



## Advancing cyanobacteria biomass estimation from hyperspectral observations: Demonstrations with HICO and PRISMA imagery

Ryan E. O'Shea<sup>a,b</sup>, Nima Pahlevan<sup>a,b,\*</sup>, Brandon Smith<sup>a,b</sup>, Mariano Bresciani<sup>c</sup>, Todd Egerton<sup>d</sup>, Claudia Giardino<sup>c</sup>, Lin Li<sup>e</sup>, Tim Moore<sup>f</sup>, Antonio Ruiz-Verdu<sup>g</sup>, Steve Ruberg<sup>h</sup>, Stefan G. H. Simis<sup>i</sup>, Richard Stumpf<sup>j</sup>, Diana Vaičiūtė<sup>k</sup>

<sup>a</sup> NASA Goddard Space Flight Center, Greenbelt, MD, USA

<sup>b</sup> Science Systems and Applications Inc. (SSAI), Lanham, MD, USA

<sup>c</sup> National Research Council of Italy (CNR), Institute for Electromagnetic Sensing of the Environment (IREA), Milan, Italy

<sup>d</sup> Virginia Department of Health, Norfolk, VA, USA

<sup>e</sup> Department of Earth Sciences, Indiana University-Purdue University, IN, USA

<sup>f</sup> Harbor Branch Oceanographic Institute, Florida Atlantic University, FL, USA

<sup>g</sup> Laboratory for Earth Observation, University of Valencia, Valencia, Spain

<sup>h</sup> Great Lakes Environmental Research Laboratory, National Oceanic and Atmospheric Administration, Ann Arbor, MI, USA

<sup>i</sup> Plymouth Marine Laboratory, Plymouth, UK

<sup>j</sup> National Oceanic and Atmospheric Administration (NOAA), National Center for Coastal Science Studies, Silver Spring, MD, USA

<sup>k</sup> Marine Research Institute, Klaipėda University, Klaipėda, Lithuania

### ARTICLE INFO

Edited by: Menghua Wang

#### Keywords:

Cyanobacteria  
Phycocyanin  
Machine learning  
Mixture density network  
Aquatic remote sensing  
cyanoHABs  
HICO  
PRISMA

### ABSTRACT

Retrieval of the phycocyanin concentration (PC), a characteristic pigment of, and proxy for, cyanobacteria biomass, from hyperspectral satellite remote sensing measurements is challenging due to uncertainties in the remote sensing reflectance ( $\Delta R_{rs}$ ) resulting from atmospheric correction and instrument radiometric noise. Although several individual algorithms have been proven to capture local variations in cyanobacteria biomass in specific regions, their performance has not been assessed on hyperspectral images from satellite sensors. Our work leverages a machine-learning model, Mixture Density Networks (MDNs), trained on a large ( $N = 939$ ) dataset of collocated *in situ* chlorophyll-*a* concentrations (Chl<sub>a</sub>), PCs, and remote sensing reflectance ( $R_{rs}$ ) measurements to estimate PC from all relevant spectral bands. The performance of the developed model is demonstrated *via* PC maps produced from select images of the Hyperspectral Imager for the Coastal Ocean (HICO) and Italian Space Agency's PRecursorore IperSpettrale della Missione Applicativa (PRISMA) using a matchup dataset. As input to the MDN, we incorporate a combination of widely used band ratios (BRs) and line heights (LHs) taken from existing multispectral algorithms, that have been proven for both Chl<sub>a</sub> and PC estimation, as well as novel BRs and LHs to increase the overall cyanobacteria biomass estimation accuracy and reduce the sensitivity to  $\Delta R_{rs}$ . When trained on a random half of the dataset, the MDN achieves uncertainties of 44.3%, which is less than half of the uncertainties of all viable *optimized* multispectral PC algorithms. The MDN is notably better than multispectral algorithms at preventing overestimation on low ( $<10 \text{ mg m}^{-3}$ ) PC. Visibly, HICO and PRISMA PC maps show the wider dynamic range that can be represented by the MDN. The available *in situ* and satellite-derived  $R_{rs}$  matchups and measured *in situ* PC demonstrate the robustness of the MDN for estimating low ( $<10 \text{ mg m}^{-3}$ ) PC and the reduced impact of  $\Delta R_{rs}$  on medium-to-high *in situ* PC ( $>10 \text{ mg m}^{-3}$ ). According to our extensive assessments, the developed model is anticipated to enable practical PC products from PRISMA and HICO, therefore the model is promising for planned hyperspectral missions, such as the Plankton Aerosol and Cloud Ecosystem (PACE). This advancement will enhance the complementary roles of hyperspectral

**Abbreviations:**  $\alpha$ , Mixing coefficient;  $\lambda$ , Wavelength;  $\mu$ , Mean;  $\sigma$ , Standard deviation; BR( $\lambda_1, \lambda_2$ ), Band ratio (wavelength 1, wavelength 2); Chl<sub>a</sub>, Chl<sub>b</sub>, Chl<sub>c</sub>, Concentrations of Chlorophyll-*a*, Chlorophyll-*b*, Chlorophyll-*c*; CDOM, Colored dissolved organic matter; cHAB, Cyanobacteria harmful algal bloom; LH( $\lambda_1, \lambda_2, \lambda_3$ ), Line Height (or baseline algorithm) centered on  $\lambda_2$ ; MDN, Mixture density network; NIR, Near-infrared; MA, Multispectral algorithm; PC, Phycocyanin concentration;  $R_{rs}(\lambda)$ , Remote-sensing reflectance;  $\Delta R_{rs}(\lambda)$ , Uncertainty in the satellite derived remote sensing reflectance.

\* Corresponding author at: NASA Goddard Space Flight Center, Greenbelt, MD, USA.

E-mail address: [nima.pahlevan@nasa.gov](mailto:nima.pahlevan@nasa.gov) (N. Pahlevan).

<https://doi.org/10.1016/j.rse.2021.112693>

Received 25 May 2021; Received in revised form 1 September 2021; Accepted 7 September 2021

Available online 16 September 2021

0034-4257/© 2021 The Authors.

Published by Elsevier Inc.

This is an open access article under the CC BY-NC-ND license

(<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

radiometry from satellite and low-altitude platforms for quantifying and monitoring cyanobacteria harmful algal blooms at both large and local spatial scales.

## 1. Introduction

Cyanobacteria can produce a variety of toxins that pose health risks, and even mortality, in wildlife, livestock, pets and humans through ingestion (the most common route of exposure), contact with skin, or inhalation (Health Canada, 2020; U.S. EPA, 2019). Although there is a wide range of toxins produced by cyanobacteria, some of the most common cyanobacterial toxins fall into the category of hepatotoxins (which harm the liver) and neurotoxins (which harm the nervous system). The toxins produce health effects in humans, ranging from mild effects, such as irritation of the eyes and ears, dermatitis, diarrhea, headaches, and abdominal pain, to more serious effects, such as muscle paralysis and kidney damage (World Health Organization, 2003). In extreme exposure events even death can occur, due to muscle paralysis inhibiting respiration. Overall, early warning and mapping of cyanobacteria harmful algal blooms (cHABs) is critical for water resource managers to keep constituents safe (Liu et al., 2020; Mishra et al., 2019; Schaeffer et al., 2018).

Water resource managers rely on a variety of different measurements to determine the safety of recreational and drinking waters. These measurements include concentrations of chlorophyll-*a* (Chl*a*) and phycocyanin (PC), cyanobacteria cell density, and microcystin (and other cyanotoxin) concentrations (U.S. EPA, 2019; World Health Organization, 2003). The United States Environmental Protection Agency (EPA) recommends that primary contact (swimming) advisories should be issued when microcystin (hepatotoxin) and cylindrospermopsin (neurotoxin) measurements exceed 8 and 15 mg m<sup>-3</sup>, respectively (US EPA, 2019). In recreational waters, the World Health Organization (WHO) recommends action be taken at a Chl*a* level of 10 mg m<sup>-3</sup> (where microcystin are generally somewhere between 2 and 4 mg m<sup>-3</sup>, assuming cyanobacteria are dominant in the waters) or at a cyanobacteria cell density of 20,000 cells mL<sup>-1</sup> (World Health Organization, 2003). Unfortunately, these metrics either assume that cyanobacteria are dominant or require *in situ* measurement and specialization to perform (Jin et al., 2018; World Health Organization, 2003).

While the *in situ* assessment techniques are useful for identifying risk at individual timepoints for specific locations, they are inadequate for early warning and monitoring of cHABs in multiple water bodies on national scales (Hunter et al., 2009), which is critical for remedial actions (Binding et al., 2021; Clark et al., 2017; Schaeffer et al., 2018). As a complement to *in situ* measurements, water resource managers can turn to optical sensors to rapidly map cyanobacteria biomass. In freshwater ecosystems, particularly those that are most eutrophic, the pigment phycocyanin is produced by cyanobacteria (in contrast to Chl*a*, which is ubiquitous to harmful and non-harmful algal taxa). Optical sensors combined with spectral algorithms can leverage the spectral characteristics of phycocyanin to remotely determine if cyanobacteria are present in potentially harmful concentrations (Dekker, 1993) making it an indicator pigment (Stumpf et al., 2016).

Over time, algorithms for retrieving PC have evolved from simple empirical band ratios and baseline algorithms to much more complex, physics-based, semi- and quasi-analytical algorithms. Most of the existing algorithms, which have been successfully trained and developed for specific regions, are multispectral and utilize one or a few band ratios and line heights even though many more bands are available (Dekker, 1993; Mishra et al., 2009; Schalles and Yacobi, 2000; Simis et al., 2005; Hunter et al., 2010; Li and Song, 2017; Stumpf et al., 2016). Further, only a limited number of cyanobacteria algorithms have been tested for *high-altitude* (aircraft) *hyperspectral* remote sensing (Pyo et al., 2020; Kudela et al., 2015; Hunter et al., 2009; Li and Song, 2017), where atmospheric effects are present. Our research aims are to advance the use

of *hyperspectral satellite* imagery for global PC retrieval by (1) developing a robust machine-learning algorithm sensitive to a wide range of PC and applicable across geographic regions, (2) demonstrating the PC mapping performance from past and current spaceborne imaging spectrometry, and (3) identifying the impact of uncertainties in the remote sensing reflectance ( $\Delta R_{rs}$ ) on PC retrieval. This research is carried out to further the use of current demonstration hyperspectral sensors (e.g., PRecursor IperSpettrale della Missione Applicativa (PRISMA)), in preparation for future hyperspectral satellite missions (e.g., the plankton, aerosol, cloud, ocean, and ecosystem (PACE) mission, the Fluorescence Explorer (FLEX)), and to inform pre-formulation studies (e.g., NASA's Surface Biology and Geology designated observable; European Space Agency's Copernicus Hyperspectral Imaging Mission, CHIME).

In this work, we present a Mixture Density Network (MDN) architecture that uses line heights and band ratios to accurately estimate PC, in the presence of  $\Delta R_{rs}$ . The MDN architecture also uses the simultaneous estimation of Chl*a* to force the model to learn the covariance with PC and its impact on the retrieval of PC. In addition, the MDN is expected to improve PC estimation by learning the nonlinear association of the line-height and band-ratio features with spectral remote sensing reflectance ( $R_{rs}$ ) to overcome spectral variability due to other pigments (e.g., chlorophyll *b* (Chl*b*)), which have been shown to inhibit PC estimation using existing multispectral algorithms (MAs) (Simis et al., 2007; Ruiz-Verdú et al., 2008). We demonstrate the accuracy benefits of our architecture for PC retrieval from regions included within the training set (which represents a wide range of aquatic regions) by splitting the overall dataset into training and testing sets using a 50/50 split. We further evaluate the accuracy of the architecture for PC retrieval from a wide range of conditions through comparison to individual MAs, by training new MDNs with the same architecture, but leaving individual regions out of the training set and then testing on those regions (*i.e.*, leave-one-out or round robin testing). The demonstration MDN, an MDN with the same architecture which has been trained on the *entire* dataset, is evaluated on images collected by the heritage Hyperspectral Imager for the Coastal Ocean (HICO) and Italian Space Agency's demonstration mission (PRecursor IperSpettrale della Missione Applicativa; PRISMA), through comparison against historic observations and *in situ* measurements. Finally, to demonstrate the stability of our architecture to  $\Delta R_{rs}$  typical of these hyperspectral sensors, we compare the stability of our model against an MDN trained solely on  $R_{rs}$ . We evaluate the stability by analyzing PC retrievals made from the two models on collocated *in situ* and atmospherically corrected radiometric measurements. The benefits of this MDN architecture for monitoring and mapping cHABs from hyperspectral satellite sensors and its relevance to pre-formulation studies is discussed, as is the future research required to develop an operational algorithm.

## 2. Background

PC is a pigment that is only present in high concentrations in cyanobacteria, and so its spectral signature can be used to differentiate cyanobacteria biomass from non-cyanobacterial algal biomass using optical remote sensing measurements. Although PC does not directly correspond to toxin levels, as the same cyanobacteria biomass could correspond to different levels of toxin concentrations, PC has been shown to be more correlated with certain toxins (e.g., microcystin) than other optical proxies, including Chl*a* (Rinta-Kanto et al., 2009; Francy et al., 2016). This provides water managers with a metric that is much more specific to toxin risk and allows for more targeted response and advisories if necessary. There is however spectral overlap between phycocyanin, which has strong absorption around 620 nm and

fluorescence around 650 nm (Dekker, 1993; Becker et al., 2002; Schalles and Yacobi, 2000), and other pigments. Chla, the most ubiquitous pigment in algae, exhibits absorption near 620 nm, strong absorption around 670 nm, and fluorescence at 683 nm (Ficek et al., 2004; Sathyendranath et al., 1987). The absorption of pigments Chlb and chlorophylls c1 and c2 (Chlc), at the 650 nm band and on either side of the 620 nm band respectively, also overlap with the spectral features of phycocyanin (Ficek et al., 2004; Sathyendranath et al., 1987; Simis et al., 2007). In nearshore coastal estuaries, mixed phytoplankton communities typically have both phycocyanin and a variety of chlorophyll pigments in varying ratios. The absorption of these pigments together with the scattering properties of cyanobacteria cells as well as those of other optically relevant components of the water column (e.g., inorganic particles) regulates the shape and magnitude of  $R_{rs}$ , defined as the ratio of water-leaving radiance to total downwelling irradiance just above the water surface (Mobley, 1999). The interplay between the absorption and scattering properties of water constituents limits the PC retrieval accuracy.

The original PC retrieval algorithms leveraged empirical band ratios and baseline algorithms (Table 1, PC Algorithms) (Ogashawara, 2020; Li, 2020), developed from direct observation of PC and  $R_{rs}$  measurements. One of the initial PC retrieval algorithms was a semi-empirical baseline algorithm that measured the impact of the absorption by phycocyanin (at 624 nm) on the subsurface irradiance reflectance relative to the 600 and 648 nm bands (Dekker, 1993). Another empirical algorithm employing band ratios between 650 nm and 625 nm was found to have a better fit than the previously developed baseline algorithm, though variability in its estimates did increase at higher PC (Schalles and Yacobi, 2000). Unfortunately, these simple algorithms did not account for the substantial impact of Chla in the absorption band, which can have deleterious effects on the PC estimation, particularly when the ratio of PC to Chla is low (Simis et al., 2005; Ogashawara et al., 2013; Simis et al., 2007; Mishra and Mishra, 2014; Mishra et al., 2013).

More complex semi-analytical algorithms, which attempt to correct for the spectral variability due to other optically active constituents through their absorption and backscattering properties while still leveraging empirical components, have also been developed for PC retrieval. One such algorithm accounted for the absorption from water and Chla at 620 nm while solving for PC, by using band ratios between the near-infrared (NIR) and Chla absorption bands (Simis et al., 2005; Lyu et al., 2013). This semi-analytical model assumed that the absorption in the NIR is due solely to water, which may lead to inaccurate estimates either at higher PC (e.g., during CHABs) or in the presence of inorganics (Babin and Stramski, 2004; Doxaran et al., 2009), that the absorption due to CDOM is negligible, and that the overlapping pigment signatures (e.g., Chlc) impacts on PC retrieval were negligible. The Simis et al. (Simis et al., 2005) algorithm performed best as compared to previous baseline and band ratio algorithms (Ruiz-Verdú et al., 2008). Further review showed that overestimation of PC by the Simis et al. (Simis et al., 2005) model typically occurs when other pigments (e.g., Chlb, Chlc) are present (Simis et al., 2007). In a similar vein as the semi-empirical algorithm, a band-ratio algorithm between the NIR (700 nm) and yellow-orange (600 nm) was developed, and found to further increase the tolerance of the model to variations in Chla in a laboratory based study (Mishra et al., 2009). Another semi-empirical algorithm also increased PC estimation accuracy by correcting for the impact of Chla and PC absorption at 620/665 nm respectively (Ogashawara and Li, 2019). A four-band semi analytical model separated the absorption due to PC (in the 620 nm band) from that due to other phytoplankton pigments, CDOM, and water, by adding three bands that capture variability due to these components (560 nm, 709 nm, and 754 nm) and solving for a single algorithm parameter (Liu et al., 2017). Overall semi-analytical models greatly improved upon empirical models by correcting for the impact of other optically relevant constituents, increasing the transferability of these models.

A quasi-analytical algorithm, which leverages the basics of radiative

**Table 1**

Band ratios used in multispectral PC algorithms and machine learning techniques.

