

# Validation practices for satellite soil moisture products: What are (the) errors?

A. Gruber, G. De Lannoy, C. Albergel, A. Al-Yaari, L. Brocca, J.-C. Calvet, A. Colliander, M. Cosh, W. Crow, W. Dorigo, C. Draper, M. Hirschi, Y. Kerr, A. Konings, W. Lahoz, K. McColl, J. Muñoz-Sabater, J. Peng, R. Reichle, P. Richaume, C. Ruediger, T. Scanlon, R. van der Schalie, W. Wagner

## Abstract

This paper presents a community effort to develop good practice guidelines for the validation of global coarse-scale satellite soil moisture products. We provide theoretical background, a review of state-of-the-art methodologies for estimating errors in soil moisture data sets, practical recommendations on data pre-processing and presentation of statistical results, and a recommended validation protocol that is supplemented with an example validation exercise focused on microwave-based surface soil moisture products. We conclude by identifying research gaps that should be addressed in the near future.

## 1 Introduction

The validation of soil moisture data sets aims to provide quantitative information about their quality by estimating systematic and random errors (*JCGM*, 2008). For satellite-derived products, this task is far from trivial because high-quality reference data are rarely available at the coarse spatial resolution of space borne microwave instruments that are predominantly used for soil moisture retrievals ( $\sim 10^1 - 10^3$  km<sup>2</sup>), and the retrieval quality is affected by numerous spatially and temporally variable factors (i.e. climatic, topographic and land cover conditions as well as instrument characteristics and the retrieval algorithm structure) (*Ochsner et al.*, 2013; *Crow et al.*, 2012; *Molero et al.*, 2018).

A host of methods exist to reconcile the distinct spatio-temporal characteristics of satellite and reference data sets (sampling and overpass times, penetration depths, representativeness errors, etc.; *Wang et al.*, 2012; *Albergel et al.*, 2008; *Gruber et al.*, 2013a; *Nicolai-Shaw et al.*, 2015; *Colliander et al.*, 2017), which is required before calculating various performance metrics (correlation coefficients, root-mean-square-differences, triple collocation analysis, etc.; *Entekhabi et al.*, 2010a; *Albergel et al.*, 2013; *Gruber et al.*, 2016a; *Loew et al.*, 2017). Given the complexity

30 of the validation problem, however, ambiguous results for the quality and ranking of satellite  
31 soil moisture products can be found in the literature (e.g., *Wagner et al.*, 2014) depending  
32 on which pre-processing and evaluation strategies were followed and which reference data were  
33 used. This paper is a community effort that addresses this issue and aims towards standardizing  
34 good practices for the validation of satellite-based near-surface soil moisture retrievals.

35 Section 2 provides a review of on-going activities regarding the standardization of satellite  
36 soil moisture validation activities. Section 3 describes the most common reference data sources  
37 used for soil moisture validation. Section 4 discusses relevant theoretical aspects and the most  
38 common methods (including data pre-processing) for assessing soil moisture data quality. Section  
39 5 presents a community-agreed validation guidance protocol with an example implementation  
40 of that protocol provided in Appendix A. Finally, Section 6 discusses research gaps that should  
41 be addressed in the near future.

## 42 **2 Towards standardized validation practices**

43 Many efforts have been made to assess and standardize validation practices across Earth obser-  
44 vation (EO) communities (*Zeng et al.*, 2015; *Loew et al.*, 2017; *Su et al.*, 2018). In this section  
45 we review activities most relevant for satellite soil moisture products.

### 46 **2.1 CEOS LPV**

47 The main authority that guides validation activities for satellite-retrieved data of biogeophys-  
48 ical variables is the Committee on Earth Observation Satellites (CEOS) Working Group on  
49 Calibration and Validation (<http://ceos.org/ourwork/workinggroups/wgcv/>; last access: 1  
50 July 2019). Activities related to soil moisture are coordinated by its Land Product Validation  
51 (LPV) subgroup (<https://lpvs.gsfc.nasa.gov/>; last access: 1 July 2019). The CEOS LPV  
52 defines four validation stages (see Table 1) that represent the level of sophistication of validation  
53 protocols employed for a particular data product. Relevant for the work presented here is that  
54 reaching validation stage 3 requires the implementation of a sophisticated validation framework,  
55 as illustrated in Figure 1. In such a framework, standardized community-agreed methods that  
56 are ideally described in a “Validation Good Practice Document” should be employed using fidu-  
57 cial reference data (see Sec. 3) to generate standardized validation reports. With this paper we  
58 aim at providing such a document. The last validation stage 4 is reached once these validation

59 reports are updated on a regular (at least annual) basis.

## 60 **2.2 Quality Assurance Frameworks**

61 The CEOS endorses the Quality Assurance Framework for Earth Observation (QA4EO; <http://qa4eo.org/>; last access: 1 July 2019) as a framework to facilitate the provision of traceable  
62 quality indicators which “shall provide sufficient information to allow all users to readily evaluate  
63 the ‘fitness for purpose’ of the data or derived product” (QA4EO, 2010). The QA4EO provides  
64 top-level guidance documents and templates that encourage the use of metrological principles  
65 (see Sec. 2.3).

67 In 2014, the Quality Assurance for Essential Climate Variables (QA4ECV; <http://www.qa4ecv.eu/>; last access: 1 July 2019) project was initiated to developed a set of guidelines for  
68 the provision of traceable quality information taking in to account the key principles of QA4EO  
69 (*Scanlon et al.*, 2017). To demonstrate how reliable and traceable quality information can  
70 be provided, quality assurance frameworks were developed for selected ECVs (not including  
71 soil moisture; e.g., *Peng et al.*, 2017). The guidelines developed by QA4EO and QA4ECV  
72 are currently embraced by the Copernicus Climate Change Service (C3S; <https://climate.copernicus.eu/>; last access: 1 July 2019) in order to build quality assured, fully traceable  
73 Climate Data Records.

76 In 2018, the Quality Assurance for Soil Moisture project (QA4SM; <https://qa4sm.eodc.eu/>; last access: 1 July 2019) was launched, specifically to create an online validation tool that  
77 employs a community-agreed validation protocol (which is described in this paper) for automat-  
78 ically and regularly generating soil moisture product validation reports, thereby addressing the  
79 CEOS validation framework requirements (see Figure 1).

## 81 **2.3 Metrology and traceability**

82 The CEOS and the QA4EO encourage the use of metrological principles for validation purposes,  
83 which are described in the “Guide to the expression of uncertainty in measurement” (GUM;  
84 *JCGM*, 2008). The GUM is a reference document of the metrological community that provides  
85 strict guidelines on how quality estimates of measurements should be obtained and reported.  
86 In essence, it states that, since they never perfectly represent the true state of the physical  
87 quantity being measured, all measurements should be complemented by uncertainty estimates  
88 that summarize their probability density function (pdf). Furthermore, it states that these

89 uncertainties should be obtained by propagating the uncertainties from all components that  
90 contribute to the measurement process in a way that is traceable back to the “International  
91 System of Units” (SI) standards, either through the standard method for the propagation of  
92 uncertainty (*Parinussa et al.*, 2011; *Merchant et al.*, 2017) or, if not possible analytically, through  
93 Monte Carlo simulations (*JCGM*, 2008).

94 However, while being relatively straightforward in a laboratory or numerical environment,  
95 the traceable propagation of uncertainties in space borne remote sensing measurements and re-  
96 trievals thereof, in particular of soil moisture, faces two particular challenges. First, footprints of  
97 current microwave instruments used for retrieving soil moisture span over tens to thousands of  
98 square kilometers, thereby covering a large variety of climatic, topographic, and land cover condi-  
99 tions. Although certain large-scale homogeneous regions are used for calibrating instruments and  
100 determining Level 1 (L1) backscatter or brightness temperature uncertainties (e.g., rainforests  
101 or polar snow fields; *Figa-Saldaña et al.*, 2002; *Macelloni et al.*, 2006), it is virtually impossible  
102 to obtain global perfectly traceable uncertainty estimates representing all possible measurement  
103 conditions. Second, uncertainty propagation assumes that the models used to propagate uncer-  
104 tainties are themselves perfect (*Parinussa et al.*, 2011). For satellite soil moisture retrievals, this  
105 is particularly problematic because uncertainties resulting from simplifications and assumptions  
106 in both the L1 processing (i.e. geometric correction and radiometric calibration) and the Level  
107 2 (L2) soil moisture retrieval algorithms cannot be accounted for. The soil moisture and other  
108 EO communities have established certain strategies to recover this broken traceability chain  
109 by validating the soil moisture estimates post retrieval against a range of reference data from  
110 various sources. Section 3 will discuss the requirements and current availability of such reference  
111 measurements suited for validation activities. Before entering those discussions, it is necessary  
112 to provide some relevant terminology.

## 113 **2.4 Terminology**

114 The CEOS and the QA4EO encourage the use of the terminology used within the metrological  
115 community as described in the “International Vocabulary of Metrology” (VIM; *JCGM*, 2012).  
116 However, there is a certain level of ambiguity in the existing EO literature, and even within  
117 the VIM and the GUM, regarding the usage of important terms such as errors, uncertainties,  
118 validation, and others. For a comprehensive summary of the most common definitions (from the  
119 VIM, the CEOS, and other sources) we refer the reader to *Loew et al.* (2017). For the purpose



120 of this paper we stress that:

- 121 • the term *error* refers to the deviation of a single measurement from the true value of the  
122 quantity being measured (which is hence always unknown), whereas the term *uncertainty*  
123 refers to the probability distribution underlying an error. For validation purposes, this  
124 probability distribution is the actual quantity of interest;
- 125 • according to the GUM, the uncertainty of a measurement generally contains both sys-  
126 tematic and random components. The laboratory environment of metrological practices  
127 typically allows for thorough measurement calibration, where it is assumed that systematic  
128 errors can be properly determined and corrected. Satellite soil moisture retrievals, how-  
129 ever, usually contain considerable systematic errors which, especially for model calibration  
130 and refinement, provide better insight when estimated separately. Therefore, we use the  
131 term *bias* to refer to systematic errors only and the term *uncertainty* to refer to random  
132 errors only, specifically to their standard deviation (or variance);
- 133 • in the EO validation literature, bias is commonly defined as the temporal mean difference  
134 between two data sets. We follow the broader statistical definition of bias as auto-correlated  
135 error, or as a property of an estimator to systematically over- or underestimate some  
136 quantity (*Dee, 2005*). For better separability of its components, we use the terms *first-*  
137 *order bias* and *second-order bias* to refer more specifically to additive and multiplicative  
138 systematic errors, respectively (see Sec. 4.4.1);
- 139 • the terms *trueness*, *precision*, and *accuracy* are popular antonyms for systematic errors,  
140 random errors, and the combined systematic plus random errors, respectively (*JCGM,*  
141 *2012*). However, trueness and precision are very rarely used in the soil moisture validation  
142 literature and the term accuracy is often ambiguously used to refer to either systematic  
143 or random errors alone;
- 144 • in Earth sciences, the term *validation* is often distinguished from the term *evaluation* such  
145 that validation is used to refer to bias or uncertainty assessment using highly accurate or  
146 at least well traceable in situ reference data (often misleadingly referred to as “ground  
147 truth”; see Sec. 4.2), whereas evaluation is used to refer to the comparison against other  
148 coarse-resolution satellite or modelled data with supposedly less well-defined uncertain-  
149 ties. However, technically, validation more specifically refers to quantitative data quality

150 assessment (*Justice et al.*, 2000) whereas evaluation more broadly refers to “the process of  
151 judging something’s quality” (*Loew et al.*, 2017). For simplicity, we use the term *validation*  
152 in this paper to refer to the process of estimating biases and uncertainties regardless of  
153 the reference data source used;

- 154 • the concept of uncertainty is closely related to the concept of confidence intervals. Both  
155 aim at describing the pdf underlying an estimate, although the term *uncertainty* is more  
156 commonly used for describing the pdf behind an estimate that results from measurement  
157 errors (see Sec. 4.1), whereas the term *confidence interval* is more commonly used for  
158 describing the pdf behind statistical parameters (such as statistical moments or validation  
159 metrics that derive from these moments) that results from finite sample sizes (see Sec.  
160 4.5).

### 161 3 Reference data

162 The term *fiducial reference measurements* is often used to refer to a suite of independent, fully  
163 characterized, and traceable measurements that meet the requirements on *reference standards*  
164 as described by QA4EO (*Fox*, 2010), which should be used to assess the quality of EO prod-  
165 ucts. However, although highly accurate in situ soil moisture measurement techniques exist and  
166 uncertainties of the measurement devices can be reliably determined through laboratory and  
167 field calibration (*Cosh et al.*, 2004, 2006; *Rüdiger et al.*, 2010), using such point-scale measure-  
168 ments for validating satellite soil moisture data sets over large areas is a very difficult task owing  
169 to the coarse resolution of space borne microwave instruments and vast heterogeneities across  
170 landscapes (*Famiglietti et al.*, 2008; *Brocca et al.*, 2010a; *Miralles et al.*, 2010; *Crow et al.*, 2012;  
171 *Nicolai-Shaw et al.*, 2015; *Molero et al.*, 2018).

172 For satellite validation purposes, numerous field and airborne campaigns have been carried  
173 out to obtain reliable satellite footprint scale reference data and to quantitatively assess the  
174 potential spatio-temporal representativeness (see Sec. 4.2) of single or small sets of in situ soil  
175 moisture stations (*De Rosnay et al.*, 2006; *Brocca et al.*, 2012; *McNairn et al.*, 2015). Addition-  
176 ally, validation activities are complemented with land surface model output and other satellite  
177 products for comparison to get as complete a picture as possible of a product’s error character-  
178 istics (*Brocca et al.*, 2010b; *Draper et al.*, 2013; *Al-Yaari et al.*, 2014; *Dorigo et al.*, 2015; *Kerr*  
179 *et al.*, 2016; *Miyaoka et al.*, 2017). The various reference data sources and their limitations are

180 discussed below. A list of publicly available reference data sources that are commonly used for  
181 satellite soil moisture validation is provided in Table 2.

### 182 **3.1 Field campaigns**

183 Field campaigns are labour-intensive studies that use highly accurate measurement techniques  
184 to obtain reliable and traceable representations of larger scale average soil moisture. Unfortu-  
185 nately, these campaigns provide only some snapshots in time, whereas the validation of satellite  
186 products requires long and consistent time series (see Sec. 4.4). Therefore, some field cam-  
187 paigns have identified and set up a limited number of permanent measurement stations (<15) at  
188 temporally stable locations (*Vachaud et al.*, 1985; *Starks et al.*, 2006) that sufficiently capture  
189 sub-pixel heterogeneities, allowing the continuous observation of satellite footprint-scale areas  
190 with sufficient and well characterized accuracy.

191 Ground measurements are often supplemented with airborne observations, which can be used  
192 to either directly validate the L1 satellite measurements or the derived soil moisture retrievals  
193 over a much larger area, allowing to evaluate spatial soil moisture variability across multiple  
194 satellite grid cells. Moreover, they can provide valuable information about soil moisture (or  
195 backscatter/brightness temperature) sub-pixel variability.

196 Early field campaigns were focused on understanding large-scale soil moisture dynamics with  
197 aircraft support such as HAPEX-MOBILHY (*Noilhan et al.*, 1991), BOREAS (*Cuenca et al.*,  
198 1997), and the Washita'92 campaigns (*Jackson et al.*, 1995), assessing the potential of soil  
199 moisture monitoring as a part of hydrologic modeling. This evolved into satellite associated  
200 field campaigns such as the 1997 Southern Great Plains Hydrology Experiment (SGP97) and  
201 the Soil Moisture Experiments (SMEX) in 2002-2004 in the United States (*Jackson et al.*, 1999,  
202 2005; *Bindlish et al.*, 2006, 2008), the National Airborne Field Experiments (NAFE) in Australia  
203 (*Panciera et al.*, 2008), the Australian Airborne Calibration/Validation Experiments for SMOS  
204 (AACES; *Peischl et al.*, 2012), the Canadian Experiment in Soil Moisture (CANEX-10; *Magagi*  
205 *et al.*, 2013), and the CAROLS airborne campaigns (*Albergel et al.*, 2011; *Zribi et al.*, 2011).  
206 These campaigns established a protocol for the synchronous collection of ground-based soil  
207 moisture measurements with airborne microwave instrumentation, which were supplemented  
208 with long-term in situ monitoring stations, thus providing long-term high density validation  
209 sites for satellites.

210 In the process of developing standardized data collection protocols, these field campaigns

211 specifically focused on the investigation of the spatial distribution of soil moisture and its evo-  
212 lution with drying or wetting, the soil moisture variability across scales, and the statistical  
213 relationship between spatial standard deviation and extent scale. These parameters drive the  
214 potential representativeness of in situ measurements for coarse soil moisture product validation  
215 and their knowledge hence allows the determination of the number of ground samples required  
216 to obtain sufficiently reliable validation reference data (*Famiglietti et al.*, 2008).

## 217 **3.2 In situ networks**

218 A large number of in situ soil moisture networks exist worldwide with different quality and  
219 spatial sampling densities as well as varying sensing depths (*Dorigo et al.*, 2011b; *Babaeian*  
220 *et al.*, 2019). For validation purposes, the soil moisture community distinguishes between dense  
221 networks, which have a large number of soil moisture stations located within single satellite  
222 footprints, and sparse networks, where footprint-scale areas usually contain only a single or very  
223 few soil moisture stations, although the quantitative cut-off between the two is not well-defined.  
224 The overall global coverage of in situ soil moisture networks (accessible and suited for satellite  
225 soil moisture validation) is unevenly distributed across the globe and particularly scarce in the  
226 tropical regions, the Southern Hemisphere and boreal regions (Fig. 2; *Ochsner et al.*, 2013).

### 227 **3.2.1 Dense networks**

228 To meet the requirements on fiducial reference data (*Fox*, 2010), the SMAP Calibration and  
229 Validation (Cal/Val) Team defined certain criteria for dense measuring networks, so-called core  
230 validation sites, ensuring that they provide a traceable representation of footprint-scale soil  
231 moisture and therefore allow for a reliable assessment of satellite soil moisture data quality.  
232 Currently, 18 densely stationed and thoroughly calibrated in situ measurement sites fulfill these  
233 requirements (*Jackson et al.*, 2012; *Colliander et al.*, 2017), operated by independent SMAP  
234 Cal/Val partners.

235 These SMAP Cal/Val partners have a diverse heritage. Some networks were deployed for  
236 Cal/Val of the AMSR-E product (*Jackson et al.*, 2010) or SMOS (*Djamai et al.*, 2015), while  
237 others evolved from hydrologic monitoring networks (*Bogena et al.*, 2018) or from some other  
238 purpose such as aircraft validation projects like AIRMOSS (*Moghaddam et al.*, 2010). During the  
239 SMAP project, several networks were selected as potential candidate sites for Cal/Val activities.  
240 The candidate networks whose accuracy versus physically collected volumetric soil moisture was

241 already demonstrated and documented in a traceable manner, were promoted to core validation  
242 sites. To date, these sites are considered to provide the best possible ground reference data for  
243 satellite footprint-scale soil moisture dynamics (*Colliander et al.*, 2017).

### 244 **3.2.2 Sparse networks**

245 A host of other operational and experimental in situ sites exist worldwide, operating soil mois-  
246 ture measurement stations that are potentially suited for soil moisture validation yet with a  
247 considerably smaller station density and often lacking information on their coarse-scale repre-  
248 sentativeness and their own inherent error characteristics (*Gruber et al.*, 2013a; *Chen et al.*,  
249 2017). Nonetheless, these sites are valuable to complement core validation sites due to their  
250 considerably larger spatial coverage across a variety of climatic regimes and biomes (see Sec. 4).

251 An important source for data from sparse networks is the International Soil Moisture Network  
252 (ISMN; *Dorigo et al.*, 2010, 2011b), which is a data hosting facility that harmonizes soil moisture  
253 measurements from in situ networks worldwide, applies automated and uniform quality control  
254 procedures to flag suspicious measurements (*Dorigo et al.*, 2013), and distributes them on their  
255 website (<http://ismn.geo.tuwien.ac.at/>; last access: 1 July 2019) on a cost-free basis in  
256 a common format. The ISMN was established by ESA in the framework of SMOS Cal/Val  
257 activities. Currently, it contains data from more than 2400 stations worldwide, operated across  
258 59 different measurement networks (see Figure 2) including historical networks that are no longer  
259 operational.

### 260 **3.3 Model simulations**

261 Due to the limited coverage and representativeness of ground reference data, validation activ-  
262 ities are complemented with soil moisture simulations from land surface models (LSMs) as an  
263 alternative reference data source (*Lahoz and De Lannoy*, 2014). Model simulations can provide  
264 spatially complete global soil moisture maps at a spatial (grid) resolution similar to that of satel-  
265 lite footprints, but they may still contain considerable representativeness errors (see Sec. 4.2)  
266 originating from simplifications of sub-grid heterogeneities, a scale-mismatch of the underlying  
267 atmospheric forcing data, errors in the model parameterization, or simply because the meaning  
268 of the modelled “soil moisture” is different. Moreover, biases and uncertainties in model simu-  
269 lations are highly variable and often also not well quantified (*Koster et al.*, 2009; *Albergel et al.*,  
270 2013), making it difficult to separate satellite retrieval errors from modelling errors in a direct

271 comparison (see Sec. 4).

272 Some examples of readily available global model-based data sets that have been used for  
273 satellite soil moisture validation activities (*Albergel et al.*, 2012; *Al-Yaari et al.*, 2014; *Kerr et al.*,  
274 2016; *Dorigo et al.*, 2017; *Gruber et al.*, 2017; *Miyaoka et al.*, 2017) include simulations from  
275 NASA’s Global Land Data Assimilation System (GLDAS; *Rodell et al.*, 2004), NASA’s Modern-  
276 Era Retrospective analysis for Research and Applications (MERRA) land data products (*Reichle*  
277 *et al.*, 2011, 2017a), and the European Center for Medium-Range Weather Forecasts (ECMWF)  
278 Land Surface Reanalysis (ERA-Interim/Land) data sets (*Balsamo et al.*, 2015).

### 279 3.4 Satellite products

280 A multitude of soil moisture products from different satellite sensors (*Babaeian et al.*, 2019)  
281 are commonly used as additional coarse resolution reference data sets for validation purposes,  
282 either for consistency assessment through direct comparison (*Al-Yaari et al.*, 2014; *Burgin et al.*,  
283 2017), or within triple collocation analysis (*Dorigo et al.*, 2010; *Draper et al.*, 2013, see Sec. 4).  
284 Like model simulations and sparse networks, they typically lack reliable and traceable bias and  
285 uncertainty characterization. Also, available satellite sensors observe at different wavelengths,  
286 polarizations, and incidence angles and have therefore a varying sensitivity to soil moisture  
287 (*Ulaby et al.*, 2014). Hence, the information gleaned from a direct comparison is limited (see  
288 Sec. 4.4.2). Furthermore, different satellite retrieval products (and model simulations) can use  
289 similar ancillary information such as temperature and/or vegetation information in a radiative  
290 transfer model, resulting in correlated errors (*Gruber et al.*, 2016b) which may complicate a fair  
291 data comparison (see Sec. 4.4.2). Comprehensive lists of commonly used and publicly available  
292 satellite soil moisture products, including some validation information where available, can be  
293 found at <https://lpvs.gsfc.nasa.gov/producers2.php?topic=SM> (last access: 1 July 2019)  
294 and in *Babaeian et al.* (2019).

## 295 4 Theory

296 This section provides the theoretical background for error characterization and how it relates to  
297 satellite soil moisture validation, including the assumptions, limitations and pre-processing steps  
298 involved. Although our main focus here is the validation of near-surface satellite soil moisture  
299 products, many of the principles discussed below can be equally applied to assess the quality

300 of soil moisture products from other sources, as well as of other biogeophysical variables (*Loew*  
301 *et al.*, 2017).

## 302 4.1 Errors

303 A measurement error  $e_x$  is defined as the deviation of a measurement  $x$ , in our case a satellite  
304 soil moisture retrieval, from the true state  $t$  of the quantity under observation (*JCGM*, 2008):

$$e_x = x - t \tag{1}$$

305 Important for understanding errors is that the “truth” is a hypothetical concept. For the case  
306 of space borne microwave measurement instruments, actual satellite footprints are overlapping  
307 elliptical areas with strong signal intensity gradients from the footprint center outwards (de-  
308 pending on the antenna gain pattern) and varying surface property dependent vertical support  
309 (*Ulaby et al.*, 2014). Horizontal footprint boundaries are commonly defined as the -3 dB region  
310 (i.e. the antenna main beam region that covers 50% of the signal’s power). Products derived  
311 thereof are typically sampled onto spatial grids with sharp boundaries between grid cells and a  
312 constant layer depth to facilitate further geospatial analysis (*Bartalis et al.*, 2006; *Brodzik et al.*,  
313 2012; *Bauer-Marschallinger et al.*, 2014). The “true” soil moisture signal that drives the mi-  
314 crowave measurement and the subsequent gridded soil moisture retrieval will therefore never be  
315 the real average soil moisture of the grid cell to which a measurement is assigned. Moreover, for  
316 validation purposes, the unknown “truth” is approximated by reference data, which themselves  
317 contain errors and may also be driven by a soil volume that is different from the satellite grid  
318 cell they are supposed to represent (see Sec. 3).

## 319 4.2 Representativeness

320 The difference between the true soil moisture that actually affects a (microwave) measurement  
321 associated with a particular grid cell and the true soil moisture within that grid cell is often  
322 referred to as representativeness error (*Gruber et al.*, 2016a). However, it is worth noting that  
323 representativeness errors have different definitions (*Van Leeuwen*, 2015). The remote sensing  
324 community mostly assigns them to the mismatch between the spatial support of a measurement  
325 and the spatial resolution of the defined sampling grid, sometimes also referred to as scaling  
326 error (*Miralles et al.*, 2010; *Crow et al.*, 2012; *Gruber et al.*, 2013a; *Molero et al.*, 2018). In

327 the modelling community, representativeness errors mostly refer to a model’s lacking ability to  
328 represent reality and, as such, to imperfections in the model structure and in parameterization  
329 (e.g., unresolved sub-grid scale processes). For the purpose of data validation, it is practical  
330 to use a definition that potentially allows us to separate representativeness errors from other  
331 error sources upon estimation. Therefore, recall that the general definition of error in Eq. (1)  
332 requires the choice of a “truth”, which is the soil moisture state within a target volume (grid  
333 cell) that one aims to estimate as accurately as possible. We define representativeness errors  
334 as those deviations of a product from that chosen “true” state, which are related to real soil  
335 moisture variations. They can occur, for example, if the actual measurement footprint of a  
336 satellite extends beyond the grid cell boundaries associated with the “truth”, if an inadequate  
337 soil parameterization in a radiative transfer model causes the soil moisture retrievals to represent  
338 deeper soil layers than the intended “truth”, or if point-scale ground measurements are used  
339 as a reference for grid cell-scale soil moisture dynamics. As such, representativeness errors of  
340 different data sets may be correlated even if the products are otherwise independent.

341 In summary, representativeness errors have important implications for validation in that  
342 they limit the information one can glean from the comparison between products, even if a  
343 chosen reference product is itself highly accurate (see Sec. 4.4.1). Since the temporal and  
344 spatial resolution and sampling of satellite and available reference measurements hardly ever  
345 match, (relative) representativeness errors will often reach considerable magnitudes (*Miralles*  
346 *et al.*, 2010; *Crow et al.*, 2012). To minimize their influence, several pre-processing steps are  
347 typically applied, which are discussed in the following section together with other pre-processing  
348 steps that are necessary before validation metrics can or should be calculated.

### 349 **4.3 Pre-processing**

350 Pre-processing steps necessary for validation aim to find match-ups in space and time between  
351 measurements that have different spatial resolutions, are sampled on to different grids, and/or  
352 are acquired at different times. Additionally, depending on the reference data choice, statis-  
353 tical rescaling methods are often applied to minimize the impact of representativeness errors.  
354 Moreover, data pre-processing typically involves the masking of unreliable satellite retrievals  
355 and reference measurements. Lastly, data sets are sometimes decomposed into different fre-  
356 quency components in order to separately assess a product’s ability of accurately representing  
357 short-term, seasonal, and inter-annual soil moisture variability (*Draper and Reichle*, 2015).



### 358 4.3.1 Data masking

359 Satellite-derived soil moisture products are typically accompanied by a set of quality flags. They  
360 can be indicators of suspected contamination of the microwave signals or problems during the  
361 retrieval. Typical examples are indicators for the probability of frozen soil, dense vegetation  
362 coverage, radio frequency interference (RFI), or urban or water contamination, to name a few  
363 (e.g., *Parinussa et al.*, 2011; *Naeimi et al.*, 2012; *Kerr et al.*, 2012; *de Nijs et al.*, 2015). The  
364 validation of a product should be based only on those retrievals that are considered “good” for  
365 later application. While masking data points using binary “use / do not use” flags is straight-  
366 forward, some quality flags require the decision of a threshold below or above which individual  
367 retrievals are masked out (e.g., the probability of RFI occurrence or the water body fraction),  
368 which implies a trade-off between data quality and measurement density. Typically, data pro-  
369 ducers provide recommendations for these thresholds. In addition to the quality flags inherent  
370 in the soil moisture products, auxiliary static and/or dynamic data from land surface models or  
371 other sources are often used to mask out retrievals that can be considered unreliable, although  
372 it should be kept in mind that these sources themselves - and hence quality flags derived thereof  
373 - are subject to errors. The most commonly used masking criteria are based on surface and/or  
374 air temperature and snow height and/or snow water equivalent estimates obtained from land  
375 surface models, or vegetation estimates from satellite sensors or models (*Al-Yaari et al.*, 2014;  
376 *Dorigo et al.*, 2015; *Gruber et al.*, 2017). Note that reference data sets, in particular in situ  
377 measurements, also often undergo quality control procedures and provide quality flags, which  
378 should be used to mask out unreliable measurements before using them to validate satellite  
379 retrievals (as is the case for example for the ISMN; *Dorigo et al.*, 2013).

