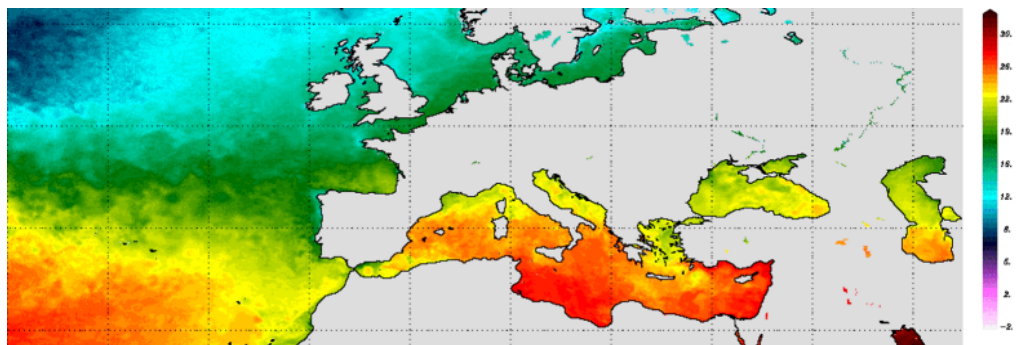


ESA CCI+ Phase 1 Sea Surface Temperature (SST)



Algorithm Theoretical Basis Document D2.1 v1

Issue Date: **02 July 2020**

ESA REF: ESA/AO/1-9322/18/I-NB

Algorithm Theoretical Basis Document D2.1 v1

Algorithm Theoretical Basis Document D2.1 v1

SIGNATURES AND COPYRIGHT

Title : ESA CCI+ Phase 1 Sea Surface Temperature (SST)

Volume : Algorithm Theoretical Basis Document D2.1 v1

Issued : 02 July 2020

Authored :



Prof Chris Merchant (UoR)

Authored :

J Hoyer

Dr Jacob Høyer (DMI)

Address : University of Reading,
Whiteknights,
Reading,
Berkshire,
RG6 6AH,
United Kingdom

Copyright : © University of Reading 2020. The Copyright of this document is the property of the University of Reading. It is supplied on the express terms that it be treated as confidential, and may not be copied, or disclosed, to any third party, except as defined in the contract, or unless authorised by the University of Reading in writing.

Algorithm Theoretical Basis Document D2.1 v1

TABLE OF CONTENTS

1. SCOPE.....	3
2. STATISTICAL RETRIEVAL USING ERA-5	4
2.1 Introduction.....	4
2.2 Input data	4
2.2.1 ESA-CCI Multi-sensor Matchup Dataset.....	4
2.2.2 AMSR-E orbital data	5
2.2.3 AMSR2 orbital data	5
2.2.4 Preprocessing	6
2.3 Algorithm description.....	6
2.3.1 WS retrieval algorithm.....	6
2.3.1.1 <i>1st-step: Global retrieval algorithm</i>	6
2.3.1.2 <i>2nd-step: Specialized WS retrieval algorithms</i>	7
2.3.2 SST retrieval algorithm	8
2.3.2.1 <i>1st-step: Specialized latitude and ascending/descending retrieval algorithms</i>	8
2.3.2.2 <i>2nd-step: Specialized SST and WS retrieval algorithms</i>	9
2.3.3 RFI filter	11
2.3.4 Uncertainty model	12
2.3.5 L2P flags and quality levels	13
2.3.6 Regression setup	15
2.4 Output data.....	16
3. OPTIMAL ESTIMATION TUNING	18
3.1 Introduction.....	18
3.2 Moderately non-linear optimal estimation.....	18
3.3 Extension to parameter estimation with reference observations	19
3.4 Bias correction parameters	22
3.5 Desroziers estimator for observation-simulation error covariance matrix	23
3.6 Desroziers estimate for prior error covariance matrix	24
3.7 Convergence	24
3.8 Iteration of parameter cycles	24
3.9 Progress towards exploitation	24
4. DESERT-DUST RELATED BIASES	26
4.1 Introduction.....	26
4.2 Notes on Data	26
4.3 Notes on methods	27
4.4 Notes on results	29
5. REFERENCES.....	34

Algorithm Theoretical Basis Document D2.1 v1

List of Acronyms

AMSR	Advanced Microwave Scanning Radiometer
ATBD	Algorithm Theoretical Basis Document
AVHRR	Advanced Very High Resolution Radiometer
BAMS	Bulletin of the American Meteorological Society
BT	brightness temperature
CAMS	Copernicus Atmospheric Monitoring Service
CCI	Climate Change Initiative
CCMP	Cross-calibrated Multi-platform
CDR	climate data record
CDS	Climate Data Store
CLW	cloud liquid water
DMI	Danish Meteorological Institute
ECMWF	European Centre for Medium-range Weather Forecasting
EDR	environmental data record
ERA	ECMWF Re-Analysis
ESA	European Space Agency
FIDUCEO	Fidelity and Uncertainty In Climate Data Records from Earth Observation
GBCS	Generalise Bayesian Cloud Screening
GDS	GHR SST Data Specification
GHR SST	Group for High Resolution SST
GT MBA	Global Tropical Moored Buoy Array
IR	infrared
MMD	multisensor matchup mataset
MMS	multisensor matchup system
MOHC	Met Office Hadley Centre
MWR	microwave radiometer
NEDT	noise equivalent differential temperature
NSIDC	National Snow and Ice Data Center
NWP	numerical weather prediction
OE	optimal estimation
OI	optimal interpolation
PMEL	Pacific Marine Environmental Laboratory
PUG	product user guide
RFI	radio frequency interference
RSS	Remote Sensing Systems
RTTOV	Radiative Transfer for TOVS
SD	standard deviation
SST	sea surface temperature
TBC	to be confirmed
TCWV	total column water vapour
UKMO	UK Met Office
WP	work package
WS	wind speed

Algorithm Theoretical Basis Document D2.1 v1

1. SCOPE

This is a report in preparation for a comprehensive Algorithm Theoretical Basis Document (ATBD) that will be completed to describe the Sea Surface Temperature Climate Change Initiative Version 3 Climate Data Record (SST CCI v3 CDR). The full ATBD will be document ATBD D2.1 v3.

The form of content of this interim ATBD report is notes arising from the algorithm development work to date within Phase 1 of SST CCI+.

The scope of work addressed in this interim ATBD report covers progress with respect to preparing passive microwave SST retrievals for potential inclusion in the CDR (in CDR v3 or v4, decision yet to be made).

A further interim ATBD report (D2.1 v2) will be prepared describing algorithm developments for AVHRR SST retrieval in the infrared, due at month 18 of the project.

The specific scope of this document is in summary:

- To describe the use of the ECMWF Re-analysis 5 (ERA-5) numerical weather prediction (NWP) fields in microwave SST
- To present developments of the theory of optimal estimation (OE) for microwave SSTs, namely how reference sensors can be used to estimate bias corrections and error covariance matrices to "tune" OE
- To present use of microwave SSTs from Phase 2 to help estimate desert-dust-related SST biases in the v2 CDR analysis in preparation for reducing biases in the v3 CDR

Algorithm Theoretical Basis Document D2.1 v1

2. STATISTICAL RETRIEVAL USING ERA-5

2.1 Introduction

This chapter describes the ATBD for the DMI regression algorithm used for retrieval of SST and wind speeds from JAXA's Advanced Microwave Scanning Radiometer - Earth Observing System (AMSR-E) and its follow-on instrument AMSR2 (Advanced Microwave Scanning Radiometer 2). The algorithm was used within the phase 2 of the European Space Agency Climate Change Initiative Sea Surface Temperature (ESA-CCI SST) project to generate a global climate data record (CDR) of level 2 SSTs with associated uncertainties (see Alerskans et al., 2020 for extensive description and validation of the algorithm). A consistent algorithm has been used for both the AMSR-E and AMSR2 observations. This version of the ATBD includes the updates that have been made for ESA-CCI+ SST phase 1, related to the use of information from ERA-5 instead of ERA-I.

2.2 Input data

The retrieval algorithm is designed to be able to use two sources of input data; a Multi-sensor Matchup Dataset (MMD), which is used for tuning, development and validation of the retrieval algorithm, and orbital AMSR-E and AMSR2 data, which is used for producing the microwave radiometer (MWR) SST climate data record.

2.2.1 ESA-CCI Multi-sensor Matchup Dataset

For tuning and development of the retrieval algorithm, as well as for assessment and validation of the performance of the algorithm, Multi-sensor Matchup Datasets (MMDs), versions MMD6c and MMD6b were used as input. The MMDs were generated using the Multi-sensor Matchup System (MMS) software which was developed within the ESA-CCI SST project and the European Union's Horizon 2020 research and innovation programme under grant agreement No 638822 (FIDUCEO project) (Block et al., 2018). The MMD6c consists of AMSR-E orbital data matched to *in situ* SST measurements and MMD6b is the corresponding matchup database for AMSR2. The *in situ* dataset contains quality controlled measurements from Global Tropical Moored Buoy Array (GTMBBA) data from NOAA PMEL (McPhaden et al., 2010), the International Comprehensive Ocean-Atmosphere Dataset (ICOADS) version 2.5.1 (Woodruff et al., 2011) and the Met Office Hadley Centre (MOHC) Ensembles dataset version 4.2.0 (EN4) (Good et al., 2013). In addition, the MMD also includes NWP data from ERA-5 (Copernicus Climate Change Service (C3S), 2017), which have been interpolated in both space and time to the matchup location. Furthermore, the Cross-Calibrated Multi-Platform (CCMP) surface vector winds

Algorithm Theoretical Basis Document D2.1 v1

(Atlas et al., 2011) were collocated with the MMDs and were used for tuning and development, as well as validation, of the wind speed (WS) retrieval algorithm.

To obtain independent results, the MMDs are divided into seven subsets each, to be used for either tuning and development or validation of the algorithm:

- WS1_TRAIN: training subset for the 1st-stage WS retrieval algorithm;
- WS1_TEST: validation subset for the 1st-stage WS retrieval algorithm;
- WS2_TRAIN: training subset for the 2nd-stage WS retrieval algorithm;
- WS2_TEST/SST_TRAIN: validation subset for the 2nd-stage WS retrieval algorithm, also used as training subset for the SST retrieval algorithm;
- SST_TEST: validation subset for the SST retrieval algorithm;
- UNCERT_TRAIN: training subset for the SST uncertainty retrieval algorithm;
and
- UNCERT_TEST: validation subset for the SST uncertainty retrieval algorithm.

2.2.2 AMSR-E orbital data

For producing a climate data record of MWR SST, the spatially resampled L2A swath data product AMSR-E V12 (Ashcroft and Wentz, 2013), produced by Remote Sensing Systems (RSS) and distributed by NASA's National Snow and Ice Data Center (NSIDC: https://nsidc.org/data/ac_l2a), is used as input. Here we used the brightness temperatures resampled to the 6.9 GHz resolutions. Hence the observations have a resolution footprint of 75 x 43 km, however, the data is distributed as a dataset with a spatial grid resolution of 10 km at all latitudes. Auxiliary data include NWP data from ERA-5. The Generalized Bayesian Cloud Screening (GBCS) software package (Merchant et al., 2008) is used to interpolate the NWP data to the satellite raster.

