

Supporting Information for

Satellite-based ~~Time Series~~ of Sea Surface Salinity designed for Ocean and Climate Studies

J. Boutin¹, N. Reul², J. Koehler³, A. Martin⁴, R. Catany⁵, S. Guimbard⁶, F. Rouffi⁷, J.L. Vergely⁷, M. Arias⁵, M. Chakroun⁷, G. Corato⁸, V. Estella-Perez^{1,*}, A. Hasson^{1,†}, S. Josey⁴, D. Khvorostyanov¹, N. Kolodziejczyk², J. Mignot¹, L. Olivier¹, G. Reverdin¹, D. Stammer³, A. Supply^{1,2}, C. Thouvenin-Masson¹, A. Turiel⁹, J. Vialard¹, P. Cipollini^{10,‡}, C. Donlon^{10,‡}, R. Sabia¹¹, S. Mecklenburg¹⁰

¹Sorbonne University, LOCEAN/IPSL Laboratory, CNRS–IRD–MNHN, Paris, France.

²University of Brest, LOPS Laboratory, IUEM, UBO–CNRS–IRD–Ifremer, Plouzané, France.

³Institut für Meereskunde, Centrum für Erdsystemwissenschaften und Nachhaltigkeit, Universität Hamburg, Germany.

⁴National Oceanography Centre, Southampton, UK.

⁵ARGANS Ltd, UK.

⁶Ocean Scope, France.

⁷ACRI-st, France.

⁸Adwaiseo, Luxemburg.

⁹Barcelona Expert Center (BEC) and Institute of Marine Sciences (ICM), CSIC, Spain.

¹⁰European Space Agency, ECSAT, Harwell, United Kingdom.

¹¹Telespazio-UK for ESA, ESRIN, Frascati, Italy.

Corresponding author: Jacqueline Boutin (jb@locean.ipsl.fr)

*Now at UL Services Spain SL

†Now at Mercator Ocean International, France

‡Now at European Space Agency, ESTEC, Noordwijk, the Netherlands.

Contents of this file

Text ~~S4~~S1, S2, S3 and S6

Introduction

This supporting information ~~gives details about~~provides information on the user survey, ~~about the differences between SSS CCI version 1 and version 2 (version 2 is described in the paper), about~~on SMOS data preprocessing, on the methodology used to build L4 SSS fields, to estimate SSS variability and level 2 SSS uncertainties used for the L4 SSS generation, ~~PIMEP validation statistics obtained for weekly SSS CCI fields and SSS anomalies computed with ISAS SSS fields.~~

All the ~~information~~material presented here ~~are~~is not essential to the comprehension of the article but ~~bring~~provides more detailed information ~~and details~~ to the reader.

S1. SSS data requirements for ocean and climate studies

To create an SSS data set that satisfies the needs of climate users, both modeling and Earth-observing scientists groups, users of satellite SSS data were consulted through various approaches: personally, via e-mail, mailing lists, or at meetings. They were invited to participate in web surveys and to specify their requirements for satellite SSS data. In our survey, we asked specific questions to find out the user's priorities (typically higher resolution or improved uncertainty estimates).

The survey (available on <https://forms.gle/BVDroYrNpVvpxFJu9>) gathered 54 answers from various countries of origin/fields (Table S1). Most responses were from the USA (28%), followed by Germany (19%) and the UK (17%).

Table S1: Percentage of respondents to the online survey.

<u>Country</u>	<u>USA</u>	<u>Germany</u>	<u>UK</u>	<u>South Korea</u>	<u>Italy</u>	<u>Spain</u>	<u>Argentina</u>	<u>Australia</u>	<u>Brazil</u>	<u>Norway</u>	<u>Japan</u>	<u>France</u>
----------------	------------	----------------	-----------	--------------------	--------------	--------------	------------------	------------------	---------------	---------------	--------------	---------------

%	28	19	17	2	4	7	2	2	2	4	2	13
---	----	----	----	---	---	---	---	---	---	---	---	----

The questionnaire had four different parts, which were 1) User profile information, 2) User dataset requirements, 3) User dataset quality information, 4) Other user requirements or suggestions.

