

Science of the Total Environment

Multi-scale trend analysis of water quality using error propagation of generalized additive models

--Manuscript Draft--

Manuscript Number:	STOTEN-D-21-10662R1
Article Type:	Research Paper
Keywords:	chlorophyll; generalized additive models; Meta-analysis; San Francisco Estuary; Trend analysis
Corresponding Author:	Marcus W Beck Tampa Bay Estuary Program UNITED STATES
First Author:	Marcus W Beck
Order of Authors:	Marcus W Beck Perry de Valpine Rebecca Murphy Ian Wren Ariella Chelsky Melissa Foley David Senn
Abstract:	<p class="FirstParagraph" style="line-height:200%">Effective stewardship of ecosystems to sustain current ecological status or mitigate impacts requires nuanced understanding of how conditions have changed over time in response to anthropogenic pressures and natural variability. Detecting and appropriately characterizing changes requires accurate and flexible trend assessment methods that can be readily applied to environmental monitoring datasets. A key requirement is complete propagation of uncertainty through the analysis. However, this is difficult when there are mismatches between sampling frequency, period of record, and trends of interest. Here, we propose a novel application of generalized additive models (GAMs) for characterizing multi-decadal changes in water quality indicators and demonstrate its utility by analyzing a 30-year record of biweekly-to-monthly chlorophyll-a concentrations in the San Francisco Estuary. GAMs have shown promise in water quality trend analysis to separate long-term (i.e., annual or decadal) trends from seasonal variation. Our proposed methods estimate seasonal averages in a response variable with GAMs, extract uncertainty measures for the seasonal estimates, and then use the uncertainty measures with mixed-effects meta-analysis regression to quantify inter-annual trends that account for full propagation of error across methods. We first demonstrate that nearly identical descriptions of temporal changes can be obtained using different smoothing spline formulations of the original time series. We then extract seasonal averages and their standard errors for an <i>a priori</i> time period within each year from the GAM results. Finally, we demonstrate how across-year trends in seasonal averages can be modeled with mixed-effects meta-analysis regression that propagates uncertainties from the GAM fits to the across-year analysis. Overall, this approach leverages GAMs to smooth data with missing observations or varying sample effort across years to estimate seasonal averages and meta-analysis to estimate trends across years. Methods are provided in the <i>wqtrends</i> R package.</p></p>
Response to Reviewers:	Please view the response to reviewers included in our submission.



July 28th, 2021

Dr. Damià Barceló, Dr. Jay Gan, Dr. Philip Hopke
Co-Editors-in-Chief
Science of the Total Environment

We are pleased to resubmit our manuscript, "Multi-scale trend analysis of water quality using error propagation of generalized additive models" to be considered as an original research paper in Science of the Total Environment. We sincerely thank the two reviewers for providing thoughtful comments on the written text and proposed methods. Our response to reviewers is included with this resubmission as a point-by-point response to each comment provided. We are confident that these changes have improved our manuscript and are hopeful that is now appropriate for publication.

Sincerely,

A handwritten signature in black ink, appearing to read "Marcus W. Beck".

Dr. Marcus W. Beck
Program Scientist
Tampa Bay Estuary Program

TAMPA BAY ESTUARY PROGRAM

263 13TH AVENUE SOUTH; SUITE 350; ST. PETERSBURG, FL 33701; (727) 893-2765; FAX (727) 893-2767; WWW.TBEP.ORG
POLICY BOARD: HILLSBOROUGH COUNTY, MANATEE COUNTY, PINELLAS COUNTY, PASCO COUNTY, CITY OF CLEARWATER, CITY OF ST. PETERSBURG,
CITY OF TAMPA, FLORIDA DEPARTMENT OF ENVIRONMENTAL PROTECTION, SOUTHWEST FLORIDA WATER MANAGEMENT DISTRICT,
U.S. ENVIRONMENTAL PROTECTION AGENCY.

1
2
3
4
5
6 **Multi-scale trend analysis of water quality using error**
7 **propagation of generalized additive models**

8
9
10 Marcus W. Beck* (mbeck@tbep.org), Tampa Bay Estuary Program, St. Petersburg, FL

11
12 Perry de Valpine (pdevalpine@berkeley.edu), University of California Berkeley, Berkeley, CA

13
14
15 Rebecca Murphy (rmurphy@chesapeakebay.net), University of Maryland Center for
16 Environmental Science, Annapolis, MD

17
18 Ian Wren (ianw@sfei.org), San Francisco Estuary Institute, Richmond, Ca

19
20 Ariella Chelsky (ariellac@sfei.org), San Francisco Estuary Institute, Richmond, Ca

21
22 Melissa Foley (melissaf@sfei.org), San Francisco Estuary Institute, Richmond, Ca

23
24 David B. Senn (davids@sfei.org), San Francisco Estuary Institute, Richmond, Ca

25
26 *Corresponding author

27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

We sincerely thank the reviewers for providing thoughtful comments on our manuscript. We have provided a point-by-point response to these comments below. We are confident that these additions and revisions improve our manuscript for readers of Science of the Total Environment.

In addition to making revisions that address points raised by the reviewers, we have also taken the opportunity to make minor improvements that we identified.

Reviewer 1

This manuscript is a welcome addition in the world of trend analyses. The authors propose combinations of GAM and mixed-effects meta-analysis regression to account for uncertainty. The methodology is sound and results are convincing. My recommendation would be even more enthusiastic if the authors had formatted their methods sections in a more traditional manner. I am of the opinion that models should be presented using clear mathematical equations and not in R code semantics. This comments and a few other can be found in the attached document.

- **Response:** We really appreciate these comments that affirm our opinion on the value of this work for trend analysis of long-term water quality data. We have changed our formulas and text to a conventional format, as opposed to R code notation. Briefly, here are the updated equations for the models described in the methods. All other instances in the main text referring to R code (except in a few locations) have been changed. Also please note that the attached pdf that was included with the reviewer's response did not include any markup comments. This is contrary to what the reviewer indicated.

$$\text{Model S: } y_i \sim \beta_0 + f_1(\text{cont_year}_i) + \epsilon_i \quad (1)$$

$$\text{Model SY: } y_i \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + \epsilon_i \quad (2)$$

$$\text{Model SYD: } y_i \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + f_2(\text{doy}_i) + \epsilon_i \quad (3)$$

$$\begin{aligned} & \text{Model SYDI: } y_i \\ & \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + f_2(\text{doy}_i) + f_3(\text{cont_year}_i, \text{doy}_i) + \epsilon_i \end{aligned} \quad (4)$$

Reviewer 2

Overview

This manuscript presents a new application of GAMs approaches to fit irregular time series data (gaps, changing frequency of measurements) to estimate trends. The paper presents a case study using Chl a data from San Francisco Bay. The manuscript is presented as a general approach with applications to other similar environmental time series data. A comparison of methods shows that the GAMs plus meta-analysis method provides some different trend conclusions from ordinary least squares and plain GAMs. The paper is generally well written and well-suited for the audience for STOTEN. I have some clarifications that would help improve/clarify the utility and readability of the manuscript. These are presented below.

- **Response:** Thank you for your comments on our manuscript. We have addressed your concerns on the utility and readability of our manuscript below.

Utility (1):

We applied the wqtrends code in R to a time series data set of X2 (salinity intrusion measure) over approximately in San Francisco Bay. Some observations:

-the code did not work with the full time series of daily X2 data. But we were able to make it work using a random sampling of 5 points per month. Could the reason for this error be explored? As time series data are concerned, we provided an input of ~100 years of daily data. A sample result for the trends are shown in the figure below. I think it would be helpful for readers/users to understand what the practical limitations of the code/algorithm are.

- **Response:** Very rarely do reviewers actually apply the proposed methods on novel datasets as a proof of concept. For that, we are very appreciative of your efforts in exploring the R package. These applications are important for testing both the method and software. As such, you exposed an important limitation of our method that was not clearly articulated in the draft. We provide two responses to this limitation.

First, the X2 dataset is at the daily scale, so one aspect of the GAM smoothing methods proposed in our manuscript for “irregularly spaced or missing” monitoring data is unnecessary. The GAMs provide two benefits: smoothing irregularly spaced data and estimating the non-independence of sequential data so that propagated uncertainty does not incorrectly assume sequential data are independent. The modelling approach is designed to provide both a continuous (e.g., daily) estimate of the long-term trend and a measure of uncertainty in the absence of more regularly collected values. So, application of our method to datasets with sampling intervals longer than the daily scale would be the most compelling use case, but they are nevertheless relevant for uncertainty propagation even with daily data. However, because some of the daily X2 values are in fact imputed, it would be appropriate to use the raw data only in our method, and possibly to enrich the GAM with whatever domain-specific model is used for smoothing X2 values, in order to fully propagate uncertainty.

Second, and more important, the practical limitations of applying a GAM with flexible smoothing to nearly 100 years of daily data is beyond the computational power of most desktop computers. This is asking too much of the smoothing splines used in the mgcv package. Given this limitation, we do feel it is important to clearly state the intended use case for these models to avoid confusion in application. We have added clarification both in the manuscript and the package vignette to clearly indicate the type of data that is appropriate for the methods. In practice, if one wants to analyze such a long series, one could break it down into multiple (possibly overlapping) shorter series for the first two stages (GAM estimation and seasonal averaging with uncertainty propagation). The only difference would be that the estimated optimal degree of smoothness might change across different time windows, but that may be appropriate and supported by having so much data.

In the manuscript, lines 194 - 196: “Approximate monthly or biweekly sampling with coverage of at least a decade is common for many long-term monitoring programs and is the motivating use case for the methods herein.”

In the package vignette, first paragraph (<https://tbep-tech.github.io/wqtrends/articles/introduction.html>): “These models are appropriate for data typically from surface water quality monitoring programs at roughly monthly or biweekly collection intervals, covering at least a decade of observations (e.g., Cloern and Schraga 2016). Daily or continuous monitoring data covering many years are not appropriate for these methods, due to computational limitations and a goal of the analysis to estimate long-term, continuous trends from irregular or discontinuous sampling.”

We have also added a paragraph in the discussion regarding limitations of the approach, one of which is related to the issues above:

Lines 532 - 541: “Several limitations of the proposed methods deserve mention. First, if sampling is so irregular that important fluctuations are missed entirely in some years, the GAM estimates and uncertainty propagation could become dubious in interpretation and usefulness. Second, estimation of GAMs for very long series can be computationally demanding. When this is an obstacle, one could do the first two analysis stages using temporal windows of the full data, with the only implication being that different degrees of smoothness may be estimated for different windows, which indeed might be justified by the data. Third, meta-analysis regression results for a very small number of years, particularly confidence intervals and associated p-values, may be inaccurate (e.g., in confidence interval coverage). In such cases, one could make alternative use of the GAM seasonal averages and standard errors, such as for pairwise comparisons among years.”

Utility (2):

From figure 7 it seems that there are very different conclusions for Station 36 depending on the method used. From Figure 8, however, it seems that Station 36 is a bit of an anomaly in this regard, and that other 8 stations in the case study have findings on trends that are similar among the three methods—and that the results are not dramatically different. Perhaps this should be the basis for added discussion in the final section—when is the added complexity of the GAMS and meta analysis method is appropriate in a real-world setting? Alternatively, it would be helpful if the authors were to examine another data set from the region that they work with, to review the relative findings of trend significance and what this depends on.

- **Response:** You are correct to note that many of the comparisons had similar results and we agree that it is useful to clarify that our approach is warranted relative to more conventional methods. If a simpler approach is used, there is no way to know for sure that an invalid result is obtained without applying our proposed methods. As such, we argue that using our approach is still the wiser option in all cases. Otherwise, a user cannot be certain the correct conclusions are obtained using simpler methods. The paragraph on lines 498 to 513 was revised to address these points:

“Incorrect conclusions on trends can have dramatic consequences for regulated parties under existing water quality compliance frameworks (Smith et al., 2001). Our examples in Figures 7 and 8 demonstrate these risks if propagation of uncertainty from raw observations across methods is unaccounted for in trend assessment. The ‘naïve’ method using OLS regression applied to seasonal averages from the raw observations fails to propagate uncertainty, similarly to averaging results within a year and applying a simple

Kendall test. In some cases the results may be similar to those from fully propagating uncertainty, but the loss of information can lead to increased Type I or II error rates depending on characteristics of the raw data and the method used for their evaluation (Shabman and Smith, 2003). Our examples demonstrated the increased potential for incorrect conclusions at specific monitoring locations and, at larger spatial scales, across all stations if simpler trend analyses are used. Even though simpler methods may produce similar results in some cases, particularly with frequent sampling and similar effort between years, the only way to confirm such an outcome would be to compare results, relying on the method with full propagation of uncertainty to be the more robust method. Use of methods that fully account for uncertainty is recommended to obtain statistically valid results in a wider range of conditions.”

We also want to mention that we recently discovered a minor bug in the seasonal feature extraction of wqtrends that prevented an estimate for years on the tail ends of the time series with incomplete data. Because of this, we have updated Figure 7 and 8, now using stations 30 and 34 as examples in the manuscript. None of the conclusions change in the manuscript.

Readability/Clarity:

The term meta-analysis seems to be an important part of the paper and is referred to multiple times, but it is never clearly explained in the manuscript. For example, at line 164. See the following sentences.

“Third, we used a mixed-effects meta-analysis to estimate trends and test hypotheses about the change in seasonal averages across years. While 166 meta-analysis methods arose from analyses of results from multiple studies, their distinguishing characteristic is propagation of uncertainty (Gasparrini et al., 2012; Sera et al., 2019). Meta-analysis uses response data that includes standard errors (uncertainties) as needed to address our questions.”

It is not clear to me what exactly is meant by meta-analysis and what is being tested in the context of the time series data. An explanation that is understandable by a general audience of potential users of this method would be helpful. It does seem likely, as that sentence points out, that most readers would be more familiar with meta-analysis methods in the context of aggregating the results of multiple studies. There are other helpful citations afterwards, but a more explicit connection between the familiar use case and the use case here of combining different seasonal features from a single regression model would have been helpful for reading this work.

- **Response:** Thank you for pointing out this ambiguity. We agree that a distinction between the more common understanding of “meta-analysis” and the one applied here needs to be made early in the manuscript. We have added content in the introduction to clarify our methods.

Lines 125 - 135: “Meta-analysis regression incorporates a known (or estimated) standard error for each response datum. Usually meta-analysis is used when each response summarizes a dataset from a separate study, along with its standard error, and the meta-analysis looks for across-study patterns in effect size (i.e., Lortie et al. 2014). Here each response summarizes one year or season of data, with standard error from the GAM, and

the meta-analysis looks for patterns across years. Thus, while meta-analysis methods are most commonly associated with combining results from multiple studies into a larger analysis, their key modeling step is propagation of uncertainty (Gasparrini et al., 2012; Sera et al., 2019). To do this, meta-analysis makes use of a known (estimated) standard error for each response datum, which is the required priority here to propagate standard errors from the GAM into a regression of seasonal averages.”

1 Abstract

2 ~~Accurate and flexible trend assessment methods are valuable tools for describing historical~~
3 ~~changes in~~Effective stewardship of ecosystems to sustain current ecological status or mitigate
4 ~~impacts requires nuanced understanding of how conditions have changed over time in response~~
5 ~~to anthropogenic pressures and natural variability. Detecting and appropriately characterizing~~
6 ~~changes requires accurate and flexible trend assessment methods that can be readily applied to~~
7 environmental monitoring datasets. A key requirement is complete propagation of uncertainty
8 through the analysis. However, this is difficult when there are mismatches between ~~time~~
9 ~~scale~~sampling frequency, period of ~~monitoring data record~~, and trends of interest. Here, we
10 propose a novel application of generalized additive models (GAMs) ~~to model seasonal and for~~
11 ~~characterizing~~ multi-decadal changes in water quality indicators and demonstrate its utility by
12 ~~analyzing a long term monitoring dataset~~30-year record of biweekly-to-monthly chlorophyll-a
13 concentrations in the San Francisco Estuary. GAMs have shown promise in water quality trend
14 analysis to separate long-term (i.e., annual or decadal) trends from seasonal variation. Our
15 proposed methods estimate seasonal averages in a response variable with GAMs, extract
16 uncertainty measures for the seasonal estimates, and then use the uncertainty measures with
17 mixed-effects meta-analysis regression to quantify inter-annual trends that account for full
18 propagation of error across methods. We first demonstrate that nearly identical descriptions of
19 temporal changes can be obtained using different smoothing ~~splines for annual or seasonal~~
20 ~~components~~spline formulations of the original time series. We then extract seasonal averages and
21 their standard errors for an *a priori* time period within each year from the GAM results. Finally,
22 we demonstrate how across-year trends in seasonal averages can be modeled with mixed-effects

23 meta-analysis regression that propagates uncertainties from the GAM fits to the across-year
24 analysis. Overall, this approach leverages GAMs to smooth data with missing observations or
25 varying sample effort across years to estimate seasonal averages and meta-analysis to estimate
26 trends across years. Methods are provided in the *wqtrends* R package.

27 *Key words:* chlorophyll, Generalized Additive Models, meta-analysis, San Francisco Estuary,
28 Trend analysis

29 **Introduction**

30 Accurate quantification of environmental trends must consider variation at different temporal
31 scales when ignoring variation at one scale could lead to incorrect conclusions about variation at
32 another scale. Many environmental monitoring programs collect temporally resolved but
33 irregular time series data to quantify trends for regulatory, management, or research purposes.
34 The mismatch between the scales of monitoring versus analysis questions or management goals
35 can present statistical challenges ([Cumming et al., 2006](#); [Forbes and Xie, 2018](#); [Urquhart et al.,](#)
36 [1998](#)). At short temporal scales typically less than a year, environmental systems exhibit
37 variability caused by multiple factors (e.g., weather events, management, or seasonal changes).
38 Such fluctuations may not be related to inter-annual trends or may not be well-suited to multi-
39 scale smoothing methods-, yet they can perturb the time-series such that the sub-annual
40 fluctuations must be addressed within the trend analysis to allow accurate quantification of inter-
41 annual trends. Many trend analysis methods lack the flexibility to evaluate one to manymultiple
42 independent variables in an extendable structure that accommodates hypothesis testing at
43 different time scales of interest.

44 In this paper, we develop methods to estimate across-year trends of within-year features, such as
45 seasonal averages, while accounting for uncertainties across analysis steps. Our overarching goal
46 for this work was to develop a flexible set of tools for accurately characterizing inter-annual
47 changes in seasonally-averaged water quality metrics that can be robustly applied to diverse
48 time-series data. At the outset, we identified several specific requirements and priorities for the
49 trend analysis methods: 1) Complete propagation of uncertainty through the analysis; 2)
50 Separation of trends on different time scales; 3) Ability to estimate linear and nonlinear
51 responses; 4) Flexibility to evaluate multiple independent variables (and random effects),
52 although the examples herein include time as the only independent variable; and 5) Robust to
53 missing observations or varying sampling effort across years, which was considered a high
54 priority to allow the methods to be applied to diverse time-series datasets and monitoring
55 programs. Existing methods that begin to address some of the above requirements and priorities
56 can be generalized into four groups: seasonal Kendall tests (and other non-parametric tests);
57 seasonal trend decomposition using loess (STL); weighted regression on time, discharge, and
58 season (WRTDS); and generalized additive models (GAMs).

59 ~~Existing methods that begin to address our objectives in water quality trend analysis can be~~
60 ~~generalized into four basic approaches: seasonal Kendall tests, seasonal trend decomposition~~
61 ~~(STL), weighted regression on time, discharge, and season (WRTDS), and generalized additive~~
62 ~~models (GAMs).~~ Seasonal Kendall tests ~~or and~~ related non-parametric approaches have been used
63 for decades in water quality trend assessments to identify inter-annual, monotonic changes ~~over~~
64 ~~several years~~ while accounting for ~~the~~ predictable patterns among seasons ([Helsel et al., 2020](#)
[Cloern et al., 2007](#); [Helsel et al., 2020](#); [Hirsch et al., 1982](#)). [Wan et al. \(2017\)](#) showed that
66 While seasonal Kendall and other non-parametric approaches have been among the most

← Formatted: First Paragraph

67 commonly used methods in long-term water quality trend analysis ~~despite critical~~
68 ~~limitations (Wan et al., 2017), they do not satisfy several of our requirements.~~ For descriptive
69 decomposition of long-term monitoring data, ~~they~~~~these approaches~~ assume seasonal patterns
70 within years do not change, ~~and~~ require regularly spaced or balanced data. ~~In addition, seasonal~~
71 ~~Kendall tests~~ do not ~~include~~~~allow for~~ additional ~~predictors~~~~independent variables~~ to explain
72 variation, ~~and~~ do not estimate a model that could be useful for other purposes. ~~Thus, while these~~
73 ~~non parametric approaches have some degree of robustness to assess magnitude (e.g.,~~
74 ~~prediction), and direction do not easily allow for propagation of trends, they apply only to narrow~~
75 ~~goals uncertainty to other trend analysis methods.~~

76 ~~The seasonal trend decomposition using loess (STL)~~ STL decomposes a time series into additive
77 components of a long-term trend, a seasonal pattern, and residuals (Cleveland et al., 1990;
78 Cloern, 2018; Cloern and Jassby, 2010; Stow et al., 2015). While useful and widely applied, this
79 method ~~also has important limitations, does not address all of our requirements.~~ STL
80 decomposition does not ~~incorporate~~~~allow for incorporating~~ explanatory variables ~~besides~~~~other~~
81 ~~than~~ time. ~~In addition,~~ it is ~~defined~~~~often characterized~~ more as an algorithm of statistical steps
82 than as a ~~coherent~~ statistical model ~~with estimated parameters~~ (e.g., Wan et al., 2017) and it
83 does not usually estimate standard errors to allow hypothesis testing (but see Hafen, 2010). STL
84 methods may also over-simplify trends into ~~stationary~~~~fixed~~ components that do not change over
85 time, e.g., a seasonal estimate that is constant across years. This limitation presents challenges
86 when addressing questions relevant to long-term water quality data, such as timing of seasonal
87 peaks that can suggest system response to changing environmental conditions (Cloern and
88 Jassby, 2010; Navarro et al., 2012).

89 The weighted regression on time, discharge, and season (WRTDS) method addresses the
90 problem of inflexibility in STL by using a more general local regression scheme (Beck et al.,
91 2018; Beck and Hagy, 2015; Hirsch et al., 2010); Hirsch et al., 2015). Designed for evaluating
92 water quality in rivers where separating the effect of discharge on constituent concentration is
93 important, WRTDS estimates a moving window regression model with components that allow
94 parameters to vary smoothly in relation to both time and discharge. This yields parameters that
95 are specific to season, year, and flow regime. The WRTDS approach is conceptually similar to
96 local kernel smoothing methods, with specific application to explanatory variables relevant for
97 water quality constituents (i.e., season, year, and discharge). Standard error estimates of
98 predictions from WRTDS are available through a block bootstrap approach applied to the model
99 results (Hirsch et al., 2015). Although a useful addition to the original method (Hirsch et al.,
100 2010), the approach requires extensive resampling using a previously fitted model. Alternative
101 methods that include standard error estimates simultaneously with model output may be
102 preferred for intensive or more iterative applications.

103 Generalized additive models (Finally, GAMs) are can satisfy the requirements and priorities
104 identified above and were adopted as a central to this paper and form the basis of component for
105 the fourth method to separate fluctuations on different time scales.trend analyses herein. GAMs
106 combine one or more smoothing splines to model patterns in data and may can be seen reasonably
107 viewed as generalizing the concepts behind STL and WRTDS (Haraguchi et al., 2015; He et al.,
108 2006; Morton and Henderson, 2008; Murphy et al., 2019; Pearce et al., 2011). The basis
109 functions used to formulate GAMs can be customized based on expected patterns in the data.
110 One example includes Examples include cyclic splines, which can be used to model seasonal
111 patterns, and low-dimensional interactions. (Wood, 2017). GAMs have added flexibility because

112 they can include both parametric (e.g., linear or quadratic) components and non-parametric
113 (spline) components. Multiple approaches have been developed to determine the optimal degree
114 of smoothness-[\(Wood, 2004; 2017\)](#). These approaches are based on optimization of out-of-
115 sample prediction error, which addresses a key concern around methods like WRTDS that do not
116 have analogs for choosing optimal degrees of smoothing. GAMs can also produce [results](#)
117 comparable ~~results similar~~ to those provided by WRTDS ([Beck and Murphy, 2017](#)) and have
118 readily obtainable uncertainty estimates. Further, GAMs have natural frequentist and Bayesian
119 interpretations, are naturally extensible to include random effects (i.e., generalized additive
120 mixed models or GAMMs), and have computationally efficient implementations ([Wood, 2017](#)).