Band ratio/Line Height	Source	Justification
Multispectral $(R(0-,600) + R(0-,648))/2 - R(0-,624)$ $R_{rs}(650)/R_{rs}(625)$ $R(0-,709)/R(0-,665)$	Dekker (1993) Schalles and Yacobi (2000) Simis et al. (2005)	Isolate phycocyanin absorption impact on $R_{rs}$ at 624 nm Phycocyanin absorption at 620 nm Chla absorption at 665 nm
$R(0-,709)/R(0-,620)$ $R_{rs}(700)/R_{rs}(600)$ $R_{rs}(725) * (1/R_{rs}(615) - 1/R_{rs}(600))$ LH(665,681,709)	Mishra et al. (2009) Hunter et al. (2010) Cyanobacteria Index (Wynne et al., 2010) Stumpf et al. (2016), Lunetta et al. (2015), Matthews and Odermatt (2015)	
LH(620,665,681)		
$R_{rs}(724) * (1/R_{rs}(629) - 1/R_{rs}(659))$ LH(654,714,754)	Mishra and Mishra (2014) Kudela et al. (2015)	
Machine learning $R_{rs}(710)/R_{rs}(665)$ $R_{rs}(715)/R_{rs}(665)$ $R_{rs}(715)/R_{rs}(690)$ $R_{rs}(710)/R_{rs}(660)$ $R_{rs}(710)/R_{rs}(690)$ $R_{rs}(715)/R_{rs}(670)$ $R_{rs}(710)/R_{rs}(620)$	Sun et al. (2012)	Band ratios that have the highest coefficient of variation with phycocyanin concentration
Added to our model (Section 4.1) (665,709,754)	Maximum chlorophyll index (665)	
(680,709,754)	Maximum chlorophyll index (680)	
(443,555,670) $(R_{rs}(709) - R_{rs}(665))/ (R_{rs}(709) + R_{rs}(665))$	Color Index NDCI	
Max_location( $R_{rs}(550-600)$ )	Schalles and Yacobi (2000)	Correlated with Chla
Max_location( $R_{rs}(694-716)$ )	Schalles and Yacobi (2000)	Correlated with PC
LH(560,620,665)		LH around phycocyanin absorption
LH(665,673,681)		LH around Chla absorption
LH(690,709,720)		LH around NIR peak associated with particle scattering
LH(620,650,670)		LH around phycocyanin fluorescence
LH(640,650,660)		LH around phycocyanin fluorescence
LH(613,620,627)		LH around phycocyanin absorption

Italic text indicates ratios used in the demonstration MDN (Section 4.1). Acronyms: the reflectance just below the surface ( $R(0^-, \lambda)$ ), the remote sensing reflectance at a wavelength ( $R_{rs}(\lambda)$ ).

transfer to solve for phycocyanin induced absorption, successfully estimated PC (from its absorption) in turbid cyanobacteria-dominated waters by accounting for the absorption by phycocyanin at 665 nm, the potential for spectrally dependent backscattering of particulates, and the absorption by CDOM at 620 nm (Mishra et al., 2013). Another quasi-analytical algorithm further increased the accuracy for low PC retrieval by separating absorption due to phycocyanin from the absorption of other pigments and CDOM (Li et al., 2015). While a range of other algorithms (e.g., Hunter et al., 2010; Li et al., 2015; Liu et al., 2017;

Li et al., 2012; Le et al., 2011) including those relying on Rayleigh-corrected reflectance spectra (Wynne et al., 2010; Binding et al., 2021; Binding et al., 2019; Matthews et al., 2012) exist, this overview of MAs demonstrates the range of wavelengths, techniques, and assumptions used for PC retrieval.

Machine learning techniques can overcome some of the limitations of individual semi-empirical, semi-analytical, and quasi-analytical PC retrieval algorithms by combining, and benefiting from, the information available from multiple optical features to estimate PC. For example, a support vector regression (SVR) trained using the seven band ratios with the highest correlation to PC, from a 92 paired  $R_{rs}$  - PC set taken from three inland turbid lakes in China, was able to achieve higher accuracy (mean absolute percent difference (MAPD) = 29.5%, and root mean square difference (RMSD) = 28.4 mg m<sup>-3</sup>) on a validation set from the same region as compared to any individual MA (e.g., the semi-analytical algorithm (Simis et al., 2005; Simis et al., 2007) had MAPD 47.5% and RMSD 46.5 mg m<sup>-3</sup>) (Sun et al., 2012). Although SVRs work well even when low amounts of data are available for training, more complex artificial neural networks can also be trained on these small datasets through input dimensionality reduction. For example, partial least squares (PLS) can be used to reduce the input dimensionality, preserving the most influential input, to an artificial neural network (forming a PLS-ANN) to estimate PC (Song et al., 2014; Song et al., 2012). The input consisted of spectral  $R_{rs}$ , derivatives of the  $R_{rs}$ , and band ratios of the  $R_{rs}$ . The PLS-ANN reduced uncertainties beneath that of traditional three band models, in data taken from two separate regions. A final machine learning method leveraged two sets of stacked autoencoders feeding artificial neural networks and support vector regressions to estimate PC from  $R_{rs}$ ; the first set atmospherically corrected the input  $R_{rs}$  and the second set calculated both PC and Chla (Pyo et al., 2020). Although these non-linear regression techniques perform well in specific regions, and outperform heritage algorithms, their sensitivity to uncertainties in the  $R_{rs}$  ( $\Delta R_{rs}$ ) and applicability to wide-area hyperspectral satellite images have not been demonstrated.

These uncertainties in the  $R_{rs}$ , which can result from both instrument noise (Moses et al., 2012) and the atmospheric correction process (Ibrahim et al., 2018), can propagate to errors in the retrieved satellite products. Modeling efforts have shown that instrument noise can result in errors as high as 50–80% in the estimated constituent concentration (e.g., Chla), depending on the aquatic conditions and algorithms used, from hyperspectral satellite sensors such as HICO (Moses et al., 2012). This error is a minimum bound, as imperfect correction of the atmospheric effects typically result in further inaccuracies in the retrieved products from hyperspectral satellite sensors (Moses et al., 2012). For multispectral sensors, the combined instrument noise and atmospheric correction was found to yield uncertainties on the order of 25–70% in certain products (Pahlevan et al., 2021a). Depending on the assumptions made for a given atmospheric correction process, accurate retrieval of blue-green wavelengths can be particularly challenging. For instance, in extremely turbid and/or eutrophic waters, classic atmospheric correction methods (Mobley et al., 2016) may fail to fully account for water-leaving radiance in the NIR, leading to particularly high errors in  $R_{rs}$  and in extreme cases can even result in (unphysical) negative  $R_{rs}$  (Ibrahim et al., 2018). Therefore, it is useful to study the impact of  $\Delta R_{rs}$  on PC retrieval algorithms developed on *in situ*  $R_{rs}$ , to inform on their limitations for PC retrieval when applied to satellite  $R_{rs}$ .

Mixture density networks (MDNs) (Bishop, 1994) are a class of machine learning algorithms particularly suited for non-unique inverse problems where limited (i.e., a low number of) training data are available. MDNs have been proven successful for a range of aquatic remote sensing tasks requiring non-unique inversion including: the estimation of Chla from Landsat-8 (Smith et al., 2021), the estimation of Chla and phytoplankton absorption ( $a_{ph}$ ) from hyperspectral data (Pahlevan et al., 2021b), the retrieval of Chla from multispectral imagers using multispectral sensors (Pahlevan et al., 2020), and the retrieval of particulate backscattering (Balasubramanian et al., 2020). The efficacy of

MDNs in optically complex coastal and inland waters demonstrates their applicability to non-unique inverse remote sensing problems.

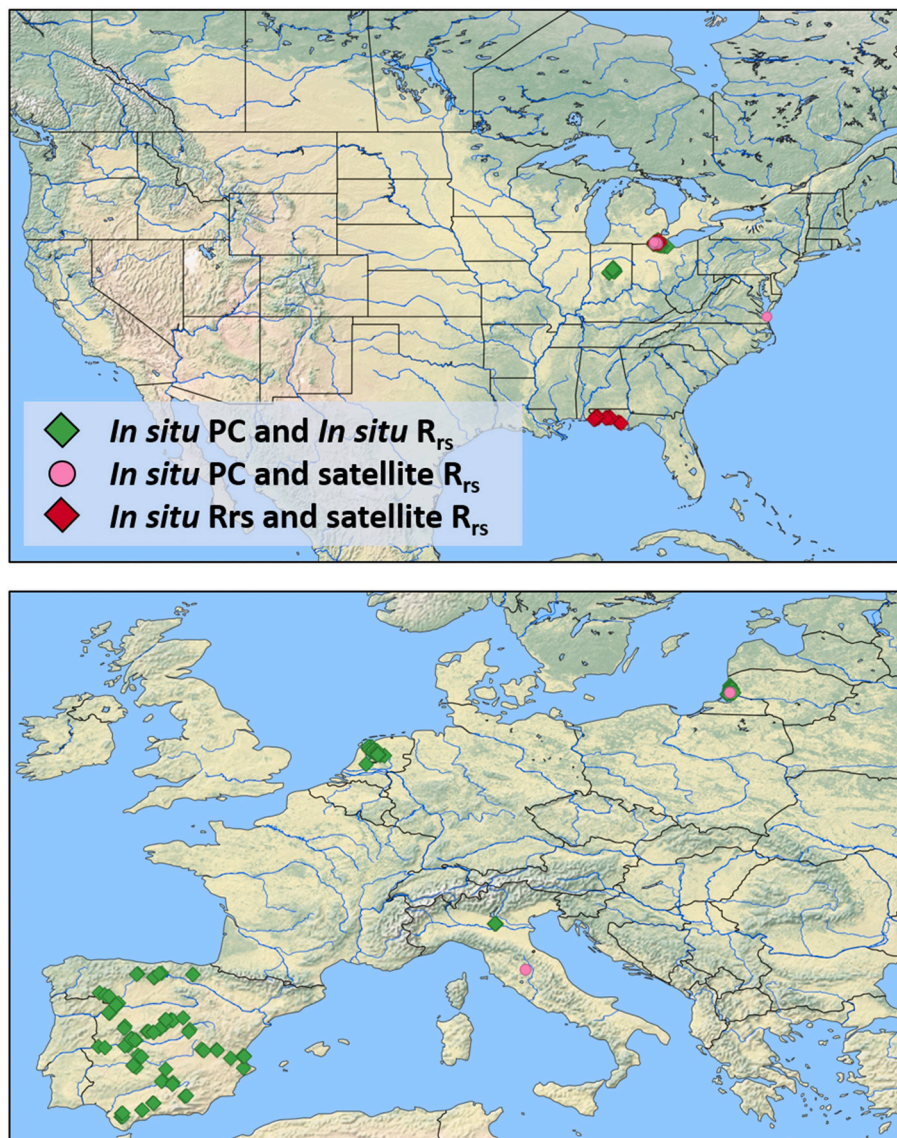
### 3. Datasets

#### 3.1. Model development data

The entire dataset ( $N = 939$ ) consists of collocated *in situ*  $R_{rs}$ , Chla, and PC measurements, the largest dataset of this kind constructed to date. The regional makeup of the *in situ* measurements is: the western basin of Lake Erie, small and large turbid lakes in the Netherlands (Simis et al., 2005), small lakes in Indiana (Li et al., 2015), lakes and reservoirs in Spain (Ruiz-Verdú et al., 2008), the Curonian Lagoon bordering Lithuania and Russia nearby the Baltic Sea, the Superior and Inferior lakes of Mantua in Italy, and Hartbeespoort, Theewaterskloof and Loskop dams in South Africa (Matthews, 2020) (Fig. 1, Table 2). Since the *in situ* dataset comes from a range of geographic regions and morphologically, hydrologically, and optically distinct waterbodies, it spans a variety of PC, Chla, and PC:Chla ranges (Fig. 2), which enables the MDN to train over a range of optical conditions. It is especially important to include data representing a broad range of PC:Chla, as Chla is the primary pigment in both toxic and non-toxic phytoplankton groups that impacts  $R_{rs}$  and therefore the recovered PC (Simis et al., 2005; Randolph et al., 2008). The dataset was pruned to only include stations with both Chla and PC measurements between 0.0 and 1000 mg m<sup>-3</sup> and  $R_{rs}$  measurements greater than 0 within the desired wavelength range (501–724 nm). The lower bound (501 nm) was chosen to avoid commonly high  $\Delta R_{rs}$  within the blue region whereas the upper bound (724 nm) was set to include the maximum number of matchups (Section 3.3). PC ranging from 0 to 10<sup>-1</sup> mg m<sup>-3</sup> were set to 10<sup>-1</sup> mg m<sup>-3</sup> because this range approaches the detection limit of *in situ* measurement, this range is similar from a water quality management standpoint and prevents the MDN from focusing on uncertainties in a largely indistinguishable range. A significant portion of the dataset consists of low PC (with an overall median concentration of 14.1 mg m<sup>-3</sup>). Finally, it is important to note that there are relatively few measurements from extreme cHAB conditions (PC >200 mg m<sup>-3</sup>).

Although laboratory extracted phycocyanin serves as our comparison to remotely retrieved values, it is worth noting that there can be large imprecisions in extracted PC from an individual extraction technique and large inaccuracies in extracted PC between techniques. For example, the error between multiple replicates, when using the freeze-thaw method (Sarada et al., 1999), can be on the order of 9% (Song et al., 2012). Additionally, on samples with low biomass (Chla < 20 mg m<sup>-3</sup>), the coefficient of variation between multiple replicates of extracted phycocyanin using the freeze-thaw methods can be on the order of 10–15% (Horváth et al., 2013). Not only are the PC extraction techniques imprecise, but they can also be inaccurate, depending on which extraction buffer is used. If a phosphate buffer serves as the extraction buffer instead of Asolectin-CHAPS, it can underestimate the PC derived from *M. aeruginosa* by ~30% (Zimba, 2012). Although these error estimates are not derived from our dataset, they serve as a general metric for an estimated lower bound for any individual remote estimation method. For a more detailed description of the dataset, we refer readers to previously published materials (Simis et al., 2005; Ruiz-Verdú et al., 2008; Li et al., 2015).

The large dataset of *in situ*  $R_{rs}$  cover a wide range of the seven optical water types (OWTs, Fig. 3) identified in Pahlevan et al. (2021a). Overall, the  $R_{rs}$  inputs are heavily dominated by OWT-5. OWT-5 has spectral features characteristic of cyanobacteria blooms, including a reflectance trough near 620 nm, likely caused by phycocyanin absorption, and a reflectance peak near 710 nm, suggesting material accumulation near the surface masking the high NIR absorption of water. This large amplitude is potentially indicative of particles consisting of gas-vacuolate cyanobacteria with a high backscatter efficiency (Moore et al., 2019). OWT-5 also has higher Chla and PC (~50 mg m<sup>-3</sup>)



**Fig. 1.** Geographic distribution of data in the United States and Europe (10 locations from South African reservoirs are not shown). Maps downloaded from <http://www.naturalearthdata.com/>.

**Table 2**

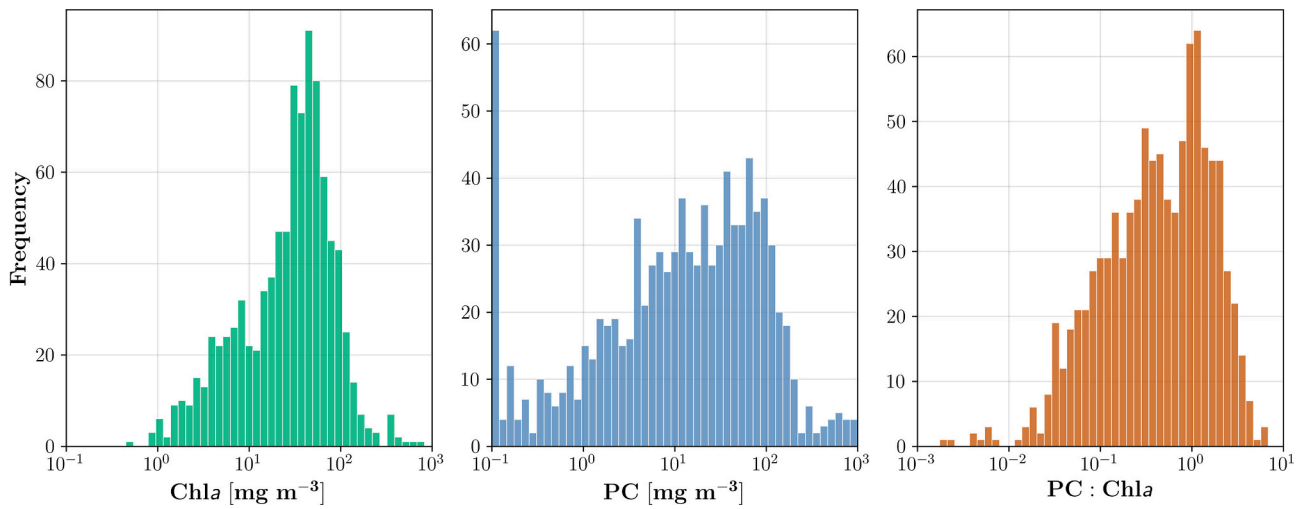
Number of samples (*N*), mean and standard deviation (SD) of PC ( $\text{mg m}^{-3}$ ), and mean and SD of Chla ( $\text{mg m}^{-3}$ ) in each region.