380 When comparing biases or uncertainties of different soil moisture products, the masking  
381 procedures applied to these data sets should be identical in order to compare the quality of  
382 retrievals from measurements that were taken under the same (or at least similar) conditions.  
383 However, if quality flags that are tailored to one data set are applied to another, some of the  
384 products may appear better or worse than they would when using only their own inherent  
385 quality control. This is especially true if the flags of one product are much more conservative  
386 than those of another. Most product comparison studies do not take this issue into account. One  
387 possible approach to address it would be to compare biases and uncertainties from collocated  
388 periods also with those in periods where only some products provide unflagged soil moisture  
389 retrievals (based on their own quality control) and to put this into perspective with the temporal

390 measurement density before and after product collocation. However, this requires the availability  
391 of appropriate reference data in collocated and non-collocated periods as well as the ability to  
392 account for possibly varying accuracy and representativeness of the reference data in these  
393 periods. Also, depending on the overall data density, it may be difficult to assess biases and  
394 uncertainties in these periods due to the presence of large statistical sampling errors (see Sec.  
395 4.5).

396 Finally, we stress that the choice of data masking criteria has a considerable impact on the  
397 overall validation results and should be carefully documented, especially for comparing different  
398 validation studies and when assessing long-term changes.

### 399 **4.3.2 Collocation**

400 Satellite sensors acquire measurements that are irregularly distributed in space and time owing  
401 to their orbiting nature and specific antenna patterns. In the soil moisture retrieval process,  
402 these measurements are typically sampled onto spatial grids (for noise reduction purposes these  
403 grids are often oversampled, i.e. the grid sampling - sometimes also referred to as grid posting -  
404 is typically higher than the antenna resolution) and sometimes also to regular time steps (e.g.,  
405 00:00 UTC) in order to generate, for example, daily global soil moisture maps and/or time  
406 series (*Kerr et al.*, 2012; *O'Neill et al.*, 2012; *H-SAF*, 2018; *Gruber et al.*, 2019a). However,  
407 neither the resolution nor the sampling of in situ reference measurements or model simulations  
408 ever perfectly match those of the satellite products being validated. Consequently, the process  
409 of finding match-ups between satellite and reference data points in space and time, commonly  
410 referred to as collocation, is essentially a resampling task (*Loew et al.*, 2017). Since the spatial  
411 resolution of the compared products can be very different (especially between in situ and satellite  
412 / modelled data), statistical rescaling methods are often additionally applied in the collocation  
413 process to minimize the impact of (especially spatial) representativeness errors on validation  
414 metrics.

#### 415 **Spatial resampling**

416 In situ measurements are point-scale measurements that sample only a few cubic centimeters  
417 of the soil (with the exception of cosmic-ray neutron sensors, which sample areas in the order  
418 of hectares; *Zreda et al.*, 2012). When used for validating satellite products, stations from  
419 sparse networks are typically sampled onto the satellite grid using a nearest-neighbour (NN)

420 search, i.e. by matching the stations to the satellite grid cells within which they are located  
421 (*Albergel et al., 2012; Dorigo et al., 2015; Chen et al., 2017*). For dense networks, commonly all  
422 stations that lie within a particular satellite grid cell are (after quality control) averaged (*Jackson*  
423 *et al., 2010; Gruber et al., 2015; Colliander et al., 2017*), either by calculating the arithmetic  
424 mean or by calculating a weighted average where higher weights are applied to stations that are  
425 expected to be more representative for the grid cell average soil moisture. Such stations can be  
426 identified, for example, via a temporal stability analysis (*Vachaud et al., 1985*), through Voronoi  
427 diagrams (*Colliander et al., 2017*), or by using landscape characteristics such as land cover or  
428 soil properties.

429 When comparing different gridded products (i.e. different satellite and/or land surface model  
430 products), one grid must be selected as the reference grid onto which the other products are  
431 resampled for collocation purposes. This is commonly done using either a NN search or inverse-  
432 distance-weighted (IDW) based approaches (*Al-Yaari et al., 2014; Gruber et al., 2017, 2019a*).  
433 However, the resampling provides mainly spatial match-ups of the data sets and can at best  
434 account for some of the spatial representativeness errors of the various data sets. How exactly  
435 these representativeness errors are affected and propagate into bias and uncertainty estimates will  
436 depend on the chosen reference grid and resampling method, and requires more research. The  
437 most common way to reduce spatial (systematic) representativeness errors is to apply statistical  
438 rescaling methods (see below).

### 439 **Temporal resampling**

440 In situ measurements and land model estimates are typically sampled more frequently than  
441 satellite soil moisture retrievals. Therefore, the reference measurements are matched in time  
442 to the irregular satellite observation times, typically by selecting the temporally closest (NN)  
443 reference measurement within a pre-defined search window (i.e. applying a maximum temporal  
444 distance threshold; *Chen et al., 2017*). Depending on the sampling interval of the reference  
445 data sets (for in situ data typically hourly and for global land surface models typically one to  
446 six hourly) and on whether or not satellite observations have been a priori resampled already  
447 (see above), this can lead to considerable differences between the actual measurement times of  
448 collocated satellite and reference data points. The issue is typically limited when using in situ  
449 or model data as reference. However, if multiple satellite products are evaluated simultaneously,  
450 their different overpass times are usually accounted for by either picking one of them as (tem-

451 poral) reference and matching the other ones against it, or by sampling all satellite products  
452 to regularized time steps (e.g., 00:00 UTC; *Gruber et al.*, 2017), which in any case favours the  
453 satellite data set whose actual measurement times are closest to the reference points. Note that  
454 the retrieval quality of satellite data sets may strongly depend on the time of observation. This  
455 is especially true for passive systems, where soil moisture retrievals are known to be strongly  
456 affected by temporal temperature fluctuations and temperature gradients in soil and vegetation  
457 cover (*Parinussa et al.*, 2015).

458 Taken together, the different measurement times of satellite and reference data sets that  
459 have been collocated will induce temporal representativeness errors, originating from the actual  
460 soil moisture changes that take place during these periods. Often these errors are assumed to be  
461 negligible or at least below the noise level of the products. In principle, one could employ more  
462 sophisticated resampling algorithms to minimize these representativeness errors, for example  
463 auto-regressive interpolation methods with or without auxiliary information such as precipita-  
464 tion, evapotranspiration, or soil texture. However, more research is needed to assess the impact  
465 of temporal interpolation approaches on validation metrics.

#### 466 **(Statistical) rescaling**

467 The resampling procedures described above provide data set match-ups in space and time which  
468 are required for statistical comparison (see Sec. 4.4). As discussed in Sec. 4.1, the measurements  
469 of the collocated products are driven by the soil moisture state of different soil volumes at  
470 different times due to the different underlying actual spatio-temporal resolution of the data  
471 sets. The latter is related to the antenna and surface properties and cannot be corrected for by  
472 common resampling methods. Therefore, a direct comparison of these products will be subject  
473 to representativeness errors, which may dominate the total soil moisture retrieval errors (*Gruber*  
474 *et al.*, 2013a; *Chen et al.*, 2017; *Molero et al.*, 2018). However, owing to the large-scale and  
475 auto-correlated nature of processes that drive soil moisture changes (*Crow et al.*, 2012), parts  
476 of these errors are systematic and can hence be corrected for by removing *relative differences*  
477 between the considered data sets (see Sec. 4.4).

478 The two most common rescaling approaches are to match either the temporal mean and  
479 standard deviation of the data sets that are to be compared (*Scipal et al.*, 2008a; *Dorigo et al.*,  
480 2010; *Albergel et al.*, 2012), or to match their complete cumulative distribution function (CDF),  
481 which additionally corrects for differences in higher statistical moments in case the products

482 are expected not to be perfectly Gaussian distributed (*Reichle and Koster, 2004; Kumar et al.,*  
483 2012). However, any rescaling approach that transforms one data set into the data space of  
484 another (without additional information) assumes the signal-to-noise ratios (SNRs) of the two  
485 involved data sets to be identical, which, since this is usually not the case, can lead to biased  
486 rescaling parameters that do not fully correct the systematic representativeness errors (see Sec.  
487 4.4.2; *Stoffelen, 1998; Yilmaz and Crow, 2013*). Alternatively, triple collocation analysis (*Stof-*  
488 *felen, 1998; Su et al., 2014; Gruber et al., 2016a*) is often employed, using a third data set to take  
489 different SNRs into account when matching the standard deviation of the underlying soil mois-  
490 ture signals, thereby potentially providing consistent rescaling parameters (*Yilmaz and Crow,*  
491 2013).

492 Note that rescaling soil moisture data sets can equally account for (systematic) represen-  
493 tativeness errors that arise from different spatial resolution and spatial and temporal mis-  
494 alignment, as well as for those arising from different vertical measurement support, i.e. wavelength-  
495 dependent penetration depths of satellites, in situ sensor placement depths, and modelled soil  
496 layer thickness (*Gruber et al., 2013a*). Also, in addition to correcting for systematic repre-  
497 sentativeness errors, rescaling can implicitly compensate for different units (provided that the  
498 used soil moisture representations are linearly related), most commonly volumetric soil moisture  
499 ( $[m^3m^{-3}]$ ) and the degree of soil saturation ( $[%]$ ) which are linked through soil porosity as a  
500 multiplicative factor (*Walker et al., 2004*). This avoids additional biases that are introduced  
501 through the use of inaccurate auxiliary data (such as soil maps) that would otherwise be needed  
502 for unit conversion.

503 After rescaling, long-term bias estimation is obviously no longer meaningful as systematic  
504 differences between the data sets, which would normally serve as proxy for biases, have been  
505 intentionally removed. However, shorter-term biases as well as random representativeness errors  
506 may remain and can considerably contribute to subsequent uncertainty estimates (see Sec. 4.4.1).

### 507 4.3.3 Signal decomposition

508 The quality of soil moisture products can vary considerably across time scales (*Su and Ryu, 2015;*  
509 *Draper and Reichle, 2015; Molero et al., 2018; Gruber et al., 2019a*). For example, some soil  
510 moisture products are better at accurately representing the seasonal cycle whereas other products  
511 more accurately capture short-term fluctuations. Therefore, products are often decomposed into  
512 different frequency components which are then validated separately (in addition to the bulk

513 time series). In Earth sciences, such decomposition is often done using moving-average windows  
514 (*Narapusetty et al.*, 2009). For soil moisture, a moving window of several weeks, centered on the  
515 measurement time, is typically used to obtain intra-annual low-frequency soil moisture dynamics  
516 (*Albergel et al.*, 2012; *Chen et al.*, 2017), referred to as seasonalities. Residuals thereof are  
517 referred to as short-term anomalies which represent higher-frequency, sub-seasonal soil moisture  
518 variations. Additionally, so-called long-term anomalies are often calculated as residuals relative  
519 to a multi-year mean seasonal cycle, referred to as the soil moisture climatology, which is typically  
520 calculated by applying a moving-average window of similar size (a few weeks) to each day-of-  
521 the-year (DOY), i.e. averaging all measurements of all years that fall inside the specified time  
522 window around a particular DOY (*Miralles et al.*, 2010; *Draper et al.*, 2013).

523 While the validation of short-term soil moisture anomalies aims at assessing a data set's  
524 capability of capturing individual drying or wetting events, uncertainties of long-term anomalies  
525 represent its performance in capturing both short-term variability and inter-annual variations  
526 such as prolonged droughts or floods as well as climate trends. However, the latter rely on a  
527 climatology estimate that requires historical data records in the order of decades (*Dorigo et al.*,  
528 2012), which are often not available, especially not at the beginning of a new mission (current  
529 microwave missions cover a time period of maximum 5-10 years). Therefore, one often has to  
530 rely on uncertainty estimates for seasonalities and short-term anomalies alone, which jointly  
531 drive uncertainties in long-term anomalies.

#### 532 4.4 Metrics

533 After satellite and reference products have been masked, collocated, and optionally decomposed  
534 and/or rescaled, validation metrics can be calculated. In this section, we summarize commonly  
535 used bias and uncertainty estimators and their underlying assumptions. Other related metrics  
536 exist (e.g., the mean absolute error, Kendall's tau, and many others), but all are derived from  
537 the same statistical moments and have therefore similar information content. Our goal here is to  
538 present the metrics that are most commonly used for soil moisture validation and are considered  
539 to provide a comprehensive picture of a product's error characteristics. These metrics also  
540 largely coincide with those used in other EO communities (*Loew et al.*, 2017). We also stress  
541 that validation specifically aims at quantitatively assessing the errors of a data set, which is  
542 different from indirectly evaluating its quality for example by investigating its skill in a particular  
543 application, e.g., drought monitoring (*Bolten et al.*, 2010). Such indirect product evaluation is

544 beyond the scope of this paper.

#### 545 4.4.1 Assumptions

546 The fundamental assumption underlying almost all satellite soil moisture validation studies is  
547 that of additive zero-mean random errors ( $\varepsilon_x$ ), and additive (first-order;  $\alpha_x$ ) and multiplicative  
548 (second-order;  $\beta_x$ ) systematic errors (*Gruber et al.*, 2016a):

$$x = \alpha_x + \beta_x t + \varepsilon_x \quad (2)$$

549 This error model applies to both the data set one aims to validate and the reference data sets.  
550 Notice that the total error  $e_x$  in Eq. (1) has now been separated into its systematic ( $\alpha_x$  and  $\beta_x$ )  
551 and random ( $\varepsilon_x$ ) components. These components contain instrument errors (i.e. noise and mis-  
552 calibration), errors in the retrieval model and parameterization, and other representativeness  
553 errors with respect to the assumed grid cell average soil moisture  $t$  (although the boundaries  
554 between the latter two are somewhat fuzzy; see Sec. 4.1).

555 To disentangle errors from different data sets and from actual soil moisture variations, all  
556 common data comparison metrics require the errors to be homoscedastic (i.e. independent from  
557 the soil moisture state, in the literature often referred to as orthogonality with respect to the  
558 truth; *Yilmaz and Crow*, 2014) and mutually uncorrelated between products. Remember, how-  
559 ever, that the *representativeness* error components of the different products may (by definition)  
560 be correlated both with the truth  $t$  and with each other, even if the products are otherwise  
561 independent (see Sec. 4.1).

562 All common validation metrics are derived from the first and second statistical moments of  
563 the data sets. This implies that soil moisture too is - even though in principle deterministic -  
564 assumed to behave as a random variable. Statistical moments are then typically estimated in  
565 the temporal domain (i.e. temporal means, variances, and covariances), assuming stationarity  
566 in soil moisture and the errors (i.e. means and variances are assumed to be constant over time),  
567 and relate to the various error components as follows:

$$\begin{aligned} \bar{x} &= \alpha_x + \beta_x \bar{t} \\ \sigma_x^2 &= \beta_x^2 \sigma_t^2 + \sigma_{\xi_x}^2 \\ \sigma_{xy} &= \beta_x \beta_y \sigma_t^2 + \sigma_{\xi_x, \xi_y} \end{aligned} \quad (3)$$

568 where the overline,  $\sigma_i^2$  and  $\sigma_{ij}$  refer to the (temporal) mean, variance, and covariance, respec-  
569 tively; and  $y$  denotes a reference data set that follows the same error model as  $x$  (Eq. (2)).  
570 Because *representativeness* errors may contain an orthogonal, a non-orthogonal, and a mutually  
571 correlated component (see above), we combine it with all other random error in the individual  
572 data set’s random error variability  $\sigma_{\xi_x}^2 = \sigma_{\varepsilon_x}^2 + 2\beta_x\sigma_{t,\varepsilon_x}$  (containing representativeness and all  
573 other random errors) and the correlated error variability  $\sigma_{\xi_x,\xi_y} = \beta_x\sigma_{t,\varepsilon_y} + \beta_y\sigma_{t,\varepsilon_x} + \sigma_{\varepsilon_x,\varepsilon_y}$  (driven  
574 by representativeness errors only), for clarity. Systematic representativeness errors are included  
575 in the  $\alpha_x$  and  $\beta_x$  coefficients.

576 The goal of validation is now to estimate  $\alpha_x$  and  $\beta_x$ , and the standard deviation of  $\varepsilon_x$  ( $\sigma_{\varepsilon_x}$ ),  
577 i.e. biases and uncertainties in the satellite data set under validation. The properties of the  
578 different reference data sets available (see Sec. 3) determine which error components will be  
579 dominant in Eq. (3), and consequently, which ones can be estimated by the available validation  
580 metrics (see Sec. 4.4.3 and 4.4.4).

#### 581 4.4.2 Relative and TCA-based metrics: opportunities and limitations

582 For discussing the various metrics we will follow the notation of fiducial reference data (see Sec.  
583 3) to refer to data sets that provide a thoroughly calibrated soil moisture proxy at the satellite  
584 scale with traceable uncertainty characteristics (i.e.  $\alpha_y \approx 0, \beta_y \approx 1$  in Eq. (2)).  $\varepsilon_y$  may be  
585 non-zero but  $\sigma_{\varepsilon_y}^2$  has to be at least well determined from laboratory experiments and field cam-  
586 paigns and could hence be corrected for in the validation metrics. As mentioned, only the core  
587 validation sites are currently considered as fiducial reference data capable of providing a reliable  
588 representation of satellite footprint-scale soil moisture (see Sec. 3.2.1). They are therefore the  
589 only reliable proxy for bias and uncertainty estimation from direct comparison, but are limited  
590 to very few regions. Non-fiducial reference data refer to coarse-resolution products such as land  
591 surface model simulations or other satellite data sets which may have non-negligible or non-  
592 traceable biases and uncertainties as well as potentially considerable representativeness errors,  
593 or to in situ data from sparse networks or not properly calibrated and validated dense networks,  
594 both of which are expected to have larger representativeness errors than coarse-resolution refer-  
595 ence data sets. Therefore, direct comparison against non-fiducial reference data can only provide  
596 information of which data set is systematically drier or wetter than the other but without rela-  
597 tion to a true grid cell average, and only lumped estimates of the uncertainty of both compared  
598 products. Nonetheless, given their larger-scale and long-term availability, sparse networks and



599 land surface models are of important complementary value for validating satellite products. In  
600 particular, one can obtain valuable information about the relative ranking of different products  
601 as well as about performance changes over time when comparing against the same reference  
602 product.

603 Introducing a second reference data set  $z$  that follows the same covariance properties (Eq.  
604 (3)) as  $y$  (commonly referred to as triple collocation analysis, TCA; *Stoffelen, 1998; Scipal et al.,*  
605 *2008b; Gruber et al., 2016a*) allows, under particular circumstances, to simultaneously estimate  
606 uncertainties of all three products and also to (partly) isolate random (relative) representative-  
607 ness errors (*Miralles et al., 2010; Gruber et al., 2013a; Chen et al., 2017*). Note, however, that  
608 the necessity of using two reference data sets instead of one may limit spatial and temporal data  
609 availability. Moreover, while non-orthogonal and mutually correlated errors are equally prob-  
610 lematic for metrics that rely on one reference data set only (see below), it may be even more  
611 difficult to find a third data set that fulfills these requirements. Commonly, any combination  
612 of in situ measurements, land surface model estimates, active-microwave-based measurements,  
613 or passive-microwave-based measurements is expected to fulfil this requirement because their  
614 sources of errors are assumed to be mostly independent (*Gruber et al., 2016a*), provided that  
615 neither of them has been used to generate another (e.g., by assimilating satellite data in to a  
616 land surface model; *Reichle et al., 2017b,c*). However, several studies suggest that mutual error  
617 correlations may exist between commonly used data set combinations (*Yilmaz and Crow, 2014;*  
618 *Pan et al., 2015*), resulting from unrecognized common data (e.g., similar vegetation or temper-  
619 ature input) or representativeness errors (e.g., if a land surface model used within TCA models  
620 a deeper layer than the sensing depth of two satellite data sets that are used in the triplet). It is  
621 therefore recommended to verify orthogonality and zero error correlation assumptions by using  
622 - where available - multiple data set triplets and checking for consistency between different TCA  
623 implementations (*Dorigo et al., 2010; Draper et al., 2013*), or by using the recently proposed  
624 TCA extension that utilizes four or more data sets to diagnose the existence, and estimate the  
625 magnitude of error correlations (*Gruber et al., 2016b; Pierdicca et al., 2017*).

626 The following sections discuss the most common bias and uncertainty metrics, either (i)  
627 based on direct comparison between two data sets, which will be referred to as relative metrics,  
628 or (ii) based on the simultaneous comparison of three products, which will be referred to as  
629 TCA-based metrics. All metrics can be equally applied to soil moisture anomaly estimates or  
630 the raw time series, except for first-order bias estimators (see below) as the anomaly calculation

631 per definition removes differences in the mean (see Sec. 4.3.3).

632 Note that none of the metrics presented below require assumptions about the shape of the  
633 pdf of the random errors or the true signal (*McColl et al.*, 2016). However, the bounded nature  
634 of soil moisture may cause violations in the orthonality assumption if cut-off values (e.g., zero  
635 and the soil porosity as lower and upper physical limit, respectively) are applied to the soil  
636 moisture estimates of a particular data sets. Especially in very dry or very wet regimes, where  
637 random errors would often cause these thresholds to be exceeded, this can result in considerable  
638 biases in all (both relative and TCA-based) uncertainty metrics.

### 639 4.4.3 Bias estimation

640 Bias estimation is only meaningful against reference data at the satellite footprint scale, i.e.  
641 without considerable representativeness errors and if no rescaling has been applied (see Sec.  
642 4.3.2).

#### 643 Temporal mean bias

644 The term bias commonly refers to the (temporal) mean difference between two data sets (*En-*  
645 *tekhabi et al.*, 2010a):

$$b_{xy} = \bar{x} - \bar{y} = \alpha_x - \alpha_y + (\beta_x - \beta_y)\bar{t} \quad (4)$$

646 Typically,  $b_{xy}$  is considered to represent first-order (additive) biases only. However, as can be  
647 seen in Eq. (4), the mean difference is also sensitive to second-order (multiplicative) biases,  
648 amplified by the actual mean soil moisture content ( $\bar{t}$ ). When using non-fiducial reference data,  
649  $b_{xy}$  provides an indication of which data set is systematically drier or wetter than the other, but  
650 without relation to the assumed true grid cell average. Moreover, a positive difference in the  
651 mean ( $\alpha_x > \alpha_y$ ) and a negative difference in variability ( $\beta_x < \beta_y$ ) can cause the same sign in  
652  $b_{xy}$  as a negative mean difference and a positive variability difference. When calculated against  
653 fiducial reference data,  $b_{xy}$  collapses to  $\alpha_x + (\beta_x - 1)\bar{t}$ . That is, it is a direct estimate for biases  
654 in the satellite retrieval, yet it is still susceptible to both first and second-order biases, and  
655 influenced by the average soil moisture conditions.

#### 656 Second-order bias

657 Most validation studies do not attempt to estimate second-order biases and neglect their impact

658 on  $b_{xy}$  and other validation metrics such as the (unbiased) Root-Mean-Square-Difference (see  
659 *Gupta et al.* (2009) and Sec. 4.4.4). TCA potentially allows for the direct estimation of second-  
660 order biases (*Gruber et al.*, 2016a) as:

$$\beta_x^y = \frac{\sigma_{xz}}{\sigma_{yz}} = \frac{\beta_x \beta_z \sigma_t^2 + \sigma_{\xi_x, \xi_z}}{\beta_y \beta_z \sigma_t^2 + \sigma_{\xi_y, \xi_z}} \approx \frac{\beta_x}{\beta_y} \quad (5)$$

661 where  $\beta_x^y$  denotes the TCA-based second-order bias estimate of  $x$  relative to  $y$  which, if  $y$  is  
662 a fiducial reference data set and if no non-orthogonal or correlated random representativeness  
663 errors exist ( $\beta_y \approx 1, \sigma_{\xi_x, \xi_z} \approx 0, \sigma_{\xi_y, \xi_z} \approx 0$ ), provides a direct estimate of the second-order bias  
664  $\beta_x$ . Notice that neither first nor second-order biases in  $z$  influence  $\beta_x^y$ . Alternatively, Eq. (5)  
665 can also be used for rescaling purposes (*Yilmaz and Crow*, 2013; *Su et al.*, 2014; *Gruber et al.*,  
666 2016a, see Sec. 4.3.2).

#### 667 4.4.4 Uncertainty estimation

668 As discussed, uncertainty estimates aim at representing the pdf of the random errors (see Sec.  
669 2), which is typically done by means of their standard deviation (or variance).

#### 670 (Unbiased) Root-Mean-Square-Difference

671 The most common relative metric for estimating uncertainty is the Root-Mean-Square-Difference  
672 (RMSD; *Entekhabi et al.*, 2010a):

$$\begin{aligned} RMSD_{xy} &= \sqrt{(x - y)^2} = \sqrt{(\bar{x} - \bar{y})^2 + \sigma_x^2 + \sigma_y^2 - 2\sigma_{xy}} \\ &= \sqrt{(\alpha_x - \alpha_y + (\beta_x - \beta_y)\bar{t})^2 + (\beta_x - \beta_y)^2 \sigma_t^2 + \sigma_{\xi_x}^2 + \sigma_{\xi_y}^2 - 2\sigma_{\xi_x, \xi_y}} \end{aligned} \quad (6)$$

673 Since the RMSD is sensitive to both systematic and random errors, the bias component is  
674 - for uncertainty estimation purposes - typically removed, resulting in the unbiased RMSD  
675 (ubRMSD):

$$\begin{aligned} ubRMSD_{xy} &= \sqrt{RMSD^2 - b_{xy}^2} = \sqrt{\sigma_x^2 + \sigma_y^2 - 2\sigma_{xy}} \\ &= \sqrt{(\beta_x - \beta_y)^2 \sigma_t^2 + \sigma_{\xi_x}^2 + \sigma_{\xi_y}^2 - 2\sigma_{\xi_x, \xi_y}} \end{aligned} \quad (7)$$

676 The common definition of the ubRMSD specifically corrects for differences between the mean of  
677 the data sets (*Entekhabi et al.*, 2010a). However, as can be seen in Eq. (7), it remains susceptible  
678 to second-order biases, which are amplified by the actual soil moisture variability ( $\sigma_t^2$ ). Moreover,

679 as was the case for  $b_{xy}$ , this second-order bias dependency in  $ubRMSD_{xy}$  persists even when  
680 calculated against fiducial reference data, in which case Eq. (7) collapses to  $\sqrt{(\beta_x - 1)^2\sigma_t^2 + \sigma_{\xi_x}^2}$ .  
681 As discussed in Sec. 4.3.2, data sets are often rescaled before calculating validation metrics to  
682 account for systematic representativeness errors, especially when validating against data from  
683 sparse networks. This is most commonly done by matching the temporal mean and the standard  
684 deviation of the data sets, or their entire cdf (i.e. also higher statistical moments). However, as  
685 can be seen from Eq. (3), this only properly corrects for relative differences in  $\beta$  if the SNRs  
686 (including random representativeness errors) of the data sets are equal, which is very unlikely.  
687 Consequently, Eq. (7) will still contain the remaining difference between  $\beta_x$  and the rescaled  $\beta_y$ ,  
688 multiplied with the actual soil moisture variability, and also random representativeness errors.

### 689 (Unbiased) Root-Mean-Square-Error

690 As mentioned in the previous section, TCA potentially allows for the estimation of relative  
691 rescaling coefficients that are independent from the SNRs of the data sets (see Eq. (5)), which  
692 would allow to fully correct for the second-order bias component in Eq. (7). Moreover, TCA  
693 allows to more directly estimate the satellite uncertainty (i.e. its error standard deviation  $\sigma_{\xi_x}$ ,  
694 commonly referred to as unbiased Root-Mean-Square-Error; ubRMSE) as:

$$\begin{aligned}
ubRMSE_x &= \sqrt{\left| \frac{(x-y)(x-z)}{\sigma_{yz}} \right|} = \sqrt{\left| \sigma_x^2 - \frac{\sigma_{xy}\sigma_{xz}}{\sigma_{yz}} \right|} \\
&= \sqrt{\left| \beta_x^2\sigma_t^2 + \sigma_{\xi_x}^2 - \frac{(\beta_x\beta_y\sigma_t^2 + \sigma_{\xi_x,\xi_y})(\beta_x\beta_z\sigma_t^2 + \sigma_{\xi_x,\xi_z})}{\beta_y\beta_z\sigma_t^2 + \sigma_{\xi_y,\xi_z}} \right|} \approx \sigma_{\xi_x}
\end{aligned} \tag{8}$$

695 Note that when calculating the ubRMSE using the cross-multiplied differences instead of the  
696 statistical moments, the data sets  $y$  and  $z$  do have to be bias-corrected with respect to  $x$  a priori  
697 using Eqs. (4) and (5). The absolute value is taken to prevent negative signs in uncertainty  
698 estimates that could occur due to sampling errors (*Gruber et al., 2018*, see Sec. 4.5). As one  
699 can see,  $ubRMSE_x$  is (as opposed to  $ubRMSD_{xy}$  in Eq. (7)) fully unbiased in that it contains  
700 neither first nor second-order biases from both the satellite and the validation data sets, and it  
701 also no longer contains the uncertainties inherent in the reference data products (*Gruber et al.,*  
702 *2016a*). However, estimates that are unbiased *with respect to the assumed true grid cell average*  
703 can only be obtained if at least one fiducial reference data set is available (*Chen et al., 2017*).  
704 Moreover,  $ubRMSE_x$  is not affected by random representativeness errors in  $y$  and  $z$  as long as  
705 they are orthogonal and not correlated. Such representativeness error correlations could occur

706 for example when applying TCA to in situ measurements together with two coarse resolution  
707 products. This case, however, provides an opportunity to estimate the representativeness of in  
708 situ stations while uncertainty estimates for the coarse resolution products remain unaffected  
709 (*Miralles et al.*, 2010; *Gruber et al.*, 2013a; *Chen et al.*, 2017). For a more detailed derivation  
710 of how representativeness errors affect the TCA-based uncertainty estimates we refer the reader  
711 to *Vogelzang and Stoffelen* (2012) and *Gruber et al.* (2016a).

712 The above described metrics are direct estimators for data set uncertainty. However, for  
713 many applications, how “good” a data set is depends on how large its uncertainties are relative  
714 to the variability of the actual soil moisture signal. Simply put, the larger the soil moisture  
715 variations one strives to observe, the more easily they can be distinguished from noise in the  
716 measurements. Therefore, some metrics aim at estimating the SNR rather than the uncertainty  
717 alone, the most important ones for soil moisture validation being discussed below.