121 GAMs have been applied previously to evaluate trends in water quality time series from long-
122 term monitoring programs ([Haraguchi et al., 2015](#); [Murphy et al., 2019](#)). For example, [Murphy et](#)
123 [al. \(2019\)](#) used GAMs to decompose water quality time series from Chesapeake Bay into long-
124 term and seasonal trends ([Murphy et al., 2019](#)) and test trend hypotheses between two points in
125 time. Other studies of environmental time series with GAMs have addressed the use of
126 transformed response data ([Yang and Moyer, 2020](#)), serial correlation in high resolution data
127 ([Morton and Henderson, 2008](#); [Yang and Moyer, 2020](#)), and quantifying time lags in
128 relationships between response and predictor variables ([Lefcheck et al., 2017](#)). The ~~study~~[method](#)
129 [development and analyses described](#) herein ~~generalizes~~[generalize](#) the approach to analyzing
130 trends of seasonal spline features, describes the relationships among alternative spline
131 formulations when spline flexibility is allowed to vary ([Wood, 2017, 2003](#)) rather than being
132 constrained *a priori* for different time scales, and prioritizes full incorporation of uncertainty.

133 ~~Our motivating problem has several characteristics that are only partially addressed by previous~~
134 ~~methods and can further build on GAMs as a starting point. Our general goal is to understand~~

135 interannual changes in seasonally averaged water quality metrics, such as chlorophyll. However,
136 the seasonal average within each year must be robust to inconsistent sampling times and
137 intervals, and any trend analysis must consider the uncertainties in seasonal averages. The
138 critical need is the ability to obtain an accurate estimate of uncertainty (e.g., a standard error) of
139 seasonal averages, even with irregular sampling and serial correlation, which is common in time
140 series data. This paper develops the use of GAMs with mixed effects meta-analysis (To
141 incorporate the uncertainty of seasonal estimates into trend analysis, we integrated GAMs with
142 mixed-effects meta-analysis (Gasparrini et al., 2012; Sera et al., 2019). In this integration, the
143 GAMs framework addressed a critical need by providing an estimate of uncertainty (e.g., a
144 standard error) of seasonal averages, even in situations with irregular sampling and serial
145 correlation, which are common in time series data. Meta-analysis regression incorporates a
146 known (or estimated) standard error for each response datum. Usually meta-analysis is used
147 when each response summarizes a dataset from a separate study, along with its standard error,
148 and the meta-analysis looks for across-study patterns in effect size (i.e., Lortie, 2014). In this
149 study, each response summarizes one year or season of data, with standard error from the GAM,
150 and the meta-analysis looks for patterns across years. Thus, while meta-analysis methods are
151 most commonly associated with combining results from multiple studies into a larger analysis,
152 their key modeling step is propagation of uncertainty (Gasparrini et al., 2012; Sera et al., 2019).³
153 to address multi-scale trend analysis questions for which seasonal Kendall tests and the more
154 complex STL and WRTDS methods are not well-suited.

155 We To do this, meta-analysis makes use of a known (estimated) standard error for each response
156 datum, which is the required priority here to propagate standard errors from the GAM into a
157 regression of seasonal averages.

158 To describe the approach and demonstrate the proposed methods by analyzing water quality
159 monitoring data from its utility, we analyze a 30 year record (1990-2019) of biweekly-to-monthly
160 chlorophyll-a concentration data, collected at 9 stations in the southern portion of the San
161 Francisco Estuary, California, USA. Approximately twice-monthly monitoring has been
162 conducted for several decades at fixed locations (stations) ~~on~~along the longitudinal axis of the
163 Bay. Analysis of these data is complicated by irregularities in timing and consistency of data
164 collection, ~~which can generate artifacts affecting simple seasonal averages of the data.~~ We were
165 interested in questions such as: Are there significant trends in spring mean chlorophyll at multi-
166 year time-scales? At what across-year ~~scale does window~~ summer-fall mean chlorophyll levels
167 change? Is there a spatial difference in chlorophyll trends? We provide examples illustrating how
168 these questions can be addressed using GAMs to estimate seasonal ~~trends~~patterns and ~~evaluated~~
169 ~~between years using~~use meta-analysis ~~methods. This approach is new to environmental trend-~~
170 ~~detection problems~~evaluate trends between years. The techniques are incorporated into an open-
171 source and ~~is provided in the wqtrends publicly available~~ R package, wqtrends, developed by the
172 authors (Beck et al., 2021, available at <https://tbep-tech.github.io/wqtrends>), including an online
173 dashboard for viewing results at [https://nutrient-](https://nutrient-data.sfei.org/apps/SFbaytrends/)
174 <https://nutrient-data.sfei.org/apps/SFbaytrends/>).

Formatted: Default Paragraph Font

175 Methods

176 Study area and data sources

177 The San Francisco Estuary (SFE) is the largest estuary on the Pacific Coast of North America.
178 Its, and its watershed covers 200 thousand km² in the US state of California. Major freshwater

179 inputs enter the system through Flows from the Sacramento-San Joaquin Delta-complex
180 upstream of Suisun, entering from the northeast, account for the vast majority of SFE-wide
181 annual-average freshwater inputs (Cloern and Jassby, 2012). Freshwater contributions to
182 southern SFE (South Bay-, Lower South Bay) come primarily from local tributaries during the
183 wet season (Nov-Apr), and from wastewater treatment plant discharges during the dry season
184 (May-Oct). Salinity ranges from 0 to 15 ppt values in the northern southern SFE subembayments
185 and (the focus of this work; Central Bay, South Bay, Lower South Bay) range from 5 to 35 ppt
186 in southern subembayments closer to the Pacific Ocean, depending and depend strongly on the
187 season (stormwater runoff), tidal cycle, and effluent discharge from wastewater treatment plants,
188 and stormwater runoff (Cloern and Jassby, 2012). An estimated 73.8 metric tons dy⁻¹
189 SFE receives 70,000 kg per day of dissolved inorganic nitrogen are discharged into the Bay,
190 primarily from wastewater (Novick and Senn, 2014). Agricultural runoff from (DIN; annual
191 average), with the upper watershed contributes 30 metric tons dy⁻¹ of nitrogen to the SFE via
192 majority of that DIN coming from wastewater treatment plant discharges. Flows from the Delta-
193 Nitrogen and phosphorus levels in deliver 30,000 kg per day of DIN to the SFE (annual average),
194 with the SFE usually exceed concentrations that cause eutrophication in other Delta's DIN load
195 varying over a 5-fold range annually (SFEI, 2014a). Based on its areal DIN loads, SFE ranks
196 among the most nutrient-enriched estuaries. However worldwide (SFEI, 2014a, b; Cloern et al.,
197 2020). Despite its nutrient-enriched status, the SFE has demonstrated resistance to
198 eutrophication, which has been generally not experienced some of the water quality impacts
199 common to other nutrient-enriched estuaries (e.g., excessive phytoplankton blooms, low
200 dissolved oxygen), with SFE's muted response attributed to high concentrations of suspended
201 sediment that reduce its highly turbid waters (reduced light penetration in within the water

202 column, low residence time caused by vigorous river flushing, and removal of primary producers
203 by); strong tidal mixing (limiting duration of water column stratification to less than several
204 days); and strong phytoplankton grazing pressure from abundant suspension feeding bivalves in
205 some regions (Alpine and Cloern, 1988; Cole and Cloern, 1984; Jassby, 2008; Kimmerer and
206 Thompson, 2014; Lehman et al., 2017).).

207 Studies over the past decade have identified changes in responses or sensitivity within SFE to
208 nutrients in deep subtidal habitats via increased phytoplankton biomass (chl-a) and gross primary
209 production (GPP) in South Bay (Cloern et al., 2007, 2010); recently documented occurrences of
210 harmful algae and their associated toxins (Sutula et al., 2017; Peacock et al., 2018); and low
211 dissolved oxygen in some tidal slough habitats (SFEI, 2021). These observations have raised
212 concerns that SFE's resistance to its high nutrient inputs could be waning (SFEI, 2014b),
213 prompting regulators to initiate the SFB Nutrient Management Strategy (SFBRWQCB, 2017).

214 The Regional Water Quality Control Board has showed renewed interest in early increases
215 in South Bay chl-a (1995-2005) were quantitatively tested using seasonal Kendall, and the signal
216 was sufficiently large and coherent it was also visually apparent in raw data. Given SFE's
217 nutrient-enriched status, there is a critical need for on-going and comprehensive characterization
218 of trends in chl-a, including the ability to examine variability at multiple time-scales and non-
219 monotonic trends, that can also be readily applied to other nutrient-related indicators (e.g.,
220 dissolved oxygen, GPP). Beyond the role these tools can play in supporting improved
221 understanding the potential for nutrient loading to negatively affect water quality for more
222 southern areas of the SFE where harmful algal blooms, elevated summer fall of system dynamics
223 in SFE, water quality managers have emphasized the importance of robust trend detection for
224 informing future nutrient management decisions.

Formatted: Font color: Custom Color(RGB(34,34,34)),
Pattern: Clear (White)

Formatted: Font color: Custom Color(RGB(34,34,34)),
Pattern: Clear (White)

225 For the trend analyses discussed below, we used near-surface (0-2 m) chlorophyll concentrations,
226 and low dissolved oxygen concentrations began around 1999 (Figure 1) (Cloern et al., 2020).
227 Although changes in the data are visually apparent, statistical analyses to quantify these changes
228 have been insufficient particularly with respect to seasonal differences between years.

229 We evaluated near surface chlorophyll (chl-a) data measured biweekly to monthly from 1990 to
230 2019 along the longitudinal axis of the SFE extending from Central Bay (stations 18-23), South
231 Bay (stations 24-32), and Lower South Bay (stations 34-36) (Table 1, Figure 2). Monitoring data
232 were obtained from the SFE Research Program of the US Geological Survey (Cloern and
233 Schraga, 2016; Schraga et al., 2020). Sampling frequency varied somewhat over time and by
234 station. Approximate monthly or biweekly sampling with coverage of at least a decade is
235 common for many long-term monitoring programs and is the motivating use case for the
236 methods herein. Every observation was included directly in the statistical models without spatial
237 or temporal binning or averaging. Log₁₀-transformed chl-a was used for all analyses to meet
238 assumptions of normally-distributed residuals. Methods for back-transformation of model results
239 are provided in the supplement.

240 GAMs with uncertainty propagation

241 We implemented our analysis in three stages. First, we used a GAM to estimate a smooth
242 temporal pattern in the raw data, along with its the uncertainty of the smoother. Second, we
243 calculated a feature of interest from the estimated GAM, along with its propagated uncertainty.

244 For this example, the examples described here, we focused on extracting seasonal averages were
245 extracted, whereas other of chl-a values. Other features could be can also be extracted using the
246 same tools, including the timing or magnitude of a seasonal peak, but those are not developed

247 here presented here (see the *wqtrends* R package, <https://tbep-tech.github.io/wqtrends>). Third, we
248 used a mixed-effects meta-analysis to estimate trends and test hypotheses about the change in
249 seasonal averages across years. While meta analysis methods arose from analyses of results from
250 multiple studies, their distinguishing characteristic is propagation of uncertainty (Gasparrini et
251 al., 2012; Sera et al., 2019). Meta analysis uses response data that includes standard errors
252 (uncertainties) as needed to address our questions.

253 First-stage analysis: GAM estimation

254 We considered four different GAMs to smooth the raw data across time. Although they, we
255 considered and tested four different GAM structures. While all four GAM structures can achieve
256 similar fits, they do so by partitioning differ in how they partition variation in the time series
257 differently (Table 2), which may unnecessarily influence understanding of temporal patterns.
258 We discuss all four to clarify their relationships and interpretations. Models are shown in the
259 notation of All models were created using the mgcv R package as formulas for the gam function
260 (R Core Team, 2020; Wood, 2017), with utility functions included in the wqtrends package
261 created by the authors (Beck et al., 2021).

262 The simplest GAM for this purpose is expressed as:

263 Model S: $y \sim s(\text{cont_year}, k = \text{num_knots_Y})$

$$264 \quad \text{Model S: } y_i \sim \beta_0 + f_1(\text{cont_year}_i) + \epsilon_i \quad (1)$$

265 where y is the time series of interest, such as measured chl-a, cont_year is an intercept, and
266 cont_year is “continuous year,” a continuous numerical date (e.g., July 1st 2019 would be
267 2019.5), $y \sim s(\dots)$ indicates that y will be explained by \dots . The $f_1()$ function is a smoothing

Formatted: First Paragraph

268 spline (in this case as a function composed of the sum of `cont_year`, and `num_knots_Y`)
269 is multiple “basis functions” multiplied by coefficients `f1(cont_year)` describes the
270 numberrelationship of knot or “connections” alongy with `cont_year` in a way that smoothly
271 follows the `splinedata` (Wood, 2017). The basis functions involve user-specified knots, a grid of
272 values on the `cont_year` axis that is discussed more below. The `e` term represents residuals
273 following a normal distribution with mean zero and constant variance.

274 Smoothing was determined using generalized cross-validation (GCV, as implemented in `mgcv`),
275 which approximately minimizes out-of-sample prediction error. GCV works by penalizing the
276 net curvature of a spline (Wood, 2004). To allow GCV (or other alternatives) to work as
277 intended, the number of knots that chosen by the analyst, which determine the maximum degrees
278 of freedom chosen by the analyst, must be sufficiently large so that the curvature penalty, rather
279 than the number of knots, determines smoothness. Results should not be sensitive to the number
280 of knots; if they are, the number of knots should be increased. In the examples below, we chose
281 the number of knots, `num_knots_Y, for f1()` as 12 times the number of years in the time series,
282 i.e., one knot per month. If the data were too sparse to fit 12 knots per year, the number of knots
283 was reduced by one knot per year until the model could be estimated (i.e., 12 * years, 11 * years,
284 etc.).

285 The next three spline formulations (Model SY, SYD, and SYDI) provide progressively
286 increasing complexity in how spline terms compose a model to smooth the raw data. Model SY
287 describes the time series using a linear trend plus a spline for `cont_year``cont_year`:

288 Model SY: `y ~ cont_year + s(cont_year, k = num_knots_Y)`
289 This model is Model SY: $y_i \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + \epsilon_i$ (2)

Formatted: Comment Reference

290 where equation (2) is the same as equation (1) with the addition of a linear term for *cont_year*
291 related to y_i by the β_1 slope parameter.

292 While Model SY contains the explicit linear trend term, $\beta_1 cont_year_i$, it is in fact
293 mathematically equivalent to model S (Table 3). The $f_1()$ spline for *cont_year* includes an
294 unpenalized linear trend, so a trend will be estimated in model S. When *cont_year* is included
295 explicitly as a linear term in model SY, `mgcv` adjusts the basis functions for the spline to exclude
296 the linear term, thereby not over-parameterizing the model. Whereas an estimated linear trend in
297 *cont_year* and its uncertainty can be extracted from the fitted spline in model S, model SY
298 provides this trend directly, giving the equivalent result. Further, package `mgcv` ~~can offer the~~
299 ~~option to~~ penalize linear trends in splines to provide a method for variable selection (option
300 `select = TRUE`), such as when numerous splines are included in the model formulation for
301 variables that may or may not be important. ~~For In the implementation of~~ our approach ~~described~~
302 ~~here~~, this option is not used and all models specify `select = FALSE`. Details in the supplement
303 explain this justification.

304 Model SYD adds an average within-year cyclic pattern as a separate spline:

305 ~~Model SYD: $y \sim cont_year + s(cont_year, k = num_knots_Y) + s(doy, bs = 'cc',$~~
306 ~~k = num_knots_D)~~

307 Model SYD: $y_i \sim \beta_0 + \beta_1 cont_year_i + f_1(cont_year_i) + f_2(doy_i) + \epsilon_i$ (3)

308 where ~~doy~~equation (3) is the same as equation (2) with the addition of a smoothing spline for
309 “day-of-year” (*doy*, i.e., Julian date, a count starting January 1 for each year)~~, $bs = 'cc'$~~
310 indicates that the spline will be). For $f_2()$, a cyclic ~~constrained to spline is specified (using $bs =$~~
311 ~~'cc'~~ in `mgcv`) to constrain the start and end at the same value), and *num_knots_D* is the ~~A user-~~

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: First Paragraph

312 specified number of knots ~~for the *doy* spline~~~~is also included in $f_2(0)$~~ . While model SYD is not
313 mathematically equivalent to models S and SY, it should produce nearly identical results. The
314 *doy* spline in model SYD gives the average within-year pattern and changes the interpretation of
315 the *cont_year* spline to represent smoothed deviations from that pattern.

Formatted: Font: Italic

Formatted: Font: Italic

316 Models S, SY, and SYD can all potentially extract a similar signal from the raw data (Table 3).
317 What differs between the models is the allocation of penalties for curvature used to determine
318 smoothness for each spline. In model SYD, there are separate penalties for the two splines, as
319 compared to S and SY that include penalties only for the *cont_year* spline. This is important
320 because variation in the response variable can be differently attributed to each spline depending
321 on the model, even while the sum of components for each model produces similar results
322 between models. Our goal is to extract seasonal averages from the fitted time series, which is that
323 are not sensitive to different allocation of penalties among the splines in each model.

Formatted: Font: Times New Roman, 12 pt, Italic

324 If the fits were to differ substantially between model SYD and models S or SY, an interpretation
325 could be difficult because the penalties for smoothing splines based on curvature are heuristic
326 (Wood, 2017). For example, if a lower AIC is achieved in one model compared to another,
327 assuming both use sufficient knots, this may just reflect the outcome of alternative penalization
328 heuristics implied by the different formulations and does not imply one model fit is better. In the
329 examples here, model SYD achieves nearly identical fits to model S or SY, where the latter by
330 definition also achieve identical fits.

331 Model SYD has the appealing feature that, if some parts of some years have limited data, model
332 SYD will impute an average seasonal pattern with the *doy* spline, thereby considering data from
333 the same period in other years in the prediction of the period with missing data. However, an

Formatted: Font: Times New Roman, 12 pt, Italic

334 interpretation of these imputations may be challenging. For example, the spring chl-a peak is a
335 notable feature every year in the SFE. If the peak occurs at the same time every year but the
336 magnitude varies, then the average within-year pattern can be interpreted as the average
337 magnitude. However, if the magnitude is the same but the timing varies across years, then the
338 magnitude of the average peak cannot be similarly interpreted ~~and instead underestimates;~~
339 instead, the average peak extracted by $f_2(\text{doy}_i)$ will underestimate the magnitude that usually
340 occurs. Moreover, the width or duration of the peak will be longer than typically occurs in a
341 given year.

342 Finally, the raw data can be smoothed using a bivariate spline representing an interaction

343 between *cont_year* and *doy*. This can be expressed as:

344 Model SYDI: $y \sim \text{cont_year} + s(\text{cont_year}, k = \text{num_knots_Y}) + s(\text{doy}, bs = "cc",$
345 $k = \text{num_knots_D}) + ti(\text{cont_year}, \text{doy}, bs = c("tp", "cc"), k =$
346 $c(\text{num_knots_Y_ti}, \text{num_knots_D_ti}))$

347 Model SYDI:

348 $y_i \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + f_2(\text{doy}_i) + f_3(\text{cont_year}_i, \text{doy}_i) + \epsilon_i \quad (4)$

349 where ~~ti()~~ specifies equation (4) is the same as equation (3) with the addition of a tensor-
350 product smoothing spline ~~for a surface(ti() in mgcv)~~ that varies smoothly as a function of both
351 ~~cont_year~~ and ~~doy~~. Both ~~cont_year~~ and ~~doy~~ include their own number of
352 knots, such that the total number of knots for the spline is the product of ~~num_knots_Y_ti~~ in the
353 ~~cont_year axis and num_knots_D_ti~~ in the ~~doy axis~~. In SYDI, the two. The need for
354 sufficient knots in SYDI can be satisfied either by ~~allowing for~~ sufficiently ~~large values for~~
355 ~~num_knots_Y_ti~~ and ~~num_knots_D_ti~~ many knots for $f_3()$ or a sufficiently ~~large value for~~

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: First Paragraph

356 knots in num_knots_Y and num_knots_D; many knots for $f_1()$ or $f_2()$, but not both, given limits
357 on the model degrees of freedom.

358 Following the rationale above, the relationship of model SYDI to model S is similar to that of
359 model SYD to model S. Model SYDI differs from model S to a greater extent than model SYD,
360 but all of the splines use the same inputs to smooth the same data. The univariate splines in

361 cont_year and doy will combined likely not capture as much variation in within model SYDI
362 compared to as captured by model S, given the fewer knots that are available to the former, $f_1()$
363 or $f_2()$ within SYDI. The ~~ti term~~ tensor-product spline represents an interaction by allowing the
364 pattern in cont_year~~cont_year~~ to vary by doy and vice-versa. The interaction term in
365 model SYDI provides an appearance that this model SYDI is fundamentally different from those
366 provided by the other models. However, models S, SY, and SYD also allow within-year
367 fluctuations to vary across years by allowing a spline to be fit through the entire time series.

368 Although model SYDI is the only model that includes an explicit interaction term, all of the
369 models support the interaction conceptually. By providing this term with sufficient knots, the
370 raw data can be fully smoothed with model SYDI to a similar degree as for the other models.

371 However, a very large number of knots in both the cont_year~~cont_year~~ spline and in both
372 dimensions of the interaction spline is impossible to achieve. The distinct aspect of model SYDI
373 is the anticipation that within-year fluctuations will vary smoothly from year to year, which.
374 However, this is unlikely an assumption about ecosystem dynamics that may not be appropriate
375 to a priori parameterize into a statistical model, including for the SFE data and the dynamics in
376 many other estuaries because where bloom size typically magnitude often varies between years.
377 Thus, the conceptual motivation for model SYDI and its practical application are not necessarily

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Times New Roman, 12 pt, Italic

378 supported in the more generalized application for which we are developing this analysisset of
379 analyses.

380 Murphy et al. (2019) used spline formulations for Chesapeake Bay water quality related analyses
381 that are comparable to those proposed here, but for different goals and with different handling of
382 smoothness. They evaluated a “gam0” with only s(a cyclic spline for doy) and linear cont_year
383 terms, a “gam1” like our SYD, and a “gam2” like our SYDI. In application, only “gam2” was
384 used, including the addition of splines as functions of hydrologic variables to account for finer-
385 scale variation. Murphy et al. Murphy et al. (2019) allowed a maximum number of knots infor
386 the s(cont_year)-term spline (f1() of 2/3 times the number of years and do not explicitly
387 consider the number of knots in the interaction spline, following an *ad hoc* allocation of variation
388 in the data to different components based on previous interpretations of water quality dynamics
389 in the system. Constraining splines with insufficient knots could inflate Type I error rates for
390 temporal changes and we seek to lower this risk by increasing the upper limit for the knots for
391 the s(cont_year)f1() term. Finally, Murphy et al. (2019) present large AIC differences
392 between their spline formulations. We instead emphasize that, given sufficient knots, the models
393 represent alternative formulations of conceptually similar explanations for the data and yield
394 similar fits (Table 3), resulting in near ties for AIC between models.

395 WeTo evaluate the theoretical and conceptual similarities and differences among the GAM
396 structures discussed above, we visually comparecompared chl-a estimates from models SYS,
397 SYD, and SYDI to emphasize that similar fits can be achieved by all of the presented models
398 (Figure 3; note that SY is identical to S and is not shown). Models S, SYD, and SYDI were fit to
399 chl-a data from station 34 using large k valuesa sufficiently high number of knots for the
400 arguments num_knots_y, num_knots_D, num_knots_Y_t1, and num_knots_D_t1respective

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Times New Roman, 12 pt, Italic

401 splines for each model. PredictionsAs expected, predictions by day of year from each model are
402 visually similar (Figure 3a) and closely follow the 1:1 line (Figure 3b). However, when
403 contrasting the estimates using only component of predictions explained by the continuous year
404 smoother (s(cont_year),f₁()) differs between models that include additional smoothers
405 (Figure 3c), offering a graphical impression of the fits differ substantially because of how each
406 model allocates variationdegree to the splineswhich data variability was partitioned to different
407 GAM components. These results are also reflected in differences in the effective degrees of
408 freedom among the additive components of each model (Table 3). Accordingly, even though the
409 models differ by which structural component describes variation in the chl-a time series, they
410 provide similar predictions.

411 For all results, model S was used with enough knots in num_knots_yf₁() to evaluate chl-a trends
412 across the monitoring stations in the SFE. This model was chosen because of the relatively faster
413 processing time to fit the model, while providing nearly identical explanatory power as compared
414 to the other models (Table 3).

415 **Second-stage analysis: Seasonal features with uncertainties**

416 In the second-stage analysis, we estimated a seasonal average, such as the mean spring chl-a
417 concentrations, along with the its associated uncertainty, in each year. We defined μ_t as the
418 seasonal average in year t , $\hat{\mu}_t$ as an estimate of μ_t , and $\hat{\sigma}_{\hat{\mu},t}$ as the estimated standard error of $\hat{\mu}_t$.
419 The season includes n days. For simplicity, the following text omits subscript t .