Region	Lake Erie	Dutch Lakes	Lakes of Indiana	Curonian Lagoon	Lakes of Italy	South African Reservoirs	Lakes of Spain
<i>N</i>	384	186	151	63	20	10	125
Mean PC	20.4	46.0	53.7	109	7.5	210.7	102.5
SD of PC	49.5	70.0	41.4	162.3	8.8	125.8	205.2
Mean Chla	34.3	46.4	52.9	69.2	32.9	106.9	54.4
SD of Chla	36.3	34.5	27.8	107.3	35.0	224.5	105.8

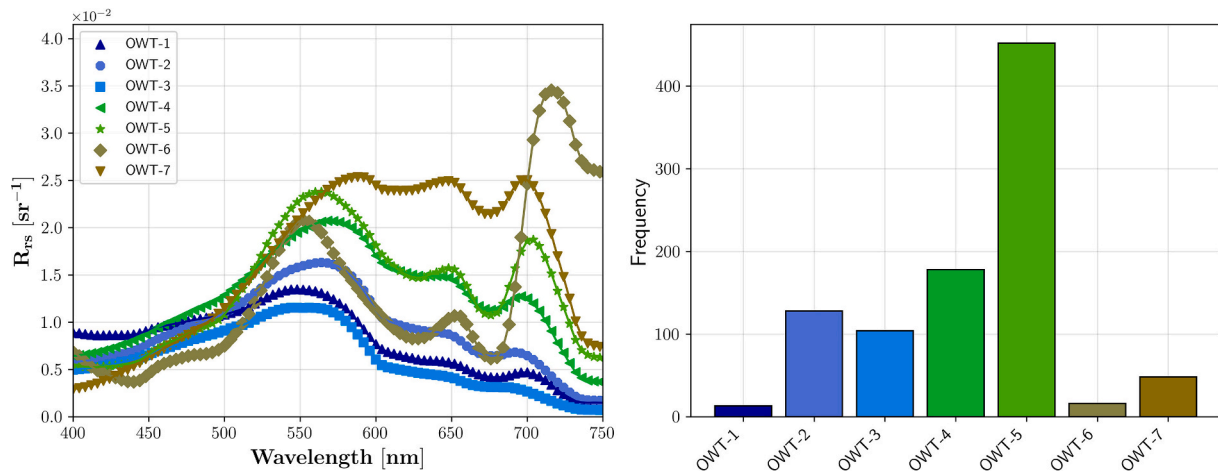
(Table 3), again characteristic of a moderate bloom. OWTs 2–4 are dominated by low PC ( $\sim < 10 \text{ mg m}^{-3}$ , Table 3), and do not have the corresponding absorption-induced dip near 620 nm (Fig. 3), though OWT-2 and OWT-4 do have a reflectance peak near 710 nm. In OWT-2, there is a slight peak in the 680 nm to 690 nm range, which may be due to Chla fluorescence (Fig. 3, left panel). There is a notably low occurrence of  $R_{rs}$  matching OWT-6, which represents intense and potentially surfacing cyanobacteria blooms (Fig. 3, left panel,  $\sim 600 \text{ mg m}^{-3}$  PC, Table 3). OWT-6 has the largest dip near 620 nm, which is characteristic of phycocyanin, and the highest NIR reflectance, which is characteristic

of extremely high cyanobacterial accumulation near the surface. OWT-7, with its relatively high and flat spectral shape in the red and NIR (Fig. 3, left panel), is representative of sediment-rich waters (Table 3). Overall, the dataset represents a broad spectrum of watercolor signals, which is useful for training general algorithms. For more information regarding the composition of the OWTs, see Table 4 in Pahlevan et al., 2021a.

The *in situ*  $R_{rs}$  are resampled to match the wavelengths and bandwidth of the prior proof-of-concept HICO (<https://oceancolor.gsfc.nasa.gov/hico/instrument/dataset-characteristics/>) and current PRISMA



**Fig. 2.** Log-scale histogram calculated for Chla, PC, and the ratio between PC and Chla (PC:Chla) within our dataset ( $N = 939$ ). The median and mean for Chla, PC, and PC:Chla are (33.15, 45.4), (14.1, 49.4), and (0.46, 0.82) respectively.



**Fig. 3.** The mean  $R_{rs}$  and frequency of each optical water type (OWT, left and right panels respectively) within our dataset. The seven OWTs are defined in [Pahlevan et al. \(2021a\)](#) from a subset of OWTs originally defined in [Spyrakos et al. \(2021\)](#). The mean Chla and PC within each OWT are included in [Table 2](#).

**Table 3**

The median (and mean in brackets) concentrations of chlorophyll-*a* (Chla), phycocyanin (PC), and phycocyanin to chlorophyll-*a* ratio (PC:Chla) for each optical water type (OWT) shown in [Fig. 3](#).

	Chla ( $\text{mg m}^{-3}$ )	PC ( $\text{mg m}^{-3}$ )	PC:Chla	$N$
OWT-1	7.4 [21.2]	0.8 [6.6]	0.11 [0.31]	13
OWT-2	13.8 [16.6]	3.9 [7.8]	0.29 [0.47]	128
OWT-3	5.6 [6.6]	0.6 [2.8]	0.11 [0.42]	104
OWT-4	25.5 [26.0]	5.7 [11.7]	0.22 [0.45]	178
OWT-5	51.2 [62.9]	47.7 [72.5]	0.93 [1.15]	452
OWT-6	337.3 [314.9]	669.8 [577.7]	1.99 [1.83]	16
OWT-7	22.5 [31.0]	4.6 [19.7]	0.2 [0.6]	48

missions by using each sensor's relative spectral response function. For HICO, the sensor's spectral response function results in center wavelengths spaced by  $\sim 5.7$  nm, each with a full-width half maximum (FWHM) of 10 nm from 400 to 745 nm and 20 nm for 746–900 nm, due to Gaussian smoothing ([Lucke et al., 2011](#); [NASA, 2021](#)). For PRISMA, the spectral response function has a FWHM of  $\sim 12$  nm, though the exact bandwidths are scene-dependent and vary as a function of the center wavelength ([Cogliati et al., 2021](#)). Most assessments within this paper leverage HICO band configurations, as the majority of the coincident *in*

*situ* observations ([Section 3.3](#)) are available for HICO.

### 3.2. Hyperspectral satellite image data

Applying the demonstration MDN to images taken from hyperspectral satellite sensors highlights the ability of the MDN (trained with *in situ* measurements) to produce realistic map products from  $R_{rs}$  despite uncertainties resulting from the atmospheric correction process as well as instrument systematic noise ([Ibrahim et al., 2018](#); [Moses et al., 2012](#)).

#### 3.2.1. HICO

Four HICO images (spatial resolution of  $\sim 90$  m) were chosen specifically because of available collocated *in situ* PC matchups (three from Lake Erie and one from the Chesapeake Bay) and an additional image (from Lake Erie) was chosen for comparison to historic observations of PC. The four images that have a total of 15 *in situ* matchups between them are from September 8<sup>th</sup>, 2014 (Lake Erie), June 16<sup>th</sup> 2014 (Lake Erie), August 19<sup>th</sup> 2013 (Lake Erie), and September 20<sup>th</sup> 2013 (Chesapeake Bay).

The HICO images were atmospherically corrected using the SeaWiFS Data Analysis System (SeaDAS v7.5.3) ([Ibrahim et al., 2018](#)) following the same procedure (using the default options) as in [Pahlevan et al. \(2021b\)](#). The atmospheric correction approach is relatively complex,

correcting for spectral transmittance effects and radiance resulting from multiple scattering by aerosols and air molecules. One of the more complex steps of the atmospheric correction process is the removal of the aerosol contribution from the top-of-atmosphere radiance. Removal of the aerosol contribution leverages a precomputed look-up-table of the optical properties of multiple different aerosol models, which are selected based on two bands in the NIR (Ibrahim et al., 2018; Gordon and Wang, 1994). The short and long wavelengths used in our atmospheric correction were 747 nm and 787 nm, respectively.

### 3.2.2. PRISMA

Application of the MDN to more recent images obtained from PRISMA (30 m spatial resolution) demonstrates the ability of the same training set and MDN architecture to be quickly adapted and applied to state-of-the-art (or future) hyperspectral missions. The MDN is retrained (Section 4) using the same *in situ* dataset resampled to the spectral response of PRISMA, but leaving out the Curonian Lagoon from the overall dataset (resulting in  $N = 929$ ) because it is used for validating one of the demonstration images. Two PRISMA images with *in situ* matchups were chosen. The images are of the shallow (mean depth of 3.8 m) fresh-to-brackish waters of the Curonian Lagoon on September 20<sup>th</sup> 2020 (Zemlys et al., 2013) and of the 124 km<sup>2</sup> large and shallow (max depth of 6.3 m) turbid waters of Lake Trasimeno (Italy), on July 25<sup>th</sup> 2020 (Ludovisi and Gaino, 2010). These two images have a total of seven *in situ* matchup locations.

PRISMA data were downloaded as L1 products (top-of-atmosphere calibrated radiance), then re-projected with a geographic lookup table (GLT) Bowtie Correction through prismaread (Busetto and Ranghetti, 2020), which also extracts ancillary information related to atmospheric correction (e.g., band centers and FWHM, Sun and view angles). The correction for atmospheric effects was performed with the Atmospheric and Topographic Correction (ATCOR v.9.3.0) (Richter and Schlpfer, 2002), which recovers reflectance spectra at the ground, from which the  $R_{rs}$  was computed by dividing ATCOR outputs by  $\pi$ .

### 3.3. Matchup data

Two different datasets are available to examine the impact of  $\Delta R_{rs}$  on satellite derived PC ( $PC^f$ ): (1) a dataset of  $PC^f$  aligned with *in situ* measured PC and (2) a dataset of  $PC^f$  aligned with PC estimated from *in situ*  $R_{rs}$  ( $PC^e$ ). For method (1), the available *in situ* PC measurements which are coaligned with HICO images ( $N = 15$ , from Lake Erie and the Chesapeake Bay) are dominated by low ( $<20 \text{ mg m}^{-3}$ ,  $N = 11$ ) and medium ( $20 < PC < 100 \text{ mg m}^{-3}$ ,  $N = 4$ ) PC. These *in situ* measurements span a PC range covering over an order of magnitude ( $0\text{--}32 \text{ mg m}^{-3}$ ) and a variety of PC:Chl $a$  ratios (0.03–0.76, ignoring the 0 PC value). Most of the reported *in situ* PC are direct measurements, except the PC from the Chesapeake Bay matchup (site label CB7.4), which is estimated from cyanobacteria biovolume following Kasinik et al., 2015. In addition to the HICO coaligned measurements, PRISMA aligned matchups ( $N = 7$ , from the Curonian Lagoon and Lake Trasimeno, with a range of PC:Chl $a$  (0.11–3.45)) offered a few higher PC estimates ( $>100 \text{ mg m}^{-3}$ ,  $N = 3$ ), where a high PC:Chl $a$  ( $\sim 3$ ) was observed.

For method (2), we used a separate dataset of  $PC^f$  coaligned with  $PC^e$  consisting of 65 near-coincident ( $\pm 3 \text{ h}$ ) *in situ* (limited to 409–724 nm) and HICO  $R_{rs}$  measurements to assess the impacts of  $\Delta R_{rs}$  on product retrieval consistency and, in extension, the quality of the retrieved PC product maps. The *in situ* and HICO  $R_{rs}$  matchups were taken from the lower Chesapeake Bay, the western basin of Lake Erie, and Florida estuaries (Keith et al., 2014; Schaeffer et al., 2015; Casey et al., 2019). Overall, the two datasets, the *in situ* measurements coaligned with hyperspectral satellite images and near-coincident  $R_{rs}$  matchups, are useful for testing the MDN architecture on a range of  $\Delta R_{rs}$  in a variety of aquatic conditions.

## 4. Methods

### 4.1. MDN architecture and hyperparameters

The MDN architecture uses novel band ratios (BRs) and line heights (LHs), as well as the BRs proven in existing MAs as inputs, replacing the  $R_{rs}$  used in previous MDN architectures for inverse aquatic remote sensing tasks (Fig. 4) (Balasubramanian et al., 2020; Pahlevan et al., 2020; Pahlevan et al., 2021b; Smith et al., 2021). All the MA BRs (Italics, Table 1) were used, while the novel BRs and LHs are chosen based on their correlation with PC. These input features were first normalized (e.g., log transformed and scaled between  $-1$  and  $1$ ), before being input to the neural network. The neural network was trained using default hyperparameters (Table 4) to produce estimates of the mean ( $\mu$ ), standard deviation ( $\sigma$ ), and mixing coefficient ( $\alpha$ ) for five Gaussians. The resultant products (PC and Chl $a$ ) were then estimated from the network by combining the Gaussians (through a combination function) and selecting the maximum likelihood estimate. The training process was repeated 10 times, with random initializations of the network, and the median output of the 10 different networks was used as the final product estimate.

#### 4.1.1. MDN input and output

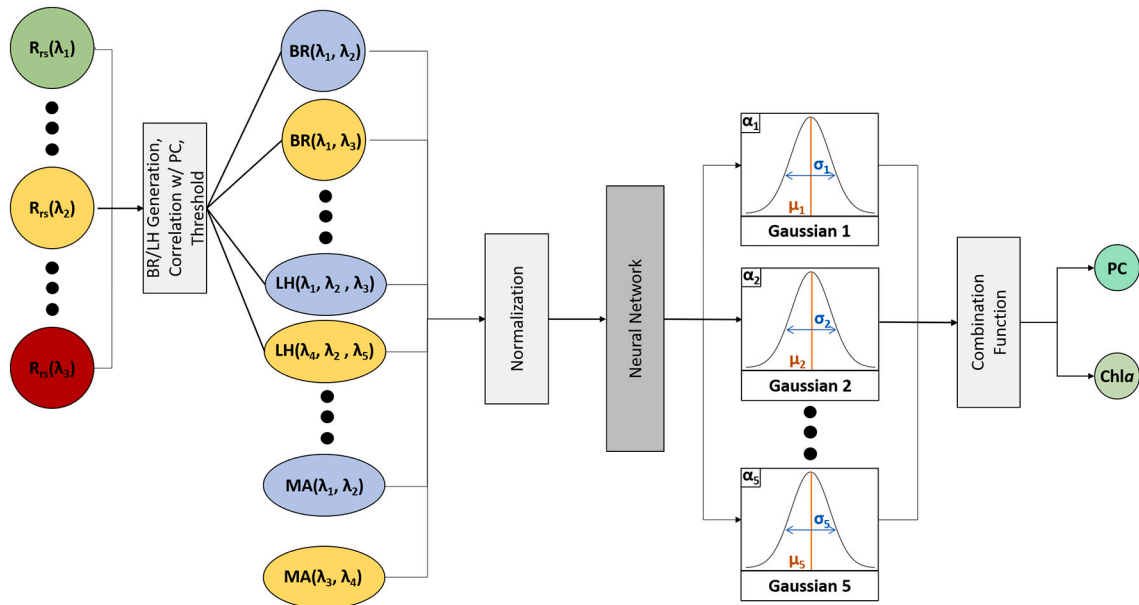
Although HICO and PRISMA offer a wide range of spectral bands from the blue to the NIR, only a specific subset of bands was considered (501–724 nm). Section 3.2 Discarding the blue bands is unlikely to reduce the accuracy of the PC estimation, as PC absorption is insignificant in the blue region (Table 1). This may have limited the ability of the MDN to correct for a strong CDOM absorption, but previous algorithms have had success without correcting for this effect while using NIR/red band ratios (Simis et al., 2007), likely due to the spectrally neutral errors in the red-NIR caused by CDOM absorption.

The MDN is trained on both spectral shape and amplitude through a combination of BRs and LHs, thus the MDN does not require the spectral  $R_{rs}$  as direct input. Band ratios were chosen as they capture relationships between separate bands (Eq. (1)), and LHs were chosen as they are relatively insensitive to spectrally neutral  $\Delta R_{rs}$  due to baseline subtraction (Eq. (2)) (Kudela et al., 2015; Wynne et al., 2010; Matthews et al., 2012).

