

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

A. Gruber¹, G. De Lannoy¹, C. Albergel², A. Al-Yaari³, L. Brocca⁴, J.-C. Calvet², A. Colliander⁵, M. Cosh⁶, W. Crow⁶, W. Dorigo⁷, C. Draper⁸, M. Hirschi⁹, Y. Kerr¹⁰, A. Konings¹¹, W. Lahoz¹², K. McColl¹³, C. Montzka¹⁴, J. Muñoz-Sabater¹⁵, J. Peng¹⁶, R. Reichle¹⁷, P. Richaume¹⁰, C. Rüdiger¹⁸, T. Scanlon⁷, R. van der Schalie¹⁹, J.-P. Wigneron²⁰,
W. Wagner⁷

¹Department of Earth and Environmental Sciences, KU Leuven, Heverlee, Belgium

²Météo-France, Toulouse, France

³Sorbonne Université, UMR 7619 METIS, Paris, France

⁴Research Institute for Geo-Hydrological Protection, National Research Council, Perugia, Italy

⁵NASA Jet Propulsion Laboratory, Pasadena, CA, USA

⁶USDA ARS, Hydrology and Remote Sensing Laboratory, Beltsville, MD, USA

⁷Department of Geodesy and Geoinformation, TU Wien, Vienna, Austria

⁸Physical Sciences Division, NOAA/Earth System Research Laboratory, Boulder, Colorado

⁹Institute for Atmospheric and Climate Science, ETH Zürich, Zürich, Switzerland

¹⁰CESBIO (UMR 5126 CNES, CNRS, UT3, IRD), Toulouse, France

¹¹Department of Earth System Science, Stanford University, Stanford, CA, United States

¹²Norwegian Institute for Air Research, 2027 Kjeller, Norway

¹³Department of Earth and Planetary Sciences, Harvard University, Cambridge, MA, USA.

¹⁴Institute of Bio- and Geosciences: Agrosphere (IBG-3), Research Center Juelich, Germany

¹⁵European Centre for Medium-Range Weather Forecasts, Shinfield Road, Reading, UK

¹⁶School of Geography and the Environment, University of Oxford, Oxford, UK

¹⁷NASA Goddard Space Flight Center, Greenbelt, Maryland, USA

¹⁸Department of Civil Engineering, Monash University, Victoria, Australia

¹⁹VanderSat B.V., Haarlem, The Netherlands

²⁰ISPA, INRA Bordeaux, Bordeaux, France

28

29 Abstract

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

37

38 **Keywords:** remote sensing, soil moisture, validation, error characterization, error estimation
39 good practice, standardisation

40 1 Introduction

41 The validation of soil moisture data sets aims to provide quantitative information about their
42 quality by estimating systematic and random errors through analytical comparison to reference
43 data, which is presumed to closely represent the truth (*Justice et al.*, 2000; *JCGM*, 2008). For
44 satellite-derived products, this task is far from trivial because high-quality reference data are
45 virtually unavailable on a global scale at the coarse spatial resolution of space borne microwave
46 instruments that are predominantly used for soil moisture retrievals ($\sim 10^1 - 10^3$ km²), and the
47 retrieval quality is affected by numerous spatially and temporally variable factors (i.e. climatic,
48 topographic and land cover conditions as well as instrument characteristics and the retrieval
49 algorithm structure) (*Ochsner et al.*, 2013; *Crow et al.*, 2012; *Molero et al.*, 2018).

50 A host of methods exists to reconcile the distinct spatio-temporal characteristics of satellite
51 and reference data sets (sampling and overpass times, penetration depths, representativeness
52 errors, etc.; *Wang et al.*, 2012; *Albergel et al.*, 2008; *Gruber et al.*, 2013a; *Nicolai-Shaw et al.*,
53 2015; *Colliander et al.*, 2017a), which is required before calculating various performance met-
54 rics (correlation coefficients, root-mean-square-differences, triple collocation-based metrics, etc.;
55 *Entekhabi et al.*, 2010a; *Albergel et al.*, 2013; *Gruber et al.*, 2016a; *Loew et al.*, 2017). Given the
56 complexity of the validation problem, however, ambiguous results for the quality and ranking of
57 satellite soil moisture products can be found in the literature (e.g., *Wagner et al.*, 2014) depend-
58 ing on which pre-processing and evaluation strategies were followed and which reference data

59 were used. This paper is a community effort that addresses this issue and aims towards stan-
60 dardizing good practices for the validation of satellite-based near-surface soil moisture retrievals,
61 building upon ongoing international activities.

62 **1.1 Towards standardized validation practices**

63 Many efforts have been made to assess and standardize validation practices across Earth obser-
64 vation (EO) communities (*Zeng et al.*, 2015; *Loew et al.*, 2017; *Su et al.*, 2018). In the following
65 we summarize activities most relevant for satellite soil moisture products.

66 **1.1.1 CEOS LPV**

67 The main authority that guides validation activities for satellite-retrieved data of biogeophys-
68 ical variables is the Committee on Earth Observation Satellites (CEOS) Working Group on
69 Calibration and Validation (<http://ceos.org/ourwork/workinggroups/wgcv/>; last access: 1
70 July 2019). Activities related to soil moisture are coordinated by its Land Product Validation
71 (LPV) subgroup (<https://lpvs.gsfc.nasa.gov/>; last access: 1 July 2019). The CEOS LPV
72 defines four validation stages (see Table 1) that represent the level of sophistication of validation
73 protocols employed for a particular data product. Relevant for the work presented here is that
74 reaching validation stage 3 requires the implementation of a sophisticated validation framework,
75 as illustrated in Figure 1. In such a framework, standardized community-agreed methods that
76 are ideally described in a “Validation Good Practice Document” should be employed using fidu-
77 cial reference data (see Sec. 2) to generate standardized validation reports. With this paper we
78 aim at providing such a document. The last validation stage 4 is reached once these validation
79 reports are updated on a regular (at least annual) basis.

80 **1.1.2 Quality Assurance Frameworks**

81 The CEOS endorses the Quality Assurance Framework for Earth Observation (QA4EO; <http://qa4eo.org/>;
82 last access: 1 July 2019) as a framework to facilitate the provision of traceable
83 quality indicators which “shall provide sufficient information to allow all users to readily evaluate
84 the ‘fitness for purpose’ of the data or derived product” (*QA4EO*, 2010). The QA4EO provides
85 top-level guidance documents and templates that encourage the use of metrological principles
86 (see Sec. 1.1.3).

87 In 2014, the Quality Assurance for Essential Climate Variables (QA4ECV; <http://www.>

88 qa4ecv.eu/; last access: 1 July 2019) project was initiated to develop a set of guidelines for
89 the provision of traceable quality information taking into account the key principles of QA4EO
90 (*Scanlon et al.*, 2017). So far, quality assurance frameworks have been developed for selected
91 ECVs, not including soil moisture (e.g., *Peng et al.*, 2017). The guidelines developed by QA4EO
92 and QA4ECV are currently embraced by the Copernicus Climate Change Service (C3S; <https://climate.copernicus.eu/>;
93 last access: 1 July 2019) in order to build quality assured, fully
94 traceable Climate Data Records.

95 In 2018, the Quality Assurance for Soil Moisture project (QA4SM; <https://qa4sm.eodc.eu/>;
96 last access: 1 July 2019) was launched, specifically to create an online validation tool that
97 employs a community-agreed validation protocol (which we aim to provide with this paper)
98 for automatically and regularly generating soil moisture product validation reports, thereby
99 addressing the CEOS validation framework requirements (see Figure 1).

100 1.1.3 Metrology and traceability

101 The CEOS and the QA4EO encourage the use of metrological principles for validation purposes,
102 which are described in the “Guide to the expression of uncertainty in measurement” (GUM;
103 *JCGM*, 2008). The GUM is a reference document of the metrological community that provides
104 strict guidelines on how quality estimates of measurements should be obtained and reported.
105 In essence, it states that, since they never perfectly represent the true state of the physical
106 quantity being measured, all measurements should be complemented by uncertainty estimates
107 that summarize their probability density function (pdf). Furthermore, it states that these
108 uncertainties should be obtained by propagating the uncertainties from all components that
109 contribute to the measurement process in a way that is traceable back to the “International
110 System of Units” (SI) standards, either through the standard method for the propagation of
111 uncertainty (*Parinussa et al.*, 2011; *Merchant et al.*, 2017) or, if not possible analytically, through
112 Monte Carlo simulations (*JCGM*, 2008).

113 However, while being relatively straightforward in a laboratory or numerical environment,
114 the traceable propagation of uncertainties in space borne remote sensing measurements and re-
115 trievals thereof, in particular of soil moisture, faces two particular challenges. First, footprints of
116 current microwave instruments used for retrieving soil moisture span over tens to thousands of
117 square kilometers, thereby covering a large variety of climatic, topographic, and land cover condi-
118 tions. Although certain large-scale homogeneous regions are used for calibrating instruments and

119 determining Level 1 (L1) backscatter or brightness temperature uncertainties (e.g., rainforests
120 or polar snow fields; *Figa-Saldaña et al.*, 2002; *Macelloni et al.*, 2006), it is virtually impossible
121 to obtain global perfectly traceable uncertainty estimates representing all possible measurement
122 conditions. Second, uncertainty propagation assumes that the models used to propagate uncer-
123 tainties are themselves perfect (*Parinussa et al.*, 2011). For satellite soil moisture retrievals, this
124 is particularly problematic because uncertainties resulting from simplifications and assumptions
125 in both the L1 processing (i.e. geometric correction and radiometric calibration) and the Level
126 2 (L2) soil moisture retrieval algorithms cannot be accounted for. Taken together, these issues
127 render the reliable and traceable propagation of uncertainties from raw measurements through
128 the whole geophysical parameter retrieval process impossible. The soil moisture and other EO
129 communities have established certain strategies to recover this broken traceability chain by eval-
130 uating the soil moisture estimates post retrieval against a range of reference data from various
131 sources. Section 2 will discuss the requirements and current availability of such reference mea-
132 surements or estimates suited for validation activities. Before entering those discussions, it is
133 necessary to provide some relevant terminology.