### 718 **Pearson correlation coefficient**

719 The most common SNR-related relative metric is the linear (Pearson) correlation coefficient,  
720 which is typically described as a measure for statistical dependency between two data sets.  
721 From the error model in Eq. (3) one can see that it is also a direct, normalized (between -1  
722 and 1) representation of the SNRs of the two data sets for which it is calculated (*Gruber et al.*,  
723 2016a):

$$\begin{aligned}
 R_{xy} &= \frac{\sigma_{ij}}{\sigma_i \sigma_j} = \frac{\beta_x \beta_y \sigma_t^2 + \sigma_{\xi_x, \xi_y}}{\sqrt{(\beta_x^2 \sigma_t^2 + \sigma_{\xi_x}^2)(\beta_y^2 \sigma_t^2 + \sigma_{\xi_y}^2)}} \\
 &\approx \text{sgn}(\sigma_{xy}) \frac{1}{\sqrt{(1 + SNR_x^{-1})(1 + SNR_y^{-1})}}
 \end{aligned}
 \tag{9}$$

724 with  $SNR_x = \frac{\beta_x^2 \sigma_t^2}{\sigma_{\xi_x}^2}$  and  $SNR_y = \frac{\beta_y^2 \sigma_t^2}{\sigma_{\xi_y}^2}$ .  $\text{sgn}(\cdot)$  denotes the signum function. When calculated  
725 against fiducial reference data,  $R_{xy}$  is a direct representation of the SNR of the satellite under  
726 validation (i.e.  $SNR_x$ ). Notice that the “signal” to which the “noise” in the SNR estimator is  
727 related is the true soil moisture variability scaled with the second-order satellite bias (i.e.  $\beta_x^2 \sigma_t^2$ ).  
728 Even if  $\beta_x$  could be estimated reliably, for example from Eq. (5), rescaling does not change  
729 the SNR as the uncertainty would be scaled as well. However, the ratio  $\frac{\beta_x^2 \sigma_t^2}{\sigma_{\xi_x}^2}$  is in fact the  
730 quantity of interest that determines how well signal variations can be distinguished from noise,  
731 regardless of whether systematic errors have been corrected for (*Gruber et al.*, 2016a), which can

732 be also interpreted as the (linear) correlation with the true soil moisture signal (*McColl et al.*,  
733 2014). When  $R_{xy}$  is calculated against non-fiducial reference data, it is additionally influenced  
734 by second-order systematic and random representativeness errors as well as the uncertainties of  
735 that reference data set.

### 736 TCA-based correlation coefficient

737 Influences of the reference data set can be again isolated using TCA (*McColl et al.*, 2014) by  
738 directly estimating  $R_x$  as:

$$\begin{aligned}
R_x &= \sqrt{\left| \frac{\sigma_{xy}\sigma_{xz}}{\sigma_x^2\sigma_{yz}} \right|} = \sqrt{\left| \frac{(\beta_x\beta_y\sigma_t^2 + \sigma_{\xi_x,\xi_y})(\beta_x\beta_z\sigma_t^2 + \sigma_{\xi_x,\xi_z})}{(\beta_x^2\sigma_t^2 + \sigma_{\xi_x}^2)(\beta_y\beta_z\sigma_t^2 + \sigma_{\xi_y,\xi_z})} \right|} \\
&\approx \sqrt{\left| \frac{\beta_x^2\sigma_t^2}{\beta_x^2\sigma_t^2 + \sigma_{\xi_x}^2} \right|} = \frac{1}{\sqrt{1 + SNR_x^{-1}}}
\end{aligned} \tag{10}$$

739 As was the case for the ubRMSE, the validity of Eq. (10) requires that there is no correlation or  
740 non-orthogonality between random representativeness errors, but their individual variance may  
741 well be non-zero. If these assumptions are respected, then  $R_x$  will be an unbiased representation  
742 of the correlation between  $x$  and the (unknown) hypothetical truth. Consequently,  $R_x$  will  
743 always be larger than  $R_{xy}$  although this difference decreases as the quality of the reference  $y$   
744 increases. Note, however, that  $R_x$  only ranges between 0 and 1, as an anti-correlation (with  
745 respect to the true signal) cannot be unambiguously inferred from the three covariances in Eq.  
746 (10).

### 747 (Logarithmic) Signal-to-Noise Ratio

748 Instead of expressing the SNR normalized between 0 and 1, it is often estimated directly and  
749 linearized by converting it into decibel (dB) units (*Gruber et al.*, 2016a):

$$SNR_x[dB] = -10 \log \left( \left| \frac{\sigma_x^2\sigma_{yz}}{\sigma_{xy}\sigma_{xz}} \right| - 1 \right) \approx 10 \log \left( \frac{\beta_x^2\sigma_t^2}{\sigma_{\xi_x}^2} \right) \tag{11}$$

750 This provides a more direct, linear representation of the ratio between soil moisture and uncer-  
751 tainty magnitude than  $R_x$ , yet the information content in both metrics is identical; it is simply  
752 a different way of presentation. Note that the  $SNR_x$  is already being used as a more coher-  
753 ent (than RMSD or RMSE based metrics) satellite data quality indicator for defining target  
754 accuracy requirements (see Sec. 4.8.2).

## 755 4.5 Statistical significance testing

756 All the above described (and also most other less common) validation metrics are based on  
757 statistical moments, sampled in time. Since these estimates are based on finite samples (i.e.  
758 the discrete soil moisture time series), they are subject to sampling errors. The most common  
759 way to deal with statistical uncertainty (i.e. sampling errors) across science communities is  
760 Null Hypothesis Significance Testing (NHST) using  $p$ -values and/or confidence intervals (*Wilks*,  
761 2011). In a validation context, typical hypotheses to be nullified are, for example, that a soil  
762 moisture product does not meet a target accuracy threshold or that one product does not  
763 exhibit higher correlation with a reference product than another. For testing such hypotheses,  
764 the sampling distribution of the statistical estimate under consideration (such as a validation  
765 metric) is constructed based on the magnitude of the estimate and the size of the sample used to  
766 draw this estimate (see below). Then, either the  $p$ -value is calculated, which is the probability  
767 of values of the sampling distribution to be equal to or below (or above, depending on which tail  
768 is considered) the pre-defined Null-value (representing the Null hypothesis), or the  $(1 - \alpha) \cdot 100\%$   
769 confidence interval is considered. A rejection of the Null-hypothesis is considered statistically  
770 significant, if the  $p$ -value is below a pre-defined significance level  $\alpha$  (typically 0.05) or if the  
771  $(1 - \alpha) \cdot 100\%$  confidence interval does not contain the Null-value. When comparing estimates  
772 of different samples (e.g., the performance of different soil moisture products), it is common to  
773 consider their relative difference as statistically significant if their confidence intervals do not  
774 overlap. Note that the term “Null-value” refers to the Null hypothesis and not to a value of zero  
775 of the test statistic (i.e. the validation metric). A common Null-value for testing soil moisture  
776 accuracy requirements, for instance, is  $0.04 \text{ m}^3\text{m}^{-3}$  ubRMSD (see Sec. 4.8.2). Hence, if the  
777  $p$ -value for  $0.04 \text{ m}^3\text{m}^{-3}$  of the sampling distribution around an estimated ubRMSD is below the  
778 defined  $\alpha$  level, the product is said to meet accuracy requirements with statistical significance.

779 However, the American Statistical Association (ASA) has recently issued a statement on sta-  
780 tistical significance and  $p$ -values (*Wasserstein and Lazar, 2016*) warning about the science-wide  
781 misuse and abuse of NHST through the replacement of scientific reasoning with a dichotomous  
782 and arbitrary classification of results into “significant” or “non-significant”. In this statement,  
783 the ASA is advocating the abandonment of statistical significance testing altogether for two main  
784 reasons. The first one is that an alarming fraction of articles in the scientific literature present  
785 unjustified inferences based on misinterpreted  $p$ -values and confidence intervals (*Greenland et al.*,  
786 2016; *Gelman and Stern, 2006*; *Wasserstein and Lazar, 2016*). For example, overlapping confi-

787 dence intervals are often wrongly considered to imply a non-significant difference between two  
788 estimates. The second and more important argument against significance testing is that  $p$ -values  
789 alone provide no grounds for meaningful decision making. While the magnitude of  $p$  itself may  
790 be informative about how consistent the data at hand are with an assumed stochastic model,  
791 “[...] a label of statistical significance does not mean or imply that an association or effect is  
792 highly probable, real, true, or important. Nor does a label of statistical nonsignificance lead to  
793 the association or effect being improbable, absent, false, or unimportant.” (*Wasserstein et al.*,  
794 2019). Therefore, no practical conclusion or decision should be based on whether  $p$ -values do or  
795 do not meet an arbitrarily defined threshold.

796 In a recent special issue of *The American Statistician* (*Wasserstein et al.*, 2019), the sta-  
797 tistical community is aiming to propose more appropriate alternatives. Their key message is  
798 that, naturally, there should not and cannot be a one-fits-all approach or threshold for statis-  
799 tical/scientific inference. Instead of strictly yet arbitrarily categorizing study results based on  
800 dichotomous significance tests, one should strive for more careful study design and more rigor-  
801 ous understanding, interpretation and reporting of the stochastic properties of the data at hand  
802 (*Greenland et al.*, 2016; *Tong*, 2019).

803 In conclusion, for soil moisture validation purposes, we follow the above guidance and  
804 recommend to avoid any statement or interpretation about statistical “significance” or “non-  
805 significance” and to instead always provide and interpret a statistical summary of calculated  
806 validation metrics in the form of confidence intervals alongside the metrics themselves. How  
807 confidence intervals can be calculated and recommendations of how they can be presented are  
808 described in the following sections.

## 809 4.6 Confidence intervals

810 In general, confidence intervals represent the pdf of the sampling errors of an estimate and  
811 are defined at a certain confidence level. A confidence level of, say, 95% means that if one  
812 would repeatedly calculate 95% confidence intervals in a series of similar experiments, then 95%  
813 of them would - on average - contain the true value, provided that all assumptions made for  
814 the stochastic model are met. Note that this is *not* the probability that the true value that  
815 is approximated by the estimate lies within the confidence interval (*Neyman*, 1937; *Greenland*  
816 *et al.*, 2016). In theory, this probability - which would indeed be more informative - could be  
817 represented by a Bayesian credible interval, but calculating it would require a priori knowledge



818 about the pdf of the parameter that is being estimated (i.e. the so-called “prior”) and this is  
 819 typically not available.

820 Estimating confidence intervals for validation metrics is not always straightforward because  
 821 the sampling error pdfs of the various estimators are often not well understood or contain  
 822 parameters that are typically unknown (*Zwieback et al., 2012*). The only validation metrics  
 823 (presented here) for which analytical solutions for confidence intervals exist are the temporal  
 824 mean bias ( $b_{xy}$ ), the unbiased RMSD ( $ubRMSD_{xy}$ ), and the Pearson correlation coefficient  
 825 ( $R_{xy}$ ). For TCA-based metrics, one has to rely on bootstrapping (*Efron and Tibshirani, 1986*)  
 826 to approximate the sampling error pdf.

#### 827 4.6.1 Analytical calculation

828 The sampling errors in  $b_{xy}$  and  $ubRMSD_{xy}$  are equivalent to the sampling errors of the popu-  
 829 lation mean and the population standard deviation of the difference series  $u = x - y$ , which are  
 830 known to follow a  $t$ -distribution and a  $\chi$ -distribution, respectively (*Gilleland, 2010; De Lannoy*  
 831 *and Reichle, 2016*):

$$\frac{\bar{u} - \mu_u}{\frac{s_u}{\sqrt{n}}} \sim t_{n-1} \quad (12)$$

832 and

$$\frac{\sqrt{n-1} s_u}{\sigma_u} \sim \chi_{n-1} \quad (13)$$

833 where  $n$  is the sample size;  $\bar{u}$  and  $s_u$  represent the sample mean and standard deviation of the  
 834 difference series ( $x - y$ ); and  $\mu_u$  and  $\sigma_u$  are their corresponding true population parameters.  
 835 The population moments of  $u$  are estimated within the  $(1 - \alpha) \cdot 100\%$  confidence intervals as  
 836 a function of the sample moments of  $u$ . Specifically, the confidence intervals ( $CI$ ) for  $b_{xy}$  and  
 837  $ubRMSD_{xy}$  can be inferred from Eqs. (12) and (13) as:

$$CI_{b_{xy}} = \left[ b_{xy} + t_{n-1}^{\alpha/2} \frac{ubRMSD_{xy}}{\sqrt{n}}, b_{xy} + t_{n-1}^{1-\alpha/2} \frac{ubRMSD_{xy}}{\sqrt{n}} \right] \quad (14)$$

838 and

$$CI_{ubRMSD_{xy}} = \left[ ubRMSD_{xy} \frac{\sqrt{n-1}}{\chi_{n-1}^{1-\alpha/2}}, ubRMSD_{xy} \frac{\sqrt{n-1}}{\chi_{n-1}^{\alpha/2}} \right] \quad (15)$$

839 No such simple direct relationships between the sampled and true values have yet been found  
840 for the other validation metrics presented here. For the Pearson correlation coefficient, it can be  
841 indirectly obtained through Fischer's  $z$ -transformation, which transforms  $R_{xy}$  into a variable that  
842 approximately follows a normal distribution with mean  $z_{xy}$  and standard deviation  $(n - 3)^{-0.5}$   
843 (*Bonett and Wright, 2000*):

$$z_{xy} = 0.5 \ln \left( \frac{1 + R_{xy}}{1 - R_{xy}} \right) \sim \mathcal{N}_{z_{xy}, (n-3)^{-0.5}} \quad (16)$$

844 The confidence interval for  $R_{xy}$  can be obtained by back-transforming  $z$  as:

$$CI_{R_{xy}} = \left[ \frac{e^{2z^{1-\alpha}} - 1}{e^{2z^{1-\alpha}} + 1}, \frac{e^{2z^\alpha} - 1}{e^{2z^\alpha} + 1} \right] \quad (17)$$

845 One major issue for calculating confidence intervals from the analytical expressions described  
846 above is the inherent assumption of independence between samples. For soil moisture time series,  
847 this assumption is often not met due to the auto-correlated nature of soil moisture governing  
848 processes. Since such auto-correlation in the data essentially causes a widening of the confidence  
849 intervals, one popular way to account for it is to reduce the degrees of freedom (sample size)  
850 of the used distribution. This is typically done by assuming a first-order auto-regressive AR(1)  
851 behaviour in the time series and using the lag-1 auto-correlation ( $\rho$ ) to calculate a correction  
852 factor for the sample size  $n$  (*Dawdy and Matalas, 1964; Draper et al., 2012*):

$$n_e = n \cdot \frac{1 - \rho}{1 + \rho} \quad (18)$$

853 where  $n_e$  is the effective sample size that is used to estimate auto-correlation corrected confidence  
854 intervals according to Eqs. (14)-(17). A combined effective value for  $\rho$ , which summarizes the  
855 possibly different lag-1 auto-correlation of the two considered time series for which the respective  
856 validation metric is calculated, can be obtained as their geometric average:

$$\rho = \sqrt{\rho_x \cdot \rho_y} \quad (19)$$

857 with  $\rho_x$  and  $\rho_y$  obtained from a fitted AR(1) model as:

$$\rho_i = e^{-\frac{dm}{\tau_i}} \quad (20)$$

858 where  $i \in [x, y]$ ,  $\tau_i$  is the fitted persistence time of the individual time series  $x$  and  $y$ , and  $d_m$   
859 is the the median distance between consecutive valid, collocated observations, i.e. the lag-1  
860 distance accounting for the typically irregular spacing between satellite measurements. Note  
861 that averaging correlation coefficients is generally not recommended (see Sec. 4.7), but required  
862 here to determine a single effective proxy of the auto-correlation of collocated data pairs with  
863 possibly deviating individual memory. Using the geometric average avoids the dominance of  
864 data sets with large auto-correlation (e.g., land surface models often have a different memory  
865 than satellite observations), which may cause excessively large confidence intervals.

866 Note that the necessity of relying on a possibly crude approximation of a lumped effective  
867 auto-correlation correction parameter for calculating confidence intervals is but one factor under-  
868 mining their ability to serve as decision basis for declaring results as significant or non-significant  
869 (see the previous section). One should always bear in mind that confidence intervals inevitably  
870 are - just as the estimates they are meant to describe - uncertain.

#### 871 **4.6.2 Bootstrapping**

872 No exact solvable analytical expressions or transformations for confidence intervals around TCA-  
873 based metrics have yet been derived. *Zwieback et al.* (2012) presented a formulation of confidence  
874 intervals for TCA-based RMSE estimates in a synthetic study which, however, required the  
875 knowledge of the true RMSE states and is therefore of limited practical use. Alternatively, several  
876 studies (e.g., *Caires and Sterl*, 2003; *Zwieback et al.*, 2012; *Draper et al.*, 2013) have suggested  
877 the use of bootstrapping as a potential non-parametric method for obtaining confidence intervals  
878 of estimators with unknown sampling distribution (*Efron and Tibshirani*, 1986).

879 Bootstrapping is a special case of Monte Carlo simulation, which uses the sample itself  
880 as approximation of the population. More specifically, it constructs an empirical probability  
881 distribution of the test statistic (in our case the validation metric) by resampling the original  
882 sample multiple times, with replacement to preserve the sample size, and repeated calculation  
883 of the test statistic from those resamples. This bootstrapped distribution then allows for the  
884 direct derivation of confidence intervals as well as other parameters of the sampling error pdf.  
885 The advantages of this method lie in its algorithmic simplicity and that it can be applied  
886 to any metric without the need to assume a particular sampling distribution (such as  $t$  or  
887  $\chi$ ). However, bootstrapping confidence intervals requires a considerable number of resamples,  
888 which may lead to large computational costs, and relies on the assumption that the sample is

indeed a reliable representation of the population, which requires a large sample size. A general recommendation for bootstrapping confidence intervals is to use a minimum of 1000 resamples (Efron and Tibshirani, 1986). However, the number of required resamples may be chosen more specifically for a given study by testing for convergence of the results with increasing sample size. For example, Draper et al. (2013) used 1000 resamples for estimating confidence intervals for TCA-based *ubRMSE* estimates, although their testing found that 500 would have been sufficient.

As was the case for the analytical expressions, bootstrapped confidence intervals are also susceptible to auto-correlation in the data. This can be accounted for by resampling blocks of data instead of single data points, referred to as block-bootstrapping (Ólafsdóttir and Mudelsee, 2014), which preserves the auto-correlation properties of the original sample. An estimate of the optimal block length ( $l_{opt}$ ) for bootstrapping CIs around TCA-based estimates can be obtained following Chen et al. (2018) as:

$$l_{opt} = \text{NINT} \left\{ \sqrt[3]{\left( \frac{\sqrt{6 \cdot n \cdot \rho}}{1 - \rho^2} \right)^2} \right\} \quad (21)$$

where  $\text{NINT}\{\cdot\}$  denotes rounding to the nearest integer. As before, a single effective value for  $\rho$  can be obtained as the geometric average of the lag-1 auto-correlations of the three data sets used to obtain the respective TCA estimate ( $\rho = \sqrt[3]{\rho_x \cdot \rho_y \cdot \rho_z}$ ). The lag-1 is the median time interval between consecutive valid, collocated data triplets. To prevent data gaps from causing an auto-correlation degradation during the resampling, we recommend to discard data blocks from the resamples if they contain less than 50% of valid data.

## 4.7 Summary statistics

Validation metrics and their confidence intervals should be calculated and assessed over a wide range of spatial locations to understand error characteristics of a soil moisture product under different climatic, topographic and land cover conditions. However, it may be practical to summarize spatially distributed skill estimates into a single combined metric (for example to obtain an overall ranking of different products or to track the performance evolution of a product over time), which requires also the aggregation of their associated confidence intervals.

915 **4.7.1 Averaging metrics**

916 The most common way of obtaining a combined skill estimate is arithmetic averaging:

$$\bar{\nu} = \mathbf{w}^\top \mathbf{v} \tag{22}$$

917 where  $\bar{\nu}$  is the average of  $k$  spatially distributed skill metrics that are summarized in the skill  
 918 vector  $\mathbf{v} = [\nu_1 \cdots \nu_k]^\top$ ; and  $\mathbf{w} = [w_1 \cdots w_k]^\top$  contains the weights that are attributed to the  
 919 individual skill estimates with  $\sum w_i = 1$ . Averaging skill metrics in a weighted fashion to  
 920 minimize the impact of sampling errors is in principle possible by deriving weights from the  
 921 sampling error magnitudes (*Aitkin, 1936*), but in most cases, an unweighted average is preferred  
 922 because validation points are typically selected to represent a wide range of varying conditions,  
 923 and areas with lower sampling errors (i.e. regions with better temporal coverage, for instance  
 924 because less data are masked out) could dominate a weighted averaged skill estimate. For such  
 925 unweighted average, the weight vector takes the form  $\mathbf{w} = [k^{-1} \cdots k^{-1}]^\top$ .

926 While many metrics can be averaged safely, it is - against common practice - not recom-  
 927 mended to average correlation coefficients (neither Pearson nor TCA-based) because they are  
 928 calculated as ratios using standard deviations (variances) and covariances or SNRs (see Eqs. (9)  
 929 and (10)). Therefore, they behave highly non-linearly and neither an average of these ratios nor  
 930 a ratio of averaged nominators / denominators would allow for a meaningful inference about  
 931 statistical properties. For example, averaging correlation coefficients of 0.1 and 0.9, which cor-  
 932 respond to a SNR of 0.01 and 4.26, respectively (in the case of Pearson correlation assuming  
 933 a random error-free reference data set), would lead to an average correlation of 0.5 with an  
 934 associated SNR of 0.33. This is far from their average SNR of 2.14 (ignoring for the moment  
 935 that this too is an average of ratios) which would correspond to a correlation coefficient of 0.83.  
 936 In contrast, correlation coefficients of 0.3 and 0.7, representing SNRs of 0.1 and 0.96, respec-  
 937 tively, would have the same average correlation yet the average of their associated SNR is 0.53,  
 938 corresponding to a correlation of 0.59. Moreover, the skewed probability distribution of the  
 939 Pearson correlation coefficient causes the arithmetic average to be systematically biased. Some  
 940 studies suggest to average Fisher-transformed  $z$ -values instead (*Corey et al., 1998*), which have  
 941 a Gaussian sampling distribution, but a back-transformed  $z$ -average is just as difficult to inter-  
 942 pret. Following the above example, averaging correlation coefficients of 0.1 and 0.9 in  $z$ -space  
 943 would lead to an average correlation (or more precisely, an inverse average- $z$ ) of 0.66 (SNR =

944 0.76), whereas when averaging  $z$ -transformed correlations of 0.3 and 0.7, it would be 0.53 (SNR  
 945 = 0.39).

946 In other words, the choice of whether to average correlation coefficients, Fisher-transformed  
 947  $z$ -values, or SNRs - albeit representing the exact same uncertainty properties - will lead to  
 948 different values / interpretations of the resulting average and this difference also depends on the  
 949 degree of variability across the estimates that are being averaged. Moreover, the resulting average  
 950 number (regardless of the approach) no longer represents an actually meaningful statistical  
 951 property. Alternatively, instead of averaging pre-calculated correlation coefficients, one may be  
 952 tempted to calculate the correlation coefficient directly over the concatenated measurements of  
 953 all available locations to obtain an overall skill estimate. However, this is not meaningful as  
 954 the effects of different populations are lumped together. As a consequence, for example, two  
 955 data sets that individually exhibit strong positive correlation in a wet and in a dry soil moisture  
 956 regime, respectively, may appear to have an overall weak anti-correlation when put together, an  
 957 effect also known as Simpson's paradox (*Blyth, 1972*). Therefore, such an approach should be  
 958 strictly avoided.

#### 959 4.7.2 Averaging confidence intervals

960 The uncertainty in the spatially averaged skill metric in Eq. (22) associated with the *sampling*  
 961 errors of the individual skill estimates can be calculated through the standard method for the  
 962 propagation of uncertainty as:

$$s_{\bar{\nu}}^2 = \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \quad (23)$$

963 where  $s_{\bar{\nu}}^2$  is the sampling uncertainty in the averaged skill  $\bar{\nu}$  (i.e. its sampling error variance);  
 964 and  $\boldsymbol{\Sigma}$  is the sampling error covariance matrix for the  $k$  individual skill estimates. The corre-  
 965 sponding aggregated confidence intervals can be derived from a Gaussian distribution (which  
 966 will generally be assured by the Central Limit Theorem for reasonably large samples) with mean  
 967  $\bar{\nu}$  and standard deviation  $s_{\bar{\nu}}$ .

968 Diagonal elements in  $\boldsymbol{\Sigma}$  are the sampling error variances of the individual skill estimates, i.e.  
 969  $\text{diag}(\boldsymbol{\Sigma}) = \mathbf{s}^2$  with  $\mathbf{s}^2 = [s_{\nu_1}^2 \cdots s_{\nu_k}^2]^\top$ . For  $b_{xy}$  and  $ubRMSE_{xy}$  estimates, they are the squared  
 970 standard errors of the sample mean and sample variance (of the difference series  $u = x - y$  at

971 each individual location), respectively:

$$\begin{aligned}
 s_{b_{xy}}^2 &= \frac{ubRMSD_{xy}^2}{n} \\
 s_{ubRMSD_{xy}}^2 &= \frac{ubRMSD_{xy}^2}{2(n-1)}
 \end{aligned}
 \tag{24}$$

972 For TCA-based metrics, the sampling error variance can be directly calculated from the boot-  
 973 strapped sampling distribution.

974 In an ideal case, the reference data used for calculating skill metrics span a wide range of  
 975 varying conditions with samples that are independent of each other in time and space. In this  
 976 case, off-diagonal elements in  $\Sigma$  would be zero. However, in many cases, differences between  
 977 soil moisture time series contain sample auto-correlation due to the large-scale auto-correlated  
 978 nature of soil moisture (*Vachaud et al.*, 1985; *Crow et al.*, 2012) and because estimation errors  
 979 (those of satellite retrievals as well as those of in situ measurements or of model predictions)  
 980 may be correlated over large distances (*Gruber et al.*, 2015, 2018). Ignoring such sampling  
 981 error auto-correlation can lead to considerably underestimated confidence intervals of spatially  
 982 averaged skill estimates. Hence, calculating off-diagonal elements in  $\Sigma$ , which represent sampling  
 983 error covariances between the skill estimates of different locations, is critical. Although these  
 984 covariances cannot be estimated directly, they can be derived from the sample auto-correlation  
 985 matrix  $\mathbf{R}$  and the sampling error *standard deviations*  $\mathbf{s}$  (see above) as:

$$\Sigma = \mathbf{R} \circ \mathbf{s}\mathbf{s}^\top
 \tag{25}$$

986 where  $\circ$  denotes the Hadamard product, i.e. element-wise matrix multiplication.  $\mathbf{R}$  differs for  
 987 the various skill metrics. For  $b_{xy}$  and  $ubRMSD_{xy}$ , it is the *spatial* auto-correlation matrix of the  
 988 difference series  $u$ , and of the squared, bias-corrected difference series  $(u - \bar{u})^2$ , respectively, at  
 989 the different locations  $u$  where skill metrics are calculated. For TCA-based metrics, the sampling  
 990 error covariance can be calculated as the covariance between the bootstrapped samples (*Gruber*  
 991 *et al.*, 2019b), provided that the order in which bootstrap-resamples are drawn is the same at  
 992 all different locations, which may be difficult when using block-bootstraps with different block-  
 993 length.

994 Earlier research (*De Lannoy and Reichle*, 2016) has proposed a clustering approach to take  
 995 possible sampling error correlations into account. This approach first calculates mean metrics  
 996 and confidence intervals per spatial cluster, assuming that the sampling errors of the spatially

997 close data sets within each cluster are perfectly correlated. Next, averaged skill metrics and con-  
998 fidence intervals from within the clusters are averaged, assuming that all clusters are completely  
999 independent. However, this approach is expected to overestimate confidence intervals because:  
1000 (i) sampling errors will never be perfectly correlated unless validation metrics are calculated  
1001 multiple times from the exact same data, and (ii) clusters are formed based on the expected  
1002 auto-correlation length of the soil moisture data sets, which will be much larger than that of the  
1003 difference series between data sets, as required in Eq. (25).

1004 Finally, although averaging of some metrics and confidence intervals is possible, we generally  
1005 recommend to retain detailed information about their spatial variability, and to leverage this  
1006 information to obtain a better understanding of product performance and its relation to land  
1007 cover, topography, climate, and other possibly important factors. If point-wise assessments are  
1008 not feasible or if simple product summaries are desired, percentile statistics such as medians  
1009 and inter-quartile-ranges (of both calculated skill estimates and their confidence intervals) are  
1010 generally more informative than spatial averages and their increasingly inaccurate averaged  
1011 confidence intervals. More specific recommendations of how validation metrics and confidence  
1012 intervals can be presented are provided in Sec. 5 and Appendix A.

## 1013 **4.8 Practical remarks**

### 1014 **4.8.1 Validating downscaled products**

1015 Currently, most space-borne microwave sensors available for soil moisture retrieval operate at  
1016 spatial resolutions of about  $25^2 - 50^2$  km<sup>2</sup> (*Gruber et al.*, 2019a). Some higher-resolution Syn-  
1017 thetic Aperture Radar (SAR) sensors exist that allow for reasonable soil moisture retrieval at  
1018 scales up to approximately 1 km<sup>2</sup> (*Pathe et al.*, 2009; *Gruber et al.*, 2013b), yet with consider-  
1019 ably lower temporal resolution and accuracy. In addition, many downscaling approaches have  
1020 been developed to improve the spatial resolution of coarse-resolution soil moisture products,  
1021 e.g., by fusing coarse-resolution radiometer or scatterometer measurements with high-resolution  
1022 SAR data (*Das et al.*, 2017; *Bauer-Marschallinger et al.*, 2018), by fusing microwave observa-  
1023 tions with optical/thermal measurements (*Chauhan et al.*, 2003), or through data assimilation  
1024 (*Reichle et al.*, 2017c). For a comprehensive review of downscaling methods see *Peng et al.*  
1025 (2017).