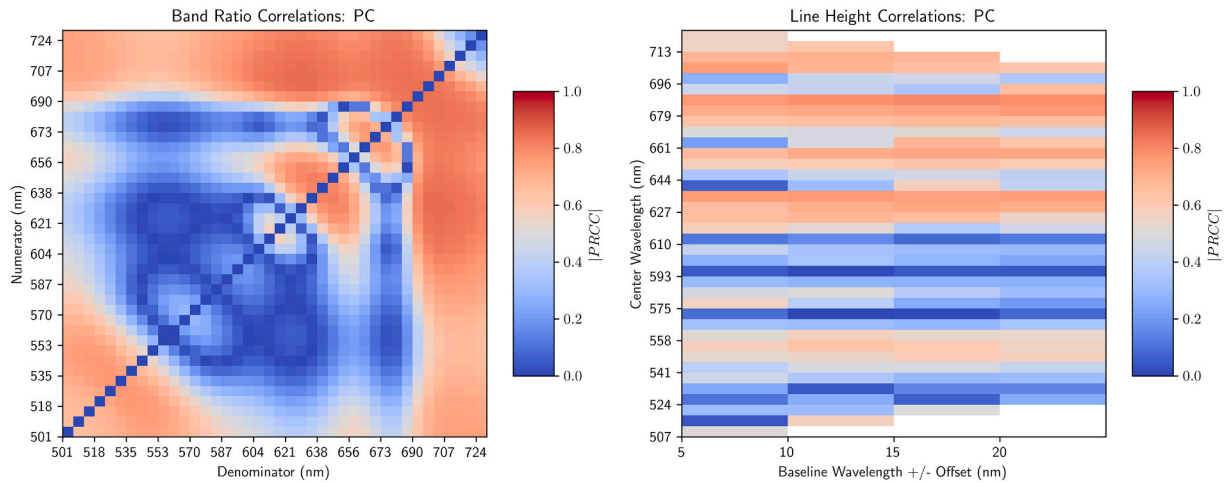
$$BR(\lambda_0, \lambda_1) = R_{rs}(\lambda_0)/R_{rs}(\lambda_1) \quad (1)$$

$$LH(\lambda_0, \lambda_1, \lambda_2) = R_{rs}(\lambda_1) - (R_{rs}(\lambda_2) + ((R_{rs}(\lambda_0) - R_{rs}(\lambda_2)) * (\lambda_2 - \lambda_1)) / (\lambda_2 - \lambda_0)) \quad (2)$$

To select the optimal BR and LH combinations, we searched the entire wavelength range used (501–724 nm). For band ratios, each wavelength combination was searched where the numerator was greater than the denominator, as it is redundant to include the inverse ratio. LHs were included for band centers between 507 and 719 nm where the baseline spanned 5, 10, 15, or 20 nm from the center wavelength. The correlations between PC and the individual BRs and LHs gave a first order estimate of how applicable each band combination was for estimating PC (Fig. 5). A threshold was applied to the absolute value of Pearson's Ranked Correlation Coefficient, at a correlation of 0.35, to exclude completely uncorrelated band ratios. In general, the spectral bands that were the most correlated with PC were 500–550, 620, 650, and 710 nm, which are the bands most typically associated with PC and Chl $a$  (Fig. 5). We used all BRs and LHs with a correlation above 0.35 (458 and 92, respectively), to allow for the MDN to fully utilize the hyperspectral data, in a format that is generally insensitive to  $\Delta R_{rs}$ . Of the small subspace of thresholds searched, this threshold was chosen as it was found to produce an MDN with high resiliency to  $\Delta R_{rs}$ . Additionally, the use of multiple features near any specific wavelength combination (e.g., the 10's of BR's available in the 709 nm/(610–630 nm) range) allow for the MDN to learn and adjust for contributions to individual features (e.g., 620 nm dip) due to absorption by other



**Fig. 4.** A block diagram representing the MDN architecture used for all (testing and demonstration) MDNs in this paper. The BRs and LHs are selected based on their individual correlation with PC from the entire dataset. The band ratios (BR( $\lambda_1, \lambda_2$ )), line heights (LH( $\lambda_1, \lambda_2, \lambda_3$ )), and operational multispectral algorithms (MA( $\lambda_1, \lambda_2$ ), Table 1) are first normalized, before being run through the neural network. The output is selected through a combination function as the maximum likelihood of the combination of five weighted probability density functions (using the mixing coefficient,  $\alpha_x$ , with their associated mean ( $\mu_x$ ) and standard deviation ( $\sigma_x$ ), to solve the non-unique inverse problem.



**Fig. 5.** The absolute value of Pearson's ranked correlation coefficient ( $|PRCC|$ ) between a variety of band ratios and line heights. The distance between the center wavelength and the evenly spaced bands on either side of the center wavelength that make up the baseline are labeled as 'Baseline Wavelength +/- Offset (nm)'. Blank cells were not calculated, because the Baseline Wavelengths exceeded the available spectral range (501–724 nm).

optically active constituents (e.g., accessory pigments such as Chlb, CDOM, sediment). In addition to this wide array of BRs and LHs, a variety of ratios and LHs from the existing MAs were also added (Table 1, *Italics text*). Overall, these features spanned and fully leveraged the available hyperspectral information in contrast to MAs that are inflexible and restricted to an individual feature.

Both PC and Chla are output from the model. While the goal of the model was to produce accurate PC estimates, including Chla during training allowed the MDNs to learn the PC-Chla covariances and increase the accuracy of PC measurements (Bishop, 1994), as PC was present in varying ratios to Chla, but affects the  $R_{rs}$  used to recover PC.

#### 4.1.2. Hyperparameters

In previous studies using MDNs for inverse problems in aquatic remote sensing, the results were found to be relatively insensitive to the

hyperparameter selection (Smith et al., 2021). Therefore, the default hyperparameters provided with the model are included in Table 4.

**Table 4**  
MDN hyperparameters.

Hyperparameter	Value
Training iterations	10,000
Number of Gaussians	5
Neurons per hidden layer	100
Hidden layers	5
Learning rate	1e-3
L2 regularization	1e-3
Epsilon	1e-3



#### 4.2. Performance assessment: 50/50 split

We first trained a theoretical MDN using a 50/50 training/testing split, which was randomly selected from all the regions in the overall dataset. By splitting the data in this manner, the generalization performance was characterized for regions represented within the training set (Section 5.1). While this version of the model was not used for demonstration throughout the rest of the paper, as it was only trained on half the data, it was used to estimate the generalization performance expected from the final demonstration model on data from within the set (e.g., within regions included in the training set).

#### 4.3. Performance assessment: leave-one-out

We applied a leave-one-out testing approach to estimate how well the theoretical MDN is expected to be transferred to regions not included within the training set (Section 5.2). In leave-one-out testing, the training dataset consisted of the entire dataset *except* one regional subset, which was used for testing. A new model was tested on each subset, until all subsets had been excluded, to assess the performance of this model MDN on out-of-training data from a variety of regions. Since our *in situ* dataset was from six distinct regions, we were able to estimate the expected accuracy of the MDN to generalize on waters with similar aquatic conditions. The estimated performance from this method likely underestimates the performance of the MDN because *in situ* measurement techniques were not fully consistent and are expected to carry various degrees of uncertainties.

#### 4.4. Benchmarking: optimization of existing multispectral algorithms

In this paper, we compared the performance of MDNs against that of example MAs (Hunter et al., 2010; Schalles and Yacobi, 2000) and more complex semi-analytical models (Simis et al., 2007; Simis et al., 2005). These algorithms have been widely cited and adequately represent both empirical and semi-analytical algorithms. Each algorithm was implemented using the nearest available HICO and/or PRISMA spectral bands (with an upper wavelength limit of 764/760 nm). The wavelengths and coefficients of MAs were originally optimized to perform well against their individual, often regional, datasets. To provide a fair comparison between the MDN (which has been trained on a mix of the underlying datasets) and these multispectral algorithms, the *coefficients* of the MAs were optimized on the same training set that was used for the MDN (Appendix A).

Gradient boosted regression trees (Friedman, 2001) were used for sequential optimization of the coefficients, with respect to the cost function. We optimized the coefficients by minimizing the median symmetric accuracy (Morley et al., 2018) between the model estimates and the known *in situ* concentration. To penalize the MAs forcing a significant fraction of the estimates to negative concentrations (which are not usable in the described optimization routine), the negative concentrations were forced to be the mean of the known concentrations ( $y$ ) if the negative (and non-finite) concentrations are larger than 25% of the data. This induced a penalty on the optimization algorithm for estimating high fractions of negative PC estimates, and therefore over-training for a specific concentration range. The bounds were set to be either four times higher/lower than the default algorithm coefficients or from  $-200$  to  $200$ , depending on the algorithm. 1500 initial points were used, using Latin Hypercube Sequence to set the initial points. Optimization is performed for 300 additional iterations. The optimized coefficients are shown in Table A1 in Appendix A.

#### 4.5. Impact of $\Delta R_{rs}$ on PC retrieval

We examined the impact of  $\Delta R_{rs}$  on PC retrieval accuracy in two ways: (1) by comparing  $PC^t$  to *in situ* measured PC and (2) by comparing  $PC^t$  to  $PC^e$  (Sections 5.3 & 5.4). The first method, which compared  $PC^t$  to

coalgined *in situ* measured PC, gave a semi-quantitative estimate of how well the demonstration MDN was expected to perform despite  $\Delta R_{rs}$  of hyperspectral satellite images. Due to a limited number of coaligned measurements, a second method, comparing  $PC^t$  to  $PC^e$ , was used to provide further qualitative evidence on the sensitivity of the algorithm to  $\Delta R_{rs}$ . The results from the second method were compared between the demonstration MDN and an MDN architecture which uses solely  $R_{rs}$  as input (MDN- $R_{rs}$ ) to determine the relative insensitivity of the demonstration MDN architectures (which uses LHs as input, Section 4.1) to  $\Delta R_{rs}$ .

### 5. Results

#### 5.1. 50/50 training/testing split: comparison to existing multispectral algorithms

In general, the MDN achieved lower median symmetric accuracy (termed hereafter as “uncertainty”;  $\epsilon$ ) and symmetric signed percentage bias (termed hereafter as bias;  $\beta$ ) (Morley et al., 2018) when compared against the optimized MAs. The retrieval uncertainties from the MDN (44.3%), when trained on a randomly selected half of the dataset, was significantly lower than MAs ( $\sim 90$ – $115\%$ ) optimized on the same training set (Fig. 6). The difference in  $\epsilon$  between the MDN and MAs was significantly higher than that between *in situ* replicates ( $\sim 10\%$ ) reported in previous studies (Song et al., 2012). The MDN uncertainties were on the order of (but higher than) the potential uncertainties in PC extraction between different extraction techniques ( $\sim 30\%$ , (Zimba, 2012)). To allow for comparisons with previous studies, a variety of alternative performance metrics (including the root mean square difference (RMSD) and mean absolute percentage difference (MAPD)) are also documented in Table A2 in Appendix A, and these metrics in general show similar results, with the MDN retrievals achieving the lowest uncertainties. Even while being re-trained on the generally lower PC contained within the training dataset, the MAs all overestimated PC, and often produced invalid (e.g., negative, or infinite) results, which falsely reduce the reported  $\epsilon$ , as the invalid concentrations are not included in the calculation of  $\epsilon$ . In general, the MAs typically overestimated low PC ( $<10 \text{ mg m}^{-3}$ ), performed particularly well in the medium-high ( $10$ – $100 \text{ mg m}^{-3}$ ) concentration range, and slightly underestimated at high PC ( $>100 \text{ mg m}^{-3}$ ). The linearity of the MDN extended over the entire concentration range, most notably on the lowest PC ( $<10 \text{ mg m}^{-3}$ ). There are a few outliers at  $10^{-1}$  due to the sheer number of estimates available at this concentration, which represents the  $0$ – $10^{-1}$  range (Fig. 2). While the optimized algorithms are shown in Fig. 6, the unoptimized coefficients produced results in the form of (%  $\epsilon$ , # invalid) for Schalles and Yacobi, 2000 (111.1%, 228 invalid), Simis et al., 2007 (95.0%, 30 invalid), Hunter et al., 2010 (170.0%, 20 invalid), showing that the optimization process typically reduced  $\epsilon$  while maintaining an equivalent number of invalid estimates (except for the Schalles and Yacobi, 2000 algorithm, where the  $\epsilon$  increased slightly but the number of invalid estimates dropped dramatically). Overall, the MDN architecture, when trained on half the dataset and tested on the other half, produced estimates with less than half the  $\epsilon$  of that of MAs, was notably more accurate at low concentrations, and did not produce any invalid results.

#### 5.2. Leave-one-out testing: MDN transferability

In the leave-one-out analysis, the MDN architecture and remaining data generalized better in regions with low-medium PC (Erie and Italy datasets Table 5) than existing MAs, while maintaining competitive accuracy in regions with medium-high PC (Dutch Lakes, Indiana, the Curonian Lagoon, Table 5). The Simis et al. algorithm (Simis et al., 2005; Simis et al., 2007) performed slightly better in the medium-high range (using its default coefficients). The Schalles and Yacobi, 2000 algorithm performed best on regions with high concentrations (South Africa and Spain, Table 5). The MDN performed notably poorly when estimating in

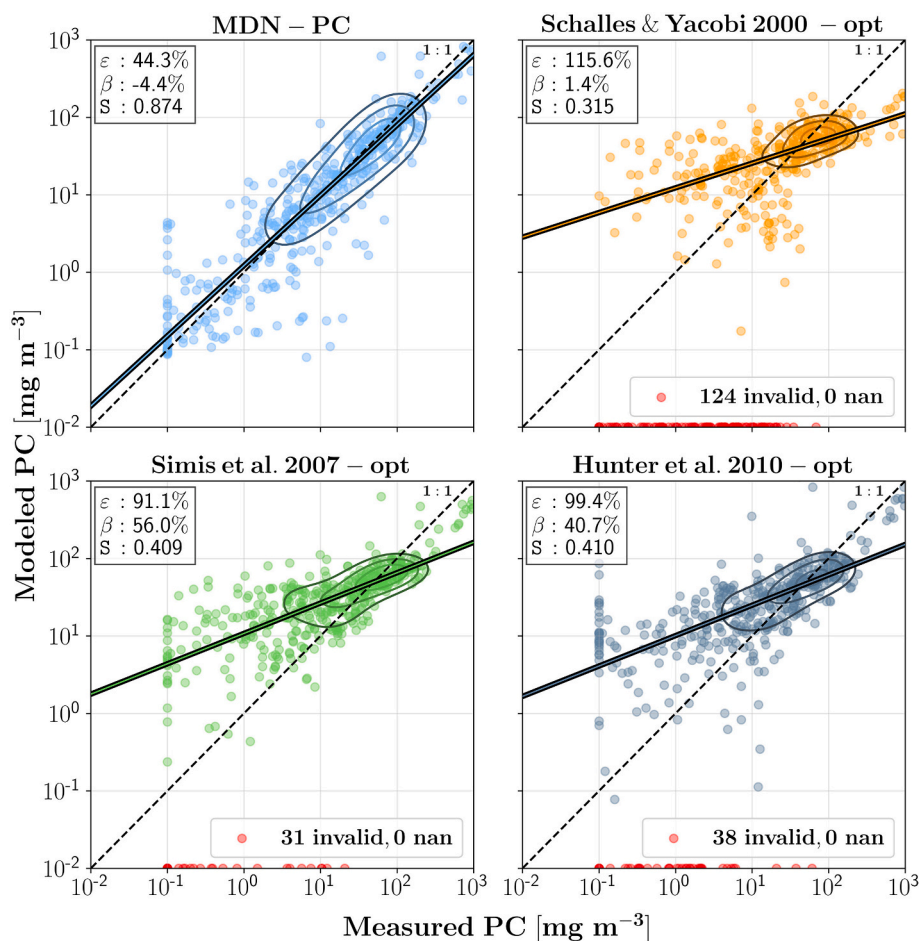


Fig. 6. Performance assessment of MDN and optimized (–opt) multispectral algorithms (Section 4.4) using *in situ* measured PC on a randomly selected half of the total set ( $N = 470$ ) with BR and LH features computed from HICO-simulated  $R_{rs}$  spectra as model input. Invalid counts (red) consist of negative, non-finite, or NaN values (NaN value counts are also shown). Reported metrics are median symmetric accuracy ( $\epsilon$ ), symmetric signed percentage bias ( $\beta$ ) (Morley et al., 2018), and the slope of linear regression (S). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 5

Results of leave-one-out testing, with % $\epsilon$  in PC estimates followed by the invalid number of estimates in parenthesis for each estimation method (rows), when each region (column) is used as the testing set.

Estimation Method	Lake Erie	Dutch Lakes	Lakes of Indiana	Curonian Lagoon	Lakes of Italy	South African Reservoirs	Lakes of Spain
MDN	<b>183.8</b>	95.4	46.9	56.1	<b>175</b>	568.4	155.1
Schalles and Yacobi (2000)	195.6 (228)	111.3 (51)	74.0 (43)	65.5 (18)	(20)	<b>67.9 (3)</b>	<b>75.6 (78)</b>
Simis et al. (2007)	187.0 (26)	<b>75.8</b>	<b>45.4</b>	<b>50.6</b>	328.8	139.8	81.5 (23)
Hunter et al. (2010)	206.0	128.7	140.1	167.5	231.3	238.7 (1)	273.0

Lowest  $\epsilon$  method for each region is bolded. See Table 2 for the number of samples and statistics.

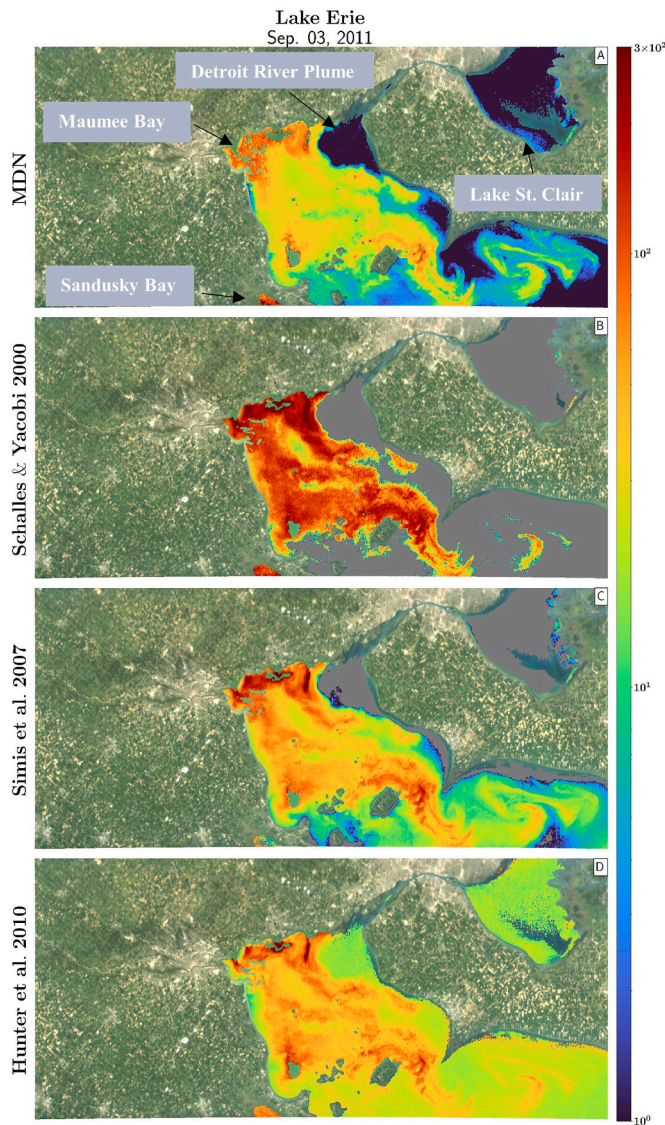
regions with the most intense blooms (South Africa), but this can be explained due to the South Africa dataset, containing a large proportion of the most intense CHABs, and with minimal similarities to datasets in other geographic regions.

