00	1.00	<i>0.12</i>	0.97	0.01	0.99	<i>0.08</i>	0.99
Seasonal								
JFM	-0.01	1.01	0.09	<i>0.92</i>	0.01	1.00	<i>0.20</i>	<i>0.84</i>
AMJ	<i>0.28</i>	<i>0.88</i>	<i>0.27</i>	<i>0.89</i>	<i>0.38</i>	<i>0.84</i>	<i>0.34</i>	<i>0.87</i>
JAS	-0.08	1.03	<i>0.34</i>	<i>0.89</i>	<i>0.30</i>	<i>0.85</i>	<i>0.39</i>	<i>0.88</i>
OND	0.02	0.98	<i>0.13</i>	<i>0.95</i>	<i>0.38</i>	<i>0.80</i>	0.03	1.00
Flow								
1 (Low)	<i>0.14</i>	<i>0.92</i>	-0.03	1.01	<i>0.46</i>	<i>0.77</i>	<i>0.16</i>	<i>0.95</i>
2	0.00	1.00	<i>0.12</i>	<i>0.96</i>	<i>0.14</i>	<i>0.94</i>	0.01	1.00
3	0.09	0.96	<i>0.21</i>	<i>0.91</i>	<i>0.12</i>	<i>0.96</i>	-0.02	1.00
4 (High)	0.09	<i>0.96</i>	0.03	<i>0.97</i>	<i>0.09</i>	<i>0.96</i>	<i>0.09</i>	<i>0.95</i>

TABLE 7: Summaries of model performance comparing observed chl-*a* with predicted values ($Chl_{obs} \sim \hat{Chl}_{obs}$) and biological chl-*a* with flow-normalized values ($Chl_{bio} \sim \hat{Chl}_{bio}$) for the three simulated time series (no flow, constant flow, and increasing flow effect). Summaries are RMSE values (deviance in parentheses) comparing results from each model (GAM, WRTDS).

Simulations	$Chl_{obs} \sim \hat{Chl}_{obs}$	$Chl_{bio} \sim \hat{Chl}_{bio}$
No flow		
GAM	0.51 (31.2)	0.53 (33.2)
WRTDS	0.50 (29.4)	0.52 (31.7)
Constant flow		
GAM	0.51 (31.2)	0.58 (39.8)
WRTDS	0.53 (32.8)	0.57 (38.9)
Increasing flow		
GAM	0.51 (31.2)	0.54 (35.0)
WRTDS	0.50 (29.7)	0.52 (31.9)

TABLE 8: Qualitative comparisons of generalized additive models and WRTDS. Qualities are grouped by ease of use and statistical considerations. Ease of use qualities are described as good, moderate, or poor and statistical qualities as yes/no.

Quality	GAM	WRTDS
ease of use		
computational requirements	good	poor
interpretation ¹	poor	moderate
software and documentation ²	moderate	good
visualization	moderate	good
statistical		
additional variables	y	n
censored data	n	y
confidence intervals	y	y
quantile fits	n	y

¹Relates to statistical foundation

²In reference to analysis of water quality trends

FIGURE 1: Patuxent River estuary with Chesapeake Bay inset. Locations monitored by the Maryland Department of Natural Resources are shown along the longitudinal axis with distance from the mouth (km). Study sites are in bold. Salinity regions in the Patuxent for the larger Chesapeake Bay area are also shown (TF = tidal fresh, OH = oligohaline, MH = mesohaline). See Table 1 for a numeric summary of station characteristics.

FIGURE 2: Annual, seasonal, and flow differences in chl-*a* trends at each monitoring station in the Patuxent River Estuary. Size and color are proportional medians of ln-chl-*a* by year, season, and flow categories. See Figure 1 for station numbers.

FIGURE 3: Predicted chl-*a* from generalized additive models (GAM) and weighted regression (WRTDS) for stations LE1.2 and TF1.6. Figure 3a shows results at monthly time steps and Figure 3b shows results averaged by year. Values in blue are model predictions and values in yellow are flow-normalized predictions.

FIGURE 4: Seasonal variation from observed and model predictions of chl-*a* by station. Predictions are points by day of year from 1986 to 2014. The blue line is a loess (locally estimated) polynomial smooth to characterize the seasonal components.

FIGURE 5: Changes in the relationship between chl-*a* and freshwater inputs (salinity decrease, flow increase) across the time series. Separate panels are shown for each station (LE1.2, TF1.6), model type (GAM, WRTDS), and chosen months. Changes over time are shown as different predictions for each year in the time series (1986 to 2014). Salinity was used as a tracer of freshwater inputs at LE1.2, whereas the flow record at Bowie, Maryland was used at TF1.6. Axes for salinity and flow are reversed for comparison. Units are proportions of the total range in the observed data.

FIGURE 6: Examples of changing relationships between chl-*a* (ln-transformed, $\mu\text{g L}^{-1}$) and flow (ln-transformed, $\text{m}^3 \text{s}^{-1}$) over time (2005–2015) for each simulated time series (Appendix B). The plots are based on August predictions from the three WRTDS models and GAMs for each time series.

FIGURE 7: Prediction errors of GAMs and WRTDS for different training and validation datasets. Datasets were created from a weekly simulated time series (Appendix B) using different split ratios (5% validation and 95% training, 10% validation and 90% training, etc.) and sampling methods of the complete dataset. Sampling to create the validation datasets varied from completely random (no blocks) to block sizes of different percentages of the total sample size (10, 50, 100%) required for the split ratio of the training and validation data. RMSE values are summarized as the median, 5th, and 95th percentile of model results for 100 repetitions of each dataset type.

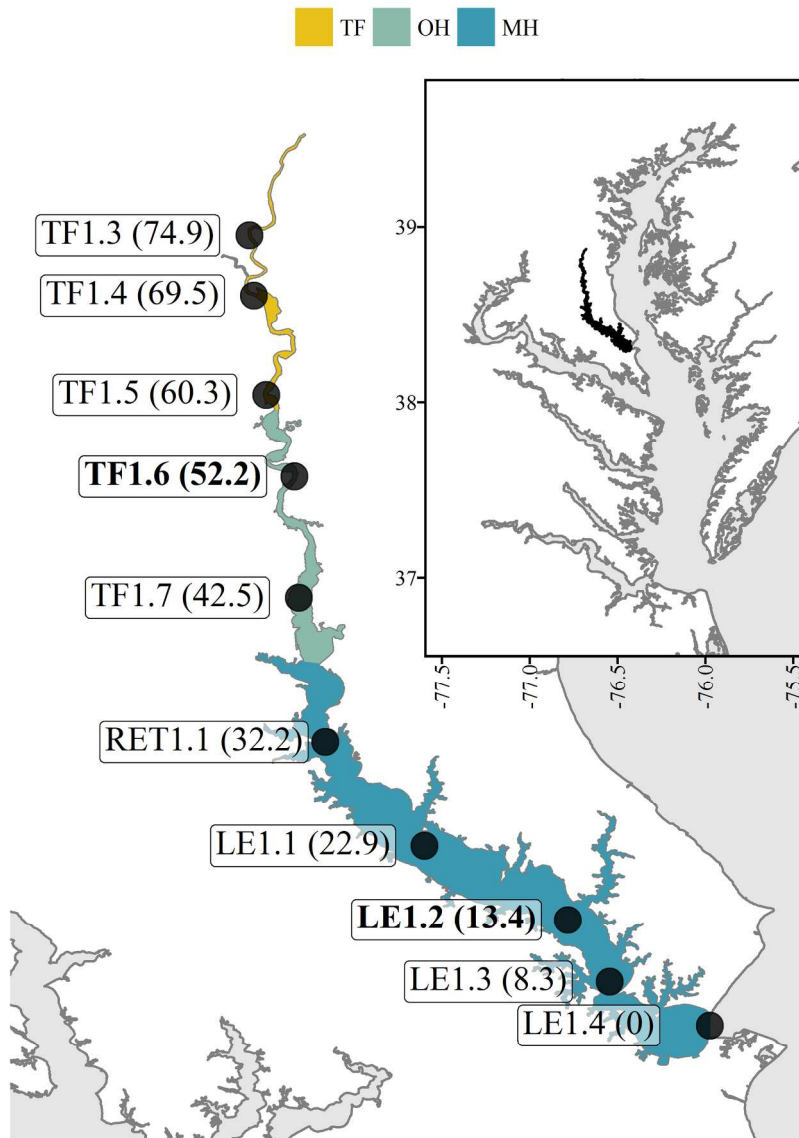


FIGURE 1: Patuxent River estuary with Chesapeake Bay inset. Locations monitored by the Maryland Department of Natural Resources are shown along the longitudinal axis with distance from the mouth (km). Study sites are in bold. Salinity regions in the Patuxent for the larger Chesapeake Bay area are also shown (TF = tidal fresh, OH = oligohaline, MH = mesohaline). See Table 1 for a numeric summary of station characteristics.

177x248mm (300 x 300 DPI)

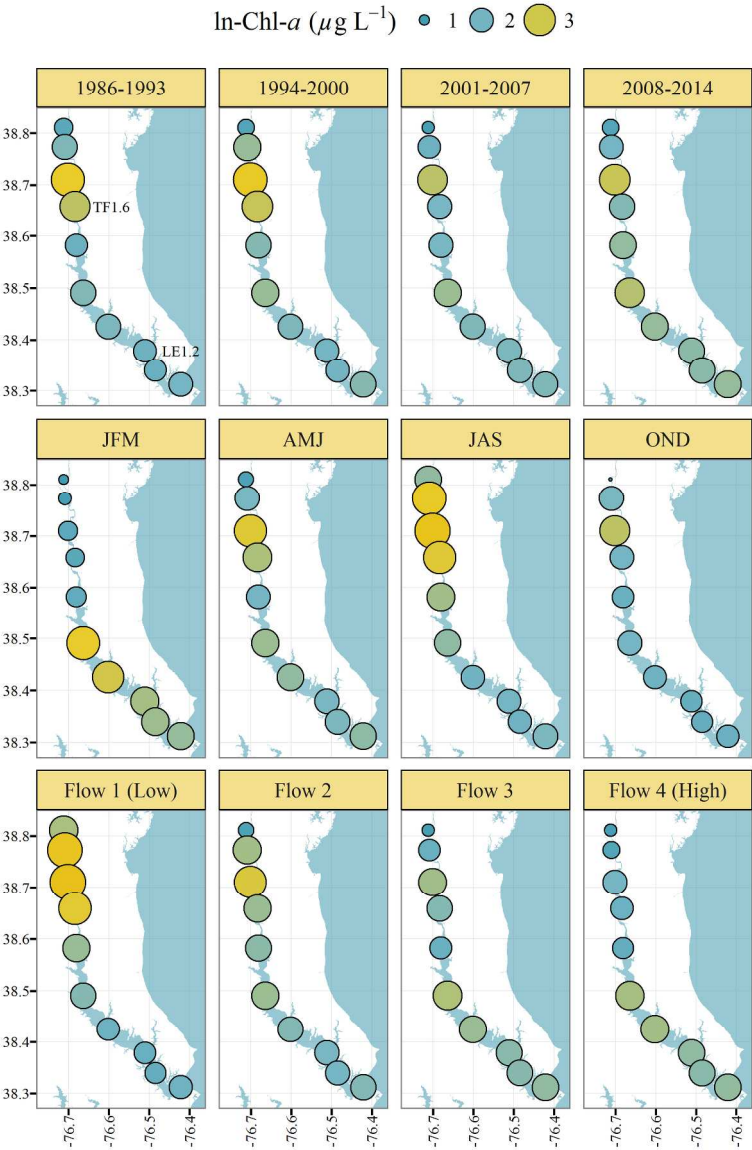


FIGURE 2: Annual, seasonal, and flow differences in chl-*a* trends at each monitoring station in the Patuxent River Estuary. Size and color are proportional medians of ln-chl-*a* by year, season, and flow categories. See Figure 1 for station numbers.

196x305mm (300 x 300 DPI)

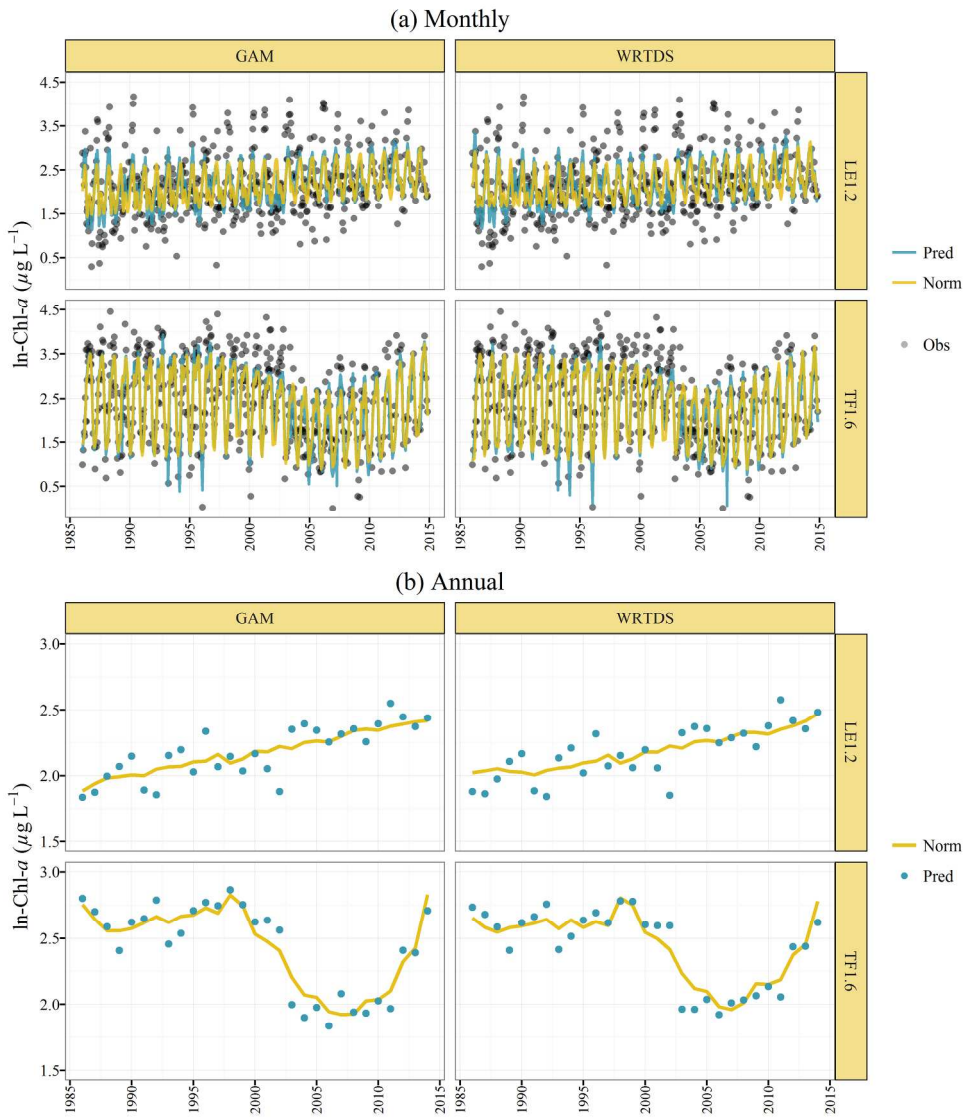


FIGURE 3: Predicted chl-a from generalized additive models (GAM) and weighted regression (WRTDS) for stations LE1.2 and TF1.6. Figure 3a shows results at monthly time steps and Figure 3b shows results averaged by year. Values in blue are model predictions and values in yellow are flow-normalized predictions.

228x257mm (300 x 300 DPI)

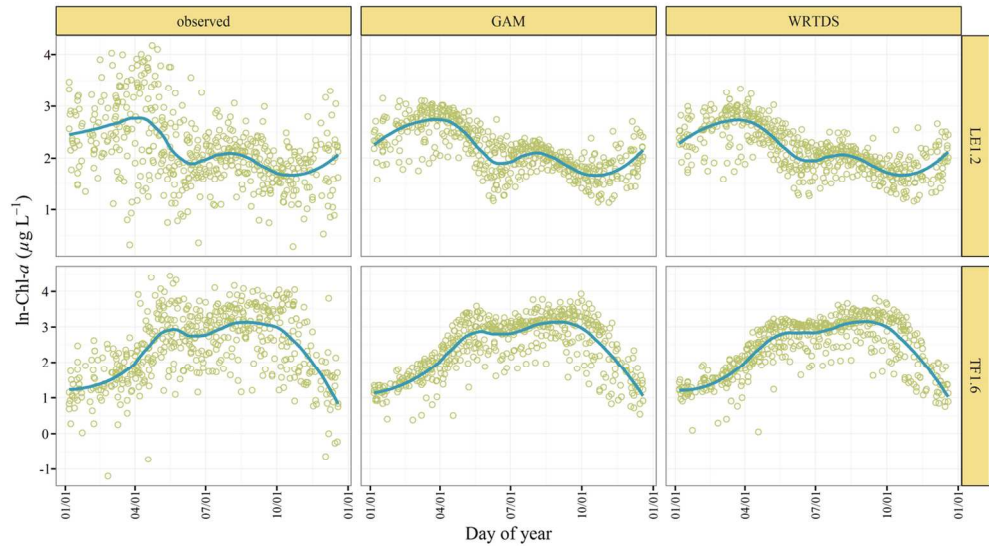


FIGURE 4: Seasonal variation from observed and model predictions of chl-a by station. Predictions are points by day of year from 1986 to 2014. The blue line is a loess (locally estimated) polynomial smooth to characterize the seasonal components.

127x70mm (300 x 300 DPI)

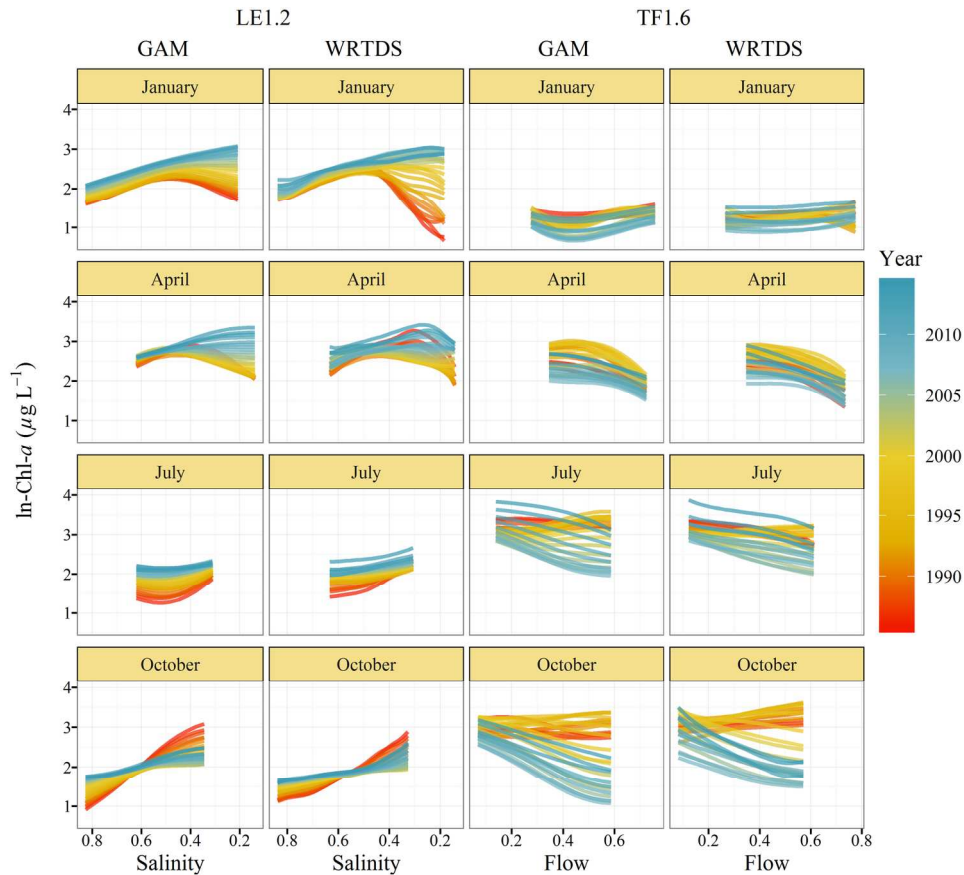


FIGURE 5: Changes in the relationship between chl-a and freshwater inputs (salinity decrease, flow increase) across the time series. Separate panels are shown for each station (LE1.2, TF1.6), model type (GAM, WRTDS), and chosen months. Changes over time are shown as different predictions for each year in the time series (1986 to 2014). Salinity was used as a tracer of freshwater inputs at LE1.2, whereas the flow record at Bowie, Maryland was used at TF1.6. Axes for salinity and flow are reversed for comparison. Units are proportions of the total range in the observed data.

165x143mm (300 x 300 DPI)

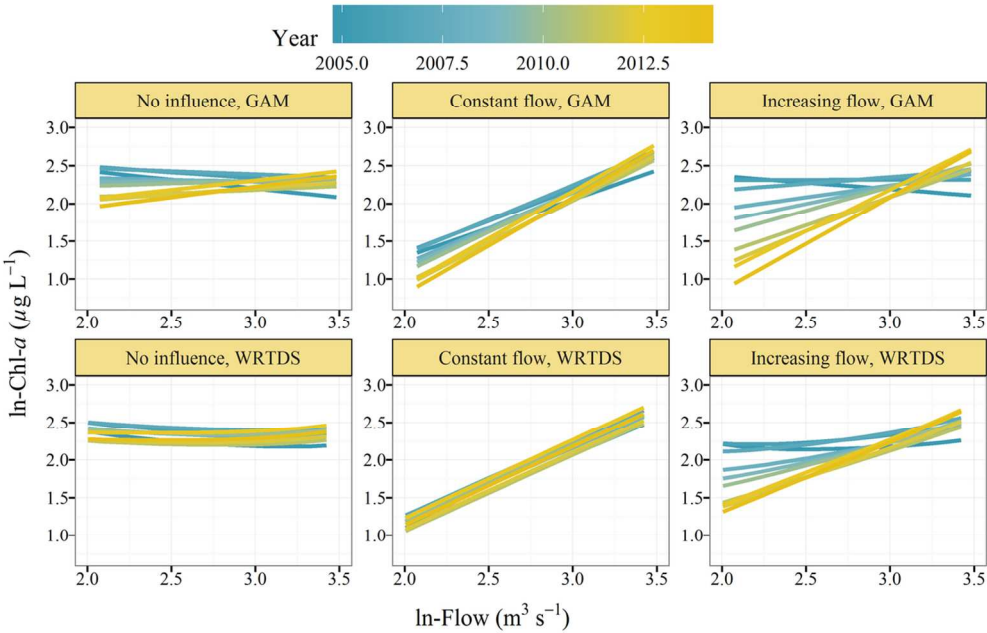


FIGURE 6: Examples of changing relationships between chl-a (\ln -transformed, $\mu\text{g L}^{-1}$) and flow (\ln -transformed, $\text{m}^3 \text{s}^{-1}$) over time (2005–2015) for each simulated time series (Appendix B). The plots are based on August predictions from the three WRTDS models and GAMs for each time series.

114x73mm (300 x 300 DPI)

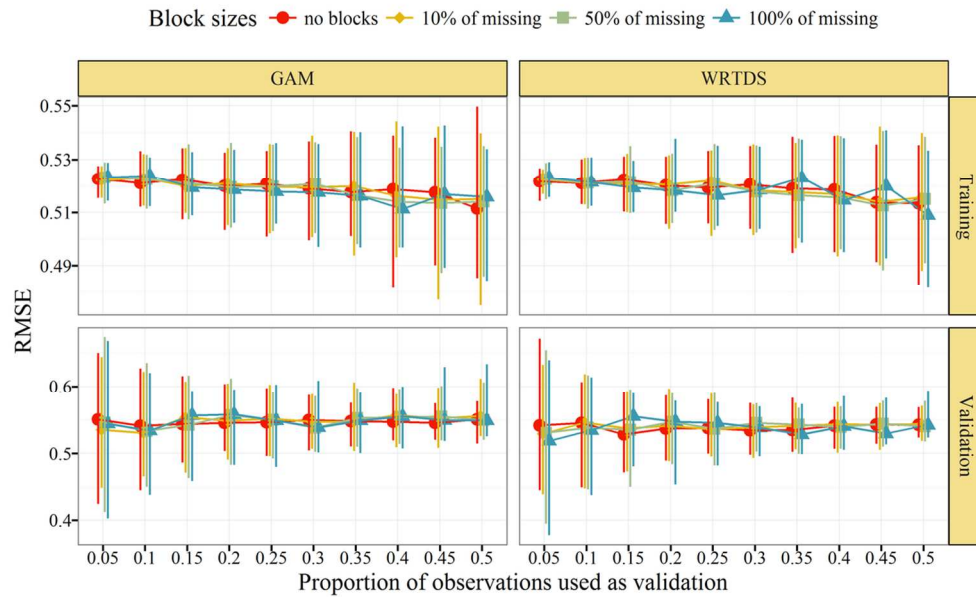


FIGURE 7: Prediction errors of GAMs and WRTDS for different training and validation datasets. Datasets were created from a weekly simulated time series (Appendix B) using different split ratios (5% validation and 95% training, 10% validation and 90% training, etc.) and sampling methods of the complete dataset. Sampling to create the validation datasets varied from completely random (no blocks) to block sizes of different percentages of the total sample size (10, 50, 100%) required for the split ratio of the training and validation data. RMSE values are summarized as the median, 5th, and 95th percentile of model results for 100 repetitions of each dataset type.

114x73mm (300 x 300 DPI)