134 1.2 Terminology

135 The CEOS and the QA4EO encourage the use of the terminology used within the metrological
136 community as described in the “International Vocabulary of Metrology” (VIM; *JCGM*, 2012).
137 However, there is a certain level of ambiguity in the existing EO literature, and even within
138 the VIM and the GUM, regarding the usage of important terms such as errors, uncertainties,
139 validation, and others. For a comprehensive summary of the most common definitions (from the
140 VIM, the CEOS, and other sources) we refer the reader to *Loew et al.* (2017). For the purpose
141 of this paper we stress that:

- 142 • in the scientific literature, the term *validation* is ubiquitous, yet its meaning and whether
143 or not anything can actually be *validated* - given the fundamental problem of an unknown
144 “truth” - has been subject to a decade-long debate (*Rykiel Jr.*, 1996). No consensus has
145 been found yet, because this is mainly a philosophical question. In the Earth sciences,
146 *validation* is used rather loosely and is often distinguished from the term *evaluation* such
147 that validation is used to refer to bias or uncertainty assessment using highly accurate or at
148 least well traceable in situ reference data (often misleadingly referred to as “ground truth”;
149 see Sec. 3.2), whereas evaluation is used to refer to the comparison against other coarse-

150 resolution satellite or modelled data with supposedly less well-defined uncertainties. How-
151 ever, ground reference data that could serve as reliable proxy for soil moisture retrievals at
152 a satellite scale are practically non-existent (with the exception of a marginally small num-
153 ber of heavily-equipped validation sites; see Sec. 2.2.1). Therefore, we more generally refer
154 to *validation* as the holistic process of gathering information from as many independent
155 sources as possible to enable a reliable quantitative judgement of the error characteristics
156 of a particular data set. This includes all, evaluation against ground measurements, com-
157 parison with estimates from land surface models, and satellite inter-comparisons. The final
158 declaration of a certain product to be *valid*, however, requires the specification of target
159 requirements for an intended use. As we will discuss later (see Sec. 3.8.2 and Sec. 5), no
160 meaningful requirements have yet been defined for satellite soil moisture applications;

- 161 • the term *measurement* refers to a quantity directly observed by a sensor (also called the
162 measurand), whereas the terms *estimate* and *retrieval* refer to a related quantity that has
163 been derived from the measurand. Accordingly, satellite sensors *measure* radiances from
164 which soil moisture or other quantities are being *estimated* or *retrieved*. Note, however,
165 that also in situ sensing technology *measures* only quantities related to water content,
166 such as dielectric constants, capacitance or weight, from which water content *estimates*
167 are derived. Notwithstanding, in situ soil moisture *estimates* are virtually always referred
168 to as *measurements*, and we will stick to this convention;
- 169 • the term *error* refers to the deviation of a single measurement (estimate) from the true
170 value of the quantity being measured (estimated), which is always unknown, whereas the
171 term *uncertainty* refers to the probability distribution underlying an error. For validation
172 purposes, this probability distribution is the actual quantity of interest;
- 173 • according to the GUM, the uncertainty of a measurement (estimate) generally contains
174 both systematic and random components. The laboratory environment of metrological
175 practices typically allows for thorough measurement calibration, where it is assumed that
176 systematic errors can be properly determined and corrected. Satellite soil moisture re-
177 trievals, however, usually contain considerable systematic errors which, especially for model
178 calibration and refinement, provide better insight when estimated separate from random
179 errors. Therefore, we use the term *bias* to refer to systematic errors only and the term
180 *uncertainty* to refer to random errors only, specifically to their standard deviation (or

- 181 variance);
- 182 • in the EO validation literature, bias is commonly estimated as the temporal mean difference
183 between two data sets. We follow the broader statistical definition of bias as auto-correlated
184 error, or as a property of an estimator to systematically over- or underestimate some
185 quantity (Dee, 2005). For better separability of its components, we use the terms *first-*
186 *order bias* and *second-order bias* to refer more specifically to additive and multiplicative
187 systematic errors, respectively (see Sec. 3.4.1);
 - 188 • the terms *trueness*, *precision*, and *accuracy* are popular antonyms for systematic errors,
189 random errors, and the combined systematic plus random errors, respectively (JCGM,
190 2012). However, trueness and precision are very rarely used in the soil moisture validation
191 literature and the term accuracy is often ambiguously used to refer to either systematic
192 or random errors alone; and
 - 193 • the concept of uncertainty is closely related to the concept of confidence intervals. Both
194 aim at describing the pdf underlying an estimate, although the term *uncertainty* is more
195 commonly used for describing the pdf behind an estimate that results from measurement
196 or retrieval errors (see Sec. 3.1), whereas the term *confidence interval* is more commonly
197 used for describing the pdf behind statistical parameters (such as statistical moments or
198 validation metrics that derive from these moments) that results from finite sample sizes
199 (see Sec. 3.5).

200 The remainder of this paper is organized as follows. Section 2 describes the most common
201 reference data sources used for soil moisture validation. Section 3 discusses relevant theoretical
202 aspects and the most common methods (including data pre-processing) for assessing soil moisture
203 data quality. Section 4 presents a validation guidance protocol that has been developed by a
204 gathering of experts across the community with an example implementation of that protocol
205 provided in Appendix A. Finally, Section 5 discusses research gaps that should be addressed in
206 the near future.

207 2 Reference data

208 The term *fiducial reference measurements* is often used to refer to a suite of independent, fully
209 characterized, and traceable measurements that meet the requirements on *reference standards*

210 as described by QA4EO (Fox, 2010), which should be used to assess the quality of EO products.
211 However, although highly accurate in situ soil moisture measurements exist and uncertainties
212 of the measurement devices can be reliably determined through laboratory and field calibration
213 activities (Cosh et al., 2005; Rüdiger et al., 2010; Caldwell et al., 2018), using such point-
214 scale measurements for evaluating satellite soil moisture data sets over large areas is a very
215 difficult task owing to the coarse resolution of space borne microwave instruments and vast
216 heterogeneities across landscapes (Cosh et al., 2004, 2006; Famiglietti et al., 1999; Brocca et al.,
217 2010a; Miralles et al., 2010; Crow et al., 2012; Nicolai-Shaw et al., 2015; Molero et al., 2018).
218 While general calibration functions can yield soil moisture measurement uncertainties in the
219 order of 0.02 to 0.03 m³m⁻³ (Seyfried et al., 2005), which can be improved to below 0.005 m³m⁻³
220 when applying a dedicated field calibration (Bogena et al., 2017), spatial representativeness
221 errors that arise when using in situ sensors to represent soil moisture variations at the satellite
222 scale (see Sec. 3.2) can easily exceed these numbers (Gruber et al., 2013a).

223 For satellite validation purposes, numerous field and airborne campaigns have been carried
224 out to obtain reliable satellite footprint scale reference data and to quantitatively assess the
225 potential spatio-temporal representativeness (see Sec. 3.2) of single or small sets of in situ soil
226 moisture stations (Famiglietti et al., 2008; Cosh et al., 2008; Brocca et al., 2012; McNairn et al.,
227 2015). Additionally, validation activities are complemented with land surface model output and
228 other satellite products for comparison to get as complete a picture as possible of a product's
229 error characteristics (Brocca et al., 2010b; Draper et al., 2013; Al-Yaari et al., 2014; Dorigo
230 et al., 2015; Kerr et al., 2016; Miyaoka et al., 2017). The various reference data sources and
231 their limitations are discussed below. Some publicly available reference data sources that are
232 commonly used for satellite soil moisture validation are listed in Table 2.

233 2.1 Field campaigns

234 Field campaigns are labor-intensive studies that use highly accurate measurement techniques
235 to obtain reliable and traceable representations of larger scale average soil moisture. Addition-
236 ally, many field campaigns collect other relevant surface properties such as soil texture, surface
237 roughness, vegetation cover, etc. The campaigns provide snapshots in time that have a set
238 of parameters characterized in detail and can answer certain specific questions related to the
239 calibration and validation of soil moisture products. However, the full validation of satellite
240 products requires long and consistent time series (see Sec. 3.4). Therefore, a number of field

241 campaigns have supported this goal by focusing on various specific aspects for improving the
242 scalability of in situ measurement networks to remote sensing footprint size. An example of
243 this is the establishment of temporally stable locations (*Vachaud et al.*, 1985; *Starks et al.*,
244 2006) that sufficiently capture sub-pixel heterogeneities, allowing the continuous observation of
245 satellite footprint-scale areas with sufficient and well-characterized accuracy. Moreover, field
246 experiment often supplement the ground measurements with airborne observations. Airborne
247 observations can be used to evaluate soil moisture retrievals over a larger area, allowing to assess
248 the spatial soil moisture (as well as brightness temperature and backscatter) variability within
249 and across multiple satellite grid cells.

250 Early field campaigns were focused on understanding large-scale soil moisture dynamics
251 with aircraft support such as the HAPEX-MOBILHY (*Noilhan et al.*, 1991), the BOREAS
252 (*Cuenca et al.*, 1997), the Washita'92 (*Jackson et al.*, 1995), and the 1997 Southern Great
253 Plains Hydrology Experiment (SGP97) campaigns (*Jackson et al.*, 1999). These experiments
254 assessed the potential of soil moisture remote sensing over larger domains as a part of hydrologic
255 research. This evolved into satellite associated field campaigns, which can be divided into pre-
256 launch and post-launch experiments based on their objectives. The Soil Moisture Experiments
257 (SMEX) in 2002-2004 in the United States (*Jackson et al.*, 2005; *Bindlish et al.*, 2006, 2008)
258 were designed in large part for the evaluation of AMSR-E soil moisture products. The National
259 Airborne Field Experiment (NAFE) in Australia (*Panciera et al.*, 2008) was designed for pre-
260 launch studies of SMOS, while the Australian Airborne Calibration/Validation Experiments
261 for SMOS (AACES; *Peischl et al.*, 2012) targeted the evaluation of SMOS retrievals. The
262 objective of the Canadian Experiment for Soil Moisture (CANEX-10; *Magagi et al.*, 2013) was
263 to contribute to the evaluation of SMOS and pre-launch activities for SMAP, and the CAROLS
264 airborne campaigns (*Albergel et al.*, 2011; *Zribi et al.*, 2011) were designed for the evaluation of
265 SMOS. The SMAP mission also carried out a dedicated pre-launch campaign in 2012 (SMAP
266 Validation Experiment 2012, SMAPVEX12; *McNairn et al.*, 2015) and post-launch validation
267 campaigns in 2015 and 2016 (*Colliander et al.*, 2017b, 2019).

268 The earlier campaigns established a protocol for the synchronous collection of ground-based
269 soil moisture measurements with airborne microwave instrumentation, which was followed in
270 most of the subsequent experiments. In the process of developing standardized data collection
271 protocols, these field campaigns specifically focused on the investigation of the spatial distri-
272 bution of soil moisture and its evolution with drying or wetting, the soil moisture variability

273 across scales, and the statistical relationship between spatial standard deviation and extent scale.
274 These parameters drive the potential representativeness of in situ measurements for coarse soil
275 moisture product evaluation and their knowledge hence allows the determination of the number
276 of ground samples required to obtain sufficiently reliable reference data. To this end, at many
277 of the experiment locations, the labor-intensive field campaign observations were supplemented
278 with long-term in situ monitoring stations, thus providing long-term high-density satellite vali-
279 dation sites.

280 **2.2 In situ networks**

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

291 **2.2.1 Dense networks**

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

299 These SMAP Cal/Val partners have a diverse heritage. Some networks were originally de-
300 ployed for Cal/Val of the AMSR-E product (*Martínez-Fernández and Ceballos, 2005; Jackson*
301 *et al., 2010*), SMOS (*Bircher et al., 2012; Smith et al., 2012; Djamaï et al., 2015*), or SMAP
302 (*Caldwell et al., 2019*), while others evolved from hydrologic monitoring networks (*Bogena et al.,*

2018) or from some other purpose such as aircraft validation projects like AIRMOSS (*Moghaddam et al.*, 2010). During the SMAP project, several networks were selected as potential candidate sites for Cal/Val activities. The candidate networks whose accuracy versus physically collected volumetric soil moisture was already demonstrated and documented in a traceable manner, were promoted to core validation sites. To date, these sites are considered to provide the best possible ground reference data for satellite footprint-scale soil moisture dynamics (*Colliander et al.*, 2017a; *Chen et al.*, 2019).