The majority of users requires global spatial coverage and temporal coverage from at least 9 years. The resolution requirements vary according to the studied phenomena. About 33% of respondents require data with a temporal resolution of 1-3 days, while, for 35% (28%) of respondents, weekly (monthly) averaged data are sufficient (Figure S1a). In terms of spatial resolution, 39% of respondents require data on a 0.25° spatial grid, while 28% of respondents require data on 1° spatial grid (Figure S1b). The majority of respondents would prefer a data product with high spatial and temporal resolution (weekly, 0.25°) on a regular latitude-longitude grid. Interestingly, a majority of users would prefer a product with high temporal and spatial resolution and a lower accuracy rather than working with a product with high accuracy but a lower resolution (Figure S1c). It was also found that the participants are aware of the data set limitations and have realistic expectations.

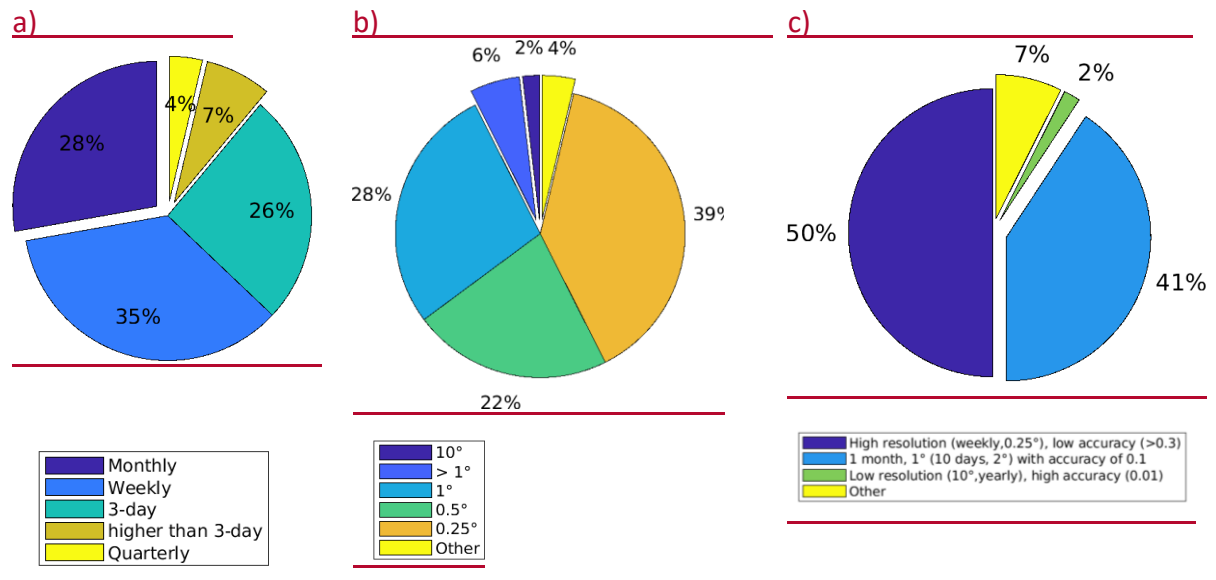


Figure S1: Percentage of required (a) temporal and (b) spatial resolution. c) Preferred SSS product based on the user's spatial and temporal resolution needs.

According to the survey, data should be combined to overcome the weaknesses of individual datasets. 50% prefer a combination of satellite and in-situ measurements, whereas 39%

require the combination of data from different satellite sensors. However, in the CCI L4 SSS products described here, information from in-situ measurements was restricted to a minimum in order to work with measurements having homogeneous spatio-temporal resolution and sampling. By making available the multiple-sensor datasets on different spatial-temporal grids, the needs of different users can be met. The most common requirement is for L4 data (43%), directly followed by requirements for L3 (37%). Some potential users, mainly modelers or scientists investigating rapid SSS changes, require L2 (20%). L3 and L2 data are already available from the original data centers. L2 and L3 datasets including the CCI+SSS systematic corrections are kept as an internal CCI+SSS product.