2.2.3 AMSR2 orbital data

The spatially resampled AMSR2 L1R version 2 swath data product: Dataset of Brightness Temperature Modified Using the Antenna Pattern Matching Technique (Maeda et al., 2016) is used for producing the MWR SST CDR. This product contains similar spatially resampled brightness temperatures to the AMSR-E dataset. NWP data from ERA-Interim are used as auxiliary data. As for the AMSR-E processing, the GBCS software package is used to interpolate the auxiliary data to the satellite raster.

Algorithm Theoretical Basis Document D2.1 v1

2.2.4 Preprocessing

Preprocessing of the input data is necessary before running the regression model. Different fields need to be calculated depending on the input data used. For AMSR-E and AMSR2, the following fields need to be calculated:

- Relative angle between satellite azimuth angle and wind direction (ϕ_{REL}), which is calculated by the following expression

$$\phi_{REL} = \phi_{SAT} - \phi_{WD}$$

where ϕ_{SAT} is the satellite azimuth angle, relative to north, and ϕ_{WD} is the wind direction, relative to north, that the wind is blowing toward.

- Sun glint angle (ϕ_{SGA}), which is calculated by the following expression

$$\phi_{SGA} = \arccos(\sin(\theta_{SOL}) \cdot \sin(\theta_{SAT}) \cdot \cos(\phi_{REL} + 180) + \cos(\theta_{SOL}) \cdot \cos(\theta_{SAT}))$$

where θ_{SOL} and θ_{SAT} are the solar zenith angle and satellite zenith angle, respectively, and ϕ_{REL} is the relative azimuth angle between solar azimuth angle (ϕ_{SOL}) and satellite azimuth angle (ϕ_{SAT}).

2.3 Algorithm description

2.3.1 WS retrieval algorithm

A regression-based retrieval algorithm is used to retrieve WS given satellite brightness temperatures and NWP fields. The WS retrieval algorithm described here is a two-step multiple linear regression model. In the first step, an initial estimate of WS is retrieved using a global retrieval algorithm, i.e. one set of regression coefficients is used for all wind speeds. In the second stage, a final WS is retrieved using specialized WS algorithms, i.e. the algorithm is trained to perform well over restricted WS domains.

2.3.1.1 1st-step: Global retrieval algorithm

An initial estimate of wind speed (WS_a) is obtained through the use of a global regression model, where the regression coefficients are obtained through training on the WSS1_TRAIN subsets. The WS retrieval algorithm is inspired by the NOAA AMSR-2 WS retrieval algorithm (Chang et al., 2015) and expresses WS in terms of brightness temperature (T_B) and incidence angle (θ_{EIA})

$$WS_a = a_0 + \sum_{i=1}^{10} (a_{1i}t_i + a_{2i}t_i^2) + a_3\theta \quad (1)$$

Algorithm Theoretical Basis Document D2.1 v1

where

$$t_i = T_{Bi} - 150, \quad \text{for all channels except the 23.6 GHz channels} \quad (2)$$

$$t_i = \ln(290 - T_{Bi}), \quad \text{for the two 23.6 GHz channels} \quad (3)$$

$$\theta = \theta_{EIA} - 55 \quad (4)$$

The coefficients a_0 , a_1 , a_2 and a_3 are regression coefficients, referred to as $\mathbf{B}_{\text{global}}$, determined through use of a training dataset, the summation index i represents the summation over 10 AMSR-E channels; 6.9, 10.7, 18.7, 23.6 and 36.5 GHz (dual polarization), and T_{Bi} denotes the brightness temperature for the i th channel.

2.3.1.2 2nd-step: Specialized WS retrieval algorithms

In the second stage, a final wind speed (WS_r) is retrieved through the use of a specialized WS regression model, using the same retrieval algorithm as in the first step. Regression coefficients are derived through training on subsets of the WS2_TRAIN subsets, defined for restricted reference WS intervals. Retrieved WS from the first stage (WS_a) is used to determine the correct WS bin, from which regression coefficients are selected to perform the WS retrieval. The specialized algorithms are derived for reference wind speeds in the interval 0 to 20 ms^{-1} , with a bin size of 1 ms^{-1} , giving a total of 20 specialized WS algorithms, which takes the form

$$WS_{rk} = b_{0k} + \sum_{i=1}^{10} (b_{1ik}t_i + b_{2ik}t_i^2) + b_{3k}\theta \quad (5)$$

where

$$t_i = T_{Bi} - 150, \quad \text{for all channels except the 23.6 GHz channels} \quad (6)$$

$$t_i = \ln(290 - T_{Bi}), \quad \text{for the two 23.6 GHz channels} \quad (7)$$

$$\theta = \theta_{EIA} - 55 \quad (8)$$

where k denotes the reference WS. The coefficients b_0 , b_1 , b_2 and b_3 are regression coefficients, referred to as \mathbf{B}_{WS} , determined through use of a training dataset. The final retrieved WS is found by performing a linear interpolation between WS_{rk} and the WS retrieved using the closest neighboring WS algorithm

Algorithm Theoretical Basis Document D2.1 v1

$$WS_r = \sum_{k=k_0}^{k_0+1} w_{k-k_0} WS_{rk} \quad (9)$$

where

$$w_0 = 1 - \alpha \quad (10)$$

$$w_1 = \alpha \quad (11)$$

$$\alpha = \frac{WS_a}{\Delta k} - k_0 \quad (12)$$

$$k_0 = \text{floor}\left(\frac{WS_r}{\Delta k}\right) \quad (13)$$

and $\Delta k = 1 \text{ ms}^{-1}$ is the WS bin size. Using constant wind intervals ensures that intervals with few observations get adequate weight.

2.3.2 SST retrieval algorithm

We use a regression model to retrieve SST given satellite brightness temperatures and NWP fields. The retrieval algorithm is inspired by the RSS AMSR-E SST retrieval algorithm (Wentz and Meissner, 2007) and expresses SST in terms of brightness temperature (T_B), incidence angle (θ_{EIA}), retrieved wind speed (WS_r) and the relative angle between satellite azimuth angle and NWP wind direction (φ_{REL}). The SST retrieval algorithm uses 12 brightness temperature channels; 6.9, 10.7, 18.7, 23.8, 36.5 and 89.0 (vertical and horizontal polarization). The SST retrieval algorithm described here is a two-step multiple linear regression model with specialized regression algorithms. With “specialized” we mean that the algorithm is trained to perform well over specialized domains. In the first stage, the algorithm is trained to perform well over restricted latitude domains and for ascending and descending orbit, respectively, whereas in the second stage, the algorithm is trained to perform well over restricted SST and WS domains.

2.3.2.1 1st-step: Specialized latitude and ascending/descending retrieval algorithms

An initial estimate of SST (SST_a) is obtained through the use of a specialized orbit and latitude regression retrieval algorithm. Regression coefficients are obtained through training on subsets of the SST_TRAIN subsets, defined for restricted reference latitude intervals and for ascending and descending orbits. Latitude and ascending or descending orbit are used to determine the correct latitude and orbit bin, from which regression coefficients are selected to perform the SST retrieval. The specialized algorithms are derived for reference

Algorithm Theoretical Basis Document D2.1 v1

latitudes in the interval -90 to 90°, with a bin size of 2°, and ascending (0) or descending (1) orbit, giving a total of 182 specialized latitude and orbit algorithms, which takes the form

$$SST_{alm} = c_{0lm} + \sum_{i=1}^{12} (c_{1ilm} t_i + c_{2ilm} t_i^2) + c_{3lm} \theta + c_{4lm} WS_r + \sum_{j=1}^2 (c_{5jlm} \cos j\varphi_{REL} + c_{6jlm} \sin j\varphi_{REL}) \quad (14)$$

where l denotes the reference latitude and m denotes the reference orbit, which ranges from 0 (descending) to 1 (ascending). The coefficients $c_0, c_1, c_2, c_3, c_4, c_5$ and c_6 are regression coefficients, referred to as $\mathbf{B}_{LAT,ORB}$, determined through use of the SST_TRAIN subset. The initial estimate of SST is found by performing a linear interpolation between SST_{alm} and the SST retrieved using the closest neighboring latitude and orbit algorithm

$$SST_a = \sum_{l=l_0}^{l_0+1} w_{l-l_0} SST_{alm} \quad (15)$$

where

$$w_0 = 1 - \alpha \quad (16)$$

$$w_1 = \alpha \quad (17)$$

$$\alpha = \frac{\phi_{LAT} - l_0}{\Delta l} \quad (18)$$

$$l_0 = \text{floor}\left(\frac{\phi_{LAT}}{\Delta l}\right) \quad (19)$$

and ϕ_{LAT} denotes latitude and $\Delta l = 2^\circ$ is the latitude bin size.

2.3.2.2 2nd-step: Specialized SST and WS retrieval algorithms

In the second stage, final SST (SST_r) is retrieved through the use of a specialized SST and WS regression model. Regression coefficients are obtained through training on subsets of the SST_TRAIN subset, defined for restricted reference SST and WS intervals. Retrieved wind speed (WS_r) and the retrieved SST from the first stage (SST_a) are used to determine the correct SST and WS bin, from which regression coefficients are selected to perform the SST retrieval. The specialized algorithms are derived for reference SSTs in the interval -2

Algorithm Theoretical Basis Document D2.1 v1

to 34°C, with a bin size of 2°C, and reference WS in the interval 0 to 20 ms⁻¹, with a bin size of 2 ms⁻¹, giving a total of 209 specialized SST and WS algorithms, which takes the form

$$\begin{aligned}
 SST_{rnp} = & d_{0np} + \sum_{i=1}^{12} (d_{1inp}t_i + d_{2inp}t_i^2) + d_{3np}\theta + d_{4np}WS_r \\
 & + \sum_{j=1}^2 (d_{5jnp} \cos j\varphi_{REL} + d_{6jnp} \sin j\varphi_{REL})
 \end{aligned} \tag{20}$$

where n denotes the reference SST and p denotes the reference WS. The coefficients d_0 , d_1 , d_2 , d_3 , d_4 , d_5 and d_6 are regression coefficients, referred to as $\mathbf{B}_{SST,WS}$, and are determined through use of subsets of the SST_TRAIN subsets. The final retrieved SST is found by performing a bi-linear interpolation between SST_{rnp} and the SSTs retrieved using the three closest neighboring SST and WS algorithms

$$SST_r = \sum_{n=n_0}^{n_0+1} \sum_{p=p_0}^{p_0+1} \omega_{n-n_0, p-p_0} SST_{rnp} \tag{21}$$

where

$$\omega_{0,0} = (1 - \beta) \cdot (1 - \gamma) \tag{22}$$

$$\omega_{1,0} = \beta \cdot (1 - \gamma) \tag{23}$$

$$\omega_{0,1} = (1 - \beta) \cdot \gamma \tag{24}$$

$$\omega_{1,1} = \beta \cdot \gamma \tag{25}$$

$$\beta = \frac{SST_a}{\Delta n} - n_0, \quad \gamma = \frac{WS_r}{\Delta p} - p_0 \tag{26}$$

$$n_0 = \text{floor}\left(\frac{SST_a}{\Delta n}\right), \quad p_0 = \text{floor}\left(\frac{WS_r}{\Delta p}\right) \tag{27}$$

and $\Delta n = 2^\circ\text{C}$ and $\Delta p = 2 \text{ ms}^{-1}$ is the SST and WS bin size, respectively.