2.2.2 Sparse networks

A host of other operational and experimental in situ sites exist worldwide, operating soil moisture measurement stations that are potentially suited for satellite soil moisture evaluation yet with a considerably smaller station density and often lacking information on their coarse-scale representativeness and their own inherent error characteristics (*Gruber et al.*, 2013a; *Chen et al.*, 2017). Nonetheless, these sites are valuable to complement core validation sites due to their considerably larger spatial coverage across a variety of climatic regimes and biomes (see Sec. 3).

An important source for data from sparse networks is the International Soil Moisture Network (ISMN; *Dorigo et al.*, 2011a,b), which is a data hosting facility that harmonizes soil moisture measurements from in situ networks worldwide, applies automated and uniform quality control procedures to flag suspicious measurements (*Dorigo et al.*, 2013), and distributes them on a cost-free basis in a common format (<http://ismn.geo.tuwien.ac.at/>; last access: 1 July 2019). The ISMN was established by ESA in the framework of SMOS Cal/Val activities. Currently, it contains data from more than 2400 stations worldwide, operated across 59 different measurement networks (see Figure 2) including historical networks that are no longer operational. In addition to soil moisture, many networks provide measurements of other variables such as precipitation or temperature as well as ancillary information such as soil texture or land cover. Note, however, that sensor technologies and data quality vary greatly across networks and measurement stations (*Dorigo et al.*, 2011b; *Babaeian et al.*, 2019).

2.3 Model simulations

Due to the limited coverage and representativeness of ground reference data, validation activities are complemented with soil moisture simulations from land surface models (LSMs) as an alternative reference data source (*Lahoz and De Lannoy*, 2014). Model simulations can provide

333 spatially complete global soil moisture maps at a spatial (grid) resolution similar to that of satel-
334 lite footprints, but they may still contain considerable representativeness errors (see Sec. 3.2)
335 originating from simplifications of sub-grid heterogeneities, a scale-mismatch of the underlying
336 atmospheric forcing data, errors in the model parameterization, or simply because the meaning
337 of the modelled “soil moisture” is different (e.g. representing a different layer depth or expressed
338 in different units). Moreover, biases and uncertainties in model simulations are highly variable
339 and often also not well quantified (*Koster et al.*, 2009; *Albergel et al.*, 2013), making it difficult
340 to separate satellite retrieval errors from modelling errors in a direct comparison (see Sec. 3).

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

348 2.4 Satellite products

349 A multitude of soil moisture products from different satellite sensors (*Babaeian et al.*, 2019)
350 are commonly used as additional coarse resolution reference data sets for validation purposes,
351 either for consistency assessment through direct comparison (*Al-Yaari et al.*, 2014; *Burgin et al.*,
352 2017), or within triple collocation analysis (*Dorigo et al.*, 2010; *Draper et al.*, 2013, see Sec. 3).
353 Like model simulations and sparse networks, they typically lack reliable and traceable bias and
354 uncertainty characterization. Also, available satellite sensors observe at different wavelengths,
355 polarizations, and incidence angles and have therefore a varying sensitivity to soil moisture
356 (*Ulaby et al.*, 2014). Hence, the information gleaned from a direct comparison is limited (see
357 Sec. 3.4.2). Furthermore, different satellite retrieval products (and model simulations) can use
358 similar ancillary information such as temperature and/or vegetation information in a radiative
359 transfer model, resulting in correlated errors (*Gruber et al.*, 2016b) which may complicate a fair
360 data comparison (see Sec. 3.4.2). Comprehensive lists of commonly used and publicly available
361 satellite soil moisture products, including some validation information where available, can be
362 found at <https://lpvs.gsfc.nasa.gov/producers2.php?topic=SM> (last access: 1 July 2019)
363 and in *Babaeian et al.* (2019).

364 **3 Theory**

365 This section provides the theoretical background for error characterization and how it relates to
366 satellite soil moisture validation, including the assumptions, limitations and pre-processing steps
367 involved. Although our main focus here is the validation of near-surface satellite soil moisture
368 products, many of the principles discussed below can be equally applied to assess the quality
369 of soil moisture products from other sources, as well as of other biogeophysical variables (*Loew*
370 *et al.*, 2017).

371 **3.1 Errors**

372 An estimation error e_x is defined as the deviation of an estimate x , in our case a satellite soil
373 moisture retrieval, from the true state t of the quantity being estimated (*JCGM*, 2008):

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

374 Important for understanding errors is that the “truth” is a hypothetical concept. For the case
375 of space borne microwave instruments, actual satellite footprints are overlapping elliptical areas
376 with strong signal intensity gradients from the footprint center outwards (depending on the
377 antenna gain pattern) and varying, surface property dependent signal penetration depth (*Ulaby*
378 *et al.*, 2014). Horizontal footprint boundaries are commonly defined as the 3 dB region, i.e.
379 the region of the antenna pattern projection on the ground where the gain is within 3 dB (50
380 %) of the peak value. Products derived thereof are typically sampled onto spatial grids with
381 sharp boundaries between grid cells and a constant layer depth to facilitate further geospatial
382 analysis (*Bartalis et al.*, 2006; *Brodzik et al.*, 2012; *Bauer-Marschallinger et al.*, 2014). The
383 “true” soil moisture signal that drives the microwave measurement and the subsequent gridded
384 soil moisture retrieval will therefore never be the real average soil moisture of the grid cell
385 to which the retrieval is assigned. Moreover, for validation purposes, the unknown “truth” is
386 approximated by reference data, which themselves contain errors and may also be driven by a
387 soil volume that is different from the satellite grid cell they are supposed to represent (see Sec.
388 2).

389 3.2 Representativeness

390 The difference between the true soil moisture that actually affects a (microwave) measurement
391 associated with a particular grid cell and the true soil moisture within that grid cell is often
392 referred to as representativeness error (*Gruber et al.*, 2016a). However, it is worth noting that
393 representativeness errors have different definitions (*Van Leeuwen*, 2015). The remote sensing
394 community mostly assigns them to the mismatch between the spatial support of a measurement
395 and the spatial resolution of the defined sampling grid, sometimes also referred to as scaling
396 error (*Miralles et al.*, 2010; *Crow et al.*, 2012; *Gruber et al.*, 2013a; *Molero et al.*, 2018). In
397 the modelling community, representativeness errors mostly refer to a model’s lacking ability to
398 represent reality and, as such, to imperfections in the model structure and in parameterization
399 (e.g., unresolved sub-grid scale processes). For the purpose of data validation, it is practical
400 to use a definition that potentially allows us to separate representativeness errors from other
401 error sources upon estimation. Therefore, recall that the general definition of error in Eq. (1)
402 requires the choice of a “truth”, which is the soil moisture state within a target volume (grid
403 cell) that one aims to estimate as accurately as possible. We define representativeness errors as
404 those deviations of a product from such chosen, unknown “true” state, which are related to real
405 soil moisture variations. They can occur, for example, if the actual measurement footprint of a
406 satellite extends beyond the grid cell boundaries associated with the chosen, unknown “truth”,
407 if an inadequate soil parameterization in a radiative transfer model causes the soil moisture
408 retrievals to represent deeper soil layers than the chosen, unknown “truth”, or if point-scale
409 ground measurements are used as a reference for grid cell-scale soil moisture dynamics. As
410 such, representativeness errors of different data sets may be correlated even if the products are
411 otherwise independent.

412 In summary, representativeness errors have important implications for validation in that
413 they limit the information one can glean from the comparison between products, even if a
414 chosen reference product is itself highly accurate (see Sec. 3.4.1). Since the temporal and
415 spatial resolution and sampling of satellite and available reference measurements or estimates
416 hardly ever match, (relative) representativeness errors will often reach considerable magnitudes
417 (*Miralles et al.*, 2010; *Crow et al.*, 2012). To minimize their influence, several pre-processing
418 steps are typically applied, which are discussed in the following section together with other
419 pre-processing steps that are necessary before validation metrics can or should be calculated.

420 **3.3 Pre-processing**

421 Pre-processing steps necessary for validation aim to find match-ups in space and time between
422 measurements and/or estimates that have different spatial resolutions, are sampled on to differ-
423 ent grids, and/or are acquired at different times. Additionally, depending on the reference data
424 choice, statistical rescaling methods are often applied to minimize the impact of representative-
425 ness errors. Moreover, data pre-processing typically involves the masking of unreliable satellite
426 retrievals and reference measurements or estimates. Lastly, data sets are sometimes decom-
427 posed into different frequency components in order to separately assess a product’s ability of
428 accurately representing short-term, seasonal, and inter-annual soil moisture variability (*Draper*
429 *and Reichle, 2015*).

430 **3.3.1 Data masking**

431 Satellite-derived soil moisture products are typically accompanied by a set of quality flags. They
432 can be indicators of suspected contamination of the microwave signals or problems during the
433 retrieval. Typical examples are indicators for the probability of frozen soil, dense vegetation
434 coverage, radio frequency interference (RFI), or urban or water contamination, to name a few
435 (e.g., *Parinussa et al., 2011; Naeimi et al., 2012; Kerr et al., 2012; de Nijs et al., 2015*).

436 The validation of a product should be based only on those retrievals that are considered
437 “good” for a given application. While masking data points using binary “use / do not use” flags
438 is straightforward, some quality flags require the decision of a threshold below or above which
439 individual retrievals are masked out (e.g., the probability of RFI occurrence or the water body
440 fraction), which implies a trade-off between data quality and measurement density. Typically,
441 data producers provide recommendations for these thresholds. In addition to the quality flags
442 inherent in the soil moisture products, auxiliary static and/or dynamic data from land surface
443 models or other sources are often used to mask out retrievals that can be considered unreliable.
444 The most commonly used masking criteria are based on surface and/or air temperature and
445 snow height and/or snow water equivalent estimates obtained from land surface models, or
446 vegetation-related estimates (such as vegetation water content or vegetation optical depth) from
447 satellite sensors or models (*Al-Yaari et al., 2014; Dorigo et al., 2015; Gruber et al., 2017*). It
448 should be kept in mind, however, that all quality flags (both provided alongside a product or
449 derived from an ancillary source) are based on data which themselves are subject to errors and
450 are therefore inherently uncertain.

451 Note that also reference data sets, in particular in situ measurements, also often undergo
452 quality control procedures and provide quality flags, which should be used to mask out unreliable
453 measurements before using them to evaluate satellite retrievals (as is the case for example for the
454 ISMN; *Dorigo et al.*, 2013). When comparing biases or uncertainties of different soil moisture
455 products, the masking procedures applied to these data sets should be identical in order to
456 compare the quality of retrievals from measurements that were taken under the same (or at
457 least similar) conditions. However, if quality flags that are tailored to one data set are applied
458 to another, some of the products may appear better or worse than they would when using only
459 their own inherent quality control. This is especially true if the flags of one product are much
460 more conservative than those of another. Most product comparison studies do not take this issue
461 into account. One possible approach to address it would be to compare biases and uncertainties
462 from common periods also with those in periods where only some products provide unflagged
463 soil moisture retrievals (based on their own quality control) and to put this into perspective with
464 the temporal measurement density before and after product collocation. However, this requires
465 the availability of appropriate reference data in collocated and non-collocated periods as well as
466 the ability to account for possibly varying accuracy and representativeness of the reference data
467 in these periods. Also, depending on the overall data density, it may be difficult to assess biases
468 and uncertainties in these periods due to the presence of large statistical sampling errors (see
469 Sec. 3.5).