Uncertainty information for each SSS grid point has to be fully characterized, including random noise and systematic uncertainties of the applied adjustments. Information about bias (systematic uncertainty) correction is most commonly required by respondents. 46% of the respondents would prefer a quality information easy to use, such as a good/bad flag or the probability that a value is good/bad.

User Requirement Survey results show the importance of contacting users and promoting communication between the users and potential users of CCI L4 SSS fields. Users will be regularly contacted to refine requirements, as well as to check their satisfaction with the CCI L4 SSS product. The recommendations regarding resolution, format, quality, and additional information derived from the user consultation are summarized in the User Requirement Document (URD available on <https://climate.esa.int/en/projects/sea-surface-salinity/key-documents/>).

S2. SMOS SSS data pre-processing

We only consider SMOS SSS satisfying the following criteria (same notations as in (Boutin et al., 2018)): normalized χ of the retrieval, $\chi_N < 3$, SSS random uncertainty, $E_{SSS\ L2} < 3$, pixel within +/-400 km from the center of the swath, with small number of Tb outliers (level 2 fg_outlier flag), uncontaminated by ice (level 2 Dg_suspect_ice=0; this flag removes pixels in cold waters ($SST < 2^\circ C$) in which at least one Tb differs by more than 20K from modeled Tb. This is a very stringent filtering that is likely to be removed in future versions), with moderate to strong RFI and ice contamination as detected using SMOS retrieved pseudo-dielectric constant, A_{card} ($|A_{card\ smos} - A_{card\ mod}| < 2$ and $A_{card} > 42$; see more in (Supply et al., 2020b)), wind speed less than 16 m/s, SSS between 2 and 45.

Deficiencies in the dielectric constant model leads us to adjust SMOS SSS with a polynomial SST function derived from comparisons between Aquarius SSS (retrieved with similar dielectric constant model and atmospheric model as SMOS processing models), and Argo SSS (blue dotted curve in Figure 16 of (Dinnat et al., 2019)). A correction for seasonal latitudinal varying biases is also applied, similar to what is described in (Boutin et al., 2018):

$$\text{SSS}_{\text{obs}}(t, \phi, \lambda, x_{\text{swath}}, x_{\text{orb}}) = \text{SSS}_{\text{smos}_{\text{ref}}}(\phi, \lambda, m) - b_{\text{lat}}(\phi, x_{\text{swath}}, x_{\text{orb}}, m)$$

where SSS_{obs} is the observed SSS, t is the time of the measurement, ϕ , and λ , are respectively the latitude and the longitude of the considered pixel over the ocean, x_{swath} corresponds to the pixel location across the swath, x_{orb} indicates the satellite orbit direction (ascending or descending), b_{lat} is a latitudinal correction that varies seasonally as a function of the month, m , and $\text{SSS}_{\text{smos}_{\text{ref}}}$ is a reference SMOS SSS taken at a given x_{swath} and x_{orb} , chosen so that $\text{SSS}_{\text{smos}_{\text{ref}}}$ interannual variability for the considered month and the corresponding pixel is the closest with that of in situ interpolated (In-situ Analysis System, ISAS) SSS after 5° latitudinal smoothing. b_{lat} is estimated through a least square minimization approach, and through a series of iterations. In order to avoid land-sea contamination, b_{lat} is derived from SMOS measurements further than 1200 km from coast except at the high northern latitudes where the distance to coast is reduced to 600km in order to get enough measurements. It is computed over 2012-2018 to avoid large RFIs in the North Atlantic in 2010 and 2011. b_{lat} is then removed from all SMOS SSS whatever their distance to coast and before estimating the land-sea contamination correction.

S1.S3. Generation of level 4 SSS **fields: detailed algorithm**

The algorithm is looking searching for solutions $\text{SSS}(t)$ and b_c that both minimize the cost function. Each grid node is processed separately. All available SSS data associated with the grid node considered are used by the algorithm. The problem is linear-, so that to minimize the cost function, a classic Raphson-Newton descent is used.