5.3. Mapping PC: demonstration on satellite observations

The retrieved product map produced by the MDN agreed with the general distribution of estimates predicted from three MAs (Fig. 7), but the magnitude of the estimates varies substantially. For example, all four algorithms produced relatively low (or sometimes, in the case of the Schalles and Yacobi, 2000 and Simis et al., 2007, no) estimates in Lake St. Clair (upper right water body), the Detroit river (connecting Lake St. Clair to Western Lake Erie), the Detroit River plume, and the central basin of Lake Erie. The very low estimates provided by the MDN agree best with historic observations, which show minimal cyanobacteria in the Detroit River plume (which is fed by Lake St. Clair) and the central basin (Binding et al., 2019). The transition region between the Detroit

River plume and the Maumee Bay, the Maumee Bay itself, and the Sandusky Bay exhibited the highest retrieved concentrations for both the MDN and the MAs, conforming to previous reports (Moore et al., 2017; Binding et al., 2019). Similar trends are also present for the three other HICO-derived maps of Lake Erie (Fig. 8). The greater extent of the September 8<sup>th</sup>, 2014 image captures more of the cHAB in the Sandusky Bay region, as well as the gradient of the cyanobacteria plume from the Sandusky Bay into the central basin of Lake Erie (Binding et al., 2019) (Fig. 8, Panel A). Overall, the algorithm matched the spatial patterns of the existing MAs in Lake Erie, while providing estimates in low-medium PC regions that agree better with historic observations.

Independent matchups of western Lake Erie (Fig. 8, Panels A, C, & E) and the Chesapeake Bay (Fig. 8, Panel G) served as a semi-quantitative metric for assessing HICO-derived estimates from the MDN (Table 6). The limited *in situ* measurements were semi-quantitatively compared against the PC retrieved from the models and MAs using both the satellite measured  $R_{rs}$  and *in situ* measured  $R_{rs}$  (where available) in Table 6. In general, the MDN retrievals agreed with very low PC measurements



**Fig. 7.** PC product maps of western Lake Erie and Lake St. Clair on September 9th, 2011 produced by the demonstration MDN and operational multispectral algorithms (Schalles and Yacobi, 2000; Hunter et al., 2010; Simis et al., 2007; Simis et al., 2005). Aquatic regions with negative  $R_{rs}$  or those flagged by the atmospheric correction do not have PC estimates (e.g., Lake St. Clair), and the underlying RGB image is displayed. Non-finite or negative PC retrievals are flagged with grey.

( $\sim 1 \text{ mg m}^{-3}$ , Fig. 8 Panel C and G), while MAs failed to generate valid (e.g., finite positive) or accurate estimates at these concentrations (Table 6). The estimates for medium PC were more accurately estimated by the Hunter et al., 2010 and Simis et al., 2007 algorithms, though the Hunter et al., 2010 algorithm commonly produced predictions within this range as seen in Fig. 7. The *in situ* measurements did not include high PC for  $\Delta R_{rs}$  sensitivity analysis.

Comparison between products retrieved from the MDN versus MAs produced from PRISMA (Fig. 9) and independent *in situ* measurements at Lake Trasimeno and the Curonian Lagoon showed similar results, but also included information on high PC matchups (Table 6). In general, the MDN estimated low PC ( $< 20 \text{ mg m}^{-3}$ ) sites (e.g., the Curonian Lagoon sites 37b & Uostadvaris as well as Lake Trasimeno sites TRS30 & TRS35) relatively well, though the Simis et al., 2007 algorithm also performed well in this range. The increased accuracy at low concentrations held despite low PC:Chla ratios ( $\sim 0.1$ – $0.2$ ). Interestingly, all test algorithms severely underestimated PC (by a factor of  $\sim 2$ ) at high ( $> 100 \text{ mg m}^{-3}$ )

PC, potentially due to a high PC:Chla ( $\sim 3$ ). Overall, these results, where the MDN performed best at low concentrations despite large differences in the PC:Chla ratio, match with those derived from HICO retrieved images.

#### 5.4. Demonstration MDN: sensitivity to $\Delta R_{rs}$

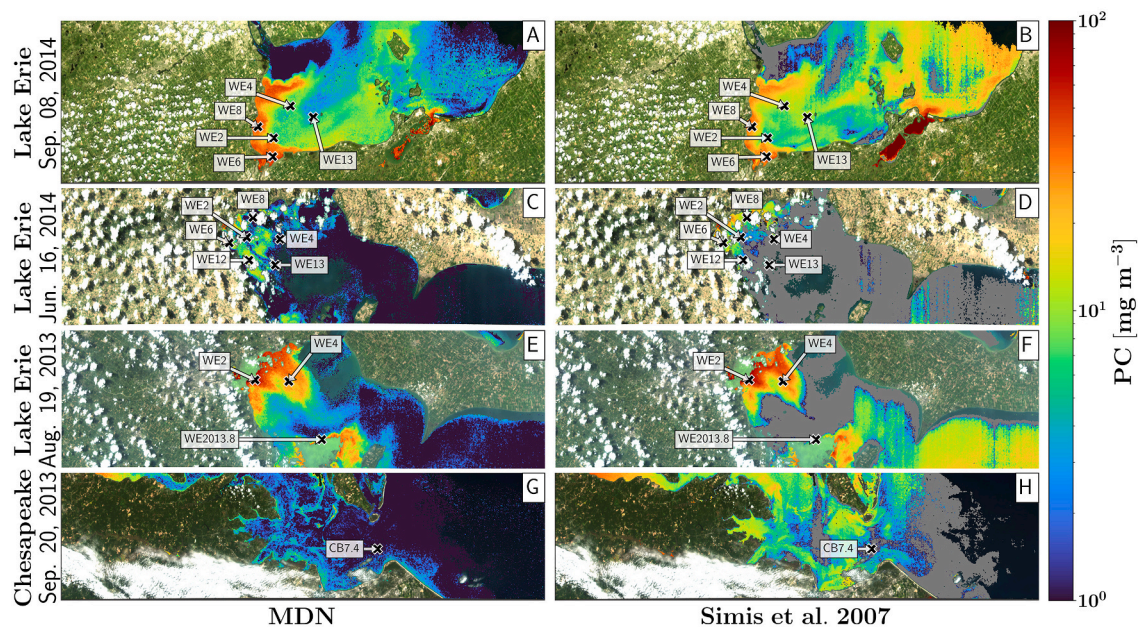
As previously shown, the demonstration MDN, which was trained on the entire ( $N = 939$ ) *in situ* dataset, is sensitive to  $\Delta R_{rs}$ . A small independent dataset of *in situ*  $R_{rs}$  and HICO measured  $R_{rs}$  (Fig. 1), covering additional sites (e.g., Florida estuaries including Pensacola Bay) (Section 3.3, Pahlevan et al., 2021b), were used to estimate PC from the MDN ( $PC^e$  and  $PC^f$ , Fig. 10, Panel A). The same dataset was used to estimate  $PC^e$  and  $PC^f$  from an MDN with the same architecture but used  $R_{rs}$  as input instead of both BRs and LHs (MDN- $R_{rs}$ , Fig. 10, Panel B). The demonstration MDN (Fig. 10, Panel A) achieved a higher Spearman's rank correlation coefficient (Spr) and a more linear slope as compared to the MDN- $R_{rs}$  (Fig. 10, Panel B). Interestingly, the  $\Delta R_{rs}$  forced the demonstration MDN estimated  $PC^e$  of  $< \sim 3 \text{ mg m}^{-3}$  to  $PC^f$  of  $< 1 \text{ mg m}^{-3}$ . While inaccurate, this sensitivity is more practical from a water quality monitoring standpoint than the alternative of MDN- $R_{rs}$ , where  $PC^e$  within the  $0.1$ – $1 \text{ mg m}^{-3}$  range were mapped to  $PC^f$  of  $6$ – $10 \text{ mg m}^{-3}$ , potentially leading to false alarms. The false alarms are apparent in product maps for the same regions produced by the demonstration MDN and MDN- $R_{rs}$  (Appendix B, Fig. B1, arrows 1–5 showing pronounced overestimation in Lake St. Clair and the Detroit River Plume). The five stations with measured PC (blue dots with black borders) lined up near to the 1:1 line for the demonstration MDN (Fig. 10, Panel A). The overall estimates of  $PC^f$  were not biased in one direction relative to *in situ* measurements (pink dots, *in situ* concentrations replace  $PC^e$ ). While no *in situ* measurements of PC were available for most of the matchups (blue dots), CHABS were not specifically noted (Keith et al., 2014), so generally very low estimates were expected. Overall, although the  $\Delta R_{rs}$  biases  $PC^e < \sim 3 \text{ mg m}^{-3}$  to even lower  $PC^f$  estimates for the demonstration MDN, the reduced false alarms, the better linearity, and the higher Spr demonstrate its overall desirable response to  $\Delta R_{rs}$  relative to MDN- $R_{rs}$ .

The PC retrievals from individual *in situ*  $R_{rs}$  and HICO retrieved  $R_{rs}$  showed the impact of  $\Delta R_{rs}$  in different optical scenarios (Fig. 11). Large offsets in the  $R_{rs}$  (stations WE4, WE6, WE13, and WE8) resulted in a difference of  $\pm 100\%$  between  $PC^f$  and  $PC^e$ . While most of the matchup sites (e.g., those not from western Lake Erie) did not have associated PC matchups, visually examining spectral differences in the remote and *in situ*  $R_{rs}$  provided some information on the expected differences in the PC estimates. For example, WE8 and WE6 exhibited substantially larger PC induced absorption dips in the remotely sensed  $R_{rs}$  near  $620 \text{ nm}$  (relative to the *in situ*  $R_{rs}$ ), and as expected have substantially larger  $PC^f$ . Finally, some of the matchups have different spectral shapes (e.g., CH01, showing that there may be significant temporal and/or spatial misalignment for certain stations). Overall, the matchups displayed substantially different  $R_{rs}$  spectra, which heavily impacts the retrieved PC estimates.

## 6. Discussion

### 6.1. Generalizability of MDN

The MDN, which used BRs and LHs spanning the hyperspectral range that were correlated with PC, achieved lower estimation uncertainties on testing datasets spanning multiple orders of magnitude relative to the existing MAs that only make use of one to two band ratios (Fig. 6, Table 1). Even when the coefficients of MAs were optimized to the training dataset, they were unable to match the linearity across the PC range that covers four orders of magnitude. The MDN continued to generalize best on low PC ( $< 20 \text{ mg m}^{-3}$ ), even when it had not been trained using data from the specific region. It is to be expected that the MDN will perform well in this range, as low PC measurements comprise



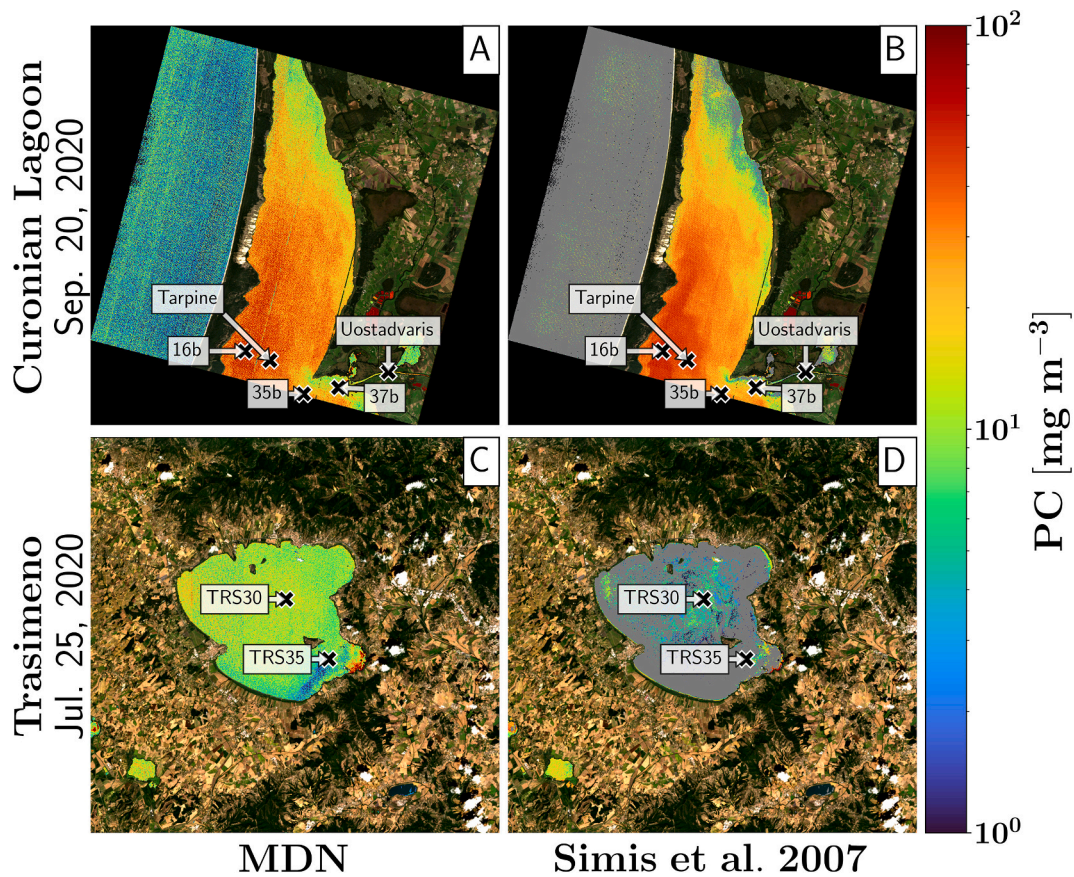
**Fig. 8.** HICO product maps retrieved using the MDN, with associated *in situ* matchups (labels, which match with *in situ* concentrations in Table 6) from Lake Erie and the Chesapeake Bay. Aquatic regions with negative (or 0)  $R_{rs}$  or those flagged by the atmospheric correction algorithm (e.g., the bloom below WE2013.8 in panels E/ F) do not have PC estimates, and the underlying RGB image is displayed. Non-finite or negative estimates are flagged with grey. These images have not been geolocated.

**Table 6**

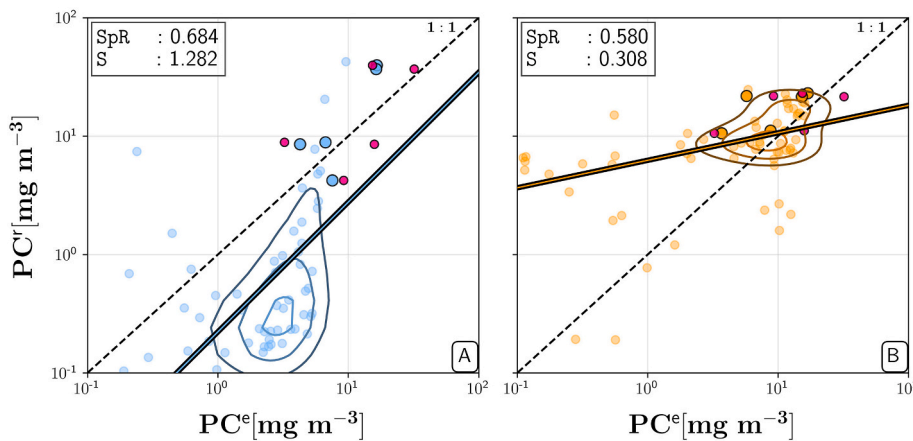
*In situ* measured PC compared against retrieved PC from *in situ* and remotely measured  $R_{rs}$  ( $PC^e$  and  $PC^f$  respectively; Section 5.3) from the labeled stations within Figs. 8 and 9.