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

473 **3.3.2 Collocation**

474 Satellite sensors acquire measurements that are irregularly distributed in space and time owing
475 to their orbiting nature and specific antenna patterns. In the soil moisture retrieval process,
476 these measurements are typically sampled onto spatial grids (for noise reduction purposes these
477 grids are often oversampled, i.e. the grid sampling - sometimes also referred to as grid posting -
478 is typically higher than the antenna resolution) and sometimes also to regular time steps (e.g.,
479 00:00 UTC) in order to generate, for example, daily global soil moisture maps and/or time
480 series (*Kerr et al.*, 2012; *O'Neill et al.*, 2012; *H-SAF*, 2018; *Gruber et al.*, 2019a). However,
481 neither the resolution nor the sampling of in situ reference measurements or model simulations

482 ever perfectly match those of the satellite products being evaluated. Consequently, the process
483 of finding match-ups between satellite and reference data points in space and time, commonly
484 referred to as collocation, is essentially a resampling task (*Loew et al.*, 2017). Since the spatial
485 resolution of the compared products can be very different (especially between in situ and satellite
486 / modelled data), statistical rescaling methods are often additionally applied in the collocation
487 process to minimize the impact of (especially spatial) representativeness errors on validation
488 metrics.

489 **Spatial resampling**

490 In situ measurements are point-scale measurements that sample only a few cubic centimeters of
491 the soil. When used for evaluating satellite products, stations from sparse networks are typically
492 sampled onto the satellite grid using a nearest-neighbour (NN) search, i.e. by matching the
493 stations to the satellite grid cells within which they are located (*Albergel et al.*, 2012; *Dorigo et al.*,
494 2015; *Chen et al.*, 2017). For dense networks, commonly all stations that lie within a particular
495 satellite grid cell are (after quality control) averaged (*Jackson et al.*, 2010; *Gruber et al.*, 2015;
496 *Colliander et al.*, 2017a), either by calculating the arithmetic mean or by calculating a weighted
497 average where higher weights are applied to stations that are expected to be more representative
498 for the grid cell average soil moisture. Such stations can be identified, for example, via a temporal
499 stability analysis (*Vachaud et al.*, 1985; *Yee et al.*, 2016), through Voronoi diagrams (*Colliander*
500 *et al.*, 2017a), or by using landscape characteristics such as land cover or soil properties.

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

511 **Temporal resampling**

512 In situ measurements and land model estimates are typically sampled more frequently than satel-

513 lite soil moisture retrievals. Therefore, the reference measurements and estimates are matched
514 in time to the irregular satellite observation times, typically by selecting the temporally closest
515 (NN) reference measurement or estimate within a pre-defined search window (i.e. applying a
516 maximum temporal distance threshold; *Chen et al.*, 2017). Depending on the sampling in-
517 terval of the reference data sets (for in situ data typically hourly and for global land surface
518 models typically one to six hourly) and on whether or not satellite observations have been a
519 priori resampled already (see above), this can lead to considerable differences between the ac-
520 tual measurement/estimation times of collocated satellite and reference data points. The issue
521 is typically limited when using in situ or model data as reference. However, if multiple satel-
522 lite products are evaluated simultaneously, their different overpass times are usually accounted
523 for by either picking one of them as (temporal) reference and matching the other ones against
524 it, or by sampling all satellite products to regularized time steps (e.g., 00:00 UTC; *Gruber*
525 *et al.*, 2017), which in any case favours the satellite data set whose actual measurement times
526 are closest to the reference points. Note that the retrieval quality of satellite data sets may
527 strongly depend on the time of observation. This is especially true for passive systems, where
528 soil moisture retrievals are known to be strongly affected by temporal temperature fluctuations
529 and temperature gradients in soil and vegetation cover (*Parinussa et al.*, 2015).

530 Taken together, the different measurement/estimation times of satellite and reference data
531 sets that have been collocated will induce temporal representativeness errors, originating from
532 the actual soil moisture changes that take place during these periods. Often these errors are
533 assumed to be negligible or at least below the noise level of the products. In principle, one could
534 employ more sophisticated resampling algorithms to minimize these representativeness errors,
535 for example auto-regressive interpolation methods with or without auxiliary information such
536 as precipitation, evapotranspiration, or soil texture. However, more research is needed to assess
537 the impact of temporal interpolation approaches on validation metrics.

538 **(Statistical) rescaling**

539 The resampling procedures described above provide data set match-ups in space and time which
540 are required for statistical comparison (see Sec. 3.4). As discussed in Sec. 3.1, the measurements
541 or estimates of the collocated products are driven by the soil moisture state of different soil
542 volumes at different times due to the different underlying actual spatio-temporal resolution of
543 the data sets. The latter is related to the antenna and surface properties and cannot be corrected

544 for by common resampling methods. Therefore, a direct comparison of these products will be
545 subject to representativeness errors, which may dominate the total soil moisture retrieval errors
546 (*Gruber et al.*, 2013a; *Chen et al.*, 2017; *Molero et al.*, 2018). However, owing to the large-scale
547 and auto-correlated nature of processes that drive soil moisture changes (*Crow et al.*, 2012), parts
548 of these errors are systematic and can hence be corrected for by removing *relative differences*
549 between the considered data sets (see Sec. 3.4).

550 The two most common rescaling approaches are to match either the temporal mean and
551 standard deviation of the data sets that are to be compared (*Scipal et al.*, 2008a; *Dorigo et al.*,
552 2010; *Albergel et al.*, 2012), or to match their complete cumulative distribution function (CDF),
553 which additionally corrects for differences in higher statistical moments in case the products
554 are expected not to be perfectly Gaussian distributed (*Reichle and Koster*, 2004; *Kumar et al.*,
555 2012). However, any rescaling approach that transforms one data set into the data space of
556 another (without additional information) assumes the signal-to-noise ratios (SNRs) of the two
557 involved data sets to be identical, which, since this is usually not the case, can lead to biased
558 rescaling parameters that do not fully correct the systematic representativeness errors (see Sec.
559 3.4.2; *Stoffelen*, 1998; *Yilmaz and Crow*, 2013). Alternatively, triple collocation analysis (*Stof-*
560 *felen*, 1998; *Su et al.*, 2014; *Gruber et al.*, 2016a) is often employed, using a third data set to take
561 different SNRs into account when matching the standard deviation of the underlying soil mois-
562 ture signals, thereby potentially providing consistent rescaling parameters (*Yilmaz and Crow*,
563 2013).

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

575 After rescaling, long-term bias estimation is obviously no longer meaningful as systematic

576 differences between the data sets, which would normally serve as proxy for biases, have been
577 intentionally removed. However, shorter-term biases as well as random representativeness errors
578 may remain and can considerably contribute to subsequent uncertainty estimates (see Sec. 3.4.1).

579 3.3.3 Signal decomposition

580 The quality of soil moisture products can vary considerably across time scales (*Su and Ryu, 2015;*
581 *Draper and Reichle, 2015; Molero et al., 2018; Gruber et al., 2019a*). For example, some soil
582 moisture products are better at accurately representing the seasonal cycle whereas other products
583 more accurately capture short-term fluctuations. Therefore, products are often decomposed into
584 different frequency components which are then evaluated separately (in addition to the bulk
585 time series). In Earth sciences, such decomposition is often done using moving-average windows
586 (*Narapusetty et al., 2009*). For soil moisture, a moving window of several weeks, centered on
587 the measurement or estimation time, is typically used to obtain intra-annual low-frequency
588 soil moisture dynamics (*Albergel et al., 2012; Chen et al., 2017*), referred to as seasonalities.
589 Residuals thereof are referred to as short-term anomalies which represent higher-frequency, sub-
590 seasonal soil moisture variations, that is, short-term drying and wetting events. Additionally,
591 so-called long-term anomalies are often calculated as residuals relative to a multi-year mean
592 seasonal cycle, referred to as the soil moisture climatology, which is typically calculated by
593 applying a moving-average window of similar size (a few weeks) to each day-of-the-year (DOY),
594 i.e. averaging all measurements or estimates of all years that fall inside the specified time window
595 around a particular DOY (*Miralles et al., 2010; Draper et al., 2013*). These long-term anomalies
596 contain information about both short-term drying and wetting events and seasonal deviations
597 from the long-term mean seasonal cycle.

598 While the evaluation of short-term soil moisture anomalies aims at assessing a data set's
599 capability of capturing individual drying or wetting events, uncertainties of long-term anomalies
600 represent its performance in capturing both short-term variability and inter-annual variations
601 such as prolonged droughts or floods as well as climate trends. However, the latter rely on a
602 climatology estimate that requires historical data records in the order of decades (*Dorigo et al.,*
603 *2012*), which are often not available, especially not at the beginning of a new mission (current
604 microwave missions cover a time period of maximum 5-10 years). Therefore, one often has to
605 rely on uncertainty estimates for seasonalities and short-term anomalies alone, which jointly
606 drive uncertainties in long-term anomalies.

607 **3.4 Metrics**

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

620 **3.4.1 Assumptions**

621 The fundamental assumption underlying almost all satellite soil moisture validation studies is
622 that of additive zero-mean random errors (ε_x), and additive (first-order; α_x) and multiplicative
623 (second-order; β_x) systematic errors (*Gruber et al.*, 2016a):

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

624 This error model applies to both the data set one aims to evaluate and the reference data sets.
625 Notice that the total error e_x in Eq. (1) has now been separated into its systematic (α_x and β_x)
626 and random (ε_x) components. These components contain instrument errors (i.e. noise and mis-
627 calibration), errors in the retrieval model and parameterization, and other representativeness
628 errors with respect to the assumed grid cell average soil moisture t (although the boundaries
629 between the latter two are somewhat fuzzy; see Sec. 3.1).

630 To disentangle errors from different data sets and from actual soil moisture variations, all
631 common data comparison metrics require the errors to be homoscedastic (i.e. independent from
632 the soil moisture state, in the literature often referred to as orthogonality with respect to the
633 truth; *Yilmaz and Crow*, 2014) and mutually uncorrelated between products. Remember, how-
634 ever, that the *representativeness* error components of the different products may (by definition)

635 be correlated both with the truth t and with each other, even if the products are otherwise
 636 independent (see Sec. 3.1).

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

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

643 where the overline, σ_i^2 and σ_{ij} refer to the (temporal) mean, variance, and covariance, respec-
 644 tively; and y denotes a reference data set that follows the same error model as x (Eq. (2)).
 645 Because *representativeness* errors may contain an orthogonal, a non-orthogonal, and a mutually
 646 correlated component (see above), we combine it with all other random error in the individual
 647 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
 648 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
 649 by representativeness errors only), for clarity. Systematic representativeness errors are included
 650 in the α_x and β_x coefficients.

651 The goal of validation is now to estimate α_x and β_x , and the standard deviation of ε_x (σ_{ε_x}),
 652 i.e. biases and uncertainties in the satellite data set under evaluation. The properties of the
 653 different reference data sets available (see Sec. 2) determine which error components will be
 654 dominant in Eq. (3), and consequently, which ones can be estimated by the available validation
 655 metrics (see Sec. 3.4.3 and 3.4.4).

656 Note, however, that α_x , β_x , and σ_{ε_x} contain lumped estimates of all systematic and random
 657 errors that accumulate in the soil moisture retrieval process, such as instrument noise, errors
 658 in the radiometric calibration, and imperfections in the retrieval model (e.g., resulting from
 659 the oversimplification and underdetermination of common radiative transfer models; *Quast and*
 660 *Wagner, 2016; Wigneron et al., 2017*), which can typically not be disentangled into its individual
 661 components.