1026 The validation of downscaled products is mostly done as for coarse-resolution products, i.e.  
1027 through time series analysis with a focus on temporal dynamics at individual locations (see



1028 Sec. 4). In doing so, it has been shown that the downscaling process often actually decreases  
1029 the temporal performance of the products, that is, the original coarse-resolution products often  
1030 correlate better with local soil moisture dynamics, even at a point scale, than their downscaled  
1031 counterparts (*Peng et al.*, 2015). While downscaled soil moisture images provide more visual  
1032 level-of-detail, only few studies have quantitatively assessed whether the obtained spatial pat-  
1033 terns actually represent real soil moisture variations (e.g., *Bauer-Marschallinger et al.*, 2018;  
1034 *Sabaghy et al.*, in review) or whether they are just mimicking spatial patterns of ancillary data  
1035 such as soil texture maps (for a comprehensive review of validation studies for downscaled prod-  
1036 ucts see *Peng et al.*, 2017).

1037 Therefore, we highly recommend that future validation studies for downscaled products  
1038 put a strong emphasis on assessing also the spatial soil moisture variations obtained from the  
1039 downscaling, e.g., by estimating spatial correlation coefficients (*Sahoo et al.*, 2013; *Kolassa*  
1040 *et al.*, 2017; *Sabaghy et al.*, in review), in addition to time series analyses. To that end, we  
1041 further encourage the setup of field campaigns / validation sites dedicated to support such  
1042 high-resolution validation activities, especially in regions where soil moisture variations are very  
1043 heterogeneous.

#### 1044 **4.8.2 Target accuracy requirements**

1045 Satellite soil moisture validation studies most commonly evaluate products against a target ac-  
1046 curacy threshold of  $0.04 \text{ m}^3\text{m}^{-3}$  ubRMSD across the globe, excluding regions of snow and ice,  
1047 frozen ground, complex topography, open water, urban areas, and vegetation with water content  
1048 greater than  $5 \text{ kg/m}^2$ . This requirement was defined by the Soil Moisture and Ocean Salinity  
1049 (SMOS; *Kerr et al.*, 2001) and the Soil Moisture Active Passive (SMAP; *Entekhabi et al.*,  
1050 2010a) missions, and by the Terrestrial Observation Panel for Climate (TOPC; *WMO*, 2016).  
1051 Alternatively, the Satellite Application Facility in Support to Operational Hydrology and Wa-  
1052 ter Management (H SAF) of the European Organisation for the Exploitation of Meteorological  
1053 Satellites (EUMETSAT) has defined (TCA-based) SNR product requirements (*H-SAF*, 2017)  
1054 for the operational soil moisture products that are retrieved from measurements of the Advanced  
1055 Scatterometer (ASCAT) onboard the MetOp satellites (*Naeimi et al.*, 2009). In particular, the  
1056 EUMETSAT H SAF defines 0, 3 and 6 dB SNR as threshold, target and optimal SNR require-  
1057 ments to make product assessment possible on a larger scale and spatially better comparable  
1058 (see Sec. 4.4).

1059 Both of these requirements are based on relatively practical, easy-to-estimate single numbers  
1060 that represent a rough average of what is currently achievable rather than being an indication  
1061 of “good” or “bad” product quality. While they provide easy means to monitor product per-  
1062 formance evolution over time and to compare products, they are not necessarily related to the  
1063 suitability of a product for specific applications. However, the actual specification of bias and  
1064 uncertainty requirements for the fitness-for-purpose for a particular application (including the  
1065 specification of the appropriate metrics) is a task of the respective user community and needs  
1066 further research (*Entekhabi et al.*, 2010b).

### 1067 **4.8.3 Reproducibility**

1068 The research community generally suffers from a lack of reproducibility in scientific studies  
1069 (*Baker*, 2016). Also in soil moisture validation studies, contradictory results for the performance  
1070 and relative ranking between different satellite products have been reported (e.g., *Wagner et al.*,  
1071 2014). These ambiguities originate from: (i) the choice of reference data and product versions;  
1072 (ii) the use of different spatial regions and time periods; (iii) different approaches used for data  
1073 preparation and pre-processing; (iv) statistical sampling errors; and (v) software implementation  
1074 errors. Note, however, that contradicting results are not necessarily caused by bad study design  
1075 but often originate from stochastic uncertainties, which are inevitably dominant in space borne  
1076 Earth observation measurements (*Greenland et al.*, 2016).

1077 Embracing statistical uncertainty and developing an in-depth understanding of soil moisture  
1078 product quality requires more comprehensive descriptions of data sets, software, and methodol-  
1079 ogy than are usually provided as well as the mandatory, additional estimation and presentation  
1080 of sampling errors. To that end, we recommend that:

- 1081 • all validation results should be accompanied by confidence intervals as measure for sam-  
1082 pling errors;
- 1083 • all methodological steps should be described with sufficient detail to be reproducible;
- 1084 • all data sets used for the study should be made publicly available and unambiguously  
1085 identifiable by providing their exact product version information and, where available,  
1086 their Digital Object Identifier (DOI);
- 1087 • all used software packages that are relevant for the exact reproduction of validation re-  
1088 sults should be referenced with their complete version number and, where available, their

1089 DOI. If not accessible via open repositories (in particular software specifically designed for  
1090 that study), we recommend to make source code publicly available, preferably on GitHub  
1091 (<https://github.com/>; last access: 1 July 2019).

1092 A list of some current publicly available software that is specifically aimed at, or closely related to  
1093 soil moisture validation is provided in Table 3. An online validation tool that is built around these  
1094 software packages and follows the good practice guidelines presented in this paper is provided by  
1095 the Quality Assurance Framework for Soil Moisture (QA4SM; <https://qa4sm.eodc.eu/>; last  
1096 access: 1 July 2019).

1097 Note that the re-distribution of in situ measurements (see the third point above) may be  
1098 particularly problematic as many networks do not operate for free. Requiring networks to  
1099 freely distribute their data will likely decrease the number of datasets available for validation  
1100 activities, which may ultimately hamper the evolution of satellite soil moisture products and  
1101 downstream products derived thereof. We therefore emphasize the tremendous value of ground  
1102 reference measurements and encourage the community to support, by any means possible, the  
1103 development and continuation of operational Cal/Val sites.

## 1104 **5 Validation Good Practice Protocol**

1105 This section provides a compilation of the theoretical considerations presented above in the form  
1106 of a validation good practice protocol for satellite soil moisture products, i.e. guidelines for:

- 1107 • the selection of reference data;
- 1108 • data pre-processing steps;
- 1109 • the selection and implementation of appropriate metrics;
- 1110 • the presentation of validation results.

1111 Figure 3 illustrates the process and Appendix A provides an example that follows these recom-  
1112 mendations. We stress that there is no one-size-fits-all approach for validating Earth observation  
1113 data. Depending on the application in question, several analyses may not be necessary. Also,  
1114 recommended thresholds may need to be adjusted depending on data quality requirements (e.g.,  
1115 more strict data masking procedures may be employed) or data availability (e.g., the allowed in  
1116 situ measurement depth may be increased if only retrievals from long wavelengths in dry and  
1117 sandy regions are used).

## 1118 **5.1 Data selection**

1119 As discussed in Sec. 3, no reference data source provides a sufficiently accurate and traceable  
1120 soil moisture proxy for reliable error assessment on a global scale. A complete and comprehen-  
1121 sive product validation therefore requires comparisons against each of the following: (i) dense  
1122 networks, in particular core validation sites; (ii) sparse networks; (iii) land surface model out-  
1123 put; and (iv) other satellite products, always making sure that the latest or most recommended  
1124 product versions are used. However, given the large number of satellite and reference prod-  
1125 ucts available, a complete analysis that considers all these data sources is typically beyond the  
1126 capacity of a single validation study. Therefore, separate studies may be conducted for dense  
1127 network validation (*Colliander et al.*, 2017), sparse network validation (*Dorigo et al.*, 2015),  
1128 or coarse-resolution product inter-comparison (*Al-Yaari et al.*, 2014) and their results compiled  
1129 together in a meta-analysis.

1130 Since satellite soil moisture retrievals represent only the top few centimeters of the soil, in  
1131 situ sensors and modelled soil layers used for validation should reach no deeper than 5-10 cm,  
1132 which is considered as the maximum sensing depth for currently available microwave wavelengths  
1133 (X-band to L-band). Information where currently publicly available reference data sets can be  
1134 accessed is provided in Table 2.

## 1135 **5.2 Pre-processing**

### 1136 **5.2.1 Masking**

1137 In situ measurements and satellite retrievals are typically accompanied by quality flags, which  
1138 must be used to mask out all measurements that are considered unreliable. If this masking  
1139 requires the decision of a threshold (for example the probability of RFI occurrence), recom-  
1140 mendations from data providers should be followed and the employed thresholds carefully docu-  
1141 mented. Ancillary information on dynamic geophysical variables, such as snow, temperature or  
1142 vegetation, should be used to screen microwave-based satellite measurements since no reliable  
1143 soil moisture retrieval is possible under frozen or snow-covered conditions and the quality of  
1144 soil moisture retrievals depends on the vegetation density. Such ancillary data can be supplied  
1145 by land surface models or complementary satellite data. Specifically, we recommend masking  
1146 out pixels classified as tropical forests, water bodies, wetlands, and inundation areas as well  
1147 as all measurements on days with non-zero snow indicators (e.g., snow height or snow-water-

1148 equivalent), or surface or soil temperature below 4°C. When biases or uncertainties of multiple  
1149 products are compared, they should be calculated from the exact same, collocated data points.  
1150 However, care should be taken that single products with poor data coverage do not distort the  
1151 overall assessment (see Sec. 6).

1152 To avoid excessively large confidence intervals that can hamper meaningful data comparison,  
1153 grid cells with less than 50-500 collocated data points may be masked out depending on data  
1154 availability (*Zwieback et al.*, 2012). Also, many studies mask out correlation coefficients based  
1155 on Student’s t-test (i.e. applying p-value thresholds for correlation coefficients), and/or bias and  
1156 uncertainty estimates based on vegetation density (e.g., vegetation water content > 5 kg/m<sup>2</sup>)  
1157 or other thresholds (e.g., open-water fraction > 0.05) (*Dorigo et al.*, 2010; *Brocca et al.*, 2011;  
1158 *Al-Yaari et al.*, 2014). However, carefully reporting and interpreting confidence intervals and  
1159 sample sizes at locations with low data coverage could indeed provide valuable additional insight  
1160 and may be more informative than masking out estimates completely (*Wasserstein et al.*, 2019).  
1161 Also, complete reporting of results prevents generating publication biases due to “cherry-picking”  
1162 which is sometimes found in the scientific literature (*Greenland et al.*, 2016).

## 1163 5.2.2 Collocation

1164 Spatial collocation requires the selection of a spatial comparison grid, which is often the grid  
1165 of the satellite product under validation. In situ measurements should be assigned to the grid  
1166 cell in which they are located. For dense networks, all stations that lie within a particular grid  
1167 cell should be averaged, if possible taking their respective spatial representativeness for that  
1168 grid cell into account. To avoid artificial jumps due to sensor drop-outs, only time steps where  
1169 all stations provide valid measurements should be considered. For the SMAP core validation  
1170 sites (see Sec. 3.2.1), a validation grid that minimizes upscaling errors has been developed as  
1171 described in *Colliander et al.* (2017).

1172 Gridded reference products (i.e. other satellite and land surface model products) should be  
1173 resampled onto the chosen comparison grid, e.g., using a Nearest Neighbor (NN) search. If the  
1174 grid resolution of the reference product is coarser than that of the comparison grid, individual  
1175 grid cells of that product may be assigned to multiple comparison grid cells. If the grid resolution  
1176 is much finer, all NNs of single comparison grid cells (in case more than one exist) should be  
1177 averaged, if possible taking spatial representativeness into account.

1178 Temporal collocation at comparison time steps should minimize the time difference between

1179 data match-ups and be based on a NN-search with a maximum time difference threshold of  
1180 1-12 hours, depending on data availability. Note that the choice of the comparison grid and  
1181 time steps may affect the presence and distribution of (spatial and temporal) representativeness  
1182 errors among the considered data sets (see Sec. 6).

### 1183 **5.2.3 Decomposition**

1184 All validation metrics should be calculated for the raw soil moisture time series (of collocated  
1185 retrievals and reference data) as well as for short-term and long-term anomalies, except for  
1186 temporal mean biases whose calculation is trivial for anomalies. Short-term anomalies should  
1187 be estimated as residuals from a seasonality that is computed by applying a 4-8 week moving  
1188 average window to the time series. Long-term anomalies should be estimated as residuals from  
1189 a climatology that is computed by averaging the measurements of all years within a 4-8 week  
1190 moving window around each DOY, but only if at least 5-10 years of data are available. To avoid  
1191 data-density related artefacts, especially in the transition periods from frozen to non-frozen  
1192 periods, moving averages should only be calculated if at least 25-50% of the maximal data pair  
1193 coverage is available within a particular time window.

### 1194 **5.2.4 Rescaling**

1195 When using fiducial reference data, units (e.g.,  $m^3m^{-3}$  and degree of saturation) should be  
1196 unified for the purpose of bias estimation using soil texture information, keeping in mind that  
1197 inaccuracy in soil information directly propagates into the bias estimates. To account for (hor-  
1198 izontal and vertical) systematic representativeness errors and different soil moisture units, the  
1199 data set under validation should be rescaled (before decomposition for validating raw time  
1200 series and after decomposition for validating anomalies) towards the reference data when esti-  
1201 mating absolute uncertainties (i.e. ubRMSDs or ubRMSEs). When calculating relative metrics,  
1202 data sets should be rescaled by matching their temporal mean and standard deviation. When  
1203 calculating TCA-based metrics, data sets should be rescaled using also TCA-based rescaling  
1204 coefficients. Note that no rescaling or unit conversion is necessary for Pearson correlation co-  
1205 efficients or TCA-based correlation and SNR estimates, since these metrics are not affected by  
1206 linear data transformation.

### 1207 **5.3 Metric calculation**

1208 Remember that all covariance-based metrics require zero error correlation. Any combination of  
1209 in situ measurements, land surface model estimates, active-microwave-based measurements, or  
1210 passive-microwave-based measurements is expected to mostly fulfil this requirement (see Sec.  
1211 4.4.2; *Gruber et al.*, 2016a). Different products from within any of these categories (except for  
1212 in situ data), on the other hand, are expected to have correlated errors (*Gruber et al.*, 2016b).  
1213 Therefore, the metrics described below should not be applied to such product combinations.  
1214 Moreover, since non-zero error correlations may exist even when using products from different  
1215 categories (see Sec. 4.4.2; *Yilmaz and Crow*, 2014; *Pan et al.*, 2015), it is strongly recommended  
1216 to verify if assumptions are met (see Sec. 5.3.2).

#### 1217 **5.3.1 Relative metrics**

1218 Temporal mean biases (Eq. (4)) should be calculated between all data sets that are expected  
1219 to be properly collocated and have comparable spatial resolution, and are hence not dominated  
1220 by spatial representativeness errors. These data sets may include dense networks, land surface  
1221 models, and other satellite data sets. It should be kept in mind, however, that the underlying  
1222 measurement resolution often considerably differs from the sampling grid resolution, which  
1223 potentially causes representativeness errors that are not directly apparent as such. Correlation  
1224 coefficients and unbiased Root-Mean-Square-Differences (Eqs. (9) and (7), respectively) should  
1225 be calculated between all data sets whose errors are not expected to be correlated (see above).

#### 1226 **5.3.2 TCA-based metrics**

1227 Second-order biases (Eq. (5)) of the validation data set should be calculated using fiducial  
1228 reference data (i.e. at the core validation sites). Unbiased Root-Mean-Square-Errors and SNRs  
1229 (Eqs. (8) and (11), respectively) should be calculated for all data sets. If more than one triplet  
1230 with independent errors is available to estimate the bias or uncertainty of a particular product,  
1231 TCA should be applied to all possible triplets and redundant estimates should be averaged  
1232 (*Gruber et al.*, 2016b). The spread between redundant estimates should be used as a diagnostic  
1233 to verify if orthogonality and zero error correlation assumptions are met (*Dorigo et al.*, 2010;  
1234 *Draper et al.*, 2013; *Chen et al.*, 2017).

### 1235 **5.3.3 Confidence intervals**

1236 For each metric, 80-95% confidence intervals should be calculated using their analytical esti-  
1237 mators (Eqs. (14)-(17)) or, if not available, block-bootstrapping. The latter should be based  
1238 on at least 1000 bootstrap samples (*Efron and Tibshirani*, 1986) or possibly less if tested for  
1239 convergence, and all confidence intervals should be corrected for sample auto-correlation.

## 1240 **5.4 Presentation**

1241 Validation metrics together with sample size and upper and lower confidence intervals/limits  
1242 should be presented for each location where they are calculated, either by means of spatial  
1243 maps or, if not meaningful (for example for core validation sites), in tabular form. Additionally,  
1244 summary statistics (representing average conditions and spatial variability) of both validation  
1245 metrics and their confidence intervals/limits should be provided, e.g., in the form of boxplots  
1246 (i.e. median, inter-quartile-range and 5th/95th percentiles). The presentation can be further  
1247 customized, for example by stratifying the summary statistics for climatological or land surface  
1248 conditions.

1249 Ratio-based metrics (i.e. Pearson and TCA-based correlation coefficients as well as SNRs)  
1250 must not be averaged. Differences between these metrics must always be related to their absolute  
1251 values and be interpreted with care (see Sec. 4.7). SNR-related properties of different products  
1252 may be compared in terms of SNR ratios or SNR differences in decibel space (Eq. (11)).

1253 Examples of how validation metrics and associated confidence intervals can be presented are  
1254 provided in Appendix A.

## 1255 **6 Towards best practices: discussion and conclusions**

1256 In this paper we have reviewed state-of-the-art validation methods, including reference data  
1257 sources and data pre-processing procedures, and provided community-agreed good practice  
1258 guidelines for the validation of satellite soil moisture products. Moreover, we have identified  
1259 several weak links that require careful attention to increase the reliability of soil moisture data  
1260 quality assessments. Specifically, the following research gaps should be addressed in the near  
1261 future:

- 1262 • On assumptions: the majority of studies assume that estimated biases and uncertainties  
1263 are stationary (i.e. constant over time) or at least that they represent the average data



1264 quality of a product. However, given the strong link between soil moisture data quality  
1265 and vegetation (*van der Schalie et al., 2018; Zwieback et al., 2018; Gruber et al., 2019a*),  
1266 retrieval accuracy can be expected to vary strongly between seasons and many applications  
1267 could greatly benefit from temporally varying quality information. Given the rapidly  
1268 growing temporal coverage of soil moisture products, efforts should be made to provide  
1269 bias and uncertainty estimates at different time scales, which also requires the use of  
1270 seasonally varying bias correction (i.e. rescaling) parameters.

- 1271 • On pre-processing: very little is known about how spatial and temporal collocation mis-  
1272 matches contribute to bias and uncertainty estimates. Using simple NN or IDW approaches  
1273 to find match-ups between measurements that sample very different soil volumes or were  
1274 taken at different times will give rise to representativeness errors that may considerably  
1275 affect the overall picture of the quality of a product. More research is needed to quantify  
1276 these representativeness errors and to develop resampling methods that more rigorously  
1277 take actual measurement resolution into account.
- 1278 • On metric calculation: most current studies neglect the impact of second-order biases on  
1279 various validation metrics such as the temporal mean difference or the ubRMSD. Several  
1280 attempts are made to mitigate their impact using rescaling methods that match the sta-  
1281 tistical moments of the data sets, yet most of these methods do not account for random  
1282 errors and therefore match the moments in an insufficient manner. More research is needed  
1283 to quantify the impact of suboptimal rescaling on second-order biases, on the impact of  
1284 uncorrected second-order biases on validation metrics, and on how such uncorrected biases  
1285 can be accounted for.
- 1286 • On reference data: validation targets are typically defined against an unknown truth.  
1287 Comparing metrics against error-prone estimates of this truth (i.e. reference data) will  
1288 be inflated by some unknown amount. Efforts should be made to obtain proper bias and  
1289 uncertainty estimates for reference data sets, which should be further used to correct over-  
1290 or underestimated validation metrics (*Miralles et al., 2010; Chen et al., 2017*).
- 1291 • On statistical uncertainty: most validation studies do not report confidence intervals,  
1292 even though they are critical for a reliable interpretation of validation results. Although  
1293 an accurate analytical calculation of confidence intervals for large-scale validation is not  
1294 trivial for all metrics, bootstrapping provides an easy and robust alternative. However,

1295 care must be taken to properly account for spatial and temporal auto-correlation in the  
1296 data.

1297 • On continuity: given the perpetual changes in the land surface character and climate as  
1298 well as progressively increasing data record lengths, sensor drifts, changing reference data  
1299 availability, and improving soil moisture retrieval algorithms, validation should be a con-  
1300 tinuous process and validation reports frequently (at least annually) updated throughout  
1301 and beyond the lifetime of the various satellite missions.

1302 • On accuracy requirements: the well-known soil moisture mission target accuracy require-  
1303 ment of  $0.04 \text{ m}^3\text{m}^{-3}$  (as specified by the Global Climate Observing System as well as  
1304 for individual products and missions), against which soil moisture products are typically  
1305 evaluated, does not relate to the fitness-for-purpose for a specific application. We there-  
1306 fore strongly encourage a closer collaboration between satellite data providers and the soil  
1307 moisture user community to determine application specific accuracy requirements that  
1308 provide deeper insight into what constitutes “good” or “bad” soil moisture data quality,  
1309 thereby fostering the development of improved satellite products.

1310 Finally, many of the discussed principles and methods are not exclusively restricted to soil  
1311 moisture. By setting this example, we hope to also nurture the development and evolution of  
1312 validation good practice guidelines in other Earth observation communities.

## 1313 7 Acknowledgements

1314 We acknowledge the support from the International Space Science Institute (ISSI; <http://www.issibern.ch/>;  
1315 last access: 1 July 2019). This publication is an outcome of the ISSI’s Team  
1316 on “Adding value to soil moisture information for climate studies” and has received funding  
1317 from the earthH2Observe project (European Union’s Seventh Framework Programme, Grant  
1318 Agreement No. 603608), from the KU Leuven C1 internal fund C14/16/045 and from the  
1319 Research Foundation Flanders (FWO-1530019N).

## 1320 Appendix

### 1321 A Validation example

1322 Sec. 5 compiles community-agreed validation good practices into a recommended validation  
1323 protocol. In this appendix, we provide an example that follows this protocol, not to actually  
1324 assess the quality of certain products, but to show an illustrative scenario that can be easily  
1325 extrapolated to more specific validation tasks that readers may face. This includes a comprehen-  
1326 sive description of the validation setup, demonstrative examples of how validation results may  
1327 be presented, and a discussion on where the currently available satellite soil moisture validation  
1328 literature often fails to comply with the good practice recommendations presented here. Results  
1329 shown in this appendix have been generated using the python programming language. All source  
1330 code is available at [https://github.com/alexgruber/validation\\_good\\_practice/](https://github.com/alexgruber/validation_good_practice/) (last ac-  
1331 cess: 1 July 2019). Metric calculation routines have been additionally translated into MATLAB.

#### 1332 A.1 Data sets and study area

1333 Select validation examples are shown for soil moisture retrievals from the Advanced SCATterom-  
1334 eter (ASCAT; *Naeimi et al.*, 2009), the Soil Moisture and Ocean Salinity (SMOS) mission (*Kerr*  
1335 *et al.*, 2010), and the Soil Moisture Active Passive (SMAP) mission (*Entekhabi et al.*, 2010a).  
1336 Reference data used are coarse-resolution model estimates from the Modern-Era Retrospective  
1337 analysis for Research and Applications, Version 2 (MERRA-2; *Gelaro et al.*, 2017). This analy-  
1338 sis is performed over the Contiguous United States (CONUS) using data from the beginning of  
1339 2015 through the end of 2018.

1340 ASCAT data used are the EUMETSAT H SAF H113 data record and its extension H114,  
1341 which are Level 2 (L2) soil moisture products that have been retrieved from inter-calibrated  
1342 backscatter measurements from identical ASCAT instruments onboard the MetOp-A and MetOp-  
1343 B satellites using the TU Wien WAtER Retrieval Package (WARP) algorithm (*Wagner et al.*,  
1344 1999; *Naeimi et al.*, 2009). ASCAT is an active C-band radar with a spatial resolution of 25 km.  
1345 Soil moisture is retrieved as the degree of saturation and sampled onto a 12.5 km discrete global  
1346 grid. Data can be obtained upon registration from [http://hsaf.meteoam.it/soil-moisture.](http://hsaf.meteoam.it/soil-moisture.php)  
1347 [php](http://hsaf.meteoam.it/soil-moisture.php) (last access: 1 July 2019).

1348 SMOS data are the reprocessed L2 soil moisture retrievals version V650, which can be ob-  
1349 tained upon registration from <https://smos-diss.eo.esa.int/> (last access: 1 July 2019; *Kerr*

1350 *et al.*, 2012). SMOS is a passive L-band interferometric radiometer with an average spatial res-  
1351 olution of 43 km. Soil moisture is retrieved in volumetric units and sampled on a 15 km discrete  
1352 global grid.

1353 SMAP data used are the 36 km L2 radiometer-only soil moisture retrievals (SPL2SMP), al-  
1354 gorithm version 5 (R16010) (*O'Neill et al.*, 2018, DOI: 10.5067/SODMLCE6LGLL). The passive  
1355 SMAP radiometer operates at L-band at a spatial resolution of 40 km. Soil moisture is retrieved  
1356 in volumetric units and sampled on the 36 km EASE grid version 2 (*Brodzik et al.*, 2012).

1357 MERRA-2 (*Gelaro et al.*, 2017) is the latest atmospheric reanalysis produced by NASA's  
1358 Global Modelling and Assimilation Office. Soil moisture is estimated on a  $0.5^\circ \times 0.625^\circ$  grid in  
1359 volumetric units as internal state variable of its land surface component, the Catchment Land  
1360 Surface Model (*Koster et al.*, 2000). Here we use soil moisture estimates of the surface layer,  
1361 which refers to the top 5 cm of the soil (*GMAO*, 2015). MERRA-2 data can be downloaded  
1362 from [https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/data\\_access/](https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/data_access/) (last access: 1 July  
1363 2019).

## 1364 A.2 Pre-processing

1365 Unreliable soil moisture retrievals of the individual satellite products are masked out following  
1366 the recommendations of the data providers. ASCAT soil moisture retrievals are masked out if  
1367 the correction flag has a value other than 0 or 4, if the confidence flag and the processing flag  
1368 have values other than 0, or if the surface state flag (*Naeimi et al.*, 2012) has a value other than  
1369 1. SMOS retrievals are masked out if the RFI probability exceeds 0.1 or if the Chi-2 probability  
1370 drops below 0.05. SMAP data are masked out if the retrieval quality flag has a value other than  
1371 0 or 8. In addition, soil moisture retrievals of all satellite products are masked out at time steps  
1372 where MERRA-2 estimates a soil temperature below  $4^\circ\text{C}$  or non-zero snow mass.

1373 ASCAT, SMOS and MERRA-2 are resampled to the 36 km EASE v2 grid that is used  
1374 for SMAP retrievals using a nearest-neighbor approach. Note that ASCAT data is, although  
1375 sampled on a 12.5 km grid, not aggregated as the actual measurement resolution (25 km) is  
1376 already close to the EASE v2 grid resolution. Data sets are collocated in time by resampling them  
1377 to fixed reference time steps with 24 hour intervals using a nearest-neighbor search. Reference  
1378 time steps are selected for each grid cell separately such that they maximize the number of  
1379 collocated time steps where all data sets provide valid soil moisture estimates. Note that the  
1380 choice of this reference time step can increase or decrease the sample size - depending on the

1381 spatial location of the grid cell - by up to a factor of two.

1382 After spatial and temporal collocation, short-term anomalies are calculated for each data set  
1383 using a 35-day moving average window. Long-term anomalies are not considered here because  
1384 the study period of four years (2015-2018) is too short to calculate reliable long-term clima-  
1385 tologies. The term “raw time series” is used to refer to the non-decomposed data, i.e. before  
1386 anomalies have been calculated. For the estimation of unbiased RMSDs, data sets (both raw  
1387 and anomaly time series) are rescaled by matching their temporal mean and standard deviation  
1388 using MERRA-2 as scaling reference for comparability.

### 1389 **A.3 Skill metrics and presentation**

#### 1390 **A.3.1 Sample size**

1391 All metrics are calculated from the same collocated data points, i.e. days where all four data  
1392 sets provide valid soil moisture estimates. The number of temporal matches at each grid cell  
1393 within our study domain is shown in Figure A.1. As discussed in Sec. 4, sample size directly  
1394 translates into statistical power, i.e. reliability (in terms of confidence intervals) of the calculated  
1395 skill metrics. Sample sizes obtained here, which range from 150 in the more mountainous areas  
1396 to up to about 300-500 in the rest of the CONUS, are typically considered high and associated  
1397 with reasonably low confidence intervals for validation purposes.

1398 However, as discussed in Sec. 4.6, confidence intervals are affected by temporal auto-  
1399 correlation. “Effective” sample sizes, corrected for auto-correlation using Eq. (18), are ad-  
1400 ditionally shown in Figure A.1 considering all data sets (for TCA metrics), and in Figure A.2  
1401 for raw soil moisture time series and Figure A.3 for soil moisture anomalies considering different  
1402 data set pairs. Effective sample sizes are considerably smaller than actual sample sizes, especially  
1403 for raw time series due to the strong auto-correlation of the seasonal soil moisture cycle. Since  
1404 auto-correlation levels vary between data sets, effective sample sizes vary when calculated for  
1405 different data set pairs (albeit only slightly), which in turn leads to differences in the confidence  
1406 intervals of relative skill metrics that are calculated between these data pairs.

1407 In the following, all analytical confidence intervals (Eqs. (14), (15), and (17)) are calculated  
1408 using these auto-correlation corrected effective sample sizes. For bootstrapped confidence in-  
1409 tervals, temporal auto-correlation is accounted for using block-bootstrapping (see Sec. 4.6.2)  
1410 where block-lengths are estimated from the same auto-correlation levels that are underlying the  
1411 calculation of effective sample sizes (see Eq. (21)).

### 1412 **A.3.2 Relative metrics**

1413 Figures A.4, A.5 and A.6 show spatial plots of relative (mean) bias, ubRMSD and  $R^2$  (coefficient  
1414 of determination or squared Pearson correlation) estimates for raw soil moisture values, respec-  
1415 tively, and Figures A.7 and A.8 show ubRMSD and  $R^2$  estimates for soil moisture anomalies,  
1416 respectively.