420 Point estimates of response values for the fitted GAM take the form $\hat{y} = \mathbf{X}\hat{\beta}$, where $\hat{\beta}$ is the
421 vector of parameter estimates and \mathbf{X} is a model matrix of explanatory variables, including spline
422 basis function values. Vector $\hat{\beta}$ includes both fixed effect parameters and spline parameters, and

423 \mathbf{X} contains columns corresponding to each. For example, using model SY, if a point estimate for
424 chl-a is needed for a single day, given as *cont_year* = r , then \mathbf{X} would have a row with 1 in the
425 first column (for the intercept parameter), r (for the linear time trend) in the second column, and
426 an evaluation of each spline basis function at r in the remaining columns. The number of spline
427 basis functions is related to the number of knots. Note that r can be any time, not necessarily the
428 time of an observation.

Formatted: Font: Times New Roman, 12 pt, Italic

Formatted: Font: Italic

429 To obtain a vector, $\hat{\mathbf{y}}$, of fitted point estimates for every day in a season, \mathbf{X} would have one row
430 for each day. Here, the seasonal averages used in our examples were calculated at the resolution
431 of days. The estimated spline yields both $\hat{\boldsymbol{\beta}}$ and $\hat{\Sigma}_{\hat{\boldsymbol{\beta}}}$, an estimate of the covariance matrix of the
432 sampling distribution of $\hat{\boldsymbol{\beta}}$. The scalar standard errors of $\hat{\boldsymbol{\beta}}$ are the square roots of the diagonal
433 elements of $\hat{\Sigma}_{\hat{\boldsymbol{\beta}}}$, whereas the off-diagonal elements are the correlations among the elements of $\hat{\boldsymbol{\beta}}$.
434 Since parameter estimates are correlated, the covariance of $\hat{\mathbf{y}}$ is $\hat{\Sigma}_{\hat{\mathbf{y}}} = \mathbf{X}\hat{\Sigma}_{\hat{\boldsymbol{\beta}}}\mathbf{X}^T$.

435 The estimated seasonal average was calculated from the vector of daily values for each of the n
436 days in the season of interest with $\hat{\mu} = A^T\hat{\mathbf{y}}$, where A^T is a row vector with all values equal to
437 $1/n$. The variance of $\hat{\mu}$ is $\hat{\sigma}_{\hat{\mu}}^2 = A^T\hat{\Sigma}_{\hat{\mathbf{y}}}A$ and its standard error is $\hat{\sigma}_{\hat{\mu}}$. Each of these estimates are
438 from the approximate multivariate normality of the sampling distribution of $\hat{\boldsymbol{\beta}}$.

439 **Third-stage analysis: Trend analysis of seasonal features with uncertainties**

440 In stage three of the analysis, we used a meta-analysis method to evaluate linear trends across
441 years of seasonal-average water quality, characterized by the within-year means ($\hat{\mu}_t$) and their
442 standard errors ($\hat{\sigma}_{\hat{\mu},t}$) that we estimated in stage two of the analysis. This analysis provided a
443 direct answer to the basis for directly answering the question: *Is there a significant linear trend*

Formatted: Font: Italic

444 across a group of years in a seasonal average², taking into account uncertainty? For example,
445 is there a trend in the spring chl-a average from 1990 to 2000? This question can also be
446 posed in a moving-window manner across a time series (e.g., spring average trend from 1990-
447 2000, 1991-2001, etc.). For all analyses, the response data of interest are $\hat{\mu}_t$, $t = 1, \dots, N$, with
448 their associated standard errors, $\hat{\sigma}_{\hat{\mu},t}$. N is the number of years of the study.

449 A mixed-effects meta-analysis model can estimate linear trends when each observation has an
450 associated standard error, which is the case with our estimates $\hat{\mu}_t$ and $\hat{\sigma}_{\hat{\mu},t}$. Differences in
451 standard errors, which may result from different monitoring effort between years, are explicitly
452 considered in the analysis. The model can be expressed using notation similar to Sera et al.
453 (2019):

454

$$\hat{\mu}_t = \beta_0 + \beta_t t + b_t + \epsilon_t \quad (45)$$

455 where β_0 is the intercept, t is the year, β_t is the slope, b_t is the random effect for year t , and ϵ_t is
456 the residual for year t . Accordingly, the seasonal average for year t is $\mu_t = \beta_0 + \beta_t t + b_t$. The
457 “residual,” ϵ_t , represents estimation error in $\hat{\mu}_t$, namely $\hat{\mu}_t - \mu_t$. The residuals are assumed to be
458 independent and normally distributed with mean 0 and variances $\hat{\sigma}_{\hat{\mu},t}^2$, where the latter is
459 estimated from the calculations above. The random effect, b_t , is the difference between μ_t and
460 $\beta_0 + \beta_t t$ and is considered the “residual” in the sense of unexplained variation not due to the
461 estimation error. The random effect follows a normal distribution with mean 0 and variance, σ_b^2 ,
462 to be estimated.

463 We estimated the model (equation (45)) using the *mixmeta* package in R (Sera et al., 2019).
464 Results from *mixmeta* have a similar interpretation as those from regression analysis, but
465 parameter estimates and their standard errors incorporate the known standard errors of the

466 response values. The default estimation method for *mixmeta*, restricted maximum likelihood
467 (REML), was used. The meta-analysis models were applied to a chosen sequence or “window”
468 of years for estimating the linear trend.

469 **Trend comparisons**

470 The above methods were applied to each station by evaluating changes in seasonal averages ~~from~~
471 ~~January to June and July to December~~ for approximately ten year moving windows from 1991 to
472 2019. ~~The choice~~The overall technique and *wgtrends* R package allow for easy and flexible
473 selection of ~~within-year~~seasonal windows. As an example focused on illustrating the approach
474 (as opposed to extracting mechanistic interpretations), we selected two broad seasons, January-
475 June and July-December, that are generally relevant to phytoplankton bloom phenology in the
476 SFE (Cloern et al., 2020). ~~The~~and combined include conditions over the full year. A moving-
477 window approach ~~applied~~encoded within *wgtrends* was used to apply the meta-analysis to each
478 decadal window (e.g., 1991-2001, 1992-2002, etc.), allowing changes in slope and its
479 significance to be interpreted as the window is shifted one year at a time. We interpret the slope
480 as representative for the central year for each block, but a predictive trend for the final year of
481 the window could also be interpreted. ~~For~~As a means for summarizing some results, we focus on
482 the windows 1991-2000, 2000-2010, and 2010-2019.

483 Finally, trend results from the meta-analysis regression method for each season and different
484 time periods were compared to “naive” across-year regressions that do not propagate uncertainty
485 to demonstrate how different and potentially misleading conclusions can be obtained. Trend
486 estimates were compared to 1) trends from ordinary least squares (OLS) regression applied to
487 seasonal averages from the raw data and 2) trends from OLS regression applied to GAM

488 seasonal averages. Select examples ~~were used~~ where differences were pronounced ~~to illustrate~~
489 ~~false positive or negatives that may occur with alternative methods.~~ ~~were used for illustration.~~
490 This analysis was then applied to all stations. The method that formally propagates uncertainty
491 should have more robust statistical properties, such as accurate confidence interval coverage,
492 than naïve methods. For this reason, even when results are similar across methods, the more
493 robust method provides the best support for those results.

494 **Results**

495 **Model performance and predictions**

496 Model predictions for chl-a trends across all stations had an average R-squared value of 71%
497 (Table 4) and ~~ranging range~~ from 59% (station 22) to 78% (station 18). GAM predictions from
498 north to south showed more pronounced annual and seasonal changes in chl-a towards the more
499 southern stations (Figures S1-S9). All the models suggested 1) increasing chl-a from 1990 until
500 2005 to 2010, followed by decreasing chl-a until the end of the record in 2019, 2) a spring chl-a
501 peak, particularly at southern stations, and 3) a fall chl-a peak that was smaller than the spring
502 peak. The magnitude of the fall peak did not vary noticeably by location (Figures S1-S9).

503 **Inter-annual trend estimates**

504 Estimates ~~of linear trends in seasonal averages~~ ~~from the seasonal trend analyses (mixed-effects~~
505 ~~meta-analysis regressions)~~ across roughly ten-year windows for different seasons are shown for
506 station 34 (Figure 4). Plots a-c show trends in January to June averages while plots d-f show
507 trends in July to December averages. ~~The seasonal trend analyses showed that~~ January to June

508 chl-a increased (\log_{10} chl-a slope $0.03 \mu\text{g L}^{-1} \text{ yr}^{-1}$, 0.01 - 0.0605 95% confidence interval) from
509 1991 to 2000, whereas July to December chl-a did not change significantly. Chl-a also increased
510 from 2000 to 2010, but only for July to December (\log_{10} slope 0.03 , 0.01 - 0.05 95% confidence
511 interval). Finally, chl-a decreased from 2010 to 2019 but again only for July to December (\log_{10}
512 chl-a slope -0.02 , -0.04 - 0 95% confidence interval). Because the trends were confined to certain
513 times of the year, the seasonal estimates provide additional information beyond coarser estimates
514 that cover the entire year.

515 Temporal changes varied among regions of the Bay and were more pronounced at southern
516 stations. Figure 5 shows results from similar analyses as those in Figure 4, but applied to all
517 stations. ~~Mixed effects meta analysis regressions applied to~~The seasonal ~~average~~trend analyses
518 showed that increases (based on $p < 0.05$) for the January to June period were observed at
519 stations 32,~~34~~, and 36~~34~~ from 1991 to 2000 ~~and station 18 from 2000 to 2010~~; decreases were
520 observed at stations 30 and 32 from 2010 to 2019. For the July to December period, increases
521 were observed at stations 24, 27, 30, and 32 from 1991 to 2000 and stations 18, 21, 22, 24, and
522 34 from 2000 to 2010, whereas decreases were observed at stations 30, 32, and 34 from 2010 to
523 2019.

524 Results from a ten-year moving window comparison of seasonal trends provided additional
525 context on when significant changes were occurring at each station (Figure 6). Trends were
526 observed at all stations that followed a general pattern of increases early in the record followed
527 by decreases later in the record. Increases and decreases were observed in both the January to
528 June and July to December seasonal periods, with some notable exceptions. In particular, the
529 most southern stations (32, 34, 36) had increasing trends prior to 2005 ~~only that were more often~~
530 observed in the July to December period. Additionally, chl-a at the more northern stations has

531 not changed in recent years for ~~both~~either seasonal ~~periods~~period. For most stations and seasonal
532 periods, a change from increasing to decreasing chl-a occurred around 2007.

533 Importance of uncertainty propagation

534 Results showing trend estimates from meta-analysis on GAM seasonal ~~estimates~~averages
535 provided different conclusions than those from either OLS regression through seasonal averages
536 from raw data (Figure 7 row 1) or OLS regression through GAM ~~estimates~~seasonal averages
537 without uncertainty propagation (Figure 7 row 2). Figure 7a shows trend estimates for station
538 ~~3630~~ for January to ~~July~~June averages from ~~1991~~2000 to ~~2000~~2010. Only the meta-analysis
539 regression results show a trend in this example (based on $p < 0.05$). The OLS regression on
540 observed estimates (top plot) and OLS regression on GAM estimates (middle plot) did not
541 identify trends. Figure 7b shows trend estimates ~~for the same at~~ station 34 for July to December
542 averages from ~~1991 to~~ 2000 ~~to~~ 2010. Unlike the first example, only the ~~top middle~~ figure shows a
543 trend, whereas the ~~top and~~ bottom ~~two~~ plots do not show trends. In both cases, only the meta-
544 analysis results give reliable conclusions because of full propagation of uncertainty across
545 methods. Even in cases where the p-value threshold is not of interest, the confidence intervals
546 from the alternative methods will be inaccurate.

547 Applying the same comparison to all stations showed that different trend analysis methods
548 provided conflicting information on the magnitude and significance of the seasonal chl-a changes
549 in each decade (Figure 8). In many cases, the slope estimates were similar in magnitude, with
550 some exceptions at the more southern stations where the OLS estimates suggested a larger trend
551 than the meta-analysis methods. More importantly, differences in the magnitude of the
552 confidence intervals between the OLS models applied to the GAM averages and the meta-

553 analyses were also observed, reflecting the ability of the latter to more accurately assess
554 significance of trends by accounting for uncertainty in the average estimates.

555 **Discussion**

556 Propagation of uncertainty from within-year features of estimated GAMs to across-year trends
557 using mixed-effects meta-analysis is a new approach that can address different questions than
558 previous methods. Our approach has several advantages over more conventional approaches for
559 analysis of water quality data from long-term monitoring programs. GAMs are capable of
560 modelling time series with missing observations or irregular sampling which can complicate
561 trend assessment and comparison of trends between locations ([Junninen et al., 2004](#); [Racault et](#)
562 [al., 2014](#)). As noted above, non-parametric approaches (i.e., seasonal Kendall tests) are by far the
563 most common trend analysis methods applied to long-term water quality data ([Helsel et al.,](#)
564 [2020](#); [Hirsch et al., 1982](#)). These methods only assess the direction and significance of
565 comparisons ~~between year pairs across years~~, and importantly, do not account for full propagation
566 of uncertainty inherent in raw observations if the raw data are aggregated to meet test
567 requirements. Aggregation of raw data, e.g., averaging of observations within a year or season to
568 comply with the requirements of Kendall tests, risks loss of information by removing variation
569 between observations at smaller time scales. The logical outcome is increased risk of incorrect
570 conclusions from test results.

571 Incorrect conclusions on trends can have dramatic consequences for regulated parties under
572 existing water quality compliance frameworks ([Smith et al., 2001](#)). Our examples in Figures 7
573 and 8 demonstrate these risks if propagation of uncertainty from raw observations across
574 methods is unaccounted for in trend assessment. ~~Our assessment of trends-The “naïve” method~~

575 using OLS regression applied to seasonal averages from the raw observations ~~is effectively~~
576 ~~like fails to propagate uncertainty, similarly to~~ averaging results within a year and applying a
577 simple Kendall test. In ~~many some~~ cases the results may be similar ~~to those from fully~~
578 ~~propagating uncertainty~~, but ~~the~~ loss of information ~~with averaging~~ can lead to increased Type I
579 or II error rates depending on characteristics of the raw data and the method used for their
580 evaluation (Shabman and Smith, 2003). Our examples demonstrated the increased potential for
581 incorrect conclusions at specific monitoring locations and, at ~~much~~ larger spatial scales, across
582 all stations if simpler trend analyses are used. Even though simpler methods may produce similar
583 results in some cases, particularly with frequent sampling and similar effort between years, the
584 only way to confirm such an outcome would be to compare results, relying on the method with
585 full propagation of uncertainty to be the more robust method. Use of methods that fully account
586 for uncertainty is recommended to obtain statistically valid results in a wider range of conditions.

587 Results here also show that GAM structure (i.e., choice of smoothing terms) was less important
588 than allowing the model sufficient freedom to fit the data. This is an important conclusion that
589 provides guidance on how GAMs could be used to model time series from long-term
590 environmental monitoring programs. Models with separate smoothers for continuous year and
591 day of year can produce nearly identical results in the predicted trends if the knots are
592 sufficiently high to allow the GAMs to be fit as intended by the methods in the mgcv package
593 (Figure 3). The approach presented here leverages the ability of GAMs to objectively estimate
594 smoothed trends across years by identifying an optimal level of smoothing using generalized
595 cross-validation to extract an underlying signal in the observed data (Wood, 2017, 2004).

596 The underlying cross-validation methods used by GAMs in the mgcv package also reduce the
597 decisions that may be necessary for the implementation of alternative trend assessment methods.

598 For example, WRTDS and similar smoothing approaches (e.g., LOESS) require decisions on
599 appropriate window widths or bandwidths to define the neighborhood of observations for
600 smoothing (Hirsch et al., 2010; Wan et al., 2017). This is especially problematic for policy
601 analysis or regulatory decisions if the results change based on arbitrary decisions of the analyst.
602 Because these decisions are not needed for GAMs, the results can be considered a more objective
603 and potentially “true” accurate signal of actual trends that are minimally influenced by process or
604 observation error present in the raw data.

605 Several limitations of the proposed methods deserve mention. First, if sampling is so irregular
606 that important fluctuations are missed entirely in some years, the GAM estimates and uncertainty
607 propagation could become dubious in interpretation and usefulness. Second, estimation of GAMs
608 for very long series can be computationally demanding. When this is an obstacle, one could do
609 the first two analysis stages using temporal windows of the full data, with the only implication
610 being that different degrees of smoothness may be estimated for different windows, which
611 indeed might be justified by the data. Third, meta-analysis regression results for a very small
612 number of years, particularly confidence intervals and associated p-values, may be inaccurate
613 (e.g. in confidence interval coverage). In such cases, one could make alternative use of the GAM
614 seasonal averages and standard errors, such as for pairwise comparisons among years.

615 Future work

616 Additional work could be conducted to further strengthen the conclusions based on trends from
617 meta-analysis regression applied to the GAM seasonal averages. Our third stage analyses require
618 *a priori* decisions on long-term time scales of interest and future work could generalize these
619 choices. Although there are undoubtedly many scenarios where years of interest can be chosen

620 objectively by the needs of an analysis (e.g., regulatory compliance periods, time since
621 management intervention), a more general question of “when” changes occur independent of
622 user decisions is also important to address. Additional methods could be developed using
623 objective criteria to identify inflection points or other important periods where changes occur
624 independent of a user choice. Assessing water quality changes beyond an evaluation of seasonal
625 averages could also be possible with our approach, such as assessing changes in the timing or
626 magnitude of a seasonal peak across years.

627 Additional explanatory variables could be identified ~~that may be associated with~~to explain trends
628 in either the trend after GAM stage or the trend has been adequately described meta-analysis
629 stage of analysis. This information ~~has~~would have obvious implications for management
630 decisions on factors that influence water quality changes, e.g., wastewater treatment upgrades,
631 large-scale climatic factors, or flow regulation practices. ~~An advantage of GAMs is their~~
632 ~~flexibility in including alternative predictors, such that the significance of a predictor or~~
633 ~~comparison of nested models with and without different predictors can provide evidence of~~
634 ~~which predictors are driving the observed trends~~Including alternative predictors in the GAMs
635 (Wood and Augustin, 2002; Zuur et al., 2009). ~~In such cases, considerations of model structure~~
636 ~~can have direct implications on conclusions given how GAMs could be used to assess different~~
637 ~~questions. Our goal was to describe chl a changes relative to time, where the predictors were~~
638 ~~variations on a general theme (e.g., season vs. year). This is a different application from using~~
639 ~~GAMs with predictors selected to explain those changes over time. Therefore, using our~~
640 ~~approach to evaluate explanatory variables will require testing of different model structures)~~
641 could reduce the uncertainties of its estimates and, if relevant, allow the influence of those
642 variables to be removed. Including alternative predictors in the meta-analysis could help in

643 explaining long-term trends of seasonal averages or other metrics obtained from the GAMs. Our
644 goal here was to describe chl-a changes relative to time, so the single predictor in both modeling
645 stages was time.

646 Finally, the evaluation of trends for alternative water quality variables in addition to chl-a is a
647 simple and logical extension of the methods proposed in this study. The long-term monitoring
648 program maintained by USGS includes multiple parameters in addition to chl-a that can provide
649 additional context into broader water quality trends in the SFE (Cloern and Schraga, 2016;
650 Schraga et al., 2020). These parameters include salinity, temperature, light attenuation, dissolved
651 oxygen, suspended particulate matter, and dissolved inorganic nutrients, which collectively can
652 be used to provide a broader understanding of potential eutrophication patterns or ecosystem
653 shifts at seasonal and multi-decadal scales. Chl-a measurements can also be used to estimate
654 gross primary production to assess process rates that may be more indicative of system function
655 (Cloern et al., 2007; Jassby et al., 2002). The open-source *wqtrends* R package (Beck et al.,
656 2021) developed for this manuscript can be used for these analyses to provide additional insight
657 into potential drivers of water quality change in the SFE and other estuarine systems.

Formatted: Font: Italic

658 Acknowledgments

659 This work was supported by funding from the San Francisco Bay Nutrient Management Strategy
660 (NMS). We thank the staff of the US Geological Survey that collect and maintain long-term
661 monitoring data in San Francisco Bay. This work benefited from discussions with the San
662 Francisco BayNMS Nutrient Technical Workgroup and Steering Committee. We thank James D.
663 Hagy III and two anonymous reviewers for reviewing an earlier draft of providing helpful
664 comments on this manuscript.

→ **Formatted:** First Paragraph, Line spacing: Double

666 **References**

- 667 Alpine, A.E., Cloern, J.E., 1988. Phytoplankton growth rates in a light-limited environment, San
668 Francisco Bay. *Marine Ecology Progress Series* 44, 167–173.
- 669 Beck, M.W., de Valpine, P., Murpy, R., Wren, I., Chelsky, A., Foley, M., Senn, D., 2021. tbep-
670 tech/wqtrends: v1.1.0 (Version v1.1.0). Zenodo. <http://doi.org/10.5281/zenodo.4509638>.
- 671 Beck, M.W., Hagy, J.D., III, 2015. Adaptation of a weighted regression approach to evaluate
672 water quality trends in an estuary. *Environmental Modelling and Assessment* 20, 637–655.
673 <https://doi.org/10.1007/s10666-015-9452-8>
- 674 Beck, M.W., Jabusch, T.W., Trowbridge, P.R., Senn, D.B., 2018. Four decades of water quality
675 change in the upper San Francisco Estuary. *Estuarine, Coastal and Shelf Science* 212, 11–22.
676 <https://doi.org/10.1016/j.ecss.2018.06.021>
- 677 Beck, M.W., Murphy, R.R., 2017. Numerical and qualitative contrasts of two statistical models
678 for water quality change in tidal waters. *Journal of the American Water Resources Association*
679 53, 197–219. <https://doi.org/10.1111/1752-1688.12489>
- 680 Cleveland, R.B., Cleveland, W.S., McRae, J.E., Terpenning, I., 1990. STL: A seasonal-trend
681 decomposition procedure based on Loess. *Journal of Official Statistics* 6, 3–73.
- 682 Cloern, J.E. 2018. Patterns, pace, and processes of water-quality variability in a long-studied
683 estuary. *Limnology and Oceanography* 64, S192-S208. <https://doi.org/10.1002/lno.10958>
- 684 Cloern, J.E., Jassby, A.D., 2012. Drivers of change in estuarine-coastal ecosystems: Discoveries
685 from four decades of study in San Francisco Bay. *Reviews of Geophysics* 50.
686 <https://doi.org/10.1029/2012RG000397>
- 687 Cloern, J.E., Jassby, A.D., 2010. Patterns and scales of phytoplankton variability in estuarine-
688 coastal ecosystems. *Estuaries and Coasts* 33, 230–241.
- 689 Cloern, J.E., Jassby, A.D., Thompson, J.K., Hieb, K.A., 2007. A cold phase of the East Pacific
690 triggers new phytoplankton blooms in San Francisco Bay. *Proceedings of the National Academy
691 of Sciences of the United States of America* 104, 18561–18565.
- 692 Cloern, J.E., Schraga, T.S., 2016. USGS measurements of water quality in San Francisco Bay
693 (CA), 1969-2015: U.S. Geological Survey data release.
694 <https://doi.org/10.5066/F7TQ5ZPR>. <https://doi.org/10.5066/F7TQ5ZPR>.
- 695 Cloern, J.E., Shraga, T.S., Nejad, E., Martin, C., 2020. Nutrient status of San Francisco Bay and
696 its management implications. *Estuaries & Coasts* 43, 1299–1317.
697 <https://doi.org/10.1007/s12237-020-00737-w>
- 698 Cole, B.E., Cloern, J.E., 1984. Significance of biomass and light availability to phytoplankton
699 productivity in San Francisco Bay. *Marine Ecology Progress Series* 17, 15–24.
- 700 Cumming, G.S., Cumming, D.H.M., Redman, C.L., 2006. Scale mismatches in social-ecological
701 systems: Causes, consequences, and solutions. *Ecology and Society* 11, 14.