Algorithm Theoretical Basis Document D2.1 v1

2.3.3 RFI filter

The baseline SST retrieval algorithm, described in Section 2.3.2, uses 12 brightness temperature channels; 6.9, 10.7, 18.7, 23.6, 36.5 and 89.0 GHz (dual polarization). For detection of RFI, two additional SST retrieval algorithms were defined; the -10GHz and -18GHz algorithms. Typically, RFI frequencies are very specific in terms of the frequency and will thus only enter one channel. Using alternative retrievals with a different channel constellation are useful for filter for RFI. These retrievals are formulated exactly as the baseline algorithm, with the exception that they use only 10 brightness temperature channels; the same as the baseline algorithm minus the 10 GHz channel (-10GHz algorithm) and minus the 18 GHz channels (-18GHz algorithm). As for the baseline retrieval algorithm, WS is first retrieved using the two-step regression model with the baseline WS algorithm and then the two-step regression model is used to retrieve SST.

A new RFI mask, based on the two additional retrieval algorithms, has been developed. A 3σ -filter on the difference between retrieved SST for the two additional algorithms, -10GHz and -18GHz, and the baseline algorithm is used to detect RFI. Data is flagged if any of the following expressions is true

$$\left| (SST_{r,\text{baseline}} - SST_{r,-10\text{GHz}}) - \mu_{-10\text{GHz}} \right| > 3\sigma_{-10\text{GHz}} \quad (28)$$

$$\left| (SST_{r,\text{baseline}} - SST_{r,-18\text{GHz}}) - \mu_{-18\text{GHz}} \right| > 3\sigma_{-18\text{GHz}} \quad (29)$$

where $SST_{r,-10\text{GHz}}$, $SST_{r,-18\text{GHz}}$ and $SST_{r,\text{baseline}}$ are the final retrieved SST using the -10GHz, -18GHz and baseline algorithms, respectively. $\mu_{-10\text{GHz}}$ and $\mu_{-18\text{GHz}}$ denote the mean of the difference $SST_{r,-10\text{GHz}} - SST_{r,\text{baseline}}$ and $SST_{r,-18\text{GHz}} - SST_{r,\text{baseline}}$, respectively, whereas $\sigma_{-10\text{GHz}}$ and $\sigma_{-18\text{GHz}}$ denote the standard deviation of the corresponding differences. The mean and standard deviation of differences used are shown in Table 1.

Table 1: Mean and standard deviation of differences for retrieved SSTs using the -10GHz and -18GHz algorithm minus baseline retrieved SST.

Sensor	$\mu_{-10\text{GHz}}$ (K)	$\mu_{-18\text{GHz}}$ (K)	$\sigma_{-10\text{GHz}}$ (K)	$\sigma_{-18\text{GHz}}$ (K)
AMSR-E	0.0024	0.0071	0.192	0.138
AMSR2	-0.0087	0.0043	0.170	0.130

Algorithm Theoretical Basis Document D2.1 v1

2.3.4 Uncertainty model

Following the approach within the ESA-CCI SST project (Rayner et al., 2015), the total uncertainty for the retrieved SST can be divided into three independent components; a random component (ε_{random}) a local systematic component (ε_{local}) and a global systematic component (ε_{global}). The total uncertainty is given by

$$\varepsilon_{SST_r} = \sqrt{\varepsilon_{random}^2 + \varepsilon_{local}^2 + \varepsilon_{global}^2} \quad (30)$$

The local systematic uncertainty component and the random uncertainty component are both retrieved through the use of a regression model. The global systematic uncertainty component is assumed to be small and is therefore populated by zeros.

To get an estimate of the local systematic and the random uncertainty components, an NEDT of 0.1 K (Wentz and Meissner, 2000) was propagated through the SST retrieval algorithm and a new set of SSTs were obtained ($SST_{r,rand}$). The data in the UNCERT_TRAIN subset was then pre-binned for retrieved SST, retrieved WS, latitude and solar zenith angle. Two standard deviations were calculated

- $\sigma_{\Delta SST_r}$: the standard deviation of the SST_r minus in situ SST difference; and
- $\sigma_{\Delta SST_r,md}$: the standard deviation of the SST_r minus $SST_{r,md}$ difference.

The first standard deviation, $\sigma_{\Delta SST_r}$, is used to represent local effects, including drifter uncertainty and sampling effect, on the total uncertainty, whereas the second standard deviation, $\sigma_{\Delta SST_r,md}$, is used to represent random and uncorrelated effects.

The same retrieval algorithm is used for both the local systematic uncertainty component and the random uncertainty component. The uncertainties are expressed in terms of baseline retrieved SST (SST_r), retrieved wind speed (WS_r), solar zenith angle (θ_{SOL}) and latitude (ϕ_{LAT})

$$\varepsilon_{SST_r} = e_0 + e_1 SST_r + e_2 SST_r^2 + e_3 WS_r + e_4 WS_r^2 + e_5 \theta_{SOL} + e_6 \theta_{SOL}^2 + \sum_{p=1}^4 \left(e_{7p} \cos \frac{\phi_{LAT}}{p} + e_{8p} \sin \frac{\phi_{LAT}}{p} \right) \quad (31)$$

Algorithm Theoretical Basis Document D2.1 v1

where the coefficients $e_0, e_1, e_2, e_3, e_4, e_5, e_6, e_7$ and e_8 are regression coefficients, denoted $\mathbf{B}_{\text{local}}$ and \mathbf{B}_{rnd} for the local and random uncertainty components, respectively, which are determined through training on the UNCERT_TRAIN subsets. p is a summation index for the harmonic function used for the latitude.

The retrieval algorithm for the random uncertainty component was regressed towards $\sigma_{\Delta\text{SST},\text{rnd}}$. To obtain only the variations due to local systematic effects, the local systematic uncertainty component was regressed towards pre-binned standard deviations where uncorrelated random effects, as well as drifter uncertainty and sampling effects, were excluded, σ_{local} . The drifter uncertainty was set to 0.2 K and the sampling effect was assumed to be mainly spatial (Høyer et al., 2012) and was therefore estimated by calculating the pixel-to-footprint variability for one year of GHRSSST Level 4 DMI_OI Global Foundation Sea Surface Temperature Analysis (DMI, 2007).

2.3.5 L2P flags and quality levels

The MWR SST CDR retrievals follow the GHRSSST GDS 2.0 data specification (GHRSSST Science Team, 2010) for L2P and each retrieval was assigned a quality level to denote the quality of the retrieval. The definition of quality levels, together with corresponding checks and thresholds, are shown in Table 2.

The quality of the individual SST retrievals is represented by a quality level, ranging from 0 to 5. Quality level 0 denotes the lowest quality indicator level, which is assigned if no data is retrieved. Quality level 1 is the lowest quality level for retrievals, indicating retrievals of bad quality which should not be used, whereas quality level 5 is the highest quality level, only given to retrievals with the best quality. Retrievals are assigned quality level 1 if the input data is of bad quality or if the retrieval is compromised, e.g. due to atmospheric and surface effects. The following criteria decide if a retrieval is of quality level 1:

AMSR-E scan quality or channel quality indicates bad satellite data.

- Any brightness temperature is outside the normal range ($0 \text{ K} < T_B < 320 \text{ K}$).
- Sea ice contamination.
- Coastal contamination.
- Contamination due to RFI (masked according to section 2.3.3).
- Rain contamination ($T_{B18V} \geq 240 \text{ K}$).
- Sun glitter contamination ($\varphi_{SGA} \leq 25^\circ$).

Algorithm Theoretical Basis Document D2.1 v1

- Cases where the atmospheric contribution exceeds the information from the surface, i.e. if the difference between the horizontal and vertical polarization brightness temperatures for channel 18-36 GHz is negative.
- The retrieved WS is outside the accepted range ($0 \text{ ms}^{-1} \leq WS_r \leq 20 \text{ ms}^{-1}$).
- The retrieved SST is outside the accepted range ($-2 \text{ }^\circ\text{C} \leq SST_r \leq 35 \text{ }^\circ\text{C}$).
- The retrieved SST deviates with more than 10 °C from a background SST.

Quality level 2, which denotes the worst-quality yet usable retrievals, is assigned to retrievals with a total uncertainty greater than 1. In addition, the proximity to sea ice and land is also used to determine if the retrieval is of quality level 2. If the distance to sea ice is less than 200 km or if the distance to land is less than 40 km, the retrieval is classified as being of quality level 2. Quality level 3 to 5 are determined based solely on the retrieved total SST uncertainty. If the SST uncertainty is in the range 0.5-1 K, the retrieval is assigned quality level 3 (low quality), if it is in the range 0.35-0.5 K, the retrieval is assigned quality level 4 (acceptable quality) and if the uncertainty is 0.35 K or smaller, the retrieval is assigned quality level 5 (best quality).

Table 2. Quality levels with corresponding checks and thresholds.