1417 Biases are only calculated for raw soil moisture time series and between soil moisture esti-  
1418 mates that are expressed in the same unit, i.e. for SMOS, SMAP, and MERRA-2 which provide  
1419 estimates of volumetric soil moisture. ASCAT estimates of the degree of saturation could be  
1420 converted into volumetric units using porosity information, but since the quality of soil texture  
1421 maps on these scales is questionable, this is not recommended for bias estimation purposes. Note  
1422 also, that the biases between the remaining three data sets also include collocation and (vertical  
1423 and horizontal) scale mismatches and should therefore be interpreted with care.

1424 Along with the skill estimates, maps of confidence intervals are shown as the difference  
1425 between the upper and lower confidence limits, chosen to be the 90th and the 10th percentile of  
1426 the sampling distribution, respectively. Important to note is that confidence intervals for  $R^2$  and  
1427 ubRMSD estimates depend on the magnitude of the respective skill estimate, and are for  $R^2$  not  
1428 centered around the skill estimate. Misinterpretations may be avoided by directly presenting  
1429 the actual confidence limits (see Sec. 4.7).

1430 We choose a confidence level of 80% because confidence intervals at the more common (yet  
1431 completely arbitrary) 95% confidence level typically become excessively large for the sample  
1432 sizes available from collocated satellite products (*Gruber et al.*, 2019a), especially when taking  
1433 temporal auto-correlation into account.

1434 Figure A.9 shows spatial summary statistics of the relative skill metrics as well as of their  
1435 upper and lower confidence limits. Hardly any skill differences would be considered significant  
1436 when tested in the common way of checking for overlap between upper and lower confidence  
1437 limits, even though Figures A.4 - A.8 show clear differences in spatial patterns.

### 1438 **A.3.3 Triple collocation metrics**

1439 As discussed in Sec. 4, TCA requires three data sets with independent random errors. Since  
1440 errors of SMAP and SMOS are expected to be correlated (see Sec. 5.3), two independent data  
1441 set triplets can be formed, i.e. ASCAT - SMOS - MERRA-2 and ASCAT - SMAP - MERRA-2.  
1442 This results in unambiguous skill estimates for SMAP and SMOS, and in two skill estimates for

1443 ASCAT, which are averaged for increased precision.

1444 Figures A.10 and A.11 show spatial plots of TCA-based ubRMSE and  $R^2$  (coefficient of  
1445 determination w.r.t. the unknown truth) estimates, respectively, and Figures A.12 and A.13  
1446 show ubRMSE and  $R^2$  estimates for short-term soil moisture anomalies, respectively. The skill  
1447 estimates represent the median of the bootstrapped sampling distribution, which are more robust  
1448 than the direct estimates, and 80 % confidence intervals (i.e. the range between the 90th and  
1449 the 10th percentile of the bootstrapped sampling distribution) are provided. Spatial summary  
1450 statistics of the TCA estimates (sampling distribution median) as well as of the upper and lower  
1451 confidence limits are shown in Figure A.14.

1452 The two degrees of freedom in TCA-based ASCAT skill estimates can not only be used for  
1453 increasing the precision of the estimates by averaging them, but also to verify if TCA assumptions  
1454 (i.e. zero error cross-correlation and error orthogonality) are met because if so, skill estimates  
1455 should be identical. To this end, Figure A.15 shows the differences between  $R^2$  and ubRMSE  
1456 estimates for ASCAT when calculated once using SMOS as third data set and once using SMAP  
1457 as third data set.

1458 On average, differences are close to zero and especially  $R^2$  estimates do not exhibit spatial  
1459 patterns of notable magnitude, which suggests that differences are mainly caused by sampling  
1460 errors and hence that the TCA assumptions are generally respected. Some positive skill biases  
1461 for raw soil moisture estimation for ASCAT are apparent in some northern and western parts  
1462 of the CONUS, with skill estimates being slightly higher when using SMOS rather than SMAP  
1463 in the triplet. These areas strongly coincide with regions of generally poor ASCAT performance  
1464 (see Figure A.11), which is more pronounced in the ubRMSD because SNR biases of a given  
1465 magnitude are associated with larger biases in error variance at low SNR levels than at high  
1466 SNR levels. (see Sec. 4.7). Poor ASCAT performance in the northern CONUS is associated  
1467 with issues in the vegetation correction of the WARP retrieval algorithm (see Sec. A.1). These  
1468 uncorrected vegetation signals are removed when using soil moisture anomalies, which results in  
1469 a considerable increase in skill metrics (see Figure A.13) and also removes the non-zero difference  
1470 in ASCAT skill estimates when using SMOS versus SMAP for TCA, i.e. spurious error cross-  
1471 correlations (see Figure A.15).

## 1472 A.4 Final remarks

1473 In this appendix, we provide an illustrative validation example that follows the good practice  
1474 guidelines presented in this paper. For brevity, we omit the presentation of ground data compar-  
1475 isons, which can be calculated and presented in the exact same way as the area-wide coarse-scale  
1476 comparisons shown above. For simplicity, results are presented in spatial maps and boxplots  
1477 that cover all of CONUS without further stratification. For summary information or if metrics  
1478 are only computed at a few locations using ground reference data, results could be further pre-  
1479 sented in tabular format. Some examples of comprehensive ground reference data comparison  
1480 including both sparse networks and core validation sites can be found in *Dorigo et al. (2015)*;  
1481 *Chen et al. (2017)*; *Colliander et al. (2017)*.

## 1482 References

- 1483 Aitkin, A. (1936), On least squares and linear combination of observations, *Proceedings of the*  
1484 *Royal Society of Edinburgh*, **55**, p. 42–48, doi:10.1017/S0370164600014346.
- 1485 Al-Yaari, A., J.-P. Wigneron, A. Ducharne, Y. Kerr, W. Wagner, G. De Lannoy, R. Reichle,  
1486 A. Al Bitar, W. Dorigo, P. Richaume, et al. (2014), Global-scale comparison of passive (SMOS)  
1487 and active (ASCAT) satellite based microwave soil moisture retrievals with soil moisture  
1488 simulations (MERRA-Land), *Remote Sensing of Environment*, **152**, p. 614–626, doi:10.1016/  
1489 j.rse.2014.07.013.
- 1490 Albergel, C., C. Ruediger, T. Pellarin, J. Calvet, N. Fritz, F. Froissard, D. Suquia, A. Pe-  
1491 titpa, B. Piguet, and E. Martin (2008), From near-surface to root-zone soil moisture us-  
1492 ing an exponential filter: an assessment of the method based on in-situ observations  
1493 and model simulations., *Hydrology and earth system sciences.*, **12**(6), p. 1323–1337, doi:  
1494 10.5194/hess-12-1323-2008.
- 1495 Albergel, C., E. Zakharova, J.-C. Calvet, M. Zribi, M. Pardé, J.-P. Wigneron, N. Novello,  
1496 Y. Kerr, A. Mialon, and N. ed Dine Fritz (2011), A first assessment of the smos data in  
1497 southwestern france using in situ and airborne soil moisture estimates: The carols airborne  
1498 campaign, *Remote Sensing of Environment*, **115**(10), p. 2718 – 2728, doi:10.1016/j.rse.2011.  
1499 06.012.
- 1500 Albergel, C., P. de Rosnay, C. Gruhier, J. Munoz-Sabater, S. Hasenauer, L. Isaksen, Y. Kerr, and



1501 W. Wagner (2012), Evaluation of remotely sensed and modelled soil moisture products using  
1502 global ground-based in situ observations, *Remote Sensing of Environment*, **118**, p. 215–226,  
1503 doi:10.1016/j.rse.2011.11.017.

1504 Albergel, C., W. Dorigo, R. Reichle, G. Balsamo, P. De Rosnay, J. Muñoz-Sabater, L. Isaksen,  
1505 R. De Jeu, and W. Wagner (2013), Skill and global trend analysis of soil moisture from  
1506 reanalyses and microwave remote sensing, *Journal of Hydrometeorology*, **14**(4), p. 1259–1277,  
1507 doi:10.1175/JHM-D-12-0161.1.

1508 Babaeian, E., M. Sadeghi, S. B. Jones, C. Montzka, H. Vereecken, and M. Tuller (2019), Ground,  
1509 proximal, and satellite remote sensing of soil moisture, *Reviews of Geophysics*, **57**, doi:10.1029/  
1510 2018RG000618.

1511 Baker, M. (2016), 1,500 scientists lift the lid on reproducibility, *Nature News*, **533**(7604), p. 452,  
1512 doi:10.1038/533452a.

1513 Balsamo, G., C. Albergel, A. Beljaars, S. Boussetta, E. Brun, H. Cloke, D. Dee, E. Dutra,  
1514 J. Muñoz-Sabater, F. Pappenberger, et al. (2015), ERA-Interim/Land: a global land surface  
1515 reanalysis data set, *Hydrology and Earth System Sciences*, **19**(1), p. 389–407, doi:10.5194/  
1516 hess-19-389-2015.

1517 Bartalis, Z., R. Kidd, and K. Scipal (2006), Development and implementation of a discrete  
1518 global grid system for soil moisture retrieval using the MetOp ASCAT scatterometer, in *1st*  
1519 *EPS/MetOp RAO Workshop*, vol. ESA SP-618, ESRIN, Frascati, Italy.

1520 Bauer-Marschallinger, B., D. Sabel, and W. Wagner (2014), Optimisation of global grids for  
1521 high-resolution remote sensing data, *Computers & Geosciences*, **72**, p. 84–93, doi:10.1016/j.  
1522 cageo.2014.07.005.

1523 Bauer-Marschallinger, B., C. Paulik, S. Hochstöger, T. Mistelbauer, S. Modanesi, L. Ciabatta,  
1524 C. Massari, L. Brocca, and W. Wagner (2018), Soil moisture from fusion of scatterometer  
1525 and sar: Closing the scale gap with temporal filtering, *Remote Sensing*, **10**(7), p. 1030, doi:  
1526 10.3390/rs10071030.

1527 Bindlish, R., T. J. Jackson, A. J. Gasiewski, M. Klein, and E. G. Njoku (2006), Soil moisture  
1528 mapping and AMSR-E validation using the PSR in SMEX02, *Remote Sensing of Environment*,  
1529 **103**(2), p. 127–139, doi:10.1016/j.rse.2005.02.003.

- 1530 Bindlish, R., T. Jackson, A. Gasiewski, B. Stankov, M. Klein, M. Cosh, I. Mladenova, C. Watts,  
1531 E. Vivoni, V. Lakshmi, et al. (2008), Aircraft based soil moisture retrievals under mixed  
1532 vegetation and topographic conditions, *Remote Sensing of Environment*, **112**(2), p. 375–390,  
1533 doi:10.1016/j.rse.2007.01.024.
- 1534 Blyth, C. R. (1972), On Simpson’s paradox and the sure-thing principle, *Journal of the American*  
1535 *Statistical Association*, **67**(338), p. 364–366, doi:10.1080/01621459.1972.10482387.
- 1536 Bogena, H., C. Montzka, J. Huisman, A. Graf, M. Schmidt, M. Stockinger, C. von Hebel,  
1537 H. Hendricks-Franssen, J. van der Kruk, W. Tappe, et al. (2018), The TERENO-Rur hydro-  
1538 logical observatory: A multiscale multi-compartment research platform for the advancement  
1539 of hydrological science, *Vadose Zone Journal*, **17**(1), doi:10.2136/vzj2018.03.0055.
- 1540 Bolten, J. D., W. T. Crow, X. Zhan, T. J. Jackson, and C. A. Reynolds (2010), Evaluating the  
1541 utility of remotely sensed soil moisture retrievals for operational agricultural drought moni-  
1542 toring, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*,  
1543 **3**(1), p. 57–66, doi:10.1109/JSTARS.2009.2037163.
- 1544 Bonett, D. G., and T. A. Wright (2000), Sample size requirements for estimating pearson, kendall  
1545 and spearman correlations, *Psychometrika*, **65**(1), p. 23–28, doi:10.1007/BF02294183.
- 1546 Brocca, L., F. Melone, T. Moramarco, and R. Morbidelli (2010a), Spatial-temporal variability of  
1547 soil moisture and its estimation across scales, *Water Resources Research*, **46**(2), doi:10.1029/  
1548 2009WR008016.
- 1549 Brocca, L., F. Melone, T. Moramarco, W. Wagner, and S. Hasenauer (2010b), ASCAT soil  
1550 wetness index validation through in situ and modeled soil moisture data in central italy,  
1551 *Remote Sensing of Environment*, **114**(11), p. 2745–2755, doi:10.1016/j.rse.2010.06.009.
- 1552 Brocca, L., S. Hasenauer, T. Lacava, F. Melone, T. Moramarco, W. Wagner, W. Dorigo, P. Mat-  
1553 gen, J. Martinez-Fernandez, P. Llorens, J. Latron, C. Martin, and M. Bittelli (2011), Soil  
1554 moisture estimation through ASCAT and AMSR-E sensors: An intercomparison and vali-  
1555 dation study across europe, *Remote Sensing of Environment*, **115**(12), p. 3390–3408, doi:  
1556 10.1016/j.rse.2011.08.003.
- 1557 Brocca, L., T. Tullo, F. Melone, T. Moramarco, and R. Morbidelli (2012), Catchment scale soil

1558 moisture spatial-temporal variability, *Journal of Hydrology*, **422-423**, p. 63–75, doi:10.1016/  
1559 j.jhydrol.2011.12.039.

1560 Brodzik, M. J., B. Billingsley, T. Haran, B. Raup, and M. H. Savoie (2012), EASE-Grid 2.0:  
1561 Incremental but significant improvements for earth-gridded data sets, *ISPRS International*  
1562 *Journal of Geo-Information*, **1**(1), p. 32–45, doi:10.3390/ijgi1010032.

1563 Burgin, M. S., A. Colliander, E. G. Njoku, S. K. Chan, F. Cabot, Y. H. Kerr, R. Bindlish, T. J.  
1564 Jackson, D. Entekhabi, and S. H. Yueh (2017), A comparative study of the SMAP passive  
1565 soil moisture product with existing satellite-based soil moisture products, *IEEE Transactions*  
1566 *on Geoscience and Remote Sensing*, **55**(5), p. 2959–2971, doi:10.1109/TGRS.2017.2656859.

1567 Caires, S., and A. Sterl (2003), Validation of ocean wind and wave data using triple collocation,  
1568 *Journal of Geophysical Research: Oceans*, **108**(C3), doi:10.1029/2002JC001491.

1569 Chauhan, N. S., S. Miller, and P. Ardanuy (2003), Spaceborne soil moisture estimation at  
1570 high resolution: a microwave-optical/ir synergistic approach, *International Journal of Remote*  
1571 *Sensing*, **24**(22), p. 4599–4622, doi:10.1080/0143116031000156837.

1572 Chen, F., W. T. Crow, A. Colliander, M. H. Cosh, T. J. Jackson, R. Bindlish, R. H. Reichle,  
1573 S. K. Chan, D. D. Bosch, P. J. Starks, et al. (2017), Application of triple collocation in ground-  
1574 based validation of soil moisture active/passive (SMAP) level 2 data products, *IEEE Journal*  
1575 *of Selected Topics in Applied Earth Observations and Remote Sensing*, **10**(2), p. 489–502,  
1576 doi:10.1109/JSTARS.2016.2569998.

1577 Chen, F., W. T. Crow, R. Bindlish, A. Colliander, M. S. Burgin, J. Asanuma, and K. Aida (2018),  
1578 Global-scale evaluation of SMAP, SMOS and ASCAT soil moisture products using triple  
1579 collocation, *Remote Sensing of Environment*, **214**, p. 1–13, doi:10.1016/j.rse.2018.05.008.

1580 Colliander, A., T. J. Jackson, R. Bindlish, S. Chan, N. Das, S. Kim, M. Cosh, R. Dunbar,  
1581 L. Dang, L. Pashaian, et al. (2017), Validation of SMAP surface soil moisture products with  
1582 core validation sites, *Remote sensing of environment*, **191**, p. 215–231, doi:10.1016/j.rse.2017.  
1583 01.021.

1584 Corey, D. M., W. P. Dunlap, and M. J. Burke (1998), Averaging correlations: Expected val-  
1585 ues and bias in combined pearson rs and fisher’s z transformations, *The Journal of general*  
1586 *psychology*, **125**(3), p. 245–261, doi:10.1080/00221309809595548.

- 1587 Cosh, M. H., T. J. Jackson, R. Bindlish, and J. H. Prueger (2004), Watershed scale temporal and  
1588 spatial stability of soil moisture and its role in validating satellite estimates, *Remote sensing*  
1589 *of Environment*, **92**(4), p. 427–435, doi:10.1016/j.rse.2004.02.016.
- 1590 Cosh, M. H., T. J. Jackson, P. Starks, and G. Heathman (2006), Temporal stability of surface  
1591 soil moisture in the little washita river watershed and its applications in satellite soil moisture  
1592 product validation, *Journal of Hydrology*, **323**(1–4), p. 168–177, doi:10.1016/j.jhydrol.2005.  
1593 08.020.
- 1594 Crow, W. T., A. A. Berg, M. H. Cosh, A. Loew, B. P. Mohanty, R. Panciera, P. de Rosnay,  
1595 D. Ryu, and J. P. Walker (2012), Upscaling sparse ground-based soil moisture observations  
1596 for the validation of coarse-resolution satellite soil moisture products, *Rev. Geophys.*, **50**(2),  
1597 p. RG2002, doi:10.1029/2011RG000372.
- 1598 Cuenca, R. H., D. E. Stangel, and S. F. Kelly (1997), Soil water balance in a boreal forest, *Journal*  
1599 *of Geophysical Research-Atmospheres*, **102**(D 24), p. 29,355–29,365, doi:10.1029/97JD02312.
- 1600 Das, N. N., D. Entekhabi, S. Kim, T. Jagdhuber, S. Dunbar, S. Yueh, and A. Colliander (2017),  
1601 High-resolution enhanced product based on smap active-passive approach using sentinel 1a  
1602 and 1b sar data, in *2017 IEEE International Geoscience and Remote Sensing Symposium*  
1603 *(IGARSS)*, p. 2543–2545, IEEE, doi:10.1109/IGARSS.2017.8127513.
- 1604 Dawdy, D., and N. Matalas (1964), *Statistical and probability analysis of hydrologic data, part*  
1605 *III: Analysis of variance, covariance and time series*, McGraw-Hill.
- 1606 De Lannoy, G. J., and R. H. Reichle (2016), Assimilation of SMOS brightness temperatures  
1607 or soil moisture retrievals into a land surface model, *Hydrology and Earth System Sciences*,  
1608 **20**(12), p. 4895–4911, doi:10.5194/hess-20-4895-2016.
- 1609 de Nijs, A. H., R. M. Parinussa, R. A. de Jeu, J. Schellekens, and T. R. Holmes (2015), A  
1610 methodology to determine radio-frequency interference in AMSR2 observations, *Geoscience*  
1611 *and Remote Sensing, IEEE Transactions on*, **53**(9), p. 5148–5159, doi:10.1109/TGRS.2015.  
1612 2417653.
- 1613 De Rosnay, P., J.-C. Calvet, Y. Kerr, J.-P. Wigneron, F. Lemaître, M. J. Escorihuela, J. M.  
1614 Sabater, K. Saleh, J. Barrié, G. Bouhours, et al. (2006), SMOSREX: A long term field cam-

1615 paign experiment for soil moisture and land surface processes remote sensing, *Remote Sensing*  
1616 *of Environment*, **102**(3-4), p. 377–389, doi:10.1016/j.rse.2006.02.021.

1617 Dee, D. P. (2005), Bias and data assimilation, *Quarterly Journal of the Royal Meteorological*  
1618 *Society*, **131**(613), p. 3323–3343, doi:10.1256/qj.05.137.

1619 Djamai, N., R. Magagi, K. Goïta, M. Hosseini, M. H. Cosh, A. Berg, and B. Toth (2015),  
1620 Evaluation of SMOS soil moisture products over the CanEx-SM10 area, *Journal of hydrology*,  
1621 **520**, p. 254–267, doi:10.1016/j.jhydrol.2014.11.026.

1622 Dorigo, W., P. van Oevelen, W. Wagner, M. Drusch, S. Mecklenburg, A. Robock, and T. Jackson  
1623 (2011a), A new international network for in situ soil moisture data, *Eos Transactions AGU*,  
1624 **92**(17), p. 141–142, doi:10.1029/2011EO170001.

1625 Dorigo, W., R. de Jeu, D. Chung, R. Parinussa, Y. Liu, W. Wagner, and D. Fernández-Prieto  
1626 (2012), Evaluating global trends (1988–2010) in harmonized multi-satellite surface soil mois-  
1627 ture, *Geophysical Research Letters*, **39**(18), doi:10.1029/2012GL052988.

1628 Dorigo, W., A. Xaver, M. Vreugdenhil, A. Gruber, H. A. A. Sanchis-Dufau, D. Zamojski,  
1629 C. Cordes, W. Wagner, and M. Drusch (2013), Global automated quality control of in situ  
1630 soil moisture data from the international soil moisture network, *Vadose Zone Journal*, **12**(3),  
1631 doi:10.2136/vzj2012.0097.

1632 Dorigo, W., A. Gruber, R. De Jeu, W. Wagner, T. Stacke, A. Loew, C. Albergel, L. Brocca,  
1633 D. Chung, R. Parinussa, et al. (2015), Evaluation of the ESA CCI soil moisture product using  
1634 ground-based observations, *Remote Sensing of Environment*, **162**, p. 380–395, doi:10.1016/j.  
1635 rse.2014.07.023.

1636 Dorigo, W., W. Wagner, C. Albergel, F. Albrecht, G. Balsamo, L. Brocca, D. Chung, M. Ertl,  
1637 M. Forkel, A. Gruber, et al. (2017), ESA CCI soil moisture for improved earth system un-  
1638 derstanding: state-of-the art and future directions, *Remote Sensing of Environment*, **203**,  
1639 p. 185–215, doi:10.1016/j.rse.2017.07.001.

1640 Dorigo, W. A., K. Scipal, R. M. Parinussa, Y. Y. Liu, W. Wagner, R. A. M. de Jeu, and  
1641 V. Naeimi (2010), Error characterisation of global active and passive microwave soil moisture  
1642 datasets, *Hydrol. Earth Syst. Sci.*, **14**(12), p. 2605–2616, doi:10.5194/hessd-7-5621-2010.

- 1643 Dorigo, W. A., W. Wagner, R. Hohensinn, S. Hahn, C. Paulik, A. Xaver, A. Gruber, M. Drusch,  
1644 S. Mecklenburg, P. van Oevelen, A. Robock, and T. Jackson (2011b), The international soil  
1645 moisture network: a data hosting facility for global in situ soil moisture measurements, *Hydrol.*  
1646 *Earth Syst. Sci.*, **15**(5), p. 1675–1698, doi:10.5194/hess-15-1675-2011.
- 1647 Draper, C., and R. Reichle (2015), The impact of near-surface soil moisture assimilation at  
1648 subseasonal, seasonal, and inter-annual timescales, *Hydrology and Earth System Sciences*,  
1649 **19**(12), p. 4831, doi:10.5194/hess-19-4831-2015.
- 1650 Draper, C., R. Reichle, G. De Lannoy, and Q. Liu (2012), Assimilation of passive and ac-  
1651 tive microwave soil moisture retrievals, *Geophysical Research Letters*, **39**(4), doi:10.1029/  
1652 2011GL050655.
- 1653 Draper, C., R. Reichle, R. de Jeu, V. Naeimi, R. Parinussa, and W. Wagner (2013), Estimating  
1654 root mean square errors in remotely sensed soil moisture over continental scale domains,  
1655 *Remote Sensing of Environment*, **137**, p. 288–298, doi:10.1016/j.rse.2013.06.013.
- 1656 Efron, B., and R. Tibshirani (1986), Bootstrap methods for standard errors, confidence intervals,  
1657 and other measures of statistical accuracy, *Statistical science*, **1**(1), p. 54–75, doi:10.1214/ss/  
1658 1177013815.
- 1659 Entekhabi, D., E. Njoku, P. O’Neill, K. Kellogg, W. Crow, W. Edelstein, J. Entin, S. Good-  
1660 man, T. Jackson, J. Johnson, J. Kimball, J. Piepmeier, R. Koster, N. Martin, K. McDonald,  
1661 M. Moghaddam, S. Moran, R. Reichle, J. Shi, M. Spencer, S. Thurman, L. Tsang, and  
1662 J. Van Zyl (2010a), The soil moisture active passive (SMAP) mission, *Proceedings of the*  
1663 *IEEE*, **98**(5), p. 704–716, doi:10.1109/JPROC.2010.2043918.
- 1664 Entekhabi, D., R. H. Reichle, R. D. Koster, and W. T. Crow (2010b), Performance metrics  
1665 for soil moisture retrievals and application requirements, *J. Hydrometeor*, **11**(3), p. 832–840,  
1666 doi:10.1175/2010JHM1223.1.
- 1667 Famiglietti, J. S., D. Ryu, A. A. Berg, M. Rodell, and T. J. Jackson (2008), Field observations  
1668 of soil moisture variability across scales, *Water Resour. Res.*, **44**(1), p. W01,423, doi:10.1029/  
1669 2006WR005804.
- 1670 Figa-Saldaña, J., J. J. Wilson, E. Attema, R. Gelsthorpe, M. Drinkwater, and A. Stoffelen  
1671 (2002), The advanced scatterometer (ASCAT) on the meteorological operational (MetOp)

1672 platform: A follow on for european wind scatterometers, *Canadian Journal of Remote Sensing*,  
1673 **28**(3), p. 404–412, doi:10.5589/m02-035.

1674 Fox, N. (2010), A guide to “reference standards” in support of quality assur-  
1675 ance requirements of GEO, *Tech. Rep. QA4EO-QAEO-GEN-DQK-003, v4.0*, QA4EO,  
1676 [http://qa4eo.org/docs/QA4EO-QAEO-GEN-DQK-003\\_v4.0.pdf](http://qa4eo.org/docs/QA4EO-QAEO-GEN-DQK-003_v4.0.pdf), last access: 1 July 2019.

1677 Gelaro, R., W. McCarty, M. J. Suárez, R. Todling, A. Molod, L. Takacs, C. A. Randles, A. Dar-  
1678 menov, M. G. Bosilovich, R. Reichle, et al. (2017), The modern-era retrospective analysis for  
1679 research and applications, version 2 (MERRA-2), *Journal of Climate*, **30**(14), p. 5419–5454,  
1680 doi:10.1175/JCLI-D-16-0758.1.

1681 Gelman, A., and H. Stern (2006), The difference between “significant” and “not significant” is  
1682 not itself statistically significant, *The American Statistician*, **60**(4), p. 328–331, doi:10.1198/  
1683 000313006X152649.

1684 Gilleland, E. (2010), Confidence intervals for forecast verification, *NCAR Technical Note*, **TN-**  
1685 **479**, doi:10.5065/D6WD3XJM.

1686 GMAO (2015), Global Modeling and Assimilation Office (GMAO), MERRA-2 tavg1\_2d\_lnd\_Nx:  
1687 2d,1-Hourly,Time-Averaged,Single-Level,Assimilation,Land Surface Diagnostics V5.12.4,  
1688 Greenbelt, MD, USA, Goddard Earth Sciences Data and Information Services Center (GES  
1689 DISC), Accessed: 1 Nov 2018, doi:10.5067/RKPHT8KC1Y1T.

1690 Greenland, S., S. J. Senn, K. J. Rothman, J. B. Carlin, C. Poole, S. N. Goodman, and D. G. Alt-  
1691 man (2016), Statistical tests, p values, confidence intervals, and power: a guide to misinterpre-  
1692 tations, *European journal of epidemiology*, **31**(4), p. 337–350, doi:10.1007/s10654-016-0149-3.

1693 Gruber, A., W. Dorigo, S. Zwieback, A. Xaver, and W. Wagner (2013a), Characterizing coarse-  
1694 scale representativeness of in situ soil moisture measurements from the international soil mois-  
1695 ture network, *Vadose Zone Journal*, **12**(2), doi:10.2136/vzj2012.0170.

1696 Gruber, A., W. Wagner, A. Hegyiova, F. Greifeneder, and S. Schläffer (2013b), Potential of  
1697 sentinel-1 for high-resolution soil moisture monitoring, in *Geoscience and Remote Sensing*  
1698 *Symposium (IGARSS), 2013 IEEE International*, p. 4030–4033, IEEE, doi:10.1109/IGARSS.  
1699 2017.8127513.

1700 Gruber, A., W. Crow, W. Dorigo, and W. Wagner (2015), The potential of 2D Kalman filtering  
1701 for soil moisture data assimilation, *Remote Sensing of Environment*, **171**, p. 137–148, doi:  
1702 10.1016/j.rse.2015.10.019.

1703 Gruber, A., C.-H. Su, S. Zwieback, W. Crow, W. Dorigo, and W. Wagner (2016a), Recent  
1704 advances in (soil moisture) triple collocation analysis, *International Journal of Applied Earth*  
1705 *Observation and Geoinformation*, **45**, p. 200–211, doi:10.1016/j.jag.2015.09.002.

1706 Gruber, A., C.-H. Su, W. Crow, S. Zwieback, W. Dorigo, and W. Wagner (2016b), Estimating  
1707 error cross-correlations in soil moisture data sets using extended collocation analysis, *Journal*  
1708 *of Geophysical Research: Atmospheres*, **121(3)**, p. 1208–1219, doi:10.1002/2015JD024027.

1709 Gruber, A., W. A. Dorigo, W. Crow, and W. Wagner (2017), Triple collocation-based merging  
1710 of satellite soil moisture retrievals, *IEEE Transactions on Geoscience and Remote Sensing*,  
1711 **55(12)**, p. 6780–6792, doi:10.1109/TGRS.2017.2734070.

1712 Gruber, A., W. Crow, and W. Dorigo (2018), Assimilation of spatially sparse in situ soil moisture  
1713 networks into a continuous model domain, *Water Resources Research*, **54(2)**, p. 1353–1367,  
1714 doi:10.1002/2017WR021277.

1715 Gruber, A., T. Scanlon, R. van der Schalie, W. Wagner, and W. Dorigo (2019a), Evolution  
1716 of the esa cci soil moisture climate data records and their underlying merging methodology,  
1717 *Earth System Science Data*, **11(2)**, p. 717–739, doi:10.5194/essd-11-717-2019.