SSS_{obs} is the observation vector that contains SMOS, SMAP and Aquarius data:

$$\text{SSS}_{\text{obs}} = \begin{pmatrix} \text{SSS}_{\text{smos}} \\ \text{SSS}_{\text{aqua}} \\ \text{SSS}_{\text{smap}} \end{pmatrix}$$

The parameter vector is written:

$$m = \begin{pmatrix} \text{SSS} \\ \text{bc_smos} \\ \text{bc_aqua} \\ \text{bc_smap} \end{pmatrix}$$

bc_smos, bc_aqua, bc_smap are the vectors that contain the biases for each type of acquisition (ascending/descending, dwell lines, etc) that can be grouped into a vector bc. We take as a priori bc=0 for all sensors and acquisition types.

The vector parameter a priori is written:

$$m_{\text{prior}} = \begin{pmatrix} \text{SSSprior} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

SSSprior is the initial SSS value, this value is constant over time. It is taken equal to the ~~SSS of the SMOS central dwell line ascending orbits when possible. Otherwise, the~~ median of the observed SSS ~~is used~~.

We call H, the matrix of partial derivatives:

$$H = \begin{bmatrix} \frac{\partial \text{SSS}_{\text{smos}}}{\partial \text{SSS}} & \frac{\partial \text{SSS}_{\text{smos}}}{\partial \text{bc_smos}} & \frac{\partial \text{SSS}_{\text{smos}}}{\partial \text{bc_aqua}} & \frac{\partial \text{SSS}_{\text{smos}}}{\partial \text{bc_smap}} \\ \frac{\partial \text{SSS}_{\text{aqua}}}{\partial \text{SSS}} & \frac{\partial \text{SSS}_{\text{aqua}}}{\partial \text{bc_smos}} & \frac{\partial \text{SSS}_{\text{aqua}}}{\partial \text{bc_aqua}} & \frac{\partial \text{SSS}_{\text{aqua}}}{\partial \text{bc_smap}} \\ \frac{\partial \text{SSS}_{\text{smap}}}{\partial \text{SSS}} & \frac{\partial \text{SSS}_{\text{smap}}}{\partial \text{bc_smos}} & \frac{\partial \text{SSS}_{\text{smap}}}{\partial \text{bc_aqua}} & \frac{\partial \text{SSS}_{\text{smap}}}{\partial \text{bc_smap}} \end{bmatrix}$$

where: $\text{SSS}_{\text{sensor}} = F(m) = \text{SSS} - \text{bc}_{\text{sensor}}$

with "sensor" = smos (SMOS), aqua (Aquarius) or smap (SMAP).

This matrix is calculated on the observation points.

The covariance matrices used are defined as follows:

- C_d the error matrix with data uncertainties derived,
- C_m the matrix of SSS variability and a priori uncertainty on bc,
- C_r the matrix of representativity uncertainties.

$$C_d = \begin{bmatrix} C_{d_smos} & 0 & 0 \\ 0 & C_{d_aqua} & 0 \\ 0 & 0 & C_{d_smap} \end{bmatrix}$$

$$C_o = \begin{bmatrix} C_{o_smos} & 0 & 0 \\ 0 & C_{o_aqua} & 0 \\ 0 & 0 & C_{o_smap} \end{bmatrix}$$

$$C_m = \begin{bmatrix} \text{CSSS} & 0 \\ 0 & \text{Cbc} \end{bmatrix}$$

CSSS is a time smoothing operator that contains the expected variability that is provided as auxiliary data. Thus, the covariance of the SSS that links two times t1 and t2 (either between two observational times or between an observational time and a sampled time of the OI) is written:

$$\text{CSSS}(t_1, t_2) = \text{sigSSS}(t_1) \text{sigSSS}(t_2) \exp\left(-\frac{(t_1 - t_2)^2}{\xi^2}\right)$$

with $\xi=25$ days and 6 days for monthly and weekly products respectively. "sigSSS" is interpolated temporally to t1 and t2 from seasonal variability.