#	Description	Checks and thresholds
0	No data	
1	Bad data	Quality controls and various atmospheric and surface effects
2	Worst-quality usable data	<ul style="list-style-type: none"> ■ $\epsilon_{SST_r} \geq 1$ ■ Proximity to sea ice ■ Proximity to land
3	Low quality	$0.50 < \epsilon_{SST_r} < 1$
4	Acceptable quality	$0.35 < \epsilon_{SST_r} \leq 0.50$
5	Best quality	$\epsilon_{SST_r} \leq 0.35$

Algorithm Theoretical Basis Document D2.1 v1

2.3.6 Regression setup

The setup of the DMI regression model with the different processes and steps is illustrated in Alerskans et al., 2020. The observation loop is started for each satellite pixel by reading in the satellite orbital data and the auxiliary data (NWP fields). The next step is to read in all regression coefficients, whereupon the retrieval process can begin. First, the 1st-stage global WS retrieval algorithm is used to retrieve an initial estimate of WS (WS_a). Thereafter, WS_a is used to select regression coefficients for the second step of the WS retrieval algorithm, \mathbf{B}_{WS} . Subsequently, final retrieved WS (WS_r) is computed using the specialized WS retrieval algorithms. Next, the two-step SST retrieval algorithm is performed for the three SST retrieval algorithms; baseline, -10GHz and -18GHz, with the algorithm loop being initialized with the baseline algorithm ($i = 0$). First, latitude and ascending/descending orbit are used to select regression coefficients for the 1st-stage SST retrieval algorithm, $\mathbf{B}_{algo_i,LAT,ORB}$, whereupon the specialized latitude and ascending/descending retrieval algorithm is used to compute an initial estimate of retrieved SST ($SST_{a,algo_i}$). For the final SST retrieval, the initially retrieved SST ($SST_{a,algo_i}$) and the final retrieved WS (WS_r) are used to select regression coefficients, $\mathbf{B}_{algo_i,SST,WS}$. In the following step, a final retrieved SST ($SST_{r,algo_i}$) is computed using the specialized SST and WS retrieval algorithm. The algorithm loop is then performed for the two additional algorithms, -10GHz ($i = 1$) and -18GHz ($i = 2$). When exiting the algorithm loop, RFI is detected and masked using the new proposed RFI mask. In the next step, regression coefficients \mathbf{B}_{local} and \mathbf{B}_{rnd} for the uncertainty retrieval algorithm are used to compute the uncertainty for the baseline-retrieved SST (ε_{SST_r}). Thereafter, the retrieval is assigned a quality level and flagged according to the quality level and L2P flagging criteria described in section 2.3.5. Finally, the baseline retrieved SST and uncertainty is saved together with the L2P flags and quality levels.

Algorithm Theoretical Basis Document D2.1 v1

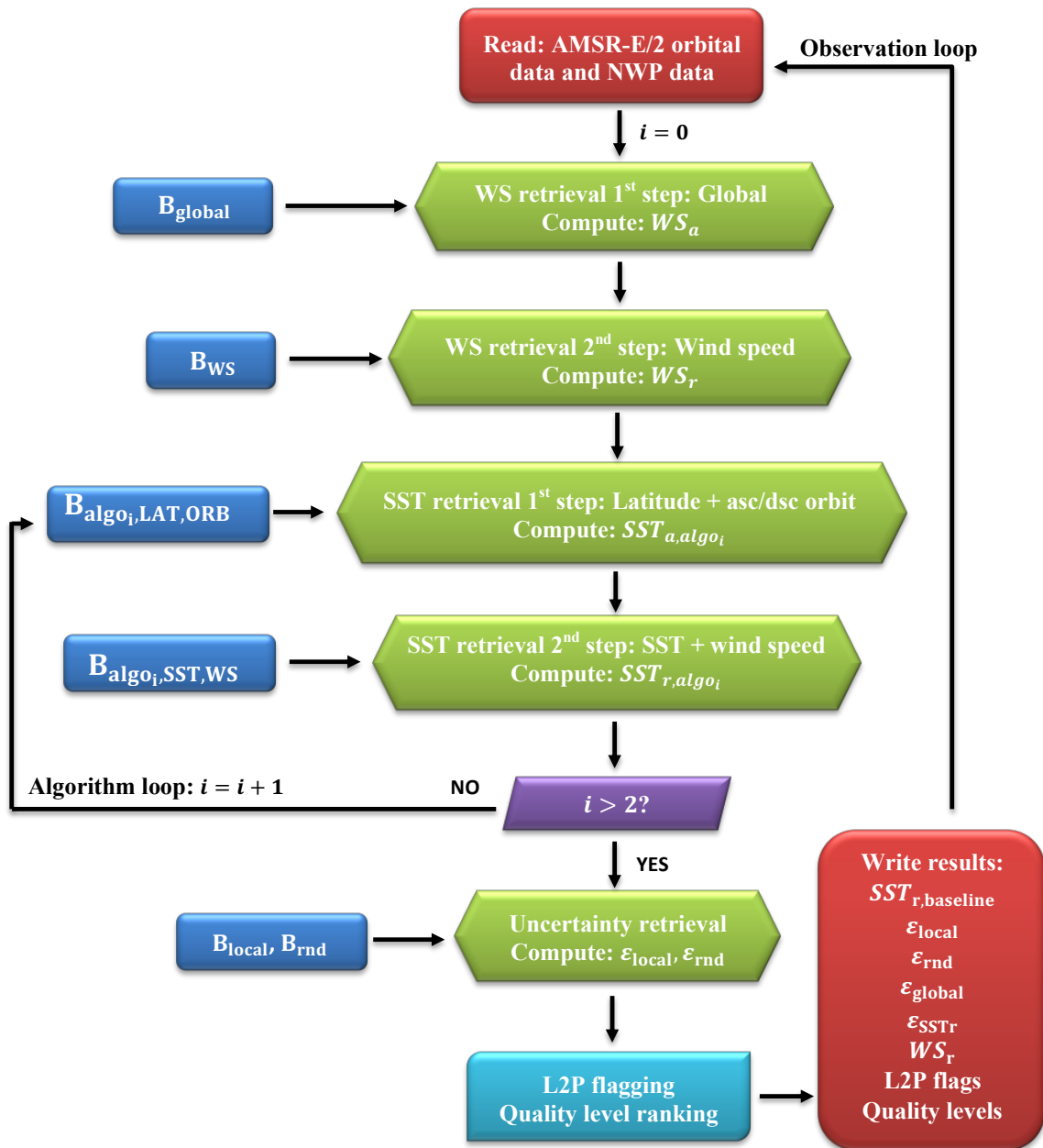


Figure 1: Setup of the DMI regression model for MWR SST retrievals using AMSR-E/2 orbital data as input. i denotes the algorithm used to retrieve SST; baseline ($i = 0$), -10GHz ($i = 1$) and -18GHz ($i = 2$).

2.4 Output data

The outputs from the regression retrieval algorithm are

- Sea surface temperature ($SST_{r,baseline}$), in Kelvin.
- Local systematic uncertainty component (ϵ_{local}), in Kelvin.
- Random uncertainty component (ϵ_{rnd}), in Kelvin.

Algorithm Theoretical Basis Document D2.1 v1

- Global systematic uncertainty component (ϵ_{global}), in Kelvin.
- Total SST uncertainty (ϵ_{SSTr}), in Kelvin.
- Wind speed (WS_r), in ms^{-1} .
- L2P flags.
- Quality levels.

Algorithm Theoretical Basis Document D2.1 v1

3. OPTIMAL ESTIMATION TUNING

3.1 Introduction

Merchant et al. (2020) have presented a theory for using reference SSTs to determine bias correction and error covariance parameters for OE in the case of infrared SST retrieval. OE SST is used for AVHRR 2- and 3-channel retrievals in SST CCI in order to achieve known and satisfactory levels of SST sensitivity, by explicitly using prior NWP to address deficits in the window-channel information about the atmospheric influence on brightness temperatures. The OE solution is obtained in one step, because the retrieval context is adequately linear.

OE for microwave SST differs from the infrared case in the following respects:

- a larger state vector is needed, comprising (at least) SST, total column water vapour (TCWV), wind speed (U) and cloud liquid water (CLW)
- a larger observation vector is needed, comprising two or more channels (minimum: ~7 GHz and ~10 GHz) in two polarisations
- a larger and non-linear dependence of the signals (BTs) on wind speed via surface emissivity effects, meaning that in the microwave OE is applied as a moderately non-linear retrieval, requiring some iterations rather than a single step

The theory developed for the infrared therefore need to be extended to the case of moderately non-linear OE. This theoretical extension is presented here.

3.2 Moderately non-linear optimal estimation

The approach presented here is obtained from Rodgers (2000). The prior for a retrieval is a reduced state vector \mathbf{z}_a . Let the k^{th} iterative retrieved state vector be \mathbf{z}_k , then Newtonian iteration to the OE solution means that

$$\mathbf{z}_k = \mathbf{z}_a + \tag{32}$$

$$\left(\mathbf{K}_{k-1}^T \mathbf{S}_\epsilon^{-1} \mathbf{K}_{k-1} + \mathbf{S}_a^{-1} \right)^{-1} \mathbf{K}_{k-1}^T \mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{z}_{k-1}) + \mathbf{K}_{k-1}(\mathbf{z}_{k-1} - \mathbf{z}_a))$$

where the iteration starts from $k = 1$ with $\mathbf{z}_{k-1} = \mathbf{z}_0 = \mathbf{z}_a$. Here $\mathbf{F}(\mathbf{z}_{k-1})$ is the forward model calculated from the full state vector corresponding to the previous reduced state vector estimate; \mathbf{K}_{k-1} is the corresponding set of partial derivatives of the forward model brightness temperature with respect to the reduced state vector elements; \mathbf{y} is the observation vector; \mathbf{S}_a^{-1} is the inverse of the prior error covariance matrix; and \mathbf{S}_ϵ^{-1} is the inverse of the error covariance in the observation-simulation difference (i.e., includes

Algorithm Theoretical Basis Document D2.1 v1

instrument errors and forward modelling errors). 0.32 is iterated until \mathbf{z}_{k-1} is within linear range of the final solution (convergence is reached).

3.3 Extension to parameter estimation with reference observations

The method of Merchant et al. (2020) is here modified to an equivalent formulation that uses reference observations as observations (as elements in the observation vector) rather than (as previously) using them to modify the prior state vector. This is a practical simplification.

The concept of OE tuning is that the reference observations are assumed, by definition, to be unbiased (although with some uncertainty in each reference observation), and that they enable inference of useful bias corrections (to observations and to the prior). Additionally, self-consistency requirements on the retrieval and the outcomes are available to improve estimation of the error covariance assumptions.

All biases and error covariance matrices can be derived as a function of relevant quantities, which is here done by binning data and assuming piecewise linear interpolation between bins. This piece-wise linear aspect is not explicitly present in the equations for simplicity: the equations describe effectively the derivation of parameters for a single bin.

The basis of the tuning is to extend the observation vector with the reference observations, $\boldsymbol{\rho}$:

$$\tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{y} \\ \boldsymbol{\rho} \end{bmatrix} \quad (33)$$

and to extend the retrieved state to include bias parameters, for the bias in the prior state ($\boldsymbol{\gamma}'$) and the observation ($\boldsymbol{\beta}$):

$$\tilde{\mathbf{z}} = \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\gamma}' \\ \boldsymbol{\beta} \end{bmatrix} \quad (34)$$

$\boldsymbol{\gamma}'$ may contain prior state bias correction values for a shorter list of elements than the full state vector (we may only feel we can bias correct some aspects), which is indicated by the prime. $\boldsymbol{\gamma}$ refers to a vector the same length as \mathbf{z}_a in which the elements not in $\boldsymbol{\gamma}'$ are padded with zeros, so that when we write " $\mathbf{z}_a + \boldsymbol{\gamma}$ " it is clear that the non-bias-corrected elements in \mathbf{z}_a are unchanged.

