



*Supplement of*

## **Sea-surface dimethylsulfide (DMS) concentration from satellite data at global and regional scales**

**Martí Galí et al.**

*Correspondence to:* Martí Galí (marti.gali.tapias@gmail.com)

The copyright of individual parts of the supplement might differ from the CC BY 4.0 License.

This file contains detailed information regarding data sources and pre-processing of the extended DMS database (sections 1-3), statistical analyses regarding the development, optimization and validation of the DMS<sub>SAT</sub> algorithm (section 4), and algorithm implementation (section 5).

## **S1 Sea surface DMS database and quality control**

In situ concentrations of DMS, DMSPt and chlorophyll a (Chl), accompanied by ancillary data (bottom depth, temperature, salinity, wind speed), were downloaded from the public sea-surface DMS database (<https://saga.pmel.noaa.gov/dms/>) on 13 April 2017. This database ( $n = 47745$  for DMS) was complemented with additional datasets obtained by the authors' teams ( $n = 403$  for DMS; Table S1) as detailed by (Galí et al., 2015). Quality control involved the deletion of DMS and DMSPt measurements potentially affected by methodological issues according to the criteria described by (Galí et al., 2015). DMS concentrations lower than 0.1 nM (0.4% of data) or higher than 100 nM (0.3% of data) were also removed. Removal of DMS < 0.1 nM is justified because standard gas chromatography methods, with a detection limit of a few pmol S, would require sample volumes of >50 mL to resolve such low concentrations, whereas most studies analyzed DMS in smaller sample volumes (Bell et al., 2012). DMS concentrations >100 nM are seldom measured in seawater and, in certain datasets obtained before the 2000s, can certainly be attributed to methodological artifacts, i.e. sparging unfiltered samples that contained *Phaeocystis* sp. (del Valle et al., 2009) or other DMS-producing phytoplankton sensitive to mechanical strain (Wolfe et al., 2002). After selecting surface data (depth  $\leq 10$  m), samples taken on the same day and within a radius of 100 m were averaged. The final dataset had 41304, 3700 and 9182 measurements for DMS, DMSPt and Chl, respectively, with 3637 DMS-DMSPt and 8141 DMS-Chl pairs.

## **S2 Satellite matchup data**

Daily and 8-day level 3-binned (L3BIN) data from the SeaWiFS and MODIS-Aqua sensors (9.28 and 4.64 km resolution, respectively) were matched to simultaneous in situ data from the DMS database (see Table S2). Matchups were done using individual pixels and the average of 3x3 and 5x5 pixel boxes centered on the in situ measurement location using SeaDAS 6.4 (Galí et al., 2015). For both sensors, the percentage of valid satellite matchups was around 10% and 40% for daily and 8-day composites, respectively. Merged satellite variables were created in order to increase the amount of data available for statistical analyses, after observing that inconsistencies between the two satellite datasets were small compared to other sources of uncertainty. The merged Chl<sub>SAT</sub>, Kd490<sub>SAT</sub>, PIC<sub>SAT</sub> and PAR<sub>SAT</sub> variables were created by averaging SeaWiFS and MODIS-Aqua match-ups with a hierarchical search procedure, i.e. prioritizing daily data over 8-day data and single-pixel data over 3x3 and 5x5 pixel box means. The resulting satellite matchups originated in 51% of cases from "quasi-simultaneous" SeaWiFS and MODIS-Aqua retrievals. The remaining 49% of observations was divided evenly between the two sensors. Daily and 8-day gridded SST (4.6 km) from the AVHRR sensor was also matched to the in situ database.

### S3 Binning of the extended sea surface DMS database

Statistical analyses were conducted using (i) non-binned data, (ii) data binned by month and 5x5 latitude-longitude bins (*M5x5*), and (iii) data binned by month and the 56 Longhurst biogeochemical provinces (*MLongh*) (Longhurst, 2010). For binned data, bins with less than 3 data counts (*M5x5*) and 5 data counts (*MLongh*) were discarded (for being poorly documented) in order to increase the robustness of regression models. These cutoff values are rather arbitrary, but similar results were obtained with slightly larger cutoff values. The statistics of data bins and the amount of bins and individual DMS measurements discarded through this procedure are shown in Table S3, showing that the amount of individual data points discarded through the binning procedure was <2.5%. The mean (median) data counts per bin were 26.5 (10) for *M5x5* binned data (see Fig. 1d) and 132.6 (57) for *MLongh* binned data.