- 702 Forbes, D.J., Xie, Z., 2018. Identifying process scales in the Indian River Lagoon, Florida using
703 wavelet transform analysis of dissolved oxygen. Ecological Complexity 36, 149–167.
704 <https://doi.org/10.1016/j.ecocom.2018.07.005>
- 705 Gasparrini, A., Armstrong, B., Kenward, M.G., 2012. Multivariate meta-analysis for non-linear
706 and other multi-parameter associations. Statistics in Medicine 31, 3821–3839.
707 <https://doi.org/10.1002/sim.5471>
- 708 Hafner, R.P., 2010. Local regression models: Advancements, applications, and new methods
709 (PhD thesis). Purdue University, West Lafayette, Indiana.
- 710 Haraguchi, L., Carstensen, J., Abreu, P.C., Odebrecht, C., 2015. Long-term changes of the
711 phytoplankton community and biomass in the subtropical shallow Patos Lagoon Estuary, Brazil.
712 Estuarine, Coastal and Shelf Science 162, 76–87.
- 713 He, S., Mazumdar, S., Arena, V.C., 2006. A comparative study of the use of GAM and GLM in
714 air pollution research. Environmetrics 17, 81–93. <https://doi.org/10.1002/env.751>
- 715 Helsel, D.R., Hirsch, R.M., Ryberg, K.R., Archfield, S.A., Gilroy, E.J., 2020. Statistical methods
716 in water resources, 2nd ed. U.S. Geological Survey Techniques; Methods, book 4, chapter A3,
717 version 1.1, Reston, Virginia.
- 718 Hirsch, R.M., Archfield, S.A., De Cicco, L.A., 2015. A bootstrap method for estimating
719 uncertainty of water quality trends. Environmental Modelling and Software 73, 148–166.
720 <https://doi.org/10.1016/j.envsoft.2015.07.017>
- 721 Hirsch, R.M., Moyer, D.L., Archfield, S.A., 2010. Weighted regressions on time, discharge, and
722 season (WRTDS), with an application to Chesapeake Bay river inputs. Journal of the American
723 Water Resources Association 46, 857–880.
- 724 Hirsch, R.M., Slack, J.R., Smith, R.A., 1982. Techniques of trend analysis for monthly water
725 quality data. Water Resources Research 18, 107–121.
- 726 Jassby, A.D., 2008. Phytoplankton in the Upper San Francisco Estuary: Recent biomass trends,
727 their causes, and their trophic significance. San Francisco Estuary and Watershed Science 6, 1–
728 24.
- 729 Jassby, A.D., Cloern, J.E., Cole, B.E., 2002. Annual primary production: Patterns and
730 mechanisms of change in a nutrient-rich tidal ecosystem. Limnology and Oceanography 47, 698–
731 712.
- 732 Junninen, H., Niska, H., Tuppurainen, K., Ruuskanen, J., Kolehmainen, M., 2004. Methods for
733 imputation of missing values in air quality data sets. Atmospheric Environment 38, 2895–2907.
734 <https://doi.org/10.1016/j.atmosenv.2004.02.026>
- 735 Kimmerer, W.J., Thompson, J.K., 2014. Phytoplankton growth balanced by clam and
736 zooplankton grazing and net transport into the low-salinity zone of the San Francisco Estuary.
737 Estuaries and Coasts 37, 1202–1218.

- 738 Lefcheck, J.S., Wilcox, D.J., Murphy, R.R., Marion, S.R., Orth, R.J., 2017. Multiple stressors
739 threaten the imperiled coastal foundation species eelgrass (*zostera marina*) in Chesapeake Bay,
740 USA. Global Change Biology 23, 3474–3483. <https://doi.org/10.1111/gcb.13623>
- 741 Lehman, P.W., Kurobe, T., Lesmeister, S., Baxa, D., Tung, A., Teh, S.J., 2017. Impacts of the
742 2014 severe drought on the *Microcystis* bloom in San Francisco Estuary. Harmful Algae 63, 94–
743 108. <https://doi.org/10.1016/j.hal.2017.01.011>
- 744 Lortie, C.J., 2014. Formalized synthesis opportunities for ecology: Systematic reviews and meta-
745 analyses. OIKOS 123, 897–902. <https://doi.org/10.1111/j.1600-0706.2013.00970.x>
- 746 Morton, R., Henderson, B.L., 2008. Estimation of nonlinear trends in water quality: An
747 improved approach using generalized additive models. Water Resources Research 44, W07420.
748 <https://doi.org/10.1029/2007WR006191>
- 749 Murphy, R.R., Perry, E., Harcum, J., Keisman, J., 2019. A Generalized Additive Model
750 Approach to evaluating water quality: Chesapeake Bay case study. Environmental Modelling &
751 Software 118, 1–13. <https://doi.org/10.1016/j.envsoft.2019.03.027>
- 752 Navarro, G., Caballero, I., Prieto, L., Vázquez, A., Flecha, S., Huertas, I.E., Ruiz, J., 2012.
753 Seasonal-to-interannual variability of chlorophyll-*a* bloom timing associated with physical
754 forcing in the Gulf of Cádiz. Advances in Space Research 50, 1164–1172.
755 <https://doi.org/10.1016/j.asr.2011.11.034>
- 756 Novick, E., Senn, D., 2014. External nutrient loads to San Francisco Bay (No. Contribution
757 Number 704). Peacock, M.B., Gibble, C.M., Senn, D.B., Cloern, J.E., Kudela, R.M. 2018.
758 Blurred lines: Multiple freshwater and marine algal toxins at the land-sea interface of San
759 Francisco Bay, California. Harmful Algae 73, 138–147.
760 <https://doi.org/10.1016/j.hal.2018.02.005>
- 761 San Francisco Estuary Institute, Richmond, CA.
- 762 Pearce, J.L., Beringer, J., Nicholls, N., Hyndman, R.J., Tapper, N.J., 2011. Quantifying the
763 influence of local meteorology on air quality using generalized additive models. Atmospheric
764 Environment 45, 1328–1336. <https://doi.org/10.1016/j.atmosenv.2010.11.051>
- 765 R Core Team, 2020. R: A language and environment for statistical computing. R Foundation for
766 Statistical Computing, R v4.0.2, Vienna, Austria.
- 767 Racault, M.F., Sathyendranath, S., Platt, T., 2014. Impact of missing data on the estimation of
768 ecological indicators from satellite ocean-colour time series. Remote Sensing of Environment
769 152, 15–28. <https://doi.org/10.1016/j.rse.2014.05.016>
- 770 Schraga, T.S., Nejad, E.S., Martin, C.A., Cloern, J.E., 2020. USGS measurements of water
771 quality in San Francisco (CA), beginning in 2016 (ver. 3.0, March 2020): U.S. Geological
772 Survey data release. <https://doi.org/10.5066/F7D21WGF>.
- 773 Sera, F., Armstrong, B., Blangiardo, M., Gasparri, A., 2019. An extended mixed-effects
774 framework for meta-analysis. Statistics in Medicine 38, 5429–5444.
775 <https://doi.org/10.1002/sim.8362>

- 776 [SFBRWQCB \(San Francisco Bay Regional Water Quality Control Board\), 2017. Water quality control plan \(basin plan\) for the San Francisco Bay basin. Prepared by California Regional Water Quality Control Board, San Francisco Bay Region, Oakland, CA.](https://www.waterboards.ca.gov/sanfranciscobay/basin_planning.html)
777
778
779 https://www.waterboards.ca.gov/sanfranciscobay/basin_planning.html
- 780 [SFEI \(San Francisco Estuary Institute\), 2014a. External nutrient loads to San Francisco Bay, SFEI Contribution Number 704. San Francisco Estuary Institute, Richmond, CA.](#)
781
- 782 [SFEI \(San Francisco Estuary Institute\), 2014b. Scientific foundation for the San Francisco Bay Nutrient Management Strategy. SFEI Contribution Number 979. San Francisco Estuary Institute, Richmond, CA.](#)
783
784
- 785 [SFEI \(San Francisco Estuary Institute\), 2021. Connections to tidal marsh and restored salt ponds drive seasonal and spatial variability in ecosystem metabolic rates in Lower South San Francisco Bay. SFEI Contribution No. 992. San Francisco Estuary Institute, Richmond, CA.](#)
786
787
- 788 Shabman, L., Smith, E., 2003. Implications of applying statistically based procedures for water quality assessment. *Journal of Water Resources Planning and Management* 129, 330–336.
789
790 [https://doi.org/10.1061/\(ASCE\)0733-9496\(2003\)129:4\(330\)](https://doi.org/10.1061/(ASCE)0733-9496(2003)129:4(330))
- 791 Smith, E.P., Ye, K., Hughes, C., Shabman, L., 2001. Statistical assessment of violations of water quality standards under section 303 (d) of the Clean Water Act. *Environmental science & technology* 35, 606–612. <https://doi.org/10.1021/es001159e>
792
793
- 794 Stow, C.A., Cha, Y., Johnson, L.T., Confesor, R., Richards, R.P., 2015. Long-term and seasonal trend decomposition of Maumee River nutrient inputs to western Lake Erie. *Environmental Science and Technology* 49, 3392–3400. <https://doi.org/10.1021/es5062648>
795
796
- 797 [Sutula, M.A., Kudela, R.M., Hagy, J.D., III, Harding, L.W., Jr., Senn, D.B., Cloern, J.E., Bricker, S., Berg, G.M., Beck, M.W. 2017. Novel analyses of long-term data provide a scientific basis for chlorophyll-a thresholds in San Francisco Bay. *Estuarine, Coastal and Shelf Science*, 197, 107-118. https://doi.org/10.1016/j.ecss.2017.07.009](#)
798
799
800
- 801 Urquhart, N.S., Paulsen, S.G., Larsen, D.P., 1998. Monitoring for policy-relevant regional trends over time. *Ecological Applications* 8, 246–257. [https://doi.org/10.1890/1051-0761\(1998\)008\[0246:MFPRRO\]2.0.CO;2](https://doi.org/10.1890/1051-0761(1998)008[0246:MFPRRO]2.0.CO;2)
802
803
- 804 Wan, Y., Wan, L., Li, Y., Doering, P., 2017. Decadal and seasonal trends of nutrient concentration and export from highly managed coastal catchments. *Water Research* 115, 180–194.
805
806
- 807 Wood, S.N., 2017. Generalized additive models: An introduction with r, 2nd ed. Chapman; Hall, CRC Press, London, United Kingdom.
808
- 809 Wood, S.N., 2004. Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association* 99, 673–686.
810
811 <https://doi.org/10.1198/016214504000000980>
- 812 Wood, S.N., 2003. Thin-plate regression splines. *Journal of the Royal Statistical Society (B)* 65, 95–114. <https://doi.org/10.1111/1467-9868.00374>
813

- 814 Wood, S.N., Augustin, N.H., 2002. GAMs with integrated model selection using penalized
815 regression splines and applications to environmental modelling. Ecological Modelling 157, 157–
816 177. [https://doi.org/10.1016/S0304-3800\(02\)00193-X](https://doi.org/10.1016/S0304-3800(02)00193-X)
- 817 Yang, G., Moyer, D.L., 2020. Estimation of nonlinear water-quality trends in high-frequency
818 monitoring data. Science of The Total Environment 715, 136686.
819 10.1016/j.scitotenv.2020.136686.
- 820 Zuur, A.F., Ieno, E.N., Walker, N.J., Saveliev, A.A., Smith, G.M., 2009. Mixed effects models
821 and extensions in ecology with r. Springer-Verlag, New York, New York.
- 822

Formatted: Normal, Space After: 0 pt

823 Figures

824
825 *Figure 1: Observed chl-a concentrations for all stations in central and south San Francisco*
826 *Estuary (18-36, Figure 2), with (a) annual summer/fall concentrations (Aug - Dec) and (b)*
827 *monthly concentrations by decade.*

828 *Figure 2: Station locations in the central and south San Francisco Estuary used for analysis. See*
829 *Table 1 for station descriptions. Full dataset described in Schraga et al. (2020).*

830 *Figure 3: GAM output of estimated chl-a at station 32 for models S, SYD, and SYDI. Model SY is*
831 *identical to S and is not shown. Plots in (a) show model predictions by day of year with separate*
832 *lines for each year. Plots in (b) show pairwise comparisons of predicted chl-a between the*
833 *models and plots in (c) show the same comparisons as in (b) but only for results from the*
834 *estimated smoother for the cont_year variable. The plots demonstrate that results between the*
835 *models are similar except for a few observations at extreme values (a, b), but they vary in how*
836 *they allocate contributions to the predictions among different additive splines (c). The 1:1 lines*
837 *are in red to facilitate comparisons.*

838 *Figure 4: Examples of seasonal averages and trend estimates in ten year blocks from meta-*
839 *analyses using results of GAM predictions for station 34. Plots (a), (b), and (c) show trend*
840 *estimates for January to June averages and (d), (e), and (f) show trend estimates for July to*
841 *December averages. The trend lines estimate the rate of change of chl-a per year, reported as*
842 *the log₁₀-slope (+/- 95 % confidence interval) in the sub-plot titles. ns: not significant at $\alpha =$*
843 *0.05, * $p < 0.05$*

844 *Figure 5: Interannual trend estimates of seasonal averages by decade for chl-a at each station.*
845 *Point type and color represent the direction and magnitude of an estimated trend as the log₁₀*
846 *slope for chl-a concentration per year. Trends with $p < 0.05$ are marked with an asterisk. All*
847 *results are from Model S.*

848 *Figure 6: Estimates of log₁₀ chl-a change per year (+/- 95% confidence interval) from applying*
849 *the meta-analysis across the seasonal averages for each station. Stations are arranged top to*
850 *bottom from north to south. Plots in (a) show estimates for seasonal averages from January to*
851 *June and plots in (b) show estimates for seasonal averages from July to December. Results are*
852 *from a ten-year, centered moving window where each point shows a linear trend estimate from*
853 *five years prior to five years after each year. Estimates prior to 1996 and after 2014 are not*
854 *available because of an incomplete ten year record for estimating the trend. Significant estimates*
855 *are shown in red.*

856 *Figure 7: Trend estimate comparisons (arithmetic scale) for three models applied to seasonal*
857 *averages of chl-a in different annual periods at station 36. (a) OLS regression applied to seasonal averages of chl-a from*
858 *the raw data, the second row shows OLS regression applied to seasonal averages of chl-a from*
859 *the GAM (without error propagation), and the third row shows meta-analysis regression applied*
860 *to the seasonal averages of chl-a from the GAM. Regressions in each plot are fit through the*
861 *seasonal estimates indicated in the plot titles for a specified year range. These examples on*
862 *selected periods of time show that slope estimates can be similar, but the confidence intervals*
863 *vary.*

865 Figure 8: Trend estimate comparisons for three models applied to seasonal averages of chl-a in
866 different annual periods at each station. The “OLS raw” trend model is based on an ordinary
867 least squares (OLS) regression fit to the seasonal averages of chl-a from the raw data, the “OLS
868 GAM” trend model is based on an OLS regression fit to the seasonal averages of chl-a from the
869 GAM model (without error propagation), and the “Meta-analysis GAM” trend model is based
870 on a meta-analysis regression fit to the seasonal averages of chl-a from the GAM model. Values
871 for each model are the \log_{10} -slope estimates (+/- 95% confidence interval) as annual change per
872 year within each season, with line style denoting trend significance.

873

Formatted: Image Caption

Formatted: Font: 12 pt, Not Bold, Font color: Auto

874 **Tables**

875 *Table 1: Station locations, sample sizes (from 1991 to 2019), and summary values (median,
876 minimum, maximum) for chl-a ($\mu\text{g L}^{-1}$). Stations are arranged from north to south.*

Station	Latitude	Longitude	n	Med.	Min.	Max.
18	37.836	-122.418	414	3.6	0.2	16.6
21	37.784	-122.351	576	4.4	0.6	40.0
22	37.752	-122.351	569	4.0	0.7	53.1
24	37.686	-122.334	595	4.2	0.7	47.3
27	37.617	-122.285	596	4.5	0.5	50.9
30	37.551	-122.184	608	5.1	0.8	112.2
32	37.517	-122.133	591	5.9	0.7	282.1
34	37.485	-122.086	544	6.5	0.6	158.3
36	37.468	-122.067	476	6.2	1.1	328.4

877

 Formatted Table

878 Table 2: Summary and details for each of the GAM structures. In practice, a sufficiently large
 879 number of knots provided to the additive terms will produce identical or comparable estimates
 880 for a response variable. The models differ in the allocation of penalties for the smoothness of
 881 each spline $s()$.

GAM	Additive components	Details
S	$s(\text{cont_year})f_1(\text{cont_year})$	A single smoother over a continuous year variable
SY	$\text{cont_year} + s(\text{cont_year})\beta_1 \text{cont_year} + f_1(\text{cont_year})$	A linear continuous year variable and a single smoother over a continuous year variable
SYD	$\text{cont_year} + s(\text{cont_year}) + s(\text{doy})\beta_1 \text{cont_year} + f_1(\text{cont_year}) + f_2(\text{doy})$	A linear continuous year variable, a smoother over a continuous year variable, and a smoother over a day of year variable
SYDI	$\text{cont_year} + s(\text{cont_year}) + s(\text{doy}) + t_i(\text{cont_year}, \text{doy})\beta_1 \text{cont_year} + f_1(\text{cont_year}) + f_2(\text{doy}) + f_3(\text{cont_year}, \text{doy})$	A linear continuous year variable, a smoother over a continuous year variable, a smoother over a day of year variable, and an interaction smoother across continuous year and day of year variables

882

← Formatted Table

883 *Table 3: Comparison of the four model structures (S, SY, SYD, SYDI) described in the first stage*
 884 *analysis of GAM estimation. The four models provide either identical or comparable ability to*
 885 *describe chl-a trends at an example station (32) in the southern end of the San Francisco*
 886 *Estuary. The models differ in additive smoothers and the amount of effective degrees of freedom*
 887 *(edf) in the smoothers (measure of wiggliness in each component), but the overall model*
 888 *predictions are similar. AIC: Akaike Information Criterion, GCV: generalized cross-validation*
 889 *score, R2: r-squared values for predictions, edf: effective degrees of freedom, F: F-statistic, p-*
 890 *val: probability value, ** p < 0.001*

model	AIC	GCV	R2	smoother	edf	F	p-val
S	-138.2	0.06	0.74	s(cont_year)	242.23	6.27	**
SY	-138.2	0.06	0.74	s(cont_year)	242.23	6.27	**
SYD	-135.66	0.06	0.74	s(cont_year) s(doy)	229.33 8.07	3.88 0.11	**
SYDI	-123.57	0.06	0.73	s(cont_year) s(doy)	136.88 9.54	3.13 0.79	**
				ti(cont_year,doy)	69.31	0.74	**

← Formatted Table

891

→ Formatted: Font: Not Italic

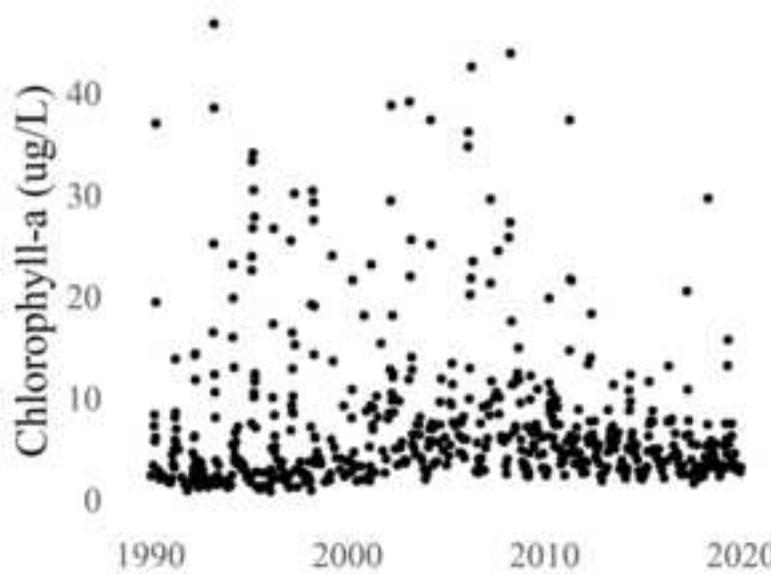
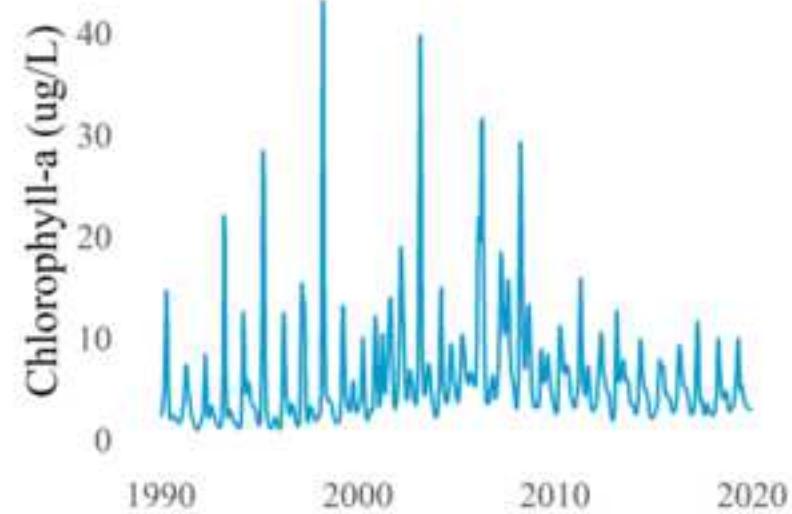
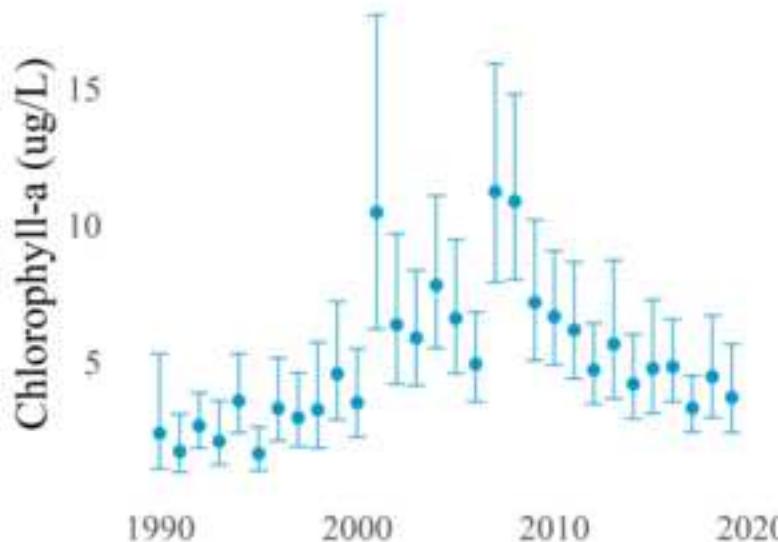
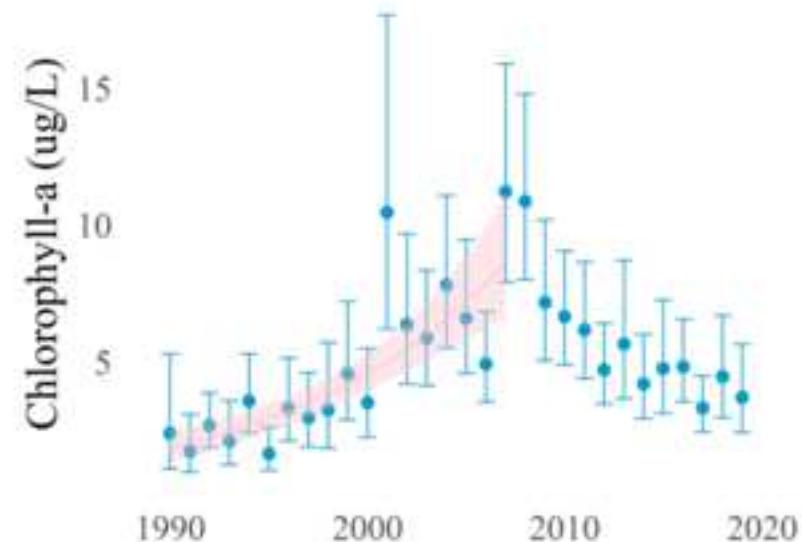
892 Table 4: Model performance statistics for each station as Akaike Information Criterion scores
893 (AIC), generalized cross-validation scores (GCV), and r-squared values.

station	AIC	GCV	R-squared
18	-430.94	0.04	0.78
21	-305.93	0.04	0.70
22	-160.07	0.05	0.59
24	-250.77	0.04	0.69
27	-188.72	0.05	0.72
30	-172.79	0.05	0.74
32	-138.20	0.06	0.74
34	-3.17	0.07	0.68
36	-22.05	0.07	0.73

894  Formatted Table

894  Formatted: Font: 12 pt, Not Bold, Font color: Auto

894  Formatted: Body Text

1) Water quality data**2) GAM estimated signal****3) GAM seasonal averages****4) Meta-analysis trend**

Highlights (for review : 3 to 5 bullet points (maximum 85 characters including spaces per bullet point)

- Trend analyses of water quality data must consider full propagation of uncertainty
- GAMs with different smoothing splines can extract nearly identical trends
- GAMs can estimate seasonal averages with uncertainty from monitoring data
- GAMs coupled with mixed-effects meta-analysis can accurately assess trends

1 Abstract

2 Effective stewardship of ecosystems to sustain current ecological status or mitigate impacts
3 requires nuanced understanding of how conditions have changed over time in response to
4 anthropogenic pressures and natural variability. Detecting and appropriately characterizing
5 changes requires accurate and flexible trend assessment methods that can be readily applied to
6 environmental monitoring datasets. A key requirement is complete propagation of uncertainty
7 through the analysis. However, this is difficult when there are mismatches between sampling
8 frequency, period of record, and trends of interest. Here, we propose a novel application of
9 generalized additive models (GAMs) for characterizing multi-decadal changes in water quality
10 indicators and demonstrate its utility by analyzing a 30-year record of biweekly-to-monthly
11 chlorophyll-a concentrations in the San Francisco Estuary. GAMs have shown promise in water
12 quality trend analysis to separate long-term (i.e., annual or decadal) trends from seasonal
13 variation. Our proposed methods estimate seasonal averages in a response variable with GAMs,
14 extract uncertainty measures for the seasonal estimates, and then use the uncertainty measures
15 with mixed-effects meta-analysis regression to quantify inter-annual trends that account for full
16 propagation of error across methods. We first demonstrate that nearly identical descriptions of
17 temporal changes can be obtained using different smoothing spline formulations of the original
18 time series. We then extract seasonal averages and their standard errors for an *a priori* time
19 period within each year from the GAM results. Finally, we demonstrate how across-year trends
20 in seasonal averages can be modeled with mixed-effects meta-analysis regression that propagates
21 uncertainties from the GAM fits to the across-year analysis. Overall, this approach leverages
22 GAMs to smooth data with missing observations or varying sample effort across years to

23 estimate seasonal averages and meta-analysis to estimate trends across years. Methods are
24 provided in the *wqtrends* R package.