Algorithm Theoretical Basis Document D2.1 v1

The assumed property of reference observations means

$$\boldsymbol{\rho} = \mathbf{z}'' + \boldsymbol{\epsilon}_\rho ; \langle \boldsymbol{\epsilon}_\rho \rangle = \mathbf{0} \quad (35)$$

where \mathbf{z}'' is the true state of a subset (indicated by the double prime) of the elements of the state vector \mathbf{z} (those elements for which we have reference observations), and $\boldsymbol{\epsilon}_\rho$ is the set of errors for a given instance of reference observations. These errors are assumed to be zero mean by our definition of reference observations (and will usually be independent of each other too, which is hereafter assumed). Unbiasedness implies that the reference observations are fully sensitive to changes in the aspect of the state they measure.

The extended forward model evaluated for the prior state, given bias correction parameters, is

$$\tilde{\mathbf{F}}(\tilde{\mathbf{z}}_a) = \begin{bmatrix} \mathbf{F}(\mathbf{z}_a + \boldsymbol{\gamma}) + \boldsymbol{\beta} \\ \mathbf{P}(\mathbf{z}_a + \boldsymbol{\gamma}) \end{bmatrix} \quad (36)$$

Here \mathbf{P} is our forward model of the reference observations given a state vector, and is simply $\mathbf{P}(\mathbf{z}) = \mathbf{z}''$.

The K matrix is "doubly extended" since it has dimensions of the extended observation vector times the extended state vector.

$$\tilde{\mathbf{K}} = \begin{bmatrix} \left. \frac{\partial \mathbf{F}}{\partial \mathbf{z}} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} & \left. \frac{\partial \mathbf{F}}{\partial \boldsymbol{\gamma}'} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} & \left. \frac{\partial \mathbf{F}}{\partial \boldsymbol{\beta}} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} \\ \left. \frac{\partial \mathbf{P}}{\partial \mathbf{z}} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} & \left. \frac{\partial \mathbf{P}}{\partial \boldsymbol{\gamma}'} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} & \left. \frac{\partial \mathbf{P}}{\partial \boldsymbol{\beta}} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} \end{bmatrix} \quad (37)$$

Where the prior state bias corrections are not too large, we can use the approximation

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{z}} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} = \left. \frac{\partial \mathbf{F}}{\partial \mathbf{z}} \right|_{\mathbf{z}_a}.$$

Given the form of the forward model, clearly $\left. \frac{\partial \mathbf{F}}{\partial \boldsymbol{\gamma}'} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} = \left. \frac{\partial \mathbf{F}}{\partial \mathbf{z}'} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}}$ and $\left. \frac{\partial \mathbf{F}}{\partial \boldsymbol{\beta}} \right|_{\mathbf{z}_a + \boldsymbol{\gamma}} = \mathbf{I}$.

Algorithm Theoretical Basis Document D2.1 v1

$\frac{\partial P}{\partial z'} \Big|_{z_a+\gamma} = I$ (the "fully sensitive" assumption), while $\frac{\partial P}{\partial z} \Big|_{z_a+\gamma}$ is similar but with additional columns of zeros corresponding to the elements that don't have reference measurements -- we write this zero-padded identity matrix as $\frac{\partial P}{\partial z} \Big|_{z_a+\gamma} = I''$.

The i - j th element of $\frac{\partial P}{\partial \gamma'} \Big|_{z_a+\gamma}$ is 1 wherever the i th element of ρ is an observation of the state vector element bias-corrected by the j th element of γ' and is zero elsewhere. It can be expected that all state vector elements with reference observations will appear in γ' (i.e., their prior will be bias corrected), but since the anchoring from the references may help de-bias some other elements too, γ' may be longer than ρ . Thus, this is another extended identity matrix, and we write this matrix as $\frac{\partial P}{\partial \gamma'} \Big|_{z_a+\gamma} = I'$.

$$\frac{\partial \rho}{\partial \beta} \Big|_{z_a+\gamma} = \mathbf{0}.$$

Thus, combining these considerations:

$$\tilde{\mathbf{K}} \cong \begin{bmatrix} \frac{\partial F}{\partial z} \Big|_{z_a} & \frac{\partial F}{\partial z'} \Big|_{z_a} & I \\ I'' & I' & \mathbf{0} \end{bmatrix} \quad (38)$$

The extended observation vector requires an extended observation error covariance matrix, given by

$$\tilde{\mathbf{S}}_{\epsilon} = \begin{bmatrix} \mathbf{S}_{\epsilon} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{\rho} \end{bmatrix} \quad (39)$$

Algorithm Theoretical Basis Document D2.1 v1

where S_p contains the known (or at least, characterised) error covariances of the reference data, each of which are independent of all other observations, and hence S_p is diagonal and the off-diagonal blocks are zero matrices.

3.4 Bias correction parameters

Bias correction parameters are obtained by doing the extended retrieval formulated above across iterations i across randomly drawn satellite-reference matches until the bias correction parameters converge. This iteration over matches is an outer loop, with the iterative retrieval for each given match as an inner loop. This is expressed below:

$$\tilde{\mathbf{z}}_k^i = \tilde{\mathbf{z}}_a + \quad (40)$$

$$\left(\tilde{\mathbf{K}}_{k-1}^T \tilde{\mathbf{S}}_\epsilon^{-1} \tilde{\mathbf{K}}_{k-1} + \tilde{\mathbf{S}}^{-1} \right)^{-1} \tilde{\mathbf{K}}_{k-1}^T \tilde{\mathbf{S}}_\epsilon^{-1} \left(\tilde{\mathbf{y}} - \tilde{\mathbf{F}}(\tilde{\mathbf{z}}_{k-1}^{i-1}) + \tilde{\mathbf{K}}_{k-1}(\tilde{\mathbf{z}}_{k-1}^{i-1} - \tilde{\mathbf{z}}_a) \right)$$

$$\tilde{\mathbf{z}}_{k-1}^{i-1} = \begin{bmatrix} \mathbf{z}_{k-1} + \boldsymbol{\gamma}_{i-1} \\ \boldsymbol{\gamma}'_{i-1} \\ \boldsymbol{\beta}_{i-1} \end{bmatrix}$$

$$\tilde{\mathbf{K}} = \begin{bmatrix} \left. \frac{\partial \mathbf{F}}{\partial \mathbf{z}} \right|_{\mathbf{z}_{k-1} + \boldsymbol{\gamma}_{i-1}} & \left. \frac{\partial \mathbf{F}}{\partial \mathbf{z}'} \right|_{\mathbf{z}_{k-1} + \boldsymbol{\gamma}_{i-1}} & \mathbf{I} \\ \mathbf{I}'' & \mathbf{I}' & \mathbf{0} \end{bmatrix}$$

$$\tilde{\mathbf{S}} = \begin{bmatrix} \mathbf{S}_a + \mathbf{S}_{\boldsymbol{\gamma}_{i-1}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{\boldsymbol{\gamma}'_{i-1}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{S}_{\boldsymbol{\beta}_{i-1}} \end{bmatrix}$$

The index i is updated only after the iterations of k conclude, and then the new parameter estimates are passed on to the next randomly selected match $i + 1$.

The error covariance matrix for the retrieval is:

Algorithm Theoretical Basis Document D2.1 v1

In implementing this with a matchup database, we get a covariance matrix for the retrieval which is

$$\mathbf{S}_{\hat{z}_i} = \left(\tilde{\mathbf{K}}^T \mathbf{S}_\epsilon^{-1} \tilde{\mathbf{K}} + \tilde{\mathbf{S}}^{-1} \right)^{-1} = \begin{bmatrix} \mathbf{S}_{z_i} & \mathbf{A} & \mathbf{B} \\ \mathbf{A}^T & \mathbf{S}_{y_i} & \mathbf{C}_i \\ \mathbf{B}^T & \mathbf{C}_i^T & \mathbf{S}_{\beta_i} \end{bmatrix} \quad (41)$$

When iterating across the matches, the matrix $\begin{bmatrix} \mathbf{S}_{y_i} & \mathbf{C}_i \\ \mathbf{C}_i^T & \mathbf{S}_{\beta_i} \end{bmatrix}$ is passed to the $\tilde{\mathbf{S}}$ of the subsequent match.

3.5 Desroziers estimator for observation-simulation error covariance matrix

Having obtained a set of bias correction parameters from the previous steps, these are fixed and used in a set of retrievals across all the matches. Again, these retrievals are made using the references as additional observations, but the state vector here is not extended, since the bias correction parameters are, for this step, fixed. Thus the tangent linear matrix is here "singly" extended, $\tilde{\mathbf{K}} = \begin{bmatrix} \frac{\partial F}{\partial z} |_{z_a+\gamma} + \frac{\partial \beta}{\partial z} |_{z_a} \\ \mathbf{I}'' \end{bmatrix} \cong \begin{bmatrix} \frac{\partial F}{\partial z} |_{z_a} \\ \mathbf{I}'' \end{bmatrix}$.

Iterative retrieval is undertaken for each match obtaining a converged solution \hat{z} . The Desroziers (Desroziers et al., 2005) estimate for the extended observation-simulation error covariance matrix, $\hat{\mathbf{S}}_\epsilon$ can be formulated as follows:

$$\hat{\mathbf{S}}_\epsilon = \begin{bmatrix} \hat{\mathbf{S}}_\epsilon & \mathbf{D} \\ \mathbf{D}^T & \hat{\mathbf{S}}_\rho \end{bmatrix} = \frac{1}{2} \langle \mathbf{d}_r^o \mathbf{d}_a^{oT} + \mathbf{d}_a^o \mathbf{d}_r^{oT} \rangle \quad (42)$$

$$\mathbf{d}_r^o = \tilde{\mathbf{y}} - \tilde{\mathbf{F}}(\hat{z}) - \langle \tilde{\mathbf{y}} - \tilde{\mathbf{F}}(\hat{z}) \rangle$$

$$\mathbf{d}_a^o = \tilde{\mathbf{y}} - (\tilde{\mathbf{F}}(\hat{z}) - \tilde{\mathbf{K}}_{\hat{z}}(\hat{z} - \mathbf{z}_a)) - \langle \tilde{\mathbf{y}} - (\tilde{\mathbf{F}}(\hat{z}) - \tilde{\mathbf{K}}_{\hat{z}}(\hat{z} - \mathbf{z}_a)) \rangle$$

where the averages are performed across all the matches to references (perhaps in bins if a piece-wise linear dependence for $\hat{\mathbf{S}}_\epsilon$ is required). $\tilde{\mathbf{K}}_{\hat{z}}$ is the tangent linear matrix evaluated at the solution state. The result should have the properties that $\mathbf{D} \sim \mathbf{0}$ and that $\hat{\mathbf{S}}_\rho$ is

Algorithm Theoretical Basis Document D2.1 v1

consistent with what we understand about the uncertainty of the references; these considerations add a check that the system is behaving as expected.