662 3.4.2 Relative and TCA-based metrics: opportunities and limitations

663 For discussing the various metrics we will follow the notation of fiducial reference data (see Sec.
664 2) to refer to data sets that provide a thoroughly calibrated soil moisture proxy at the satellite
665 scale with traceable uncertainty characteristics (i.e. $\alpha_y \approx 0, \beta_y \approx 1$ in Eq. (2)). ε_y may be
666 non-zero but $\sigma_{\varepsilon_y}^2$ has to be at least well determined from laboratory experiments and field cam-
667 paigns and could hence be corrected for in the validation metrics. As mentioned, only the core
668 validation sites are currently considered as fiducial reference data capable of providing a reliable
669 representation of satellite footprint-scale soil moisture (see Sec. 2.2.1). They are therefore the
670 only reliable proxy for bias and uncertainty estimation from direct comparison, but are limited
671 to very few regions. Non-fiducial reference data refer to coarse-resolution products such as land
672 surface model simulations or other satellite data sets which may have non-negligible or non-
673 traceable biases and uncertainties as well as potentially considerable representativeness errors,
674 or to in situ data from sparse networks or not properly calibrated and validated dense networks,
675 both of which are expected to have larger representativeness errors than coarse-resolution refer-
676 ence data sets. Therefore, direct comparison against non-fiducial reference data can only provide
677 information of which data set is systematically drier or wetter than the other but without rela-
678 tion to a true grid cell average, and only lumped estimates of the uncertainty of both compared
679 products. Nonetheless, given their larger-scale and long-term availability, sparse networks and
680 land surface models are of important complementary value for validating satellite products. In
681 particular, one can obtain valuable information about the relative ranking of different products
682 as well as about performance changes over time when comparing against the same reference
683 product.

684 Introducing a second reference data set z that follows the same covariance properties (Eq.
685 (3)) as y (commonly referred to as triple collocation analysis, TCA; *Stoffelen, 1998; Scipal*
686 *et al., 2008b; Gruber et al., 2016a*) allows, under particular circumstances, simultaneous esti-
687 mation of the uncertainty of all three products and also (partly) isolation of random (relative)
688 representativeness errors (*Miralles et al., 2010; Gruber et al., 2013a; Chen et al., 2017*). Note,
689 however, that the necessity of using two reference data sets instead of one may limit spatial
690 and temporal data availability. Moreover, while non-orthogonal and mutually correlated er-
691 rors are equally problematic for metrics that rely on one reference data set only (see below), it
692 may be even more difficult to find a third data set that fulfills these requirements. Commonly,
693 any combination of in situ measurements, land surface model estimates, active-microwave-based

694 retrievals, or passive-microwave-based retrievals is expected to fulfil this requirement because
695 their sources of errors are assumed to be mostly independent (*Gruber et al.*, 2016a), provided
696 that neither of them has been used to generate another (e.g., by assimilating satellite data in
697 to a land surface model; *Reichle et al.*, 2017a,b). However, several studies suggest that mutual
698 error correlations may exist between commonly used data set combinations (*Yilmaz and Crow*,
699 2014; *Pan et al.*, 2015), resulting from representativeness errors (e.g., if a land surface model
700 used within TCA models a deeper layer than the sensing depth of two satellite data sets that
701 are used in the triplet) or from unrecognized common data. Examples for the latter can be
702 found in some SMOS and SMAP products, which use modelled temperature estimates from
703 ECMWF’s Integrated Forecast System (IFS) and NASA’s Goddard Earth Observing System
704 Model, version 5 (GEOS-5), respectively, as input to the soil moisture retrieval algorithm (*Kerr*
705 *et al.*, 2012; *O’Neill et al.*, 2018). Research is needed to quantify the degree to which that af-
706 fects inter-comparisons between the satellite soil moisture retrievals and soil moisture estimates
707 from models that rely on the same temperature input (such as MERRA2, ERA-Interim/Land,
708 or others; e.g. *Chen et al.*, 2018). It is therefore recommended to verify orthogonality and
709 zero error correlation assumptions by using - where available - multiple data set triplets and
710 checking for consistency between different TCA implementations (*Dorigo et al.*, 2010; *Draper*
711 *et al.*, 2013), or by using the recently proposed TCA extension that utilizes four or more data
712 sets to diagnose the existence, and estimate the magnitude of error correlations (*Gruber et al.*,
713 2016b; *Pierdicca et al.*, 2017).

714 The following sections discuss the most common bias and uncertainty metrics, either (i)
715 based on direct comparison between two data sets, which will be referred to as relative metrics,
716 or (ii) based on the simultaneous comparison of three products, which will be referred to as
717 TCA-based metrics. All metrics can be equally applied to soil moisture anomaly estimates or
718 the raw time series, except for first-order bias estimators (see below) as the anomaly calculation
719 per definition removes differences in the mean (see Sec. 3.3.3).

720 Note that none of the metrics presented below require assumptions about the shape of the
721 pdf of the random errors or the true signal (*McColl et al.*, 2016). However, the bounded nature
722 of soil moisture may cause violations in the orthogonality assumption if cut-off values (e.g., zero
723 and the soil porosity as lower and upper physical limit, respectively) are applied to the soil
724 moisture estimates of a particular data sets. Especially in very dry or very wet regimes, where
725 random errors would often cause these thresholds to be exceeded, this can result in considerable

726 biases in all (both relative and TCA-based) uncertainty metrics.

727 **3.4.3 Bias estimation**

728 Bias estimation is only meaningful against reference data at the satellite footprint scale, i.e.
729 without considerable representativeness errors and if no rescaling has been applied (see Sec.
730 3.3.2).

731 **Temporal mean bias**

732 Bias estimates are commonly based on the (temporal) mean difference between two data sets
733 (*Entekhabi et al., 2010a*):

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

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

744 **Second-order bias**

745 Most validation studies do not attempt to estimate second-order biases and neglect their impact
746 on b_{xy} and other validation metrics such as the (unbiased) Root-Mean-Square-Difference (see
747 *Gupta et al. (2009)* and Sec. 3.4.4). TCA potentially allows for the direct estimation of second-
748 order biases (*Gruber et al., 2016a*) as:

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

749 where β_x^y denotes the TCA-based second-order bias estimate of x relative to y which, if y is
750 a fiducial reference data set and if no non-orthogonal or correlated random representativeness
751 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
752 β_x . Notice that neither first nor second-order biases in z influence β_x^y . Alternatively, Eq. (5)
753 can also be used for rescaling purposes (*Yilmaz and Crow, 2013; Su et al., 2014; Gruber et al.,*
754 *2016a*, see Sec. 3.3.2).

755 3.4.4 Uncertainty estimation

756 As discussed, uncertainty estimates aim at representing the pdf of the random errors (see Sec.
757 1.1), which is typically done by means of their standard deviation (or variance).

758 (Unbiased) Root-Mean-Square-Difference

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

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

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

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

764 The common definition of the ubRMSD specifically corrects for differences between the mean of
765 the data sets (*Entekhabi et al., 2010a*). However, as can be seen in Eq. (7), it remains susceptible
766 to second-order biases, which are amplified by the actual soil moisture variability (σ_t^2). Moreover,
767 as was the case for b_{xy} , this second-order bias dependency in $ubRMSD_{xy}$ persists even when
768 calculated against fiducial reference data, in which case Eq. (7) collapses to $\sqrt{(\beta_x - 1)^2\sigma_t^2 + \sigma_{\xi_x}^2}$.
769 As discussed in Sec. 3.3.2, data sets are often rescaled before calculating validation metrics to
770 account for systematic representativeness errors, especially when evaluating against data from
771 sparse networks. This is most commonly done by matching the temporal mean and the standard
772 deviation of the data sets, or their entire cdf (i.e. also higher statistical moments). However, as

773 can be seen from Eq. (3), this only properly corrects for relative differences in β if the SNRs
774 (including random representativeness errors) of the data sets are equal, which is very unlikely.
775 Consequently, Eq. (7) will still contain the remaining difference between β_x and the rescaled β_y ,
776 multiplied with the actual soil moisture variability, and also random representativeness errors.

777 (Unbiased) Root-Mean-Square-Error

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

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

783 Note that when calculating the ubRMSE using the cross-multiplied differences instead of the
784 statistical moments, the data sets y and z do have to be bias-corrected with respect to x a priori
785 using Eqs. (4) and (5). The absolute value is taken to prevent negative signs in uncertainty
786 estimates that could occur due to sampling errors (*Gruber et al., 2018*, see Sec. 3.5). As one
787 can see, $ubRMSE_x$ is (as opposed to $ubRMSE_{xy}$ in Eq. (7)) fully unbiased in that it contains
788 neither first nor second-order biases from both the satellite and the reference data sets, and it
789 also no longer contains the uncertainties inherent in the reference data products (*Gruber et al.,*
790 *2016a*). However, estimates that are unbiased *with respect to the assumed true grid cell average*
791 can only be obtained if at least one fiducial reference data set is available (*Chen et al., 2017*).
792 Moreover, $ubRMSE_x$ is not affected by random representativeness errors in y and z as long as
793 they are orthogonal and not correlated. Such representativeness error correlations could occur
794 for example when applying TCA to in situ measurements together with two coarse resolution
795 products. This case, however, provides an opportunity to estimate the representativeness of in
796 situ stations while uncertainty estimates for the coarse resolution products remain unaffected
797 (*Miralles et al., 2010; Gruber et al., 2013a; Chen et al., 2017*). For a more detailed derivation
798 of how representativeness errors affect the TCA-based uncertainty estimates we refer the reader
799 to *Vogelzang and Stoffelen (2012)* and *Gruber et al. (2016a)*.

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

806 **Pearson correlation coefficient**

807 The most common SNR-related relative metric is the linear (Pearson) correlation coefficient,
 808 which is typically described as a measure for statistical dependency between two data sets.
 809 From the error model in Eq. (3) one can see that it is also a direct, normalized (between -1
 810 and 1) representation of the SNRs of the two data sets for which it is calculated (*Gruber et al.*,
 811 2016a):

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

812 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
 813 against fiducial reference data, R_{xy} is a direct representation of the SNR of the satellite under
 814 evaluation (i.e. SNR_x). Notice that the “signal” to which the “noise” in the SNR estimator is
 815 related is the true soil moisture variability scaled with the second-order satellite bias (i.e. $\beta_x^2 \sigma_t^2$).
 816 Even if β_x could be estimated reliably, for example from Eq. (5), rescaling does not change
 817 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
 818 quantity of interest that determines how well signal variations can be distinguished from noise,
 819 regardless of whether systematic errors have been corrected for (*Gruber et al.*, 2016a), which can
 820 be also interpreted as the (linear) correlation with the true soil moisture signal (*McCull et al.*,
 821 2014). When R_{xy} is calculated against non-fiducial reference data, it is additionally influenced
 822 by second-order systematic and random representativeness errors as well as the uncertainties
 823 of that reference data set. Note that the Pearson correlation coefficient is sometimes presented
 824 squared (R_{xy}^2), referred to as coefficient of determination and interpreted as “percentage of
 825 variance explained”, which provides a slightly more intuitive link to the SNR and may hence

826 be preferable, even though the information content is identical.