25 *Key words:* chlorophyll, Generalized Additive Models, meta-analysis, San Francisco Estuary,
26 Trend analysis

27 **Introduction**

28 Accurate quantification of environmental trends must consider variation at different temporal
29 scales when ignoring variation at one scale could lead to incorrect conclusions about variation at
30 another scale. Many environmental monitoring programs collect temporally resolved but
31 irregular time series data to quantify trends for regulatory, management, or research purposes.

32 The mismatch between the scales of monitoring versus analysis questions or management goals
33 can present statistical challenges ([Cumming et al., 2006](#); [Forbes and Xie, 2018](#); [Urquhart et al.,](#)
34 [1998](#)). At short temporal scales typically less than a year, environmental systems exhibit
35 variability caused by multiple factors (e.g., weather events, management, or seasonal changes).

36 Such fluctuations may not be related to inter-annual trends or may not be well-suited to multi-
37 scale smoothing methods, yet they can perturb the time-series such that the sub-annual
38 fluctuations must be addressed within the trend analysis to allow accurate quantification of inter-
39 annual trends. Many trend analysis methods lack the flexibility to evaluate multiple independent
40 variables in an extendable structure that accommodates hypothesis testing at different time scales
41 of interest.

42 In this paper, we develop methods to estimate across-year trends of within-year features, such as
43 seasonal averages, while accounting for uncertainties across analysis steps. Our overarching goal

44 for this work was to develop a flexible set of tools for accurately characterizing inter-annual
45 changes in seasonally-averaged water quality metrics that can be robustly applied to diverse
46 time-series data. At the outset, we identified several specific requirements and priorities for the
47 trend analysis methods: 1) Complete propagation of uncertainty through the analysis; 2)
48 Separation of trends on different time scales; 3) Ability to estimate linear and nonlinear
49 responses; 4) Flexibility to evaluate multiple independent variables (and random effects),
50 although the examples herein include time as the only independent variable; and 5) Robust to
51 missing observations or varying sampling effort across years, which was considered a high
52 priority to allow the methods to be applied to diverse time-series datasets and monitoring
53 programs. Existing methods that begin to address some of the above requirements and priorities
54 can be generalized into four groups: seasonal Kendall tests (and other non-parametric tests);
55 seasonal trend decomposition using loess (STL); weighted regression on time, discharge, and
56 season (WRTDS); and generalized additive models (GAMs).

57 Seasonal Kendall tests and related non-parametric approaches have been used for decades in
58 water quality trend assessments to identify inter-annual, monotonic changes while accounting for
59 predictable patterns among seasons ([Cloern et al., 2007](#); [Helsel et al., 2020](#); [Hirsch et al., 1982](#)).
60 While seasonal Kendall and other non-parametric approaches have been among the most
61 commonly used methods in long-term water quality trend analysis ([Wan et al., 2017](#)), they do not
62 satisfy several of our requirements. For descriptive decomposition of long-term monitoring data,
63 these approaches assume seasonal patterns within years do not change and require regularly
64 spaced or balanced data. In addition, seasonal Kendall tests do not allow for additional
65 independent variables to explain variation, do not estimate a model that could be useful for other

66 purposes (e.g., prediction), and do not easily allow for propagation of uncertainty to other trend
67 analysis methods.

68 STL decomposes a time series into additive components of a long-term trend, a seasonal pattern,
69 and residuals (Cleveland et al., 1990; Cloern, 2018; Cloern and Jassby, 2010; Stow et al., 2015).

70 While useful and widely applied, this method does not address all of our requirements. STL
71 decomposition does not allow for incorporating explanatory variables other than time. In
72 addition, it is often characterized more as an algorithm of statistical steps than as a statistical
73 model with estimated parameters (e.g., Wan et al., 2017) and it does not usually estimate
74 standard errors to allow hypothesis testing (but see Hafen, 2010). STL methods may also over-
75 simplify trends into fixed components that do not change over time, e.g., a seasonal estimate that
76 is constant across years. This limitation presents challenges when addressing questions relevant
77 to long-term water quality data, such as timing of seasonal peaks that can suggest system
78 response to changing environmental conditions (Cloern and Jassby, 2010; Navarro et al., 2012).

79 The weighted regression on time, discharge, and season (WRTDS) method addresses the
80 problem of inflexibility in STL by using a more general local regression scheme (Beck et al.,
81 2018; Beck and Hagy, 2015; Hirsch et al., 2010; Hirsch et al., 2015). Designed for evaluating
82 water quality in rivers where separating the effect of discharge on constituent concentration is
83 important, WRTDS estimates a moving window regression model with components that allow
84 parameters to vary smoothly in relation to both time and discharge. This yields parameters that
85 are specific to season, year, and flow regime. The WRTDS approach is conceptually similar to
86 local kernel smoothing methods, with specific application to explanatory variables relevant for
87 water quality constituents (i.e., season, year, and discharge). Standard error estimates of
88 predictions from WRTDS are available through a block bootstrap approach applied to the model

89 results (Hirsch et al., 2015). Although a useful addition to the original method (Hirsch et al.,
90 2010), the approach requires extensive resampling using a previously fitted model. Alternative
91 methods that include standard error estimates simultaneously with model output may be
92 preferred for intensive or more iterative applications.

93 Finally, GAMs can satisfy the requirements and priorities identified above and were adopted as a
94 central component for the trend analyses herein. GAMs combine one or more smoothing splines
95 to model patterns in data and can be reasonably viewed as generalizing the concepts behind STL
96 and WRTDS (Haraguchi et al., 2015; He et al., 2006; Morton and Henderson, 2008; Murphy et
97 al., 2019; Pearce et al., 2011). The basis functions used to formulate GAMs can be customized
98 based on expected patterns in the data. Examples include cyclic splines, which can be used to
99 model seasonal patterns, and low-dimensional interactions (Wood, 2017). GAMs have added
100 flexibility because they can include both parametric (e.g., linear or quadratic) components and
101 non-parametric (spline) components. Multiple approaches have been developed to determine the
102 optimal degree of smoothness (Wood, 2004; 2017). These approaches are based on optimization
103 of out-of-sample prediction error, which addresses a key concern around methods like WRTDS
104 that do not have analogs for choosing optimal degrees of smoothing. GAMs can also produce
105 results comparable to those provided by WRTDS (Beck and Murphy, 2017) and have readily
106 obtainable uncertainty estimates. Further, GAMs have natural frequentist and Bayesian
107 interpretations, are naturally extensible to include random effects (i.e., generalized additive
108 mixed models or GAMMs), and have computationally efficient implementations (Wood, 2017).

109 GAMs have been applied previously to evaluate trends in water quality time series from long-
110 term monitoring programs (Haraguchi et al., 2015; Murphy et al., 2019). For example, Murphy et
111 al. (2019) used GAMs to decompose water quality time series from Chesapeake Bay into long-

112 term and seasonal trends (Murphy et al., 2019) and test trend hypotheses between two points in
113 time. Other studies of environmental time series with GAMs have addressed the use of
114 transformed response data (Yang and Moyer, 2020), serial correlation in high resolution data
115 (Morton and Henderson, 2008; Yang and Moyer, 2020), and quantifying time lags in
116 relationships between response and predictor variables (Lefcheck et al., 2017). The method
117 development and analyses described herein generalize the approach to analyzing trends of
118 seasonal spline features, describes the relationships among alternative spline formulations when
119 spline flexibility is allowed to vary (Wood, 2017, 2003) rather than being constrained *a priori* for
120 different time scales, and prioritizes full incorporation of uncertainty.

121 To incorporate the uncertainty of seasonal estimates into trend analysis, we integrated GAMs
122 with mixed-effects meta-analysis (Gasparrini et al., 2012; Sera et al., 2019). In this integration,
123 the GAMs framework addressed a critical need by providing an estimate of uncertainty (e.g., a
124 standard error) of seasonal averages, even in situations with irregular sampling and serial
125 correlation, which are common in time series data. Meta-analysis regression incorporates a
126 known (or estimated) standard error for each response datum. Usually meta-analysis is used
127 when each response summarizes a dataset from a separate study, along with its standard error,
128 and the meta-analysis looks for across-study patterns in effect size (i.e., Lortie, 2014). In this
129 study, each response summarizes one year or season of data, with standard error from the GAM,
130 and the meta-analysis looks for patterns across years. Thus, while meta-analysis methods are
131 most commonly associated with combining results from multiple studies into a larger analysis,
132 their key modeling step is propagation of uncertainty (Gasparrini et al., 2012; Sera et al., 2019).
133 To do this, meta-analysis makes use of a known (estimated) standard error for each response

134 datum, which is the required priority here to propagate standard errors from the GAM into a
135 regression of seasonal averages.

136 To describe the approach and demonstrate its utility, we analyze a 30 year record (1990-2019) of
137 biweekly-to-monthly chlorophyll-a concentration data, collected at 9 stations in the southern
138 portion of the San Francisco Estuary, California, USA. Approximately twice-monthly
139 monitoring has been conducted for several decades at fixed locations (stations) along the
140 longitudinal axis of the Bay. Analysis of these data is complicated by irregularities in timing and
141 consistency of data collection. We were interested in questions such as: Are there significant
142 trends in spring mean chlorophyll at multi-year time-scales? At what across-year window does
143 summer-fall mean chlorophyll levels change? Is there a spatial difference in chlorophyll trends?
144 We provide examples illustrating how these questions can be addressed using GAMs to estimate
145 seasonal patterns and use meta-analysis to evaluate trends between years. The techniques are
146 incorporated into an open-source and publicly available R package, *wqtrends*, developed by the
147 authors (Beck et al., 2021, available at <https://tbep-tech.github.io/wqtrends>, including an online
148 dashboard for viewing results at <https://nutrient-data.sfei.org/apps/SFbaytrends/>).

149 **Methods**

150 **Study area and data sources**

151 The San Francisco Estuary (SFE) is the largest estuary on the Pacific Coast of North America,
152 and its watershed covers 200 thousand km² in the US state of California. Flows from the
153 Sacramento-San Joaquin Delta, entering from the northeast, account for the vast majority of
154 SFE-wide annual-average freshwater inputs (Cloern and Jassby, 2012). Freshwater contributions

155 to southern SFE (South Bay, Lower South Bay) come primarily from local tributaries during the
156 wet season (Nov-Apr), and from wastewater treatment plant discharges during the dry season
157 (May-Oct). Salinity values in southern SFE subembayments (the focus of this work; Central
158 Bay, South Bay, Lower South Bay) range from 5 to 35 ppt and depend strongly on season
159 (stormwater runoff), tidal cycle, and effluent discharge from wastewater treatment plants ([Cloern
160 and Jassby, 2012](#)).

161 SFE receives 70,000 kg per day of dissolved inorganic nitrogen (DIN; annual average), with the
162 majority of that DIN coming from wastewater treatment plant discharges. Flows from the Delta
163 deliver 30,000 kg per day of DIN to the SFE (annual average), with the Delta's DIN load varying
164 over a 5-fold range annually ([SFEI, 2014a](#)). Based on its areal DIN loads, SFE ranks among the
165 most nutrient-enriched estuaries worldwide ([SFEI, 2014a, b](#); [Cloern et al., 2020](#)). Despite its
166 nutrient-enriched status, the SFE has generally not experienced some of the water quality
167 impacts common to other nutrient-enriched estuaries (e.g., excessive phytoplankton blooms, low
168 dissolved oxygen), with SFE's muted response attributed to its highly turbid waters (reduced
169 light penetration within the water column); strong tidal mixing (limiting duration of water
170 column stratification to less than several days); and strong phytoplankton grazing pressure from
171 abundant suspension feeding bivalves in some regions ([Alpine and Cloern, 1988](#); [Cole and
172 Cloern, 1984](#); [Jassby, 2008](#); [Kimmerer and Thompson, 2014](#)).

173 Studies over the past decade have identified changes in responses or sensitivity within SFE to
174 nutrients in deep subtidal habitats via increased phytoplankton biomass (chl-a) and gross primary
175 production (GPP) in South Bay ([Cloern et al., 2007, 2010](#)); recently documented occurrences of
176 harmful algae and their associated toxins ([Sutula et al., 2017](#); [Peacock et al., 2018](#)); and low
177 dissolved oxygen in some tidal slough habitats ([SFEI, 2021](#)). These observations have raised

178 concerns that SFE's resistance to its high nutrient inputs could be waning ([SFEI, 2014b](#)),
179 prompting regulators to initiate the SFB Nutrient Management Strategy ([SFBRWQCB, 2017](#)).
180 The early increases in South Bay chl-a (1995-2005) were quantitatively tested using seasonal
181 Kendall, and the signal was sufficiently large and coherent it was also visually apparent in raw
182 data. Given SFE's nutrient-enriched status, there is a critical need for on-going and
183 comprehensive characterization of trends in chl-a, including the ability to examine variability at
184 multiple time-scales and non-monotonic trends, that can also be readily applied to other nutrient-
185 related indicators (e.g., dissolved oxygen, GPP). Beyond the role these tools can play in
186 supporting improved understanding of system dynamics in SFE, water quality managers have
187 emphasized the importance of robust trend detection for informing future nutrient management
188 decisions.

189 For the trend analyses discussed below, we used near-surface (0-2 m) chlorophyll concentrations
190 ([Figure 1](#)) measured biweekly to monthly from 1990 to 2019 along the longitudinal axis of the
191 SFE extending from Central Bay (stations 18-23), South Bay (stations 24-32), and Lower South
192 Bay (stations 34-36) ([Table 1](#), [Figure 2](#)). Monitoring data were obtained from the SFE Research
193 Program of the US Geological Survey ([Cloern and Schraga, 2016](#); [Schraga et al., 2020](#)).
194 Sampling frequency varied somewhat over time and by station. Approximate monthly or
195 biweekly sampling with coverage of at least a decade is common for many long-term monitoring
196 programs and is the motivating use case for the methods herein. Every observation was included
197 directly in the statistical models without spatial or temporal binning or averaging. Log₁₀-
198 transformed chl-a was used for all analyses to meet assumptions of normally-distributed
199 residuals. Methods for back-transformation of model results are provided in the supplement.

200 **GAMs with uncertainty propagation**

201 We implemented our analysis in three stages. First, we used a GAM to estimate a smooth
202 temporal pattern in the raw data, along with the uncertainty of the smoother. Second, we
203 calculated a feature of interest from the estimated GAM, along with its propagated uncertainty.
204 For the examples described here, we focused on extracting seasonal averages of chl-a values.
205 Other features can also be extracted using the same tools, including the timing or magnitude of a
206 seasonal peak, but those are not presented here (see the *wqtrends* R package, [https://tbep-
207 tech.github.io/wqtrends](https://tbep-tech.github.io/wqtrends)). Third, we used a mixed-effects meta-analysis to estimate trends and
208 test hypotheses about the change in seasonal averages across years.

209 **First-stage analysis: GAM estimation**

210 To smooth the raw data across time, we considered and tested four different GAM structures.
211 While all four GAM structures can achieve similar fits, they differ in how they partition variation
212 in the time series (Table 2), which may unnecessarily influence understanding of temporal
213 patterns. We discuss all four to clarify their relationships and interpretations. All models were
214 created using the *mgcv* R package (R Core Team, 2020; Wood, 2017), with utility functions
215 included in the *wqtrends* package created by the authors (Beck et al., 2021).

216 The simplest GAM for this purpose is expressed as:

217
$$\text{Model S: } y_i \sim \beta_0 + f_1(\text{cont_year}_i) + \epsilon_i \quad (1)$$

218 where y is measured chl-a, β_0 is an intercept, and cont_year is “continuous year,” a continuous
219 numerical date (e.g., July 1st 2019 would be 2019.5). The $f_1()$ function is a smoothing spline
220 composed of the sum of multiple “basis functions” multiplied by coefficients. $f_1(\text{cont_year})$

221 describes the relationship of y with *cont_year* in a way that smoothly follows the data (Wood,
222 2017). The basis functions involve user-specified knots, a grid of values on the *cont_year* axis
223 that is discussed more below. The ϵ term represents residuals following a normal distribution
224 with mean zero and constant variance.

225 Smoothing was determined using generalized cross-validation (GCV, as implemented in `mgcv`),
226 which minimizes out-of-sample prediction error. GCV works by penalizing the net curvature of a
227 spline (Wood, 2004). To allow GCV (or other alternatives) to work as intended, the number of
228 knots chosen by the analyst, which determine the maximum degrees of freedom, must be
229 sufficiently large so that the curvature penalty, rather than the number of knots, determines
230 smoothness. Results should not be sensitive to the number of knots; if they are, the number of
231 knots should be increased. In the examples below, we chose the number of knots for $f_1()$ as 12
232 times the number of years in the time series, i.e., one knot per month. If the data were too sparse
233 to fit 12 knots per year, the number of knots was reduced by one knot per year until the model
234 could be estimated (i.e., 12 * years, 11 * years, etc.).

235 The next three spline formulations (Model SY, SYD, and SYDI) provide progressively
236 increasing complexity in how spline terms compose a model to smooth the raw data. Model SY
237 describes the time series using a linear trend plus a spline for *cont_year*:

238
$$\text{Model SY: } y_i \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + \epsilon_i \quad (2)$$

239 where equation (2) is the same as equation (1) with the addition of a linear term for *cont_year*
240 related to y_i by the β_1 slope parameter.

241 While Model SY contains the explicit linear trend term, $\beta_1 \text{cont_year}_i$, it is in fact
242 mathematically equivalent to model S (Table 3). The $f_1()$ spline for *cont_year* includes an

243 unpenalized linear trend, so a trend will be estimated in model S. When *cont_year* is included
244 explicitly as a linear term in model SY, `mgcv` adjusts the basis functions for the spline to exclude
245 the linear term, thereby not over-parameterizing the model. Whereas an estimated linear trend in
246 *cont_year* and its uncertainty can be extracted from the fitted spline in model S, model SY
247 provides this trend directly, giving the equivalent result. Further, package `mgcv` offers the option
248 to penalize linear trends in splines to provide a method for variable selection (option `select` =
249 `TRUE`), such as when numerous splines are included in the model formulation for variables that
250 may or may not be important. In the implementation of our approach described here, this option
251 is not used and all models specify `select` = `FALSE`. Details in the supplement explain this
252 justification.

253 Model SYD adds an average within-year cyclic pattern as a separate spline:

254 Model SYD: $y_i \sim \beta_0 + \beta_1 \text{cont_year}_i + f_1(\text{cont_year}_i) + f_2(\text{doy}_i) + \epsilon_i$ (3)

255 where equation (3) is the same as equation (2) with the addition of a smoothing spline for “day-
256 of-year” (*doy*, i.e., Julian date, a count starting January 1 for each year). For $f_2()$, a cyclic spline
257 is specified (using `bs` = '`cc`' in `mgcv`) to constrain the start and end at the same value. A user-
258 specified number of knots is also included in $f_2()$. While model SYD is not mathematically
259 equivalent to models S and SY, it should produce nearly identical results. The *doy* spline in
260 model SYD gives the average within-year pattern and changes the interpretation of the *cont_year*
261 spline to represent smoothed deviations from that pattern.

262 Models S, SY, and SYD can all potentially extract a similar signal from the raw data (Table 3).
263 What differs between the models is the allocation of penalties for curvature used to determine
264 smoothness for each spline. In model SYD, there are separate penalties for the two splines, as

265 compared to S and SY that include penalties only for the *cont_year* spline. This is important
266 because variation in the response variable can be differently attributed to each spline depending
267 on the model, even while the sum of components for each model produces similar results
268 between models. Our goal is to extract seasonal averages from the fitted time series that are not
269 sensitive to different allocation of penalties among the splines in each model.

270 If the fits were to differ substantially between model SYD and models S or SY, an interpretation
271 could be difficult because the penalties for smoothing splines based on curvature are heuristic
272 ([Wood, 2017](#)). For example, if a lower AIC is achieved in one model compared to another,
273 assuming both use sufficient knots, this may just reflect the outcome of alternative penalization
274 heuristics implied by the different formulations and does not imply one model fit is better. In the
275 examples here, model SYD achieves nearly identical fits to model S or SY, where the latter by
276 definition also achieve identical fits.

277 Model SYD has the appealing feature that, if some parts of some years have limited data, model
278 SYD will impute an average seasonal pattern with the *doy* spline, thereby considering data from
279 the same period in other years in the prediction of the period with missing data. However, an
280 interpretation of these imputations may be challenging. For example, the spring chl-a peak is a
281 notable feature every year in the SFE. If the peak occurs at the same time every year but the
282 magnitude varies, then the average within-year pattern can be interpreted as the average
283 magnitude. However, if the magnitude is the same but the timing varies across years, then the
284 magnitude of the average peak cannot be similarly interpreted; instead, the average peak
285 extracted by $f_2(doy_i)$ will underestimate the magnitude that usually occurs. Moreover, the width
286 or duration of the peak will be longer than typically occurs in a given year.

287 Finally, the raw data can be smoothed using a bivariate spline representing an interaction
288 between *cont_year* and *doy*. This can be expressed as:

289 Model SYDI:

290 $y_i \sim \beta_0 + \beta_1 cont_year_i + f_1(cont_year_i) + f_2(doy_i) + f_3(cont_year_i, doy_i) + \epsilon_i$ (4)

291 where equation (4) is the same as equation (3) with the addition of a tensor-product smoothing
292 spline (`ti()` in `mgcv`) that varies smoothly as a function of both *cont_year* and *doy*. Both
293 *cont_year* and *doy* include their own number of knots, such that the total number of knots for
294 the spline is the product of the two. The need for sufficient knots in SYDI can be satisfied either
295 by allowing for sufficiently many knots for $f_3()$ or sufficiently many knots for $f_1()$ or $f_2()$, but
296 not both, given limits on the model degrees of freedom.

297 Following the rationale above, the relationship of model SYDI to model S is similar to that of
298 model SYD to model S. Model SYDI differs from model S to a greater extent than model SYD,
299 but all of the splines use the same inputs to smooth the same data. The univariate splines in
300 *cont_year* and *doy* will combined likely not capture as much variation within model SYDI as
301 captured by model S, given the fewer knots that are available to $f_1()$ or $f_2()$ within SYDI. The
302 tensor-product spline represents an interaction by allowing the pattern in *cont_year* to vary by
303 *doy* and vice-versa. The interaction term provides an appearance that model SYDI is
304 fundamentally different from those provided by the other models. However, models S, SY, and
305 SYD also allow within-year fluctuations to vary across years by allowing a spline to be fit
306 through the entire time series. Although model SYDI is the only model that includes an explicit
307 interaction term, all of the models support the interaction conceptually. By providing this term
308 with sufficient knots, the raw data can be fully smoothed with model SYDI to a similar degree as

309 for the other models. However, a very large number of knots in both the *cont_year* spline *and* in
310 both dimensions of the interaction spline is impossible to achieve. The distinct aspect of model
311 SYDI is the anticipation that within-year fluctuations will vary smoothly from year to year.
312 However, this is an assumption about ecosystem dynamics that may not be appropriate to *a*
313 *priori* parameterize into a statistical model, including for SFE and many other estuaries where
314 bloom magnitude often varies between years. Thus, the conceptual motivation for model SYDI
315 and its practical application are not necessarily supported in the more generalized application for
316 which we are developing this set of analyses.

317 [Murphy et al. \(2019\)](#) used spline formulations for Chesapeake Bay water quality analyses that
318 are comparable to those proposed here, but for different goals and with different handling of
319 smoothness. They evaluated a “`gam0`” with only a cyclic spline for *doy* and linear *cont_year*
320 terms, a “`gam1`” like our SYD, and a “`gam2`” like our SYDI. In application, only “`gam2`” was
321 used, including the addition of splines as functions of hydrologic variables to account for finer-
322 scale variation. [Murphy et al. \(2019\)](#) allowed a maximum number of knots for the *cont_year*
323 spline ($f_1()$) of 2/3 times the number of years and do not explicitly consider the number of knots
324 in the interaction spline, following an *ad hoc* allocation of variation in the data to different
325 components based on previous interpretations of water quality dynamics in the system.
326 Constraining splines with insufficient knots could inflate Type I error rates for temporal changes
327 and we seek to lower this risk by increasing the upper limit for the knots for the $f_1()$ term.
328 Finally, [Murphy et al. \(2019\)](#) present large AIC differences between their spline formulations.
329 We instead emphasize that, given sufficient knots, the models represent alternative formulations
330 of conceptually similar explanations for the data and yield similar fits (Table 3), resulting in near
331 ties for AIC between models.

