

Extracting information from ocean color

B. B. Cael¹, Kelsey M Bisson², Emmanuel S Boss³, and Zachary K Erickson⁴

¹National Oceanography Centre

²Oregon State University

³University of Maine

⁴NASA Goddard Space Flight Center

December 8, 2022

Abstract

Products derived from remote sensing reflectances ($R_{rs}(\lambda)$), e.g. chlorophyll, phytoplankton carbon, euphotic depth, or particle size, are widely used in oceanography. Problematically, $R_{rs}(\lambda)$ may have fewer degrees of freedom (DoF) than measured wavebands or derived products. A global sea surface hyperspectral $R_{rs}(\lambda)$ dataset has DoF=4. MODIS-like multispectral equivalent data also have DoF=4, while their SeaWiFS equivalent has DoF=3. Both multispectral-equivalent datasets predict individual hyperspectral wavelengths' $R_{rs}(\lambda)$ within nominal uncertainties. Remotely sensed climatological multispectral $R_{rs}(\lambda)$ have DoF=2, as information is lost by atmospheric correction, shifting to larger spatiotemporal scales, and/or more open-ocean measurements, but suites of $R_{rs}(\lambda)$ -derived products have DoF=1. These results suggest that remote sensing products based on existing satellites' $R_{rs}(\lambda)$ are not independent and should not be treated as such, that existing $R_{rs}(\lambda)$ measurements hold unutilized information, and that future multi- or especially hyper-spectral algorithms must rigorously consider correlations between $R_{rs}(\lambda)$ wavebands.

Extracting information from ocean color

B. B. Cael¹, Kelsey Bisson², Emmanuel Boss³, and Zachary K. Erickson⁴

¹National Oceanography Centre, Southampton, UK

²Oregon State University, Corvallis, OR, USA

³University of Maine, Orono, ME, USA

⁴NOAA Pacific Marine Environmental Laboratory, Seattle, WA, USA

Key Points:

- In situ hyperspectral $R_{rs}(400-700\text{nm})$ have 4 degrees of freedom & are predicted within uncertainties by MODIS & SeaWiFS wavebands.
- Degrees of freedom are lost upscaling to global satellite climatologies and again to $R_{rs}(\lambda)$ -derived products like chlorophyll.
- Information exists in satellite $R_{rs}(\lambda)$ that's underutilized by products' algorithms. Future algorithms must consider correlations carefully.

Corresponding author: B. B. Cael, cael@noc.ac.uk

Abstract

Products derived from remote sensing reflectances ($R_{rs}(\lambda)$), e.g. chlorophyll, phytoplankton carbon, euphotic depth, or particle size, are widely used in oceanography. Problematically, $R_{rs}(\lambda)$ may have fewer degrees of freedom (DoF) than measured wavebands or derived products. A global sea surface hyperspectral $R_{rs}(\lambda)$ dataset has DoF=4. MODIS-like multispectral equivalent data also have DoF=4, while their SeaWiFS equivalent has DoF=3. Both multispectral-equivalent datasets predict individual hyperspectral wavelengths' $R_{rs}(\lambda)$ within nominal uncertainties. Remotely sensed climatological multispectral $R_{rs}(\lambda)$ have DoF=2, as information is lost by atmospheric correction, shifting to larger spatiotemporal scales, and/or more open-ocean measurements, but suites of $R_{rs}(\lambda)$ -derived products have DoF=1. These results suggest that remote sensing products based on existing satellites' $R_{rs}(\lambda)$ are not independent and should not be treated as such, that existing $R_{rs}(\lambda)$ measurements hold unutilized information, and that future multi- or especially hyper-spectral algorithms must rigorously consider correlations between $R_{rs}(\lambda)$ wavebands.

Plain Language Summary

The reflectance of sunlight from the ocean can be observed from satellites and is used to derive many biologically-relevant parameters, such as the concentration of chlorophyll in the upper ocean. Reflectances are currently observed at about ten different wavelengths, but this will soon be expanded to hundreds with the upcoming launch of a new ocean color satellite, PACE, in early 2024. Many new algorithms are being proposed to make use of the wealth of ocean color data which will be provided. However, there are strong correlations between reflectances at different wavelengths; these correlations mean there will be far fewer products that can be independently derived than there will be reflectance wavelengths observed. Here we use a ship-based measurements similar to what will be provided from PACE to suggest that, on a global scale, only a few independent variables can be calculated from hundreds of reflectance wavelengths. Current and past satellites provide a similar amount of independent data to what is projected from PACE. We then show that, on a global scale, a set of six derived parameters only contains one independent piece of information, suggesting that more information exists in ocean color data than is being currently used.

1 Introduction

Ocean color satellites have revolutionized the study of ocean ecology and biogeochemistry in recent decades by providing a near-continuous global picture of surface ocean properties (Hovis et al., 1980; O'Reilly et al., 1998). Satellites measure the spectral radiance emanating from the ocean and atmosphere. Remote sensing reflectance ($R_{rs}(\lambda)$) is obtained following the removal of the contribution of atmospheric and surface effects and normalization to downwelling solar irradiance. Algorithms have been developed to estimate numerous biogeochemically-relevant surface variables from $R_{rs}(\lambda)$, such as chlorophyll concentration (Chl, [$\mu\text{g/L}$]) (O'Reilly et al., 1998; Hu et al., 2012), the spectral slope of the particle size distribution (ξ) (Kostadinov et al., 2009), the concentrations of phytoplankton and particulate organic and inorganic carbon (C_{phyto} , POC, and PIC, [$\mu\text{g/L}$]) (Graff et al., 2015; Evers-King et al., 2017; Mitchell et al., 2017), euphotic layer depth (Z_{eu} [m]) (Lee et al., 2007), and, using additional input variables, net primary production (NPP, [$\text{mg/m}^2\text{d}$]) (Behrenfeld & Falkowski, 1997; Silsbe et al., 2016; Westberry et al., 2008). Such products are used in a wide variety of applications, such as validation of complex ocean ecosystem and biogeochemistry models (Dutkiewicz et al., 2020; Cael et al., 2021) or as inputs for simpler models that predict other variables such as vertical particulate organic carbon fluxes from ocean color (Siegel et al., 2014; Cael et al., 2017; DeVries & Weber, 2017; Nowicki et al., 2022; Bisson et al., 2020).

Existing $R_{rs}(\lambda)$ data are multispectral, meaning they are measured within several individually determined wavebands. Derived products generally rely only on a subset of these wavebands and are commonly expressed as functions of band ratios between just two wavelengths (e.g. Hu et al., 2012). Some algorithms attempt to simultaneously estimate multiple products to match the full spectrum of $R_{rs}(\lambda)$; for example, the Generalized Inherent Optical Properties (GIOP) approach (Werdell et al., 2013) uses known and assumed spectral shapes of backscattering and absorption from different optical constituents to estimate the suite of products that best represents the observed $R_{rs}(\lambda)$. However, the most widely used products, such as for Chl and POC, treat all outputs as independent quantities and are fully empirical.

Correlations between $R_{rs}(\lambda)$ at different wavebands can be quite strong (Huot & Antoine, 2016), depending also on the spatiotemporal scales considered (see §3). This presents multiple potential issues for both users and developers of ocean color derived products. If multiple products are used simultaneously and treated as independent when they are in fact not, this can lead to overconfidence in model skill or miscalculation of uncertainties. An unintended consequence of treating satellite products independently within models is a functional limit on model complexity. Adding different (yet correlated) satellite products to a model can result in model output redundancy (Bisson et al., 2020). These issues will only be exacerbated by the hyperspectral resolution of the next generation of ocean color satellites, namely the Plankton, Clouds, Aerosols and Ecosystems (PACE) satellite scheduled to launch January 2024 (Werdell et al., 2019). In addition to the common suite of multispectral products, PACE also plans to move beyond chlorophyll and enable characterizations of phytoplankton communities (e.g. Chase et al., 2017), substantially increasing the number of products available from $R_{rs}(\lambda)$.

The strong correlations among $R_{rs}(\lambda)$ wavelengths can be framed in terms of the degrees of freedom (DoF) of $R_{rs}(\lambda)$ measurements and suites of derived products. DoF represents the effective number of dimensions of a dataset after accounting for correlations and uncertainties between variables and is in essence the number of independent variables in that dataset. It has been shown that the DoF of globally distributed near-surface measured hyperspectral absorption spectra is about five (Cael et al., 2020). This could be considered a possible upper limit for the DoF of satellite-measured $R_{rs}(\lambda)$ given higher uncertainties on satellite measurements – particularly associated with atmospheric correction (Bisson et al., 2021; Cael et al., 2020). The DoF of PACE’s hyperspectral measurements might then be expected to be much lower than the number of wavelengths for which it will measure $R_{rs}(\lambda)$, which will appreciably affect how hyperspectral satellite $R_{rs}(\lambda)$ products should be constructed. For both existing and future satellite $R_{rs}(\lambda)$, in other words, understanding the DoF of $R_{rs}(\lambda)$ measurements and derived products is crucial for appropriate usage and optimal construction of such products.

Here we investigate the DoF of $R_{rs}(\lambda)$. We first find that a global sea surface hyperspectral $R_{rs}(\lambda)$ database has four DoF. Coarsening hyperspectral $R_{rs}(\lambda)$ to their MODIS (Moderate Resolution Imaging Spectrometer) equivalent retains four DoF, though the SeaWiFS (Sea-viewing Wide Field of view Sensor) equivalent only has three DoF. Both of these multispectral equivalents can, however, predict individual hyperspectral $R_{rs}(\lambda)$ wavelengths within nominal uncertainties for satellite sensors. We then consider climatological $R_{rs}(\lambda)$ and derived products. We find that both MODIS-Aqua and SeaWiFS $R_{rs}(\lambda)$ have two DoF at the climatological scale, suggesting that $R_{rs}(\lambda)$ complexity is lost either through atmospheric correction, relatively more inclusion of open-ocean data, or averaging over larger scales in space and time. Suites of derived products, however, only retain one DoF. This latter result suggests that derived products should not be treated as independent by users. We close by discussing the substantial implications these findings have for the construction and use of ocean color products, from both existing and future $R_{rs}(\lambda)$.

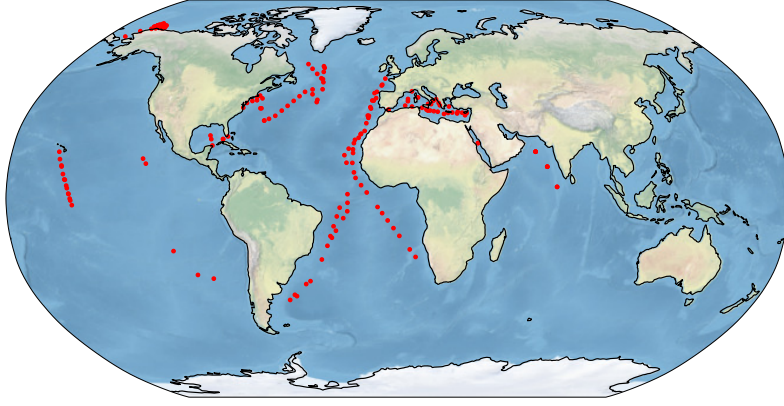


Figure 1. Locations of the 191 stations considered in this study (red dots).