~~"sigSSS" is interpolated temporally to the acquisition times from seasonal variability.~~

"Cbc" is a diagonal matrix that contains the a priori standard deviation of biases. This standard deviation is set at 4 pss.

The Cr matrix corresponds to representativity uncertainties:

$$C_r = \begin{bmatrix} C_{r_smos} & 0 & 0 \\ 0 & C_{r_aqua} & 0 \\ 0 & 0 & C_{r_smap} \end{bmatrix}$$

With, in CCI v1 and v2, "Cr_smos" and "Cr_smap" set to 0.

In addition to measurement uncertainties, representativity uncertainties are added:

$$C_t = C_d C_o + C_r$$

Representativity uncertainties are reported monthly. They are interpolated temporally to the acquisition times.

In this formalism the cost function is written for each grid node:

$$C(\text{SSS}, bc) = \langle \text{SSS}_{\text{obs}} - F(m) | C_t^{-1} \cdot (\text{SSS}_{\text{obs}} - F(m)) \rangle + \langle m - m_{\text{prior}} | C_m^{-1} \cdot (m - m_{\text{prior}}) \rangle$$

$$C(\text{SSS}, bc) = (\text{SSS}_{\text{obs}} - F(m))^T C_t^{-1} \cdot (\text{SSS}_{\text{obs}} - F(m)) + (m - m_{\text{prior}})^T C_m^{-1} \cdot (m - m_{\text{prior}})$$

with:

$$F(m) = \text{SSS} - bc$$

We look for SSS_est and bc_est that minimizes C (SSS, bc). The solution of minimization is written:

$$m_{\text{est}} = m_{\text{prior}} + C_m \cdot H^T \cdot (H \cdot C_m \cdot H^T + C_t)^{-1} C_m^{-1} \cdot H^T \cdot (H \cdot C_m \cdot H^T + H^T + C_t)^{-1} \cdot (\text{SSS}_{\text{obs}} - F(m_{\text{prior}}))$$

where "T" indicates the transpose operator, Cm is the matrix of variability operating in the observational space and Cm' the matrix of variability between observational time and regular sampled time of the OI.

Estimation of monthly SSS

In order to estimate the monthly SSS, we proceed in 3 steps:

1) a first estimation of the biases and time series of SSS, is performed spatial grid node by spatial grid node ~~is performed~~,

2) a 3-sigma filtering of the observed SSS in comparison with the estimated SSS is done.

The aim here is to identify any outliers against the returned SSS field. Outliers can be linked to intermittent RFIs. It is ~~considered~~assumed here that stable RFI contamination can be corrected.

3) a second estimate of SSS biases and time series after removing outliers.

The relative biases used to derive monthly SSS are estimated taking the averaged SSS from the SMOS central across swath location as a priori.

Estimation of weekly SSS

To estimate the weekly SSS, the biases calculated at the monthly SSS generation step are frozen (it is assumed that the biases will not be better estimated from a weekly smoothing). We start from the monthly SSS as a priori value. We estimate the weekly fluctuations around this a priori, taking into account the acceptable SSS variability between weekly and monthly fields that was derived as a monthly climatology from the Mercator model. A 3-sigma filter is used in order to eliminate outliers that deviate too far from what is expected. Here, $\sigma =$

$\sqrt{\text{error_L2}^2 + \text{variability}^2}$. The weekly SSS field estimate is done in a single step.

Absolute SSS correction

At the end of the inter-sensor bias correction step, the salinities ~~obtained are set on average of the SSS of all sensors. However, the SMOS ascending central across swath location that was taken as reference can itself be affected by a bias~~ are obtained in relative values, i.e. they are known within one additive constant. This is corrected by adjusting a quantile of the CCI and ISAS SSS statistical distributions in each grid node over the period considered. The dynamics of ~~the~~ SSS

variability are not affected by this adjustment as only one constant value, grid node per grid node, is added for the entire period. In regions where SSS variability is low, we assume that high frequency variability better sampled by CCI than by ISAS does not affect significantly the median of the SSS and we therefore adjust ~~both~~the SSS median (50% quantile). In regions with larger variability, given that intermittent freshening is much more frequent than intermittent over-salting, we expect the high part of the SSS distribution to be less affected by the higher frequency sampling by satellite than by ISAS. Hence in case of high weekly variability, we perform the calibration of CCI SSS on ISAS SSS, not by using the median but a highhigher quantile, in order to promote the calibration on the high SSS values. A high quantile is not used everywhere as in case the SSS error is greater than the variability, the high quantile of the satellite SSS is expected to differ (be higher) from the one of ISAS.