3.6 Desroziers estimate for prior error covariance matrix

Using the fixed bias correction parameters and the updated estimate for the observation-simulation error covariance matrix, the (iterative) retrieval is undertaken again for all the satellite-reference matches. The set of new results is used to evaluate an estimate for the prior error covariance matrix:

$$\hat{\mathbf{S}}_a = \frac{1}{2} \langle (\tilde{\mathbf{K}}_z^T \tilde{\mathbf{K}}_z)^{-1} \tilde{\mathbf{K}}_z^T (\mathbf{d}_a^r \mathbf{d}_a^{oT} + \mathbf{d}_a^o \mathbf{d}_a^{rT}) \tilde{\mathbf{K}}_z (\tilde{\mathbf{K}}_z^T \tilde{\mathbf{K}}_z)^{-1} \rangle \quad (43)$$

$$\mathbf{d}_a^r = \tilde{\mathbf{K}}_z(\hat{\mathbf{z}} - \mathbf{z}_a) - \langle \tilde{\mathbf{K}}_z(\hat{\mathbf{z}} - \mathbf{z}_a) \rangle$$

Where the error covariance matrix has a functional dependence, this is evaluated per bin.

3.7 Convergence

Overall convergence of the derived set of OE parameters can be monitored by ensuring that the statistic

$$\text{su} \left(\langle \hat{\mathbf{S}}_\epsilon + \tilde{\mathbf{K}}_z \hat{\mathbf{S}}_a \tilde{\mathbf{K}}_z^T \rangle^{-1} \langle \mathbf{d}_a^o \mathbf{d}_a^{oT} \rangle - \mathbf{I} \right) \quad (44)$$

decreases as estimate of the OE parameters are iterated. However, it is difficult to identify a value of the metric that would indicate convergence for our purposes. A practical alternative is to monitor also the change in SST as OE parameters are improved: once the SST change is very small (e.g., < 0.01 K) the system is adequately converged.

3.8 Iteration of parameter cycles

The steps outlined in sections 3.4 to 3.7 constitute a single cycle of parameter estimation. A few cycles of parameter estimation are needed to obtain bias corrections parameters and covariance parameters that are fully tuned (meeting the overall convergence criterion) because the bias correction estimates respond to the covariances and vice versa.

3.9 Progress towards exploitation

This theoretical development has been undertaken within the context of SST CCI Phase 3 WP 20. Preliminary evaluation of the approach has been promising, indicating that a useful approach is to simultaneously use skin-adjusted Argo or drifting buoy SSTs and the best

Algorithm Theoretical Basis Document D2.1 v1

available source of reference wind speeds as simultaneous reference observations. However, the preliminary results (not shown, being too preliminary) are not yet mature enough for implementation in the v3 CDR. Given the general interest of the approach, ad hoc further research will be continued with a view to having a plausible route to implementation ready for Phase 4 and the v4 CDR.

Algorithm Theoretical Basis Document D2.1 v1

4. DESERT-DUST RELATED BIASES

4.1 Introduction

As explained in (Merchant et al., 2019), there are unsatisfactory levels of negative bias in the v2 CDR related to infrared retrieval in the presence of desert-dust aerosol, particularly in the single-view retrievals obtained using OE applied to AVHRR instruments. Two steps are necessary to improve this situation: infrared algorithm development so as to better deal with desert-dust aerosol, such as enabling a measure of such aerosol in OE; and ensuring the prior SST used for the OE has low bias in respect of such aerosol. This section addresses the latter point.

The intention for the v3 CDR is to use the v2 analysis SST as prior SST, if that analysis can be adequately corrected for desert-dust related SST bias. One approach to this is to use v2 microwave SSTs, since these are not biased by desert dust, although they have their own complex spatio-temporal bias characteristics. This approach has been pursued within SST CCI Phase 3 WP 20, as explained below.

4.2 Notes on Data

For a self-consistent and spatially completed field of desert-dust aerosol, we use outputs of the Copernicus Atmosphere Monitoring Service (CAMS) reanalysis (Inness et al., 2019). Data were obtained for 2003 to 2017 inclusive, spanning the period for which CCI v2 SSTs are available from microwave radiometers (MWRs). The data were obtained on a grid matching the MWR data, namely, daily, 0.5° latitude-longitude resolution. The possible variables from the re-analysis describing desert dust are: aerosol optical depth at 550 nm; and vertically integrated mass of dust aerosol in three size classes, the coarsest of which covers 9 to 20 μm and is most likely to interact with split-window channel brightness temperatures. Reasoning that desert-dust related biases in SST CCI data are most likely to arise in split-window-based retrievals, the coarse-mode dust mass is used to characterise desert dust prevalence in this analysis.

SSTs estimates that are assumed unbiased by desert dust are provided by the MWR SSTs of CCI v2, which are available on a daily, 0.5° grid. The MWR SSTs come from AMSR-E and AMSR-2 retrievals (Alerskans, 2020). While the MWR SSTs are, given the negligible interaction of microwaves with aerosols of order 10 μm in size, unbiased by desert dust, the retrievals are prone to a variety of seasonally and locally systematic error effects and overall uncertainty ~ 0.5 K.

Algorithm Theoretical Basis Document D2.1 v1

Lastly, the dataset to be corrected is the CCI v2.1 SST analysis, which is a daily 0.05° product spanning 1981 to 2019 (at time of writing). Corrections based on CAMS/MWR data can be derived only on a spatial scale that is x10 the analysis resolution and only on a climatological basis for times outside the period of the CAMS re-analysis.

4.3 Notes on methods

Investigations into the spatio-temporal resolution with which statistically reliable SST adjustments with respect to desert dust could be estimated concluded that the seasonally and locally variable SST biases in the MWR SST product preclude time and space resolved determination of the sensitivity of analysis SST, \hat{x} , with respect to CAMS re-analysis dust mass, m , $s = \frac{\partial \hat{x}}{\partial m}$. Therefore, a global all-year average sensitivity is calculated using all MWR data (both AMSR-E and AMSR-2) when available.

A further factor to consider when estimating s is that the impact of desert dust on analysis SSTs is non-linear because of the interaction of dust loading with cloud detection. The heaviest dust events trigger cloud detection and therefore under these conditions, IR SSTs entering the SST analysis system are missing over such events rather than being present-but-biased, and the bias impact on the analysis will be damped therefore for the most extreme dust events. For this reason, it is appropriate to estimate s by a method that is driven more usual (not the extreme) aerosol loadings.

The procedure used is therefore as follows:

- for each ocean grid location, $p = (i, j)$, we calculate the 5, 25, 75 and 95 percentile levels of dust mass ($m_5, m_{25}, m_{75}, m_{95}$ respectively, each implicitly a function of p), across all days for which MWR SSTs are available
- for each grid location we identify the set $T_{lo}(p) = \{t | m_5 < m(t, p) < m_{25}\}$ of times, t , for which the dust mass $m(t, p)$ satisfies $m_5 < m(t, p) < m_{25}$; these constitute the “low-aerosol” baseline for each location
- similarly, identify $T_{hi}(p) = \{t | m_{75} < m(t, p) < m_{95}\}$, the “high-aerosol” cases
- calculate $m_{lo}(p) = \langle m(t, p) \rangle_{T_{lo}(p)}$, which is the arithmetic average of dust mass for the low-aerosol baseline
- similarly calculate

$$m_{hi}(p) = \langle m(t, p) \rangle_{T_{hi}(p)}$$

$$\delta_{lo}(p) = \langle x_{L4}(t, p) - x_{MWR}(t, p) \rangle_{T_{lo}(p)}$$

$$\delta_{hi}(p) = \langle x_{L4}(t, p) - x_{MWR}(t, p) \rangle_{T_{hi}(p)}$$

Algorithm Theoretical Basis Document D2.1 v1

- not all areas experience a significant variation in dust mass, so identify the set of locations $P = \{p | m_{hi}(p) - m_{lo}(p) > m_{th}\}$; the area of ocean over which dust is routinely significant at some point in the annual cycle is ~5% of the total area, so choose m_{th} to be the 95 percentile of $m_{75}(p)$.
- calculate the typical dust mass sensitivity of the SST analysis as $s = \left\langle \frac{\delta_{hi} - \delta_{lo}}{m_{hi} - m_{lo}} \right\rangle_P$, where a robust mean (median) is used to avoid undue influence of outliers
- use the variability, σ , of this ratio across P as a measure of the uncertainty in applying s to any given place and time, again measuring variability with a robust metric (scaled median absolute deviation from the median, also called “robust standard deviation”)

The sensitivity, s , is then used to transform the CAMS dust mass for a given location and day to an SST correction to apply to the SST analysis, as follows:

$$\delta_x(t, p) = s \cdot m(t, p)$$

Given the definition of s , this correction is added to the existing analysis value of SST.

If the reported SST analysis uncertainty (which doesn't account for the desert dust bias) is $u_x(t, p)$, then the revised uncertainty in the SST analysis thus corrected is taken to be

$$\sqrt{u_x^2 + \sigma^2 m^2}$$

A limitation of this estimate being that the uncertainty in the re-analysis value of dust mass is not explicitly present, although dust mass uncertainty does contribute to the value obtained for σ and it therefore this uncertainty is at least partially accounted for.

For the periods outside of the CAMS re-analysis data (prior to 2003), the re-analysis dust mass, m , must be replaced by the climatological value of the dust mass, $\overline{m}(d, p)$ where d refers to the day of year. The climatology calculation uses the years 2003 and 2017 and a five-day rolling window centred on the target day of year. A simple average is used, which is conservative (in the sense of tending to over-correct) because the dust mass tends to be log-normally distributed. The variability (as a standard deviation) around the climatological value is also calculated and is designated $\sigma_{\overline{m}}$. The uncertainty in the corrected analysis SST in the case that the correction is climatological is

$$\sqrt{u_x^2 + \sigma^2 \overline{m}^2 + s^2 \sigma_{\overline{m}}^2}$$

where the additional term accounts for the additional uncertainty in the correction arising from having to use the climatological rather than specific value.

Algorithm Theoretical Basis Document D2.1 v1

4.4 Notes on results

As an example, the climatological dust mass distribution for day of year 182 is shown in

Figure 2, along with the variability around the climatological value.