Station	PC		MDN	Hunter 2010		Simis 2007	MDN		Hunter 2010		Simis 2007
	<i>insitu</i>	Chla		$PC^e$			$PC^f$				
Erie 09/08/2014											
WE2	3.24	20.77	<b>6.7</b>	8.6	13.9	11.2 [8.9]	16.4	<b>8.1</b>			
WE4	15.81	26.24	4.3	12.3	<b>13.25</b>	8.7 [8.5]	20.0	<b>11.5</b>			
WE6	15.32	53.38	<b>16.5</b>	13.1	31.9	36.9 [39.7]	<b>19.7</b>	29.4			
WE8	31.94	41.82	16.3	12.2	<b>27.1</b>	39.4 [36.9]	21.4	<b>28.5</b>			
WE13	9.21	13.09	7.5	<b>10.2</b>	12.1	4.7 [4.2]	20.6	<b>10.5</b>			
Erie 08/19/2013											
WE2	24.82	89.08	N/A	N/A	N/A	59.7	<b>44.2</b>	56.0			
WE4	30.82	62.18	N/A	N/A	N/A	22.7	<b>31.2</b>	22.4			
WE2013.8	29.9	57.27	N/A	N/A	N/A	20.1	<b>32.6</b>	4.2 (6.3)			
Erie 06/16/2014											
WE2	0.31	10.72	N/A	N/A	N/A	4.4	8.5	<b>0.9</b>			
WE4	0.44	12.26	N/A	N/A	N/A	<b>(1.8)</b>	(28.8)	0			
WE6	1.82	23.04	N/A	N/A	N/A	12.1	<b>11.8</b>	17.7			
WE8	1.45	5.74	N/A	N/A	N/A	<b>(1.9)</b>	(11.6)	(6.2)			
WE12	1.06	34.37	N/A	N/A	N/A	(13.9)	(6.8)	<b>(2.5)</b>			
WE13	0.21	5.09	N/A	N/A	N/A	<b>1.7 (0.66)</b>	7.8 (9.7)	0			
Chesapeake 09/20/2013											
CB7.4	0	4.38	N/A	N/A	N/A	<b>0.9</b>	19.1	0			
Curonian 09/20/2020											
35b	156.74	50.35	<b>81.4</b>	29.67	76.9	<b>43.0 (28.9)</b>	23.5	34.7			
16b	152.67	44.21	64.8	31.1	<b>87.6</b>	53.4	<b>54.5 (34.9)</b>	51.7			
Tarpine	141.90	46.68	65.7	30.7	<b>81.8</b>	31.5 (45.9)	32.6 (30.5)	<b>44.6</b>			
37b	6.89	60.39	<b>10.1</b>	17.5	34.5	<b>12.7 (11.3)</b>	28.3 (24.2)	15.6			
Uostadvaris	17.39	82.13	13.8	<b>19.7</b>	40.5	<b>17.0 (14.0)</b>	22.8 (16.0)	2.9 (10.7)			
Trasimeno 07/25/2020											
TRS30	3.5	18.7	N/A	N/A	N/A	7.1 (10.7)	19.8	<b>3.8</b>			
TRS35	2.2	13.2	N/A	N/A	N/A	<b>4.4 (4.6)</b>	21.6 (10.7)	<b>(0)</b>			

When the  $R_{rs}$  were masked by the atmospheric correction algorithm or the *in situ*  $R_{rs}$  did not exist, the results are reported as N/A. When the algorithm failed to produce results, the cells are filled with empty parenthesis (e.g., ()). Parenthesis are filled with median concentrations from a  $7 \times 7$  window for PRISMA and a  $3 \times 3$  window for HICO are shown within parenthesis. The highest accuracy  $PC^e$  and  $PC^f$  products from the MDN and MAs are bolded. Units are  $mg\ m^{-3}$ . Brackets ([ ]) surrounding the concentrations from Erie on 09/08/2014 denote the estimates from Fig. 11 (MDN estimates for 09/08/2014 are not the exact same as Fig. 11 due to changes in the geospatial alignment of the images). To save space, the Schalles and Yacobi, 2000 algorithm has been dropped due to a low number of valid results (7/22). Text highlighted in red would be incorrectly classified as low/medium/high risk (with 20 and  $95\ mg\ m^{-3}$  PC serving as the minimum values for medium and high risk respectively).



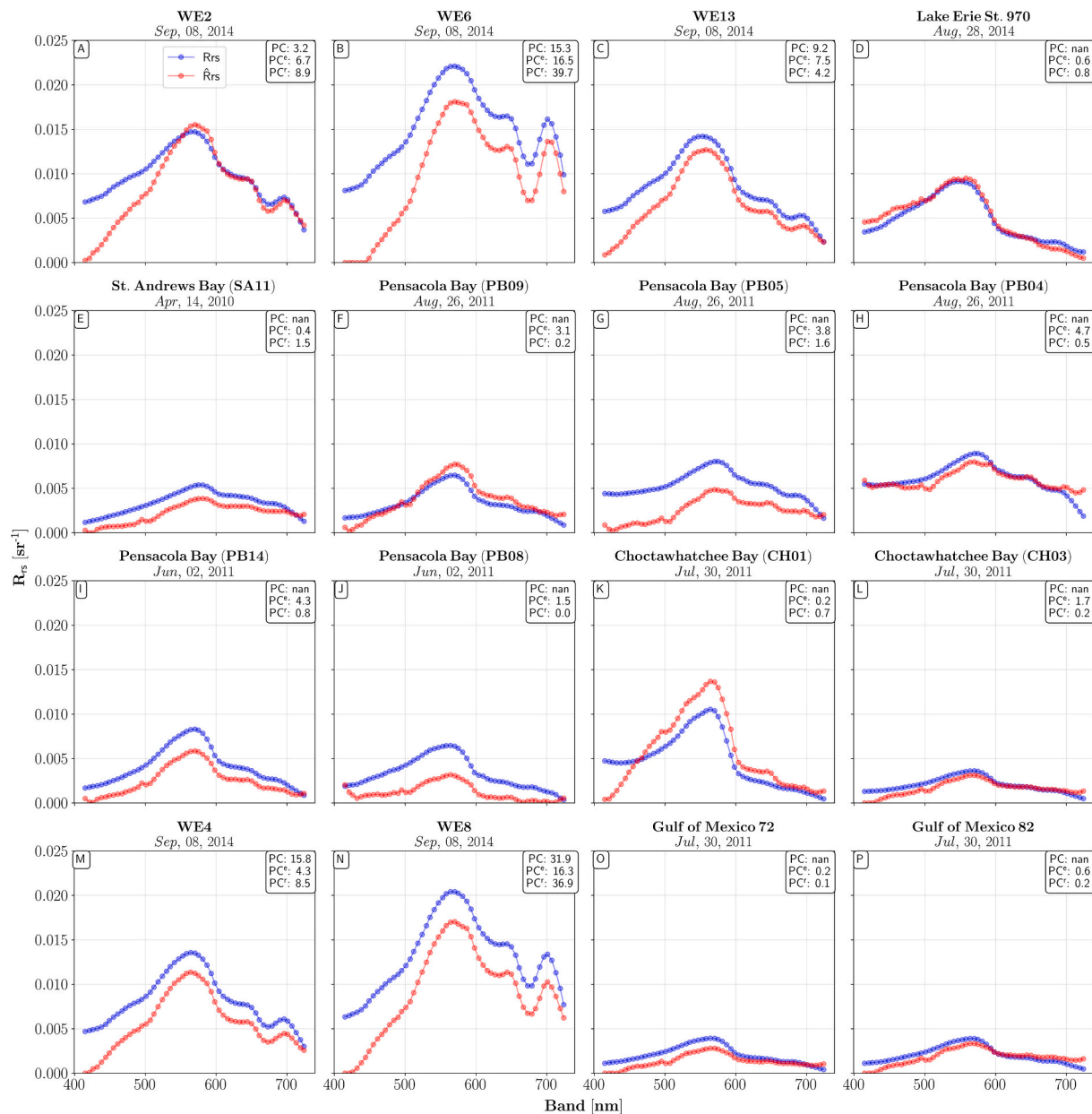
**Fig. 9.** PC maps of the Curonian Lagoon and Lake Trasimeno produced from PRISMA images, shown on the underlying RGB image, with *in situ* matchups (labels match with *in situ* concentrations in Table 6). Non-finite or negative concentrations are shown in grey. The normalized difference water index (NDWI) is used with a threshold of 0.1 to identify the water regions.



**Fig. 10.** PC estimated from *in situ* measured  $R_{rs}$  ( $PC^e$ ) and HICO derived  $R_{rs}$  ( $PC^c$ ) from the demonstrations MDN (blue dots, Panel A) and MDN- $R_{rs}$  (orange dots, Panel B) demonstrate the impact  $\Delta R_{rs}$  (Ibrahim et al., 2018; Moses et al., 2012) has on estimated PC ( $N = 65$ ). The five *in situ* matchups which have *in situ* PC measurements associated with them (blue dots with black borders) occur at higher concentrations. The  $PC^c$  estimates plotted against their *in situ* measurements (pink dots) demonstrate the accuracy of each MDN against *in situ* measured PC. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

the bulk of the training data (Figs. 2 & 3, Tables 2 and 3). Estimating low PC in diverse phytoplankton communities has been particularly challenging for the multispectral PC algorithms, due to the presence of other optically relevant constituents such as accessory pigments (e.g., Chlb) and CDOM (Simis et al., 2007; Ruiz-Verdú et al., 2008). The use of multiple bands near features of interest (e.g., 709/(610–630 nm) BRs, Fig. 4) may provide the additional information necessary to overcome the impact of accessory pigments and other optically relevant constituents on individual features of interest. While the MAs performed best on regions with higher average concentrations (Table 5), this came at the cost of invalid estimates, and significant over-estimates at lower PC.

The MDN may have particularly inaccurate results in regions with very high ( $> 200 \text{ mg m}^{-3}$ ) PC (e.g., South Africa, Table 5) where there is a non-linear relationship between PC and absorption at 620 nm, which may be due to a lower pigment-specific absorption coefficient (i.e., the package effect). The non-linear relationship at higher concentrations has led to underestimation of PC using MAs (Ruiz-Verdú et al., 2008). While machine learning techniques easily learn non-linear relationships given sufficient data, it is possible that the dearth of information at very high concentrations in the rest of the dataset restricts the ability of the MDN to learn the complex nonlinear relationship between band ratios and PC (Fig. 2, Table 5). Overall, the combination of the MDN architecture, its



**Fig. 11.** Site by site comparison of *in situ* measured PC versus the estimated PC using both *in situ* and remotely measured  $R_{rs}$  as input to the model ( $PC^e$  and  $PC^r$ , respectively). *In situ* PC measurements were only available for stations WE2, WE6, WE13, WE4, WE8 (panels A-C, M, and N). The atmospheric correction had significant impacts on the  $R_{rs}$ . Based on visual analysis of the Chla fluorescence band, some of the matchups may not be perfectly spatially/temporally aligned (e.g., CH01), but the locations were not removed to match (Pahlevan et al., 2021b) for ease of comparison.

inherent skill in learning multimodal distribution within the target space, and training dataset result in better generalization on low-medium ( $< 20 \text{ mg m}^{-3}$ ) PC conditions, without sacrificing substantial accuracy in the medium-high ( $\sim 50 \text{ mg m}^{-3}$ ) range (Table 5).

Comparison to historic observations and *in situ* matchups shows that by leveraging BRs and LHs the demonstration MDN does in fact better estimate low PC regions than existing multispectral algorithms, despite  $\Delta R_{rs}$ . First, the demonstration MDN produced generally low  $PC^r$  estimates despite the assumed  $\Delta R_{rs}$ , in regions without significant PC being reported, while MDN- $R_{rs}$  overestimates in these regions (Figs. 10 & 11). Second, the demonstration MDN produced maps with very low PC estimates in regions which have been observed to have low PC, such as Lake St. Clair, the Detroit River, and the central basin of Lake Erie (Figs. 7 & 8) (Binding et al., 2019; Moore et al., 2017). Third, the demonstration MDN produced maps that aligned with physical oceanographic characteristics. For example, in the Rappahannock River

(Fig. 8, Panel G, top left), the highest PC estimates were in the tidal brackish water region, and the PC estimates decreased in intensity as the river's salinity moved to the higher saline Chesapeake Bay waters, where low PC estimates were supported by *in situ* measurement (Table 6). Finally, the higher accuracy at lower concentrations produced by the combination of the architecture and training dataset is also supported quantitatively through comparison to *in situ* matchups from four regions, Lake Erie, the Chesapeake Bay, Lake Trasimeno, and the Curonian Lagoon, even with very low PC:Chla (Figs. 8 & 9, Table 6). This is particularly remarkable in the lower Chesapeake Bay where diverse phytoplankton groups (e.g., diatoms, dinoflagellates) with other accessory pigments are present (Harding et al., 2016). Interestingly, although the  $PC^r$  estimates typically underestimated the  $PC^e$  estimates at  $< 10 \text{ mg m}^{-3}$  (Fig. 10), the  $PC^r$  tended to overestimate relative to the *in situ* PC measurements (Table 6). Overall, the demonstration MDN produces high-accuracy results for regions with low PC despite  $\Delta R_{rs}$ , large

variations in the PC:Chla, and the use of two different sensors.

In a limited assessment ( $N = 22$ , Table 6) comparing the risk categories of cHABs *via* individual matchups, defined using the *in situ* PC, the MDN PC<sup>r</sup> retrievals can more accurately classify risk levels than *in situ* Chla measurements. The risk classification groups are defined as: 20,000 cells mL<sup>-1</sup> cyanobacteria or 10 mg m<sup>-3</sup> Chla in cyanobacteria dominated waters for low risk, 20,000–100,000 cells mL<sup>-1</sup> of cyanobacteria or 10–50 mg m<sup>-3</sup> Chla for moderate risk, and higher levels represent high risk (World Health Organization, 2003). Of the 15 stations with *in situ* PC measurements <~20 mg m<sup>-3</sup> (approximate cyanobacteria concentration of <20,000 cells mL<sup>-1</sup> (Bastien et al., 2011)), 12 stations would be flagged as moderate or high risk using the *in situ* Chla (>10 mg m<sup>-3</sup>) (Table 6). Of those 12 stations, three (Erie 09/08/14 WE6 & the Curonian Lagoon 37b & *Ustadvaris*) would be flagged as high-risk regions. Alternatively, if we use the PC<sup>r</sup> retrievals to determine whether a site has dangerous levels of cyanobacteria, only one of the 15 stations would be mislabeled as moderate risk (with an *in situ* PC concentrations of 15.32 mg m<sup>-3</sup>). Of the four moderate risk stations (with *in situ* PC ~20–95 mg m<sup>-3</sup>), none would be mislabeled by PC<sup>r</sup>, while three would be mislabeled as high-risk using *in situ* Chla. Finally, of the high-risk stations (PC > 95 mg m<sup>-3</sup>), all three stations in the Curonian Lagoon would be mislabeled as moderate-risk using PC<sup>r</sup> and two would be mislabeled using *in situ* Chla. This brief comparison demonstrates that the MDN produced less false positive estimates of low-risk regions using PC<sup>r</sup> from HICO instead of *in situ* Chla (*i.e.*, Chla may not be a reliable indicator for presence of cyanobacteria in all global regions that experience cHABs). Due to the limited dataset, additional matchups are necessary to determine if these results are applicable to a wider range of regions and their associated aquatic conditions, though the limited initial results are promising for determining the extent and impact of different risk categories.

## 6.2. Limitations of available dataset for training a generalized MDN

A couple of factors limit the accuracy and transferability of data from global regions. First, despite having a large dataset ( $N = 939$ ) and spatial extent (six regions from different parts of the globe) relative to previous studies, this is still a limited range on global scales, with minimal redundancy and obvious bias towards waterbodies known to be periodically affected by cyanobacteria. Most of the samples originate from only a few regions, and overall, the samples do not cover the wide range of aquatic conditions that exist, which may limit its generalization capability. Since certain cyanobacteria species have different vertical distributions, and the same species can exhibit different vertical distributions (Moore et al., 2019), the resultant  $R_{rs}$  may vary markedly for similar *in situ* PC measurements taken at a given depth. Second, the accuracy and generalizability are limited by the wide range of collection techniques for both  $R_{rs}$  and PC that are employed during data collection, which add uncertainty to the truth concentrations used during training. A larger dataset consisting of data collected using a common approach (Dierssen et al., 2020) would enable use of more sophisticated, and therefore potentially more accurate, algorithms for PC retrieval.

Spatial and temporal offsets may explain some of the inaccuracy between PC<sup>e</sup> and PC<sup>r</sup> retrievals (Figs. 10 & 11). First, without geolocating the images, HICO spatial accuracy can be severely limited (Pahlevan et al., 2021b), so PC<sup>r</sup> may be calculated from different  $R_{rs}$  than PC<sup>e</sup>. Second, variations in the surface cyanobacterial concentration (and in extension the PC) can occur between the time that the *in situ* measurements are taken and when the remote sensing images are captured (Kutser, 2004). The vertical structure of the cyanobacterial bloom and dominant species can vary both spatially and temporally, which affects the red and NIR signal used for PC retrieval (Moore et al., 2019). Changes in these bands, which comprise a substantial portion of the bands used by the MDN, would drastically alter the PC retrievals. Overall, relatively small spatial scale and temporal misalignment can result in large differences in the *in situ* measured PC and remotely

estimated PC.

A limitation of the current dataset is the lack of very high PC (>~200 mg m<sup>-3</sup>) which comprise the highest intensity blooms (Figs. 2 & 3). Both the model development dataset and the matchups suffer from a lack of very high PC measurements. The lack of high PC measurements in the training dataset reduces the ability for the MDN to learn the relationship between high PC waters and the hyperspectral band ratios. While high accuracy at higher concentrations is of interest scientifically, increasing the accuracy at higher concentrations is not of interest from a water quality management perspective, as estimates of high concentrations (Fig. 6) indicate the potential risk (Stumpf et al., 2016). The dearth of high PC measurement with associated satellite measurements is concerning for estimating the efficacy of these algorithms from satellite measurements because these regions may suffer most from  $\Delta R_{rs}$ , as cHABs may hinder atmospheric correction due to their high NIR reflectance (Fig. 3, OWT-6) (Ogashawara et al., 2013). While at such high concentrations, the blooms may become visually apparent in the RGB imagery (*e.g.*, Fig. 8, Panel E, below station ID WE2013.8), quantification of cHABs *via* visual cues is subjective and discouraged. If the bloom forms scum, the atmospheric correction methods often fail, hence models might consider alternative methods that do not depend on  $R_{rs}$  (Smith et al., 2021). Overall, substantially higher PC measurements are required for both the training dataset and satellite matchup dataset to increase and demonstrate estimation accuracy at these concentration levels from hyperspectral satellite images.