2 Sea surface R_{rs} : hyperspectral versus multispectral

We first analyze a global sea surface hyperspectral $R_{rs}(\lambda)$ dataset to determine its DoF and how the DoF depends on spectral resolution (Chase et al., 2017; Kramer et al., 2022). The dataset includes $R_{rs}(\lambda)$ data at 191 locations at an effective 3.35 nm resolution (Chase et al., 2017) from 400–800 nm, linearly interpolated to 1 nm (Figure 1). We trimmed spectra to 700nm due to the large fraction of missing values $>700\text{nm}$; note that most of the non-empty values $>700\text{nm}$ are zeros and the non-zero-non-empty values, with a median of $<4 \times 10^{-5} \text{ sr}^{-1}$, have very small signal-to-noise ratios. The dataset includes measurements taken from 2004 to 2018 evenly distributed across months of the year, and from all major ocean basins ranging in latitude from 41°S to 74°N . We also compare these data to their MODIS-Aqua and SeaWiFS multispectral equivalents by convoluting the hyperspectral $R_{rs}(\lambda)$ with the MODIS-Aqua and SeaWiFS spectral response functions (available at https://oceancolor.gsfc.nasa.gov/docs/rsr/HMODISA_RSRs.txt and https://oceancolor.gsfc.nasa.gov/docs/rsr/SeaWiFS_RSRs.txt) to generate 10-waveband and 6-waveband datasets which correspond to what each instrument would have measured from the same optical input that the radiometer received when generating the hyperspectral $R_{rs}(\lambda)$ data.

We then apply principal component analysis (PCA) (Wold et al., 1987) to these 301-, 10- and 6-dimensional $R_{rs}(\lambda)$ datasets. PCA is a widely used method to reduce the dimensionality of datasets by identifying orthogonal vectors that explain the most variance in the data. PCA is linear in nature, which may result in an overestimation of effective dimensions by poorly approximating non-linear relationships between variables (e.g. a PCA performed on the pair (x, y) where $y = x^2$ will yield two DoF). Nonlinear generalizations do exist (Weinberger et al., 2004; Scholz et al., 2008), though these are less widely applied due to their additional complexity and computational requirements that make interpretation challenging. One may therefore consider the DoF we report to be upper bounds. We perform a PCA on each $R_{rs}(\lambda)$ dataset, standardizing each first by subtracting from each waveband its mean and then dividing by its standard deviation. This results in a percentage of total variance explained by each component. We use the broken-stick rule to choose the DoF, which states that the DoF is equal to the number of components that explain more variance than would be expected by randomly distributed data; this method was shown to be more consistent than a suite of others in a comparison (Jackson, 1993). These results can be shown visually as a ‘scree’ plot, which plots the percentage of variance explained by each component and for randomly distributed data; the DoF is the number of components with a higher percentage of variance explained than would be expected for randomly distributed data. Our figures also visibly demon-

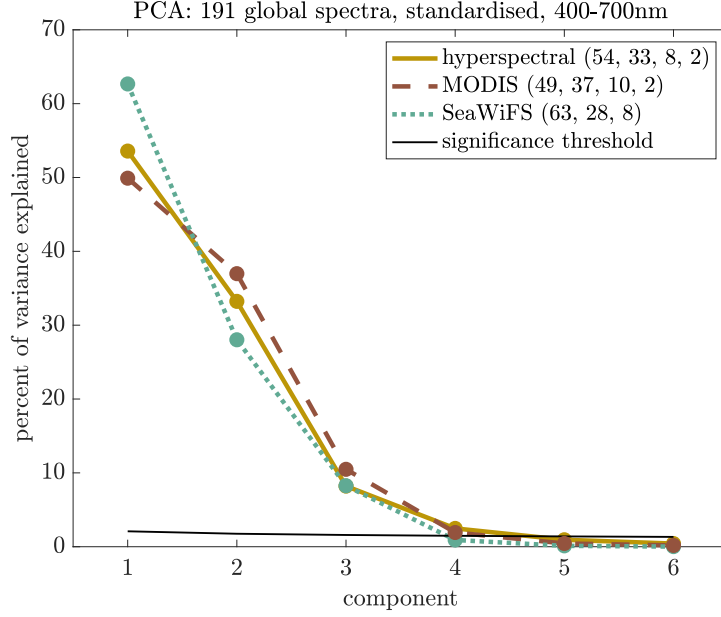


Figure 2. Scree plot of percent variance explained versus component for hyperspectral $R_{rs}(\lambda)$ dataset and MODIS-Aqua and SeaWiFS equivalents calculated from their spectral response functions. Black line indicates broken-stick significance threshold for hyperspectral data; numbers in legend give percent variance explained for each mode above this threshold in each case.

strate that one would get the same results from using the scree plot rule, which states that the DoF is equal to the number of components not sitting on the straight line made by the higher-order components, and was found to consistently capture the correct DoF plus one when the first point on this straight line was included (Jackson, 1993).

PCA analysis reveals that the hyperspectral in situ $R_{rs}(\lambda)$ dataset has four DoF (Figure 2); the first four components explain 54%, 33%, 8%, and 2%, totalling 97%, of the variance. The first four MODIS-Aqua equivalent $R_{rs}(\lambda)$ principal components have very similar percentages of variance explained: 49%, 37%, 10%, and 2%, totalling 99% of the total variance. In contrast, the first three SeaWiFS equivalent $R_{rs}(\lambda)$ principal components explain 63%, 28%, and 8%, totalling 99%, of the variance. This suggests that the hyperspectral $R_{rs}(\lambda)$ have four DoF, or four independent variables within the data, and that these four variables are effectively captured when reducing spectral resolution to the ten MODIS-Aqua wavebands, but not to the six SeaWiFS wavebands.

The ability of coarsened, MODIS-equivalent data to obtain the same number of DoF as the hyperspectral dataset is further supported by predictions of hyperspectral $R_{rs}(\lambda)$ from multispectral equivalents. To illustrate this, for each hyperspectral wavelength we perform a multivariate linear regression of $R_{rs}(\lambda)$ at that wavelength regressed against $R_{rs}(\lambda)$ at each waveband of both the MODIS-Aqua and SeaWiFS equivalent $R_{rs}(\lambda)$. We then calculate the root-mean-square-error (RMSE) of these regressions. For all wavelengths below 578 nm in the SeaWiFS case and 582 nm in the MODIS-Aqua case, the RMSE is smaller – and for many, much smaller – than 5% of the mean $R_{rs}(\lambda)$ at that wavelength, where 5% is a nominal relative uncertainty for satellite $R_{rs}(\lambda)$ (Figure 3). Even for wavelengths greater than this, the RMSE is still very small in absolute terms, $<0.00007 \text{ sr}^{-1}$, far smaller than the nominal 0.0003 sr^{-1} absolute error for 1km-by-1km pixels for PACE (Gordon & Wang, 1994). These small errors in predicting hyperspectral $R_{rs}(\lambda)$ from its multispectral equivalent underscore the extent to which different wavelengths’ $R_{rs}(\lambda)$

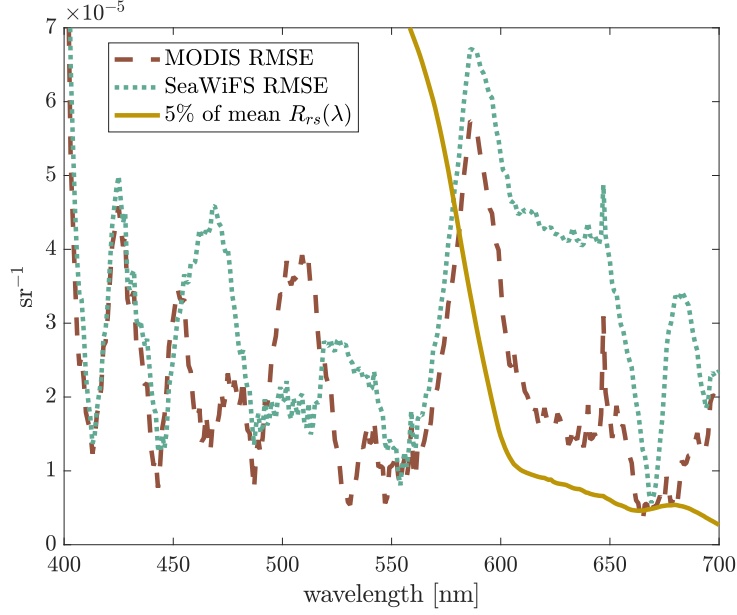


Figure 3. Root-mean-square-error of multivariate linear regressions of each hyperspectral wavelength versus the MODIS-Aqua and SeaWiFS equivalent $R_{rs}(\lambda)$. Solid line is 5% of the mean of each wavelength’s hyperspectral $R_{rs}(\lambda)$.

are correlated and demonstrate the ability of MODIS-Aqua equivalent multispectral $R_{rs}(\lambda)$ to preserve the dimensionality of hyperspectral $R_{rs}(\lambda)$. The fact that SeaWiFS-like $R_{rs}(\lambda)$ can accurately predict hyperspectral $R_{rs}(\lambda)$ to within PACE uncertainties but has fewer DoF than the in situ hyperspectral dataset is a reflection of the lower uncertainty on the in situ dataset than the expected PACE $R_{rs}(\lambda)$, and suggests that PACE $R_{rs}(\lambda)$ may have fewer DoF than the in situ hyperspectral dataset.

We also note that excluding wavelengths 651–700nm affects the DoF numbers presented here but not our conclusions. A choice of an upper limit of 650nm would be based on the fact that for all wavelengths above 648nm, >95% of measurements are below 0.0003 sr^{-1} , the nominal uncertainty of a 1km-by-1km pixel for PACE (Gordon & Wang, 1994). Repeating this analysis over 400–650nm results in hyperspectral and MODIS-Aqua-equivalent $R_{rs}(\lambda)$ data having three DoF, and SeaWiFS-equivalent $R_{rs}(\lambda)$ data having two DoF. This suggests that there is one DoF in the 651–700nm range that is picked up by hyperspectral and multispectral $R_{rs}(\lambda)$ alike; however, the $R_{rs}(\lambda)$ values are small enough (mean and median both $<1.2 \times 10^{-4} \text{ sr}^{-1}$ for all wavelengths 651–700nm) compared to the nominal 1km-by-1km pixel uncertainty $3 \times 10^{-4} \text{ sr}^{-1}$) that this DoF may not be useful for satellite applications, which we are interested in here. This is corroborated by the DoF < 3 in the next section, despite incorporating the full wavebands of both MODIS-Aqua and SeaWiFS. Note that when estimating the MODIS-Aqua- and SeaWiFS-equivalent data from 400–650nm hyperspectral data, the contribution of hyperspectral $R_{rs} > 650\text{nm}$ is not included; while both MODIS-Aqua and SeaWiFS have wavebands centered at $>650\text{nm}$, these wavebands’ spectral response functions are nonzero for some wavelengths in the range 400–650nm, and it is only the influence of these hyperspectral wavelengths on all wavebands that is considered. In other words, $R_{rs}(\lambda)$ is effectively set to zero for all hyperspectral wavelengths $>650\text{nm}$ when calculating the multispectral equivalent datasets in this case.

3 Climatologies: R_{rs} versus products

The analysis in Section 2 is based on instantaneous, local-scale $R_{rs}(\lambda)$ values measured in situ at the sea surface. The power of satellite $R_{rs}(\lambda)$ and derived products, however, lies in their near-continuous global spatial coverage, and many users are primarily interested in climatological data, which is near the coarsest spatial and temporal scales. In this section we therefore analyze climatological $R_{rs}(\lambda)$ and derived products, again via PCA to determine DoF.