827 TCA-based correlation coefficient

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

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

830 As was the case for the ubRMSE, the validity of Eq. (10) requires that there is no correlation or
831 non-orthogonality between random representativeness errors, but their individual variance may
832 well be non-zero. If these assumptions are respected, then R_x will be an unbiased representation
833 of the correlation between x and the (unknown) hypothetical truth. Consequently, R_x will
834 always be larger than R_{xy} although this difference decreases as the quality of the reference y
835 increases. Note, however, that R_x only ranges between 0 and 1, as an anti-correlation (with
836 respect to the true signal) cannot be unambiguously inferred from the three covariances in Eq.
837 (10). To provide a more intuitive link to the SNR, R_x may also be presented squared (i.e. as
838 TCA-based coefficient of determination; R_x^2).

839 (Logarithmic) Signal-to-Noise Ratio

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

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

842 This provides a more direct, linear representation of the ratio between soil moisture and uncer-
843 tainty magnitude than R_x , yet the information content in both metrics is identical; it is simply
844 a different way of presentation. Note that the SNR_x is already being used as a more coher-
845 ent (than RMSD or RMSE based metrics) satellite data quality indicator for defining target
846 accuracy requirements (see Sec. 3.8.2).

847 3.5 Statistical significance testing

848 All the above described (and also most other less common) validation metrics are based on
849 statistical moments, sampled in time. Since these estimates are based on finite samples (i.e.
850 the discrete soil moisture time series), they are subject to sampling errors. The most common
851 way to deal with statistical uncertainty (i.e. sampling errors) across science communities is
852 Null Hypothesis Significance Testing (NHST) using p -values and/or confidence intervals (*Wilks*,
853 2011). In a validation context, typical hypotheses to be nullified are, for example, that a soil
854 moisture product does not meet a target accuracy threshold or that one product does not
855 exhibit higher correlation with a reference product than another. For testing such hypotheses,
856 the sampling distribution of the statistical estimate under consideration (such as a validation
857 metric) is constructed based on the magnitude of the estimate and the size of the sample used to
858 draw this estimate (see below). Then, either the p -value is calculated, which is the probability
859 of values of the sampling distribution to be equal to or below (or above, depending on which tail
860 is considered) the pre-defined Null-value (representing the Null hypothesis), or the $(1 - \alpha) \cdot 100\%$
861 confidence interval is considered. A rejection of the Null-hypothesis is considered statistically
862 significant, if the p -value is below a pre-defined significance level α (typically 0.05) or if the
863 $(1 - \alpha) \cdot 100\%$ confidence interval does not contain the Null-value. When comparing estimates
864 of different samples (e.g., the performance of different soil moisture products), it is common to
865 consider their relative difference as statistically significant if their confidence intervals do not
866 overlap. Note that the term “Null-value” refers to the Null hypothesis and not to a value of zero
867 of the test statistic (i.e. the validation metric). A common (yet inappropriate; see Sec. 3.8.2)
868 Null-value for testing soil moisture accuracy requirements, for instance, is $0.04 \text{ m}^3\text{m}^{-3}$ ubRMSD
869 . Hence, if the p -value for $0.04 \text{ m}^3\text{m}^{-3}$ of the sampling distribution around an estimated ubRMSD
870 is below the defined α level, the product is said to meet accuracy requirements with statistical
871 significance.

872 However, the American Statistical Association (ASA) has recently issued a statement on sta-
873 tistical significance and p -values (*Wasserstein and Lazar, 2016*) warning about the science-wide
874 misuse and abuse of NHST through the replacement of scientific reasoning with a dichotomous
875 and arbitrary classification of results into “significant” or “non-significant”. In this statement,
876 the ASA is advocating the abandonment of statistical significance testing altogether for two
877 main reasons. The first one is that an alarming fraction of articles in the scientific literature
878 present unjustified inferences based on misinterpreted p -values and confidence intervals (*Green-*

879 *land et al.*, 2016; *Gelman and Stern*, 2006; *Wasserstein and Lazar*, 2016). The second and more
880 important argument is that p -values alone provide no grounds for meaningful decision making.
881 While the magnitude of p itself can be informative about how consistent the data at hand are
882 with an assumed stochastic model, “[...] a label of statistical significance does not mean or imply
883 that an association or effect is highly probable, real, true, or important. Nor does a label of
884 statistical nonsignificance lead to the association or effect being improbable, absent, false, or
885 unimportant.” (*Wasserstein et al.*, 2019). Therefore, no practical conclusion or decision should
886 be based on whether p -values do or do not meet an arbitrarily defined threshold. Instead of
887 strictly yet arbitrarily categorizing study results based on dichotomous significance tests, one
888 should strive for more careful study design and more rigorous understanding, interpretation
889 and reporting of the stochastic properties of the data at hand (*Greenland et al.*, 2016; *Tong*,
890 2019). Note that the same can be said for an arbitrarily defined target accuracy threshold of
891 $0.04 \text{ m}^3\text{m}^{-3}$, which is often used to declare a product - without any solid grounds - as “valid”
892 or “invalid” (see Sec. 3.8.2 and Sec. 5).

893 In conclusion, for soil moisture validation purposes, we follow the guidance of the ASA and
894 recommend to avoid any statement or interpretation about statistical “significance” or “non-
895 significance”, and to instead always provide and interpret a statistical summary of calculated
896 validation metrics in the form of confidence intervals alongside the metrics themselves. How
897 confidence intervals can be calculated and recommendations of how they can be presented are
898 provided in the following sections.

899 **3.6 Confidence intervals**

900 In general, confidence intervals represent the pdf of the sampling errors of an estimate and
901 are defined at a certain confidence level. A confidence level of, say, 95% means that if one
902 would repeatedly calculate 95% confidence intervals in a series of similar experiments, then 95%
903 of them would - on average - contain the true value, provided that all assumptions made for
904 the stochastic model are met. Note that this is *not* the probability that the true value that
905 is approximated by the estimate lies within the confidence interval (*Neyman*, 1937; *Greenland*
906 *et al.*, 2016). In theory, this probability - which would indeed be more informative - could be
907 represented by a Bayesian credible interval, but calculating it would require a priori knowledge
908 about the pdf of the parameter that is being estimated (i.e. the so-called “prior”) and this is
909 typically not available.

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

917 3.6.1 Analytical calculation

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

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

922 and

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

923 where n is the sample size; \bar{u} and s_u represent the sample mean and standard deviation of the
 924 difference series ($x - y$); and μ_u and σ_u are their corresponding true population parameters.
 925 The population moments of u are estimated within the $(1 - \alpha) \cdot 100\%$ confidence intervals as
 926 a function of the sample moments of u . Specifically, the confidence intervals (CI) for b_{xy} and
 927 $ubRMSD_{xy}$ can be inferred from Eqs. (12) and (13) as:

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

928 and

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

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

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

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

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

935 The confidence interval for the coefficient of determination (R_{xy}^2) can be derived by simply
 936 squaring the confidence interval of R_{xy} in Eq. (17).

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

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

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

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

949 with ρ_x and ρ_y obtained from a fitted AR(1) model as:

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

950 where $i \in [x, y]$, τ_i is the fitted persistence time of the individual time series x and y , i.e. the time
951 lag at which the auto-correlation drops below $1/e$, and d_m is the median time distance between
952 consecutive valid, collocated observations, i.e. the lag-1 distance accounting for the typically
953 irregular spacing between satellite retrievals. Note that averaging correlation coefficients is
954 generally not recommended (see Sec. 3.7), but required here to determine a single effective
955 proxy of the auto-correlation of collocated data pairs with possibly deviating individual memory.
956 Using the geometric average avoids the dominance of data sets with large auto-correlation (e.g.,
957 land surface models often have a different memory than satellite observations), which may cause
958 excessively large confidence intervals.

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

964 **3.6.2 Bootstrapping**

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

972 Bootstrapping is a special case of Monte Carlo simulation, which uses the sample itself
973 as approximation of the population. More specifically, it constructs an empirical probability
974 distribution of the test statistic (in our case the validation metric) by resampling the original
975 sample multiple times, with replacement to preserve the sample size, and repeated calculation
976 of the test statistic from those resamples. This bootstrapped distribution then allows for the

977 direct derivation of confidence intervals as well as other parameters of the sampling error pdf.
978 The advantages of this method lie in its algorithmic simplicity and that it can be applied
979 to any metric without the need to assume a particular sampling distribution (such as t or
980 χ). However, bootstrapping confidence intervals requires a considerable number of resamples,
981 which may lead to large computational costs, and relies on the assumption that the sample is
982 indeed a reliable representation of the population, which requires a large sample size. A general
983 recommendation for bootstrapping confidence intervals is to use a minimum of 1000 resamples
984 (*Efron and Tibshirani, 1986*). However, the number of required resamples may be chosen more
985 specifically for a given study by testing for convergence of the results with increasing sample
986 size. For example, *Draper et al. (2013)* used 1000 resamples for estimating confidence intervals
987 for TCA-based *ubRMSE* estimates, although their testing found that 500 would have been
988 sufficient.

989 As was the case for the analytical expressions, bootstrapped confidence intervals are also
990 susceptible to auto-correlation in the data. This can be accounted for by resampling blocks of
991 data instead of single data points, referred to as block-bootstrapping (*Ólafsdóttir and Mudelsee,*
992 2014), which preserves the auto-correlation properties of the original sample. An estimate of the
993 optimal block length (l_{opt}) for bootstrapping CIs around TCA-based estimates can be obtained
994 following *Chen et al. (2018)* as:

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

995 where $\text{NINT}\{\cdot\}$ denotes rounding to the nearest integer. As before, a single effective value for
996 ρ can be obtained as the geometric average of the lag-1 auto-correlations of the three data sets
997 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
998 interval between consecutive valid, collocated data triplets. To prevent data gaps from causing
999 an auto-correlation degradation during the resampling, we recommend to discard data blocks
1000 from the resamples if they contain less than 50% of valid data.

1001 3.7 Summary statistics

1002 Validation metrics and their confidence intervals should be calculated and assessed over a wide
1003 range of spatial locations to understand error characteristics of a soil moisture product under
1004 different climatic, topographic and land cover conditions. However, it may be practical to

1005 summarize spatially distributed skill estimates into a single combined metric (for example to
 1006 obtain an overall ranking of different products or to track the performance evolution of a product
 1007 over time), which requires also the aggregation of their associated confidence intervals.

1008 3.7.1 Averaging metrics

1009 The most common way of obtaining a combined skill estimate is arithmetic averaging:

$$\bar{\nu} = \mathbf{w}^\top \mathbf{v} \quad (22)$$

1010 where $\bar{\nu}$ is the average of k spatially distributed skill metrics that are summarized in the skill
 1011 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
 1012 individual skill estimates with $\sum w_i = 1$. Averaging skill metrics in a weighted fashion to
 1013 minimize the impact of sampling errors is in principle possible by deriving weights from the
 1014 sampling error magnitudes (*Aitkin*, 1936), but in most cases, an unweighted average is preferred
 1015 because validation points are typically selected to represent a wide range of varying conditions,
 1016 and areas with lower sampling errors (i.e. regions with better temporal coverage, for instance
 1017 because less data are masked out) could dominate a weighted averaged skill estimate. For such
 1018 unweighted average, the weight vector takes the form $\mathbf{w} = [k^{-1} \cdots k^{-1}]^\top$.