Further, although our strategy to minimize the sensitivity of MDN to  $\Delta R_{rs}$  proved to improve magnitudes of retrievals (Figs. 10 and B1), in general, the BR, LH and other derived features tend to amplify noise (Pahlevan et al., 2020). This is particularly pronounced for instruments with low radiometric fidelity, *i.e.*, low signal-to-noise ratios and/or systematic noise. For example, one may notice the speckled noise in both PRISMA- (Fig. 9A & C) and HICO-derived (Fig. B1A, C, E, and G) maps compared to maps derived from MAs or MDN- $R_{rs}$ . Under these scenarios, the model may overfit to some of the critical BR and LH features that contain large  $\Delta R_{rs}$ , leading to unnatural local variability. From a water resource management perspective, however, a median spatial filter can suppress some of the noise to provide smoother map products.

## 6.3. Future directions

A variety of different techniques could be attempted to increase the accuracy of the MDN for estimation of high PC waters from satellite sensors using the currently available dataset. Probably the most effective technique would be to use known blooms that have been mapped effectively and routinely, such as those in Lake Erie (Stumpf et al., 2016) or Lake Winnipeg (Binding et al., 2018), and assume a 1:1 relationship of PC to Chla (Table 4; Stumpf et al., 2016), which is adequate given the uncertainties in PC measurement. Since the remote retrieval of PC critically depends on reductions in  $\Delta R_{rs}$  from the atmospheric correction process (*e.g.*, Figs. 10 & 11), alternative techniques could focus on overcoming uncertainties in the atmospheric correction process. First, a subset of the independent matchups could be used during training as a validation dataset, to choose the model iteration which has the highest accuracy on the real-world satellite measurements. Second, a forward optical model, such as MODTRAN (Berk and Bernstein, 1989), could be used to generate simulated satellite  $R_{rs}$ , by simulating TOA radiance and then applying typical glint correction approaches, to better train the network on expected  $\Delta R_{rs}$  from the atmospheric correction process (Kravitz et al., 2021). Third, and finally, alternative atmospheric correction settings (*e.g.*, using ultraviolet spectral information for aerosol model selection) and approaches should be further explored for estimation in high PC regions (Frouin et al., 2019). Overall, while these techniques may slightly increase the accuracy for satellite remote sensing of high PC waters, the largest gains would likely come simply from a larger *in situ* training dataset containing more high PC concentrations.

While we focus on the application of the MDN architecture to *hyperspectral satellite* sensors, this architecture could be extended to other platforms and imagers. We assessed the performance of the MDN on satellite instruments, whose observations suffer from atmospheric absorption and scattering, but due to its high accuracy relative to existing multispectral algorithms, the MDN could also be useful for PC retrieval from low-altitude platforms. For example, the MDN could be applied to hyperspectral measurements from handheld sensors, tower platforms (Vansteenkewegen et al., 2019; O'Shea et al., 2020), or drones (Kwon et al., 2020), which would suffer from minimal atmospheric effects. The MDN architecture could also be adapted for multispectral instruments, though it would be challenging due to the lack of spectral bands available to accurately estimate PC, despite impacts on the  $R_{rs}$  from other optically relevant constituents in the water column. A potential benefit of using current operational multispectral sensors is that the  $\Delta R_{rs}$  may be less prone to image artifacts (e.g., striping) found in proof-of-concept or demonstration hyperspectral measurements. We currently plan to develop this approach for use on Sentinel-3's Ocean and Land Color Imager (OLCI). Overall, while we focus on the application of the MDN to the specific task of PC retrieval from *hyperspectral satellite* measurements, the MDN architecture is applicable to alternative platforms and could be retrained for sensors with limited spectral content.

## 7. Conclusion

This study was the first to develop Mixture Density Networks for the non-unique inversion problem of estimating PC from hyperspectral  $R_{rs}$ , apply them to images of proof-of-concept satellite missions, and assess their sensitivity to uncertainties in  $R_{rs}$ . The large dataset ( $N = 939$ ), which included a substantial number of lower PC than previously used for model training, increased the ability of the MDN to make predictions in areas with low PC relative to existing multispectral algorithms. The model performance was evaluated via both a commonly used training-testing data split and a leave-one-out approach to inform users of the range of model uncertainties anticipated from model predictions. The multispectral algorithms overestimated low PC, or produced invalid estimates, more frequently than the MDN, on both the testing dataset and atmospherically corrected satellite images. This was further corroborated via visual assessments of PC maps produced from multiple instances of HICO and PRISMA images over lakes and estuaries. We also demonstrated that the band-ratio and line-height features accessible via hyperspectral data diminish the model sensitivity to uncertainties in  $R_{rs}$ . With further widespread assessments using PRISMA imagery, the MDN

may be utilized for producing PC maps to complement and support water resource management practices, as the combination of more valid results and the reduced overestimation on low PC regions allows resources to be focused on regions with the highest probability of cHAB formation. A larger dataset, particularly with additional high PC measurements, will yet aid in quantifying intense cHABs. While these results are useful for post-processing of HICO data, and application to the current demonstration satellite mission, PRISMA, the MDNs could readily be retrained for the spectral configurations of next-generation hyperspectral imagers planned onboard PACE and FLEX. While this article focused on the utility of the MDN to hyperspectral sensors deployed on satellites, the presented MDN architecture could be extended to viable multispectral spaceborne sensors and/or to sensors deployed on low-altitude platforms (e.g., towers or drones).

## Code availability

The codes and pretrained models for retrieving PC from HICO or PRISMA data are available via <https://github.com/STREAM-RS/STREAM-RS>.

## Declaration of Competing Interest

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.

## Acknowledgements

We acknowledge the public *in situ* phycocyanin and/or phytoplankton databases prepared and curated by NOAA's Great Lakes monitoring program, the Chesapeake Bay Program, and their partners. We are thankful to Reagan Errera with NOAA GLERL for data collection and curation of Lake Erie dataset. This work was mainly supported through the NASA ROSES grants #80NSSC20M0235, PACE Science and Applications Team, and #80NSSC21K0499, Ocean Biology and Biogeochemistry (OBB) program, as well as the United States Geological Survey Landsat Science Team Award #140G0118C0011. Support through the PRISCAV project (ASI grant #2019-5-HH.0) and EU Horizon 2020 projects PRIMEWATER (Grant Agreement #870497) and CERTO (Grant Agreement #870349) is also acknowledged. Project carried out using PRISMA Products© of the Italian Space Agency (ASI) was delivered under an ASI License to use.

## Appendix A

**Table A1**

Optimized coefficients of MAs, default coefficients shown in parentheses.

Algorithm	Coefficient 1	Coefficient 2
Schalles and Yacobi (2000)	0.911091083942539 (0.97)	0.0027912243467980704 (0.000912)
Simis et al. (2007)	0.25175305920006175 (0.24)	
Hunter et al. (2010)	-19.353972440084295 (-4.96)	625.8629863442636 (266)

**Table A2**

Alternative performance metrics provided for direct comparison to previous studies.

Algorithm	$\epsilon$ [%]	$\beta$ [%]	Slope []	RMSD [ $\text{mg m}^{-3}$ ]	RMSLD []	MAD [ $\text{mg m}^{-3}$ ]	MAPD [%]	$\leq 0$   NaN
Original								
Schalles and Yacobi (2000)	111.059	54.201	0.38	119.54	1.558	74.769	70.104	228
Simis et al. (2007)	95.029	66.42	0.409	97.8	1.712	35.925	69.379	30
Hunter et al. (2010)	170.036	-22.253	0.304	99.057	1.722	39.746	69.731	20
Optimized								

(continued on next page)



Table A2 (continued)

Algorithm	$\epsilon$ [%]	$\beta$ [%]	Slope []	RMSD [ $\text{mg m}^{-3}$ ]	RMSLD []	MAD [ $\text{mg m}^{-3}$ ]	MAPD [%]	$\leq 0$   NaN
Schalles and Yacobi (2000)	115.586	<b>1.429</b>	0.315	106.592	1.542	47.927	61.913	124
Simis et al. (2007)	91.069	55.951	0.409	96.247	1.687	35.932	61.698	31
Hunter et al. (2010)	99.397	40.674	0.41	85.066	1.772	33.186	66.717	38
<b>MDN</b>	<b>44.339</b>	<b>-4.384</b>	<b>0.874</b>	<b>84.735</b>	<b>1.003</b>	<b>22.303</b>	<b>36.546</b>	<b>0</b>

Rows containing best performers within each column are shown in bold. Linear metrics (RMSD and MAPD) are not recommended for future performance assessments. Additional performance metrics include the root mean square logarithmic difference (RMSLD), mean absolute difference (MAD), and median absolute percentage difference (MAPD). Performance metrics are calculated on finite and positive estimates.

## Appendix B

Additional product maps have been included (Fig. B1) as a further means of visually comparing between the demonstration MDN, which uses LHs and BRs as inputs, and MDN- $R_{rs}$  which directly uses the  $R_{rs}$  as input. MDN- $R_{rs}$  often overestimates in regions with historically low PC (Fig. B1, arrows), which was expected due to overestimation due to  $\Delta R_{rs}$  seen in the matchups (Fig. 10).

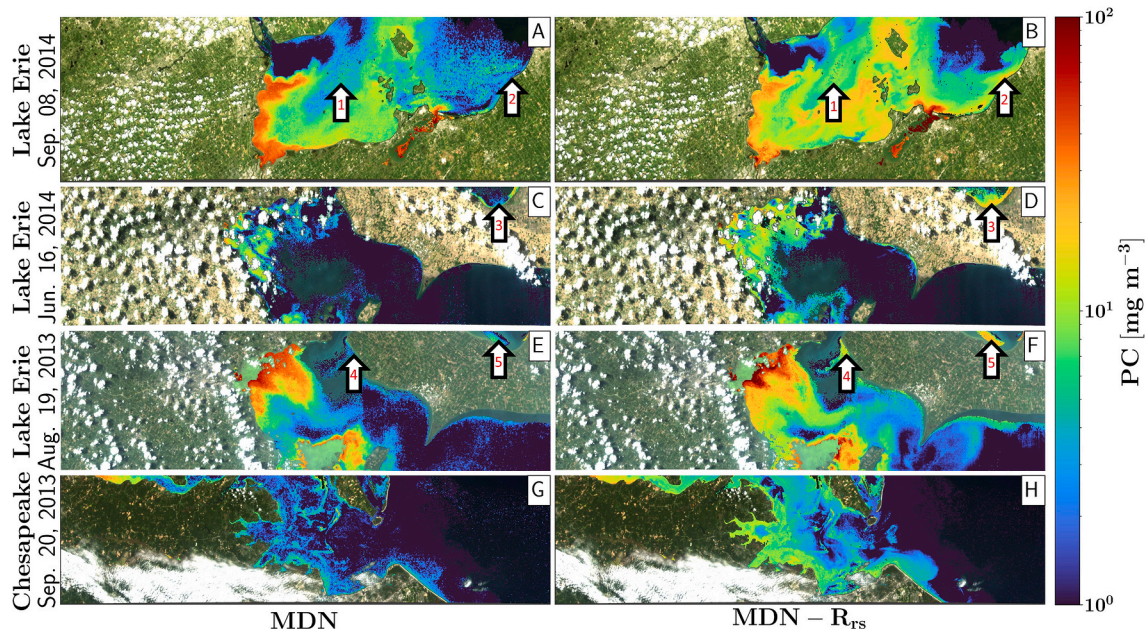


Fig. B1. Product Maps for Lake Erie and the Chesapeake Bay created using MDN with only  $R_{rs}$  as input (which produced results shown in Fig. 10, panel B). Notable differences with the demonstration MDN occur in low estimated regions from the MDN with LH and BR, Fig. 10, particularly for the Chesapeake Bay (Panel G), the Central Basin of Lake Erie, and Lake St. Clair (upper right section of panels C/E).

## References

- Babin, M., Stramski, D., 2004. Variations in the mass-specific absorption coefficient of mineral particles suspended in water. *Limnol. Oceanogr.* 49, 756–767.
- Balasubramanian, S.V., Pahlevan, N., Smith, B., Binding, C., Schalles, J., Loisel, H., Boss, E., 2020. Robust algorithm for estimating total suspended solids (TSS) in inland and nearshore coastal waters. *Remote Sens. Environ.* 246 (111768), 1–16.
- Bastien, C., Cadin, R., Veilleux, E., Deblois, C., Warren, A., Laurion, I., 2011. Performance evaluation of phycocyanin probes for the monitoring of cyanobacteria. *J. Environ. Monit.* 13, 110–118.
- Becker, A., Meister, A., Wilhelm, C., 2002. Flow cytometric discrimination of various phycobilin-containing phytoplankton groups in a hypertrophic reservoir. *Cytometry* 48, 45–57.
- Berk, A., Bernstein, L.S., 1989. MODTRAN: A Moderate Resolution Model for LOWTRAN7. Report GL-TR-89-0122. US Air Force Geophysical Laboratory, Hanscom, Massachusetts.
- Binding, C.E., Greenberg, T.A., McCullough, G., Watson, S.B., Page, E., 2018. An analysis of satellite-derived chlorophyll and algal bloom indices on Lake Winnipeg. *J. Great Lakes Res.* 44 (3), 436–446.
- Binding, C.E., Zastepa, A., Zeng, C., 2019. The impact of phytoplankton community composition on optical properties and satellite observations of the 2017 western Lake Erie algal bloom. *J. Great Lakes Res.* 45, 573–586.
- Binding, C.E., Pizzolato, L., Zeng, C., 2021. EOLakeWatch; delivering a comprehensive suite of remote sensing algal bloom indices for enhanced monitoring of Canadian eutrophic lakes. *Ecol. Indic.* 121 (106999), 1–13.
- Bishop, C.M., 1994. Mixture Density Networks: NCRG/94/004. Aston University, Aston, UK.
- Busetto, L., Ranghetti, L., . Prismaread: A Tool for Facilitating Access and Analysis of PRISMA L1/L2 Hyperspectral Imagery v1.0.0. URL: <https://lbusett.github.io/prismaread/> <https://doi.org/10.5281/zenodo.4019081>.
- Casey, K.A., Rousseaux, C.S., Gregg, W.W., Boss, E., Chase, A.P., Craig, S.E., Maritorena, S., 2019. A global compilation of in situ aquatic high spectral resolution inherent and apparent optical property data for remote sensing applications. *Earth Syst. Data Sci.* 1–29.
- Clark, J.M., Schaeffer, B.A., Darling, J.A., Urquhart, E.A., Johnston, J.M., Ignatius, A.R., Stumpf, R.P., 2017. Satellite Monitoring of cyanobacterial harmful algal bloom frequency in recreational waters and drinking water sources. *Ecol. Indic.* 80, 84–95.
- Cogliati, S., Sarti, F., Chiarantini, L., Cosi, M., Lorusso, R., Lopinto, E., Miglietta, F., Genesio, L., Guanter, L., Damm, A., 2021. The PRISMA imaging spectroscopy mission: overview and first performance analysis. *Remote Sens. Environ.* 262, 112499 accepted.
- Dekker, A.G., 1993. Detection of Optical water quality parameters for eutrophic waters by high resolution remote sensing. *Free Universit.* 1–222.
- Dierssen, H., Bracher, A., Brando, V.E., Loisel, H., Ruddick, K.G., 2020. Data needs for hyperspectral detection of algal diversity across the globe. *Oceanography* 33.
- Doxaran, D., Ruddick, K., McKee, D., Gentili, B., Tailiez, D., Chami, M., Babin, M., 2009. Spectral variations of light scattering by marine particles in coastal waters, from visible to near infrared. *Limnol. Oceanogr.* 54 (4), 1257–1271.
- Ficek, D., Kaczmarek, S., Stoń-Egiert, J., Woźniak, B., Majchrowski, R., Dera, J., 2004. Spectra of light absorption by phytoplankton pigments in the Baltic; conclusions to be drawn from a Gaussian analysis of empirical data. *Oceanologia* 46 (4), 533–555.
- Francy, D.S., Brady, A.M., Ecker, C.D., Graham, J.L., Stelzer, E.A., Struffolino, P., Loftin, K.A., 2016. Estimating microcystin levels at recreational sites in western Lake Erie and Ohio. *Harmful Algae* 58, 23–34.
- Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *Ann. Stat.* 29 (5), 1189–1232.