We generated a $1^\circ \times 1^\circ$ climatology for each month using $R_{rs}(\lambda)$ data from SeaWiFS spanning 1997–2008, excluding the final 2 years of the mission due to known instrument issues (Siegel et al., 2014), using data downloaded from <https://oceancolor.gsfc.nasa.gov/>. We did the same for MODIS-Aqua, spanning the time period from July 2002 – June 2022. We generated analogous climatologies for derived products from each satellite over the same period and at the same spatial and temporal resolution, namely the extensive (i.e. mass-dependent) variables Chl, C_{phyto} , POC, PIC, and the intensive (i.e. mass-independent) variables Z_{eu} , ξ , the fraction of biovolume in the microplankton size class f_{micro} calculated from ξ as described in (Kostadinov et al., 2009), the particulate backscatter to chlorophyll ratio $b_{bp} : \text{Chl}$, and NPP as estimated by the CAFE (Silsbe et al., 2016) and CbPMv2 (Westberry et al., 2008) models. Chl, POC, and PIC were downloaded from <https://oceancolor.gsfc.nasa.gov/>, as was b_{bp} to calculate C_{phyto} according to (Graff et al., 2015) and $b_{bp} : \text{Chl}$ and the diffuse attenuation coefficient at 490nm to calculate Z_{eu} according to (Lee et al., 2007); SeaWiFS ξ and f_{micro} were derived as in (Kostadinov et al., 2009); and NPP products were downloaded from <http://sites.science.oregonstate.edu/ocean.productivity/index.php>. In total we then have climatologies for MODIS-Aqua, SeaWiFS $R_{rs}(\lambda)$, and ten derived products. We consider the six products Chl, C_{phyto} , POC, PIC, ξ , and Z_{eu} , to be core products and f_{micro} , $b_{bp} : \text{Chl}$, CAFE NPP, and CbPMv2 NPP to be ancillary products as these are either derived from the core products or rely on ancillary data other than $R_{rs}(\lambda)$.

We note that a PCA on the MODIS-Aqua climatologies of $R_{rs}(\lambda)$ and products other than ξ and f_{micro} yields the same results as those for SeaWiFS below, so we focus here only on the SeaWiFS climatologies because ξ and f_{micro} are not readily available for MODIS-Aqua. We find two DoF for SeaWiFS $R_{rs}(\lambda)$, but only one for the products (Figure 4). This result is not sensitive to which combination of products is used – for instance, including all the ancillary products as well still results in one DoF for the products. This result is also not sensitive to log-transformations of the variables that are log-normally (e.g. Chl, POC, PIC, C_{phyto} (Campbell, 1995)) or log-skew-normally (e.g. NPP, (Cael et al., 2018; Cael, 2021)) distributed, or removal of outliers, zeros, or negative values.

That MODIS-Aqua $R_{rs}(\lambda)$ have three DoF for the data in the previous section but two DoF from satellite-derived climatologies suggests that some reduction of complexity of the data occurs via some combination of increased sensor noise relative to ship-based data, atmospheric correction, or averaging over large space and time scales (Scott & Werdell, 2019). (Note (Scott & Werdell, 2019) also point out the difference between averaging $R_{rs}(\lambda)$ versus taking the ratio of averaged water-leaving radiance Lw and downwelling irradiance, which may introduce a slight bias but is unlikely to affect our results here.) Two DoF remain in satellite climatological $R_{rs}(\lambda)$ for both SeaWiFS and MODIS-Aqua, indicating the possibility of generating two independent products from these data. The suite of products tested above, however, has one fewer DoF than the $R_{rs}(\lambda)$. This is likely due to derived products’ appreciable uncertainties and/or strong correlations with chlorophyll. POC, ξ , and Z_{eu} , for instance, have Spearman rank correlations (across all months and 1° grid cells) of >0.9 with Chl. C_{phyto} ’s rank correlation with Chl is still fairly high, at 0.61, and is low largely due to small fluctuations when both are small; a simple spline fit of $\log(C_{phyto})$ against $\log(\text{Chl})$ yields an r^2 of 0.7.

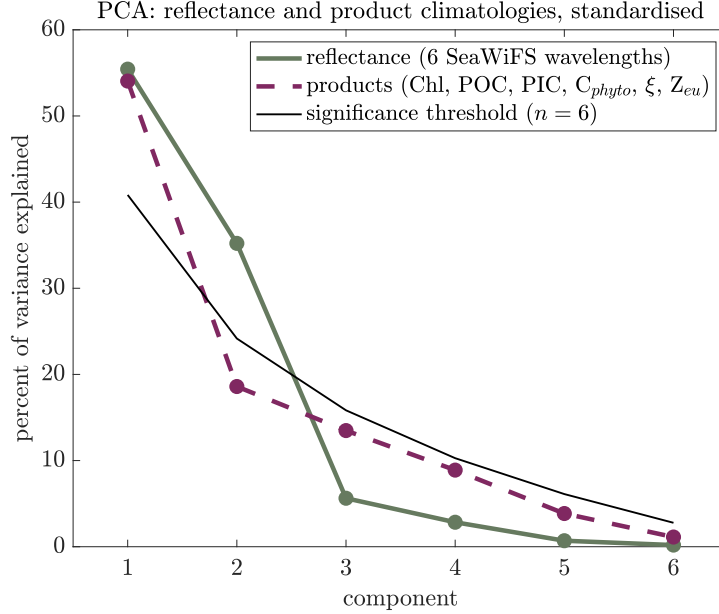


Figure 4. Scree plot of percent variance explained versus component for climatologies of SeaWiFS $R_{rs}(\lambda)$ and of six SeaWiFS- $R_{rs}(\lambda)$ -derived products. Black line indicates broken-stick significance threshold for six-dimensional data.

The exception is PIC, which has a rank correlation with Chl of 0.11. PIC, however, is highly sensitive to small variations in $R_{rs}(\lambda)$ for typical $R_{rs}(\lambda)$ values. To substantiate this, we performed a simple sensitivity analysis with the standard two-band PIC algorithm used by NASA for all but the most optically bright waters (see <https://oceancolor.gsfc.nasa.gov/atbd/pic/>). We calculated PIC for the climatological median $R_{rs}(\lambda)$ at 443 nm and 555 nm and for 5% variations, converting to normalized water-leaving radiance by multiplying by the global mean extraterrestrial solar irradiance. We then perturbed these $R_{rs}(\lambda)$ values with Gaussian noise at the 5% level, corresponding to the nominal uncertainty in $R_{rs}(\lambda)$. This noise at 443 nm results in 68% noise in PIC. By contrast, POC only varies 5% with these 5% variations in $R_{rs}(\lambda)$ at either wavelength. This indicates that in the bulk of cases, satellite-derived PIC is highly uncertain, on the order of 70% (and note the PIC uncertainty will be magnified more when considering documented uncertainties for $R_{rs}(\lambda)$ of 15-40% in some regions (Bisson et al., 2021)). In contrast, for relatively bright waters, the same exercise resulted in PIC variations of <10%, indicating that this algorithm performs well in instances when PIC values are high. Nonetheless, the high sensitivity to typical uncertainty in $R_{rs}(\lambda)$ for median waters explains why we find one DoF for the products even though PIC and Chl are not strongly correlated: derived PIC is noisy most of the time.

These results have two key implications. One is that there is additional information in climatological $R_{rs}(\lambda)$ that is not included in current derived products. This implies that existing products do not utilize the full set of $R_{rs}(\lambda)$ wavelengths. The other implication is that these products are not at all independent, and should not be treated as such when using them simultaneously. In other words, there are more products than there are DoF in satellite data. A numerical ecosystem model that reproduces the satellite-derived climatology of chlorophyll and of the particle size distribution's spectral slope should not be considered to be capturing two independent properties of the Earth sys-

tem. When using satellite products as inputs to other models, these products and their propagated uncertainties must be treated simultaneously rather than independently.

The results presented here are appropriate for global ocean analyses. The open ocean represents the largest area, and is composed primarily of Case 1 waters; that is, waters in which optical variability is dominated by chlorophyll (Morel & Prieur, 1977). In this context, it is in a sense unsurprising that the suite of $R_{rs}(\lambda)$ -derived products produced only one DoF. More optically complex waters, such as coastal regions and inland waters, have optical variability that is influenced by other constituents, such as colored dissolved organic material (CDOM), inorganic particles, and other pigments in addition to chlorophyll (e.g. Brown et al., 2008; Nelson & Siegel, 2013)). Analyses focused on these waters is likely to reveal a higher number of DoF from both $R_{rs}(\lambda)$ and derived products. Indeed, algorithms to derive concentrations of cyanobacteria and suspended particulate (Wang et al., 2016)) or distinguish between different phytoplankton species (Erickson et al., 2020) can be successful in such waters. However, we note that the in situ dataset used here (Figure 1) represents waters with $R_{rs}(\lambda)$ variability similar to that of the ocean as a whole, which can be seen by comparing the variation in $R_{rs}(\lambda)$ at each MODIS-Aqua wavelength from global satellite data with the same satellite data sub-sampled to the locations with in situ measurements (or the closest non-cloudy location). Sub-sampled satellite measurements have similar, and slightly lower, $R_{rs}(\lambda)$ in bluer wavelengths, indicating that the in situ dataset is oriented more towards optically complex coastal waters with substantial CDOM. This suggests that part of the explanation for the drop in DoF in satellite-derived climatologies comes from the fact that the in situ dataset sampled, as a whole, more optically complex waters.

We find that both $R_{rs}(\lambda)$ and variables derived from $R_{rs}(\lambda)$ are highly inter-correlated, reducing the number of DoF associated with each, with a greater reduction in DoF in the derived products. This becomes a problem when products are derived using empirical relationships with $R_{rs}(\lambda)$, and especially when the same wavelengths are used for the products that are assumed to be independent of each other; for example, over much of the ocean PIC, POC, and chlorophyll all are functions only of $R_{rs}(\lambda)$ at two wavelengths, at (or near, depending on the sensor) 443 and 555 nm. Certain combinations of PIC, POC, and chlorophyll, which may occur in the surface ocean, are therefore impossible to find using these algorithms. This is distinct from algorithms, typically called “quasi-analytical” or “semi-empirical”, that use known or assumed spectral shapes for absorption and scattering properties of optical constituents that can be related to the same derived products, such as PIC, POC, and chlorophyll (Werdell et al., 2013). These approaches may result in similar correlations and DoF between derived products, but do not inherently have the same problems as empirical approaches. We note that PACE will have, in addition to hyperspectral visible bands, UV bands from 350nm as well as spectral polarized bands. These measurements are expected to both improve the atmospheric correction (hence reduce the $R_{rs}(\lambda)$ uncertainties) as well as provide their own ocean signals, both of which may increase the DoF compared to those found here. In addition, it has been shown that adding other environmental variables such as SST can add useful information to inversions of phytoplankton groups (e.g. Chase et al., 2022) and thus another approach to increase DoF for inversions by adding relevant and independent information (e.g. mixed-layer depth and nutrients from BGC-Argo assimilating models).

4 Conclusion

The results presented here highlight the high degree of co-dependence between remote sensing reflectances at different wavelengths and of the products derived from these reflectances. For users of products based on existing reflectances, this primarily means factoring in the relationships between products when using more than one simultaneously. For the algorithms that generate these products from existing reflectances, these

results indicate a potential to improve the suite of available products to be more accurate and precise, and to account for the relationships between products and $R_{rs}(\lambda)$ wavebands. One way to do this, consistent with the findings above, would be to derive a single product such as chlorophyll as a function of all reflectance wavebands, derive an anomaly from chlorophyll-based expectations of a secondary product (e.g., phytoplankton community composition, size, POC, PIC, and so forth), then specify all other products explicitly as a function of these two, along the lines of Alvain et al. (2005). Ancillary and independent information can also be added to algorithms, as is currently done with net primary production models via temperature and mixed layer depth.

These findings are most relevant for algorithms that will generate products from hyperspectral reflectances in the future. The small number of degrees of freedom in hyperspectral reflectances indicates that only a few quantities can be estimated independently, and that different wavelengths' reflectances as measured from space will be strongly correlated. Complex algorithms that utilize the full spectrum of reflectance will need to factor in these correlations in order to generate reliable products. Crucially, if more than a few products are generated from hyperspectral reflectances, as is likely the case, such algorithms will also need to output the covariance information encoding the uncertainty in each product and the relationships between them. This can be achieved by some, but not all, machine learning techniques, on which this new generation of algorithms are likely to be based. The fact that hyperspectral reflectances can be predicted within nominal uncertainties by their multispectral equivalents suggests that hyperspectral resolution can play a role in improving ocean color products, but that it will be challenging to provide a substantially finer-grained picture of surface ocean ecosystems and biogeochemical cycles. Here by relying on principal component analysis we have focused on broad, first-order variations, but where such resolution may be most useful and generate novel insights is in investigating outliers and rare events, such as blooms or binning data over coherent features like eddies, where e.g. monospecific signatures may be resolved with spectral precision.

Open Research

Remote sensing data were downloaded from <https://oceancolor.gsfc.nasa.gov/> and <http://sites.science.oregonstate.edu/ocean.productivity/index.php>. All data and code are available at github.com/bbcael/eifoc for review purposes and will be given a Zenodo DOI should this manuscript be accepted for publication.

Acknowledgments

It is a pleasure to thank the many scientists whose collective work has generated the data on which this work relies. Cael acknowledges support from the National Environmental Research Council through Enhancing Climate Observations, Models and Data, and the European Union's Horizon 2020 Research and Innovation Programme under grant agreement No. 820989 (project COMFORT). The work reflects only the authors' view; the European Commission and their executive agency are not responsible for any use that may be made of the information the work contains. KB acknowledges support from NASA grant 80NSSC18K0957. EB acknowledges support from NASA grant 80NSSC20M0203. Cael lead and all other authors assisted with all aspects of this work. The authors have no competing interests to declare. This is PMEL contribution number 5445.

References

- Alvain, S., Moulin, C., Dandonneau, Y., & Breon, F.-M. (2005). Remote sensing of phytoplankton groups in case 1 waters from global seawifs imagery. *Deep Sea Research Part I: Oceanographic Research Papers*, 52(11), 1989–2004.

- Behrenfeld, M. J., & Falkowski, P. G. (1997). Photosynthetic rates derived from satellite-based chlorophyll concentration. *Limnology and oceanography*, 42(1), 1–20.
- Bisson, K., Boss, E., Werdell, P. J., Ibrahim, A., Frouin, R., & Behrenfeld, M. (2021). Seasonal bias in global ocean color observations. *Applied optics*, 60(23), 6978–6988.
- Bisson, K., Siegel, D. A., & DeVries, T. (2020). Diagnosing mechanisms of ocean carbon export in a satellite-based food web model. *Frontiers in Marine Science*, 7, 505.
- Brown, C. A., Huot, Y., Werdell, P. J., Gentili, B., & Claustre, H. (2008). The origin and global distribution of second order variability in satellite ocean color and its potential applications to algorithm development. *Remote Sensing of Environment*, 112(12), 4186–4203. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0034425708002162> doi: <https://doi.org/10.1016/j.rse.2008.06.008>
- Cael, B. (2021). Variability-based constraint on ocean primary production models. *Limnology and Oceanography Letters*, 6(5), 262–269.
- Cael, B., Bisson, K., & Follett, C. L. (2018). Can rates of ocean primary production and biological carbon export be related through their probability distributions? *Global biogeochemical cycles*, 32(6), 954–970.
- Cael, B., Bisson, K., & Follows, M. J. (2017). How have recent temperature changes affected the efficiency of ocean biological carbon export? *Limnology and Oceanography Letters*, 2(4), 113–118.
- Cael, B., Chase, A., & Boss, E. (2020). Information content of absorption spectra and implications for ocean color inversion. *Applied Optics*, 59(13), 3971–3984.
- Cael, B., Dutkiewicz, S., & Henson, S. (2021). Abrupt shifts in 21st-century plankton communities. *Science advances*, 7(44), eabf8593.
- Campbell, J. W. (1995). The lognormal distribution as a model for bio-optical variability in the sea. *Journal of Geophysical Research: Oceans*, 100(C7), 13237–13254.
- Chase, A. P., Boss, E., Cetinić, I., & Slade, W. (2017). Estimation of phytoplankton accessory pigments from hyperspectral reflectance spectra: toward a global algorithm. *Journal of Geophysical Research: Oceans*, 122(12), 9725–9743.
- Chase, A. P., Boss, E. S., Haëntjens, N., Culhane, E., Roesler, C., & Karp-Boss, L. (2022). Plankton imagery data inform satellite-based estimates of diatom carbon. *Geophysical Research Letters*, 49(13), e2022GL098076. Retrieved from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2022GL098076> (e2022GL098076 2022GL098076) doi: <https://doi.org/10.1029/2022GL098076>
- DeVries, T., & Weber, T. (2017). The export and fate of organic matter in the ocean: New constraints from combining satellite and oceanographic tracer observations. *Global Biogeochemical Cycles*, 31(3), 535–555.
- Dutkiewicz, S., Cermenó, P., Jahn, O., Follows, M. J., Hickman, A. E., Taniguchi, D. A., & Ward, B. A. (2020). Dimensions of marine phytoplankton diversity. *Biogeosciences*, 17(3), 609–634.
- Erickson, Z. K., Werdell, P. J., & Cetinić, I. (2020). Bayesian retrieval of optically relevant properties from hyperspectral water-leaving reflectances. *Applied Optics*, 59(23), 6902–6917. doi: 10.1364/AO.398043
- Evers-King, H., Martínez-Vicente, V., Brewin, R. J., Dall’Olmo, G., Hickman, A. E., Jackson, T., ... others (2017). Validation and intercomparison of ocean color algorithms for estimating particulate organic carbon in the oceans. *Frontiers in Marine Science*, 251.
- Gordon, H. R., & Wang, M. (1994). Retrieval of water-leaving radiance and aerosol optical thickness over the oceans with seawifs: a preliminary algorithm. *Applied optics*, 33(3), 443–452.

- Graff, J. R., Westberry, T. K., Milligan, A. J., Brown, M. B., Dall’Olmo, G., van Dongen-Vogels, V., ... Behrenfeld, M. J. (2015). Analytical phytoplankton carbon measurements spanning diverse ecosystems. *Deep Sea Research Part I: Oceanographic Research Papers*, 102, 16–25.
- Hovis, W. A., Clark, D., Anderson, F., Austin, R., Wilson, W., Baker, E., ... others (1980). Nimbus-7 coastal zone color scanner: system description and initial imagery. *Science*, 210(4465), 60–63.
- Hu, C., Lee, Z., & Franz, B. (2012). Chlorophyll algorithms for oligotrophic oceans: A novel approach based on three-band reflectance difference. *Journal of Geophysical Research: Oceans*, 117(C1).
- Huot, Y., & Antoine, D. (2016). Remote sensing reflectance anomalies in the ocean. *Remote Sensing of Environment*, 184, 101–111. doi: 10.1016/j.rse.2016.06.002
- Jackson, D. A. (1993). Stopping rules in principal components analysis: a comparison of heuristical and statistical approaches. *Ecology*, 74(8), 2204–2214.
- Kostadinov, T., Siegel, D., & Maritorena, S. (2009). Retrieval of the particle size distribution from satellite ocean color observations. *Journal of Geophysical Research: Oceans*, 114(C9).
- Kramer, S. J., Siegel, D. A., Maritorena, S., & Catlett, D. (2022). Modeling surface ocean phytoplankton pigments from hyperspectral remote sensing reflectance on global scales. *Remote Sensing of Environment*, 270, 112879.
- Lee, Z., Weidemann, A., Kindle, J., Arnone, R., Carder, K. L., & Davis, C. (2007). Euphotic zone depth: Its derivation and implication to ocean-color remote sensing. *Journal of Geophysical Research: Oceans*, 112(C3).
- Mitchell, C., Hu, C., Bowler, B., Drapeau, D., & Balch, W. (2017). Estimating particulate inorganic carbon concentrations of the global ocean from ocean color measurements using a reflectance difference approach. *Journal of Geophysical Research: Oceans*, 122(11), 8707–8720.
- Morel, A., & Prieur, L. (1977). Analysis of variations in ocean color 1. *Limnology and oceanography*, 22(4), 709–722.
- Nelson, N. B., & Siegel, D. A. (2013). The global distribution and dynamics of chromophoric dissolved organic matter. *Annual review of marine science*, 5, 447–476.
- Nowicki, M., DeVries, T., & Siegel, D. A. (2022). Quantifying the carbon export and sequestration pathways of the ocean’s biological carbon pump. *Global Biogeochemical Cycles*, 36(3), e2021GB007083.
- O’Reilly, J. E., Maritorena, S., Mitchell, B. G., Siegel, D. A., Carder, K. L., Garver, S. A., ... McClain, C. (1998). Ocean color chlorophyll algorithms for seawifs. *Journal of Geophysical Research: Oceans*, 103(C11), 24937–24953.
- Scholz, M., Fraunholz, M., & Selbig, J. (2008). Nonlinear principal component analysis: neural network models and applications. In *Principal manifolds for data visualization and dimension reduction* (pp. 44–67). Springer.
- Scott, J. P., & Werdell, P. J. (2019). Comparing level-2 and level-3 satellite ocean color retrieval validation methodologies. *Optics Express*, 27(21), 30140–30157.
- Siegel, D., Buesseler, K., Doney, S. C., Sailley, S., Behrenfeld, M. J., & Boyd, P. (2014). Global assessment of ocean carbon export by combining satellite observations and food-web models. *Global Biogeochemical Cycles*, 28(3), 181–196.
- Silsbe, G. M., Behrenfeld, M. J., Halsey, K. H., Milligan, A. J., & Westberry, T. K. (2016). The cafe model: A net production model for global ocean phytoplankton. *Global Biogeochemical Cycles*, 30(12), 1756–1777.
- Wang, G., Lee, Z., Mishra, D. R., & Ma, R. (2016). Retrieving absorption coefficients of multiple phytoplankton pigments from hyperspectral remote sensing reflectance measured over cyanobacteria bloom waters. *Limnology and Oceanography: Methods*, 14(7), 432–447. doi: 10.1002/lom3.10102
- Weinberger, K. Q., Sha, F., & Saul, L. K. (2004). Learning a kernel matrix for nonlinear dimensionality reduction. In *Proceedings of the twenty-first international*

- 492 *conference on machine learning* (p. 106).
- 493 Werdell, P. J., Behrenfeld, M. J., Bontempi, P. S., Boss, E., Cairns, B., Davis,
 494 G. T., . . . others (2019). The plankton, aerosol, cloud, ocean ecosystem
 495 mission: status, science, advances. *Bulletin of the American Meteorological*
 496 *Society*, 100(9), 1775–1794.
- 497 Werdell, P. J., Franz, B. A., Bailey, S. W., Feldman, G. C., Boss, E., Brando, V. E.,
 498 . . . others (2013). Generalized ocean color inversion model for retrieving
 499 marine inherent optical properties. *Applied Optics*, 52(10), 2019–2037.
- 500 Westberry, T., Behrenfeld, M., Siegel, D., & Boss, E. (2008). Carbon-based primary
 501 productivity modeling with vertically resolved photoacclimation. *Global Bio-*
 502 *geochemical Cycles*, 22(2).
- 503 Wold, S., Esbensen, K., & Geladi, P. (1987). Principal component analysis. *Chemo-*
 504 *metrics and intelligent laboratory systems*, 2(1-3), 37–52.

Extracting information from ocean color

B. B. Cael¹, Kelsey Bisson², Emmanuel Boss³, and Zachary K. Erickson⁴

¹National Oceanography Centre, Southampton, UK

²Oregon State University, Corvallis, OR, USA

³University of Maine, Orono, ME, USA

⁴NOAA Pacific Marine Environmental Laboratory, Seattle, WA, USA

Key Points:

- In situ hyperspectral $R_{rs}(400-700\text{nm})$ have 4 degrees of freedom & are predicted within uncertainties by MODIS & SeaWiFS wavebands.
- Degrees of freedom are lost upscaling to global satellite climatologies and again to $R_{rs}(\lambda)$ -derived products like chlorophyll.
- Information exists in satellite $R_{rs}(\lambda)$ that's underutilized by products' algorithms. Future algorithms must consider correlations carefully.

Corresponding author: B. B. Cael, cael@noc.ac.uk

Abstract

Products derived from remote sensing reflectances ($R_{rs}(\lambda)$), e.g. chlorophyll, phytoplankton carbon, euphotic depth, or particle size, are widely used in oceanography. Problematically, $R_{rs}(\lambda)$ may have fewer degrees of freedom (DoF) than measured wavebands or derived products. A global sea surface hyperspectral $R_{rs}(\lambda)$ dataset has DoF=4. MODIS-like multispectral equivalent data also have DoF=4, while their SeaWiFS equivalent has DoF=3. Both multispectral-equivalent datasets predict individual hyperspectral wavelengths' $R_{rs}(\lambda)$ within nominal uncertainties. Remotely sensed climatological multispectral $R_{rs}(\lambda)$ have DoF=2, as information is lost by atmospheric correction, shifting to larger spatiotemporal scales, and/or more open-ocean measurements, but suites of $R_{rs}(\lambda)$ -derived products have DoF=1. These results suggest that remote sensing products based on existing satellites' $R_{rs}(\lambda)$ are not independent and should not be treated as such, that existing $R_{rs}(\lambda)$ measurements hold unutilized information, and that future multi- or especially hyper-spectral algorithms must rigorously consider correlations between $R_{rs}(\lambda)$ wavebands.

Plain Language Summary

The reflectance of sunlight from the ocean can be observed from satellites and is used to derive many biologically-relevant parameters, such as the concentration of chlorophyll in the upper ocean. Reflectances are currently observed at about ten different wavelengths, but this will soon be expanded to hundreds with the upcoming launch of a new ocean color satellite, PACE, in early 2024. Many new algorithms are being proposed to make use of the wealth of ocean color data which will be provided. However, there are strong correlations between reflectances at different wavelengths; these correlations mean there will be far fewer products that can be independently derived than there will be reflectance wavelengths observed. Here we use a ship-based measurements similar to what will be provided from PACE to suggest that, on a global scale, only a few independent variables can be calculated from hundreds of reflectance wavelengths. Current and past satellites provide a similar amount of independent data to what is projected from PACE. We then show that, on a global scale, a set of six derived parameters only contains one independent piece of information, suggesting that more information exists in ocean color data than is being currently used.

1 Introduction

Ocean color satellites have revolutionized the study of ocean ecology and biogeochemistry in recent decades by providing a near-continuous global picture of surface ocean properties (Hovis et al., 1980; O'Reilly et al., 1998). Satellites measure the spectral radiance emanating from the ocean and atmosphere. Remote sensing reflectance ($R_{rs}(\lambda)$) is obtained following the removal of the contribution of atmospheric and surface effects and normalization to downwelling solar irradiance. Algorithms have been developed to estimate numerous biogeochemically-relevant surface variables from $R_{rs}(\lambda)$, such as chlorophyll concentration (Chl, [$\mu\text{g/L}$]) (O'Reilly et al., 1998; Hu et al., 2012), the spectral slope of the particle size distribution (ξ) (Kostadinov et al., 2009), the concentrations of phytoplankton and particulate organic and inorganic carbon (C_{phyto} , POC, and PIC, [$\mu\text{g/L}$]) (Graff et al., 2015; Evers-King et al., 2017; Mitchell et al., 2017), euphotic layer depth (Z_{eu} [m]) (Lee et al., 2007), and, using additional input variables, net primary production (NPP, [$\text{mg/m}^2\text{d}$]) (Behrenfeld & Falkowski, 1997; Silsbe et al., 2016; Westberry et al., 2008). Such products are used in a wide variety of applications, such as validation of complex ocean ecosystem and biogeochemistry models (Dutkiewicz et al., 2020; Cael et al., 2021) or as inputs for simpler models that predict other variables such as vertical particulate organic carbon fluxes from ocean color (Siegel et al., 2014; Cael et al., 2017; DeVries & Weber, 2017; Nowicki et al., 2022; Bisson et al., 2020).

Existing $R_{rs}(\lambda)$ data are multispectral, meaning they are measured within several individually determined wavebands. Derived products generally rely only on a subset of these wavebands and are commonly expressed as functions of band ratios between just two wavelengths (e.g. Hu et al., 2012). Some algorithms attempt to simultaneously estimate multiple products to match the full spectrum of $R_{rs}(\lambda)$; for example, the Generalized Inherent Optical Properties (GIOP) approach (Werdell et al., 2013) uses known and assumed spectral shapes of backscattering and absorption from different optical constituents to estimate the suite of products that best represents the observed $R_{rs}(\lambda)$. However, the most widely used products, such as for Chl and POC, treat all outputs as independent quantities and are fully empirical.

Correlations between $R_{rs}(\lambda)$ at different wavebands can be quite strong (Huot & Antoine, 2016), depending also on the spatiotemporal scales considered (see §3). This presents multiple potential issues for both users and developers of ocean color derived products. If multiple products are used simultaneously and treated as independent when they are in fact not, this can lead to overconfidence in model skill or miscalculation of uncertainties. An unintended consequence of treating satellite products independently within models is a functional limit on model complexity. Adding different (yet correlated) satellite products to a model can result in model output redundancy (Bisson et al., 2020). These issues will only be exacerbated by the hyperspectral resolution of the next generation of ocean color satellites, namely the Plankton, Clouds, Aerosols and Ecosystems (PACE) satellite scheduled to launch January 2024 (Werdell et al., 2019). In addition to the common suite of multispectral products, PACE also plans to move beyond chlorophyll and enable characterizations of phytoplankton communities (e.g. Chase et al., 2017), substantially increasing the number of products available from $R_{rs}(\lambda)$.

The strong correlations among $R_{rs}(\lambda)$ wavelengths can be framed in terms of the degrees of freedom (DoF) of $R_{rs}(\lambda)$ measurements and suites of derived products. DoF represents the effective number of dimensions of a dataset after accounting for correlations and uncertainties between variables and is in essence the number of independent variables in that dataset. It has been shown that the DoF of globally distributed near-surface measured hyperspectral absorption spectra is about five (Cael et al., 2020). This could be considered a possible upper limit for the DoF of satellite-measured $R_{rs}(\lambda)$ given higher uncertainties on satellite measurements – particularly associated with atmospheric correction (Bisson et al., 2021; Cael et al., 2020). The DoF of PACE’s hyperspectral measurements might then be expected to be much lower than the number of wavelengths for which it will measure $R_{rs}(\lambda)$, which will appreciably affect how hyperspectral satellite $R_{rs}(\lambda)$ products should be constructed. For both existing and future satellite $R_{rs}(\lambda)$, in other words, understanding the DoF of $R_{rs}(\lambda)$ measurements and derived products is crucial for appropriate usage and optimal construction of such products.

Here we investigate the DoF of $R_{rs}(\lambda)$. We first find that a global sea surface hyperspectral $R_{rs}(\lambda)$ database has four DoF. Coarsening hyperspectral $R_{rs}(\lambda)$ to their MODIS (Moderate Resolution Imaging Spectrometer) equivalent retains four DoF, though the SeaWiFS (Sea-viewing Wide Field of view Sensor) equivalent only has three DoF. Both of these multispectral equivalents can, however, predict individual hyperspectral $R_{rs}(\lambda)$ wavelengths within nominal uncertainties for satellite sensors. We then consider climatological $R_{rs}(\lambda)$ and derived products. We find that both MODIS-Aqua and SeaWiFS $R_{rs}(\lambda)$ have two DoF at the climatological scale, suggesting that $R_{rs}(\lambda)$ complexity is lost either through atmospheric correction, relatively more inclusion of open-ocean data, or averaging over larger scales in space and time. Suites of derived products, however, only retain one DoF. This latter result suggests that derived products should not be treated as independent by users. We close by discussing the substantial implications these findings have for the construction and use of ocean color products, from both existing and future $R_{rs}(\lambda)$.

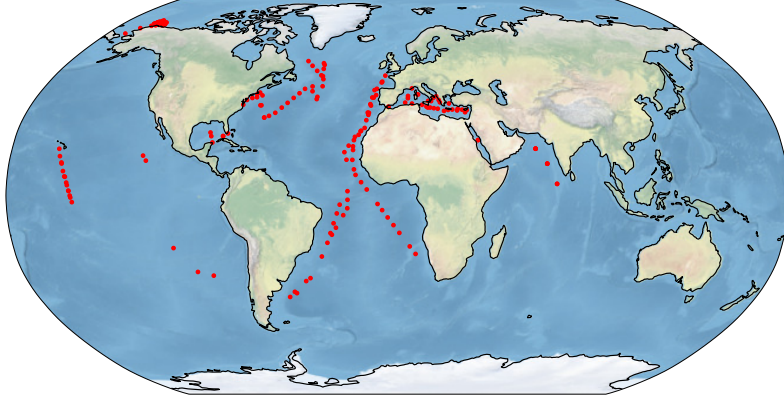


Figure 1. Locations of the 191 stations considered in this study (red dots).

2 Sea surface R_{rs} : hyperspectral versus multispectral

We first analyze a global sea surface hyperspectral $R_{rs}(\lambda)$ dataset to determine its DoF and how the DoF depends on spectral resolution (Chase et al., 2017; Kramer et al., 2022). The dataset includes $R_{rs}(\lambda)$ data at 191 locations at an effective 3.35 nm resolution (Chase et al., 2017) from 400–800 nm, linearly interpolated to 1 nm (Figure 1). We trimmed spectra to 700nm due to the large fraction of missing values $>700\text{nm}$; note that most of the non-empty values $>700\text{nm}$ are zeros and the non-zero-non-empty values, with a median of $<4 \times 10^{-5} \text{ sr}^{-1}$, have very small signal-to-noise ratios. The dataset includes measurements taken from 2004 to 2018 evenly distributed across months of the year, and from all major ocean basins ranging in latitude from 41°S to 74°N . We also compare these data to their MODIS-Aqua and SeaWiFS multispectral equivalents by convoluting the hyperspectral $R_{rs}(\lambda)$ with the MODIS-Aqua and SeaWiFS spectral response functions (available at https://oceancolor.gsfc.nasa.gov/docs/rsr/HMODISA_RSRs.txt and https://oceancolor.gsfc.nasa.gov/docs/rsr/SeaWiFS_RSRs.txt) to generate 10-waveband and 6-waveband datasets which correspond to what each instrument would have measured from the same optical input that the radiometer received when generating the hyperspectral $R_{rs}(\lambda)$ data.

We then apply principal component analysis (PCA) (Wold et al., 1987) to these 301-, 10- and 6-dimensional $R_{rs}(\lambda)$ datasets. PCA is a widely used method to reduce the dimensionality of datasets by identifying orthogonal vectors that explain the most variance in the data. PCA is linear in nature, which may result in an overestimation of effective dimensions by poorly approximating non-linear relationships between variables (e.g. a PCA performed on the pair (x, y) where $y = x^2$ will yield two DoF). Nonlinear generalizations do exist (Weinberger et al., 2004; Scholz et al., 2008), though these are less widely applied due to their additional complexity and computational requirements that make interpretation challenging. One may therefore consider the DoF we report to be upper bounds. We perform a PCA on each $R_{rs}(\lambda)$ dataset, standardizing each first by subtracting from each waveband its mean and then dividing by its standard deviation. This results in a percentage of total variance explained by each component. We use the broken-stick rule to choose the DoF, which states that the DoF is equal to the number of components that explain more variance than would be expected by randomly distributed data; this method was shown to be more consistent than a suite of others in a comparison (Jackson, 1993). These results can be shown visually as a ‘scree’ plot, which plots the percentage of variance explained by each component and for randomly distributed data; the DoF is the number of components with a higher percentage of variance explained than would be expected for randomly distributed data. Our figures also visibly demon-

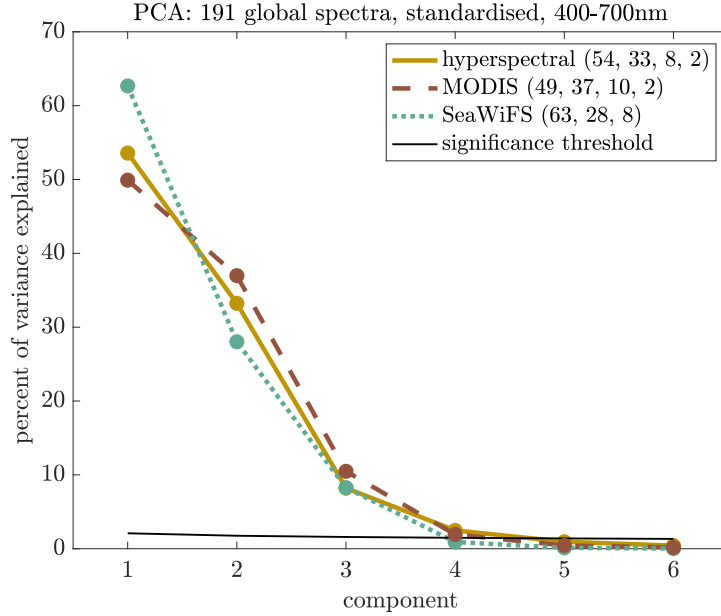


Figure 2. Scree plot of percent variance explained versus component for hyperspectral $R_{rs}(\lambda)$ dataset and MODIS-Aqua and SeaWiFS equivalents calculated from their spectral response functions. Black line indicates broken-stick significance threshold for hyperspectral data; numbers in legend give percent variance explained for each mode above this threshold in each case.

strate that one would get the same results from using the scree plot rule, which states that the DoF is equal to the number of components not sitting on the straight line made by the higher-order components, and was found to consistently capture the correct DoF plus one when the first point on this straight line was included (Jackson, 1993).

PCA analysis reveals that the hyperspectral in situ $R_{rs}(\lambda)$ dataset has four DoF (Figure 2); the first four components explain 54%, 33%, 8%, and 2%, totalling 97%, of the variance. The first four MODIS-Aqua equivalent $R_{rs}(\lambda)$ principal components have very similar percentages of variance explained: 49%, 37%, 10%, and 2%, totalling 99% of the total variance. In contrast, the first three SeaWiFS equivalent $R_{rs}(\lambda)$ principal components explain 63%, 28%, and 8%, totalling 99%, of the variance. This suggests that the hyperspectral $R_{rs}(\lambda)$ have four DoF, or four independent variables within the data, and that these four variables are effectively captured when reducing spectral resolution to the ten MODIS-Aqua wavebands, but not to the six SeaWiFS wavebands.

The ability of coarsened, MODIS-equivalent data to obtain the same number of DoF as the hyperspectral dataset is further supported by predictions of hyperspectral $R_{rs}(\lambda)$ from multispectral equivalents. To illustrate this, for each hyperspectral wavelength we perform a multivariate linear regression of $R_{rs}(\lambda)$ at that wavelength regressed against $R_{rs}(\lambda)$ at each waveband of both the MODIS-Aqua and SeaWiFS equivalent $R_{rs}(\lambda)$. We then calculate the root-mean-square-error (RMSE) of these regressions. For all wavelengths below 578 nm in the SeaWiFS case and 582 nm in the MODIS-Aqua case, the RMSE is smaller – and for many, much smaller – than 5% of the mean $R_{rs}(\lambda)$ at that wavelength, where 5% is a nominal relative uncertainty for satellite $R_{rs}(\lambda)$ (Figure 3). Even for wavelengths greater than this, the RMSE is still very small in absolute terms, $<0.00007 \text{ sr}^{-1}$, far smaller than the nominal 0.0003 sr^{-1} absolute error for 1km-by-1km pixels for PACE (Gordon & Wang, 1994). These small errors in predicting hyperspectral $R_{rs}(\lambda)$ from its multispectral equivalent underscore the extent to which different wavelengths’ $R_{rs}(\lambda)$

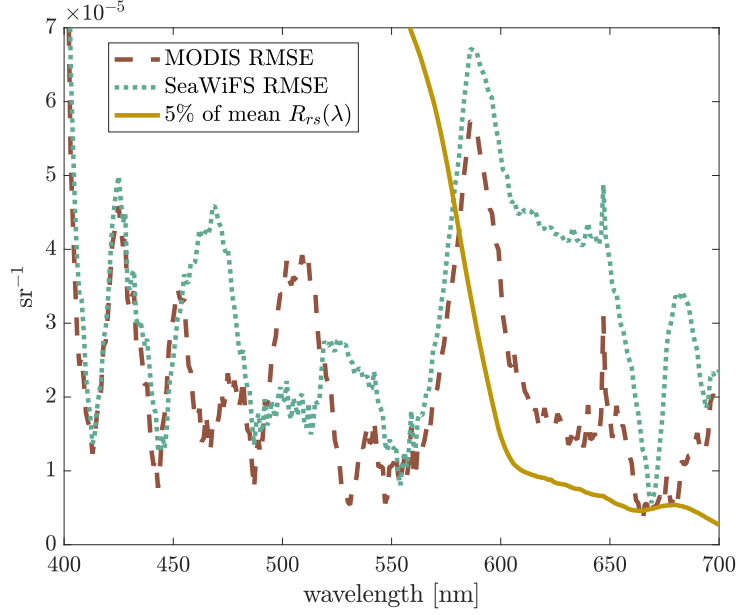


Figure 3. Root-mean-square-error of multivariate linear regressions of each hyperspectral wavelength versus the MODIS-Aqua and SeaWiFS equivalent $R_{rs}(\lambda)$. Solid line is 5% of the mean of each wavelength’s hyperspectral $R_{rs}(\lambda)$.

are correlated and demonstrate the ability of MODIS-Aqua equivalent multispectral $R_{rs}(\lambda)$ to preserve the dimensionality of hyperspectral $R_{rs}(\lambda)$. The fact that SeaWiFS-like $R_{rs}(\lambda)$ can accurately predict hyperspectral $R_{rs}(\lambda)$ to within PACE uncertainties but has fewer DoF than the in situ hyperspectral dataset is a reflection of the lower uncertainty on the in situ dataset than the expected PACE $R_{rs}(\lambda)$, and suggests that PACE $R_{rs}(\lambda)$ may have fewer DoF than the in situ hyperspectral dataset.

We also note that excluding wavelengths 651–700nm affects the DoF numbers presented here but not our conclusions. A choice of an upper limit of 650nm would be based on the fact that for all wavelengths above 648nm, >95% of measurements are below 0.0003 sr^{-1} , the nominal uncertainty of a 1km-by-1km pixel for PACE (Gordon & Wang, 1994). Repeating this analysis over 400–650nm results in hyperspectral and MODIS-Aqua-equivalent $R_{rs}(\lambda)$ data having three DoF, and SeaWiFS-equivalent $R_{rs}(\lambda)$ data having two DoF. This suggests that there is one DoF in the 651–700nm range that is picked up by hyperspectral and multispectral $R_{rs}(\lambda)$ alike; however, the $R_{rs}(\lambda)$ values are small enough (mean and median both $<1.2 \times 10^{-4} \text{ sr}^{-1}$ for all wavelengths 651–700nm) compared to the nominal 1km-by-1km pixel uncertainty $3 \times 10^{-4} \text{ sr}^{-1}$) that this DoF may not be useful for satellite applications, which we are interested in here. This is corroborated by the DoF < 3 in the next section, despite incorporating the full wavebands of both MODIS-Aqua and SeaWiFS. Note that when estimating the MODIS-Aqua- and SeaWiFS-equivalent data from 400–650nm hyperspectral data, the contribution of hyperspectral $R_{rs} > 650\text{nm}$ is not included; while both MODIS-Aqua and SeaWiFS have wavebands centered at $>650\text{nm}$, these wavebands’ spectral response functions are nonzero for some wavelengths in the range 400–650nm, and it is only the influence of these hyperspectral wavelengths on all wavebands that is considered. In other words, $R_{rs}(\lambda)$ is effectively set to zero for all hyperspectral wavelengths $>650\text{nm}$ when calculating the multispectral equivalent datasets in this case.

3 Climatologies: R_{rs} versus products

The analysis in Section 2 is based on instantaneous, local-scale $R_{rs}(\lambda)$ values measured in situ at the sea surface. The power of satellite $R_{rs}(\lambda)$ and derived products, however, lies in their near-continuous global spatial coverage, and many users are primarily interested in climatological data, which is near the coarsest spatial and temporal scales. In this section we therefore analyze climatological $R_{rs}(\lambda)$ and derived products, again via PCA to determine DoF.

We generated a $1^\circ \times 1^\circ$ climatology for each month using $R_{rs}(\lambda)$ data from SeaWiFS spanning 1997–2008, excluding the final 2 years of the mission due to known instrument issues (Siegel et al., 2014), using data downloaded from <https://oceancolor.gsfc.nasa.gov/>. We did the same for MODIS-Aqua, spanning the time period from July 2002 – June 2022. We generated analogous climatologies for derived products from each satellite over the same period and at the same spatial and temporal resolution, namely the extensive (i.e. mass-dependent) variables Chl, C_{phyto} , POC, PIC, and the intensive (i.e. mass-independent) variables Z_{eu} , ξ , the fraction of biovolume in the microplankton size class f_{micro} calculated from ξ as described in (Kostadinov et al., 2009), the particulate backscatter to chlorophyll ratio $b_{bp} : \text{Chl}$, and NPP as estimated by the CAFE (Silsbe et al., 2016) and CbPMv2 (Westberry et al., 2008) models. Chl, POC, and PIC were downloaded from <https://oceancolor.gsfc.nasa.gov/>, as was b_{bp} to calculate C_{phyto} according to (Graff et al., 2015) and $b_{bp} : \text{Chl}$ and the diffuse attenuation coefficient at 490nm to calculate Z_{eu} according to (Lee et al., 2007); SeaWiFS ξ and f_{micro} were derived as in (Kostadinov et al., 2009); and NPP products were downloaded from <http://sites.science.oregonstate.edu/ocean.productivity/index.php>. In total we then have climatologies for MODIS-Aqua, SeaWiFS $R_{rs}(\lambda)$, and ten derived products. We consider the six products Chl, C_{phyto} , POC, PIC, ξ , and Z_{eu} , to be core products and f_{micro} , $b_{bp} : \text{Chl}$, CAFE NPP, and CbPMv2 NPP to be ancillary products as these are either derived from the core products or rely on ancillary data other than $R_{rs}(\lambda)$.

We note that a PCA on the MODIS-Aqua climatologies of $R_{rs}(\lambda)$ and products other than ξ and f_{micro} yields the same results as those for SeaWiFS below, so we focus here only on the SeaWiFS climatologies because ξ and f_{micro} are not readily available for MODIS-Aqua. We find two DoF for SeaWiFS $R_{rs}(\lambda)$, but only one for the products (Figure 4). This result is not sensitive to which combination of products is used – for instance, including all the ancillary products as well still results in one DoF for the products. This result is also not sensitive to log-transformations of the variables that are log-normally (e.g. Chl, POC, PIC, C_{phyto} (Campbell, 1995)) or log-skew-normally (e.g. NPP, (Cael et al., 2018; Cael, 2021)) distributed, or removal of outliers, zeros, or negative values.

That MODIS-Aqua $R_{rs}(\lambda)$ have three DoF for the data in the previous section but two DoF from satellite-derived climatologies suggests that some reduction of complexity of the data occurs via some combination of increased sensor noise relative to ship-based data, atmospheric correction, or averaging over large space and time scales (Scott & Werdell, 2019). (Note (Scott & Werdell, 2019) also point out the difference between averaging $R_{rs}(\lambda)$ versus taking the ratio of averaged water-leaving radiance Lw and downwelling irradiance, which may introduce a slight bias but is unlikely to affect our results here.) Two DoF remain in satellite climatological $R_{rs}(\lambda)$ for both SeaWiFS and MODIS-Aqua, indicating the possibility of generating two independent products from these data. The suite of products tested above, however, has one fewer DoF than the $R_{rs}(\lambda)$. This is likely due to derived products’ appreciable uncertainties and/or strong correlations with chlorophyll. POC, ξ , and Z_{eu} , for instance, have Spearman rank correlations (across all months and 1° grid cells) of >0.9 with Chl. C_{phyto} ’s rank correlation with Chl is still fairly high, at 0.61, and is low largely due to small fluctuations when both are small; a simple spline fit of $\log(C_{phyto})$ against $\log(\text{Chl})$ yields an r^2 of 0.7.

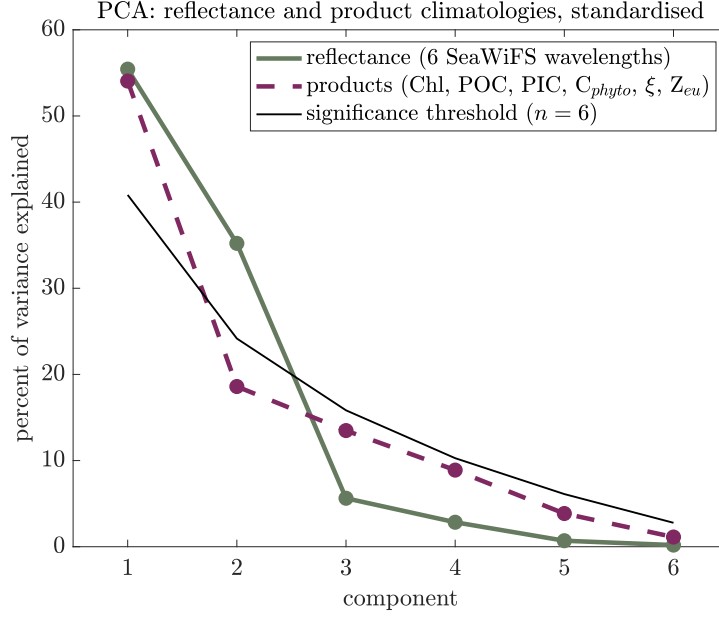


Figure 4. Scree plot of percent variance explained versus component for climatologies of SeaWiFS $R_{rs}(\lambda)$ and of six SeaWiFS- $R_{rs}(\lambda)$ -derived products. Black line indicates broken-stick significance threshold for six-dimensional data.

The exception is PIC, which has a rank correlation with Chl of 0.11. PIC, however, is highly sensitive to small variations in $R_{rs}(\lambda)$ for typical $R_{rs}(\lambda)$ values. To substantiate this, we performed a simple sensitivity analysis with the standard two-band PIC algorithm used by NASA for all but the most optically bright waters (see <https://oceancolor.gsfc.nasa.gov/atbd/pic/>). We calculated PIC for the climatological median $R_{rs}(\lambda)$ at 443 nm and 555 nm and for 5% variations, converting to normalized water-leaving radiance by multiplying by the global mean extraterrestrial solar irradiance. We then perturbed these $R_{rs}(\lambda)$ values with Gaussian noise at the 5% level, corresponding to the nominal uncertainty in $R_{rs}(\lambda)$. This noise at 443 nm results in 68% noise in PIC. By contrast, POC only varies 5% with these 5% variations in $R_{rs}(\lambda)$ at either wavelength. This indicates that in the bulk of cases, satellite-derived PIC is highly uncertain, on the order of 70% (and note the PIC uncertainty will be magnified more when considering documented uncertainties for $R_{rs}(\lambda)$ of 15-40% in some regions (Bisson et al., 2021)). In contrast, for relatively bright waters, the same exercise resulted in PIC variations of <10%, indicating that this algorithm performs well in instances when PIC values are high. Nonetheless, the high sensitivity to typical uncertainty in $R_{rs}(\lambda)$ for median waters explains why we find one DoF for the products even though PIC and Chl are not strongly correlated: derived PIC is noisy most of the time.

These results have two key implications. One is that there is additional information in climatological $R_{rs}(\lambda)$ that is not included in current derived products. This implies that existing products do not utilize the full set of $R_{rs}(\lambda)$ wavelengths. The other implication is that these products are not at all independent, and should not be treated as such when using them simultaneously. In other words, there are more products than there are DoF in satellite data. A numerical ecosystem model that reproduces the satellite-derived climatology of chlorophyll and of the particle size distribution's spectral slope should not be considered to be capturing two independent properties of the Earth sys-

tem. When using satellite products as inputs to other models, these products and their propagated uncertainties must be treated simultaneously rather than independently.

The results presented here are appropriate for global ocean analyses. The open ocean represents the largest area, and is composed primarily of Case 1 waters; that is, waters in which optical variability is dominated by chlorophyll (Morel & Prieur, 1977). In this context, it is in a sense unsurprising that the suite of $R_{rs}(\lambda)$ -derived products produced only one DoF. More optically complex waters, such as coastal regions and inland waters, have optical variability that is influenced by other constituents, such as colored dissolved organic material (CDOM), inorganic particles, and other pigments in addition to chlorophyll (e.g. Brown et al., 2008; Nelson & Siegel, 2013)). Analyses focused on these waters is likely to reveal a higher number of DoF from both $R_{rs}(\lambda)$ and derived products. Indeed, algorithms to derive concentrations of cyanobacteria and suspended particulate (Wang et al., 2016)) or distinguish between different phytoplankton species (Erickson et al., 2020) can be successful in such waters. However, we note that the in situ dataset used here (Figure 1) represents waters with $R_{rs}(\lambda)$ variability similar to that of the ocean as a whole, which can be seen by comparing the variation in $R_{rs}(\lambda)$ at each MODIS-Aqua wavelength from global satellite data with the same satellite data sub-sampled to the locations with in situ measurements (or the closest non-cloudy location). Sub-sampled satellite measurements have similar, and slightly lower, $R_{rs}(\lambda)$ in bluer wavelengths, indicating that the in situ dataset is oriented more towards optically complex coastal waters with substantial CDOM. This suggests that part of the explanation for the drop in DoF in satellite-derived climatologies comes from the fact that the in situ dataset sampled, as a whole, more optically complex waters.

We find that both $R_{rs}(\lambda)$ and variables derived from $R_{rs}(\lambda)$ are highly inter-correlated, reducing the number of DoF associated with each, with a greater reduction in DoF in the derived products. This becomes a problem when products are derived using empirical relationships with $R_{rs}(\lambda)$, and especially when the same wavelengths are used for the products that are assumed to be independent of each other; for example, over much of the ocean PIC, POC, and chlorophyll all are functions only of $R_{rs}(\lambda)$ at two wavelengths, at (or near, depending on the sensor) 443 and 555 nm. Certain combinations of PIC, POC, and chlorophyll, which may occur in the surface ocean, are therefore impossible to find using these algorithms. This is distinct from algorithms, typically called “quasi-analytical” or “semi-empirical”, that use known or assumed spectral shapes for absorption and scattering properties of optical constituents that can be related to the same derived products, such as PIC, POC, and chlorophyll (Werdell et al., 2013). These approaches may result in similar correlations and DoF between derived products, but do not inherently have the same problems as empirical approaches. We note that PACE will have, in addition to hyperspectral visible bands, UV bands from 350nm as well as spectral polarized bands. These measurements are expected to both improve the atmospheric correction (hence reduce the $R_{rs}(\lambda)$ uncertainties) as well as provide their own ocean signals, both of which may increase the DoF compared to those found here. In addition, it has been shown that adding other environmental variables such as SST can add useful information to inversions of phytoplankton groups (e.g. Chase et al., 2022) and thus another approach to increase DoF for inversions by adding relevant and independent information (e.g. mixed-layer depth and nutrients from BGC-Argo assimilating models).

4 Conclusion

The results presented here highlight the high degree of co-dependence between remote sensing reflectances at different wavelengths and of the products derived from these reflectances. For users of products based on existing reflectances, this primarily means factoring in the relationships between products when using more than one simultaneously. For the algorithms that generate these products from existing reflectances, these

results indicate a potential to improve the suite of available products to be more accurate and precise, and to account for the relationships between products and $R_{rs}(\lambda)$ wavebands. One way to do this, consistent with the findings above, would be to derive a single product such as chlorophyll as a function of all reflectance wavebands, derive an anomaly from chlorophyll-based expectations of a secondary product (e.g., phytoplankton community composition, size, POC, PIC, and so forth), then specify all other products explicitly as a function of these two, along the lines of Alvain et al. (2005). Ancillary and independent information can also be added to algorithms, as is currently done with net primary production models via temperature and mixed layer depth.

These findings are most relevant for algorithms that will generate products from hyperspectral reflectances in the future. The small number of degrees of freedom in hyperspectral reflectances indicates that only a few quantities can be estimated independently, and that different wavelengths' reflectances as measured from space will be strongly correlated. Complex algorithms that utilize the full spectrum of reflectance will need to factor in these correlations in order to generate reliable products. Crucially, if more than a few products are generated from hyperspectral reflectances, as is likely the case, such algorithms will also need to output the covariance information encoding the uncertainty in each product and the relationships between them. This can be achieved by some, but not all, machine learning techniques, on which this new generation of algorithms are likely to be based. The fact that hyperspectral reflectances can be predicted within nominal uncertainties by their multispectral equivalents suggests that hyperspectral resolution can play a role in improving ocean color products, but that it will be challenging to provide a substantially finer-grained picture of surface ocean ecosystems and biogeochemical cycles. Here by relying on principal component analysis we have focused on broad, first-order variations, but where such resolution may be most useful and generate novel insights is in investigating outliers and rare events, such as blooms or binning data over coherent features like eddies, where e.g. monospecific signatures may be resolved with spectral precision.

Open Research

Remote sensing data were downloaded from <https://oceancolor.gsfc.nasa.gov/> and <http://sites.science.oregonstate.edu/ocean.productivity/index.php>. All data and code are available at github.com/bbcael/eifoc for review purposes and will be given a Zenodo DOI should this manuscript be accepted for publication.

Acknowledgments

It is a pleasure to thank the many scientists whose collective work has generated the data on which this work relies. Cael acknowledges support from the National Environmental Research Council through Enhancing Climate Observations, Models and Data, and the European Union's Horizon 2020 Research and Innovation Programme under grant agreement No. 820989 (project COMFORT). The work reflects only the authors' view; the European Commission and their executive agency are not responsible for any use that may be made of the information the work contains. KB acknowledges support from NASA grant 80NSSC18K0957. EB acknowledges support from NASA grant 80NSSC20M0203. Cael lead and all other authors assisted with all aspects of this work. The authors have no competing interests to declare. This is PMEL contribution number 5445.

References

- Alvain, S., Moulin, C., Dandonneau, Y., & Breon, F.-M. (2005). Remote sensing of phytoplankton groups in case 1 waters from global seawifs imagery. *Deep Sea Research Part I: Oceanographic Research Papers*, 52(11), 1989–2004.

- Behrenfeld, M. J., & Falkowski, P. G. (1997). Photosynthetic rates derived from satellite-based chlorophyll concentration. *Limnology and oceanography*, 42(1), 1–20.
- Bisson, K., Boss, E., Werdell, P. J., Ibrahim, A., Frouin, R., & Behrenfeld, M. (2021). Seasonal bias in global ocean color observations. *Applied optics*, 60(23), 6978–6988.
- Bisson, K., Siegel, D. A., & DeVries, T. (2020). Diagnosing mechanisms of ocean carbon export in a satellite-based food web model. *Frontiers in Marine Science*, 7, 505.
- Brown, C. A., Huot, Y., Werdell, P. J., Gentili, B., & Claustre, H. (2008). The origin and global distribution of second order variability in satellite ocean color and its potential applications to algorithm development. *Remote Sensing of Environment*, 112(12), 4186–4203. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0034425708002162> doi: <https://doi.org/10.1016/j.rse.2008.06.008>
- Cael, B. (2021). Variability-based constraint on ocean primary production models. *Limnology and Oceanography Letters*, 6(5), 262–269.
- Cael, B., Bisson, K., & Follett, C. L. (2018). Can rates of ocean primary production and biological carbon export be related through their probability distributions? *Global biogeochemical cycles*, 32(6), 954–970.
- Cael, B., Bisson, K., & Follows, M. J. (2017). How have recent temperature changes affected the efficiency of ocean biological carbon export? *Limnology and Oceanography Letters*, 2(4), 113–118.
- Cael, B., Chase, A., & Boss, E. (2020). Information content of absorption spectra and implications for ocean color inversion. *Applied Optics*, 59(13), 3971–3984.
- Cael, B., Dutkiewicz, S., & Henson, S. (2021). Abrupt shifts in 21st-century plankton communities. *Science advances*, 7(44), eabf8593.
- Campbell, J. W. (1995). The lognormal distribution as a model for bio-optical variability in the sea. *Journal of Geophysical Research: Oceans*, 100(C7), 13237–13254.
- Chase, A. P., Boss, E., Cetinić, I., & Slade, W. (2017). Estimation of phytoplankton accessory pigments from hyperspectral reflectance spectra: toward a global algorithm. *Journal of Geophysical Research: Oceans*, 122(12), 9725–9743.
- Chase, A. P., Boss, E. S., Haëntjens, N., Culhane, E., Roesler, C., & Karp-Boss, L. (2022). Plankton imagery data inform satellite-based estimates of diatom carbon. *Geophysical Research Letters*, 49(13), e2022GL098076. Retrieved from <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2022GL098076> (e2022GL098076 2022GL098076) doi: <https://doi.org/10.1029/2022GL098076>
- DeVries, T., & Weber, T. (2017). The export and fate of organic matter in the ocean: New constraints from combining satellite and oceanographic tracer observations. *Global Biogeochemical Cycles*, 31(3), 535–555.
- Dutkiewicz, S., Cermenó, P., Jahn, O., Follows, M. J., Hickman, A. E., Taniguchi, D. A., & Ward, B. A. (2020). Dimensions of marine phytoplankton diversity. *Biogeosciences*, 17(3), 609–634.
- Erickson, Z. K., Werdell, P. J., & Cetinić, I. (2020). Bayesian retrieval of optically relevant properties from hyperspectral water-leaving reflectances. *Applied Optics*, 59(23), 6902–6917. doi: 10.1364/AO.398043
- Evers-King, H., Martínez-Vicente, V., Brewin, R. J., Dall’Olmo, G., Hickman, A. E., Jackson, T., ... others (2017). Validation and intercomparison of ocean color algorithms for estimating particulate organic carbon in the oceans. *Frontiers in Marine Science*, 251.
- Gordon, H. R., & Wang, M. (1994). Retrieval of water-leaving radiance and aerosol optical thickness over the oceans with seawifs: a preliminary algorithm. *Applied optics*, 33(3), 443–452.

- Graff, J. R., Westberry, T. K., Milligan, A. J., Brown, M. B., Dall’Olmo, G., van Dongen-Vogels, V., ... Behrenfeld, M. J. (2015). Analytical phytoplankton carbon measurements spanning diverse ecosystems. *Deep Sea Research Part I: Oceanographic Research Papers*, 102, 16–25.
- Hovis, W. A., Clark, D., Anderson, F., Austin, R., Wilson, W., Baker, E., ... others (1980). Nimbus-7 coastal zone color scanner: system description and initial imagery. *Science*, 210(4465), 60–63.
- Hu, C., Lee, Z., & Franz, B. (2012). Chlorophyll algorithms for oligotrophic oceans: A novel approach based on three-band reflectance difference. *Journal of Geophysical Research: Oceans*, 117(C1).
- Huot, Y., & Antoine, D. (2016). Remote sensing reflectance anomalies in the ocean. *Remote Sensing of Environment*, 184, 101–111. doi: 10.1016/j.rse.2016.06.002
- Jackson, D. A. (1993). Stopping rules in principal components analysis: a comparison of heuristical and statistical approaches. *Ecology*, 74(8), 2204–2214.
- Kostadinov, T., Siegel, D., & Maritorena, S. (2009). Retrieval of the particle size distribution from satellite ocean color observations. *Journal of Geophysical Research: Oceans*, 114(C9).
- Kramer, S. J., Siegel, D. A., Maritorena, S., & Catlett, D. (2022). Modeling surface ocean phytoplankton pigments from hyperspectral remote sensing reflectance on global scales. *Remote Sensing of Environment*, 270, 112879.
- Lee, Z., Weidemann, A., Kindle, J., Arnone, R., Carder, K. L., & Davis, C. (2007). Euphotic zone depth: Its derivation and implication to ocean-color remote sensing. *Journal of Geophysical Research: Oceans*, 112(C3).
- Mitchell, C., Hu, C., Bowler, B., Drapeau, D., & Balch, W. (2017). Estimating particulate inorganic carbon concentrations of the global ocean from ocean color measurements using a reflectance difference approach. *Journal of Geophysical Research: Oceans*, 122(11), 8707–8720.
- Morel, A., & Prieur, L. (1977). Analysis of variations in ocean color 1. *Limnology and oceanography*, 22(4), 709–722.
- Nelson, N. B., & Siegel, D. A. (2013). The global distribution and dynamics of chromophoric dissolved organic matter. *Annual review of marine science*, 5, 447–476.
- Nowicki, M., DeVries, T., & Siegel, D. A. (2022). Quantifying the carbon export and sequestration pathways of the ocean’s biological carbon pump. *Global Biogeochemical Cycles*, 36(3), e2021GB007083.
- O’Reilly, J. E., Maritorena, S., Mitchell, B. G., Siegel, D. A., Carder, K. L., Garver, S. A., ... McClain, C. (1998). Ocean color chlorophyll algorithms for seawifs. *Journal of Geophysical Research: Oceans*, 103(C11), 24937–24953.
- Scholz, M., Fraunholz, M., & Selbig, J. (2008). Nonlinear principal component analysis: neural network models and applications. In *Principal manifolds for data visualization and dimension reduction* (pp. 44–67). Springer.
- Scott, J. P., & Werdell, P. J. (2019). Comparing level-2 and level-3 satellite ocean color retrieval validation methodologies. *Optics Express*, 27(21), 30140–30157.
- Siegel, D., Buesseler, K., Doney, S. C., Sailley, S., Behrenfeld, M. J., & Boyd, P. (2014). Global assessment of ocean carbon export by combining satellite observations and food-web models. *Global Biogeochemical Cycles*, 28(3), 181–196.
- Silsbe, G. M., Behrenfeld, M. J., Halsey, K. H., Milligan, A. J., & Westberry, T. K. (2016). The cafe model: A net production model for global ocean phytoplankton. *Global Biogeochemical Cycles*, 30(12), 1756–1777.
- Wang, G., Lee, Z., Mishra, D. R., & Ma, R. (2016). Retrieving absorption coefficients of multiple phytoplankton pigments from hyperspectral remote sensing reflectance measured over cyanobacteria bloom waters. *Limnology and Oceanography: Methods*, 14(7), 432–447. doi: 10.1002/lom3.10102
- Weinberger, K. Q., Sha, F., & Saul, L. K. (2004). Learning a kernel matrix for nonlinear dimensionality reduction. In *Proceedings of the twenty-first international*

- 492 *conference on machine learning* (p. 106).
- 493 Werdell, P. J., Behrenfeld, M. J., Bontempi, P. S., Boss, E., Cairns, B., Davis,
 494 G. T., . . . others (2019). The plankton, aerosol, cloud, ocean ecosystem
 495 mission: status, science, advances. *Bulletin of the American Meteorological*
 496 *Society*, 100(9), 1775–1794.
- 497 Werdell, P. J., Franz, B. A., Bailey, S. W., Feldman, G. C., Boss, E., Brando, V. E.,
 498 . . . others (2013). Generalized ocean color inversion model for retrieving
 499 marine inherent optical properties. *Applied Optics*, 52(10), 2019–2037.
- 500 Westberry, T., Behrenfeld, M., Siegel, D., & Boss, E. (2008). Carbon-based primary
 501 productivity modeling with vertically resolved photoacclimation. *Global Bio-*
 502 *geochemical Cycles*, 22(2).
- 503 Wold, S., Esbensen, K., & Geladi, P. (1987). Principal component analysis. *Chemo-*
 504 *metrics and intelligent laboratory systems*, 2(1-3), 37–52.