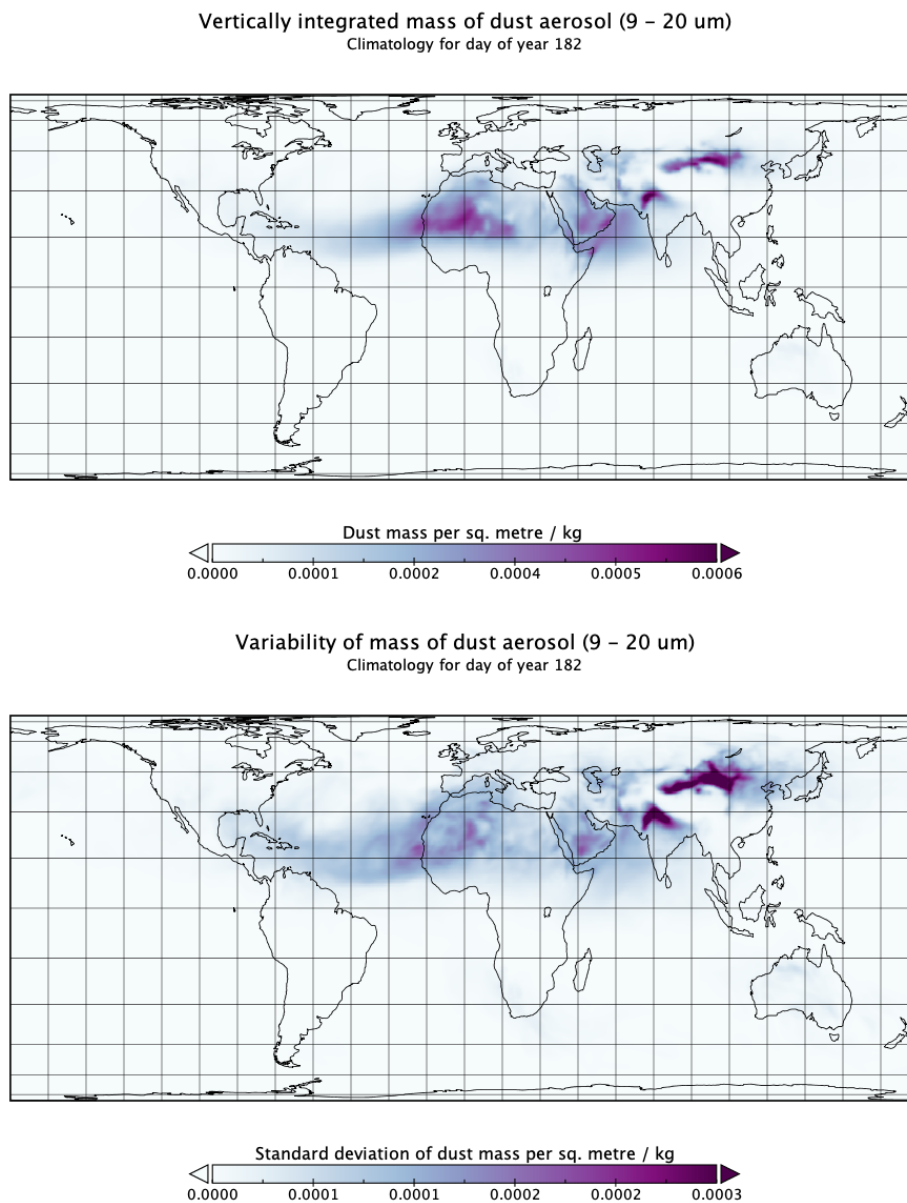


Figure 2: (Upper) Climatological coarse dust mass for day of year 182, in kg m^{-2} , based on CAMS re-analysis for 2003 to 2017. (Lower) Variability (as standard deviation) of coarse dust mass about the climatological value. To calculate the climatological statistics, a centred five-day window over all available years is used.

The highest over-ocean mass loadings are off north west Africa, but dust is also significant in the Arabian Sea. Over areas where dust is climatologically significant, the variability around the climatological value is typically $\sim 50\%$. However, around the areas where dust

Algorithm Theoretical Basis Document D2.1 v1

is climatologically prevalent, elevated variability extends further, indicating these areas are subject to occasional dust events.

The histogram of s obtained is widely spread (Figure 3), with a median value of $-2.96 \text{ K g}^{-1} \text{ m}^2$ and spread (robust standard deviation) of $0.71 \text{ K g}^{-1} \text{ m}^2$. The spread arises partly from differential impacts of a given dust mass on temperature, because of variations in aerosol height in particular; but a significant contribution to the spread arises also from seasonally and locally correlated MWR SST errors. The ocean areas retained to make this distribution, and their corresponding local estimates of s , are shown as Figure 4. (As noted in methods, the limiting of the area for calculating the sensitivity to the locations P shown in the figure is to ensure the ratio $\frac{\delta_{hi}-\delta_{lo}}{m_{hi}-m_{lo}}$ is not ill-conditioned because of a small denominator.) There is coherent geographical variation in the estimate of s , but the degree to which this reflects differences in dust impact or coherent MWR SST errors is unknown; for this reason, the median rather than local estimates of s will be used as the scaling to correct the analysis SST.

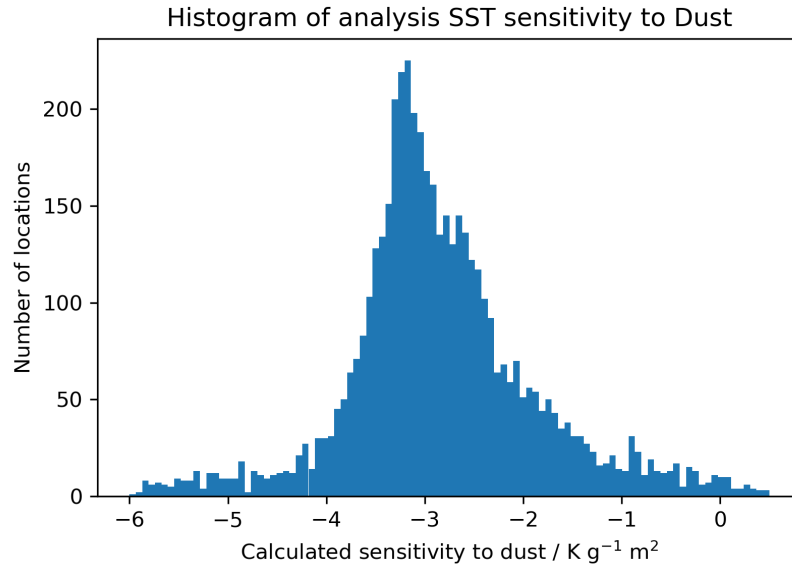


Figure 3: Distribution of estimates of $s = \frac{\partial \hat{x}}{\partial m}$.

Algorithm Theoretical Basis Document D2.1 v1

Estimates of dust sensitivity

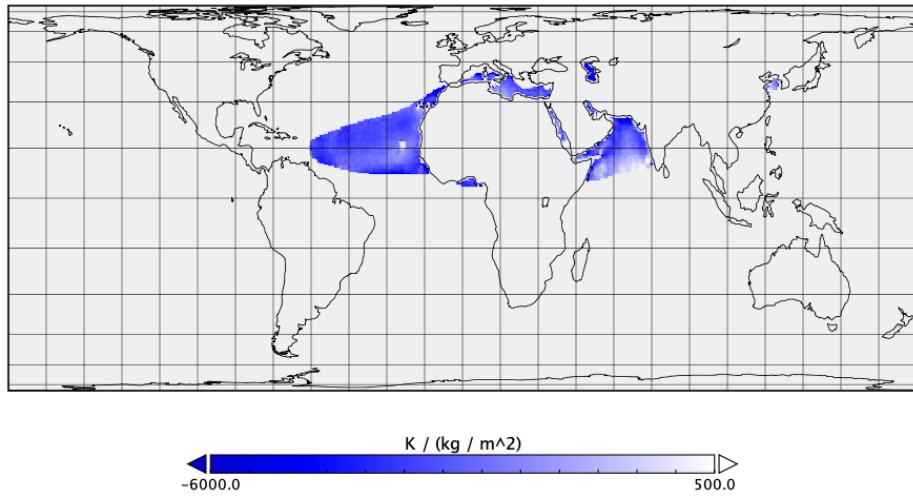


Figure 4: Locations and result of dust sensitivity calculation per 0.5 degree grid cell.

The average climatological correction that will be applied to the analysis SST using this approach is 0.02 K, and the maximum is 1.43 K. The geographical distribution of the largest climatological correction during the annual cycle is show in Figure 5. The size of corrections are as expected.

Annual maximum climatological correction to SST analysis

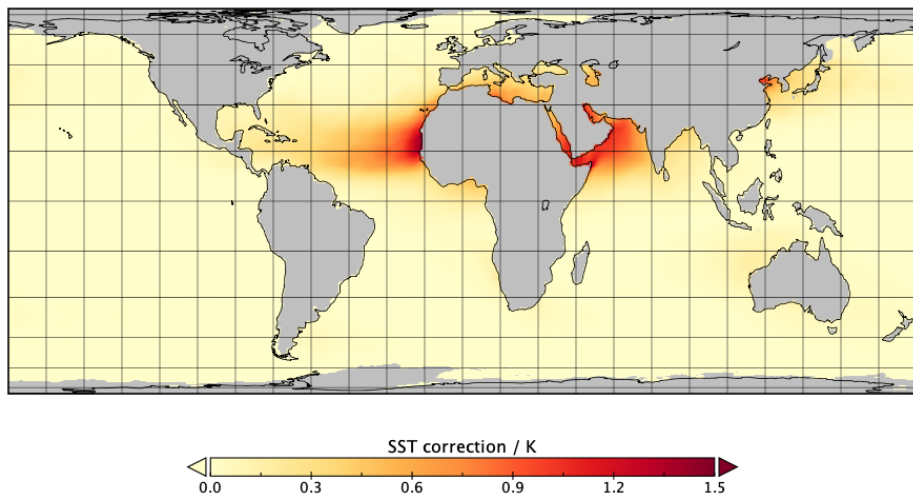


Figure 5: Annual maximum of the climatological correction to analysis SSTs for the influence of dust aerosol

Algorithm Theoretical Basis Document D2.1 v1

The climatological correction was developed on data during the 2000s, when a mix of dual-view reference sensors and single-view meteorological IR sensors were informing the SST analysis. The climatological correction must be applied prior to the CAMS re-analysis period, particularly in the 1980s, when only older, less-well calibrated, single-view sensors are available. To check it is valid then, we consider the comparison of the analysis SST against an in situ only analysis, HadSST3 (Kennedy et al., 2011a and 2011b). This data product is gappy, and estimated on monthly 5 degree scales, which determines the scale at which comparison is possible.

Absolute HadSST3 SST estimates were obtained by adding the HadSST3 ensemble median anomaly to the ensemble median climatology. The CCI SST analysis was averaged to 5 degree cells. The average for 1982 to 1990 inclusive was created for both datasets, per month.