332 To evaluate the theoretical and conceptual similarities and differences among the GAM
333 structures discussed above, we visually compared chl-a estimates from models S, SYD, and
334 SYDI (Figure 3; note that SY is identical to S and is not shown). Models S, SYD, and SYDI
335 were fit to chl-a data from station 34 using a sufficiently high number of knots for the respective
336 splines for each model. As expected, predictions by day of year from each model are visually
337 similar (Figure 3a) and closely follow the 1:1 line (Figure 3b). However, the component of
338 predictions explained by the continuous year smoother ($f_1()$) differs between models that include
339 additional smoothers (Figure 3c), offering a graphical impression of the degree to which data
340 variability was partitioned to different GAM components. These results are also reflected in
341 differences in the effective degrees of freedom among the additive components of each model
342 (Table 3).

343 For all results, model S was used with enough knots in $f_1()$ to evaluate chl-a trends across the
344 monitoring stations in the SFE. This model was chosen because of the relatively faster
345 processing time to fit the model, while providing nearly identical explanatory power as compared
346 to the other models (Table 3).

347 **Second-stage analysis: Seasonal features with uncertainties**

348 In the second-stage analysis, we estimated a seasonal average, such as the mean spring chl-a
349 concentrations, along with its associated uncertainty, in each year. We defined μ_t as the seasonal
350 average in year t , $\hat{\mu}_t$ as an estimate of μ_t , and $\hat{\sigma}_{\hat{\mu},t}$ as the estimated standard error of $\hat{\mu}_t$. The
351 season includes n days. For simplicity, the following text omits subscript t .

352 Point estimates of response values for the fitted GAM take the form $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$, where $\hat{\boldsymbol{\beta}}$ is the
353 vector of parameter estimates and \mathbf{X} is a model matrix of explanatory variables, including spline

354 basis function values. Vector $\hat{\beta}$ includes both fixed effect parameters and spline parameters, and
355 \mathbf{X} contains columns corresponding to each. For example, using model SY, if a point estimate for
356 chl-a is needed for a single day, given as $cont_year = r$, then \mathbf{X} would have a row with 1 in the
357 first column (for the intercept parameter), r (for the linear time trend) in the second column, and
358 an evaluation of each spline basis function at r in the remaining columns. The number of spline
359 basis functions is related to the number of knots. Note that r can be any time, not necessarily the
360 time of an observation.

361 To obtain a vector, $\hat{\mathbf{y}}$, of fitted point estimates for every day in a season, \mathbf{X} would have one row
362 for each day. Here, the seasonal averages were calculated at the resolution of days. The estimated
363 spline yields both $\hat{\beta}$ and $\hat{\Sigma}_{\hat{\beta}}$, an estimate of the covariance matrix of the sampling distribution of
364 $\hat{\beta}$. The scalar standard errors of $\hat{\beta}$ are the square roots of the diagonal elements of $\hat{\Sigma}_{\hat{\beta}}$, whereas
365 the off-diagonal elements are the correlations among the elements of $\hat{\beta}$. Since parameter
366 estimates are correlated, the covariance of $\hat{\mathbf{y}}$ is $\hat{\Sigma}_{\hat{\mathbf{y}}} = \mathbf{X}\hat{\Sigma}_{\hat{\beta}}\mathbf{X}^T$.

367 The estimated seasonal average was calculated from the vector of daily values for each of the n
368 days in the season of interest with $\hat{\mu} = A^T\hat{\mathbf{y}}$, where A^T is a row vector with all values equal to
369 $1/n$. The variance of $\hat{\mu}$ is $\hat{\sigma}_{\hat{\mu}}^2 = A^T\hat{\Sigma}_{\hat{\mathbf{y}}}A$, and its standard error is $\hat{\sigma}_{\hat{\mu}}$. Each of these estimates are
370 from the approximate multivariate normality of the sampling distribution of $\hat{\beta}$.

371 **Third-stage analysis: Trend analysis of seasonal features with uncertainties**

372 In stage three of the analysis, we used a meta-analysis method to evaluate linear trends across
373 years of seasonal-average water quality, characterized by the within-year means ($\hat{\mu}_t$) and their
374 standard errors ($\hat{\sigma}_{\hat{\mu},t}$) that we estimated in stage two of the analysis. This analysis provided the

375 basis for directly answering the question *Is there a significant linear trend across a group of*
376 *years in a seasonal average, taking into account uncertainty?* For example, was there a trend in
377 the spring chl-a average from 1990 to 2000? This question can also be posed in a moving-
378 window manner across a time series (e.g., spring average trend from 1990-2000, 1991-2001,
379 etc.). For all analyses, the response data of interest are $\hat{\mu}_t$, $t = 1, \dots, N$, with their associated
380 standard errors, $\hat{\sigma}_{\hat{\mu},t}$. N is the number of years of the study.

381 A mixed-effects meta-analysis model can estimate linear trends when each observation has an
382 associated standard error, which is the case with our estimates $\hat{\mu}_t$ and $\hat{\sigma}_{\hat{\mu},t}$. Differences in
383 standard errors, which may result from different monitoring effort between years, are explicitly
384 considered in the analysis. The model can be expressed using notation similar to [Sera et al.](#)
385 ([2019](#)):

386
$$\hat{\mu}_t = \beta_0 + \beta_t t + b_t + \epsilon_t \quad (5)$$

387 where β_0 is the intercept, t is the year, β_t is the slope, b_t is the random effect for year t , and ϵ_t is
388 the residual for year t . Accordingly, the seasonal average for year t is $\mu_t = \beta_0 + \beta_t t + b_t$. The
389 “residual,” ϵ_t , represents estimation error in $\hat{\mu}_t$, namely $\hat{\mu}_t - \mu_t$. The residuals are assumed to be
390 independent and normally distributed with mean 0 and variances $\hat{\sigma}_{\hat{\mu},t}^2$, where the latter is
391 estimated from the calculations above. The random effect, b_t , is the difference between μ_t and
392 $\beta_0 + \beta_t t$ and is considered the “residual” in the sense of unexplained variation not due to the
393 estimation error. The random effect follows a normal distribution with mean 0 and variance, σ_b^2 ,
394 to be estimated.

395 We estimated the model (equation (5)) using the *mixmeta* package in R ([Sera et al., 2019](#)).

396 Results from *mixmeta* have a similar interpretation as those from regression analysis, but

397 parameter estimates and their standard errors incorporate the known standard errors of the
398 response values. The default estimation method for *mixmeta*, restricted maximum likelihood
399 (REML), was used. The meta-analysis models were applied to a chosen sequence or “window”
400 of years for estimating the linear trend.

401 **Trend comparisons**

402 The above methods were applied to each station by evaluating changes in seasonal averages for
403 approximately ten year moving windows from 1991 to 2019. The overall technique and *wqtrends*
404 R package allow for easy and flexible selection of seasonal windows. As an example focused on
405 illustrating the approach (as opposed to extracting mechanistic interpretations), we selected two
406 broad seasons, January-June and July-December, that are generally relevant to phytoplankton
407 bloom phenology in the SFE ([Cloern et al., 2020](#)), and combined include conditions over the full
408 year. A moving-window approach encoded within *wqtrends* was used to apply the meta-analysis
409 to each decadal window (e.g., 1991-2001, 1992-2002, etc.), allowing changes in slope and its
410 significance to be interpreted as the window is shifted one year at a time. We interpret the slope
411 as representative for the central year for each block, but a predictive trend for the final year of
412 the window could also be interpreted. As a means for summarizing some results, we focus on the
413 windows 1991-2000, 2000-2010, and 2010-2019.

414 Finally, trend results from the meta-analysis regression method for each season and different
415 time periods were compared to “naive” across-year regressions that do not propagate uncertainty
416 to demonstrate how different and potentially misleading conclusions can be obtained. Trend
417 estimates were compared to 1) trends from ordinary least squares (OLS) regression applied to
418 seasonal averages from the raw data and 2) trends from OLS regression applied to GAM

419 seasonal averages. Select examples where differences were pronounced were used for
420 illustration. This analysis was then applied to all stations. The method that formally propagates
421 uncertainty should have more robust statistical properties, such as accurate confidence interval
422 coverage, than naïve methods. For this reason, even when results are similar across methods, the
423 more robust method provides the best support for those results.

424 **Results**

425 **Model performance and predictions**

426 Model predictions for chl-a trends across all stations had an average R-squared value of 71%
427 (Table 4) and range from 59% (station 22) to 78% (station 18). GAM predictions from north to
428 south showed more pronounced annual and seasonal changes in chl-a towards the more southern
429 stations (Figures S1-S9). All the models suggested 1) increasing chl-a from 1990 until 2005 to
430 2010, followed by decreasing chl-a until the end of the record in 2019, 2) a spring chl-a peak,
431 particularly at southern stations, and 3) a fall chl-a peak that was smaller than the spring peak.
432 The magnitude of the fall peak did not vary noticeably by location (Figures S1-S9).

433 **Inter-annual trend estimates**

434 Estimates from the seasonal trend analyses (mixed-effects meta-analysis regressions) across
435 roughly ten-year windows for different seasons are shown for station 34 (Figure 4). Plots a-c
436 show trends in January to June averages while plots d-f show trends in July to December
437 averages. January to June chl-a increased (\log_{10} chl-a slope $0.03 \mu\text{g L}^{-1} \text{yr}^{-1}$, $0.01-0.05$ 95%
438 confidence interval) from 1991 to 2000, whereas July to December chl-a did not change

439 significantly. Chl-a also increased from 2000 to 2010, but only for July to December (\log_{10} slope
440 0.03, 0-0.05 95% confidence interval). Finally, chl-a decreased from 2010 to 2019 but again only
441 for July to December (\log_{10} chl-a slope -0.02, -0.04-0 95% confidence interval). Because the
442 trends were confined to certain times of the year, the seasonal estimates provide additional
443 information beyond coarser estimates that cover the entire year.

444 Temporal changes varied among regions of the Bay and were more pronounced at southern
445 stations. Figure 5 shows results from similar analyses as those in Figure 4, but applied to all
446 stations. The seasonal trend analyses showed that increases (based on $p < 0.05$) for the January
447 to June period were observed at stations 32 and 34 from 1991 to 2000; decreases were observed
448 at stations 30 and 32 from 2010 to 2019. For the July to December period, increases were
449 observed at stations 24, 27, 30, and 32 from 1991 to 2000 and stations 18, 21, 22, 24, and 34
450 from 2000 to 2010, whereas decreases were observed at stations 30, 32, and 34 from 2010 to
451 2019.

452 Results from a ten-year moving window comparison of seasonal trends provided additional
453 context on when significant changes were occurring at each station (Figure 6). Trends were
454 observed at all stations that followed a general pattern of increases early in the record followed
455 by decreases later in the record. Increases and decreases were observed in both the January to
456 June and July to December seasonal periods, with some notable exceptions. In particular, the
457 most southern stations (32, 34, 36) had increasing trends prior to 2005 that were more often
458 observed in the July to December period. Additionally, chl-a at the more northern stations has
459 not changed in recent years for either seasonal period. For most stations and seasonal periods, a
460 change from increasing to decreasing chl-a occurred around 2007.

461 **Importance of uncertainty propagation**

462 Results showing trend estimates from meta-analysis on GAM seasonal averages provided
463 different conclusions than those from either OLS regression through seasonal averages from raw
464 data (Figure 7 row 1) or OLS regression through GAM seasonal averages without uncertainty
465 propagation (Figure 7 row 2). Figure 7a shows trend estimates for station 30 for January to June
466 averages from 2000 to 2010. Only the meta-analysis regression results show a trend in this
467 example (based on $p < 0.05$). The OLS regression on observed estimates (top plot) and OLS
468 regression on GAM estimates (middle plot) did not identify trends. Figure 7b shows trend
469 estimates at station 34 for July to December averages from 1991 to 2000. Unlike the first
470 example, only the middle figure shows a trend, whereas the top and bottom plots do not show
471 trends. In both cases, only the meta-analysis results give reliable conclusions because of full
472 propagation of uncertainty across methods. Even in cases where the p-value threshold is not of
473 interest, the confidence intervals from the alternative methods will be inaccurate.

474 Applying the same comparison to all stations showed that different trend analysis methods
475 provided conflicting information on the magnitude and significance of the seasonal chl-a changes
476 in each decade (Figure 8). In many cases, the slope estimates were similar in magnitude, with
477 some exceptions at the more southern stations where the OLS estimates suggested a larger trend
478 than the meta-analysis methods. More importantly, differences in the magnitude of the
479 confidence intervals between the OLS models applied to the GAM averages and the meta-
480 analyses were also observed, reflecting the ability of the latter to more accurately assess
481 significance of trends by accounting for uncertainty in the average estimates.

482 **Discussion**

483 Propagation of uncertainty from within-year features of estimated GAMs to across-year trends
484 using mixed-effects meta-analysis is a new approach that can address different questions than
485 previous methods. Our approach has several advantages over more conventional approaches for
486 analysis of water quality data from long-term monitoring programs. GAMs are capable of
487 modelling time series with missing observations or irregular sampling which can complicate
488 trend assessment and comparison of trends between locations ([Junninen et al., 2004](#); [Racault et](#)
489 [al., 2014](#)). As noted above, non-parametric approaches (i.e., seasonal Kendall tests) are by far the
490 most common trend analysis methods applied to long-term water quality data ([Helsel et al.,](#)
491 [2020](#); [Hirsch et al., 1982](#)). These methods only assess the direction and significance of
492 comparisons across years, and importantly, do not account for full propagation of uncertainty
493 inherent in raw observations if the raw data are aggregated to meet test requirements.
494 Aggregation of raw data, e.g., averaging of observations within a year or season to comply with
495 the requirements of Kendall tests, risks loss of information by removing variation between
496 observations at smaller time scales. The logical outcome is increased risk of incorrect
497 conclusions from test results.
498 Incorrect conclusions on trends can have dramatic consequences for regulated parties under
499 existing water quality compliance frameworks ([Smith et al., 2001](#)). Our examples in Figures 7
500 and 8 demonstrate these risks if propagation of uncertainty from raw observations across
501 methods is unaccounted for in trend assessment. The “naïve” method using OLS regression
502 applied to seasonal averages from the raw observations fails to propagate uncertainty, similarly
503 to averaging results within a year and applying a simple Kendall test. In some cases the results

504 may be similar to those from fully propagating uncertainty, but the loss of information can lead
505 to increased Type I or II error rates depending on characteristics of the raw data and the method
506 used for their evaluation ([Shabman and Smith, 2003](#)). Our examples demonstrated the increased
507 potential for incorrect conclusions at specific monitoring locations and, at larger spatial scales,
508 across all stations if simpler trend analyses are used. Even though simpler methods may produce
509 similar results in some cases, particularly with frequent sampling and similar effort between
510 years, the only way to confirm such an outcome would be to compare results, relying on the
511 method with full propagation of uncertainty to be the more robust method. Use of methods that
512 fully account for uncertainty is recommended to obtain statistically valid results in a wider range
513 of conditions.

514 Results here also show that GAM structure (i.e., choice of smoothing terms) was less important
515 than allowing the model sufficient freedom to fit the data. This is an important conclusion that
516 provides guidance on how GAMs could be used to model time series from long-term
517 environmental monitoring programs. Models with separate smoothers for continuous year and
518 day of year can produce nearly identical results in the predicted trends if the knots are
519 sufficiently high to allow the GAMs to be fit as intended by the methods in the mgcv package
520 (Figure 3). The approach presented here leverages the ability of GAMs to objectively estimate
521 smoothed trends across years by identifying an optimal level of smoothing using generalized
522 cross-validation to extract an underlying signal in the observed data ([Wood, 2017, 2004](#)).

523 The underlying cross-validation methods used by GAMs in the mgcv package also reduce the
524 decisions that may be necessary for the implementation of alternative trend assessment methods.
525 For example, WRTDS and similar smoothing approaches (e.g., LOESS) require decisions on
526 appropriate window widths or bandwidths to define the neighborhood of observations for

527 smoothing (Hirsch et al., 2010; Wan et al., 2017). This is especially problematic for policy
528 analysis or regulatory decisions if the results change based on arbitrary decisions of the analyst.
529 Because these decisions are not needed for GAMs, the results can be considered a more objective
530 and potentially accurate signal of actual trends that are minimally influenced by process or
531 observation error present in the raw data.

532 Several limitations of the proposed methods deserve mention. First, if sampling is so irregular
533 that important fluctuations are missed entirely in some years, the GAM estimates and uncertainty
534 propagation could become dubious in interpretation and usefulness. Second, estimation of GAMs
535 for very long series can be computationally demanding. When this is an obstacle, one could do
536 the first two analysis stages using temporal windows of the full data, with the only implication
537 being that different degrees of smoothness may be estimated for different windows, which
538 indeed might be justified by the data. Third, meta-analysis regression results for a very small
539 number of years, particularly confidence intervals and associated p-values, may be inaccurate
540 (e.g. in confidence interval coverage). In such cases, one could make alternative use of the GAM
541 seasonal averages and standard errors, such as for pairwise comparisons among years.

542 **Future work**

543 Additional work could be conducted to further strengthen the conclusions based on trends from
544 meta-analysis regression applied to the GAM seasonal averages. Our third stage analyses require
545 *a priori* decisions on long-term time scales of interest and future work could generalize these
546 choices. Although there are undoubtedly many scenarios where years of interest can be chosen
547 objectively by the needs of an analysis (e.g., regulatory compliance periods, time since
548 management intervention), a more general question of when changes occur independent of user

549 decisions is also important to address. Additional methods could be developed using objective
550 criteria to identify inflection points or other important periods where changes occur independent
551 of a user choice. Assessing water quality changes beyond an evaluation of seasonal averages
552 could also be possible with our approach, such as assessing changes in the timing or magnitude
553 of a seasonal peak across years.

554 Additional explanatory variables could be identified to explain trends in either the GAM stage or
555 the meta-analysis stage of analysis. This information would have obvious implications for
556 management decisions on factors that influence water quality changes, e.g., wastewater treatment
557 upgrades, large-scale climatic factors, or flow regulation practices. Including alternative
558 predictors in the GAMs ([Wood and Augustin, 2002](#); [Zuur et al., 2009](#)) could reduce the
559 uncertainties of its estimates and, if relevant, allow the influence of those variables to be
560 removed. Including alternative predictors in the meta-analysis could help in explaining long-term
561 trends of seasonal averages or other metrics obtained from the GAMs. Our goal here was to
562 describe chl-a changes relative to time, so the single predictor in both modeling stages was time.

563 Finally, the evaluation of trends for alternative water quality variables in addition to chl-a is a
564 simple and logical extension of the methods proposed in this study. The long-term monitoring
565 program maintained by USGS includes multiple parameters in addition to chl-a that can provide
566 additional context into broader water quality trends in the SFE ([Cloern and Schraga, 2016](#);
567 [Schraga et al., 2020](#)). These parameters include salinity, temperature, light attenuation, dissolved
568 oxygen, suspended particulate matter, and dissolved inorganic nutrients, which collectively can
569 be used to provide a broader understanding of potential eutrophication patterns or ecosystem
570 shifts at seasonal and multi-decadal scales. Chl-a measurements can also be used to estimate
571 gross primary production to assess process rates that may be more indicative of system function

572 (Cloern et al., 2007; Jassby et al., 2002). The open-source *wqtrends* R package (Beck et al.,
573 2021) developed for this manuscript can be used for these analyses to provide additional insight
574 into potential drivers of water quality change in the SFE and other estuarine systems.

575 **Acknowledgments**

576 This work was supported by funding from the San Francisco Bay Nutrient Management Strategy
577 (NMS). We thank the staff of the US Geological Survey that collect and maintain long-term
578 monitoring data in San Francisco Bay. This work benefited from discussions with the NMS
579 Nutrient Technical Workgroup and Steering Committee. We thank James D. Hagy III and two
580 anonymous reviewers for providing helpful comments on this manuscript.

581 **References**

- 582 Alpine, A.E., Cloern, J.E., 1988. Phytoplankton growth rates in a light-limited environment, San
583 Francisco Bay. *Marine Ecology Progress Series* 44, 167–173.
- 584 Beck, M.W., de Valpine, P., Murpy, R., Wren, I., Chelsky, A., Foley, M., Senn, D., 2021. tbep-
585 tech/wqtrends: v1.1.0 (Version v1.1.0). Zenodo. <http://doi.org/10.5281/zenodo.4509638>.
- 586 Beck, M.W., Hagy, J.D., III, 2015. Adaptation of a weighted regression approach to evaluate
587 water quality trends in an estuary. *Environmental Modelling and Assessment* 20, 637–655.
588 <https://doi.org/10.1007/s10666-015-9452-8>
- 589 Beck, M.W., Jabusch, T.W., Trowbridge, P.R., Senn, D.B., 2018. Four decades of water quality
590 change in the upper San Francisco Estuary. *Estuarine, Coastal and Shelf Science* 212, 11–22.
591 <https://doi.org/10.1016/j.ecss.2018.06.021>
- 592 Beck, M.W., Murphy, R.R., 2017. Numerical and qualitative contrasts of two statistical models
593 for water quality change in tidal waters. *Journal of the American Water Resources Association*
594 53, 197–219. <https://doi.org/10.1111/1752-1688.12489>
- 595 Cleveland, R.B., Cleveland, W.S., McRae, J.E., Terpenning, I., 1990. STL: A seasonal-trend
596 decomposition procedure based on Loess. *Journal of Official Statistics* 6, 3–73.
- 597 Cloern, J.E. 2018. Patterns, pace, and processes of water-quality variability in a long-studied
598 estuary. *Limnology and Oceanography* 64, S192-S208. <https://doi.org/10.1002/lno.10958>
- 599 Cloern, J.E., Jassby, A.D., 2012. Drivers of change in estuarine-coastal ecosystems: Discoveries
600 from four decades of study in San Francisco Bay. *Reviews of Geophysics* 50.
601 <https://doi.org/10.1029/2012RG000397>
- 602 Cloern, J.E., Jassby, A.D., 2010. Patterns and scales of phytoplankton variability in estuarine-
603 coastal ecosystems. *Estuaries and Coasts* 33, 230–241.
- 604 Cloern, J.E., Jassby, A.D., Thompson, J.K., Hieb, K.A., 2007. A cold phase of the East Pacific
605 triggers new phytoplankton blooms in San Francisco Bay. *Proceedings of the National Academy
606 of Sciences of the United States of America* 104, 18561–18565.
- 607 Cloern, J.E., Schraga, T.S., 2016. USGS measurements of water quality in San Francisco Bay
608 (CA), 1969-2015: U.S. Geological Survey data release. <https://doi.org/10.5066/F7TQ5ZPR>.
- 609 Cloern, J.E., Shcraga, T.S., Nejad, E., Martin, C., 2020. Nutrient status of San Francisco Bay and
610 its management implications. *Estuaries & Coasts* 43, 1299–1317.
611 <https://doi.org/10.1007/s12237-020-00737-w>
- 612 Cole, B.E., Cloern, J.E., 1984. Significance of biomass and light availability to phytoplankton
613 productivity in San Francisco Bay. *Marine Ecology Progress Series* 17, 15–24.
- 614 Cumming, G.S., Cumming, D.H.M., Redman, C.L., 2006. Scale mismatches in social-ecological
615 systems: Causes, consequences, and solutions. *Ecology and Society* 11, 14.