If the variability is greater than 0.8, the quantile is taken as 80%. If the variability is between 0.6 and 0.8, we take a quantile intermediate between 50% and 80% that varies linearly with the SSS variability. The map of quantiles used for the absolute calibration of the SSS is given in Figure S4.

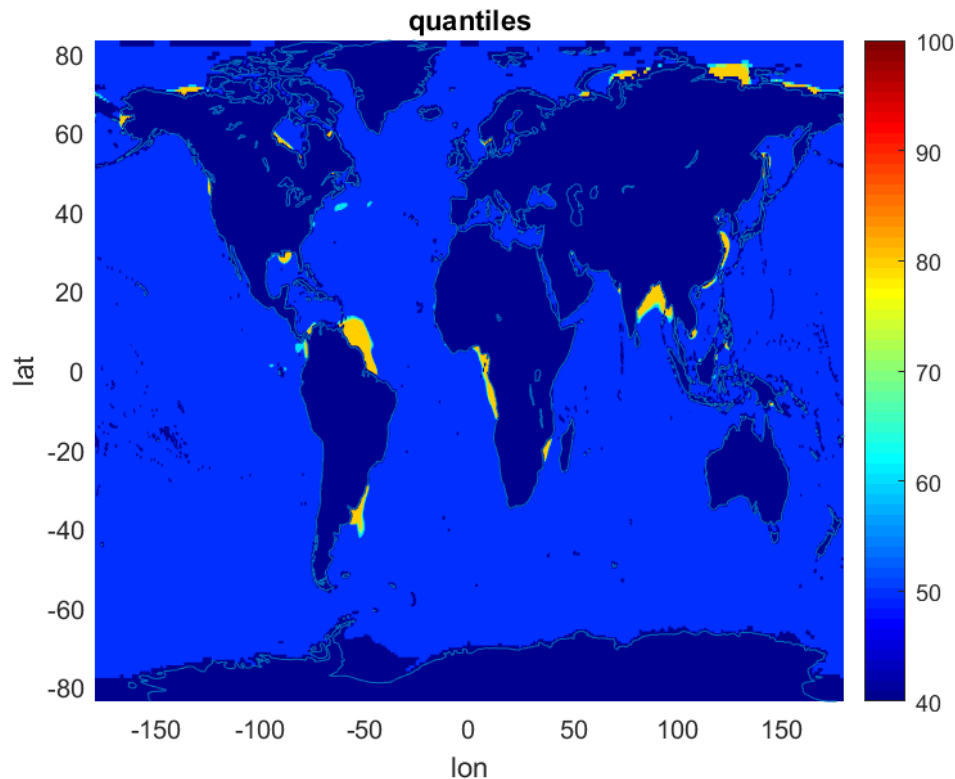


Figure S4S3: Quantile map used for the SSS absolute calibration. x and y axis units in pixel number for longitude and latitude respectively.

Estimation of SSS uncertainties

The computation of theoretical uncertainties is obtained directly from the pseudo Hessian matrix.

$$C_{post} = C_m - C_m \cdot H^T \cdot (H \cdot C_m \cdot H^T + C_t)^{-1} \cdot C_m'' - C_m' \cdot H^T \cdot (H \cdot C_m \cdot H^T + C_t)^{-1} \cdot H \cdot C_m \cdot C_m'^T$$

Where single apostrophe (') indicates covariance defined between observational time and regular sampled time of the OI. Double apostrophe (") indicates covariance acting over regular sampled time of the OI.

The problem ~~turns~~ becomes into inverting the " $H \cdot C_m \cdot H^T + C_t$ " matrix over the entire period, which is rather computationally heavy. We therefore prefer to make take a sliding window over a large time interval and invert the matrix on this time domain (the computation being similar to the one we could perform over the entire period).

Note that the *a posteriori* uncertainty is necessarily lower than the variability introduced in the operator C_m . In the monthly case, this variability corresponds to the expected monthly fluctuations with respect to the whole time series. In the weekly case, the variability is calculated relative to the monthly field. The latter is generally lower than the monthly variability. The *a posteriori* uncertainty obtained on the weekly fields should therefore be lower than that obtained on the monthly fields. However, this is only true if the weekly fields are derived from noise-corrected monthly fields, which is not the case. The propagation of uncertainties on the weekly fields must therefore take into account uncertainties on the monthly field. Thus, for the monthly fields, we have:

~~$$C_{post_month} = C_{m_month} - C_{m_month} \cdot H^T \cdot (H \cdot C_{m_month} \cdot H^T + C_t)^{-1} \cdot H \cdot C_{m_month}$$~~

$$C_{post_month} = C_{m_month}'' - C_{m_month}' \cdot H^T \cdot (H \cdot C_{m_month} \cdot H^T + C_t)^{-1} \cdot H \cdot C_{m_month}'^T$$

and for the weekly fields:

~~$$C_{post_week} = C_{post_month} + C_{m_week} - C_{m_week} \cdot H^T \cdot (H \cdot C_{m_week} \cdot H^T + C_t)^{-1} \cdot H \cdot C_{m_week}$$~~

$$C_{post_week} = C_{post_month} + C_{m_week}'' - C_{m_week}' \cdot H^T \cdot (H \cdot C_{m_week} \cdot H^T + C_t)^{-1} \cdot H \cdot C_{m_week}'^T$$

with " $C_{m_month} C_{m_month}$ ", the monthly variability and " $C_{m_week} C_{m_week}$ " the weekly variability relative to the monthly variability.

The *a posteriori* uncertainties on the monthly and weekly fields are therefore obtained as follows:

~~$$\sigma_{SSS_month} = \sqrt{\text{diag}(C_{post_month})}$$~~

~~$$\sigma_{SSS_week} = \sqrt{\text{diag}(C_{post_week})}$$~~

~~$$\sigma_{SSS_month} = \sqrt{\text{diag}(C_{post_month})}$$~~

$$\sigma_{SSS_{week}} = \sqrt{\text{diag}(C_{post_{week}})}$$

The number of outliers is also calculated on this same basis as well as the number of data available. The window sizes used are respectively +/- 30 days and +/- 10 days for monthly and weekly products respectively.

S6. SSS random uncertainties of L2 satellite SSS

The SMAP random uncertainties are derived from the std of the difference between SSS retrieved from fore and aft acquisitions (Figure S6.1, red). They are very close to a modeled error with a 0.45K radiometric noise.

The Aquarius random uncertainties are derived from comparisons of SSS sampled at successive (7-day interval days apart) Aquarius SSS measurements and are fitted with an SST dependency. (Figure S6.1, blue).

The SMOS random uncertainties are taken from the theoretical error multiplied by the Chi of the retrieval provided in SMOS L2 files which are found to give provide a reasonable estimate (Figure S6.2).

In the above estimates, only pixels further than 800 km from coast have been considered in order to avoid land-sea contamination and very large representativity uncertainties.

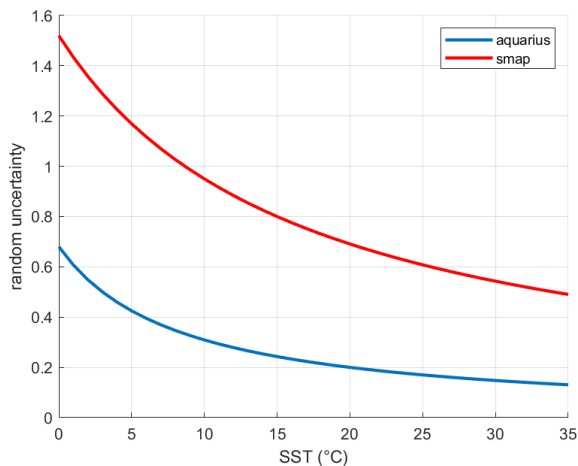


Figure S6.1. SSS uncertainties derived for SMAP (red) and Aquarius (blue) as a function of SST.