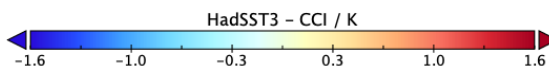
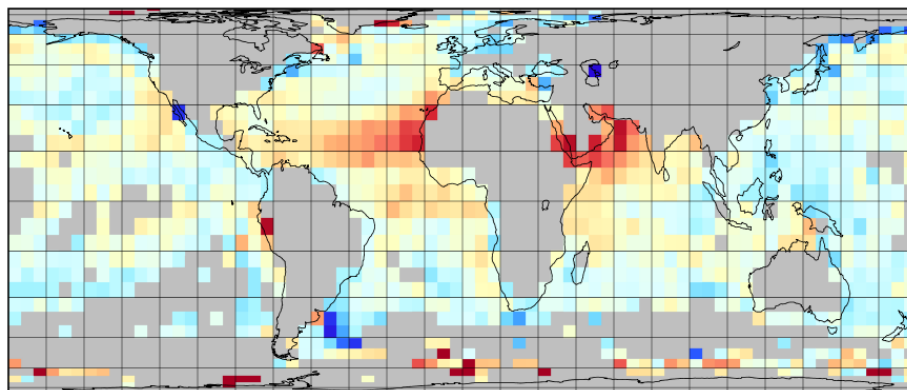
The comparison of the panels in Figure 6, for the Julys between 1982 and 1990) shows clearly that the dust aerosol correction is climatologically beneficial. The same is true for other months of the year. This is shown in the table below, where we show the mean and standard deviation of the monthly 1982-1990-mean differences across cells in the main dust-affected region, namely, 75°W to 75°E, 5°N to 40°N. Although there remains an offset between the SST analysis and HadSST3 (the former being cooler) it is less variable around the year than before correction. Moreover, the standard deviation across cells is considerably improved in the high dust months (March to September). The correction is thus inferred to be beneficial to the CCI SST analysis in reducing mean dust-related biases in the 1980s.

<i>Month</i>	Uncorrected Mean / K	Uncorrected SD / K	Corrected Mean / K	Corrected SD / K
<i>January</i>	0.20	0.22	0.16	0.20
<i>February</i>	0.24	0.23	0.17	0.21
<i>March</i>	0.29	0.23	0.19	0.21
<i>April</i>	0.29	0.23	0.16	0.21
<i>May</i>	0.31	0.28	0.14	0.21
<i>June</i>	0.42	0.43	0.15	0.26
<i>July</i>	0.42	0.45	0.11	0.26
<i>August</i>	0.41	0.39	0.18	0.24
<i>September</i>	0.32	0.36	0.17	0.27
<i>October</i>	0.14	0.25	0.04	0.21
<i>November</i>	-0.03	0.22	-0.08	0.21
<i>December</i>	0.11	0.21	0.07	0.19

Table 3: Uncorrected and corrected SST difference statistics over main dust-affected area of ocean

Algorithm Theoretical Basis Document D2.1 v1

No correction for dust aerosol
 1982 - 1990 average, July



Corrected for dust aerosol
 1982 - 1990 average, July

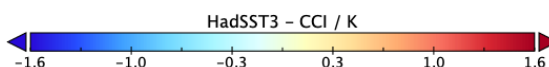
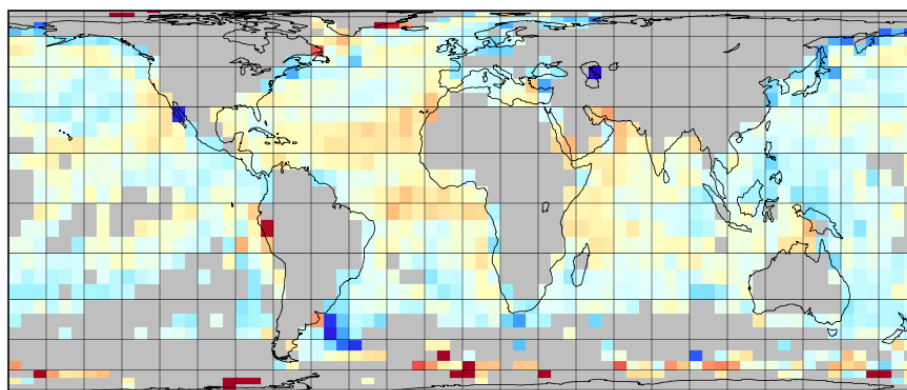


Figure 6: Comparison of HadSST3 and CCI analysis SST in the 1980s, with and without climatological correction for dust aerosol (based on dust influence during 2003 – 2017).

Algorithm Theoretical Basis Document D2.1 v1

5. REFERENCES

Alerskans, E., Høyer, J. L., Gentemann, C. L., Pedersen, L. T., Nielsen-Englyst, P., & Donlon, C. 2020. Construction of a climate data record of sea surface temperature from passive microwave measurements. *Remote Sensing of Environment*, 236, 111485.

Ashcroft, P., Wentz, F.J., 2013. AMSR-E/Aqua L2A Global Swath Spatially-Resampled Brightness Temperatures, Version 3.

Atlas, R., Hoffman, R.N., Ardizzone, J., Leidner, S.M., Jusem, J.C., Smith, D.K., Gombos, D., 2011. A cross-calibrated, multiplatform ocean surface wind velocity product for meteorological and oceanographic applications. *Bull. Am. Meteorol. Soc.* 92, 157–174. <https://doi.org/10.1175/2010BAMS2946.1>

Block, T., Embacher, S., Merchant, C.J., Donlon, C., 2018. High-performance software framework for the calculation of satellite-to-satellite data matchups (MMS version 1.2). *Geosci. Model Dev.* 11, 2419–2427. <https://doi.org/10.5194/gmd-11-2419-2018>

Chang, P.S., Jelenak, Z., Alswiss, S., 2015. Algorithm Theoretical Basis Document: GCOM-W1/AMSR2 Day-1 EDR version 1.0.

Copernicus Climate Change Service (C3S), 2017. ERA5: Fifth generation of ECMWF atmospheric reanalyses of the global climate. Copernicus Climate Change Service Climate Data Store (CDS), December 2019. <https://cds.climate.copernicus.eu/cdsapp#!/home>

Dee, D.P., Uppala, S.M., Simmons, A.J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M.A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A.C.M., van de Berg, L., Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A.J., Haimberger, L., Healy, S.B., Hersbach, H., Hólm, E. V., Isaksen, L., Kållberg, P., Köhler, M., Matricardi, M., McNally, A.P., Monge-Sanz, B.M., Morcrette, J.J., Park, B.K., Peubey, C., de Rosnay, P., Tavolato, C., Thépaut, J.N., Vitart, F., 2011. The ERA-Interim reanalysis: Configuration and performance of the data assimilation system. *Q. J. R. Meteorol. Soc.* 137, 553–597. <https://doi.org/10.1002/qj.828>

Desroziers G, Berre L, Chapnik B, Poli P. Diagnosis of observation, background and analysis-error statistics in observation space. *Quarterly Journal of the Royal Meteorological Society* 2005; 131: 3385-3396.

DMI, 2007. GHRSSST Level 4 DMI_OI Global Foundation Sea Surface Temperature Analysis (GDS version 2). Ver. 1.0.

GHRSSST Science Team, 2010. The recommended GHRSSST data specification (GDS) 2.0, document revision 5.

Algorithm Theoretical Basis Document D2.1 v1

Good, S.A., Martin, M.J., Rayner, N.A., 2013. EN4 : Quality controlled ocean temperature and salinity profiles and monthly objective analyses with uncertainty estimates. *J. Geophys. Res. Ocean.* 118, 6704–6716. <https://doi.org/10.1002/2013JC009067>

Høyer, J.L., Karagali, I., Dybkjær, G., Tonboe, R., 2012. Multi sensor validation and error characteristics of Arctic satellite sea surface temperature observations. *Remote Sens. Environ.* 121, 335–346. <https://doi.org/10.1016/j.rse.2012.01.013>

Inness, A., Ades, M., Agustí-Panareda, A., Barré, J., Benedictow, A., Blechschmidt, A.-M., Dominguez, J. J., Engelen, R., Eskes, H., Flemming, J., Huijnen, V., Jones, L., Kipling, Z., Massart, S., Parrington, M., Peuch, V.-H., Razinger, M., Remy, S., Schulz, M., and Suttie, M.: The CAMS reanalysis of atmospheric composition, *Atmos. Chem. Phys.*, 19, 3515–3556, <https://doi.org/10.5194/acp-19-3515-2019>, 2019.

Kennedy J.J., Rayner, N.A., Smith, R.O., Saunby, M. and Parker, D.E. (2011b). Reassessing biases and other uncertainties in sea-surface temperature observations since 1850 part 1: measurement and sampling errors. *J. Geophys. Res.*, 116, D14103, doi:10.1029/2010JD015218

Kennedy J.J., Rayner, N.A., Smith, R.O., Saunby, M. and Parker, D.E. (2011c). Reassessing biases and other uncertainties in sea-surface temperature observations since 1850 part 2: biases and homogenisation. *J. Geophys. Res.*, 116, D14104, doi:10.1029/2010JD015220

Maeda, T., Taniguchi, Y., Imaoka, K., 2016. GCOM-W1 AMSR2 Level 1R Product: Dataset of Brightness Temperature Modified Using the Antenna Pattern Matching Technique. *IEEE Trans. Geosci. Remote Sens.* 54, 770–782. <https://doi.org/10.1109/TGRS.2015.2465170>

McPhaden, M.J., Ando, K., Boulès, B., Freitag, H.P., Lumpkin, R., Masumoto, Y., Murty, V.S.N., Nobre, P., Ravichandran, M., Vialard, J., Vousdem, D., Yu, W., 2010. The Global Tropical Moored Buoy Array.pdf. *Proc. Ocean.* 9, 668–682.

Merchant, C.J., Old, C.P., Embury, O., MacCallum, S.N., 2008. Generalized Bayesian Cloud Screening.

Merchant CJ, Embury O, Bulgin CE, Block T, Corlett GK, Fiedler E, et al. Satellite-based time-series of sea-surface temperature since 1981 for climate applications. *Scientific Data* 2019; 6.

Merchant, C., Saux-Picart, S. and Waller, J. (2020) Bias correction and covariance parameters for optimal estimation by exploiting matched in-situ references. *Remote Sensing of Environment*, 237. 111590. ISSN 0034-4257 doi: <https://doi.org/10.1016/j.rse.2019.111590>

Rayner, N., Good, S., Block, T., 2015. SST CCI Product User Guide, Project Document, SST_CCI-PUG-UKMO-201. <http://www.esa-sst-cci.org/PUG/documents.htm>.

Algorithm Theoretical Basis Document D2.1 v1

Wentz, F.J., Meissner, T., 2007. Supplement 1 Algorithm Theoretical Basis Document for AMSR-E Ocean Algorithms, RSS Technical Report 051707. Remote Sensing Systems, Santa Rosa, CA.

Wentz, F.J., Meissner, T., 2000. Algorithm Theoretical Basis Document (ATBD): AMSR Ocean Algorithm (Version 2), RSS Tech. Proposal 121599A-1. Remote Sensing Systems, Santa Rosa, CA.

Woodruff, S.D., Worley, S.J., Lubker, S.J., Ji, Z., Eric Freeman, J., Berry, D.I., Brohan, P., Kent, E.C., Reynolds, R.W., Smith, S.R., Wilkinson, C., 2011. ICOADS Release 2.5: extensions and enhancements to the surface marine meteorological archive. *Int. journal Climatol.* 31, 951–967. <https://doi.org/10.1002/joc.2103>