- 616 Forbes, D.J., Xie, Z., 2018. Identifying process scales in the Indian River Lagoon, Florida using
617 wavelet transform analysis of dissolved oxygen. Ecological Complexity 36, 149–167.
618 <https://doi.org/10.1016/j.ecocom.2018.07.005>
- 619 Gasparrini, A., Armstrong, B., Kenward, M.G., 2012. Multivariate meta-analysis for non-linear
620 and other multi-parameter associations. Statistics in Medicine 31, 3821–3839.
621 <https://doi.org/10.1002/sim.5471>
- 622 Hafen, R.P., 2010. Local regression models: Advancements, applications, and new methods
623 (PhD thesis). Purdue University, West Lafayette, Indiana.
- 624 Haraguchi, L., Carstensen, J., Abreu, P.C., Odebrecht, C., 2015. Long-term changes of the
625 phytoplankton community and biomass in the subtropical shallow Patos Lagoon Estuary, Brazil.
626 Estuarine, Coastal and Shelf Science 162, 76–87.
- 627 He, S., Mazumdar, S., Arena, V.C., 2006. A comparative study of the use of GAM and GLM in
628 air pollution research. Environmetrics 17, 81–93. <https://doi.org/10.1002/env.751>
- 629 Helsel, D.R., Hirsch, R.M., Ryberg, K.R., Archfield, S.A., Gilroy, E.J., 2020. Statistical methods
630 in water resources, 2nd ed. U.S. Geological Survey Techniques; Methods, book 4, chapter A3,
631 version 1.1, Reston, Virginia.
- 632 Hirsch, R.M., Archfield, S.A., De Cicco, L.A., 2015. A bootstrap method for estimating
633 uncertainty of water quality trends. Environmental Modelling and Software 73, 148–166.
634 <https://doi.org/10.1016/j.envsoft.2015.07.017>
- 635 Hirsch, R.M., Moyer, D.L., Archfield, S.A., 2010. Weighted regressions on time, discharge, and
636 season (WRTDS), with an application to Chesapeake Bay river inputs. Journal of the American
637 Water Resources Association 46, 857–880.
- 638 Hirsch, R.M., Slack, J.R., Smith, R.A., 1982. Techniques of trend analysis for monthly water
639 quality data. Water Resources Research 18, 107–121.
- 640 Jassby, A.D., 2008. Phytoplankton in the Upper San Francisco Estuary: Recent biomass trends,
641 their causes, and their trophic significance. San Francisco Estuary and Watershed Science 6, 1–
642 24.
- 643 Jassby, A.D., Cloern, J.E., Cole, B.E., 2002. Annual primary production: Patterns and
644 mechanisms of change in a nutrient-rich tidal ecosystem. Limnology and Oceanography 47, 698–
645 712.
- 646 Junninen, H., Niska, H., Tuppurainen, K., Ruuskanen, J., Kolehmainen, M., 2004. Methods for
647 imputation of missing values in air quality data sets. Atmospheric Environment 38, 2895–2907.
648 <https://doi.org/10.1016/j.atmosenv.2004.02.026>
- 649 Kimmerer, W.J., Thompson, J.K., 2014. Phytoplankton growth balanced by clam and
650 zooplankton grazing and net transport into the low-salinity zone of the San Francisco Estuary.
651 Estuaries and Coasts 37, 1202–1218.

- 652 Lefcheck, J.S., Wilcox, D.J., Murphy, R.R., Marion, S.R., Orth, R.J., 2017. Multiple stressors
653 threaten the imperiled coastal foundation species eelgrass (*zostera marina*) in Chesapeake Bay,
654 USA. Global Change Biology 23, 3474–3483. <https://doi.org/10.1111/gcb.13623>
- 655 Lortie, C.J., 2014. Formalized synthesis opportunities for ecology: Systematic reviews and meta-
656 analyses. OIKOS 123, 897–902. <https://doi.org/10.1111/j.1600-0706.2013.00970.x>
- 657 Morton, R., Henderson, B.L., 2008. Estimation of nonlinear trends in water quality: An
658 improved approach using generalized additive models. Water Resources Research 44, W07420.
659 <https://doi.org/10.1029/2007WR006191>
- 660 Murphy, R.R., Perry, E., Harcum, J., Keisman, J., 2019. A Generalized Additive Model
661 Approach to evaluating water quality: Chesapeake Bay case study. Environmental Modelling &
662 Software 118, 1–13. <https://doi.org/10.1016/j.envsoft.2019.03.027>
- 663 Navarro, G., Caballero, I., Prieto, L., Vázquez, A., Flecha, S., Huertas, I.E., Ruiz, J., 2012.
664 Seasonal-to-interannual variability of chlorophyll-*a* bloom timing associated with physical
665 forcing in the Gulf of Cádiz. Advances in Space Research 50, 1164–1172.
666 <https://doi.org/10.1016/j.asr.2011.11.034>
- 667 Peacock, M.B., Gibble, C.M., Senn, D.B., Cloern, J.E., Kudela, R.M. 2018. Blurred lines:
668 Multiple freshwater and marine algal toxins at the land-sea interface of San Francisco Bay,
669 California. Harmful Algae 73, 138–147. <https://doi.org/10.1016/j.hal.2018.02.005>
- 670 Pearce, J.L., Beringer, J., Nicholls, N., Hyndman, R.J., Tapper, N.J., 2011. Quantifying the
671 influence of local meteorology on air quality using generalized additive models. Atmospheric
672 Environment 45, 1328–1336. <https://doi.org/10.1016/j.atmosenv.2010.11.051>
- 673 R Core Team, 2020. R: A language and environment for statistical computing. R Foundation for
674 Statistical Computing, R v4.0.2, Vienna, Austria.
- 675 Racault, M.F., Sathyendranath, S., Platt, T., 2014. Impact of missing data on the estimation of
676 ecological indicators from satellite ocean-colour time series. Remote Sensing of Environment
677 152, 15–28. <https://doi.org/10.1016/j.rse.2014.05.016>
- 678 Schraga, T.S., Nejad, E.S., Martin, C.A., Cloern, J.E., 2020. USGS measurements of water
679 quality in San Francisco (CA), beginning in 2016 (ver. 3.0, March 2020): U.S. Geological
680 Survey data release. <https://doi.org/10.5066/F7D21WGF>.
- 681 Sera, F., Armstrong, B., Blangiardo, M., Gasparini, A., 2019. An extended mixed-effects
682 framework for meta-analysis. Statistics in Medicine 38, 5429–5444.
683 <https://doi.org/10.1002/sim.8362>
- 684 SFBRWQCB (San Francisco Bay Regional Water Quality Control Board), 2017. Water quality
685 control plan (basin plan) for the San Francisco Bay basin. Prepared by California Regional Water
686 Quality Control Board, San Francisco Bay Region, Oakland, CA.
687 https://www.waterboards.ca.gov/sanfranciscobay/basin_planning.html
- 688 SFEI (San Francisco Estuary Institute), 2014a. External nutrient loads to San Francisco Bay.
689 SFEI Contribution Number 704. San Francisco Estuary Institute, Richmond, CA.

- 690 SFEI (San Francisco Estuary Institute), 2014b. Scientific foundation for the San Francisco Bay
691 Nutrient Management Strategy. SFEI Contribution Number 979. San Francisco Estuary Institute,
692 Richmond, CA.
- 693 SFEI (San Francisco Estuary Institute), 2021. Connections to tidal marsh and restored salt ponds
694 drive seasonal and spatial variability in ecosystem metabolic rates in Lower South San Francisco
695 Bay. SFEI Contribution No. 992. San Francisco Estuary Institute, Richmond, CA.
- 696 Shabman, L., Smith, E., 2003. Implications of applying statistically based procedures for water
697 quality assessment. *Journal of Water Resources Planning and Management* 129, 330–336.
698 [https://doi.org/10.1061/\(ASCE\)0733-9496\(2003\)129:4\(330\)](https://doi.org/10.1061/(ASCE)0733-9496(2003)129:4(330))
- 699 Smith, E.P., Ye, K., Hughes, C., Shabman, L., 2001. Statistical assessment of violations of water
700 quality standards under section 303 (d) of the Clean Water Act. *Environmental science &*
701 *technology* 35, 606–612. <https://doi.org/10.1021/es001159e>
- 702 Stow, C.A., Cha, Y., Johnson, L.T., Confesor, R., Richards, R.P., 2015. Long-term and seasonal
703 trend decomposition of Maumee River nutrient inputs to western Lake Erie. *Environmental*
704 *Science and Technology* 49, 3392–3400. <https://doi.org/10.1021/es5062648>
- 705 Sutula, M.A., Kudela, R.M., Hagy, J.D., III, Harding, L.W., Jr., Senn, D.B., Cloern, J.E.,
706 Bricker, S., Berg, G.M., Beck, M.W. 2017. Novel analyses of long-term data provide a scientific
707 basis for chlorophyll-a thresholds in San Francisco Bay. *Estuarine, Coastal and Shelf Science*.
708 197, 107-118. <https://doi.org/10.1016/j.ecss.2017.07.009>
- 709 Urquhart, N.S., Paulsen, S.G., Larsen, D.P., 1998. Monitoring for policy-relevant regional trends
710 over time. *Ecological Applications* 8, 246–257. [https://doi.org/10.1890/1051-0761\(1998\)008\[0246:MFPRRO\]2.0.CO;2](https://doi.org/10.1890/1051-0761(1998)008[0246:MFPRRO]2.0.CO;2)
- 712 Wan, Y., Wan, L., Li, Y., Doering, P., 2017. Decadal and seasonal trends of nutrient
713 concentration and export from highly managed coastal catchments. *Water Research* 115, 180–
714 194.
- 715 Wood, S.N., 2017. Generalized additive models: An introduction with r, 2nd ed. Chapman; Hall,
716 CRC Press, London, United Kingdom.
- 717 Wood, S.N., 2004. Stable and efficient multiple smoothing parameter estimation for generalized
718 additive models. *Journal of the American Statistical Association* 99, 673–686.
719 <https://doi.org/10.1198/016214504000000980>
- 720 Wood, S.N., 2003. Thin-plate regression splines. *Journal of the Royal Statistical Society (B)* 65,
721 95–114. <https://doi.org/10.1111/1467-9868.00374>
- 722 Wood, S.N., Augustin, N.H., 2002. GAMs with integrated model selection using penalized
723 regression splines and applications to environmental modelling. *Ecological Modelling* 157, 157–
724 177. [https://doi.org/10.1016/S0304-3800\(02\)00193-X](https://doi.org/10.1016/S0304-3800(02)00193-X)
- 725 Yang, G., Moyer, D.L., 2020. Estimation of nonlinear water-quality trends in high-frequency
726 monitoring data. *Science of The Total Environment* 715, 10.1016/j.scitotenv.2020.136686.

727 Zuur, A.F., Ieno, E.N., Walker, N.J., Saveliev, A.A., Smith, G.M., 2009. Mixed effects models
728 and extensions in ecology with r. Springer-Verlag, New York, New York.

729

730 **Figures**
731

732 *Figure 1: Observed chl-a concentrations for all stations in central and south San Francisco
733 Estuary (18-36, Figure 2), with (a) annual summer/fall concentrations (Aug - Dec) and (b)
734 monthly concentrations by decade.*

735 *Figure 2: Station locations in the central and south San Francisco Estuary used for analysis. See
736 Table 1 for station descriptions. Full dataset described in Schraga et al. (2020).*

737 *Figure 3: GAM output of estimated chl-a at station 32 for models S, SYD, and SYDI. Model SY is
738 identical to S and is not shown. Plots in (a) show model predictions by day of year with separate
739 lines for each year. Plots in (b) show pairwise comparisons of predicted chl-a between the
740 models and plots in (c) show the same comparisons as in (b) but only for results from the
741 estimated smoother for the cont_year variable. The plots demonstrate that results between the
742 models are similar except for a few observations at extreme values (a, b), but they vary in how
743 they allocate contributions to the predictions among different additive splines (c). The 1:1 lines
744 are in red to facilitate comparisons.*

745 *Figure 4: Examples of seasonal averages and trend estimates in ten year blocks from meta-
746 analyses using results of GAM predictions for station 34. Plots (a), (b), and (c) show trend
747 estimates for January to June averages and (d), (e), and (f) show trend estimates for July to
748 December averages. The trend lines estimate the rate of change of chl-a per year, reported as
749 the log₁₀-slope (+/- 95 % confidence interval) in the sub-plot titles. ns: not significant at $\alpha =$
750 0.05, * $p < 0.05$*

751 *Figure 5: Interannual trend estimates of seasonal averages by decade for chl-a at each station.
752 Point type and color represent the direction and magnitude of an estimated trend as the log₁₀
753 slope for chl-a concentration per year. Trends with $p < 0.05$ are marked with an asterisk. All
754 results are from Model S.*

755 *Figure 6: Estimates of log₁₀ chl-a change per year (+/- 95% confidence interval) from applying
756 the meta-analysis across the seasonal averages for each station. Stations are arranged top to
757 bottom from north to south. Plots in (a) show estimates for seasonal averages from January to
758 June and plots in (b) show estimates for seasonal averages from July to December. Results are
759 from a ten-year, centered moving window where each point shows a linear trend estimate from
760 five years prior to five years after each year. Estimates prior to 1996 and after 2014 are not
761 available because of an incomplete ten year record for estimating the trend. Significant estimates
762 are shown in red.*

763 *Figure 7: Trend estimate comparisons (arithmetic scale) for three models applied to seasonal
764 averages of chl-a in different annual periods at stations (a) 30 and (b) 34. The first row shows
765 OLS (ordinary least squares) regression applied to seasonal averages of chl-a from the raw
766 data, the second row shows OLS regression applied to seasonal averages of chl-a from the GAM
767 (without error propagation), and the third row shows meta-analysis regression applied to the
768 seasonal averages of chl-a from the GAM. Regressions in each plot are fit through the seasonal
769 estimates indicated in the plot titles for a specified year range.*

770 *Figure 8: Trend estimate comparisons for three models applied to seasonal averages of chl-a in
771 different annual periods at each station. The “OLS raw” trend model is based on an ordinary*

772 least squares (OLS) regression fit to the seasonal averages of chl-a from the raw data, the “OLS
773 GAM” trend model is based on an OLS regression fit to the seasonal averages of chl-a from the
774 GAM model (without error propagation), and the “Meta-analysis GAM” trend model is based
775 on a meta-analysis regression fit to the seasonal averages of chl-a from the GAM model. Values
776 for each model are the \log_{10} -slope estimates (+/- 95% confidence interval) as annual change per
777 year within each season, with line style denoting trend significance.

778 **Tables**779 *Table 1: Station locations, sample sizes (from 1991 to 2019), and summary values (median,*
780 *minimum, maximum) for chl-a ($\mu\text{g L}^{-1}$). Stations are arranged from north to south.*

Station	Latitude	Longitude	n	Med.	Min.	Max.
18	37.836	-122.418	414	3.6	0.2	16.6
21	37.784	-122.351	576	4.4	0.6	40.0
22	37.752	-122.351	569	4.0	0.7	53.1
24	37.686	-122.334	595	4.2	0.7	47.3
27	37.617	-122.285	596	4.5	0.5	50.9
30	37.551	-122.184	608	5.1	0.8	112.2
32	37.517	-122.133	591	5.9	0.7	282.1
34	37.485	-122.086	544	6.5	0.6	158.3
36	37.468	-122.067	476	6.2	1.1	328.4

781

782 *Table 2: Summary and details for each of the GAM structures. In practice, a sufficiently large*
 783 *number of knots provided to the additive terms will produce identical or comparable estimates*
 784 *for a response variable. The models differ in the allocation of penalties for the smoothness of*
 785 *each spline.*

GAM	Additive components	Details
S	$f_1(\text{cont_year})$	A single smoother over a continuous year variable
SY	$\beta_1 \text{cont_year} + f_1(\text{cont_year})$	A linear continuous year variable and a single smoother over a continuous year variable
SYD	$\beta_1 \text{cont_year} + f_1(\text{cont_year}) + f_2(\text{doy})$	A linear continuous year variable, a smoother over a continuous year variable, and a smoother over a day of year variable
SYDI	$\beta_1 \text{cont_year} + f_1(\text{cont_year}) + f_2(\text{doy}) + f_3(\text{cont_year}, \text{doy})$	A linear continuous year variable, a smoother over a continuous year variable, a smoother over a day of year variable, and an interaction smoother across continuous year and day of year variables

786

787 *Table 3: Comparison of the four model structures (S, SY, SYD, SYDI) described in the first stage*
 788 *analysis of GAM estimation. The four models provide either identical or comparable ability to*
 789 *describe chl-a trends at an example station (32) in the southern end of the San Francisco*
 790 *Estuary. The models differ in additive smoothers and the amount of effective degrees of freedom*
 791 *(edf) in the smoothers (measure of wiggliness in each component), but the overall model*
 792 *predictions are similar. AIC: Akaike Information Criterion, GCV: generalized cross-validation*
 793 *score, R2: r-squared values for predictions, edf: effective degrees of freedom, F: F-statistic, p-*
 794 *val: probability value, ** p < 0.001*

model	AIC	GCV	R2	smoother	edf	F	p-val
S	-138.2	0.06	0.74	s(cont_year)	242.23	6.27	**
SY	-138.2	0.06	0.74	s(cont_year)	242.23	6.27	**
SYD	-135.66	0.06	0.74	s(cont_year) s(doy)	229.33 8.07	3.88 0.11	**
SYDI	-123.57	0.06	0.73	s(cont_year) s(doy)	136.88 9.54	3.13 0.79	**
				ti(cont_year,doy)	69.31	0.74	**

795

796 *Table 4: Model performance statistics for each station as Akaike Information Criterion scores*
797 *(AIC), generalized cross-validation scores (GCV), and r-squared values.*

station	AIC	GCV	R-squared
18	-430.94	0.04	0.78
21	-305.93	0.04	0.70
22	-160.07	0.05	0.59
24	-250.77	0.04	0.69
27	-188.72	0.05	0.72
30	-172.79	0.05	0.74
32	-138.20	0.06	0.74
34	-3.17	0.07	0.68
36	-22.05	0.07	0.73

798

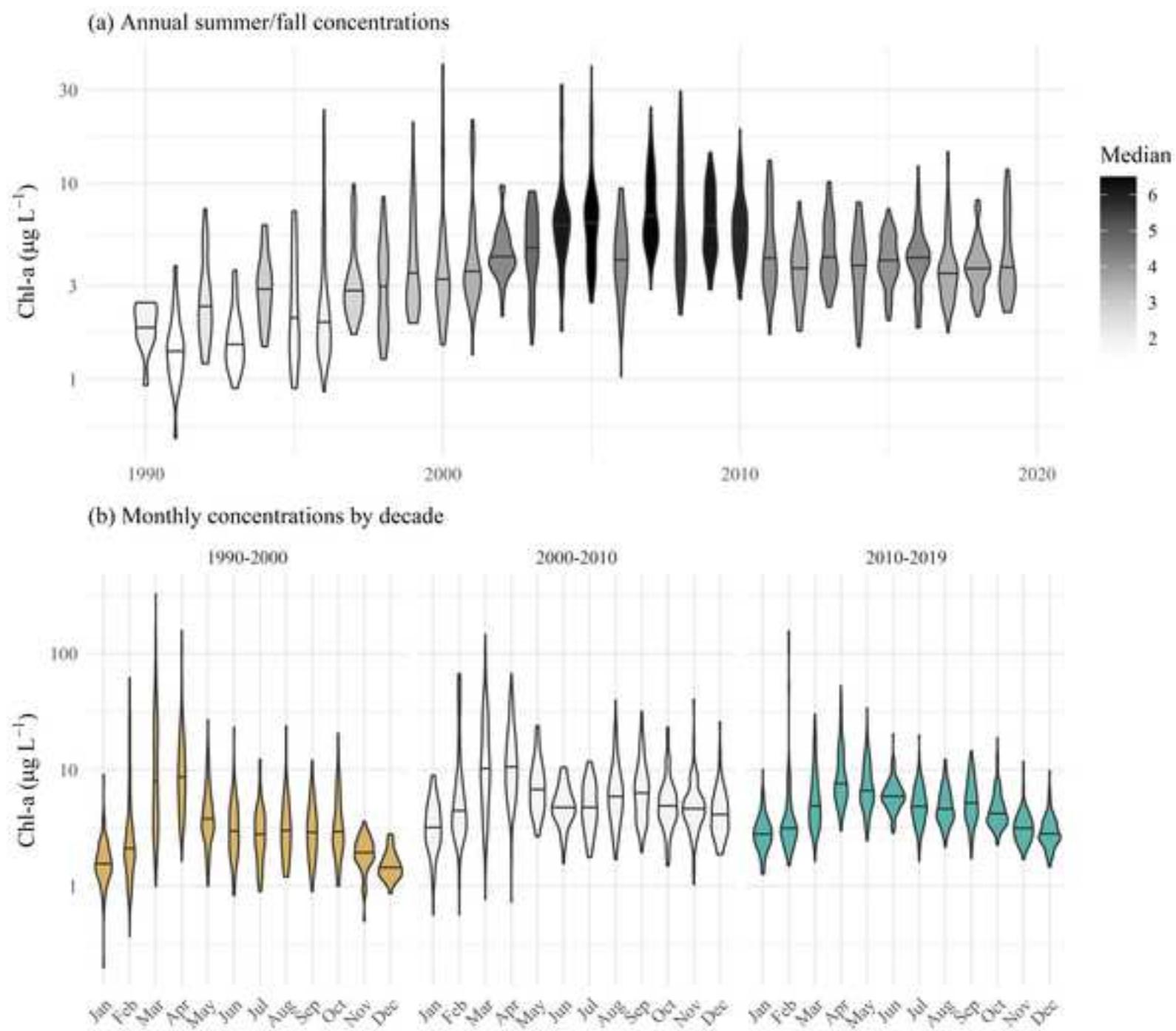
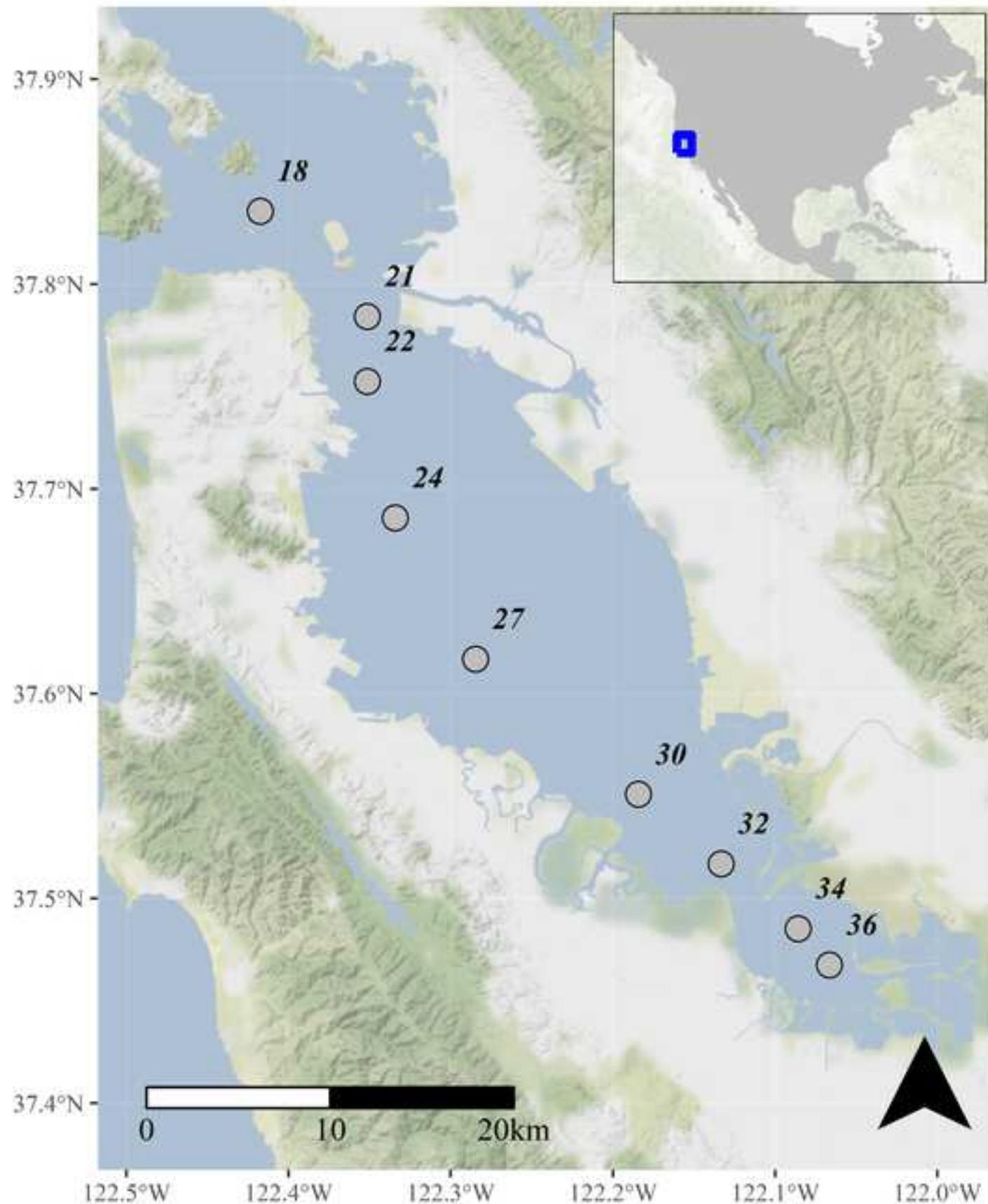
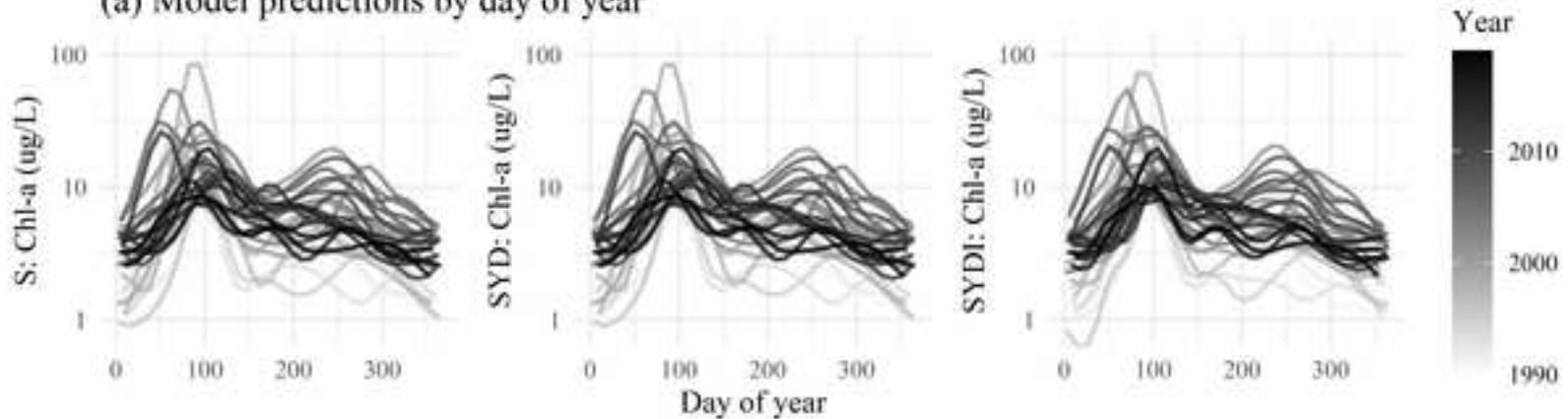