### S4 Algorithm coefficients: uncertainty and optimization

To assess the uncertainty in fitted eq. 2 coefficients, we used the bootstrap method to produce  $10^5$  sets of regression coefficients for eq. 2 using the *MLongh* binned dataset. Fig. S2 shows the nonrandom relationships among eq. 2 coefficients.

Regression-derived coefficients were further optimized for global and regional scales using a constrained nonlinear optimization approach developed for this study.

The optimal coefficients of eq. 1 were obtained by minimizing a cost function different from RMSE (which is by definition the *cost function* minimized by least-squares regression). The best model was obtained with a cost function J defined as:

$$J = \text{RMSE} + \text{abs}(1 - R^2) + \text{abs}(1 - \text{Slope}_{\text{MA}}) \quad (\text{eq. S1}),$$

where  $\text{Slope}_{\text{MA}}$  is the major axis regression between observed and predicted fields. This cost function rewards the model coefficients that predict DMS with  $R^2$  and  $\text{Slope}_{\text{MA}}$  closest to 1. The goodness-of-fit statistics used in eq. S1 were calculated in  $\log_{10}$  space using the same *MLongh* binned dataset. To obtain realistic solutions, we constrained the optimization to the 99% confidence intervals of the  $10^5$  bootstrapped regression coefficients shown in Fig. S2 (*MLongh* binned dataset). The resulting optimal model (eq. 2f) had higher DMSPt ( $\beta$ ) and PAR ( $\gamma$ ) coefficients and a smaller y-intercept than eq. 2e, and moved the modeled DMS concentration closer to the 1:1 agreement line without degrading neither RMSE nor  $R^2$  (Table 2). The optimized model coefficients were validated using an independent dataset as described in section 3.1.3 of the main text (see Fig. 4).

The same approach was used to optimize the eq. 2 coefficients for the Bermuda Atlantic Time Series (BATS) site. In this case we used the 3 years of monthly measurements (upper mixed layer means) to obtain the regionally tuned coefficients (eq. 2h), which were not further validated using independent datasets given that they were only used to demonstrate the portability of the algorithm.

## S5 Algorithm implementation: data sources and processing chain

The full DMS<sub>SAT</sub> algorithm (Fig. 2) was implemented to produce (i) a monthly global DMS<sub>SAT</sub> climatology based on SeaWiFS climatological 1997-2010 data; and (ii) regional time series with 8-day resolution for the period 2003-2016 using MODIS-Aqua data. In both cases we used reprocessing 2014.0. The data sources are summarized in Table S2.

Global DMS<sub>SAT</sub> fields were computed using ocean color data from the SeaWiFS 1/12° gridded monthly climatology (1997-2010) in combination with the 1/2° gridded monthly MLD climatology from MIMOC. The input Chl<sub>SAT</sub> product was either the band-ratio algorithm OC4-OCI (the current standard NASA Chl algorithm) or the semi-analytical GSM algorithm (Maritorena et al., 2002). The euphotic layer depth (Zeu) was computed as either the 1% penetration depth of 490 nm radiation ( $Zeu = 4.6/Kd490$ ) or the semi-analytical Zeu from (Lee et al., 2007).

Regional DMS<sub>SAT</sub> time series between 2003 and 2016 for latitudes  $>45^{\circ}\text{N}$  were computed using daily MODIS-Aqua data (4.64 km) combined with the MIMOC MLD climatology (linearly interpolated onto the MODIS-Aqua grid and a daily period). Unlike the global implementation, which used Chl and Kd490 directly downloaded from the NASA Ocean Color website, in this case we used remote sensing reflectance spectra (Rrs) to compute bio-optical variables. Chl<sub>SAT</sub> was computed using either the OC3 band-ratio algorithm (MODIS version of the OC algorithm (O'Reilly et al., 1998)) or the GSM algorithm (Maritorena et al., 2002). The diffuse attenuation coefficient at 488 nm, Kd488, which is nearly equivalent to the SeaWiFS Kd centered on 490 nm, was computed using the semi-analytical algorithm of Lee et al., (2005) and used to estimate  $Zeu = 4.6/Kd488$ .