- Frouin, R.J., Franz, B.A., Ibrahim, A., Knobelspiese, K., Ahmad, Z., Cairns, B., Chowdhary, J., Dierssen, H.M., Tan, J., Dubovik, O., 2019. Atmospheric correction of satellite ocean-color imagery during the PACE era. *Front. Earth Sci.* 7, 145.
- Gordon, H.R., Wang, M., 1994. Retrieval of water-leaving radiance and aerosol optical thickness over the oceans with SeaWiFS: a preliminary algorithm. *Appl. Opt.* 33 (3), 443–452.
- Harding Jr., L.W., Mallonee, M.E., Perry, E.S., Miller, W.D., Adolf, J.E., Gallegos, C.L., Paerl, H.W., 2016. Variable climatic conditions dominate recent phytoplankton dynamics in Chesapeake Bay. *Scient. Rep.* 6, 1–16.
- Health Canada, 2020. Guidelines for Canadian Recreational Water Quality Cyanobacteria and their Toxins.
- Horváth, H., Kovács, A.W., Riddick, C., Présing, M., 2013. Extraction methods for phycocyanin determination in freshwater filamentous cyanobacteria and their application in a shallow lake. *Eur. J. Phycol.* 48 (3), 278–286.
- Hunter, P.D., Tyler, A.N., Gilvear, D.J., Willby, N.J., 2009. Using remote sensing to aid the assessment of human health risks from blooms of potentially toxic cyanobacteria. *Environ. Sci. Technol.* 43 (7), 2627–2633.
- Hunter, P.D., Tyler, A.N., Carvalho, L., Codd, G.A., Maberly, S.C., 2010. Hyperspectral remote sensing of cyanobacterial pigments as indicators for cell populations and toxins in eutrophic lakes. *Remote Sens. Environ.* 114, 2705–2718.
- Ibrahim, A., Franz, B., Ahmad, Z., Healy, R., Knobelspiese, K., Gao, B.-C., Proctor, C., Zhai, P.-W., 2018. Atmospheric correction for hyperspectral ocean color retrieval with application to the Hyperspectral Imager for the Coastal Ocean (HICO). *Remote Sens. Environ.* 204, 60–75.
- Jin, C., Mesquita, M.M., Deglinc, J.L., Emelko, M.B., Wong, A., 2018. Quantification of cyanobacterial cells via a novel imaging-driven technique with an integrated fluorescence signature. *Scient. Rep.* 8, 1–12.
- Kasnik, J.M., Chislock, M.F., Wilson, A.E., 2015. Benchtop fluorometry of phycocyanin as a rapid approach for estimating cyanobacterial biovolume. *J. Plankton Res.* 37 (1), 248–257.
- Keith, D.J., Schaeffer, B.A., Lunetta, R.S., Gould Jr., R.W., Rocha, K., Cobb, D.J., 2014. Remote sensing of selected water-quality indicators with the hyperspectral imager for the coastal ocean (HICO) sensor. *Int. J. Remote Sens.* 35 (9), 2927–2962.
- Kravitz, J., Matthews, M., Lain, L., Fawcett, S., Bernard, S., 2021. Potential for high fidelity global mapping of common inland water quality products at high spatial and temporal resolutions based on a synthetic data and machine learning approach. *Front. Environ. Sci.* 9.
- Kudela, R.M., Palacios, S.L., Austerberry, D.C., Accorsi, E.K., Guild, L.S., Torres-Perez, J., 2015. Application of hyperspectral remote sensing to cyanobacterial blooms in inland waters. *Remote Sens. Environ.* 167, 196–205.
- Kutser, T., 2004. Quantitative detection of chlorophyll in cyanobacterial blooms by satellite remote sensing. *Limnol. Oceanogr.* 49 (6), 2179–2189.
- Kwon, Y.S., JongCheol, P., Kwon, Y.-H., Duan, H., Cho, K.H., Park, Y., 2020. Drone-based hyperspectral remote sensing of cyanobacteria using vertical cumulative pigment concentration in a deep reservoir. *Remote Sens. Environ.* 236 (111517), 1–16.
- Le, C., Li, Y., Zha, Y., Wang, Q., Zhang, H., Yin, B., 2011. Remote sensing of phycocyanin pigment in highly turbid inland waters in Lake Taihu, China. *Int. J. Remote Sens.* 32, 8253–8269.
- Li, L., Sep, 25, 2020. Cyanobacteria in inland waters: remote sensing. In: *Fresh Water and Watersheds*, 2nd edition.
- Li, L., Song, K., 2017. Bio-optical modeling of phycocyanin. In: *Bio-optical Modeling and Remote Sensing of Inland Waters*. Elsevier, pp. 233–262.
- Li, L.H., Li, L., Shi, K., Li, Z.C., Song, K.S., 2012. A semi-analytical algorithm for remote estimation of phycocyanin in inland waters. *Sci. Total Environ.* 435&436, 141–150.
- Li, L., Li, L., Song, K., 2015. Remote sensing of freshwater cyanobacteria: an extended IOP Inversion Model of Inland Waters (IIMIWI) for partitioning absorption coefficient and estimating phycocyanin. *Remote Sens. Environ.* 157, 9–23.
- Liu, G., Simis, S.G.H., Li, L., Wang, Q., Li, Y.M., Song, K.S., Lyu, H., Zheng, Z., Shi, K., 2017. A four-band semi-analytical model for estimating phycocyanin in inland waters from simulated MERIS and OLCI data. *IEEE Trans. Geosci. Remote Sens.* 99, 1–12.
- Liu, Q., Rowe, M.D., Anderson, E.J., Stow, C.A., Stumpf, R.P., Johengen, T.H., 2020. Probabilistic forecast of microcystin toxin using satellite remote sensing, in situ observations and numerical modeling. *Environ. Model. Softw.* 104705, 1–12.
- Lucke, R.L., Corson, M., McGlothlin, N.R., Butcher, S.D., Wood, D.L., Korwan, D.R., Chen, D.T., 2011. Hyperspectral Imager for the Coastal Ocean: instrument description and first images. *Appl. Opt.* 50 (11), 1501–1516.
- Ludovisi, A., Gaino, E., 2010. Meteorological and water quality changes in Lake Trasimeno (Umbria, Italy) during the last fifty years. *J. Limnol.* 69 (1), 174–188.
- Lunetta, R.S., Schaeffer, B.A., Stumpf, R.P., Keith, D., Jacobs, S.A., Murphy, M.S., 2015. Evaluation of cyanobacteria cell count detection derived from MERIS imagery across the eastern USA. *Remote Sens. Environ.* 157, 24–34.
- Lyu, H., Wang, Q., Wu, C., Zhu, L., Yin, B., Li, Y.M., Huang, J.Z., 2013. Retrieval of phycocyanin concentration from remote-sensing reflectance using a semi-analytic model in eutrophic lakes. *Ecol. Inform.* 18, 178–187.
- Matthews, M., 2020. Data for: distinguishing cyanobacteria from algae in optically complex inland waters using a radiative transfer inversion algorithm. *Mendeley Data VI*. <https://doi.org/10.17632/msf535b5cyc.1>
- Matthews, M.W., Odermatt, D., 2015. Improved algorithm for routine monitoring of cyanobacteria and eutrophication in inland and near-coastal waters. *Remote Sens. Environ.* 156, 374–382.
- Matthews, M.W., Bernard, S., Robertson, L., 2012. An algorithm for detecting trophic status (chlorophyll-a), cyanobacterial-dominance, surface scums and floating vegetation in inland and coastal waters. *Remote Sens. Environ.* 124, 637–652.
- Mishra, S., Mishra, D.R., 2014. A novel remote sensing algorithm to quantify phycocyanin in cyanobacterial algal blooms. *Environ. Res. Lett.* 9 (114003), 1–9.
- Mishra, S., Mishra, D.R., Schluchter, W.M., 2009. A novel algorithm for predicting phycocyanin concentrations in cyanobacteria: a proximal hyperspectral remote sensing approach. *Remote Sens.* 1, 758–775.
- Mishra, S., Mishra, D.R., Lee, Z., Tucker, C.S., 2013. Quantifying cyanobacterial phycocyanin concentration in turbid productive waters: a quasi-analytical approach. *Remote Sens. Environ.* 133, 141–151.
- Mishra, S., Stumpf, R.P., Schaeffer, B.A., Werdell, J.P., Loftin, K.A., Meredith, A., 2019. Measurement of cyanobacterial bloom magnitude using satellite remote sensing. *Sci. Rep.* 9 (18310), 1–17.
- Mobley, C.D., Werdell, J., Franz, B., Ahmad, Z., Bailey, S., 2016. Atmospheric correction for satellite ocean color radiometry. NASA/TM-2016-217551, GSFC-E-DAATN35509.
- Moore, T.S., Mouw, C.B., Sullivan, J.M., Twardowski, M.S., Burtner, A.M., Ciochetto, A. B., Weidemann, A., 2017. Bio-optical properties of cyanobacteria blooms in Western Lake Erie. *Front. Mar. Sci.* 4, 1–20.
- Mobley, C.D., 1999. Estimation of the remote-sensing reflectance from above-surface measurements. *Appl. Opt.* 38, 7442–7455.
- Moore, T.S., Churnside, J.H., Sullivan, J.M., Twardowski, M.S., Nayak, A.R., McFarland, M.N., Ruberg, S.A., 2019. Vertical distributions of blooming cyanobacteria populations in freshwater lake from LIDAR observations. *Remote Sens. Environ.* 225, 347–367.
- Morley, S.K., Brito, T.V., Welling, D.T., 2018. Measures of model performance based on the log accuracy ratio. *Space Weather* 16, 69–88.
- Moses, W.J., Bowles, J.H., Lucke, R.L., Corsor, M.R., 2012. Impact of signal-to-noise ratio in a hyperspectral sensor on the accuracy of biophysical parameter estimation in case II waters. *Opt. Express* 20, 4309–4330.
- NASA, 2021. *Ocean Color Web*. <https://oceancolor.gsfc.nasa.gov/>.
- Ogashawara, I., 2020. Determination of phycocyanin from space - a bibliometric analysis. *Remote Sens.* 12 (567), 1–16.
- Ogashawara, I., Mishra, D.R., Mishra, S., Curtarelli, M.P., Stech, J.L., 2013, July. A performance review of reflectance based algorithms for predicting phycocyanin concentrations in inland waters. *Remote Sens.* 5, 4774–4798.
- Ogashawara, I., Li, L., 2019. Removal of chlorophyll spectral interference to improve phycocyanin estimation from space. *Remote Sens.* 11 (1764), 1–19.
- O'Shea, R., Laney, S., Lee, Z., 2020. Evaluation of glint correction approaches for fine-scale ocean color measurements by lightweight hyperspectral imaging spectrometers. *Appl. Opt.* 59 (7), B18–B34.
- Pahlevan, N., Smith, B., Schalles, J., Binding, C., Cao, Z., Ma, R., Stumpf, R., 2020, April. Seamless retrievals of chlorophyll-a from Sentinel-2 (MSI) and Sentinel-3 (OLCI) in inland and coastal waters: a machine-learning approach. *Remote Sens. Environ.* 240 (111604), 1–21.
- Pahlevan, N., Mangin, A., Balasubramanian, S.V., Smith, B., Alikas, K., Arai, K., Warren, M., 2021a. ACIX-aqua: a global assessment of atmospheric correction methods for Landsat-8 and Sentinel-2 over lakes, rivers, and coastal waters. *Remote Sens. Environ.* 258 (112366), 1–22.
- Pahlevan, N., Smith, B.B., Gurlin, D., Li, L., Bresciani, M., Giardino, C., 2021b. Hyperspectral retrievals of phytoplankton absorption and chlorophyll-a in inland and nearshore coastal waters. *Remote Sens. Environ.* 253 (112200), 1–15.
- Pyo, J., Duan, H., Ligaray, M., Kim, M., Baek, S., Kwon, Y.S., Cho, K.H., 2020. An integrative remote sensing application of stacked autoencoder for atmospheric correction and cyanobacteria estimation using hyperspectral imagery. *Remote Sens.* 12 (1073), 1–23.
- Randolph, K., Wilson, J., Tedesco, L., Li, L., Pascual, D.L., Soyeux, E., 2008. Hyperspectral remote sensing of cyanobacteria in turbid productive water using optically active pigments, chlorophyll a and phycocyanin. *Remote Sens. Environ.* 112, 4009–4019.
- Richter, R., Schlpfer, D., 2002. Geo-atmospheric processing of airborne imaging spectrometry data part.2: atmospheric/topographic correction. *Int. J. Remote Sens.* 23 (13), 2631–2649.
- Rinta-Kanto, J.M., Konopko, E.A., DeBruyn, J.M., Bourbonniere, R.A., Boyer, G.L., Wilhelm, S.W., 2009. Lake Erie Microcystis: relationship between microcystin production, dynamics of genotypes and environmental parameters in a large lake. *Harmful Algae* 8, 665–673.
- Ruiz-Verdú, A., Simis, S.G., de Hoyos, C., Gons, H.J., Peña-Martínez, R., 2008. An evaluation of algorithms for the remote sensing of cyanobacterial biomass. *Remote Sens. Environ.* 112, 3996–4008.
- Sarada, R., Pillai, M.G., Ravishankar, G.A., 1999. Phycocyanin from *Spirulina* sp: influence of processing of biomass on phycocyanin yield, analysis of efficacy of extraction methods and stability studies on phycocyanin. *Process Biochem.* 34, 795–801.
- Sathyendranath, S., Lazara, L., Prieur, L., 1987. Variations in the spectral values of specific absorption of phytoplankton. *Limnol. Oceanogr.* 32 (2), 403–415.
- Schaeffer, B.A., Conmy, R.N., Duffy, A.E., Aukamp, J., Yates, D.F., Craven, G., 2015. Northern Gulf of Mexico estuarine coloured dissolved organic matter derived from MODIS data. *Int. J. Remote Sens.* 36 (8), 2219–2237.
- Schaeffer, B.A., Bailey, S.W., Conmy, R.N., Galvin, M., Inatius, A.R., Johnston, J.M., Wolfe, K., 2018. Mobile device application for monitoring cyanobacteria harmful algal blooms using Sentinel-3 satellite Ocean and Land Color Instruments. *Environ. Model. Softw.* 109, 93–103.
- Schalles, J., Yacobi, Y.Z., 2000, February. Remote detection and seasonal patterns of phycocyanin, carotenoid and chlorophyll pigments in eutrophic waters. *Arch. Hydrobiol. Spec. Issues Advanc. Limnol.* 153–169.
- Simis, S.G., Peters, S.W., Gons, H.J., 2005. Remote sensing of the cyanobacterial pigment phycocyanin in turbid inland water. *Limnol. Oceanogr.* 50 (1), 237–245.

- Simis, S.G., Ruiz-Verdú, A., Domínguez-Gómez, J.A., Peña-Martínez, R., Peters, S.W., Gons, H.J., 2007. Influence of phytoplankton pigment composition on remote sensing of cyanobacterial biomass. *Remote Sens. Environ.* 106, 414–427.
- Smith, B., Pahlevan, N., Schalles, J., Ruberg, S., Errera, R., Ma, R., Kangaro, K., 2021. A chlorophyll-a algorithm for Landsat-8 based on mixture density networks. *Front. Remote Sens.* 1 (623678), 1–14.
- Song, K., Li, L., Li, S., Tedesco, L., Hall, B., Li, Z., 2012. Hyperspectral retrieval of phycocyanin in potable water sources using genetic algorithm–partial least squares (GA–PLS) modeling. *Int. J. Appl. Earth Obs. Geoinf.* 18, 368–385.
- Song, K., Li, L., Tedesco, L.P., Li, S., Hall, B.E., Du, J., 2014. Remote quantification of phycocyanin in potable water sources through an adaptive model. *ISPRS J. Photogramm. Remote Sens.* 95, 68–80.
- Spyrakos, E., O'Donnell, R., Hunter, P.D., Miller, C., Scott, M., Simis, S.G., Tyler, A.N., 2021. Optical types of inland and coastal waters. *Limnol. Oceanogr.* 63 (2), 846–870.
- Stumpf, R.P., Davis, T.W., Wynna, T.T., Graham, J.L., Loftin, K.A., Johengen, T.H., Gossiaux, Duane, Palladino, D., Burtner, A., 2016. Challenges for mapping cyanotoxin patterns from remote sensing of cyanobacteria. *Harmful Algae* 54, 160–173.
- Sun, D., Li, Y., Wang, Q., Le, C., Lv, H., Huang, C., Gong, S., 2012. A novel support vector regression model to estimate the phycocyanin concentration in turbid inland waters from hyperspectral reflectance. *Hydrobiologia* 680, 199–217.
- U.S. Environmental Protection Agency Office of Water (4304T) Health and Ecological Criteria Division, 2019. Recommended Human Health Recreational Ambient Water Quality Criteria or Swimming Advisories for Microcystins and Cylindrospermopsin.
- Vansteenkoven, D., Ruddick, K., Cattijsee, A., Vanhellemont, Q., Beck, M., 2019. The Pan-and-Tilt Hyperspectral Radiometer System (PANTHYR) for autonomous satellite validation measurements – prototype design and testing. *Remote Sens.* 11 (1360).
- World Health Organization, 2003. Guidelines for Safe Recreational Water Environments: Coastal and Fresh Waters.
- Wynne, T.T., Stumpf, R.P., Tomlinson, M.C., Dyble, J., 2010. Characterizing a cyanobacterial bloom in western Lake Erie using satellite imagery and meteorological data. *Limnol. Oceanogr.* 55 (5), 2025–2036.
- Zemlys, P., Ferrarin, C., Umgiesser, G., Gulbinskas, S., Bellafiore, D., 2013. Investigation of saline water intrusions into the Curonian Lagoon (Lithuania) and two-layer flow in the Klaipėda Strait using finite element hydrodynamic model. *Ocean Sci.* 9, 573–584.
- Zimba, P.V., 2012. An Improved phycobilin extraction method. *Harmful Algae* 17, 35–39.