1718 Gruber, A., G. D. Lannoy, and W. Crow (2019b), A monte carlo based adaptive kalman filtering  
1719 framework for soil moisture data assimilation, *Remote Sensing of Environment*, **228**, p. 105  
1720 – 114, doi:10.1016/j.rse.2019.04.003.

1721 Gupta, H. V., H. Kling, K. K. Yilmaz, and G. F. Martinez (2009), Decomposition of the mean  
1722 squared error and NSE performance criteria: Implications for improving hydrological mod-  
1723 elling, *Journal of Hydrology*, **377(1)**, p. 80–91, doi:10.1016/j.jhydrol.2009.08.003.

1724 H-SAF (2017), Product validation report (PVR) h111 metop ASCAT soil mois-  
1725 ture, *Tech. Rep. SAF/HSAF/CDOP3/PVR/H111, v0.3*, EUMETSAT H SAF reports,  
1726 [http://hsaf.meteoam.it/documents/PVR/H111\\_ASCAT\\_SSM\\_CDR\\_PVR\\_v0.3.pdf](http://hsaf.meteoam.it/documents/PVR/H111_ASCAT_SSM_CDR_PVR_v0.3.pdf) (last ac-  
1727 cess: 1 July 2019).



1728 H-SAF (2018), Algorithm theoretical baseline document (ATBD) soil mois-  
1729 ture data records, metop ASCAT soil moisture time series, *Tech.*  
1730 *Rep. SAF/HSAF/CDOP3/ATBD, v0.7*, EUMETSAT H SAF reports,  
1731 [http://hsaf.meteoam.it/documents/ATDD/ASCAT\\_SSM\\_CDR\\_ATBD\\_v0.7.pdf](http://hsaf.meteoam.it/documents/ATDD/ASCAT_SSM_CDR_ATBD_v0.7.pdf) (last ac-  
1732 cess: 1 July 2019).

1733 Jackson, T., M. Cosh, R. Bindlish, P. Starks, D. Bosch, M. Seyfried, D. Goodrich, M. Moran, and  
1734 J. Du (2010), Validation of advanced microwave scanning radiometer soil moisture products,  
1735 *Geoscience and Remote Sensing, IEEE Transactions on*, **48**(12), p. 4256–4272, doi:10.1109/  
1736 TGRS.2010.2051035.

1737 Jackson, T., A. Colliander, J. Kimball, R. Reichle, W. Crow, D. Entekhabi, and P. Neill (2012),  
1738 Science data calibration and validation plan, *SMAP Mission, NASA Jet Propuls. Lab.*

1739 Jackson, T. J., D. M. Le Vine, C. T. Swift, T. J. Schmugge, and F. R. Schiebe (1995), Large  
1740 area mapping of soil moisture using the ESTAR passive microwave radiometer in washita’92,  
1741 *Remote sensing of Environment*, **54**(1), p. 27–37, doi:10.1016/0034-4257(95)00084-E.

1742 Jackson, T. J., D. M. Le Vine, A. Y. Hsu, A. Oldak, P. J. Starks, C. T. Swift, J. D. Isham,  
1743 and M. Haken (1999), Soil moisture mapping at regional scales using microwave radiometry:  
1744 The southern great plains hydrology experiment, *IEEE transactions on geoscience and remote*  
1745 *sensing*, **37**(5), p. 2136–2151, doi:10.1109/36.789610.

1746 Jackson, T. J., R. Bindlish, A. J. Gasiewski, B. Stankov, M. Klein, E. G. Njoku, D. Bosch,  
1747 T. L. Coleman, C. A. Laymon, and P. Starks (2005), Polarimetric scanning radiometer C-  
1748 and X-band microwave observations during SMEX03, *IEEE Transactions on Geoscience and*  
1749 *Remote Sensing*, **43**(11), p. 2418–2430, doi:10.1109/TGRS.2005.857625.

1750 JCGM (2008), Evaluation of measurement data—guide to the expression of uncer-  
1751 tainty in measurement (GUM), *Tech. Rep. JCGM 100:2008*, Bureau International des  
1752 Poids et Mesures (BIPM), Joint Committee for Guides in Metrology (JCGM), URL:  
1753 <https://www.bipm.org/en/publications/guides/gum.html>, last access: 1 July 2019.

1754 JCGM (2012), International vocabulary of metrology—basic and general concepts and asso-  
1755 ciated terms (VIM 3rd edition), *Tech. Rep. JCGM 200:2012*, Bureau International des  
1756 Poids et Mesures (BIPM), Joint Committee for Guides in Metrology (JCGM), URL:  
1757 <https://www.bipm.org/en/publications/guides/vim.html>, last access: 1 July 2019.

1758 Justice, C., A. Belward, J. Morisette, P. Lewis, J. Privette, and F. Baret (2000), Developments  
1759 in the 'validation' of satellite sensor products for the study of the land surface, *International*  
1760 *Journal of Remote Sensing*, **21**(17), p. 3383–3390, doi:10.1080/014311600750020000.

1761 Kerr, Y., P. Waldteufel, J.-P. Wigneron, S. Delwart, F. Cabot, J. Boutin, M. Escorihuela,  
1762 J. Font, N. Reul, C. Gruhier, S. Juglea, M. Drinkwater, A. Hahne, M. Martin-Neira, and  
1763 S. Mecklenburg (2010), The SMOS mission: New tool for monitoring key elements of the global  
1764 water cycle, *Proceedings of the IEEE*, **98**(5), p. 666–687, doi:10.1109/JPROC.2010.2043032.

1765 Kerr, Y. H., P. Waldteufel, J.-P. Wigneron, J. Martinuzzi, J. Font, and M. Berger (2001), Soil  
1766 moisture retrieval from space: The soil moisture and ocean salinity (SMOS) mission, *IEEE*  
1767 *transactions on Geoscience and remote sensing*, **39**(8), p. 1729–1735, doi:10.1109/36.942551.

1768 Kerr, Y. H., P. Waldteufel, P. Richaume, J. P. Wigneron, P. Ferrazzoli, A. Mahmoodi,  
1769 A. Al Bitar, F. Cabot, C. Gruhier, S. E. Juglea, et al. (2012), The SMOS soil moisture  
1770 retrieval algorithm, *IEEE Transactions on Geoscience and Remote Sensing*, **50**(5), p. 1384–  
1771 1403, doi:10.1109/TGRS.2012.2184548.

1772 Kerr, Y. H., A. Al-Yaari, N. Rodriguez-Fernandez, M. Parrens, B. Molero, D. Leroux, S. Bircher,  
1773 A. Mahmoodi, A. Mialon, P. Richaume, et al. (2016), Overview of SMOS performance in terms  
1774 of global soil moisture monitoring after six years in operation, *Remote Sensing of Environment*,  
1775 **180**, p. 40–63, doi:10.1016/j.rse.2016.02.042.

1776 Kolassa, J., P. Gentine, C. Prigent, F. Aires, and S. Alemohammad (2017), Soil moisture retrieval  
1777 from amsr-e and ascat microwave observation synergy. part 2: Product evaluation, *Remote*  
1778 *Sensing of Environment*, **195**, p. 202 – 217, doi:https://doi.org/10.1016/j.rse.2017.04.020.

1779 Koster, R. D., M. J. Suarez, A. Ducharne, M. Stieglitz, and P. Kumar (2000), A catchment-  
1780 based approach to modeling land surface processes in a general circulation model: 1. model  
1781 structure, *Journal of Geophysical Research: Atmospheres*, **105**(D20), p. 24,809–24,822, doi:  
1782 10.1029/2000JD900327.

1783 Koster, R. D., Z. Guo, R. Yang, P. A. Dirmeyer, K. Mitchell, and M. J. Puma (2009), On  
1784 the nature of soil moisture in land surface models, *Journal of Climate*, **22**(16), p. 4322–4335,  
1785 doi:10.1175/2009JCLI2832.1.

- 1786 Kumar, S. V., R. H. Reichle, K. W. Harrison, C. D. Peters-Lidard, S. Yatheendradas, and J. A.  
1787 Santanello (2012), A comparison of methods for a priori bias correction in soil moisture data  
1788 assimilation, *Water Resour. Res.*, **48**(3), p. W03,515, doi:10.1029/2010WR010261.
- 1789 Lahoz, W. A., and G. J. De Lannoy (2014), Closing the gaps in our knowledge of the hydrological  
1790 cycle over land: Conceptual problems, *Surveys in Geophysics*, **35**(3), p. 623–660, doi:10.1007/  
1791 s10712-013-9221-7.
- 1792 Loew, A., W. Bell, L. Brocca, C. E. Bulgin, J. Burdanowitz, X. Calbet, R. V. Donner,  
1793 D. Ghent, A. Gruber, T. Kaminski, et al. (2017), Validation practices for satellite-based  
1794 earth observation data across communities, *Reviews of Geophysics*, **55**(3), p. 779–817, doi:  
1795 10.1002/2017RG000562.
- 1796 Macelloni, G., M. Brogioni, P. Pampaloni, A. Cagnati, and M. R. Drinkwater (2006), DOMEX  
1797 2004: An experimental campaign at Dome-C antarctica for the calibration of spaceborne  
1798 low-frequency microwave radiometers, *IEEE transactions on geoscience and remote sensing*,  
1799 **44**(10), p. 2642–2653, doi:10.1109/TGRS.2006.882801.
- 1800 Magagi, R., A. A. Berg, K. Goïta, S. Bélair, T. J. Jackson, B. Toth, A. Walker, H. McNairn,  
1801 P. E. O’Neill, M. Moghaddam, et al. (2013), Canadian experiment for soil moisture in 2010  
1802 (CanEx-SM10): Overview and preliminary results, *IEEE Transactions on Geoscience and*  
1803 *Remote Sensing*, **51**(1), p. 347–363, doi:10.1109/TGRS.2012.2198920.
- 1804 McColl, K. A., J. Vogelzang, A. G. Konings, D. Entekhabi, M. Piles, and A. Stoffelen  
1805 (2014), Extended triple collocation: Estimating errors and correlation coefficients with  
1806 respect to an unknown target, *Geophysical Research Letters*, **41**(17), p. 6229–6236, doi:  
1807 10.1002/2014GL061322.
- 1808 McColl, K. A., A. Roy, C. Derksen, A. G. Konings, S. H. Alemohammed, and D. En-  
1809 tekhabi (2016), Triple collocation for binary and categorical variables: Application to val-  
1810 idating landscape freeze/thaw retrievals, *Remote Sensing of Environment*, **176**, p. 31–42,  
1811 doi:10.1016/j.rse.2016.01.010.
- 1812 McNairn, H., T. J. Jackson, G. Wiseman, S. Belair, A. Berg, P. Bullock, A. Colliander, M. H.  
1813 Cosh, S.-B. Kim, R. Magagi, et al. (2015), The soil moisture active passive validation experi-  
1814 ment 2012 (SMAPVEX12): Prelaunch calibration and validation of the SMAP soil moisture

1815 algorithms, *IEEE Transactions on Geoscience and Remote Sensing*, **53**(5), p. 2784–2801, doi:  
1816 10.1109/TGRS.2014.2364913.

1817 Merchant, C. J., F. Paul, T. Popp, M. Ablain, S. Bontemps, P. Defourny, R. Hollmann,  
1818 T. Lavergne, A. Laeng, G. d. Leeuw, et al. (2017), Uncertainty information in climate  
1819 data records from earth observation, *Earth System Science Data*, **9**(2), p. 511–527, doi:  
1820 10.5194/essd-9-511-2017.

1821 Miralles, D. G., W. T. Crow, and M. H. Cosh (2010), Estimating spatial sampling errors in  
1822 coarse-scale soil moisture estimates derived from point-scale observations, *J. Hydrometeor.*,  
1823 **11**(6), p. 1423–1429, doi:10.1175/2010JHM1285.1.

1824 Miyaoka, K., A. Gruber, F. Ticconi, S. Hahn, W. Wagner, J. Figa-Saldana, and C. Anderson  
1825 (2017), Triple collocation analysis of soil moisture from Metop-A ASCAT and SMOS against  
1826 JRA-55 and ERA-Interim, *IEEE Journal of Selected Topics in Applied Earth Observations*  
1827 *and Remote Sensing*, **10**(5), p. 2274–2284, doi:10.1109/JSTARS.2016.2632306.

1828 Moghaddam, M., D. Entekhabi, Y. Goykhman, K. Li, M. Liu, A. Mahajan, A. Nayyar,  
1829 D. Shuman, and D. Teneketzis (2010), A wireless soil moisture smart sensor web using  
1830 physics-based optimal control: Concept and initial demonstrations, *IEEE Journal of Se-*  
1831 *lected Topics in Applied Earth Observations and Remote Sensing*, **3**(4), p. 522–535, doi:  
1832 10.1109/JSTARS.2010.2052918.

1833 Molero, B., D. Leroux, P. Richaume, Y. Kerr, O. Merlin, M. Cosh, and R. Bindlish (2018),  
1834 Multi-timescale analysis of the spatial representativeness of in situ soil moisture data within  
1835 satellite footprints, *Journal of Geophysical Research: Atmospheres*, **123**(1), p. 3–21, doi:10.  
1836 1002/2017JD027478.

1837 Naeimi, V., K. Scipal, Z. Bartalis, S. Hasenauer, and W. Wagner (2009), An improved soil  
1838 moisture retrieval algorithm for ERS and METOP scatterometer observations, *Geoscience*  
1839 *and Remote Sensing, IEEE Transactions on*, **47**(7), p. 1999–2013, doi:10.1109/TGRS.2008.  
1840 2011617.

1841 Naeimi, V., C. Paulik, A. Bartsch, W. Wagner, R. Kidd, S.-E. Park, K. Elger, and J. Boike  
1842 (2012), ASCAT surface state flag (SSF): Extracting information on surface freeze/thaw con-  
1843 ditions from backscatter data using an empirical threshold-analysis algorithm, *Geoscience*

- 1844 *and Remote Sensing, IEEE Transactions on*, **50**(7), p. 2566–2582, doi:10.1109/TGRS.2011.  
1845 2177667.
- 1846 Narapusetty, B., T. DelSole, and M. K. Tippett (2009), Optimal estimation of the climatological  
1847 mean, *Journal of Climate*, **22**(18), p. 4845–4859, doi:10.1175/2009JCLI2944.1.
- 1848 Neyman, J. (1937), X—outline of a theory of statistical estimation based on the classical theory  
1849 of probability, *Philosophical Transactions of the Royal Society of London. Series A, Mathe-*  
1850 *matical and Physical Sciences*, **236**(767), p. 333–380, doi:10.1098/rsta.1937.0005.
- 1851 Nicolai-Shaw, N., M. Hirschi, H. Mittelbach, and S. I. Seneviratne (2015), Spatial representa-  
1852 tiveness of soil moisture using in situ, remote sensing, and land reanalysis data, *Journal of*  
1853 *Geophysical Research: Atmospheres*, **120**(19), p. 9955–9964, doi:10.1002/2015JD023305.
- 1854 Noilhan, J., P. Lacarrère, and P. Bougeault (1991), An experiment with an advanced surface pa-  
1855 rameterization in a mesobeta-scale model. part III: Comparison with the HAPEX-MOBILHY  
1856 dataset, *Monthly weather review*, **119**(10), p. 2393–2413, doi:10.1175/1520-0493(1991)  
1857 119(2393:AEWAAS)2.0.CO;2.
- 1858 Ochsner, T. E., M. H. Cosh, R. H. Cuenca, W. A. Dorigo, C. S. Draper, Y. Hagimoto, Y. H.  
1859 Kerr, E. G. Njoku, E. E. Small, M. Zreda, et al. (2013), State of the art in large-scale  
1860 soil moisture monitoring, *Soil Science Society of America Journal*, **77**(6), p. 1888–1919, doi:  
1861 10.2136/sssaj2013.03.0093.
- 1862 Ólafsdóttir, K., and M. Mudelsee (2014), More accurate, calibrated bootstrap confidence inter-  
1863 vals for estimating the correlation between two time series, *Mathematical Geosciences*, **46**(4),  
1864 p. 411–427, doi:10.1007/s11004-014-9523-4.
- 1865 O’Neill, P., S. Chan, E. Njoku, T. Jackson, and R. Bindlish (2012), SMAP level 2 & 3 soil  
1866 moisture (passive) algorithm theoretical basis document (ATBD), *Initial Release, version*, **1**.
- 1867 O’Neill, P., S. Chan, E. Njoku, T. Jackson, and R. Bindlish (2018), SMAP L2 radiometer half-  
1868 orbit 36 km EASE-grid soil moisture, version 5, *Boulder, Colorado USA. NASA National*  
1869 *Snow and ice Data Center Distributed Active Archive Center*, doi:https://doi.org/10.5067/  
1870 SODMLCE6LGLL.
- 1871 Pan, M., C. K. Fisher, N. W. Chaney, W. Zhan, W. T. Crow, F. Aires, D. Entekhabi, and E. F.  
1872 Wood (2015), Triple collocation: Beyond three estimates and separation of structural/non-

1873 structural errors, *Remote Sensing of Environment*, **171**, p. 299–310, doi:doi.org/10.1016/j.rse.  
1874 2015.10.028.

1875 Panciera, R., J. P. Walker, J. D. Kalma, E. J. Kim, J. M. Hacker, O. Merlin, M. Berger, and  
1876 N. Skou (2008), The NAFE'05/CoSMOS data set: Toward SMOS soil moisture retrieval,  
1877 downscaling, and assimilation, *IEEE Transactions on Geoscience and Remote Sensing*, **46**(3),  
1878 p. 736–745, doi:10.1109/TGRS.2007.915403.

1879 Parinussa, R. M., A. G. Meesters, Y. Y. Liu, W. Dorigo, W. Wagner, and R. A. De Jeu (2011),  
1880 Error estimates for near-real-time satellite soil moisture as derived from the land parameter  
1881 retrieval model, *Geoscience and Remote Sensing Letters, IEEE*, **8**(4), p. 779–783, doi:10.1109/  
1882 LGRS.2011.2114872.

1883 Parinussa, R. M., T. R. Holmes, N. Wanders, W. A. Dorigo, and R. A. de Jeu (2015), A  
1884 preliminary study toward consistent soil moisture from AMSR2, *Journal of Hydrometeorology*,  
1885 **16**(2), p. 932–947, doi:10.1175/JHM-D-13-0200.1.

1886 Pathe, C., W. Wagner, D. Sabel, M. Doubkova, and J. B. Basara (2009), Using envisat asar  
1887 global mode data for surface soil moisture retrieval over oklahoma, usa, *IEEE Transactions*  
1888 *on Geoscience and Remote Sensing*, **47**(2), p. 468–480, doi:10.1109/TGRS.2008.2004711.

1889 Peischl, S., J. P. Walker, C. Rüdiger, N. Ye, Y. H. Kerr, E. Kim, R. Bandara, and M. Al-  
1890 lahmoradi (2012), The AACES field experiments: SMOS calibration and validation across  
1891 the murrumbidgee river catchment., *Hydrology & Earth System Sciences Discussions*, **9**(3),  
1892 doi:10.5194/hessd-9-2763-2012.

1893 Peng, J., A. Loew, S. Zhang, J. Wang, and J. Niesel (2015), Spatial downscaling of satellite  
1894 soil moisture data using a vegetation temperature condition index, *IEEE Transactions on*  
1895 *Geoscience and Remote Sensing*, **54**(1), p. 558–566, doi:10.1109/TGRS.2015.2462074.

1896 Peng, J., A. Loew, O. Merlin, and N. E. Verhoest (2017), A review of spatial downscaling  
1897 of satellite remotely sensed soil moisture, *Reviews of Geophysics*, **55**(2), p. 341–366, doi:  
1898 10.1002/2016RG000543.

1899 Pierdicca, N., F. Fascetti, L. Pulvirenti, and R. Crapolicchio (2017), Error characterization of  
1900 soil moisture satellite products: Retrieving error cross-correlation through extended quadru-

1901 ple collocation, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote*  
1902 *Sensing*, **10**(10), p. 4522–4530, doi:10.1109/JSTARS.2017.2714025.

1903 QA4EO (2010), *A Quality Assurance Framework for Earth Observation: Principles*, version 4.0  
1904 ed.

1905 Reichle, R. H., and R. D. Koster (2004), Bias reduction in short records of satellite soil moisture,  
1906 *Geophys. Res. Lett.*, **31**(19), p. L19,501, doi:10.1029/2004GL020938.

1907 Reichle, R. H., R. D. Koster, G. J. De Lannoy, B. A. Forman, Q. Liu, S. P. Mahanama, and  
1908 A. Touré (2011), Assessment and enhancement of MERRA land surface hydrology estimates,  
1909 *Journal of climate*, **24**(24), p. 6322–6338, doi:10.1175/JCLI-D-10-05033.1.

1910 Reichle, R. H., C. S. Draper, Q. Liu, M. Girotto, S. P. Mahanama, R. D. Koster, and G. J.  
1911 De Lannoy (2017a), Assessment of MERRA-2 land surface hydrology estimates, *Journal of*  
1912 *Climate*, **30**(8), p. 2937–2960, doi:10.1175/JCLI-D-16-0720.1.

1913 Reichle, R. H., G. J. De Lannoy, Q. Liu, R. D. Koster, J. S. Kimball, W. T. Crow, J. V.  
1914 Ardizzone, P. Chakraborty, D. W. Collins, A. L. Conaty, et al. (2017b), Global assessment of  
1915 the SMAP level-4 surface and root-zone soil moisture product using assimilation diagnostics,  
1916 *Journal of Hydrometeorology*, **18**(12), p. 3217–3237, doi:10.1175/JHM-D-17-0130.1.

1917 Reichle, R. H., G. J. De Lannoy, Q. Liu, J. V. Ardizzone, A. Colliander, A. Conaty, W. Crow,  
1918 T. J. Jackson, L. A. Jones, J. S. Kimball, et al. (2017c), Assessment of the SMAP level-4  
1919 surface and root-zone soil moisture product using in situ measurements, *Journal of hydrome-*  
1920 *teorology*, **18**(10), p. 2621–2645, doi:10.1175/JHM-D-17-0063.1.

1921 Rodell, M., P. Houser, U. e. a. Jambor, J. Gottschalck, K. Mitchell, C. Meng, K. Arse-  
1922 nault, B. Cosgrove, J. Radakovich, M. Bosilovich, et al. (2004), The global land data as-  
1923 simulation system, *Bulletin of the American Meteorological Society*, **85**(3), p. 381–394, doi:  
1924 10.1175/BAMS-85-3-381.

1925 Rüdiger, C., A. W. Western, J. P. Walker, A. B. Smith, J. D. Kalma, and G. R. Willgoose (2010),  
1926 Towards a general equation for frequency domain reflectometers, *Journal of hydrology*, **383**(3-  
1927 4), p. 319–329, doi:10.1016/j.jhydrol.2009.12.046.

1928 Sabaghy, S., J. Walker, L. Renzullo, R. Akbar, S. Chan, J. Chaubell, N. Das, R. Dunbar,  
1929 D. Entekhabi, A. Gevaert, T. Jackson, A. Loew, O. Merlin, M. Moghaddam, J. Peng, J. Peng,

1930 J. Piepmeier, C. Rüdiger, V. Stefan, X. Wu, N. Ye, and S. Yueh (in review), Comprehensive  
1931 analysis of alternative downscaled soil moisture products, *Remote Sensing of Environment*.

1932 Sahoo, A. K., G. J. D. Lannoy, R. H. Reichle, and P. R. Houser (2013), Assimilation and  
1933 downscaling of satellite observed soil moisture over the little river experimental watershed in  
1934 georgia, usa, *Advances in Water Resources*, **52**, p. 19 – 33, doi:10.1016/j.advwatres.2012.08.  
1935 007.

1936 Scanlon, T., J. Nightingale, F. Boersma, J.-P. Muller, C. Farquhar, S. Compernelle, and J.-  
1937 C. Lambert (2017), Outline of QA4ECV quality assurance service (version 2.0), *Tech. rep.*,  
1938 QA4ECV, <http://www.qa4ecv.eu/qa-system>, last access: 1 July 2019.

1939 Scipal, K., M. Drusch, and W. Wagner (2008a), Assimilation of a ERS scatterometer derived  
1940 soil moisture index in the ECMWF numerical weather prediction system, *Advances in water  
1941 resources*, **31**(8), p. 1101–1112, doi:10.1016/j.advwatres.2008.04.013.

1942 Scipal, K., T. Holmes, R. de Jeu, V. Naeimi, and W. Wagner (2008b), A possible solution for  
1943 the problem of estimating the error structure of global soil moisture data sets, *Geophys. Res.  
1944 Lett.*, **35**(24), p. L24,403, doi:10.1029/2008GL035599.

1945 Starks, P. J., G. C. Heathman, T. J. Jackson, and M. H. Cosh (2006), Temporal stability of soil  
1946 moisture profile, *Journal of Hydrology*, **324**, p. 400–411, doi:10.1016/j.jhydrol.2005.09.024.

1947 Stoffelen, A. (1998), Toward the true near-surface wind speed: Error modeling and calibration  
1948 using triple collocation, *J. Geophys. Res.*, **103**(C4), p. 7755–7766, doi:10.1029/97JC03180.

1949 Su, C.-H., and D. Ryu (2015), Multi-scale analysis of bias correction of soil moisture, *Hydrology  
1950 and Earth System Sciences*, **19**(1), p. 17–31, doi:10.5194/hess-19-17-2015.

1951 Su, C.-H., D. Ryu, W. T. Crow, and A. W. Western (2014), Beyond triple collocation: Appli-  
1952 cations to soil moisture monitoring, *Journal of Geophysical Research: Atmospheres*, **119**(11),  
1953 p. 6419–6439, doi:10.1002/2013JD021043.

1954 Su, Z., W. Timmermans, Y. Zeng, J. Schulz, V. O. John, R. A. Roebeling, P. Poli, D. Tan,  
1955 F. Kaspar, A. K. Kaiser-Weiss, E. Swinnen, C. Toté, H. Gregow, T. Manninen, A. Riihelä,  
1956 J.-C. Calvet, Y. Ma, and J. Wen (2018), An overview of european efforts in generating climate  
1957 data records, *Bulletin of the American Meteorological Society*, **99**(2), p. 349–359, doi:10.1175/  
1958 BAMS-D-16-0074.1.



1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987

Tong, C. (2019), Statistical inference enables bad science; statistical thinking enables good science, *The American Statistician*, **73**(sup1), p. 246–261, doi:10.1080/00031305.2018.1518264.

Ulaby, F. T., D. G. Long, W. J. Blackwell, C. Elachi, A. K. Fung, C. Ruf, K. Sarabandi, H. A. Zebker, and J. Van Zyl (2014), *Microwave radar and radiometric remote sensing*, vol. 4, University of Michigan Press Ann Arbor.

Vachaud, G., A. Passerat De Silans, P. Balabanis, and M. Vauclin (1985), Temporal stability of spatially measured soil water probability density function, *Soil Sci. Soc. Am. J.*, **49**(4), p. 822–828, doi:10.2136/sssaj1985.03615995004900040006x.

van der Schalie, R., R. de Jeu, R. Parinussa, N. Rodríguez-Fernández, Y. Kerr, A. Al-Yaari, J.-P. Wigneron, and M. Drusch (2018), The effect of three different data fusion approaches on the quality of soil moisture retrievals from multiple passive microwave sensors, *Remote Sensing*, **10**(1), p. 107, doi:10.3390/rs10010107.

Van Leeuwen, P. J. (2015), Representation errors and retrievals in linear and nonlinear data assimilation, *Quarterly Journal of the Royal Meteorological Society*, **141**(690), p. 1612–1623.

Vogelzang, J., and A. Stoffelen (2012), Triple collocation, *EUMETSAT Report*. Available at [http://research.metoffice.gov.uk/research/interproj/nwpsaf/scatterometer/TripleCollocation\\_NWPSAF\\_TR.L](http://research.metoffice.gov.uk/research/interproj/nwpsaf/scatterometer/TripleCollocation_NWPSAF_TR.L) last access: 1 July 2019.

Wagner, W., G. Lemoine, and H. Rott (1999), A method for estimating soil moisture from ERS scatterometer and soil data, *Remote Sensing of Environment*, **70**(2), p. 191–207, doi:10.1016/S0034-4257(99)00036-X.

Wagner, W., L. Brocca, V. Naeimi, R. Reichle, C. Draper, R. de Jeu, D. Ryu, C.-H. Su, A. Western, J.-C. Calvet, et al. (2014), Clarifications on the “comparison between SMOS, VUA, ASCAT, and ECMWF soil moisture products over four watersheds in US”, *IEEE Transactions on Geoscience and Remote Sensing*, **52**(3), p. 1901–1906, doi:10.1109/TGRS.2013.2282172.

Walker, J. P., G. R. Willgoose, and J. D. Kalma (2004), In situ measurement of soil moisture: a comparison of techniques, *Journal of Hydrology*, **293**, p. 85–99, doi:10.1016/j.jhydrol.2004.01.008.

Wang, G., D. Garcia, Y. Liu, R. De Jeu, and A. J. Dolman (2012), A three-dimensional gap filling method for large geophysical datasets: Application to global satellite soil moisture

1988 observations, *Environmental Modelling & Software*, **30**, p. 139–142, doi:10.1016/j.envsoft.  
1989 2011.10.015.

1990 Wasserstein, R. L., and N. A. Lazar (2016), The ASA’s statement on p-values: context, pro-  
1991 cess, and purpose, *The American Statistician*, **70**(2), p. 129–133, doi:10.1080/00031305.2016.  
1992 1154108.

1993 Wasserstein, R. L., A. L. Schirm, and N. A. Lazar (2019), Moving to a world beyond “p<0.05”,  
1994 *The American Statistician*, **73**(sup1), p. 1–19, doi:10.1080/00031305.2019.1583913.

1995 Wilks, D. S. (2011), *Statistical Methods in the Atmospheric Sciences*, vol. 100, 3rd ed., Academic  
1996 Press.

1997 WMO (2016), The global observing system for climate: Implementation needs, *Implementation*  
1998 *Plan GCOS-200*, World Meteorological Organization.

1999 Yilmaz, M. T., and W. T. Crow (2013), The optimality of potential rescaling approaches in land  
2000 data assimilation., *Journal of Hydrometeorology*, **14**(2), doi:10.1175/JHM-D-12-052.1.

2001 Yilmaz, M. T., and W. T. Crow (2014), Evaluation of assumptions in soil moisture triple collocation  
2002 analysis, *Journal of Hydrometeorology*, **15**(3), p. 1293–1302, doi:10.1175/JHM-D-13-0158.  
2003 1.

2004 Zeng, Y., Z. Su, J.-C. Calvet, T. Manninen, E. Swinnen, J. Schulz, R. Roebeling, P. Poli, D. Tan,  
2005 A. Riihelä, C.-M. Tanis, A.-N. Arslan, A. Obregon, A. Kaiser-Weiss, V. John, W. Timmer-  
2006 mans, J. Timmermans, F. Kaspar, H. Gregow, A.-L. Barbu, D. Fairbairn, E. Gelati, and  
2007 C. Meurey (2015), Analysis of current validation practices in europe for space-based climate  
2008 data records of essential climate variables, *International Journal of Applied Earth Observation*  
2009 *and Geoinformation*, **42**, p. 150 – 161, doi:https://doi.org/10.1016/j.jag.2015.06.006.

2010 Zreda, M., W. Shuttleworth, X. Zeng, C. Zweck, D. Desilets, T. Franz, and R. Rosolem (2012),  
2011 COSMOS: the cosmic-ray soil moisture observing system, *Hydrology and Earth System Sci-*  
2012 *ences*, **16**(11), p. 4079–4099, doi:10.5194/hess-16-4079-2012.

2013 Zribi, M., M. Pardé, J. Boutin, P. Fanise, D. Hauser, M. Dechambre, Y. Kerr, M. Leduc-  
2014 Leballeur, G. Reverdin, N. Skou, S. Søbjaerg, C. Albergel, J. C. Calvet, J. P. Wigneron,  
2015 E. Lopez-Baeza, A. Rius, and J. Tenerelli (2011), Carols: A new airborne l-band radiometer  
2016 for ocean surface and land observations, *Sensors*, **11**(1), p. 719–742, doi:10.3390/s110100719.

2017 Zwieback, S., K. Scipal, W. Dorigo, and W. Wagner (2012), Structural and statistical properties  
2018 of the collocation technique for error characterization, *Nonlin. Processes Geophys.*, **19**(1),  
2019 p. 69–80, doi:10.5194/npg-19-69-2012.

2020 Zwieback, S., A. Colliander, M. H. Cosh, J. Martínez-Fernández, H. McNairn, P. J. Starks,  
2021 M. Thibeault, and A. Berg (2018), Estimating time-dependent vegetation biases in the SMAP  
2022 soil moisture product, *Hydrology and Earth System Sciences*, **22**(8), p. 4473–4489, doi:10.  
2023 5194/hess-22-4473-2018.

Table 1: Validation stages as defined by CEOS (modified from <https://lpvs.gsfc.nasa.gov/>; last access: 1 July 2019).

Validation Stage	Definition
0	No validation. Product accuracy has not been assessed. Product considered beta.
1	Product accuracy is assessed from a small (typically <30) set of locations and time periods by comparison with in situ or other suitable reference data.
2	Product accuracy is estimated over a considerable set of locations and time periods by comparison with reference in situ or other suitable reference data. Spatial and temporal consistency of the product and consistency with similar products has been evaluated over globally representative locations and time periods. Results are published in the peer-reviewed literature.
3	Uncertainties in the product and its associated structure are well quantified from comparison with reference in situ or other suitable reference data. Uncertainties are characterized in a statistically rigorous way over multiple locations and time periods representing global conditions. Spatial and temporal consistency of the product and with similar products has been evaluated over globally representative locations and periods. Results are published in the peer-reviewed literature.
4	Validation results for stage 3 are systematically updated when new product versions are released and as the time-series expands.

Table 2: Summary of publicly available reference data sources commonly used for satellite soil moisture validation (links last accessed: 1 July 2019).

Name	Description	Reference
ISMN	Data hosting facility for sparse soil moisture networks	<a href="http://ismn.geo.tuwien.ac.at/">http://ismn.geo.tuwien.ac.at/</a> ( <i>Dorigo et al., 2011a,b</i> )
CVS	Openly available Core Validation Site (CVS) data that have been specifically processed for SMAP validation.	<a href="https://nsidc.org/data/nsidc-0712">https://nsidc.org/data/nsidc-0712</a>
GLDAS	NASA’s global modelling and data assimilation system	<a href="https://ldas.gsfc.nasa.gov/gldas/">https://ldas.gsfc.nasa.gov/gldas/</a>
MERRA	NASA’s global reanalysis data sets	<a href="https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/">https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/</a>
ERA	ECMWF’s global reanalysis data sets	<a href="https://www.ecmwf.int/en/forecasts/datasets/browse-reanalysis-datasets/">https://www.ecmwf.int/en/forecasts/datasets/browse-reanalysis-datasets/</a>

Table 3: Open-source software that can be used for satellite soil moisture validation (links last accessed: last access: 1 July 2019).

Name	Description	Language	Reference
	Source code used to produce validation examples in this publication in Appendix A	python, MATLAB	<a href="https://github.com/alexgruber/validation_good_practice/">https://github.com/alexgruber/validation_good_practice/</a>
pytesmo	Geospatial time series validation toolbox	python	<a href="https://doi.org/10.5281/zenodo.1215760/">https://doi.org/10.5281/zenodo.1215760/</a>
poets	Geospatial image resampling toolbox	python	<a href="https://pypi.org/project/poets/">https://pypi.org/project/poets/</a>

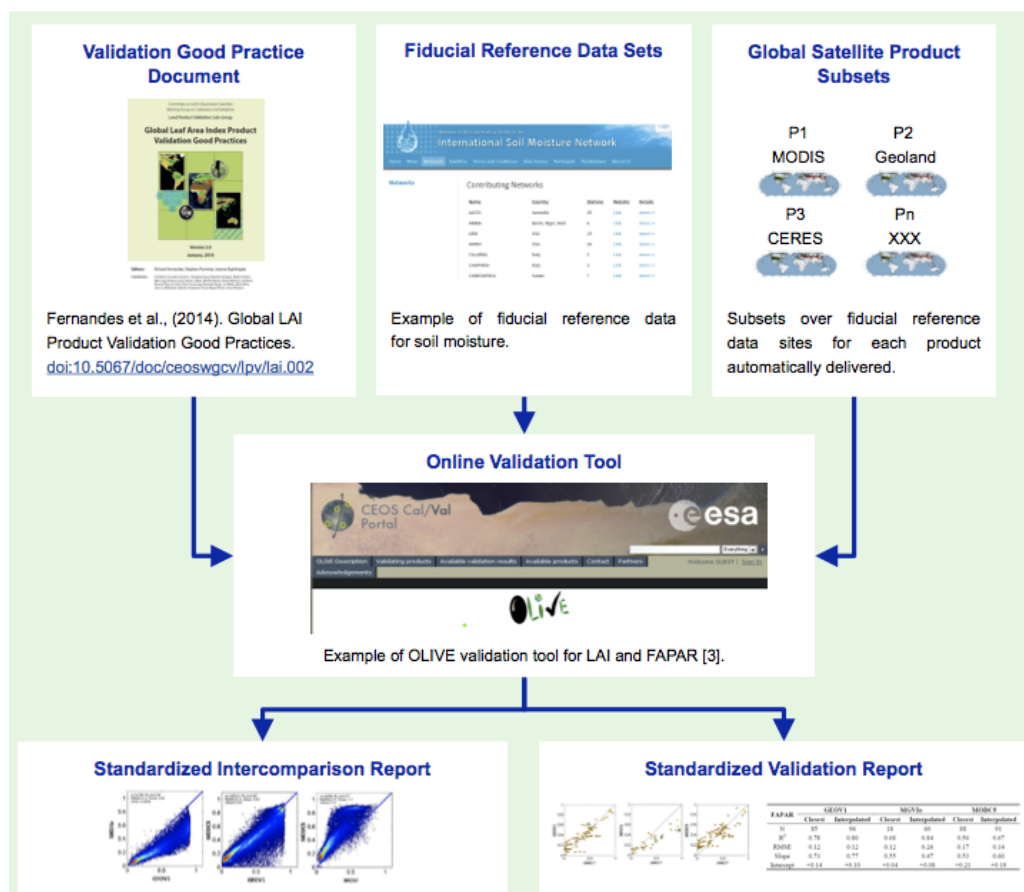


Figure 1: Validation framework as defined by CEOS (from <https://lpvs.gsfc.nasa.gov/>; last access: 1 July 2019).

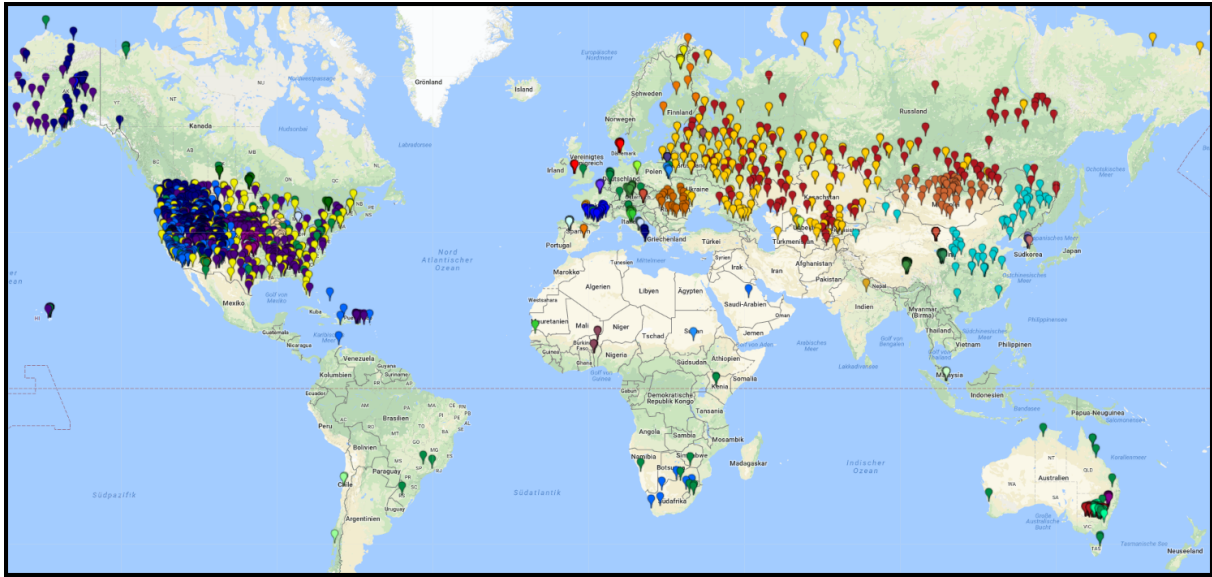


Figure 2: Currently available stations from sparse networks hosted by the ISMN. Colors represent different station hosting networks.

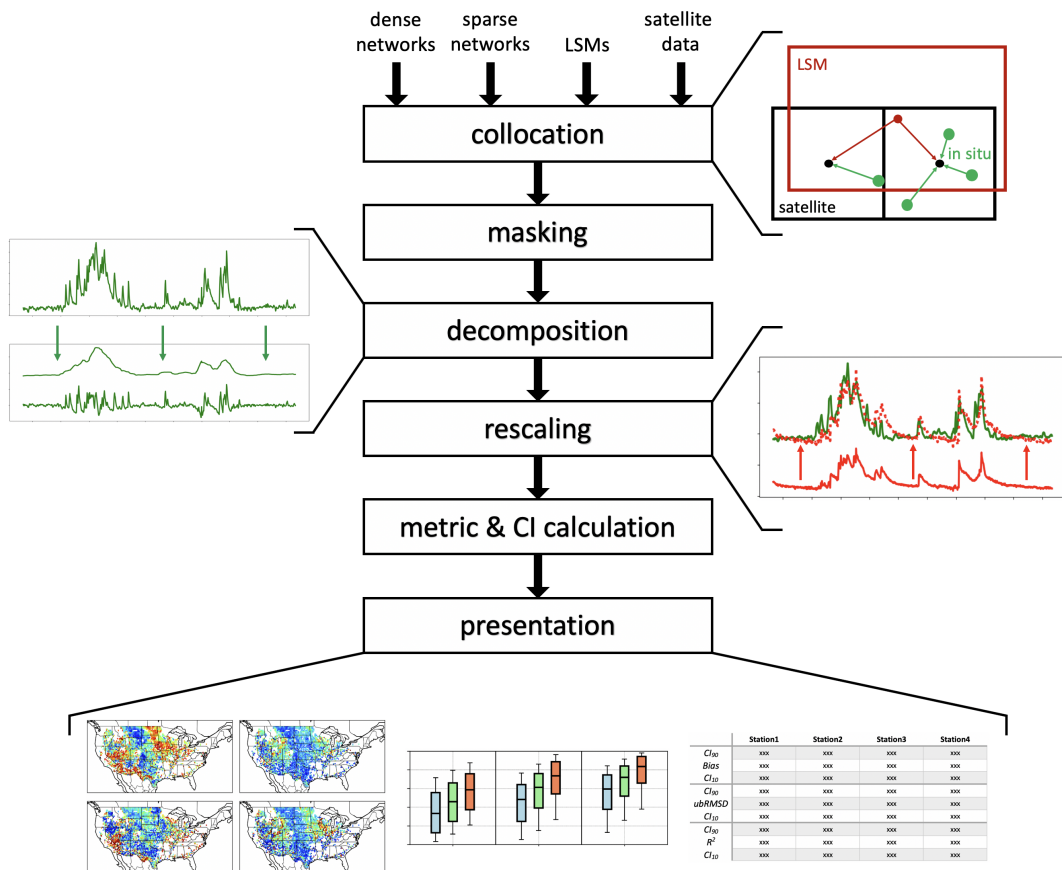


Figure 3: Validation good practice protocol illustration.

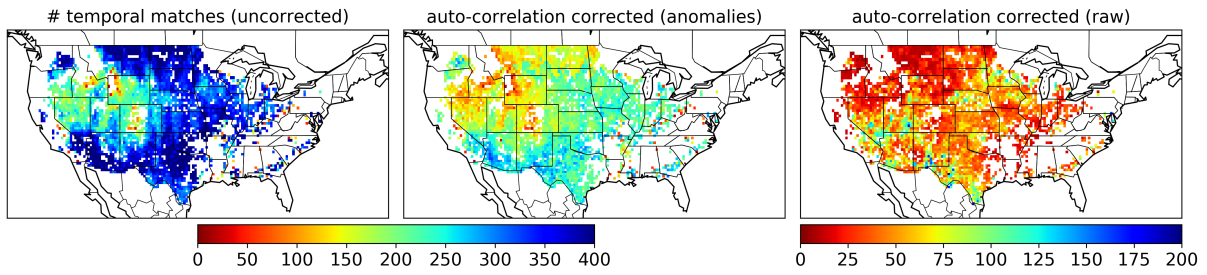


Figure A.1: Sample size for temporal matches between ASCAT, SMOS, SMAP and MERRA-2 between 2015 and 2018 (left), effective sample size when correcting for anomaly auto-correlation (middle), and effective sample size when correcting for auto-correlation in the raw time series (right).

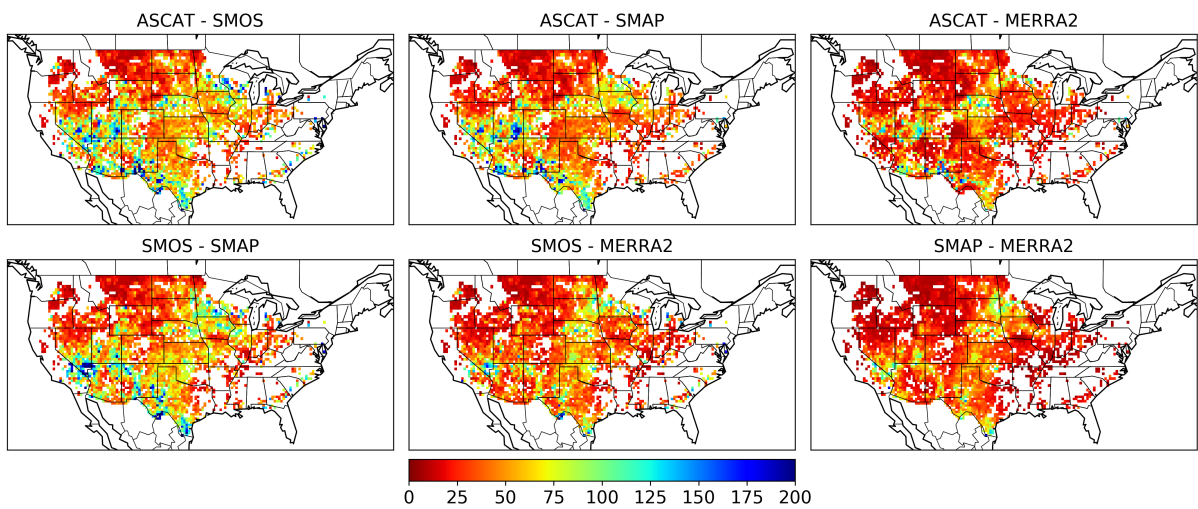


Figure A.2: Effective, raw time series auto-correlation corrected sample size for different data set combinations.

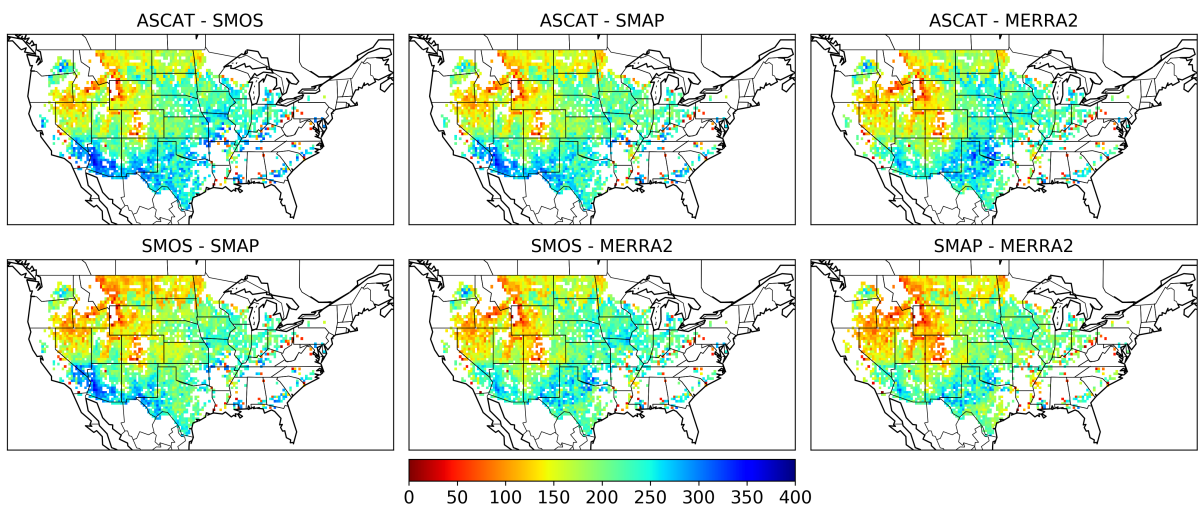


Figure A.3: Effective, anomaly auto-correlation corrected sample size for different data set combinations.

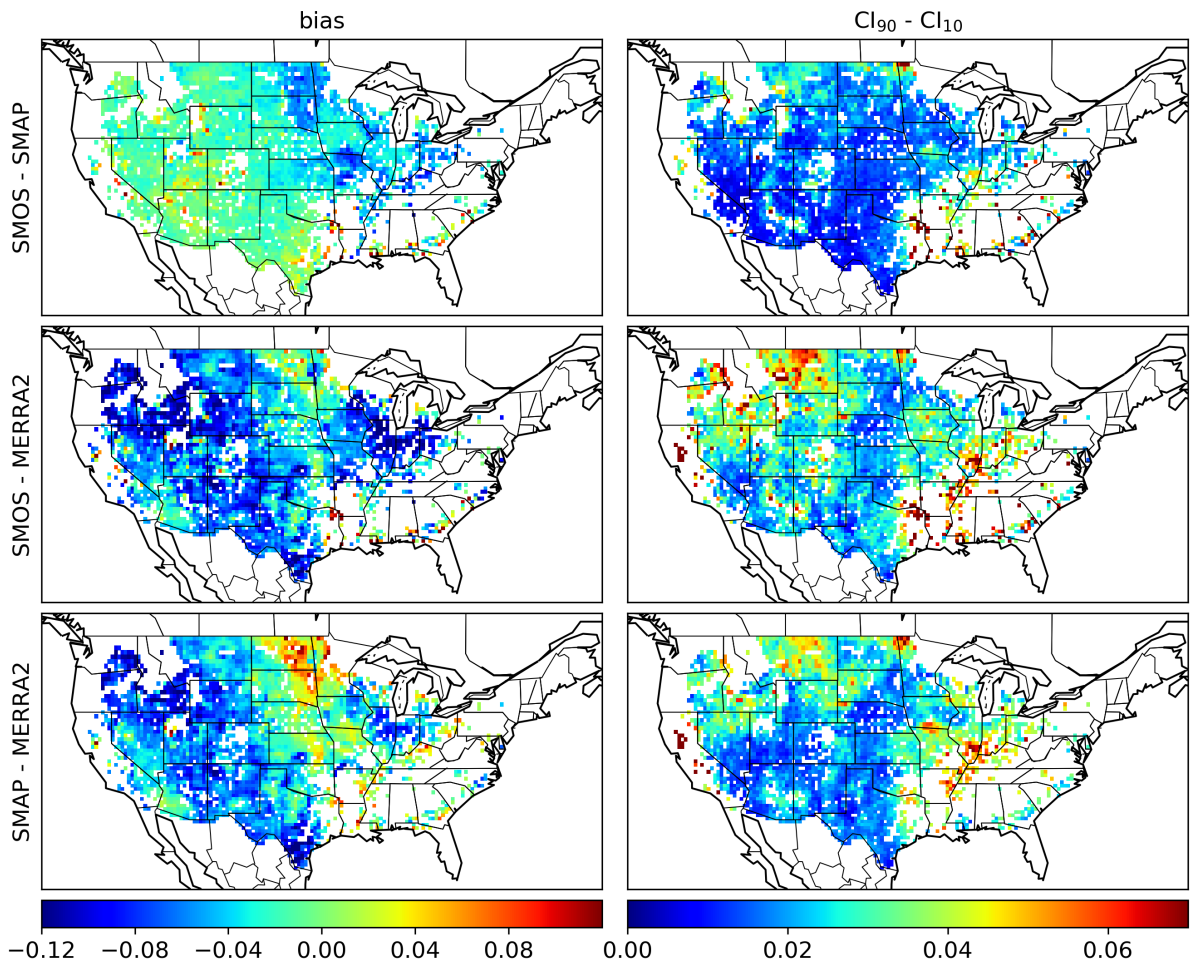


Figure A.4: Temporal mean biases [ $m^3m^{-3}$ ] (left) and associated 80% confidence intervals (right) between raw soil moisture estimates of SMOS, SMAP and MERRA-2.



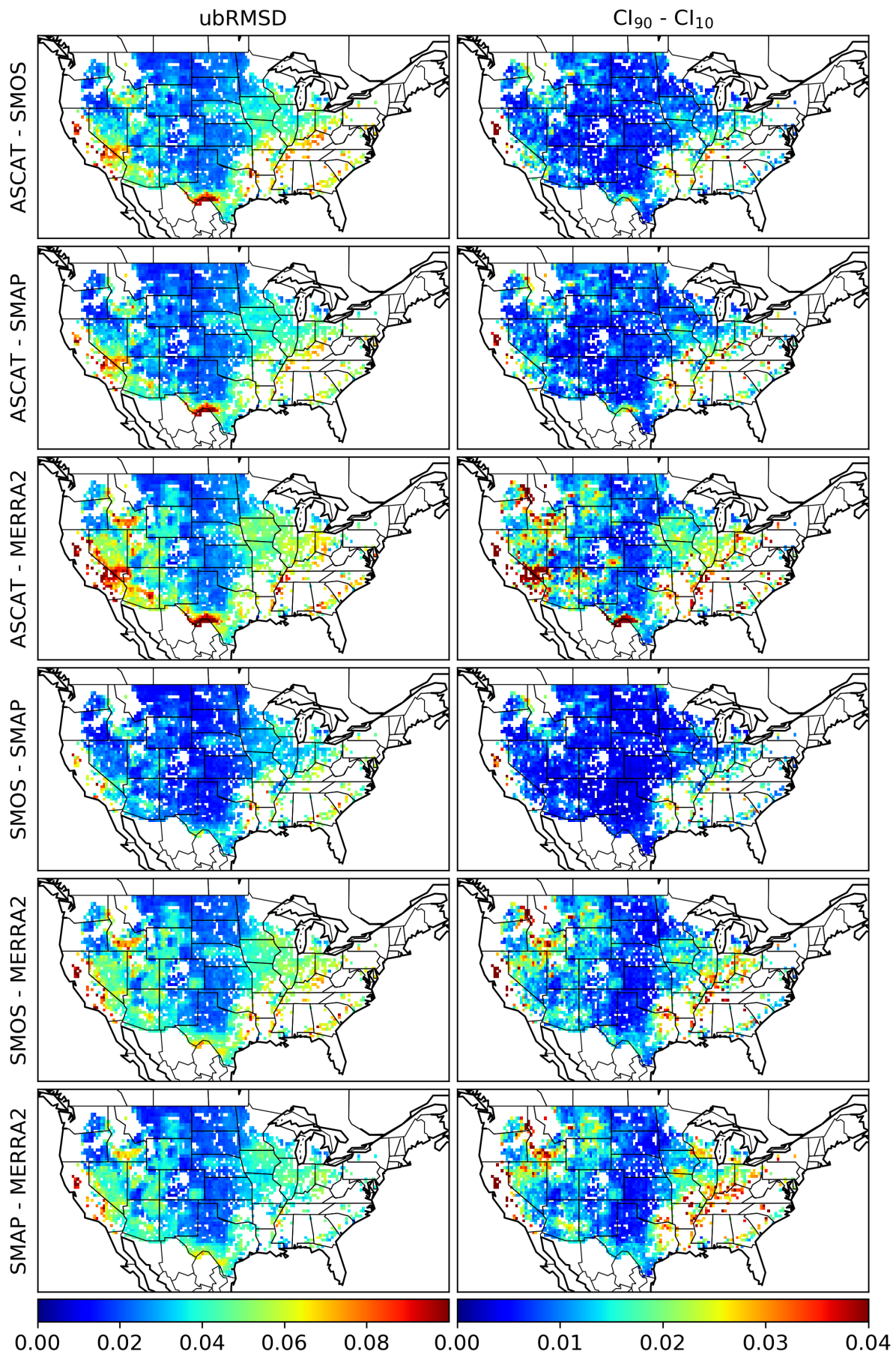


Figure A.5: Unbiased (in mean and standard deviation) root-mean-square-differences [ $m^3m^{-3}$ ] (left) and associated 80% confidence intervals (right) between raw soil moisture estimates of ASCAT, SMOS, SMAP and MERRA-2. 77

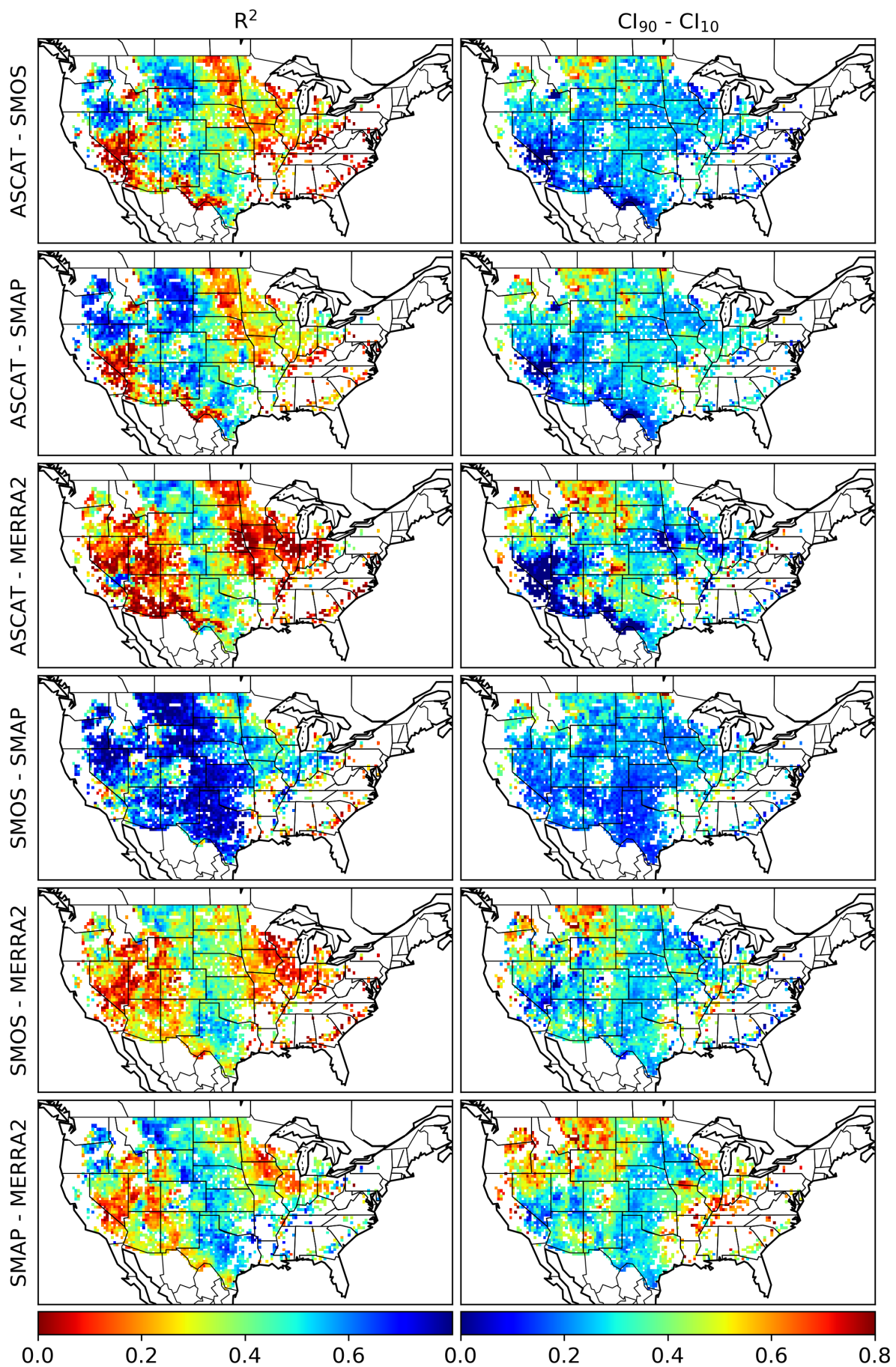


Figure A.6: Coefficients of determination [-] (left) and associated 80% confidence intervals (right) between raw soil moisture estimates of ASCAT, SMOS, SMAP and MERRA-2.