Since non-climatological satellite data frequently contain data gaps caused by cloudiness, we applied a binning and gap-filling procedure to obtain full coverage. First, we calculated DMSPt<sub>SAT</sub> (Galí et al., 2015) using daily 4.64 km MODIS-Aqua data (native L3bin resolution). Daily 4.64 km data were then averaged into 6x6 pixel boxes (27.84 km macropixels) and 8-day periods. The remaining gaps (10% pixels) were successively filled with 8-day (9%) and monthly (1%) climatologies of each variable. Finally, DMS<sub>SAT</sub> was calculated from 8-day 27.84 km DMSPt<sub>SAT</sub> and PAR<sub>SAT</sub>.

We performed a further sensitivity test to analyze the effect of different mixed layer depth (MLD) products on DMSPt<sub>SAT</sub>. Besides our standard MLD obtained from the MIMOC climatology (Schmidtko et al., 2013), we tested the monthly MLD time series diagnosed by the Global Ocean Data Assimilation System (GODAS) ocean circulation model (which covers latitudes  $<65^{\circ}$  between 1980 and present). Eight-day DMSPt<sub>SAT</sub> time series between 2003 and 2015 were generated for two areas located in the North Atlantic (Iceland Basin) and Pacific (Bering Sea), which show contrasting wintertime MLD —owing to differences in salinity stratification. To verify the accuracy of MIMOC and GODAS MLD, we compared them to MLD derived from collocated ARGO float profiles. As shown in Fig. S1, GODAS overestimated wintertime MLD in the subpolar Atlantic. This translated, in some years, in slightly lower DMSPt<sub>SAT</sub> with GODAS because the DMSPt sub-algorithm switched at a later date from the 'mixed'

( $Z_{eu}/MLD < 1$ ) to the 'stratified' ( $Z_{eu}/MLD > 1$ ) waters equation. Differences were almost absent in the Bering Sea. Overall, the use of climatological or model-derived MLD had a negligible effect on diagnosed  $DMSPt_{SAT}$ .

## Tables

**Table S1.** Compilation of studies added to the sea surface DMS database used for algorithm development and validation.

Reference	Region	Dates
Levasseur et al. 2006	NW Pacific	Jul 2002
Matrai et al. 2007	Barents Sea	1998, 1999, 2001
Royer et al. 2010	NW Pacific	Jul 2007
Lizotte et al. 2012	NW Atlantic	2003
Luce et al. 2011	Canadian Arctic	Late summer 2007, 2008
Royer et al. 2015	Tropical Atlantic, Pacific and Indian oceans	Dec 2010—Jul 2011
Royer et al. 2016	NW Mediterranean	Sep 2011, May 2012

**Table S2.** Summary of datasets used to complement the sea surface DMS database (for algorithm development and validation) and to implement the DMS<sub>SAT</sub> algorithm at global and regional scales. R: reprocessing.

Data type and source	Use, data version and type (when applicable).
SeaWiFS and MODIS-Aqua ocean color data: <a href="https://oceancolor.gsfc.nasa.gov/">https://oceancolor.gsfc.nasa.gov/</a>	<b>Matchups:</b> SeaWiFS R2010.0 and MODIS-Aqua R2013.1. DAY and 8D L3BIN data. <b>Global</b> ocean implementation of DMS <sub>SAT</sub> , SD02 and VS07 algorithms: MODIS-Aqua R2014.0. MONTH L3SMI data. <b>Regional</b> DMS <sub>SAT</sub> implementation: MODIS-Aqua R2014.0. DAY L3SMI data. <i>Note: MODIS-Aqua nighttime SST was used between 2003-2016 instead of AVHRR SST.</i>
AVHRR sea surface temperature (SST): <a href="https://www.nodc.noaa.gov/SatelliteData/pathfinder4km/">https://www.nodc.noaa.gov/SatelliteData/pathfinder4km/</a>	<b>Matchups:</b> Pathfinder v5.2. <b>Global</b> DMS <sub>SAT</sub> implementation: 1/12 degree monthly climatology based on Pathfinder v5.2. <b>Regional</b> DMS <sub>SAT</sub> implementations of DMS <sub>SAT</sub> : Pathfinder v5.3. <i>Note: MODIS-Aqua nighttime SST was used between 2003-2016 instead of AVHRR SST.</i>
Monthly Isopycnal & Mixed-layer Ocean Climatology (MIMOC): <a href="http://www.pmel.noaa.gov/mimoc/">http://www.pmel.noaa.gov/mimoc/</a>	<b>Matchups:</b> Native 0.5 degrees monthly resolution. <b>Global</b> ocean implementation of DMS <sub>SAT</sub> , SD02 and VS07 algorithms: Native 0.5 degrees monthly resolution. <b>Regional</b> DMS <sub>SAT</sub> implementation: Reprojected onto MODIS-Aqua 4.64 sinusoidal grid and interpolated to 1-day resolution.
NCEP Global Ocean Data Assimilation System (GODAS): <a href="https://www.esrl.noaa.gov/psd/data/gridded/data.godas.html">https://www.esrl.noaa.gov/psd/data/gridded/data.godas.html</a>	<b>Regional</b> DMS <sub>SAT</sub> implementation sensitivity tests: Reprojected onto MODIS-Aqua 4.64 sinusoidal grid and interpolated to 1-day resolution.
ARGO float MLD: <a href="http://mixedlayer.ucsd.edu/">http://mixedlayer.ucsd.edu/</a>	<b>Regional</b> DMS <sub>SAT</sub> implementation sensitivity tests: matched to 27.84 8-day resolution of DMS <sub>SAT</sub> dataset.

World Ocean Atlas 2009 nutrient and salinity data: <a href="http://coastwatch.pfeg.noaa.gov/erddap/griddap/nodcWoa09sea1n.html">http://coastwatch.pfeg.noaa.gov/erddap/griddap/nodcWoa09sea1n.html</a>	<b>Matchups:</b> 1-degree monthly resolution.
General Bathymetric Chart of the Oceans <a href="https://www.gebco.net/data_and_products/gridded_bathymetry_data/gebco_one_minute_grid/">https://www.gebco.net/data_and_products/gridded_bathymetry_data/gebco_one_minute_grid/</a>	<b>Matchups:</b> GEBCO One Minute Grid (November 2008)

**Table S3.** Binning statistics for the extended sea surface DMS database. V = valid, D = discarded.

Binning scheme	Minimum N per bin	Bin counts			Individual data counts		
		V	D	% D	V	D	% D
M5x5	N ≥ 3	1562	742	32.2%	40326	978	2.4%
Mlongh	N ≥ 5	322	66	17.0%	41156	148	0.4%

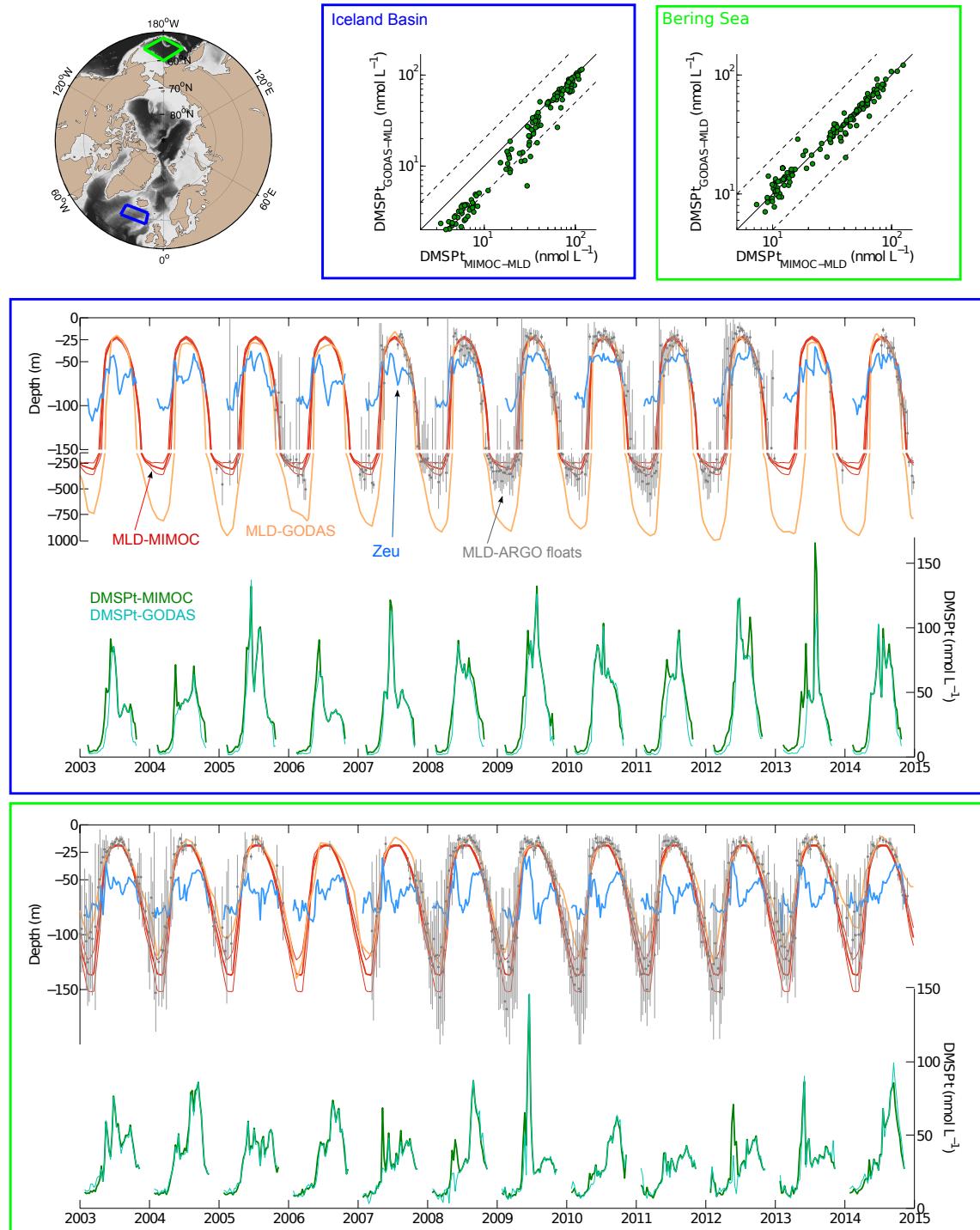
**Table S4.** Summary of the stepwise regression analysis. Only models with significant coefficients are shown. WS: wind speed.

Model	R <sup>2</sup> adj	RMSE	AIC	N
<i>Non-binned data</i>				
$\log_{10}\text{DMS} = -1.21 + 0.67 \log_{10}\text{DMSPt} + 0.0136 \text{ PAR}$	0.50	0.35	2743	3620
$\log_{10}\text{DMS} = -1.01 + 0.75 \log_{10}\text{DMSPt} + 0.0111 \text{ PAR}_{\text{MLD}}$	0.47	0.37	2992	3595
$\log_{10}\text{DMS} = -1.17 + 0.64 \log_{10}\text{DMSPt} + 0.0151 \text{ PAR} - 0.0038 \text{ SST}$	0.51	0.35	2722	3615
$\log_{10}\text{DMS} = -1.22 + 0.63 \log_{10}\text{DMSPt} + 0.0152 \text{ PAR} + 0.055 \log_{10}[\text{NO}_3]$	0.51	0.34	2357	3418
$\log_{10}\text{DMS} = -1.17 + 0.64 \log_{10}\text{DMSPt} + 0.0147 \text{ PAR} - 0.040 \log_{10}[\text{N-cline}]$	0.51	0.35	2562	3418
$\log_{10}\text{DMS} = -1.79 + 0.70 \log_{10}\text{DMSPt} + 0.0159 \text{ PAR} + 0.012 \text{ Salinity}$	0.52	0.31	1027	1911
$\log_{10}\text{DMS} = -1.04 + 0.59 \log_{10}\text{DMSPt} + 0.0118 \text{ PAR} - 0.0082 \text{ WS}$	0.48	0.32	764	1442
$\log_{10}\text{DMS} = -0.63 + 0.45 \log_{10}\text{DMSPt} + 0.0129 \text{ PAR} + 0.098 \log_{10}\text{PIC}_{\text{SAT}}$	0.35	0.29	388	1123
<i>MLongh binned data (bin medians)</i>				
$\log_{10}\text{DMS} = -1.02 + 0.45 \log_{10}\text{DMSPt} + 0.0163 \text{ PAR}$	0.57	0.21	-31.0	118
$\log_{10}\text{DMS} = -0.94 + 0.69 \log_{10}\text{DMSPt} + 0.0172 \text{ PAR}_{\text{MLD}}$	0.52	0.25	15.7	118
$\log_{10}\text{DMS} = -0.91 + 0.48 \log_{10}\text{DMSPt} + 0.0189 \text{ PAR} - 0.0087 \text{ SST}$	0.59	0.24	0.9	118
[NO <sub>3</sub> ] coefficient non-significant (p = 0.26)				118
[N-cline] coefficient non-significant (p = 0.71)				118
Salinity coefficient non-significant (p = 0.35)				102
WS coefficient non-significant (p = 0.86)				97
$\log_{10}\text{DMS} = -0.64 + 0.29 \log_{10}\text{DMSPt} + 0.0118 \text{ PAR} + 0.033 \log_{10}\text{PIC}_{\text{SAT}}$	0.52	0.21	-14.8	86