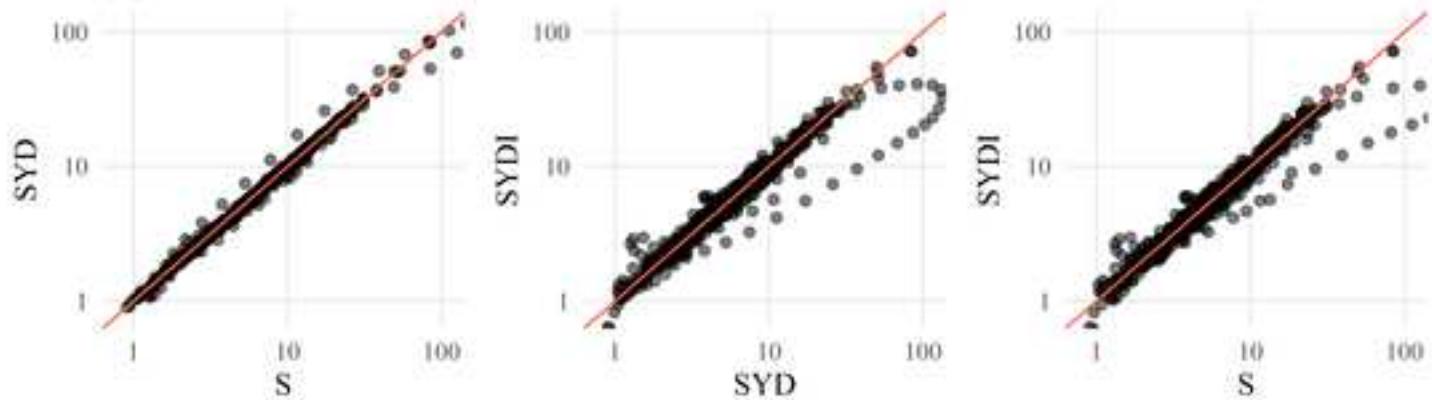
Figure 2

[Click here to access/download/](#)[Figure;Fig2.tif](#)

(a) Model predictions by day of year



(b) Estimated chl-a between models



(c) Estimated smoother for continuous year between models

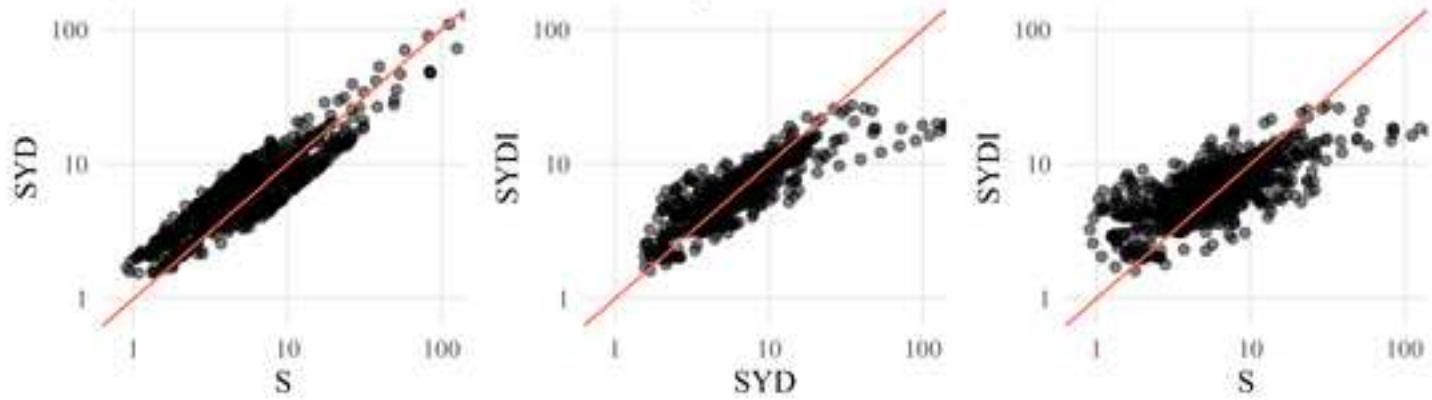


Figure 4

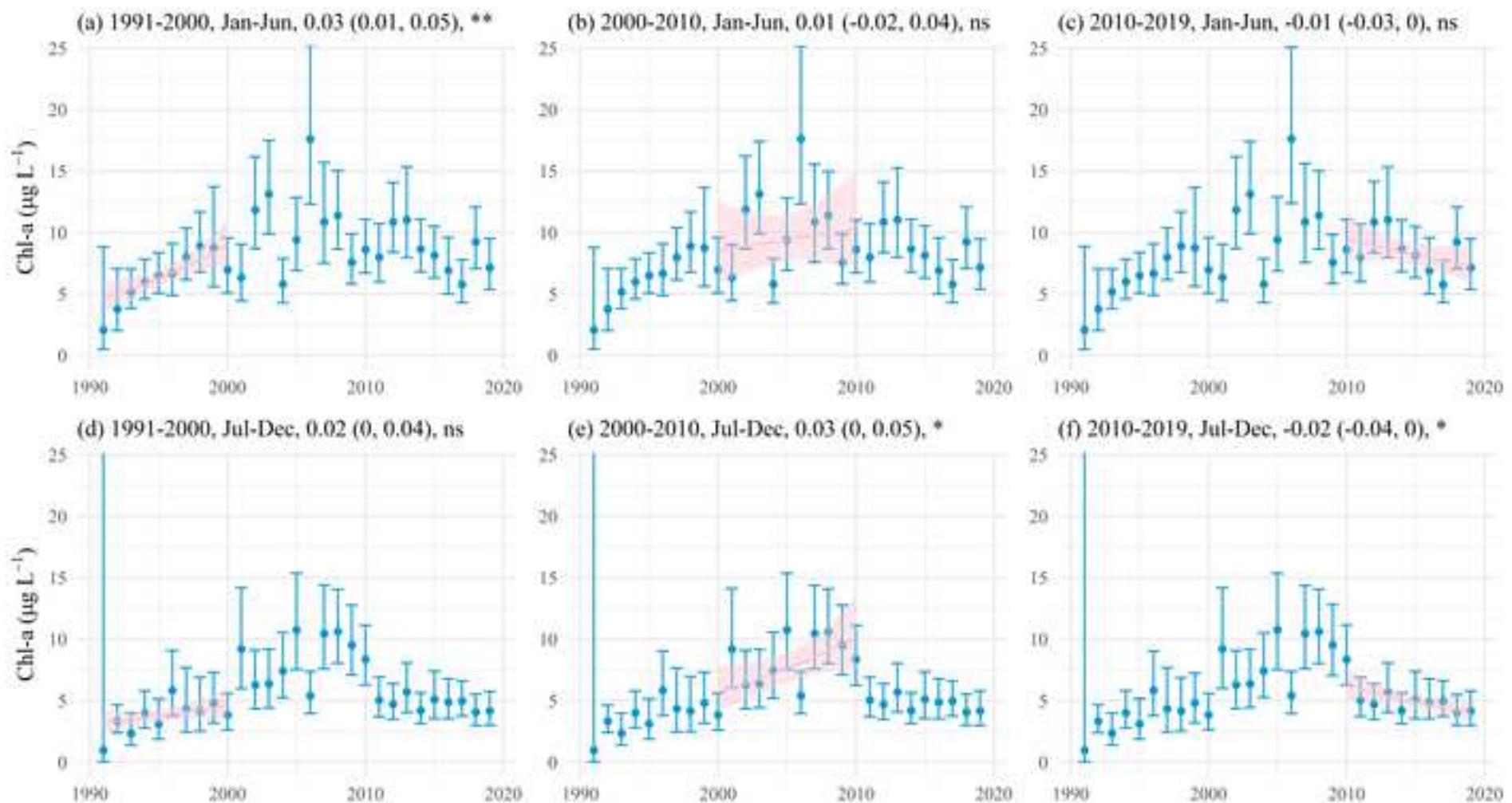
[Click here to access/download/](#)[Figure;Fig4.tiff](#)

Figure 5

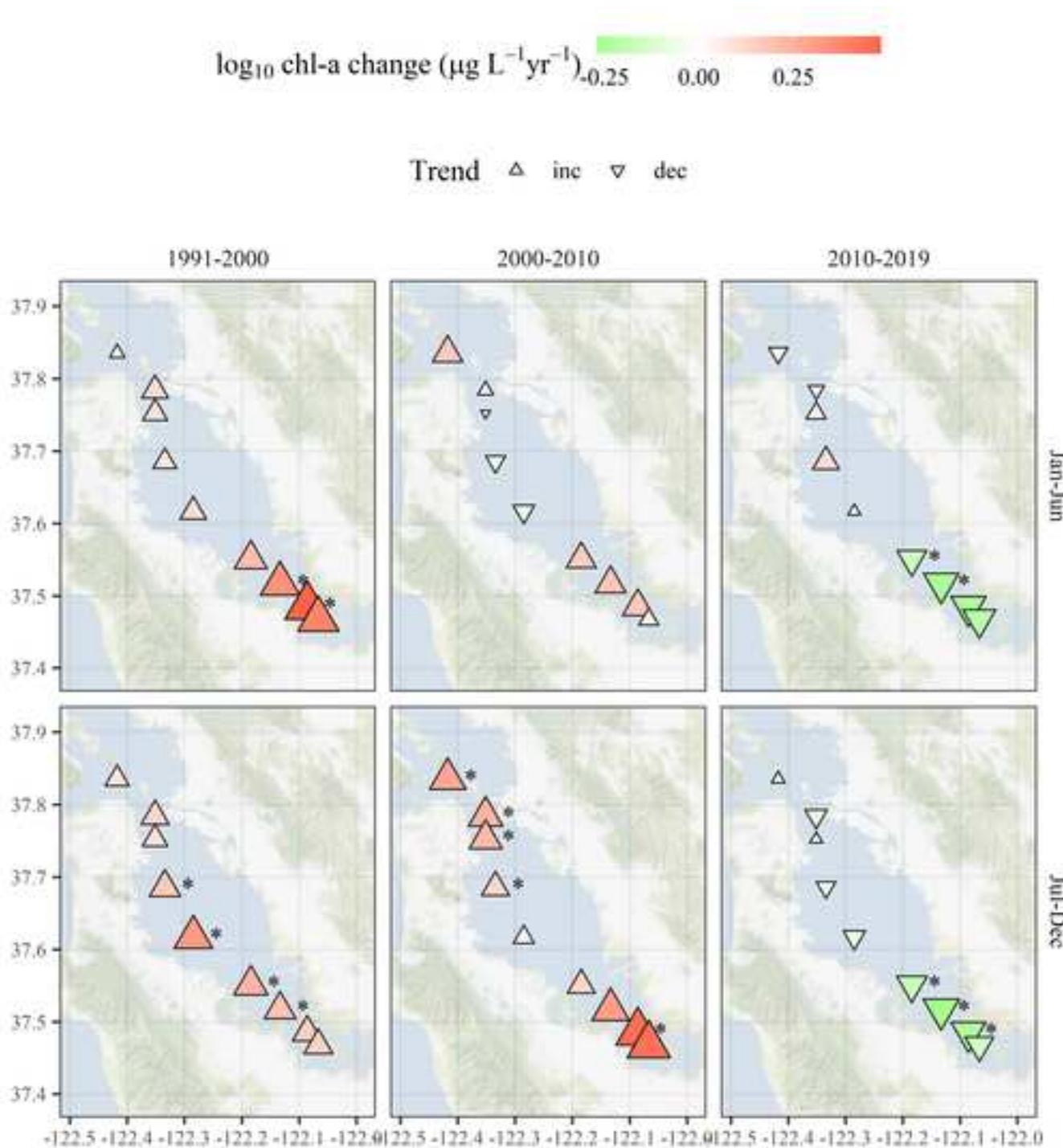
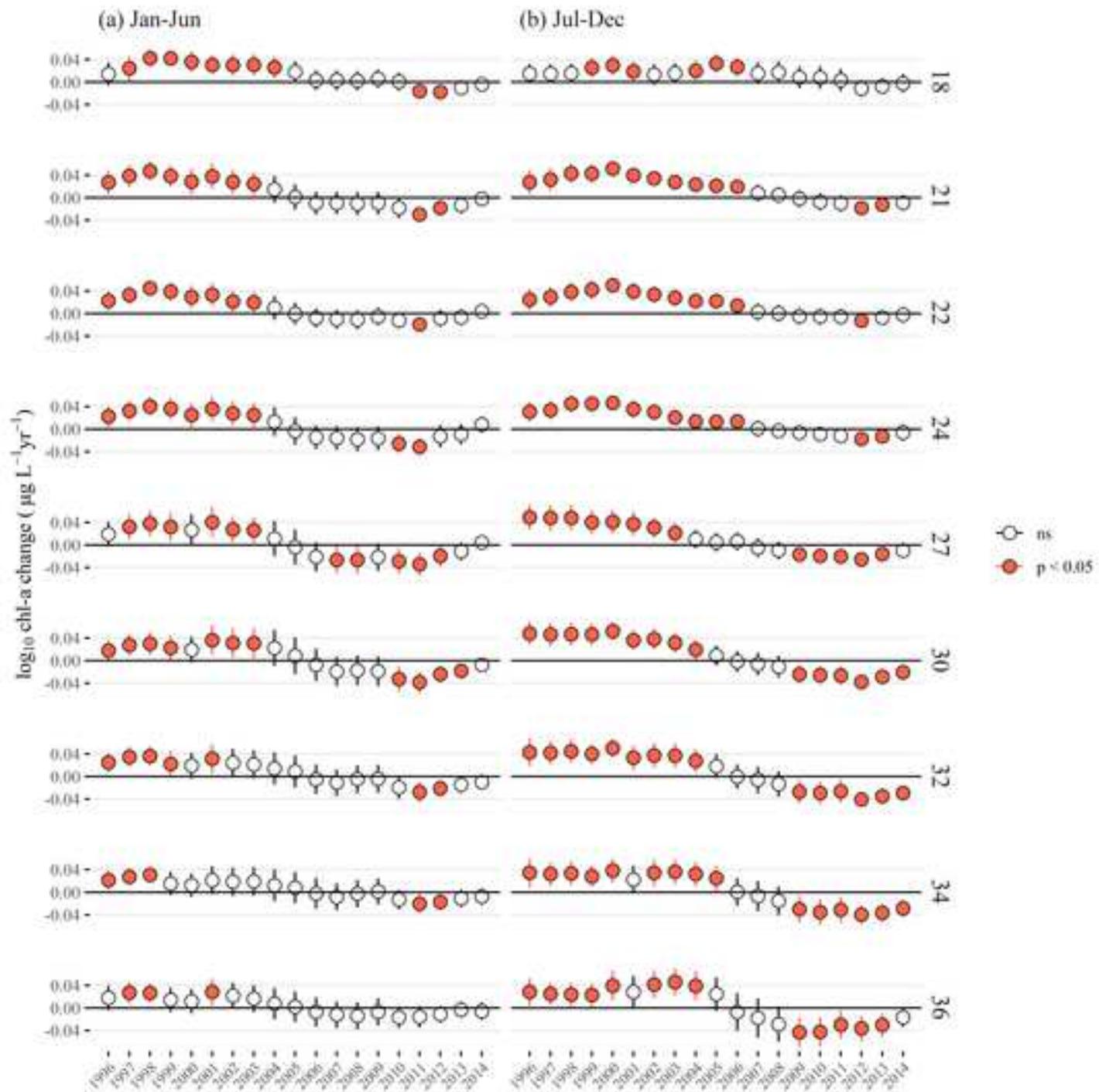
[Click here to access/download/](#)[Figure;Fig5.tif](#)

Figure 6

[Click here to access/download/](#)[Figure;Fig6.tif](#)


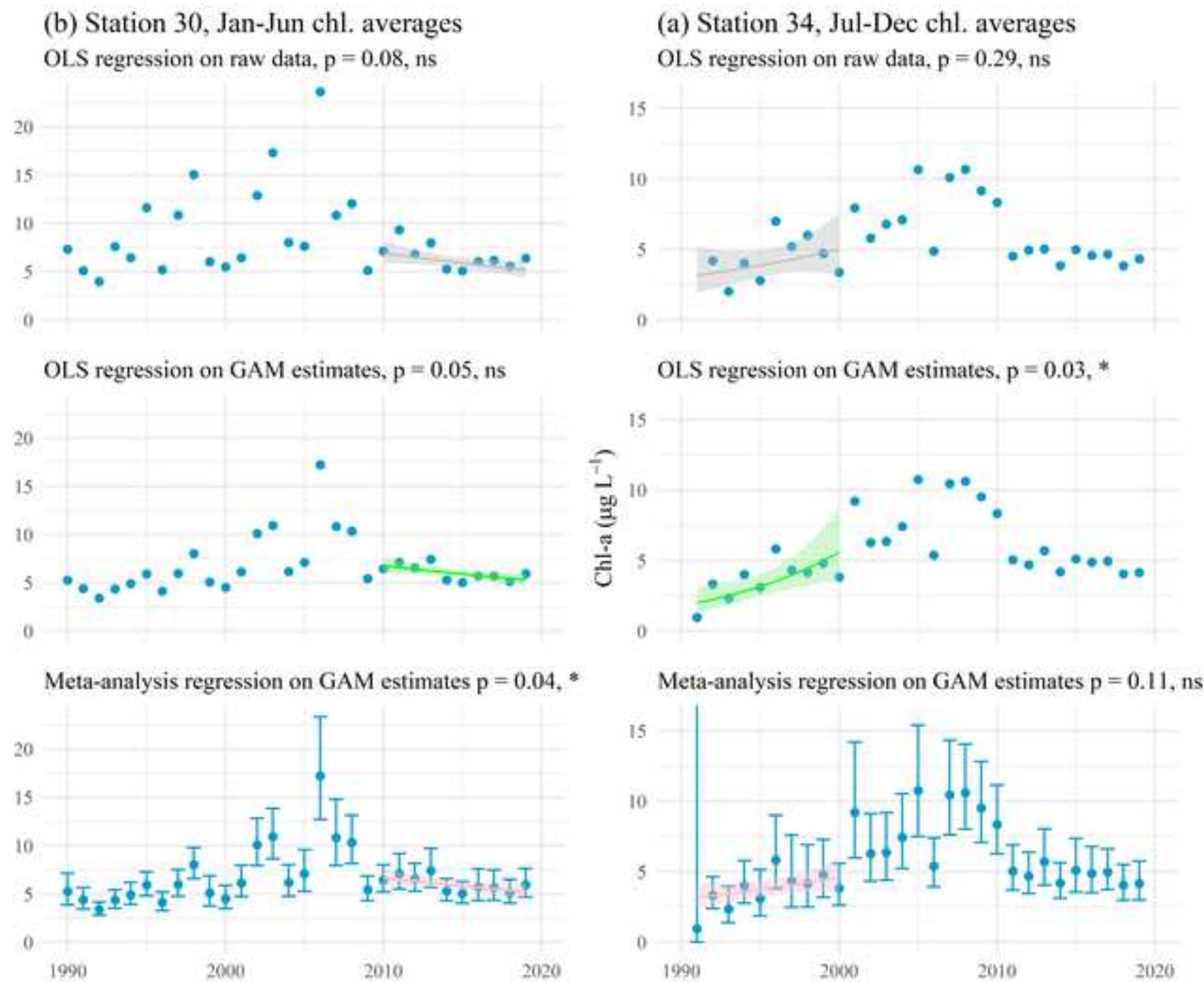
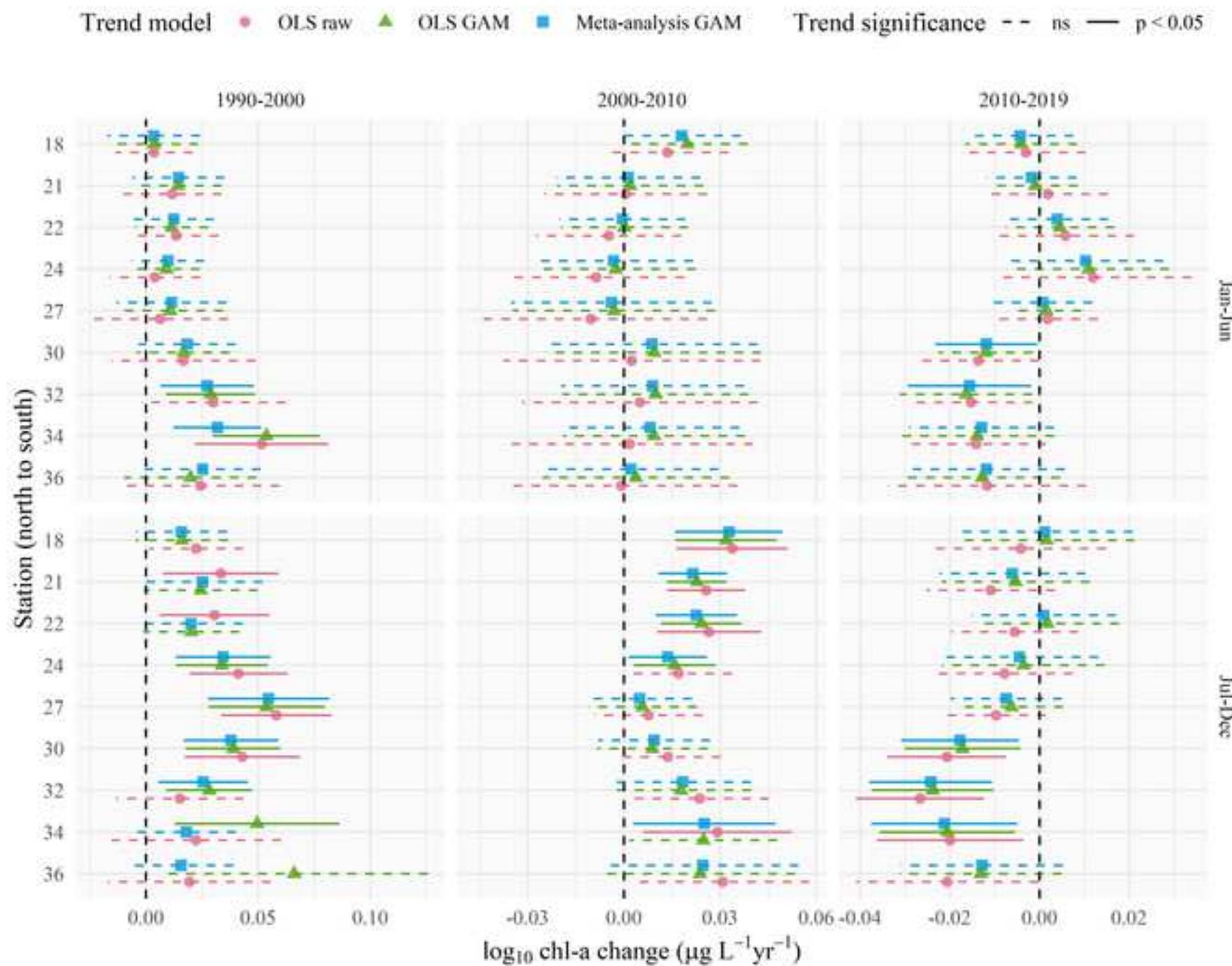


Figure 8

[Click here to access/download/](#)[Figure; Fig8.tif](#)



Click here to access/download

Supplementary material for on-line publication only
supplement.docx

Marcus W. Beck: Conceptualization, Data Curation, Formal analysis, Methodology, Software, Writing – Original Draft, **Perry de Valpine:** Conceptualization, Methodology, Software, Writing – Original Draft, **Rebecca Murphy:** Conceptualization, Writing – Review & Editing, **Ian Wren:** Conceptualization, Writing – Review & Editing, **Ariella Chelsky:** Conceptualization, Writing – Review & Editing, **Melissa Foley:** Conceptualization, Project administration, Writing – Review & Editing, **David B. Senn:** Conceptualization, Data Curation, Project administration, Writing – Review & Editing

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: