

1 **Comparison of weighted regression and additive models for**
2 **trend evaluation of water quality in tidal waters**

3 **Marcus W. Beck¹, Rebecca Murphy²**

¹ORISE Research Participation Program

USEPA National Health and Environmental Effects Research Laboratory

Gulf Ecology Division, 1 Sabine Island Drive, Gulf Breeze, FL 32561

Phone: 850-934-2480, Fax: 850-934-2401, Email: beck.marcus@epa.gov

²UMCES at Chesapeake Bay Program

410 Severn Avenue, Suite 112, Annapolis, MD 21403

Phone: 410-267-9837, Fax: 410-267-5777, Email: rmurphy@chesapeakebay.net

Version Date: Tue Oct 27 17:04:29 2015 -0500

Abstract

Long-term monitoring datasets can provide information to interpret the effects of environmental changes or management actions on ecosystem condition. The ability to link causal effects with potential changes from observed data can partly depend on the chosen method of trend analysis.

Two statistical approaches, weighted regression on time, discharge, and season (WRTDS) and generalized additive models (GAMs), have recently been used to evaluate long-term trends in chlorophyll time series in estuarine systems. Both models provide a similar approach to trend analysis by using context-dependent parameters or smoothing functions that vary continuously and have the potential to identify multiple drivers of change. However, the quantitative

capabilities of each model, including descriptions of observed and flow-normalized trends, have not been rigorously compared to determine most appropriate use of each model. We evaluated WRTDS and GAMs using thirty years of data for a monthly time series of chlorophyll in the

Patuxent River Estuary, a well-studied tributary to Chesapeake Bay. Each model was evaluated based on predictive capabilities of the observed data and ability to reproduce flow-normalized trends with simulated data that had statistical properties comparable to the original dataset. Each

model was also evaluated based on concordance of conclusions of water quality changes, and causes thereof, in different time periods. For all examples, WRTDS had lower prediction error than GAMs, although differences were minor compared to computational requirements of each

model. Comparisons of flow-normalized trends revealed distinct differences in temporal variation in chlorophyll *a* (chl-*a*) from the upper to lower Patuxent estuary. Mainstem influences of the Chesapeake Bay were apparent with a slight increase in chl-*a* trends over time in the lower

estuary, whereas flow-normalized predictions for the upper estuary showed declines in chl-*a* followed by an increase in recent years. Both models had comparable abilities to remove flow effects in simulated time series of chl-*a*, although flow-normalized predictions to actual data

suggested GAMs results were more stable with minimal observations. This information will provide researchers with valuable guidance for using statistical models in trend evaluation, with particular relevance for computational requirements, desired products, and future data needs.

Key words: chlorophyll, estuary, generalized additive models, Patuxent River Estuary, trend analysis, weighted regression

1 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 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 chlorophyll 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 chlorophyll 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 chlorophyll 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 chlorophyll 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 contextual flow inputs 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. 2013), Lake Champlain (Medalie et al. 2012), and the Mississippi River (Sprague et al. 2011). 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). Similarly, GAMs can be used to describe variation in a response variable as a sum of smoothing functions for different predictors (Hastie and Tibshirani 1990, Wood 2006). In applications to water quality time series, GAMs are similar to WRTDS in that variable effects through time can be described in relation to seasonal or annual changes. Application of GAMs to describe eutrophication endpoints in tidal waters have not been as extensive as WRTDS, although exploratory analyses have suggested that results are comparable. Moreover, GAMs are particularly appealing because they are less computationally intense and provide more accessible estimates of model uncertainty than WRTDS. Despite the potential for both approaches to characterize system dynamics, the relative merits of each have not been evaluated. Quantitative 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 description of the relative abilities of WRTDS and GAMs to describe long-term changes in time series of eutrophication response endpoints in tidal waters. A thirty year time series of monthly chlorophyll observations from the Patuxent River Estuary is 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 models'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 computational requirements, uncertainty assessment, and potential needs for additional monitoring data.

2 *Methods*

2.1 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 (Fig. 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 vertically-mixed in the upper estuary. 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 km² watershed that is 49% forest, 28% grassland, 12% developed, and 10% cropland (Jordan et al. 2003). The US Geological Survey (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 maintains 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/>, Fig. 1 and Table 1). Water quality samples have been collected since 1985 at monthly or bimonthly intervals and include salinity, temperature, chlorophyll *a* (chl-*a*), dissolved oxygen, and additional dissolved or particulate nutrients and organic carbon. Seasonal variation in chl-*a* is observed across the stations with spring and summer blooms

occurring in the upper, oligohaline section, whereas primary production is generally higher in the lower estuary during winter months (Fig. 2). Chlorophyll concentrations are generally lowest for all stations in late fall and early winter. Periods of low flow are associated with higher chlorophyll concentrations in the upper estuary, whereas the opposite is observed for high flow. Stations TF1.6 and LE1.2 were chosen as representative time series from different salinity regions to evaluate the water quality models. Observations at each station capture a longitudinal gradient of watershed influences at TF1.6 to mainstem influences from the Chesapeake Bay at LE1.2. Long-term changes in chlorophyll have also been related to historical reductions in nutrient inputs following a statewide ban on phosphorus-based detergents in 1984 and wastewater treatment improvements in the early 1990s that reduced point sources of nitrogen (Lung and Bai 2003, Testa et al. 2008). Therefore, the chosen stations provide unique datasets to evaluate the predictive and flow-normalization abilities of each model given the differing contributions of landward and seaward influences on water quality.

Thirty years of monthly chlorophyll and salinity data from 1986 to 2014 were obtained for stations TF1.6 and LE1.2 from the Chesapeake Bay Program data hub (<http://www.chesapeakebay.net/data>). All data were vertically integrated throughout the water column for each date to create a representative sample of water quality. The integration averaged all values after interpolating from the surface to the maximum depth. Observations at the most shallow and deepest sampling depth were repeated for zero depth and maximum depths, respectively, to bound the interpolations within the range of the data. Daily flow data were also obtained from the USGS stream gage station at Bowie, Maryland and merged with the nearest date in the chlorophyll 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 chlorophyll data at TF1.6, whereas the continuous salinity record was used as a tracer of discharge at LE1.2. Both chlorophyll and discharge data were log-transformed. Censored data were not present in any of the data sets. Initial quality assurance checks for all monitoring data were conducted following standard protocols adopted by the Chesapeake Bay Program.

2.2 Model descriptions

2.2.1 Weighted Regression on Time, Discharge, and Season

The WRTDS method relates a response endpoint, 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 chlorophyll concentration to salinity and time (Beck and Hagy III 2015), where salinity is a tracer of freshwater inputs or tidal changes. The functional form of the model is a simple regression that relates the natural log of chlorophyll (Chl) to decimal time (T) and salinity (Sal) on a sinusoidal annual time scale (i.e., cyclical variation by year).

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

The tidal adaptation of WRTDS uses quantile regression models (Cade and Noon 2003) to characterize trends in different conditional distributions of chlorophyll, e.g., the median or 90th percentile. For comparison to GAMs, the original WRTDS model in Hirsch et al. (2010) that characterizes the conditional mean of the response was used. Mean models require an estimation of the back-transformation bias parameter for response variables in log-space. This is achieved using the standard error of residuals for each observation along the time series during back-transformation (Hirsch et al. 2010). Additionally, the WRTDS model uses survival regression as a variation of the weighted Tobit model (Tobin 1958) to account for censored observations beyond the detection limit (Hirsch and De Cicco 2014).

The WRTDS approach obtains fitted values of the response variable by estimating regression parameters for each unique observation. Specifically, a unique regression model is estimated for each point in the period of observation. Each model is weighted by month, year, and salinity (or flow) such that a unique set of regression parameters for each observation is obtained. For example, a weighted regression centered on a single observation weights other observations in the same year, month, and similar salinity with higher importance, whereas observations for different months, years, or salinities receive lower importance. This weighting approach allows estimation of regression parameters that vary in relation to observed conditions throughout the period of record (Hirsch et al. 2010). Optimal window widths can be identified using

cross-validation, described below, that evaluates the ability of the model to generalize results with novel datasets.

Predicted values are based on an interpolation matrix from the unique regressions at each time step. A sequence of salinity or flow values based on the minimum and maximum values for the data are used to predict chlorophyll using the observed month and year based on the parameters fit to the observation. Model predictions are based on a bilinear interpolation from the grid using the salinity (flow) and date values closest to observed. Salinity- or flow-normalized values are also obtained from the prediction grid that allow an interpretation of chlorophyll trend that is independent of variation related to freshwater inputs. Normalized predictions are obtained for each observation by collecting the sample of observed salinity or flow values that occur for the same month throughout all years in the dataset. These values are assumed to be equally likely to occur across the time series at that particular month. A normalized value for each point in the time series is the average of the predicted values from each specific model based on the salinity or flow values that are expected to occur for each month.

2.2.2 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). The relationships between log-chl-*a* and the covariates were modeled with thin plate regression splines (Wood 2006) as the smooth functions using the ‘mgcv’ package in R. To allow for interaction between the model covariates (e.g., seasonal differences in the long-term chl-*a* pattern), a tensor product basis between all three covariates was constructed. The tensor product basis allows for the smooth construct to be a function of any number of covariates, without an isotropy constraint. The GAM implementation in ‘mgcv’ does not require the selection of knots for a spline basis, but instead a reasonable upper limit on the flexibility of the function is set, and a ‘wiggleness’ penalty is added to create a penalized regression spline structure. The balance between model fit and smoothness is achieved by selecting a smoothness parameter that minimizes the generalized cross-validation score (Wood 2006).

Model predictions with GAMs are straightforward to obtain after the model parameters

are selected, and can be obtained along with standard errors which are based on the Bayesian posterior covariance matrix (Wood 2006). For this comparison, salinity- or flow-normalized GAM predictions were obtained in a manner for consistency with 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.

2.2.3 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 2006). Overfitting a model with excessively small window widths or smoothing parameter will minimize prediction error but prevent extrapolation of results to different datasets. Similarly, underfitting a model with large window widths or smoothing parameter will reduce precision but will improve the ability to generalize results to different datasets. From a statistical perspective, the optimal model parameters provide a balance between over- and under-fitting. Both models use a form of cross-validation to identify model parameters that maximize the precision of model predictions with a novel dataset.

The basic premise of cross-validation is to identify the optimal set of model parameters that minimize prediction error on a dataset that was not used to develop the model. For the GAM approach, generalized cross-validation is 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 varies with the smoothing parameter (Wood 2006). Similarly, the tidal adaptation of WRTDS used k-fold cross-validation to identify the optimal half-window widths. For a given set of half-window widths, the dataset was separated into ten disjoint sets, such that ten models were evaluated for every combination of k - 1 training and remaining test datasets. That is, the training dataset for each fold was all k - 1 folds and the test dataset was the remaining fold, repeated k times. The average prediction error of the test datasets across k folds provided an indication of model performance for the given combination of

half-window widths. The optimum window widths were those that provided minimum errors on the test data. Evaluating multiple combinations of window-widths can be computationally intensive. An optimization function was implemented in R (Byrd et al. 1995, RDCT (R Development Core Team) 2015) to more efficiently evaluate model parameters using a search algorithm. Window widths were searched using the limited-memory modification of the BFGS quasi-Newton method that imposes upper and lower bounds for each parameter. The chosen parameters were based on a selected convergence tolerance for the error minimization of the search algorithm.

2.3 Comparison of modelled trends

Separate WRTDS and GAMs were created using the above methods for the chlorophyll time series at TF1.6 and LE1.2. Initial analyses indicated that model performance could be improved using the flow record to model chl-*a* at TF1.6 and the salinity record to model chl-*a* at LE1.2. For each model and station, a predicted and flow-normalized (hereafter flow-normalized refers to both flow and salinity) time series was obtained for comparison. The results were compared using several summary statistics that evaluated both the predictive performance to describe observed chlorophyll and direct comparisons between the models. Emphasis was on agreement between observed and predicted values, rather than uncertainty associated with parameter estimates or model results. As of writing, methods for estimating confidence intervals of WRTDS have been developed for the original model (Hirsch et al. 2015), but have not been fully developed for application to WRTDS in tidal waters. In addition to simple visual evaluation of trends over time, summary statistics used to compare model predictions to observed chl-*a* included root mean square error (RMSE) and average differences. For all comparisons, RMSE comparing each model's predictions to observed chl-*a* (fit) was defined as:

$$RMSE_{fit} = \sqrt{\frac{\sum_{i=1}^n (Chl_i - \widehat{Chl}_i)^2}{n}} \quad (2)$$

where n is the number of observations for a given evaluation, Chl_i is the observed value of chl-*a* for observation i , and \widehat{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 using

261 RMSE are similar, such that:

$$RMSE_{btw} = \sqrt{\frac{\sum_{i=1}^n (\widehat{Chl}_{WRTDS,i} - \widehat{Chl}_{GAM,i})^2}{n}} \quad (3) \quad \{\text{rmse_fun}\}$$

262 where the estimated chl-*a* values for each model, $\widehat{Chl}_{i,WRTDS}$ and $\widehat{Chl}_{i,GAM}$, are compared
 263 directly. Similarly, average differences (or bias) of predictions between models as a percentage
 264 was defined as:

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

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

273 The statistical comparisons described above were conducted for the entire time series at
 274 each station to evaluate overall performance. Different time periods were also evaluated to
 275 identify potential temporal variation in results, which included a comparison of results by annual
 276 and seasonal aggregations and periods with different levels of flow using the discharge record at
 277 Bowie, Maryland. Annual and seasonal aggregations shown in Fig. 2 were evaluated between the
 278 models, in addition to evaluating the models at different levels of flow defined by the quartile
 279 distributions (min–25%, 25%–median, median–75%, and 75%–max). Flow-normalized time
 280 series were compared similarly but only between the models because the true flow-independent
 281 component of the observed data is not known and can only be empirically estimated. As
 282 described below, an evaluation of flow-normalized data for each model was accomplished using
 283 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.

2.4 Comparison of flow-normalized trends

The relative abilities of each model to characterize flow-normalized trends in chlorophyll were evaluated using simulated datasets with known components. This approach was used because the flow-independent component of chlorophyll is typically not observed in raw data such that the true signal must be empirically estimated. Accordingly, the ability of each model to isolate the flow-normalized trend cannot be evaluated with reasonable certainty unless the true signal is known. Simulated time series of observed chlorophyll (Chl_{obs}) were created as additive components related to flow (Chl_{flo}) and a flow-independent biological component of chlorophyll (Chl_{bio}):

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

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 chlorophyll in the water column that is directly linked to primary production. The biological component of chlorophyll is comparable to an observation in a closed 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 chlorophyll records were obtained 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. Daily chlorophyll concentrations were estimated from fluorescence values that did not include blue-green algae blooms. Our primary concern was simulating chlorophyll concentrations at monthly or bimonthly timesteps such that taxa-specific concentrations on a daily time step were not relevant.

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 chlorophyll independent of discharge as a stationary seasonal component (\widehat{Chl}_{seas}), and 4) simulated error structure from the residuals of the seasonal chlorophyll model ($\varepsilon_{Chl, sim}$). Unless otherwise noted, chlorophyll and discharge are in ln-transformed units. Each of the four components was used to simulate the components in eq. (5):

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

$$Chl_{bio} = \widehat{Chl}_{seas} + \sigma \cdot \varepsilon_{Chl, sim} \quad (7) \quad \{chl_{bio}\}$$

The estimated flow time series within the parentheses, $\hat{Q}_{seas} + \sigma \cdot \varepsilon_{Q, sim}$, 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 chlorophyll. Moreover, the vector I (where $0 \leq I \leq 1$) can be manually changed to represent an independent effect of flow based on the desired simulation. 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. Methods for estimating each of the components in eqs. (6) and (7) are described in detail below.

First, a model for simulating flow-related chlorophyll (eq. (6)) was estimated from the stream gage data as the additive combination of a stationary seasonal component and serially-correlated errors:

$$Q_{seas} = \beta_0 + \beta_1 \sin(2\pi T) + \beta_2 \cos(2\pi T) \quad (8) \quad \{qseas\}$$

$$\varepsilon_Q = Q_{seas} - \hat{Q}_{seas} \quad (9) \quad \{qerr\}$$

A seasonal model of flow was estimated using linear regression for time, T , on an annual sinusoidal period (eq. (8)). The residuals from this regression, ε_Q (eq. (9)), were used to estimate the structure of the error distribution for simulating the stochastic component of flow. The error distribution was characterized using an Autoregressive Moving Average (ARMA) model to identify appropriate p and q coefficients (Hyndman and Khandakar 2008). The parameters were chosen using stepwise estimation for nonseasonal univariate time series that minimized Akaike Information Criterion (AIC). The resulting coefficients were used to generate random errors from a standard normal distribution for the length of the original time series, $\varepsilon_{Q,sim}$. These stochastic errors were multiplied by the standard deviation of the residuals in eq. (9) and added to the seasonal component in eq. (8) to create a simulated, daily time series of the flow-component for chlorophyll, Chl_{flo} (eq. (6)).

The chlorophyll time series was created using a similar approach. The first step estimated the stationary seasonal component of the chlorophyll time series by fitting a WRTDS model (Hirsch et al. 2010) that explicitly included discharge from the gaged station using one year of data from the whole time series:

$$Chl_{seas} = \beta_0 + \beta_1 T + \beta_2 Q + \beta_3 \sin(2\pi T) + \beta_4 \cos(2\pi T) \quad (10) \quad \{chlseas\}$$

$$\varepsilon_{Chl} = Chl_{seas} - \widehat{Chl}_{seas} \quad (11) \quad \{chlerr\}$$

This approach was used to isolate an error structure for simulation that was independent of flow

and biology, where the seasonal component (as time T on a sinusoidal annual period) was assumed to be related to biological processes. The error distribution was then estimated from the residuals (eq. (11)) as before using an ARMA estimate of the residual parameters, p and q . Standard error estimates from the regression used at each point in the one-year time series were also retained for each residual. Errors were simulated ($\varepsilon_{Chl, sim}$, eq. (7)) for the entire year using the estimated auto-regressive structure and multiplied by the corresponding standard error estimate from the regression. The entire year was repeated for every year in the observed time series. All simulated errors were rescaled to the range of the original residuals that were used to estimate the distribution. Finally, the simulated flow-component, Chl_{flo} , was added to the simulated biological model, Chl_{bio} , to create the final chlorophyll-flow time series, Chl_{obs} , in eq. (5).

A daily time series for the entire period of record was simulated using the above methods and then used to compare the relative abilities of WRTDS and GAMs to characterize flow-normalized trends. Three time series with monthly sampling frequencies and varying contributions of the flow component (Chl_{flo} in eq. (5)) were created from the daily time series (Fig. 7). One day in each month for each year was randomly sampled and used as the monthly time step for each time series. Varying effects of the flow component on observed chlorophyll were created by multiplying Chl_{flo} by different indicator vectors (I in eq. (6)). The contribution of the flow component varied from non-existent, constant, and steadily increasing. Respectively, the vector of coefficients applied to each flow component was a constant vector of zeroes, a constant vector of ones, and a linear increase starting at zero and ending at 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 (Fig. 7). Results were evaluated by first comparing the predicted ($\widehat{Chl_{obs}}$) and observed (Chl_{obs}) chlorophyll values for each simulation, following by comparing the flow-normalized results ($\widehat{Chl_{bio}}$) from each model to the original biological chlorophyll (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.

3 Results

3.1 Observed trends and relative fit

The optimal half-window widths and degrees of freedom for smoothing varied for WRTDS and GAMs, respectively, at each station. For WRTDS, optimal half-window widths identified by generalized cross-validation were 0.25 as a proportion of each year, 13.59 years, and 0.25 as a proportion of the total range of salinity 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 (annual proportion) and flow component. For GAMs, the optimal smoothing procedure resulted in a smoother model at LE1.2 than TF1.6 with effective degrees of freedom of 32.2 and 47.1, respectively. The tensor product smooth construct does not split apart the effective degrees of freedom among the three interacting parameters.

The predicted chl-*a* from each model generally followed patterns in observed chl-*a* from 1986 to 2014 (Fig. 3). At LE1.2, each model showed seasonal minimum typically in November, whereas maximum chl-*a* was observed in a spring bloom, typically March or April (Fig. 4). A secondary, smaller seasonal peak was also observed in late summer from bottom-layer regeneration and upward nutrient transport (Testa et al. 2008). Visual comparison of results between models showed that the secondary seasonal maxima at each site were more pronounced for WRTDS estimates. 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 (Fig. 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 averaging $3.5 \mu\text{g L}^{-1}$ of chl-*a* at LE1.2 and $9.5 \mu\text{g L}^{-1}$ 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* (Fig. 3a). Normalized predictions for both models were visually distinct from observed predictions such that seasonal minima and maxima and extreme predictions were not common

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 predicted values 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 0.51 at LE1.2 for WRTDS predictions and a maximum of 0.54 at TF1.6 for GAM predictions (Table 2). Overall, both models performed similarly, although WRTDS had slightly lower RMSE for the overall time series and all time periods (annual, seasonal, and flow) regardless of site (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 0.63 for GAMs, RMSE 0.60 for WRTDS), the April-May-June (AMJ) seasonal periods (0.69 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 (0.58 for GAMs, 0.58 for WRTDS) and the AMJ seasonal period (0.60 for GAMs, 0.58 for WRTDS). Error rates between models were comparable for all flow periods at TF1.6, with the exception of lower error rates during low flow (0.48 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.

Results as annual aggregations suggested that chl-*a* patterns between years have not been constant and are considerably different between sites (Fig. 3b). Both models showed a gradual and consistent increase in chl-*a* at LE1.2, with values increasing by approximately $1.5 \mu\text{g L}^{-1}$ 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 $12 \mu\text{g L}^{-1}$ to $6 \mu\text{g L}^{-1}$ from 2000 to 2006 was observed. By 2014, chlorophyll returned to values similar to those prior to the initial decrease. Flow-normalized predictions that were averaged annually 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 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 (20% 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 2008–2014 period (Table 4).

Seasonal changes were especially pronounced during the JFM and October-November-December (OND) periods where both models showed an increase and decrease, respectively, with differences between the two (JFM period, 5.7% for GAMs, 32.7% for WRTDS; OND period, –13.6% for GAMs, –17.5% for WRTDS). Percent changes by flow period 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.

3.2 Comparison of model predictions

The following describes comparisons of model results, whereas the previous section emphasized results relative to trends over time and fit to the observed data. Accordingly, direct comparisons between model predictions provided a means of identifying instances when models results were systematically different from each other. Table 5 compares average differences and RMSE 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 only one RMSE value greater than 0.20 (RMSE 0.22 for JFM at TF1.6). The greatest differences between the models were observed in the seasonal comparisons, with the most extreme differences observed at LE1.2 for the JAS (5.03%, WRTDS greater than GAMs) and OND (–5.25, WRTDS less than GAMs) periods. Differences between stations were not apparent.

Regressions comparing model results provided additional information about overall differences (significantly different intercept) and differences between the models that varied for different values (significantly different slope) (Table 6, Fig. 5). No significant differences were observed for the entire time series such that estimated intercepts and slopes were similar to zero and one, respectively, for both stations and model predictions (observed and flow-normalized). However, differences were 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 JAS period (Table 6). All site and model prediction comparisons also had significantly different intercept and slope estimates for the low flow time period. Visual comparisons of results in Fig. 5 confirm those in Table 6, particularly differences in the seasonal aggregations. For example, comparisons of WRTDS and GAMs results for the JFM period at TF1.6 showed an intercept greater than zero and a tendency for GAM estimates to be higher than WRTDS at the right-side of the distribution (i.e., slope less than one with the regression line below 1:1).

3.3 Changes in chlorophyll response to flow over time

Both models described chlorophyll 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 primary production, 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 such 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 novel data.

Fig. 6 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 an 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, particularly for January at LE1.2 where WRTDS provided a wider range, or potentially less stable response of chl-*a* to salinity changes in the earlier years.

3.4 Flow-normalization with simulated data

WRTDS and GAMs were fit to each dataset creating six models to evaluate the general fit of observed to predicted ($Chl_{obs} \sim \widehat{Chl}_{obs}$) and biological to flow-normalized chl-*a* ($Chl_{bio} \sim \widehat{Chl}_{bio}$). 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. Fig. 8 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 confirm those in Fig. 7 by showing 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 salinity is apparent in the third panel of Fig. 8 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 for the GAMs.

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 ($Chl_{obs} \sim \widehat{Chl}_{obs}$, constant flow simulation), although differences in performance were minor (Table 7). Visual comparison of results suggested that both models provided comparable information for predictions of observed values and flow-normalized predictions (Fig. 9). Additionally, the varying effect of flow on each time series was apparent in comparisons of predicted with flow-normalized results, such that \widehat{Chl}_{bio} was increasingly different from \widehat{Chl}_{obs} from no effect to constant effect of the flow component (top row, Fig. 9). 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 \widehat{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 (\widehat{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 (\widehat{Chl}_{bio}) with relatively minor differences between the models.

4 Discussion

4.1 Model comparisons and considerations

Methods for comparing alternative statistical models can be defined by the type of information desired from each technique. As applied to estuaries, both models have the general objective of describing trends from long-term monitoring datasets, whereas more specific applications of each model (e.g., hypothesis testing) will be defined by future management or research objectives. Accordingly, our comparison methods were chosen to provide a comprehensive overview of each model based on the exploratory objectives of the analysis and that each technique provides a new approach to trend assessment. We evaluated predictive performance of both observed and flow-normalized trends, comparison between models for potential bias and indications of trend, and descriptions of canonical variation related to temporal or flow effects. The variety of methods for comparing models can provide different information depending on the desired application. An improvement in predictive performance using RMSE, for example, may suggest one model is more advantageous over another if the goal is to reproduce trends, whereas this information may be irrelevant for hypothesis testing. Inferior performance for one metric does not necessarily invalidate an analysis method for all potential applications. An interpretation of our specific results should consider that the intent was to provide an overview with several techniques given that the purpose of each model will likely be better defined by future applications, although a general focus for the current assessment is on trend evaluation.

A general conclusion from our results is that both models provide similar information, both in predictive performance and trends over time in the Patuxent. The comparisons of RMSE consistently showed that WRTDS had lower prediction error than GAMs for all scenarios. A naive interpretation may suggest that WRTDS is a superior model given these results, but we emphasize that the improvement in performance is trivial both in the relative values and differences in the theoretical foundations of each model. Moreover, prediction error for GAMs

could easily be improved over WRTDS with only a slight adjustment of model parameters. This highlights a potential pitfall of using prediction error as a performance metric because the values are sensitive to tuning parameters and the statistical characteristics of a training dataset. To address this issue, comparable methods for model development were implemented to legitimize the comparisons of predictive performance. Both WRTDS and GAMs used a form of cross-validation to identify an optimal parameter space that minimized the bias-variance tradeoff on separate training and test datasets. A more generic interpretation of the benefits of cross-validation is that model development is not biased by analyst intervention as the parameters are chosen with predefined heuristics. Although further development of the technique is needed, this paper presents the first application of a statistical method of selecting optimal window widths for WRTDS. This method allowed for a valid comparison model performance and will likely facilitate future applications of the model to alternative datasets.

The comparisons of predictive performance for each model should also be interpreted relative to the statistical foundations of each model. The smoothing process in GAMs, although mathematically involved, rapidly converges to a solution as compared to the fitting process for WRTDS that requires a unique regression estimate for every point in the time series. From a practical perspective, the comparable error estimates for each model's predictions suggests that GAMs are advantageous because there appears to be little benefit of the added computational time of WRTDS. More surprisingly, the characterization of dynamic changes in flow relationships over time appear comparable for each model. For example, Fig. 6 shows similar information for each model although different methods defined by the models were used to describe chl-*a* changes over time with variation in salinity or flow. A simple grid of explanatory variables spanning the distribution space of the observed variables was used as input for the fitted GAMs, whereas WRTDS results in the same figure used each model's fitted interpolation grid. As such, novel insight into trends over time would be expected with the added computational time required for the WRTDS interpolation grids. Conventional modelling techniques that have a predefined and limited parameter space have been described as 'statistical straightjackets' that mold the data to the model. WRTDS is meant to provide a contrasting approach where the data mold the results, as compared to GAMs that could be considered overconstrained in this context by following a predefined model structure. The results do not provide a compelling contrast between GAMs and

WRTDS despite the statistical differences between the models.

The interpretation that WRTDS does not provide additional insight with the added computational time may be misguided. A lack of clear differences between each model could have been related to characteristics of the datasets. The use of data-driven models to identify emergent patterns in the data necessarily requires that these characteristics are real phenomena that are not readily observed in the raw data. A logical expectation for trend evaluation is that different methods would lead to similar conclusions for datasets that lack dynamic variation, such as differences between flow relationships and response endpoints over time. Similarity in results for WRTDS and GAMs could simply mean that relationships between time, season, and flow in the Patuxent are adequately described by the statistical theories of each approach. A priori site-selection of TF1.6 and LE1.2 was meant to capture a gradient of watershed to mainstem influences on productivity at each location. The known historical changes from management practices (e.g., wastewater treatment, banning of phosphorus-based detergents) and natural events (e.g., storm events, seagrass recovery) that have affected the Patuxent have also provided a unique context for the time series. The assumption that these characteristics of the datasets will translate to differences in the model results may have been misguided given that we were not able to show clear differences. Generalizations of the relative merits of each model should be made sparingly until additional assessments with alternative datasets are made. Finally, comparability between the two methods may also suggest that both models were equally bad at describing dynamic relationships in the dataset. An assumption throughout the analysis that is not entirely unreasonable is that both models were equally ‘good’, although the possibility that both were equally inadequate should also be considered as a potential explanation. Alternative drivers of chlorophyll response that are not explicitly included in each model could limit explanatory power if time, season, and discharge were not the dominant predictors of production.

Although our results generally indicated that comparable information was provided by both models, instances were also observed when different information was provided. Descriptions of canonical variation of seasonal trends from the model predictions suggested that WRTDS more adequately characterized seasonal minima and maxima in chlorophyll production. Although both similarly described the general structure of seasonal peaks at each station, WRTDS provided greater distinction of the secondary seasonal maxima (Fig. 4, September/October peak at LE1.2,

May/June peak at TF1.6). A post-hoc assessment of seasonal variation at each station showed that GAMs poorly characterized the secondary seasonal peak during the last ten years of the time series despite the peak being present in the observed data. A quantitative comparison of predictions by seasonal time periods also suggested that differences between the models were more often observed during specific months within each year (Tables 5 and 6 and Fig. 5), as compared to variation between different annual or flow periods. The potential for WRTDS to more adequately characterize seasonal variation in chl-*a* patterns deserves additional attention, particularly if changes in observed chl-*a* for specific months of the year have occurred at a monitoring station. Cloern and Jassby (2010) emphasize the importance and potential difficulty of characterizing different components of chl-*a* variability. WRTDS applied to tidal waters could prove a valuable tool to develop a more comprehensive understanding of spatio-temporal variation in phytoplankton variability.

Differences between the models may also indicate potential shortcomings, in addition to advantages. The computationally intensive fitting method for WRTDS was expected to allow additional insight into variation of chl-*a* in response to the primary drivers of time, discharge, and season. This novel information was investigated several ways, including comparisons of model performance, evaluations of time periods or flow-normalized trends, or a simple visual distinction between model results. As noted above, seasonal differences were observed between the models, whereas the additional comparisons suggested both methods provided comparable information. However, initial assessment of Fig. 6 suggested that WRTDS provided a more dynamic description of chl-*a* response to changes in flow or salinity over time for specific locations in the record. For example, chl-*a* response over time to salinity changes during January at LE1.2 shows WRTDS describing greater variation than GAMs, particularly for lower salinity values. Additional investigation suggested that these ‘novel’ descriptions were related to low sample size for the specific location in the record causing instability in the model predictions. Accordingly, WRTDS descriptions may be unstable at extreme or uncommon locations in the data domain where the number of observations with non-zero weights may be limited. Methods for WRTDS have been developed to address this issue (i.e., automated window width increases with low sample sizes), although they were not implemented for the current analysis due to complications during model fitting and interpretations of trend comparisons with GAMs.

4.2 Patuxent trends

A secondary benefit of applying WRTDS and GAMs to evaluate chl-*a* time series in the Patuxent is a detailed description of water quality changes in this highly visible and well-studied estuary. Although our primary objective was to evaluate the quantitative capabilities of each model, several interesting trends were described that deserve additional investigation.

4.3 Conclusions

References

- Beck MW, Hagy III JD. 2015. Adaptation of a weighted regression approach to evaluate water quality trends in an estuary. *Environmental Modelling and Assessment*, pages 1–19.
- Boicourt WC, Sanford LP. 1998. A hydrodynamic study of the Patuxent River estuary. Technical report, Maryland Department of the Environment, Baltimore, Maryland.
- Borsuk ME, Stow CA, Reckhow KH. 2004. Confounding effect of flow on estuarine response to nitrogen loading. *Journal of Environmental Engineering-ASCE*, 130(6):605–614.
- Byrd RH, Lu P, and C. Zhu JN. 1995. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing*, 16(5):1190–1208.
- Cade BS, Noon BR. 2003. A gentle introduction to quantile regression for ecologists. *Frontiers in Ecology and the Environment*, 1(8):412–420.
- Cloern JE. 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 JE, Jassby AD. 2010. Patterns and scales of phytoplankton variability in estuarine-coastal ecosystems. *Estuaries and Coasts*, 33(2):230–241.
- Cronin WB, Pritchard DW. 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.
- Hagy JD. 1996. Residence times and net ecosystem processes in Patuxent River estuary. Master's thesis, University of Maryland, College Park, Maryland.
- Harding LW. 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.
- Hastie T, Tibshirani R. 1990. *Generalized Additive Models*. Chapman and Hall, London, New York.
- Hirsch RM. 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 RM, Archfield SA, De Cicco LA. 2015. A bootstrap method for estimating uncertainty of water quality trends. *Environmental Modelling and Software*, 73:148–166.
- Hirsch RM, De Cicco L. 2014. User guide to Exploration and Graphics for RivEr Trends (EGRET) and dataRetrieval: R packages for hydrologic data. Technical Report Techniques and Methods book 4, ch. A10, US Geological Survey, Reston, Virginia.
<http://pubs.usgs.gov/tm/04/a10/>.

- Hirsch RM, Moyer DL, Archfield SA. 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.
- Hyndman RJ, Khandakar Y. 2008. Automatic time series forecasting: The forecast package for R. *Journal of Statistical Software*, 26(3):1–22.
- Jordan TE, Weller DE, Correll DL. 2003. Sources of nutrient inputs to the Patuxent River estuary. *Estuaries and Coasts*, 26(2A):226–243.
- Lung W, Bai S. 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, Hirsch RM, Archfield SA. 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 DL, Hirsch RM, Hyer KE. 2012. Comparison of two regression-based approaches for determining nutrient and sediment fluxes and trends in the Chesapeake Bay watershed. Technical Report Scientific Investigations Report 2012-544, US Geological Survey, US Department of the Interior, Reston, Virginia.
- Nixon SW. 1995. Coastal marine eutrophication: A definition, social causes, and future concerns. *Ophelia*, 41:199–219.
- Paerl HW. 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 HW, Hall NS, Peierls BL, Rossignol KL. 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.
- Powers SP, Peterson CH, Christian RR, Sullivan E, Powers MJ, Bishop MJ, Buzzelli CP. 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>.
- Sprague LA, Hirsch RM, Aulenbach BT. 2011. Nitrate in the Mississippi River and its tributaries, 1980 to 2008: Are we making progress? *Environmental Science and Technology*, 45(17):7209–7216.

- 735 Steward JS, Green WC. 2007. Setting load limits for nutrients and suspended solids based upon
736 seagrass depth-limit targets. *Estuaries and Coasts*, 30(4):657–670.
- 737 TBEP (Tampa Bay Estuary Program). 2011. Tampa Bay Water Atlas.
738 <http://www.tampabay.wateratlas.usf.edu/>. (Accessed October, 2013).
- 739 Testa JM, Kemp WM, Boynton WR, Hagy JD. 2008. Long-term changes in water quality and
740 productivity in the Patuxent River Estuary: 1985 to 2003. *Estuaries and Coasts*,
741 31(6):1021–1037.
- 742 Tobin J. 1958. Estimation of relationships for limited dependent variables. *Econometrica*,
743 26(1):24–36.
- 744 Wood SN. 2006. *Generalized Additive Models: An Introduction with R*. Chapman and Hall,
745 CRC Press, London, United Kingdom.
- 746 Zhang Q, Brady DC, Ball WP. 2013. Long-term seasonal trends of nitrogen, phosphorus, and
747 suspended sediment load from the non-tidal Susquehanna River Basin to Chesapeake Bay.
748 *Science of the Total Environment*, 452-453:208–221.

Table 1: Summary characteristics of monitoring stations on the Patuxent River estuary. Chlorophyll 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). See Fig. 1 for site locations.

Station	Lat	Long	Segment	Distance (km)	Depth (m)	ln-Chl ($\mu\text{g L}^{-1}$)	Sal (ppt)
TF1.3	38.81	-76.71	TF	74.90	2.9	1.52	0.00
TF1.4	38.77	-76.71	TF	69.50	2.0	2.31	0.02
TF1.5	38.71	-76.70	TF	60.30	10.6	2.88	0.27
TF1.6	38.66	-76.68	OH	52.20	6.2	2.44	0.90
TF1.7	38.58	-76.68	OH	42.50	3.0	2.09	4.09
RET1.1	38.49	-76.66	MH	32.20	11.2	2.47	10.25
LE1.1	38.43	-76.60	MH	22.90	12.1	2.31	12.04
LE1.2	38.38	-76.51	MH	13.40	17.1	2.16	12.73
LE1.3	38.34	-76.48	MH	8.30	23.4	2.12	12.89
LE1.4	38.31	-76.42	MH	0.00	15.4	2.21	13.46

Table 2: Summaries of model performance using RMSE of observed to predicted ln-chlorophyll for each station (LE1.2 and TF1.6). Deviance for each model as the sum of squared residuals is shown in parentheses. 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 by monthly quarters (middle), and flow periods based on quantile distributions of discharge.

Period	LE1.2		TF1.6	
	GAM	WRTDS	GAM	WRTDS
All	0.54 (153.6)	0.51 (135.1)	0.54 (150.4)	0.52 (138.6)
Annual				
1986-1993	0.54 (47.5)	0.50 (40.9)	0.53 (45.7)	0.49 (39.1)
1994-2000	0.52 (37.1)	0.50 (33.2)	0.58 (44.5)	0.58 (44.9)
2001-2007	0.63 (55.1)	0.60 (49.6)	0.54 (38.6)	0.53 (37.5)
2008-2014	0.39 (14.0)	0.36 (11.4)	0.49 (21.6)	0.44 (17.1)
Seasonal				
JFM	0.61 (38.8)	0.58 (35.3)	0.53 (28.2)	0.49 (23.8)
AMJ	0.69 (74.7)	0.64 (65.3)	0.60 (55.4)	0.58 (51.9)
JAS	0.38 (22.1)	0.35 (18.6)	0.48 (35.0)	0.46 (32.2)
OND	0.41 (18.1)	0.38 (15.9)	0.55 (31.7)	0.54 (30.7)
Flow				
1 (Low)	0.40 (21.5)	0.36 (16.7)	0.48 (29.5)	0.46 (27.7)
2	0.47 (29.1)	0.42 (23.5)	0.56 (40.3)	0.54 (37.8)
3	0.61 (48.5)	0.57 (42.9)	0.56 (39.6)	0.52 (35.4)
4 (High)	0.64 (54.6)	0.63 (52.0)	0.56 (41.0)	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-chlorophyll ($\mu\text{g L}^{-1}$) based on annual means within each category. For example, summary values for high flow for a given model and are based on instances of high flow across years. 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 by monthly quarters (middle), and flow periods based on quantile distributions of discharge.

Period	GAM		WRTDS	
	Ave.	% Change	Ave.	% Change
All	2.19	20.40	2.18	18.85
Annual				
1986-1993	2.02	3.86	2.03	1.75
1994-2000	2.12	5.78	2.12	5.50
2001-2007	2.25	4.77	2.24	5.35
2008-2014	2.38	4.08	2.37	6.07
Seasonal				
JFM	2.51	15.20	2.58	14.04
AMJ	2.41	28.00	2.33	22.47
JAS	1.92	18.06	2.01	19.91
OND	1.92	15.57	1.83	15.14
Flow				
Flow 1 (Low)	1.91	13.04	1.93	16.77
Flow 2	2.12	10.04	2.11	7.73
Flow 3	2.30	12.22	2.29	9.24
Flow 4 (High)	2.34	20.00	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-chlorophyll ($\mu\text{g L}^{-1}$) based on annual means within each category. For example, summary values for high flow for a given model and are based on instances of high flow across years. 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 by monthly quarters (middle), and flow periods based on quantile distributions of discharge.

Period	GAM		WRTDS	
	Ave.	% Change	Ave.	% Change
All	2.44	-4.36	2.44	-2.28
Annual				
1986-1993	2.63	-2.40	2.60	-3.06
1994-2000	2.70	-5.95	2.65	-3.55
2001-2007	2.16	-22.05	2.19	-21.51
2008-2014	2.24	45.59	2.30	38.35
Seasonal				
JFM	1.56	5.70	1.48	32.72
AMJ	2.59	7.51	2.62	5.14
JAS	3.08	-2.18	3.08	0.79
OND	2.15	-13.57	2.20	-17.55
Flow				
Flow 1 (Low)	2.89	-5.99	2.93	-0.42
Flow 2	2.41	17.54	2.43	20.31
Flow 3	2.28	7.50	2.27	15.20
Flow 4 (High)	2.22	-8.51	2.21	-11.27

Table 5: Comparison of predicted results between WRTDS and GAMs using average differences (%) and RMSE values at each station. 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 by monthly quarters (middle), and flow periods based on quantile distributions of discharge. Negative percentages indicate WRTDS predictions were lower than GAM predictions (eq. (4)).

Period	LE1.2		TF1.6	
	Ave. diff.	RMSE	Ave. diff.	RMSE
All	-0.11	0.15	0.01	0.17
Annual				
1986-1993	0.18	0.16	-0.78	0.17
1994-2000	0.53	0.15	-1.09	0.19
2001-2007	-0.95	0.14	0.48	0.14
2008-2014	-0.18	0.14	3.12	0.18
Seasonal				
JFM	2.91	0.14	-5.02	0.22
AMJ	-3.42	0.17	0.93	0.14
JAS	5.03	0.14	-0.10	0.17
OND	-5.25	0.14	2.08	0.17
Flow				
Flow 1 (Low)	0.19	0.16	-0.09	0.12
Flow 2	-0.83	0.16	0.73	0.15
Flow 3	0.19	0.15	0.84	0.20
Flow 4 (High)	0.03	0.13	-1.62	0.20

Table 6: Regression fits comparing predicted (*pred*) and flow-normalized (*norm*) results for WRTDS and GAMs at each station. 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 by monthly quarters (middle), and flow periods based on quantile distributions of discharge. See Fig. 5 for a graphical summary.

Period	$\beta_{0,pred}$		$\beta_{1,pred}$		$\beta_{0,norm}$		$\beta_{1,norm}$	
	LE1.2	TF1.6	LE1.2	TF1.6	LE1.2	TF1.6	LE1.2	TF1.6
All	0.01	0.04	0.99	0.98	0.02	-0.02	0.99	1.01
Annual								
1986-1993	-0.03	-0.06	1.02	1.01	0.03	-0.19	0.99	1.06
1994-2000	0.05	-0.09	0.98	1.01	-0.04	-0.18	1.02	1.05
2001-2007	0.00	0.11	0.99	0.95	-0.01	0.08	1.00	0.98
2008-2014	0.05	0.09	0.98	0.99	0.04	0.04	0.98	1.01
Seasonal								
JFM	0.10	0.33	0.99	0.74	0.01	0.50	1.02	0.62
AMJ	-0.40	0.06	1.13	0.99	-0.27	0.03	1.08	1.00
JAS	0.35	0.48	0.86	0.83	0.66	0.62	0.70	0.79
OND	0.00	0.01	0.95	1.02	-0.34	-0.16	1.13	1.09
Flow								
Flow 1 (Low)	0.22	-0.19	0.87	1.06	0.48	0.20	0.75	0.94
Flow 2	0.11	0.06	0.94	0.98	0.10	0.01	0.95	1.00
Flow 3	-0.03	0.20	1.02	0.91	-0.08	-0.10	1.04	1.03
Flow 4 (High)	-0.14	0.00	1.06	0.98	-0.13	0.00	1.05	0.99

Table 7: Summaries of model performance comparing observed chlorophyll with predicted values ($Chl_{obs} \sim \widehat{Chl}_{obs}$) and biological chlorophyll with flow-normalized values ($Chl_{bio} \sim \widehat{Chl}_{bio}$) for the three simulated time series (no flow, constant flow, and increasing flow effect). Summaries are RMSE values comparing results from each model (GAM, WRTDS) in the bottom two rows of panels in Fig. 9. Deviance for each model as the sum of squared residuals is shown in parentheses.

Simulations	$Chl_{obs} \sim \widehat{Chl}_{obs}$	$Chl_{bio} \sim \widehat{Chl}_{bio}$
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)

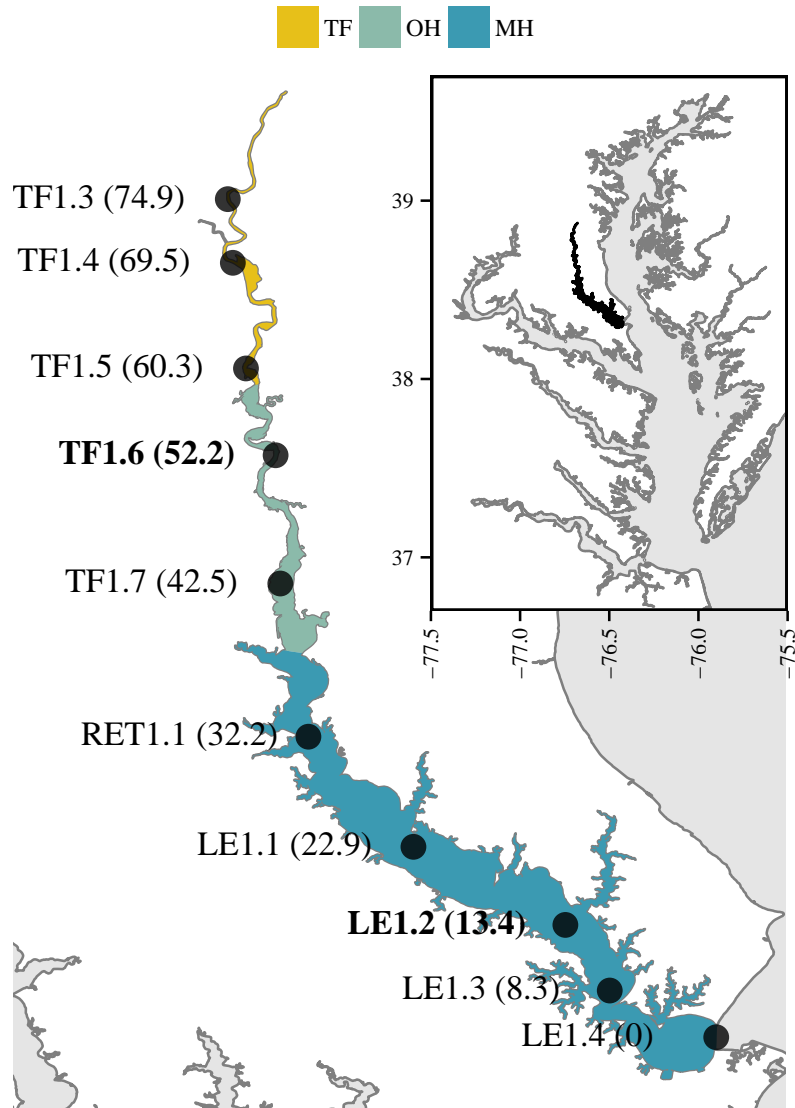


Fig. 1: Patuxent River estuary with Chesapeake Bay inset. Fixed locations monitored by the Chesapeake Bay Program at monthly frequencies are shown along the longitudinal axis with distance from the mouth (km). Study sites are in bold. Salinity regions in the Patuxent for the larger Chesapeake Bay area are also shown (TF = tidal fresh, OH = oligohaline, MH = mesohaline). See Table 1 for a numeric summary of station characteristics.

{fig:map}

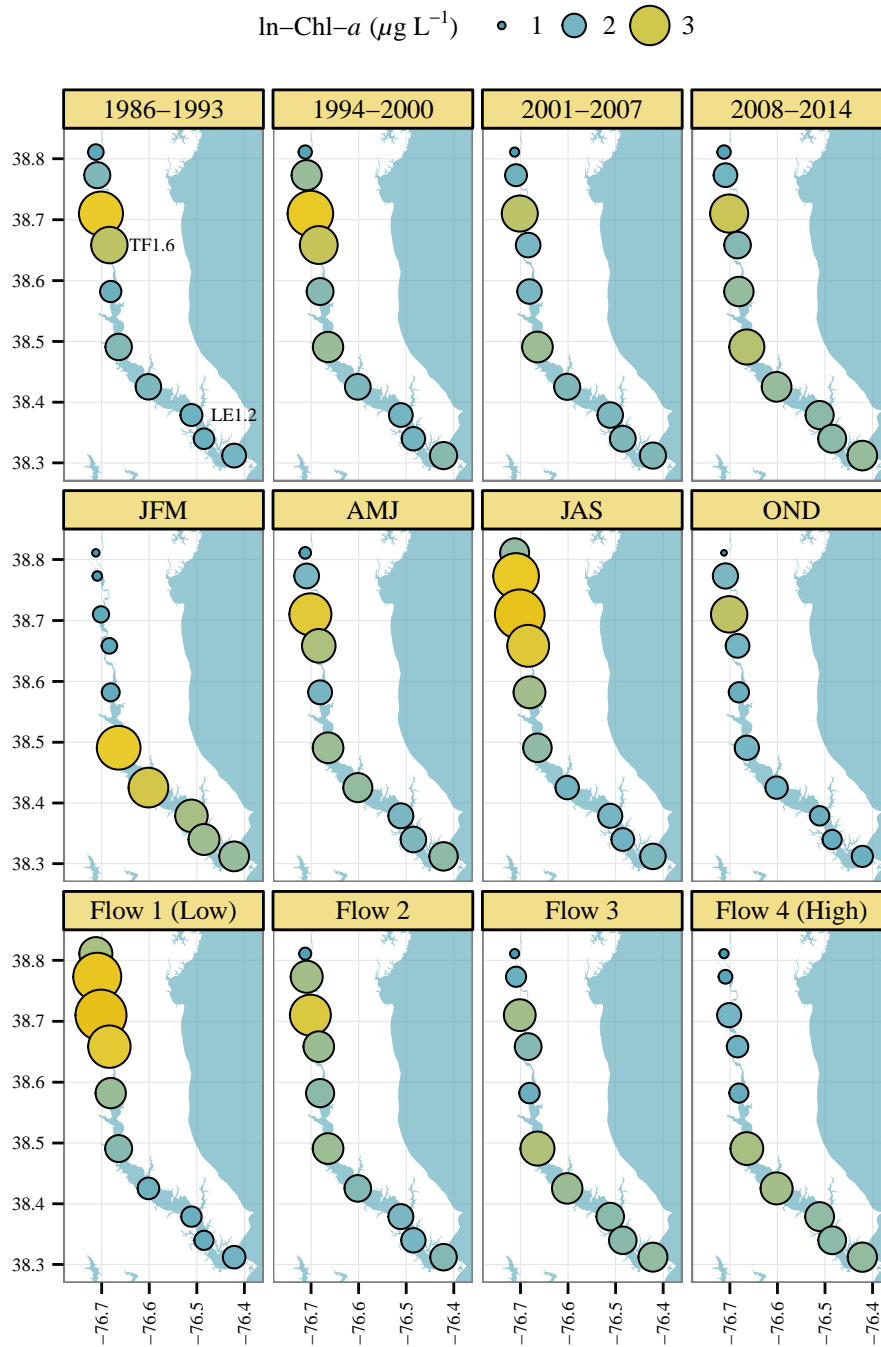


Fig. 2: Annual, seasonal, and flow differences in chlorophyll trends at each monitoring station in the Patuxent River Estuary. Size and color are proportional medians of $\ln\text{-chlorophyll-}a$ by year, season, and flow categories. See Fig. 1 for station numbers.

{fig:chly}

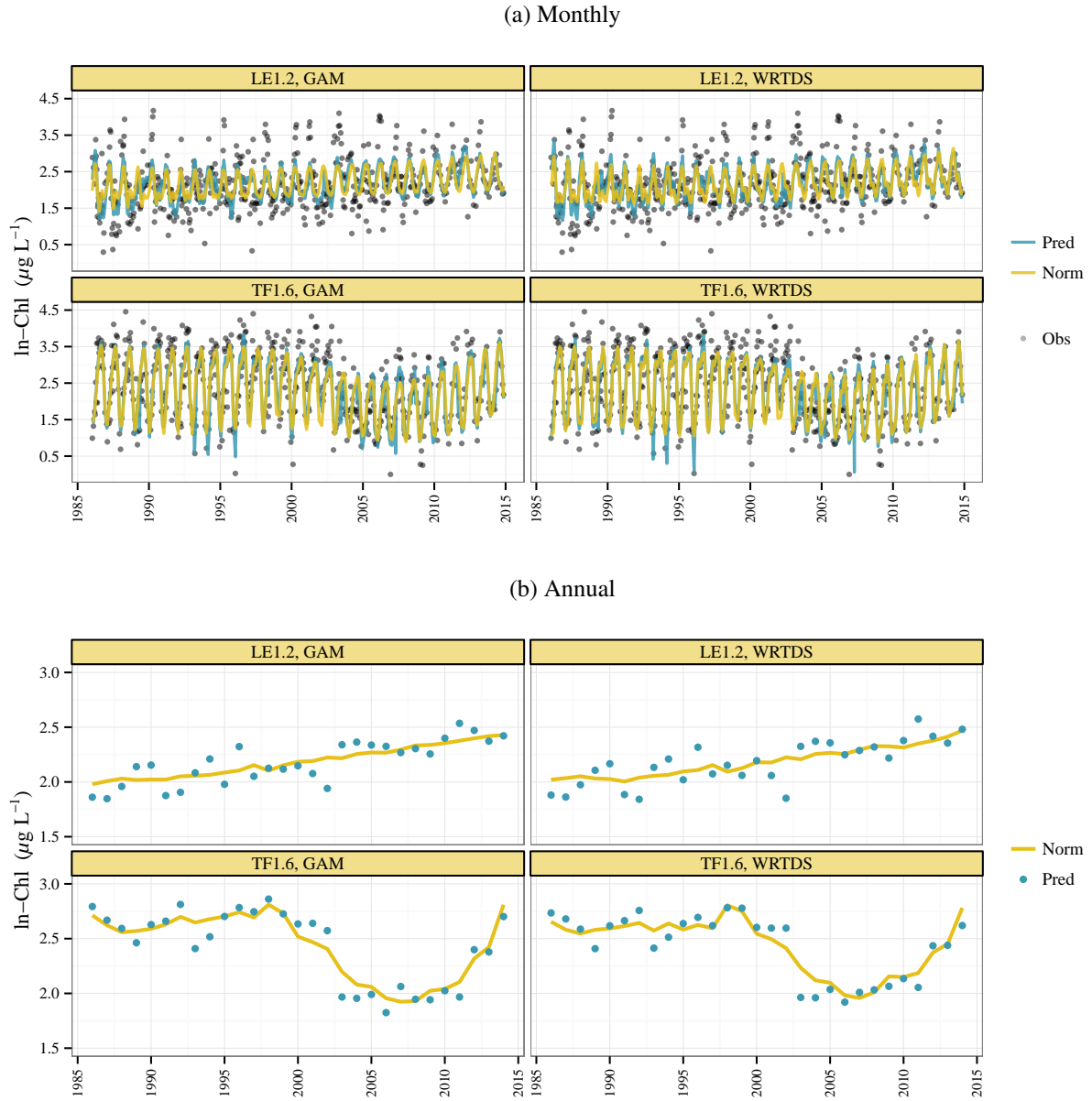


Fig. 3: Predicted chlorophyll from generalized additive models (GAM) and weighted regression (WRTDS) for LE1.2 and TF1.6 stations on the Patuxent River estuary. Fig. 3a shows results at monthly time steps and Fig. 3b shows results averaged by year. Values in blue are model predictions and values in yellow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

{fig:pred}

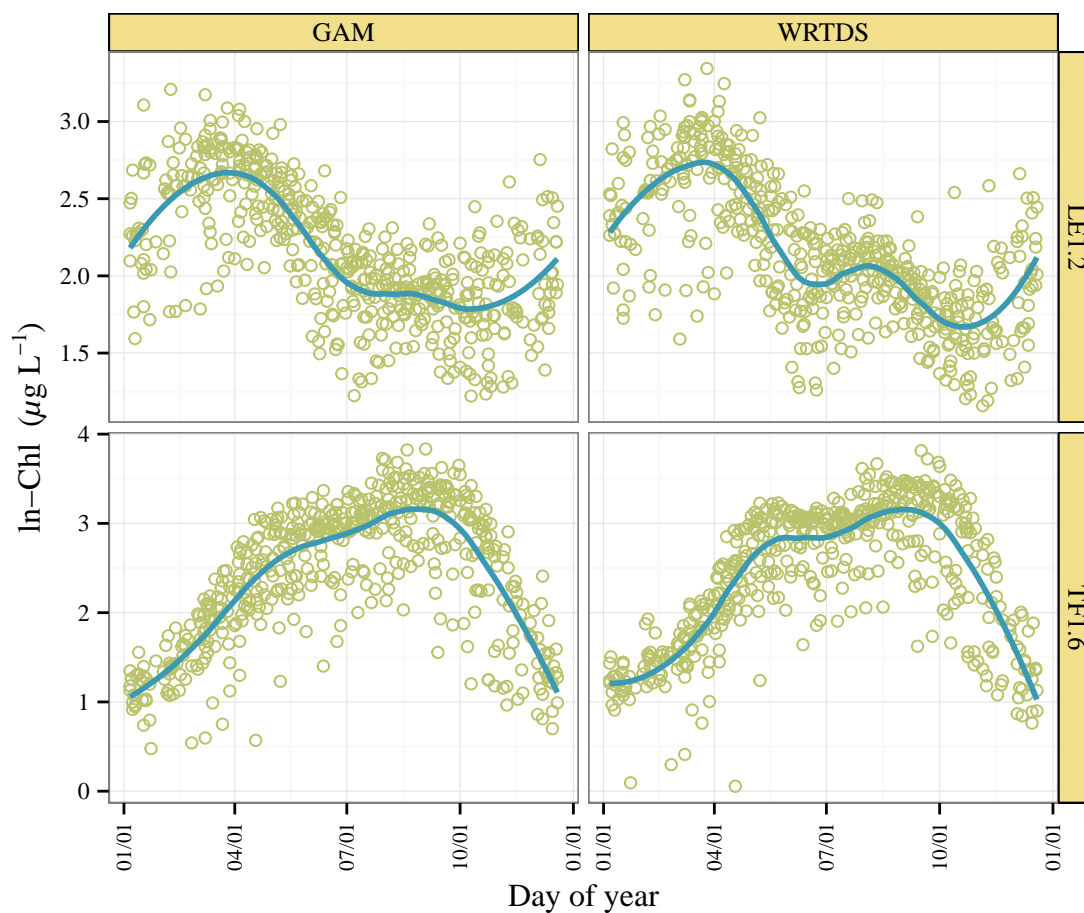


Fig. 4: Seasonal variation of predictions by station and model. Points are all model predictions by day of year from 1986 to 2014. The blue line is a loess (locally estimated) polynomial smooth to characterize the seasonal components.

{fig:seas}

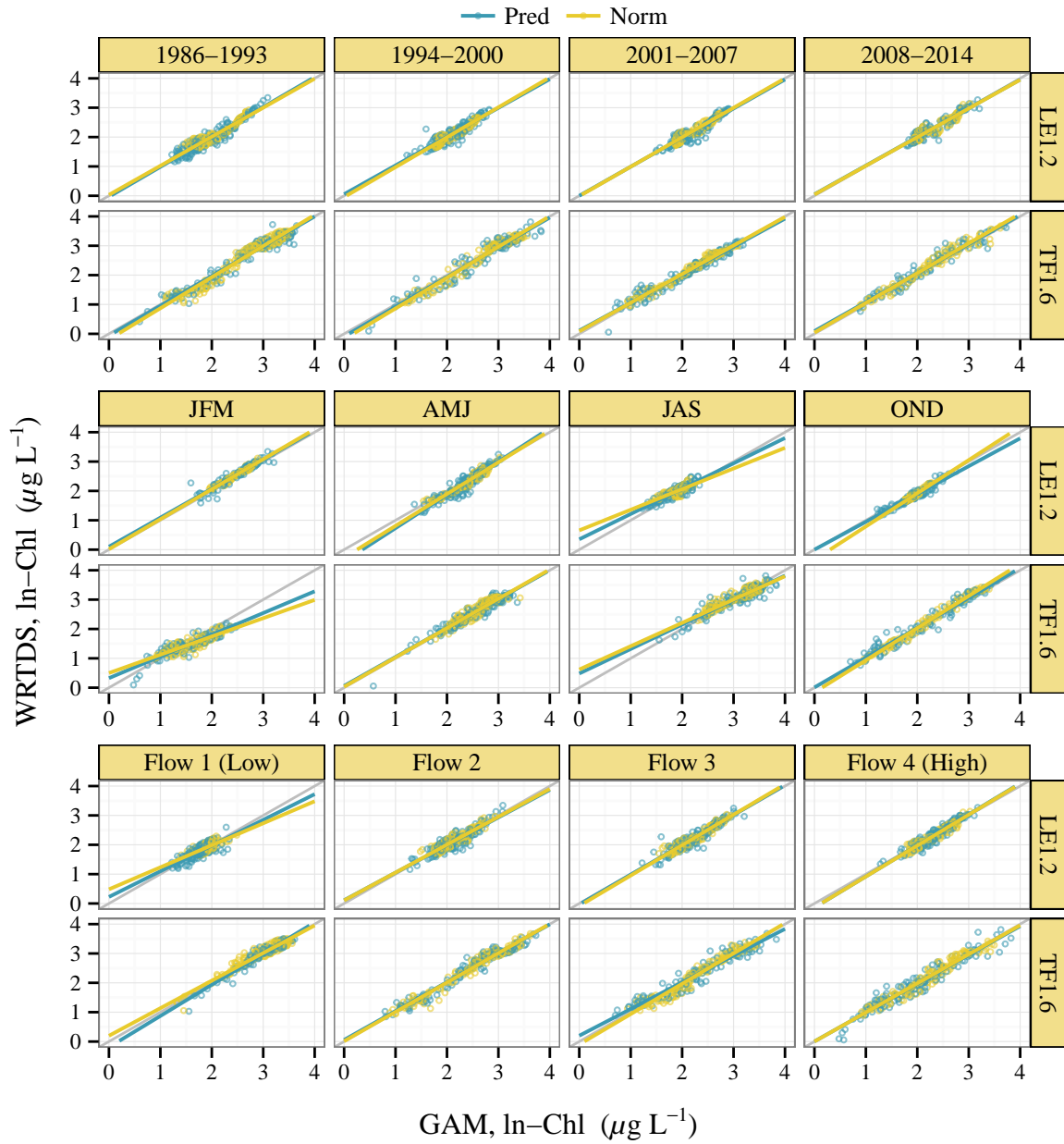


Fig. 5: Comparison of WRTDS and GAMs results at each station (LE1.2, TF1.6) and different time periods. Predicted and flow-normalized results are shown. Time periods are annual groupings every seven years (top), seasonal groupings by monthly quarters (middle), and flow periods based on quantile distributions from the discharge record (low). Regression lines for each model result and 1:1 replacement lines (thin grey) are also shown. See Table 6 for parameter estimates of regression comparisons.

{fig:regp}

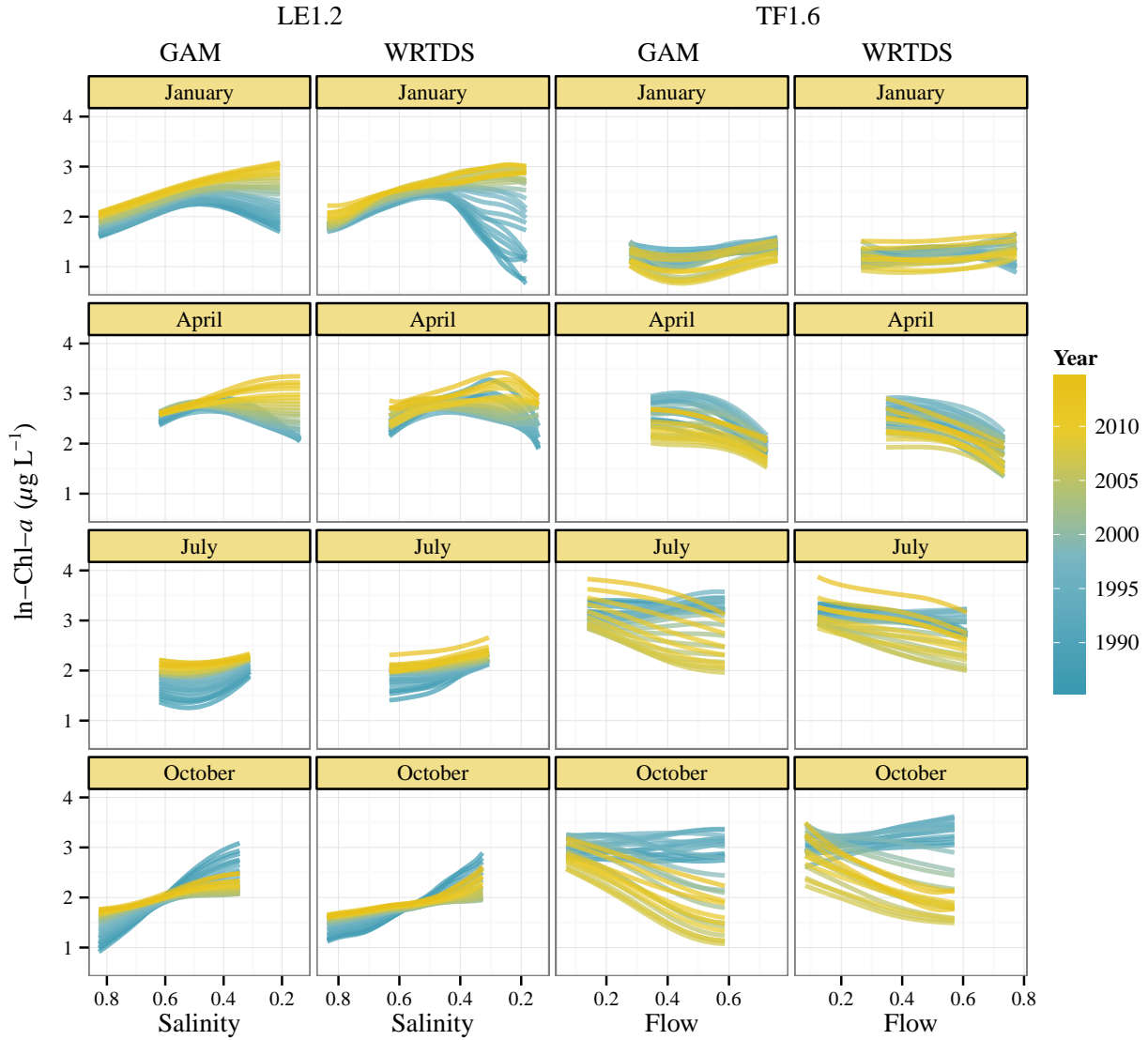


Fig. 6: Changes in the relationship between chl-*a* and freshwater inputs (salinity decrease, flow increase) across the time series. Separate panels are shown for each station (LE1.2, TF1.6), model type (GAM, WRTDS), and chosen months. Changes over time are shown as different predictions for each year 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. The scales of salinity and flow are reversed for comparison of trends. Units are proportions of the total range in the observed data with values in each plot truncated by the monthly 5th and 95th percentiles.

{fig:dyna

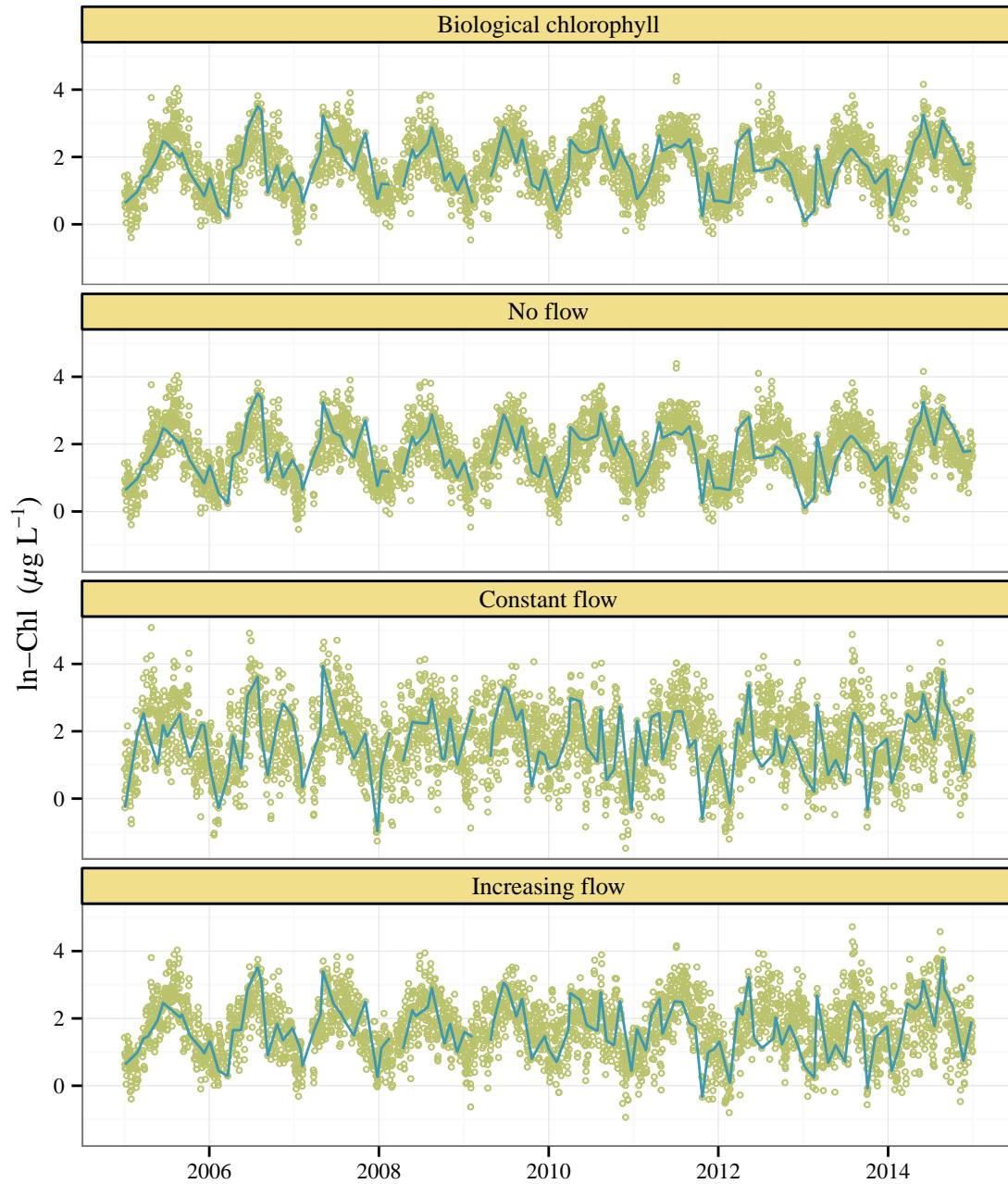


Fig. 7: Examples of simulated time series for evaluating flow-normalized results from WRTDS and GAMs. The plots show the simulated daily time series (points) and monthly samples (lines) from the daily time series used to evaluate the flow-normalized predictions from WRTDS and GAMs. From top to bottom, the time series show the biological chl-*a* independent of flow and the three simulated datasets that represent different effects of flow: none, constant, and increasing effect. The flow-normalized results for the simulated monthly time series from each model were compared to the first time series (biological chlorophyll) that was independent of flow.

{fig:sime2}

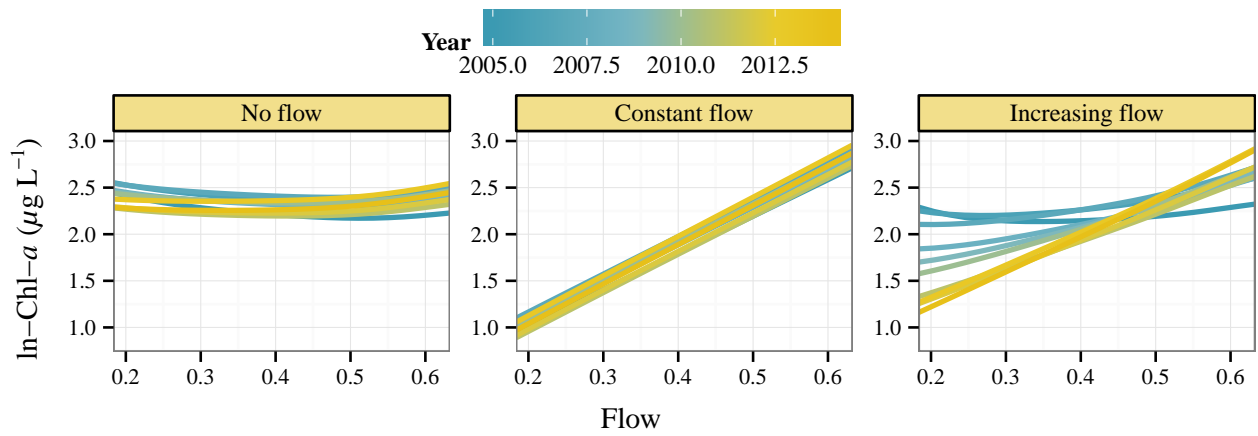


Fig. 8: Examples of changing relationships between chl-*a* ($\mu\text{g L}^{-1}$) and flow (as proportion of the total range) over time (2005–2015) for each simulated time series in Fig. 7. The plots are based on August predictions from three WRTDS models for each time series to illustrate the simulated relationships between flow and chlorophyll.

{fig:dyna

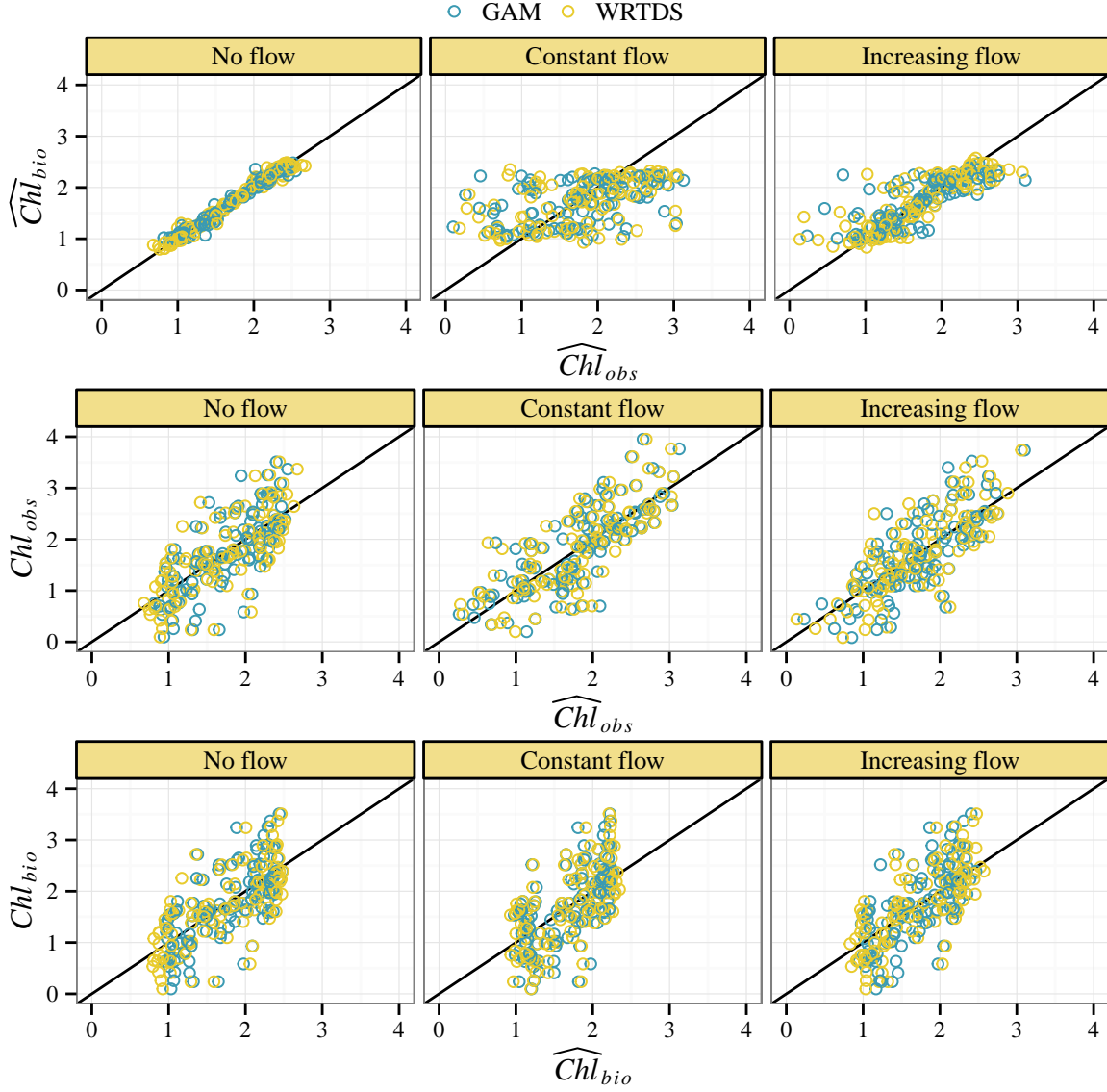


Fig. 9: Model predictions for three simulated datasets with different flow contributions (none, constant, increasing). Estimated variables (e.g., \widehat{Chl}_{bio}) are compared to simulated variables (e.g., Chl_{bio}) to evaluate the ability of each model (GAMs and WRTDS) to recreate the flow-normalized time series of chlorophyll (i.e., bottom plot, \widehat{Chl}_{bio} vs Chl_{bio}) after removing a simulated flow component from the observed chlorophyll time series (Chl_{obs}).

{fig:simr