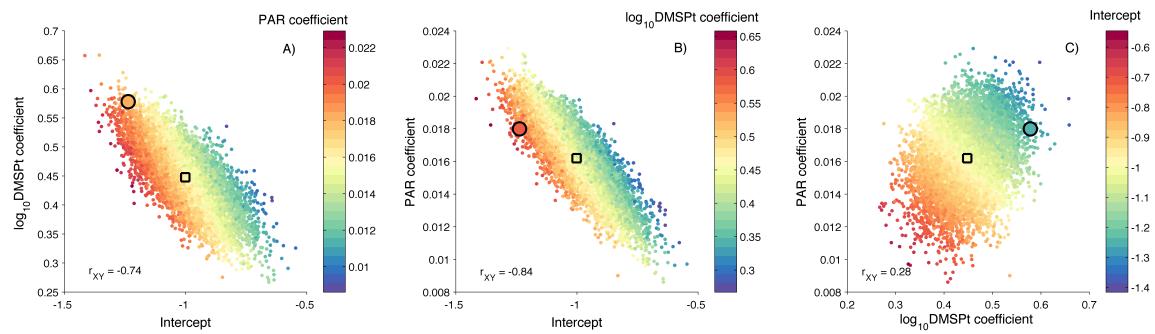
**Table S5.** DMS<sub>SAT</sub> validation statistics for constrained and *unconstrained (in italics)* Chl<sub>SAT</sub> error. N increases from 86 to 1293 as the tolerated Chl<sub>SAT</sub> error increases. N is 14677 for unconstrained Chl error.

Algorithm	R <sup>2</sup> (log <sub>10</sub> space)	RMSE (log <sub>10</sub> space)	MAPE (linear space)	Mean bias (linear space)
DMS <sub>SAT</sub> (eq. 2f)	0.40–0.53	0.25–0.30	54–66%	11–36%
	<i>0.29</i>	<i>0.38</i>	<i>108%</i>	<i>-9%</i>
SD02 (eq. 3)	0.22–0.31	0.26–0.29	50–59%	-14 – -7%
	<i>0.24</i>	<i>0.40</i>	<i>138%</i>	<i>2%</i>
VS07 (eq. 4)	0.10–0.15	0.30–0.33	78–82%	-2–10%
	<i>0.02</i>	<i>0.45</i>	<i>89%</i>	<i>-41%</i>

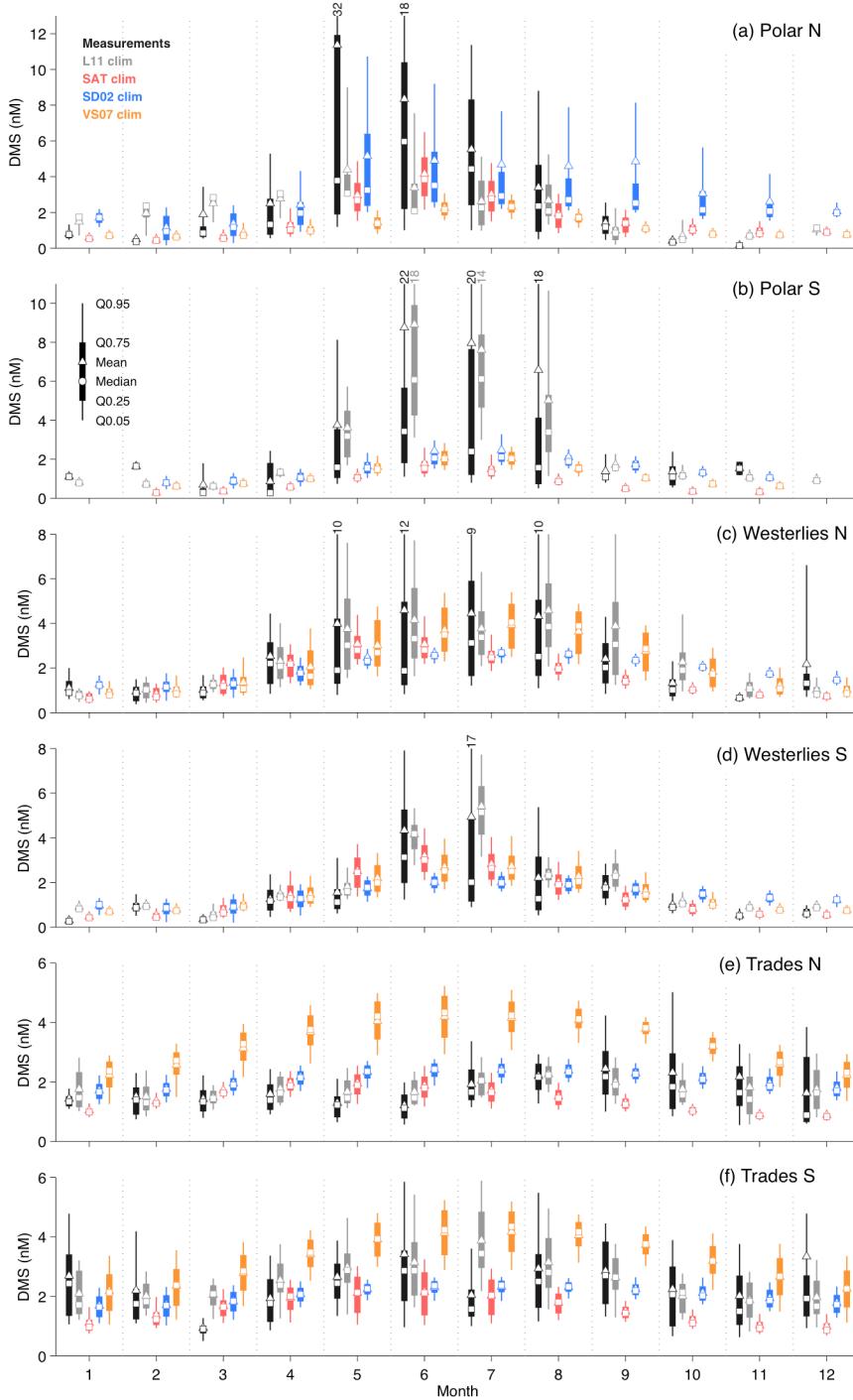
## Figures



**Figure S1: Sensitivity of the DMSP<sub>T</sub><sub>SAT</sub> sub-algorithm to different MLD input data.** Top: scatterplots comparing DMSP<sub>T</sub><sub>SAT</sub> calculated with the MIMOC and GODAS MLD products. Middle and bottom: 2003–2015 8-day resolution time series of MIMOC and GODAS MLD, ARGO float MLD, and satellite derived euphotic layer depth (Zeu). The scatterplots and time series plots show the regional means of each variable for the two regions highlighted in the map: Iceland Basin and Bering Sea.



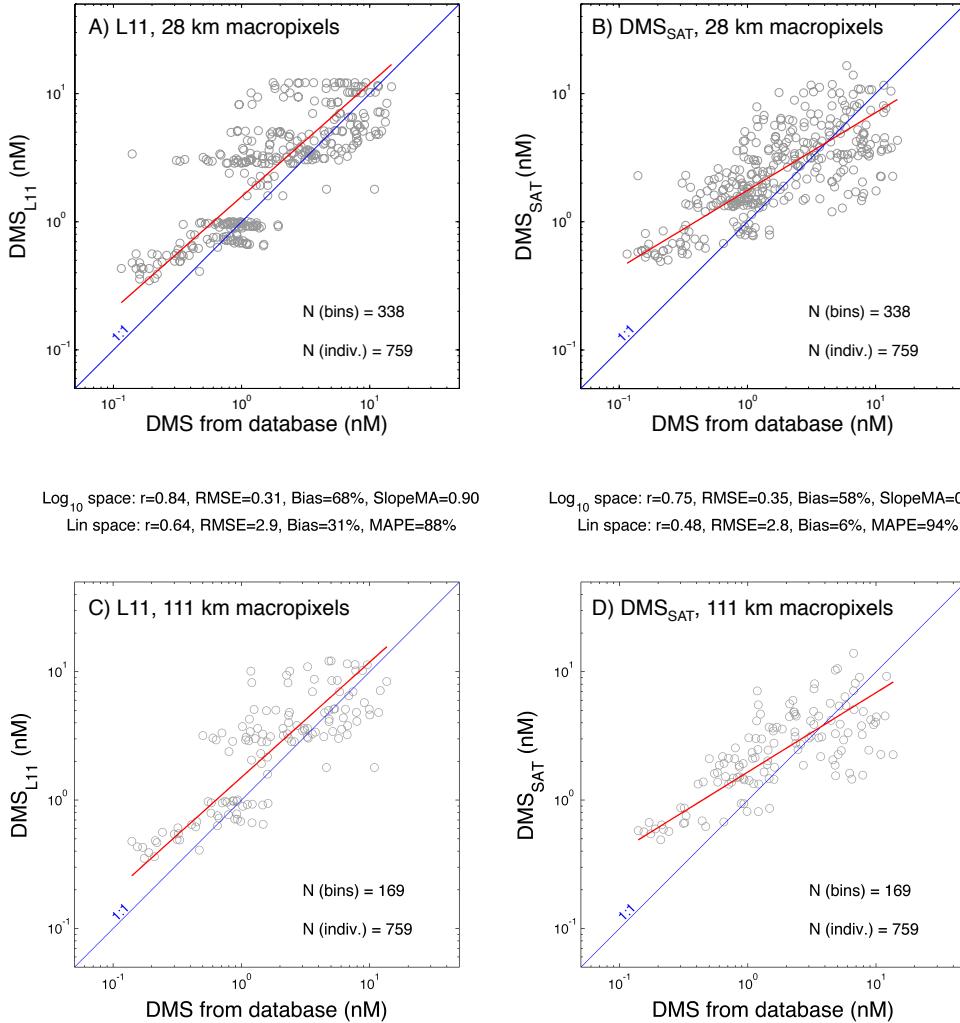
**Figure S2.** Bootstrapped regression model coefficients obtained for the equation:  $\log_{10}(\text{DMS}) = \alpha + \beta \log_{10}(\text{DMSPt}) + \gamma \text{PAR}$ , with  $n = 10000$ . Black squares show the mean coefficients, which are equivalent to those obtained through regular multiple regression (eq. 2e). The uncertainty envelopes defined by the 10000 bootstrapped coefficients were used as the bounds for a constrained optimization procedure (see text). Filled circles represent the resulting optimized coefficients (eq. 2f).



**Figure S3. DMS seasonal cycles by biomes.** The monthly means, medians, interquartile range and 5%-95% percentiles are shown for the in situ database, the L11 climatology, and remote sensing climatologies derived from the  $\text{DMS}_{\text{SAT}}$ , SD02 and VS07 algorithms. The temporal axis has been shifted by 6 months in the Southern hemisphere, i.e., July is the 1st month and June the 12th.

$\text{Log}_{10}$  space:  $r=0.78$ , RMSE=0.35, Bias=75%, SlopeMA=0.88  
 Lin space:  $r=0.61$ , RMSE=3.0, Bias=31%, MAPE=113%

$\text{Log}_{10}$  space:  $r=0.72$ , RMSE=0.37, Bias=69%, SlopeMA=0.61  
 Lin space:  $r=0.49$ , RMSE=2.9, Bias=10%, MAPE=110%



**Figure S4.  $\text{DMS}_{\text{SAT}}$  validation statistics for the subpolar North Atlantic.** In situ data points correspond to the area delimited by latitudes between  $45^{\circ}\text{N}$ – $60^{\circ}\text{N}$  and longitudes between  $55^{\circ}\text{W}$ – $15^{\circ}\text{E}$ , marked in Fig. 9 of the paper. Left-hand plots compare the L11 climatology to the in situ measurements (on which it is based). Right-hand plots compare  $\text{DMS}_{\text{SAT}}$  to the same in situ measurements. Top and bottom plots differ in the degree of spatial binning applied: 28 km and 111 km. The temporal binning period is 8 days in all cases. On top of each plot we provide both  $\text{Log}_{10}$  and linear space statistics. The perfect agreement line is shown in blue and the model-data linear regression in red ( $\text{Log}_{10}$  space).