We check the reasonable behavior of estimated random errors, σ_{SSS} , by ~~looking at~~ considering the statistical distribution of the centered reduced SSS, $SSSc$:

$$SSSc = (SSS_{obs} - SSS_{ref}) / \sigma_{SSS} \quad (1)$$

where SSS_{obs} is the retrieved SSS possibly corrected from systematic uncertainties, SSS_{ref} is a reference SSS. Figure S6.2 shows an example obtained with SMOS data further than 800 km from coast compared with a 20-day SSS average. In that case $SSSc$ is quite close to the expected Gaussian law. The slightly large value of $std(SSSc)$ (1.25 instead of 1) is partly due to the presence of outliers.

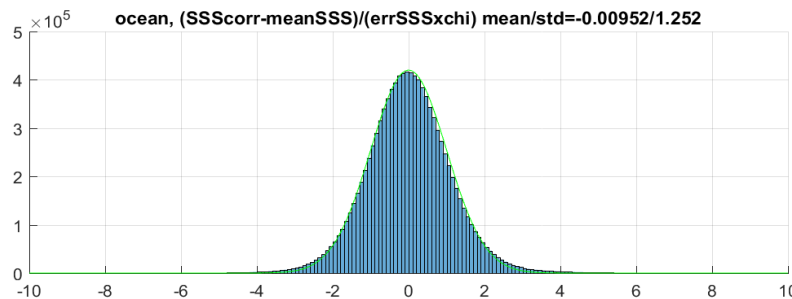


Figure S6.2. Example of the distribution of the centered reduced SSS for grid points in open ocean (further than 800 km from coast) in March 2012. $SSScorr$ represents SMOS L2 SSS corrected from systematic uncertainties. $meanSSS$ is an estimation of the true SSS obtained by averaging $SSScorr$ over a 20day period.

Closer to the coast, $std(SSSc)$ deviates more significantly from 1. Part of this difference can be associated with the variability of ~~the~~ salinity. In order to verify this, we sought to quantify $std(SSSc)$ in regions with low variability. For grid nodes with variability lower than 0.2 on SMOS-CCI SSS rmsd, we compute a robust std, and we observe it to increase towards the coast. We apply the same process to SMAP and Aquarius data. We then derive multiplicative factors which are function to the distance to the coast, $f(d_{coast})$, that will be applied to σ_{SSS} of each instrument (Figure S6.3), so that the reduced random variables normalized with the L2 random uncertainties multiplied by these factors, have a standard deviation equal to 1.

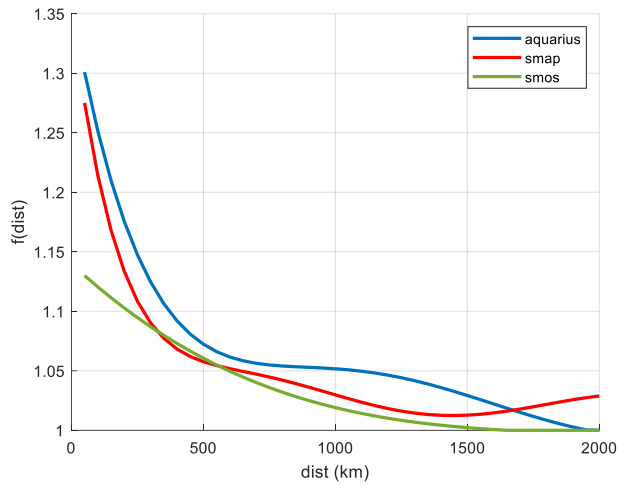


Figure S6.3: Multiplicative factor applied to the σ_{SSS} as a function to the distance to the coast.