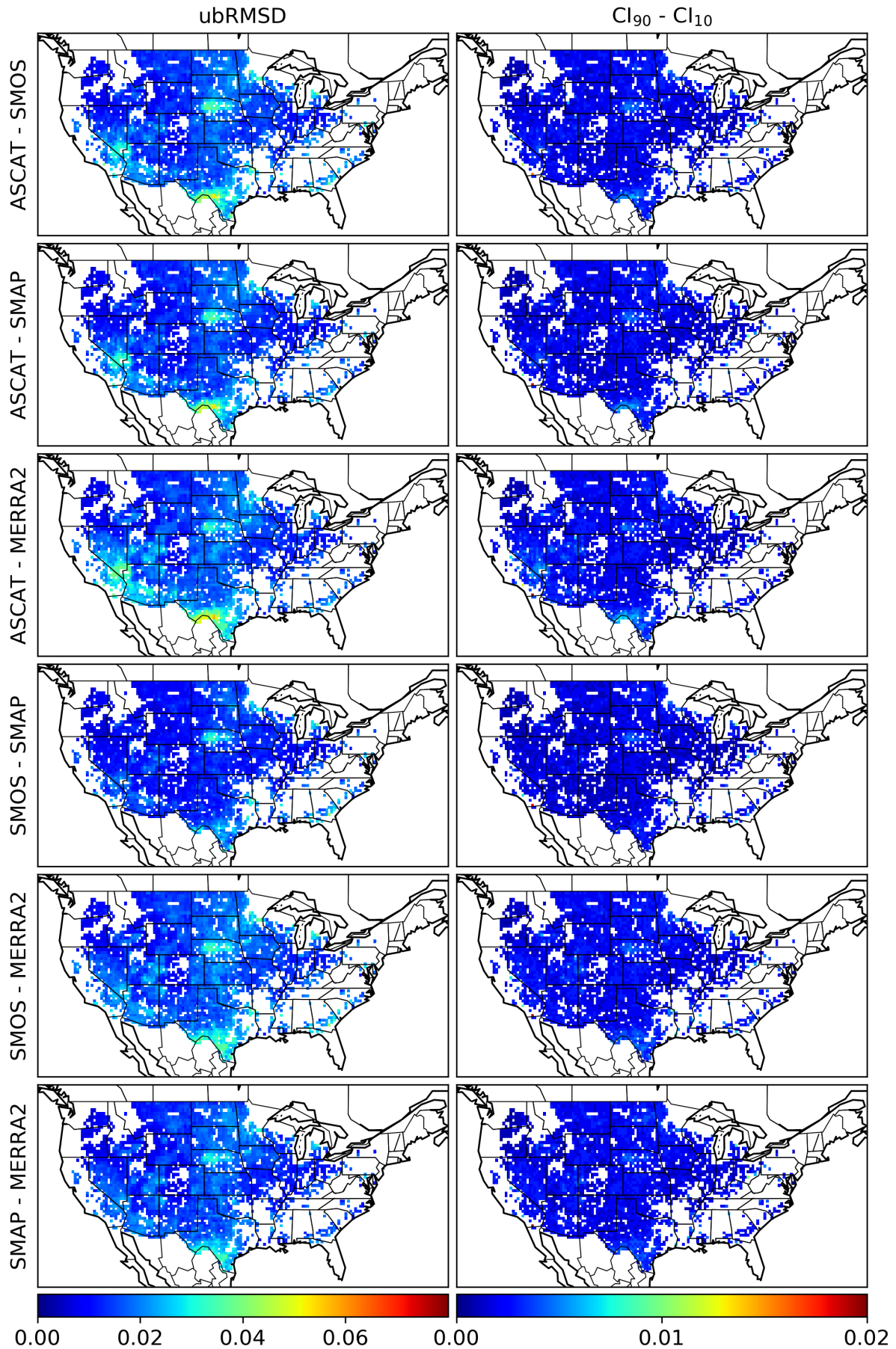


Figure A.7: Unbiased (in mean and standard deviation)  $[m^3m^{-3}]$  root-mean-square-differences (left) and associated 80% confidence intervals (right) between soil moisture anomaly estimates of ASCAT, SMOS, SMAP and MERRA-2. 79



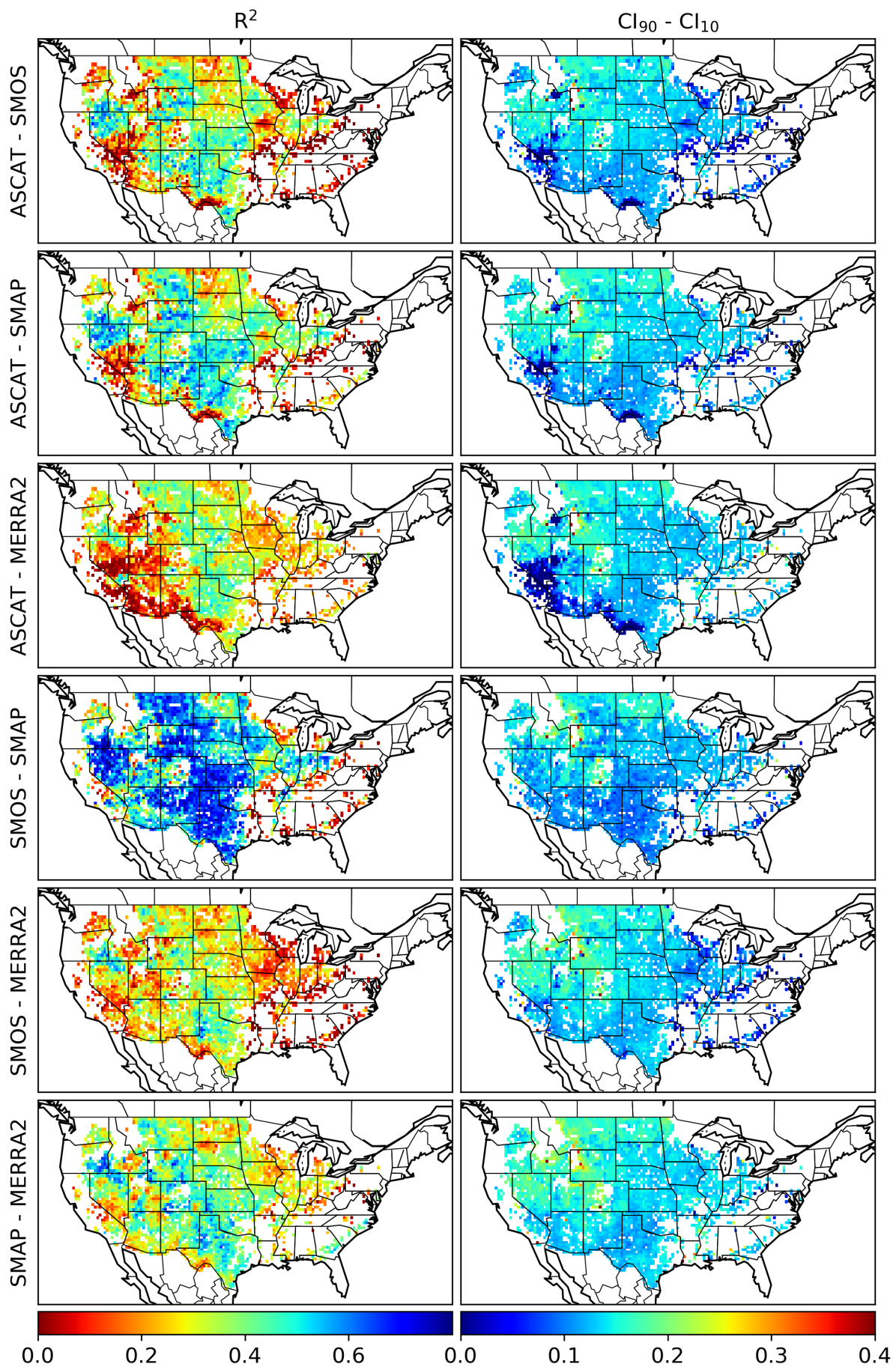


Figure A.8: Coefficients of determination [-] (left) and associated 80% confidence intervals (right) between soil moisture anomaly estimates of ASCAT, SMOS, SMAP and MERRA-2.



Figure A.9: Spatial summary statistics of biases [ $m^3m^{-3}$ ], ubRMSDs [ $m^3m^{-3}$ ], and coefficients of determination [-] and their 10% and 90% confidence limits, respectively, for raw soil moisture estimates and soil moisture anomalies of ASCAT, SMOS, SMAP and MERRA-2. Boxes represent the (spatial) median and inter-quartile-range and whiskers represent the 5 and 95 percentiles.

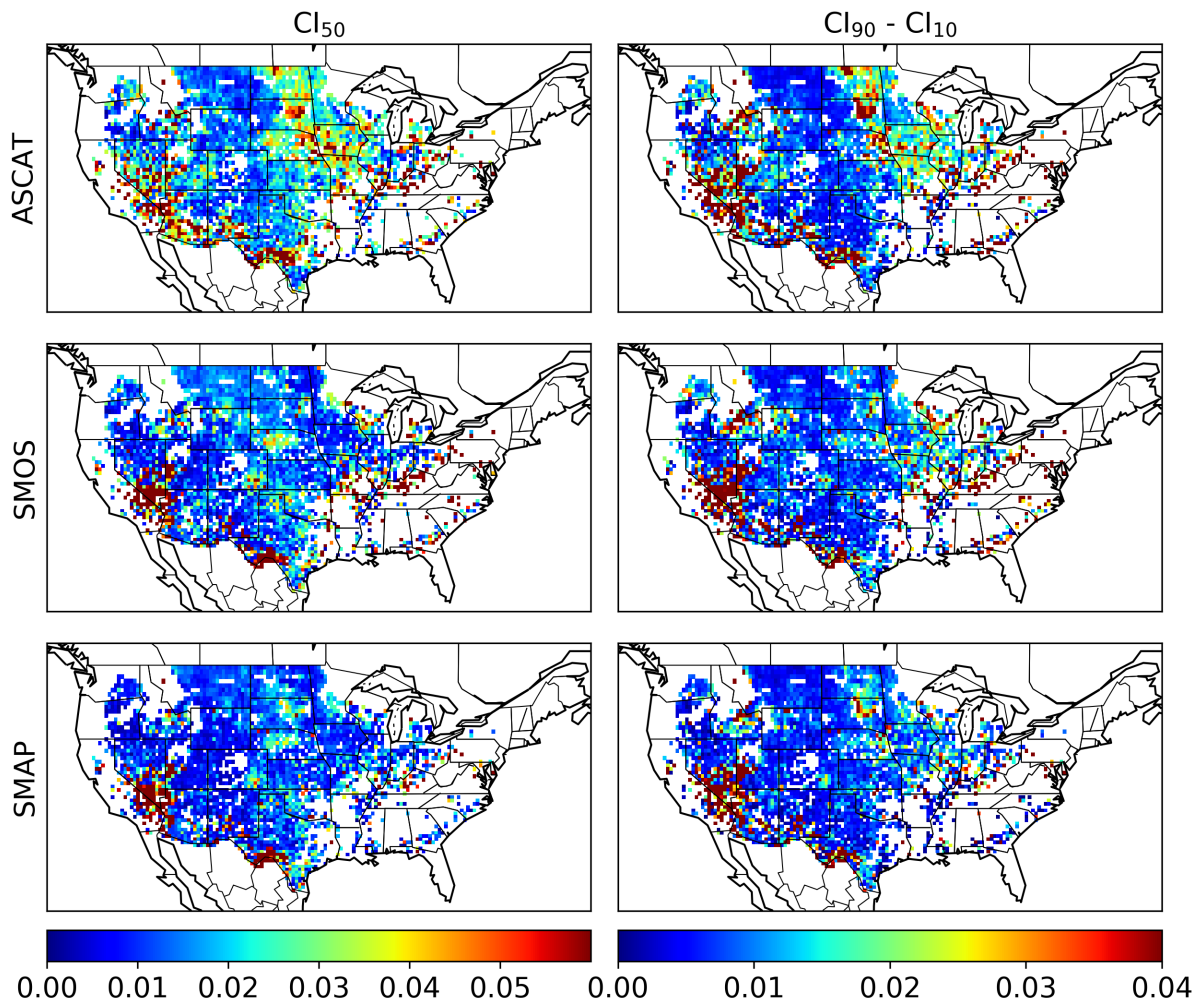


Figure A.10: Median of the bootstrapped TCA-based ubRMSEs [ $m^3m^{-3}$ ] (left) and associated 80% confidence intervals (right) of raw soil moisture estimates of ASCAT, SMOS, and SMAP.

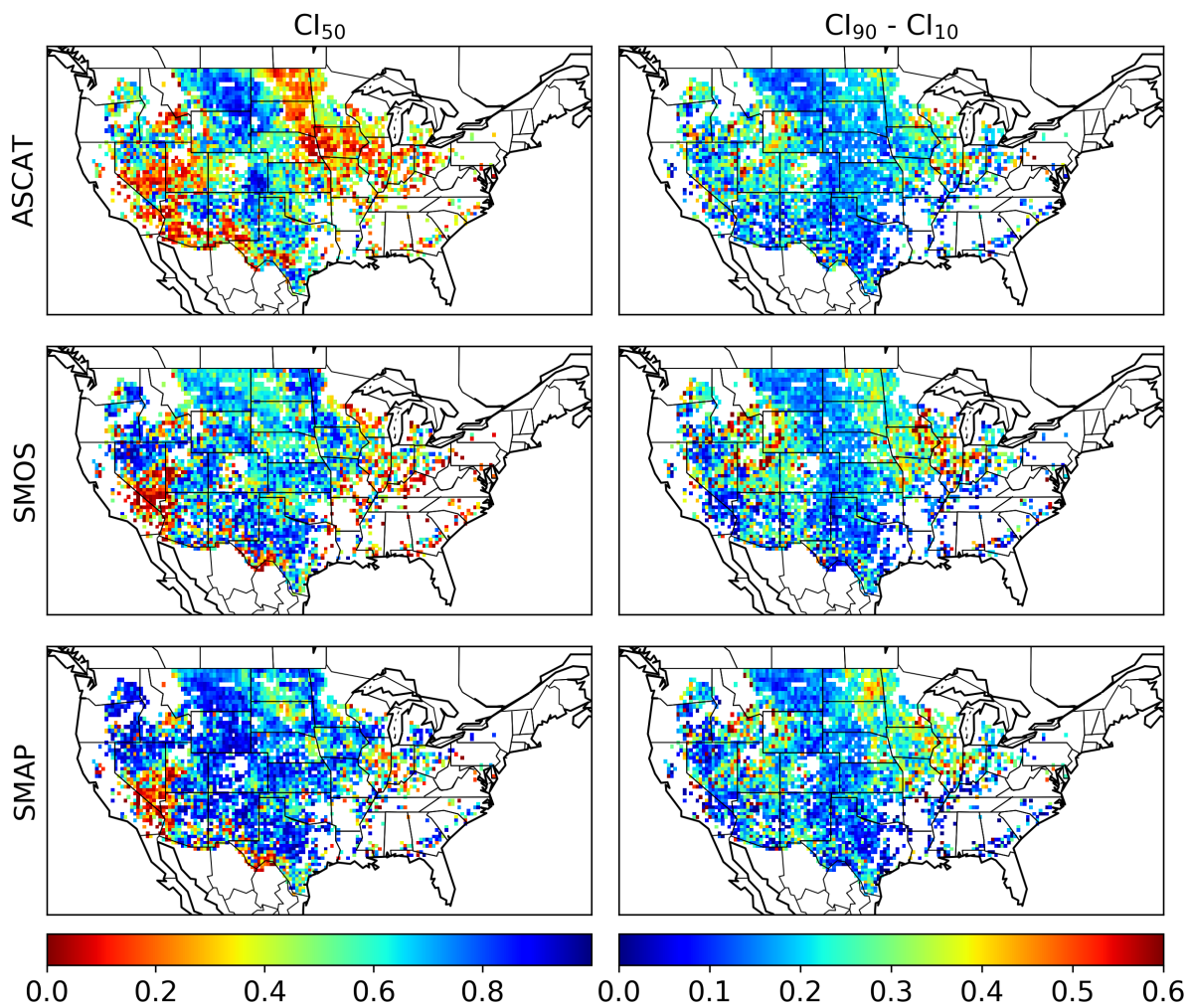


Figure A.11: Median of the bootstrapped TCA-based  $R^2$  estimates [-] (left) and associated 80% confidence intervals (right) of raw soil moisture estimates of ASCAT, SMOS, and SMAP.

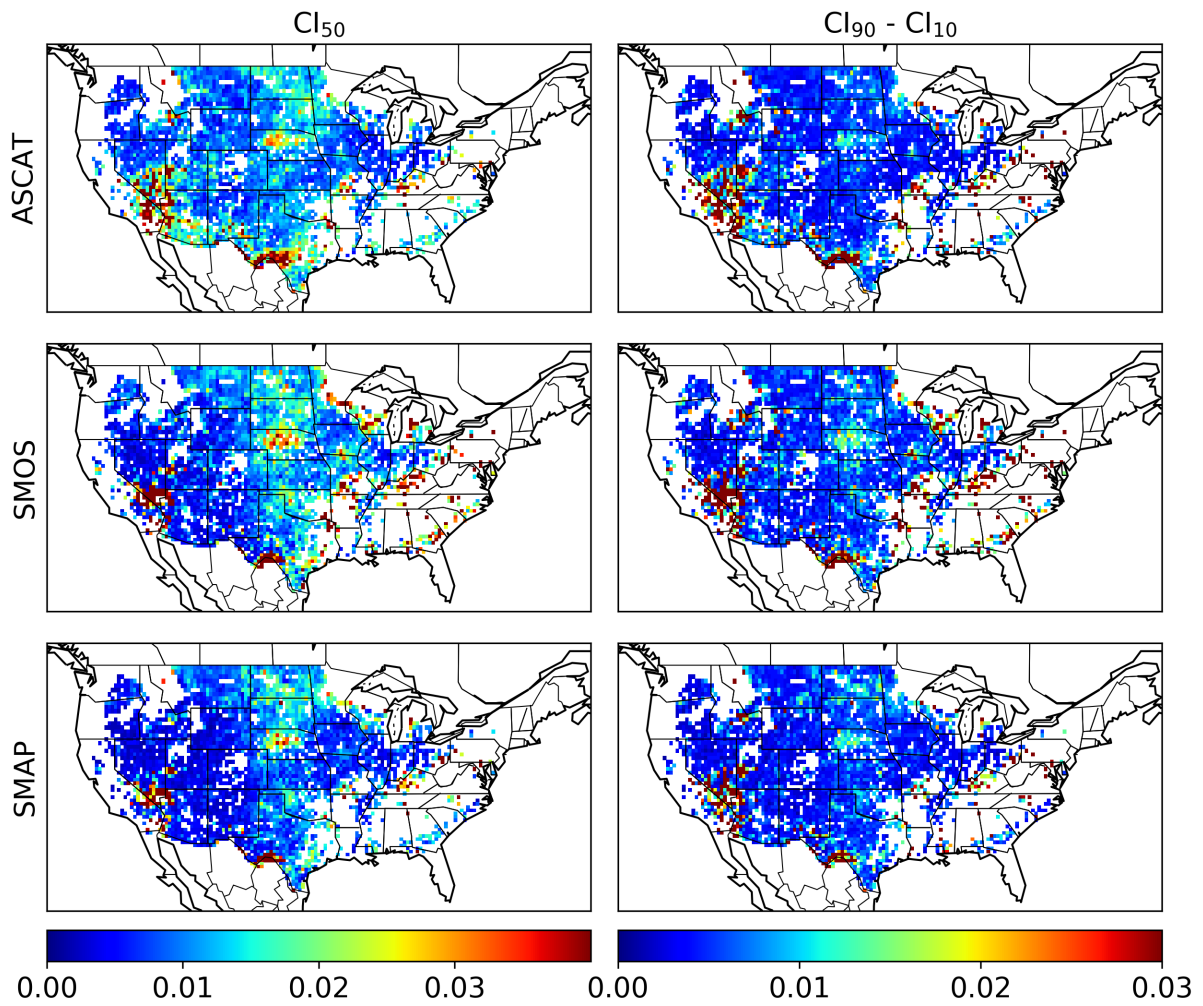


Figure A.12: Median of the bootstrapped TCA-based ubRMSEs [ $m^3m^{-3}$ ] (left) and associated 80% confidence intervals (right) of soil moisture anomaly estimates of ASCAT, SMOS, and SMAP.



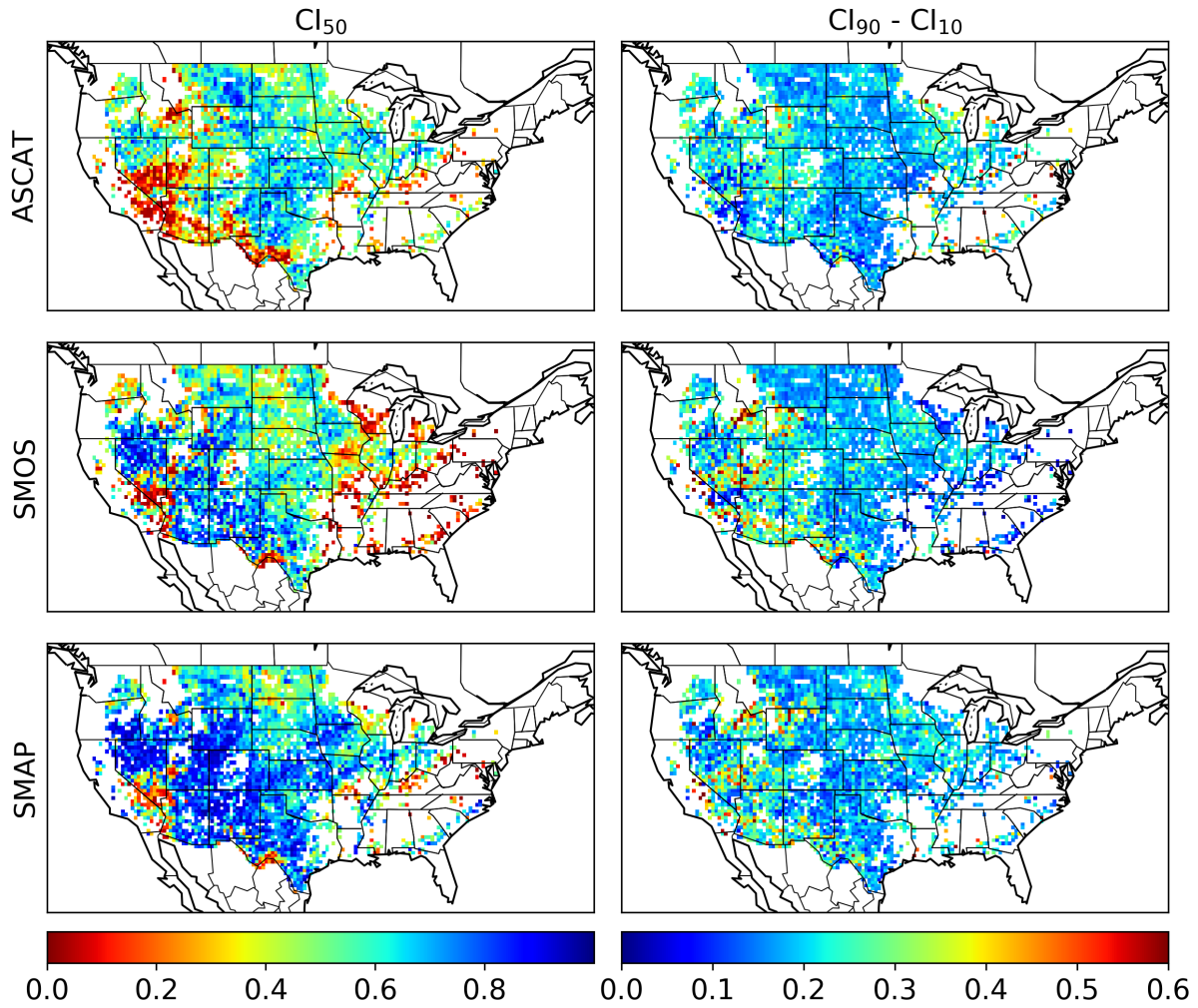


Figure A.13: Median of the bootstrapped TCA-based  $R^2$  estimates [-] (left) and associated 80% confidence intervals (right) of soil moisture anomaly estimates of ASCAT, SMOS, and SMAP.

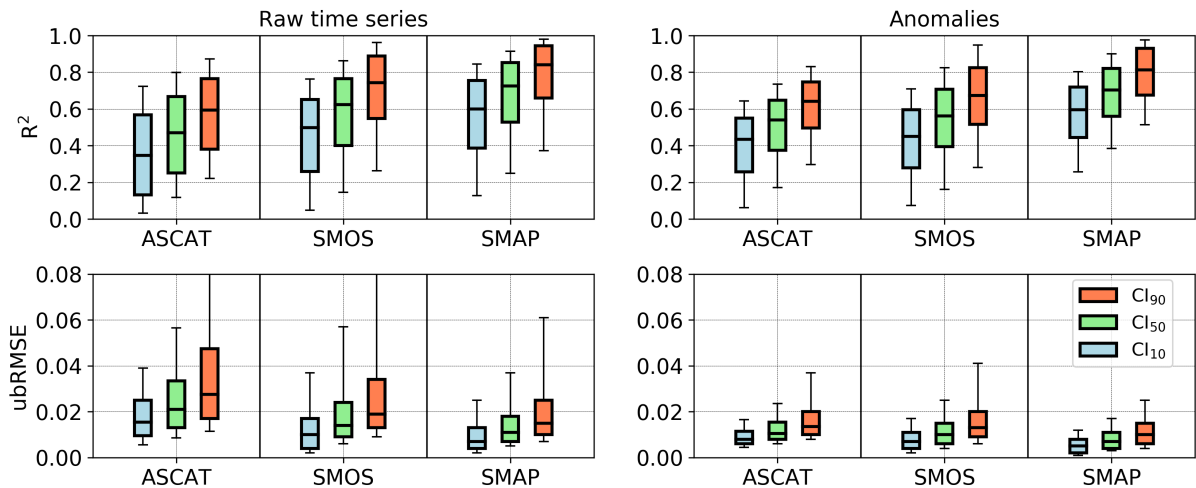


Figure A.14: Spatial summary statistics of the median of the bootstrapped TCA-based ubRMSEs [ $m^3m^{-3}$ ], and  $R^2$  estimates [-] and their 10% and 90% confidence limits, respectively, for raw soil moisture estimates and soil moisture anomalies of ASCAT, SMOS, and SMAP. Boxes represent the (spatial) median and inter-quartile-range and whiskers represent the 5 and 95 percentiles.

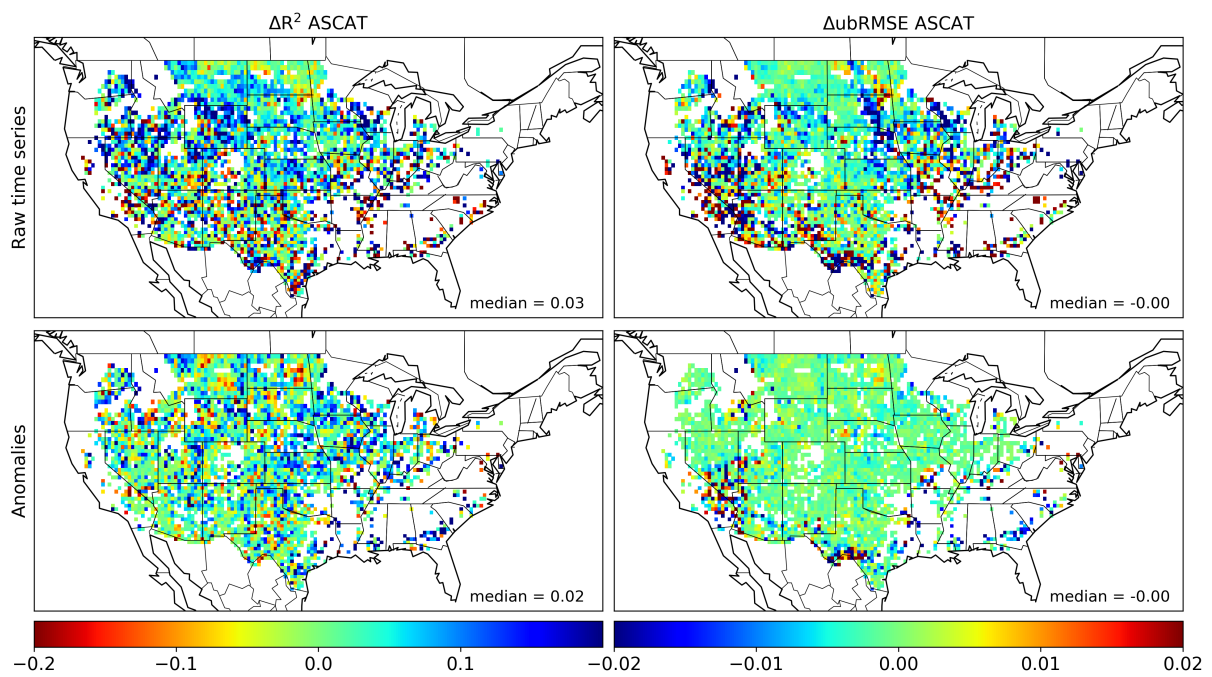


Figure A.15: Difference in TCA-based ubRMSE [ $m^3m^{-3}$ ] and  $R^2$  estimates [-] for raw soil moisture estimates (top) and soil moisture anomaly estimates (bottom) of ASCAT when using SMOS as third data set minus when using SMAP as third data set in the triplet.