1019 While many metrics can be averaged safely, it is - against common practice - not recom-
 1020 mended to average correlation coefficients (neither Pearson nor TCA-based) because they are
 1021 calculated as ratios using standard deviations (variances) and covariances or SNRs (see Eqs. (9)
 1022 and (10)). Therefore, they behave highly non-linearly and neither an average of these ratios nor
 1023 a ratio of averaged numerators / denominators would allow for a meaningful inference about
 1024 statistical properties. For example, averaging correlation coefficients of 0.1 and 0.9, which cor-
 1025 respond to a SNR of 0.01 and 4.26, respectively (in the case of Pearson correlation assuming
 1026 a random error-free reference data set), would lead to an average correlation of 0.5 with an
 1027 associated SNR of 0.33. This is far from their average SNR of 2.14 (ignoring for the moment
 1028 that this too is an average of ratios) which would correspond to a correlation coefficient of 0.83.
 1029 In contrast, correlation coefficients of 0.3 and 0.7, representing SNRs of 0.1 and 0.96, respec-
 1030 tively, would have the same average correlation yet the average of their associated SNR is 0.53,
 1031 corresponding to a correlation of 0.59. Moreover, the skewed probability distribution of the
 1032 Pearson correlation coefficient causes the arithmetic average to be systematically biased. Some

1033 studies suggest to average Fisher-transformed z -values instead (*Corey et al.*, 1998), which have
 1034 a Gaussian sampling distribution, but a back-transformed z -average is just as difficult to inter-
 1035 pret. Following the above example, averaging correlation coefficients of 0.1 and 0.9 in z -space
 1036 would lead to an average correlation (or more precisely, an inverse average- z) of 0.66 (SNR =
 1037 0.76), whereas when averaging z -transformed correlations of 0.3 and 0.7, it would be 0.53 (SNR
 1038 = 0.39).

1039 In other words, the choice of whether to average correlation coefficients, Fisher-transformed
 1040 z -values, or SNRs - albeit representing the exact same uncertainty properties - will lead to dif-
 1041 ferent values and hence interpretations of the resulting average and this difference also depends
 1042 on the degree of variability across the estimates that are being averaged. Moreover, the resulting
 1043 average number (regardless of the approach) no longer represents an actually meaningful sta-
 1044 tistical property. Alternatively, instead of averaging pre-calculated correlation coefficients, one
 1045 may be tempted to calculate the correlation coefficient directly over the concatenated measure-
 1046 ments or estimates of all available locations to obtain an overall skill estimate. However, this is
 1047 not meaningful as the effects of different populations are lumped together. As a consequence,
 1048 for example, two data sets that individually exhibit strong positive correlation in a wet and in
 1049 a dry soil moisture regime, respectively, may appear to have an overall weak anti-correlation
 1050 when put together, an effect also known as Simpson’s paradox (*Blyth*, 1972). Therefore, such
 1051 an approach should be strictly avoided.

1052 3.7.2 Averaging confidence intervals

1053 The uncertainty in the spatially averaged skill metric in Eq. (22) associated with the *sampling*
 1054 errors of the individual skill estimates can be calculated through the standard method for the
 1055 propagation of uncertainty as:

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

1056 where $s_{\bar{v}}^2$ is the sampling uncertainty in the averaged skill \bar{v} (i.e. its sampling error variance);
 1057 and $\mathbf{\Sigma}$ is the sampling error covariance matrix for the k individual skill estimates. The corre-
 1058 sponding aggregated confidence intervals can be derived from a Gaussian distribution (which
 1059 will generally be assured by the Central Limit Theorem for reasonably large samples) with mean
 1060 \bar{v} and standard deviation $s_{\bar{v}}$.

1061 Diagonal elements in Σ are the sampling error variances of the individual skill estimates, i.e.
 1062 $diag(\Sigma) = \mathbf{s}^2$ 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
 1063 standard errors of the sample mean and sample variance (of the difference series $u = x - y$ at
 1064 each individual location), respectively:

$$s_{b_{xy}}^2 = \frac{ubRMSD_{xy}^2}{n} \tag{24}$$

$$s_{ubRMSE_{xy}}^2 = \frac{ubRMSD_{xy}^2}{2(n-1)}$$

1065 For TCA-based metrics, the sampling error variance can be directly calculated from the boot-
 1066 strapped sampling distribution.

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

1067 where \circ denotes the Hadamard product, i.e. element-wise matrix multiplication. \mathbf{R} differs for
 1068 the various skill metrics. For b_{xy} and $ubRMSE_{xy}$, it is the *spatial* auto-correlation matrix of the
 1069 difference series u , and of the squared, bias-corrected difference series $(u - \bar{u})^2$, respectively, at
 1070 the different locations u where skill metrics are calculated. For TCA-based metrics, the sampling
 1071 error covariance can be calculated as the covariance between the bootstrapped samples (*Gruber*
 1072 *et al.*, 2019b), provided that the order in which bootstrap-resamples are drawn is the same at
 1073 all different locations, which may be difficult when using block-bootstraps with different block-
 1074 length.

1075 Earlier research (*De Lannoy and Reichle*, 2016) has proposed a clustering approach to take
 1076 possible sampling error correlations into account. This approach first calculates mean metrics
 1077 and confidence intervals per spatial cluster, assuming that the sampling errors of the spatially
 1078 close data sets within each cluster are perfectly correlated. Next, averaged skill metrics and con-
 1079 fidence intervals from within the clusters are averaged, assuming that all clusters are completely
 1080 independent. However, this approach is expected to overestimate confidence intervals because:
 1081 (i) sampling errors will never be perfectly correlated unless validation metrics are calculated
 1082 multiple times from the exact same data, and (ii) clusters are formed based on the expected
 1083 auto-correlation length of the soil moisture data sets, which will be much larger than that of the
 1084 difference series between data sets, as required in Eq. (25).

1085 Finally, although averaging of some metrics and confidence intervals is possible, we generally
 1086 recommend to retain detailed information about their spatial variability, and to leverage this

1087 information to obtain a better understanding of product performance and its relation to land
1088 cover, topography, climate, and other possibly important factors. If point-wise assessments are
1089 not feasible or if simple product summaries are desired, percentile statistics such as medians
1090 and inter-quartile-ranges (of both calculated skill estimates and their confidence intervals) are
1091 generally more informative than spatial averages and their increasingly inaccurate averaged
1092 confidence intervals. More specific recommendations of how validation metrics and confidence
1093 intervals can be presented are provided in Sec. 4 and Appendix A.

1094 **3.8 Practical remarks**

1095 **3.8.1 Validating downscaled products**

1096 Currently, most space-borne microwave sensors available for soil moisture retrieval operate at
1097 spatial resolutions of about $25^2 - 50^2$ km² (*Gruber et al.*, 2019a). Some higher-resolution Syn-
1098 thetic Aperture Radar (SAR) sensors exist that allow for reasonable soil moisture retrieval at
1099 scales up to approximately 1 km² (*Pathe et al.*, 2009; *Gruber et al.*, 2013b), yet with consider-
1100 ably lower temporal resolution and accuracy. In addition, many downscaling approaches have
1101 been developed to improve the spatial resolution of coarse-resolution soil moisture products,
1102 e.g., by fusing coarse-resolution radiometer or scatterometer measurements with high-resolution
1103 SAR data (*Das et al.*, 2017; *Bauer-Marschallinger et al.*, 2018), by fusing microwave observa-
1104 tions with optical/thermal measurements (*Chauhan et al.*, 2003), or through data assimilation
1105 (*Reichle et al.*, 2017b). For a comprehensive review of downscaling methods see *Peng et al.*
1106 (2017).

1107 The validation of downscaled products is mostly done as for coarse-resolution products, i.e.
1108 through time series analysis with a focus on temporal dynamics at individual locations (see
1109 Sec. 3). In doing so, it has been shown that the downscaling process often actually decreases
1110 the temporal performance of the products, that is, the original coarse-resolution products often
1111 correlate better with local soil moisture dynamics, even at a point scale, than their downscaled
1112 counterparts (*Peng et al.*, 2015). While downscaled soil moisture images provide more visual
1113 level-of-detail, only few studies have quantitatively assessed whether the obtained spatial pat-
1114 terns actually represent real soil moisture variations (e.g., *Bauer-Marschallinger et al.*, 2018;
1115 *Sabaghy et al.*, in review) or whether they are just mimicking spatial patterns of ancillary data
1116 such as soil texture maps (for a comprehensive review of validation studies for downscaled prod-
1117 ucts see *Peng et al.*, 2017).

1118 Therefore, we highly recommend that future validation studies for downscaled products
1119 put a strong emphasis on assessing also the spatial soil moisture variations obtained from the
1120 downscaling, e.g., by estimating spatial correlation coefficients (*Sahoo et al.*, 2013; *Kolassa et al.*,
1121 2017; *Sabaghy et al.*, in review), in addition to time series analyses. To that end, we further
1122 encourage the setup of field campaigns and validation sites dedicated to support such high-
1123 resolution validation activities, especially in regions where soil moisture variations are very
1124 heterogeneous.

1125 **3.8.2 Target accuracy requirements**

1126 Satellite soil moisture validation studies most commonly evaluate products against a target ac-
1127 curacy threshold of $0.04 \text{ m}^3\text{m}^{-3}$ ubRMSD across the globe, excluding regions of snow and ice,
1128 frozen ground, complex topography, open water, urban areas, and vegetation with water content
1129 greater than 5 kg/m^2 . This requirement was defined by the Soil Moisture and Ocean Salinity
1130 (SMOS; *Kerr et al.*, 2001) and the Soil Moisture Active Passive (SMAP; *Entekhabi et al.*,
1131 2010a) missions, and by the Terrestrial Observation Panel for Climate (TOPC; *WMO*, 2016).
1132 Alternatively, the Satellite Application Facility in Support to Operational Hydrology and Wa-
1133 ter Management (H SAF) of the European Organisation for the Exploitation of Meteorological
1134 Satellites (EUMETSAT) has defined (TCA-based) SNR product requirements (*H-SAF*, 2017)
1135 for the operational soil moisture products that are retrieved from measurements of the Advanced
1136 Scatterometer (ASCAT) onboard the MetOp satellites (*Naeimi et al.*, 2009). In particular, the
1137 EUMETSAT H SAF defines 0, 3 and 6 dB SNR as threshold, target and optimal SNR require-
1138 ments to make product assessment possible on a larger scale and spatially better comparable
1139 (see Sec. 3.4).

1140 Both of these requirements are based on relatively practical, easy-to-estimate single numbers
1141 that represent a rough estimate of what is currently achievable rather than being an indication
1142 of “good” or “bad” product quality. While they provide easy means to monitor product perfor-
1143 mance evolution over time and to compare products, they are entirely unrelated to the suitability
1144 of a product for specific applications. However, the actual specification of bias and uncertainty
1145 requirements for the fitness-for-purpose for a particular application (including the specification
1146 of the appropriate metrics) is a task of the respective user community and urgently requires
1147 further research (*Entekhabi et al.*, 2010b), because no data set can be declared “valid” if no
1148 validity requirements are available.

1149 3.8.3 Reproducibility

1150 The research community generally suffers from a lack of reproducibility in scientific studies
1151 (*Baker, 2016*). Also in soil moisture validation studies, contradictory results for the performance
1152 and relative ranking between different satellite products have been reported (e.g., *Wagner et al.,*
1153 2014). These ambiguities originate from: (i) the choice of reference data and product versions;
1154 (ii) the use of different spatial regions and time periods; (iii) different approaches used for data
1155 preparation and pre-processing; (iv) statistical sampling errors; and (v) software implementation
1156 errors. Note, however, that contradicting results are not necessarily caused by bad study design
1157 but often originate from stochastic uncertainties, which are inevitably dominant in space borne
1158 Earth observation measurements and retrieval algorithms (*Greenland et al., 2016*).

1159 Embracing statistical uncertainty and developing an in-depth understanding of soil moisture
1160 product quality requires more comprehensive descriptions of data sets, software, and methodol-
1161 ogy than are usually provided as well as the mandatory, additional estimation and presentation
1162 of sampling errors. To that end, we recommend that:

- 1163 • all validation results should be accompanied by confidence intervals as measure for sam-
1164 pling errors;
- 1165 • all methodological steps should be described with sufficient detail to be reproducible;
- 1166 • all data sets used for the study should be made publicly available and unambiguously
1167 identifiable by providing their exact product version information and, where available,
1168 their Digital Object Identifier (DOI);
- 1169 • all used software packages that are relevant for the exact reproduction of validation results
1170 should be referenced with their complete version number and, where available, their DOI.
1171 If not accessible via open repositories (in particular software specifically designed for that
1172 study), we recommend to make source code publicly available, for example on GitHub
1173 (<https://github.com/>; last access: 1 July 2019).

1174 A list of some current publicly available software that is specifically aimed at, or closely related to
1175 soil moisture validation is provided in Table 3. An online validation tool that is built around these
1176 software packages and follows the good practice guidelines presented in this paper is provided by
1177 the Quality Assurance Framework for Soil Moisture (QA4SM; <https://qa4sm.eodc.eu/>; last
1178 access: 1 July 2019).

1179 Note that the re-distribution of in situ measurements (see the third point above) may be
1180 particularly problematic as many networks do not operate for free. Requiring networks to
1181 freely distribute their data will likely decrease the number of datasets available for validation
1182 activities, which may ultimately hamper the evolution of satellite soil moisture products and
1183 downstream products derived thereof. We therefore emphasize the tremendous value of ground
1184 reference measurements and encourage the community to support, by any means possible, the
1185 development and continuation of operational Cal/Val sites.

1186 4 Validation Good Practice Protocol

1187 This section provides a compilation of the theoretical considerations presented above in the form
1188 of a validation good practice protocol for satellite soil moisture products, i.e. guidelines for:

- 1189 • the selection of reference data;
- 1190 • data pre-processing steps;
- 1191 • the selection and implementation of appropriate metrics; and
- 1192 • the presentation of validation results.

1193 Figure 3 illustrates the process and Appendix A provides an example that follows these recom-
1194 mendations. We stress that there is no one-size-fits-all approach for validating Earth observation
1195 data. Depending on the application in question, several analyses may not be necessary. Also,
1196 recommended thresholds may need to be adjusted depending on data quality requirements (e.g.,
1197 more strict data masking procedures may be employed) or data availability (e.g., the allowed in
1198 situ measurement depth may be increased if only retrievals from long wavelengths in dry and
1199 sandy regions are used).

1200 4.1 Data selection

1201 As discussed in Sec. 2, no reference data source provides a sufficiently accurate and traceable soil
1202 moisture proxy for reliable error assessment on a global scale. A complete and comprehensive
1203 product validation therefore requires comparisons against each of the following (*Jackson et al.*,
1204 2012): (i) dense networks, in particular core validation sites; (ii) sparse networks; (iii) land
1205 surface model output; and (iv) other satellite products, always making sure that the latest or
1206 most recommended product versions are used. However, given the large number of satellite and

1207 reference products available, a complete analysis that considers all these data sources is typically
1208 beyond the capacity of a single validation study. Therefore, separate studies may be conducted
1209 for dense network evaluation (*Colliander et al., 2017a*), sparse network evaluation (*Dorigo et al.,*
1210 *2015; Chen et al., 2017*), or coarse-resolution product inter-comparison (*Al-Yaari et al., 2014;*
1211 *Burgin et al., 2017; Chen et al., 2018*) and their results compiled together.

1212 Since satellite soil moisture retrievals represent only the top few centimeters of the soil, in
1213 situ sensors and modelled soil layers used for validation should reach no deeper than 5-10 cm,
1214 which is considered as the maximum sensing depth for currently available microwave wavelengths
1215 (X-band to L-band). Information where currently publicly available reference data sets can be
1216 accessed is provided in Table 2.

1217 4.2 Pre-processing

1218 4.2.1 Masking

1219 In situ measurements and satellite retrievals should be masked out when considered unreliable.
1220 Recommendations from data providers regarding product inherent quality flags should be fol-
1221 lowed and the employed thresholds carefully documented. Additionally, we recommend using
1222 ancillary data to mask out pixels classified as tropical forests, water bodies, wetlands, and unin-
1223 dation areas as well as all measurements on days with non-zero snow indicators (e.g., snow height
1224 or snow-water-equivalent), or surface or soil temperature below 4°C. Such ancillary data can be
1225 supplied by land surface models or complementary satellite data. When biases or uncertainties
1226 of multiple products are compared, they should be calculated from the exact same, collocated
1227 data points. However, care should be taken that single products with poor data coverage do not
1228 distort the overall assessment (see Sec. 5).

1229 To avoid excessively large confidence intervals that can hamper meaningful data comparison,
1230 grid cells with less than 50-500 collocated data points may be masked out depending on data
1231 availability (*Zwieback et al., 2012*). Also, many studies mask out correlation coefficients based
1232 on Student's t-test (i.e. applying p-value thresholds for correlation coefficients), and/or bias and
1233 uncertainty estimates based on vegetation density (e.g., vegetation water content $> 5 \text{ kg/m}^2$)
1234 or other thresholds (e.g., open-water fraction > 0.05) (*Dorigo et al., 2010; Brocca et al., 2011;*
1235 *Al-Yaari et al., 2014*). However, carefully reporting and interpreting confidence intervals and
1236 sample sizes at locations with low data coverage could indeed provide valuable additional insight
1237 and may be more informative than masking out estimates completely (*Wasserstein et al., 2019*).

1238 Also, complete reporting of results prevents generating publication biases due to “cherry-picking”
1239 which is sometimes found in the scientific literature (*Greenland et al.*, 2016).

1240 **4.2.2 Collocation**

1241 Spatial collocation requires the selection of a spatial comparison grid, which is often the grid
1242 of the satellite product under validation. In situ measurements should be assigned to the grid
1243 cell in which they are located. For dense networks, all stations that lie within a particular grid
1244 cell should be averaged, if possible taking their respective spatial representativeness for that
1245 grid cell into account. To avoid artificial jumps due to sensor drop-outs, only time steps where
1246 all stations provide valid measurements should be considered. For the SMAP core validation
1247 sites (see Sec. 2.2.1), a validation grid that minimizes upscaling errors has been developed as
1248 described in *Colliander et al.* (2017a).

1249 Gridded reference products (i.e. other satellite and land surface model products) should be
1250 resampled onto the chosen comparison grid, e.g., using a Nearest Neighbor (NN) search. If the
1251 grid resolution of the reference product is coarser than that of the comparison grid, individual
1252 grid cells of that product may be assigned to multiple comparison grid cells. If the grid resolution
1253 is much finer, all NNs of single comparison grid cells (in case more than one exist) should be
1254 averaged, if possible taking spatial representativeness into account.

1255 Temporal collocation at comparison time steps should minimize the time difference between
1256 data match-ups and be based on a NN-search with a maximum time difference threshold of
1257 1-12 hours, depending on data availability. Note that the choice of the comparison grid and
1258 time steps may affect the presence and distribution of (spatial and temporal) representativeness
1259 errors among the considered data sets (see Sec. 5).

1260 **4.2.3 Decomposition**

1261 All validation metrics should be calculated for the raw soil moisture time series (of collocated
1262 retrievals and reference data) as well as for short-term and long-term anomalies, except for
1263 temporal mean biases whose calculation is trivial for anomalies. Short-term anomalies should
1264 be estimated as residuals from a seasonality that is computed by applying a 4-8 week moving
1265 average window to the time series. Long-term anomalies should be estimated as residuals from
1266 a climatology that is computed by averaging the measurements or estimates of all years within a
1267 4-8 week moving window around each DOY, but only if at least 5-10 years of data are available.

1268 To avoid data-density related artefacts, especially in the transition periods from frozen to non-
1269 frozen periods, moving averages should only be calculated if at least 25-50% of the maximal
1270 data pair coverage is available within a particular time window.

1271 4.2.4 Rescaling

1272 When using fiducial reference data, units (e.g., m^3m^{-3} and degree of saturation) should be
1273 unified for the purpose of bias estimation using soil texture information, keeping in mind that
1274 inaccuracy in soil information directly propagates into the bias estimates. To account for (hor-
1275 izontal and vertical) systematic representativeness errors and different soil moisture units, the
1276 data set under validation should be rescaled (before decomposition for evaluating raw time
1277 series and after decomposition for evaluating anomalies) towards the reference data when esti-
1278 mating absolute uncertainties (i.e. ubRMSDs or ubRMSEs). When calculating relative metrics,
1279 data sets should be rescaled by matching their temporal mean and standard deviation. When
1280 calculating TCA-based metrics, data sets should be rescaled using also TCA-based rescaling
1281 coefficients. Note that no rescaling or unit conversion is necessary for Pearson correlation co-
1282 efficients or TCA-based correlation and SNR estimates, since these metrics are not affected by
1283 linear data transformation.

1284 4.3 Metric calculation

1285 Remember that all covariance-based metrics require zero error correlation. Any combination
1286 of in situ measurements, land surface model estimates, active-microwave-based retrievals, or
1287 passive-microwave-based retrievals is expected to mostly fulfil this requirement (see Sec. 3.4.2;
1288 *Gruber et al.*, 2016a). Different products from within any of these categories (except for in
1289 situ data), on the other hand, are expected to have correlated errors (*Gruber et al.*, 2016b).
1290 Therefore, the metrics described below should not be applied to such product combinations.
1291 Moreover, since non-zero error correlations may exist even when using products from different
1292 categories (see Sec. 3.4.2; *Yilmaz and Crow*, 2014; *Pan et al.*, 2015), it is strongly recommended
1293 to verify if assumptions are met (see Sec. 4.3.2).

1294 4.3.1 Relative metrics

1295 Temporal mean biases (Eq. (4)) should be calculated between all data sets that are expected
1296 to be properly collocated and have comparable spatial resolution, and are hence not dominated

1297 by spatial representativeness errors. These data sets may include dense networks, land surface
1298 models, and other satellite data sets. It should be kept in mind, however, that the underly-
1299 ing measurement resolution often considerably differs from the sampling grid resolution, which
1300 potentially causes representativeness errors that are not directly apparent as such. Correlation
1301 coefficients and unbiased Root-Mean-Square-Differences (Eqs. (9) and (7), respectively) should
1302 be calculated between all data sets whose errors are not expected to be correlated (see above).

1303 **4.3.2 TCA-based metrics**

1304 Second-order biases (Eq. (5)) of the validation data set should be calculated using fiducial
1305 reference data (i.e. at the core validation sites). Unbiased Root-Mean-Square-Errors and SNRs
1306 (Eqs. (8) and (11), respectively) should be calculated for all data sets. If more than one triplet
1307 with independent errors is available to estimate the bias or uncertainty of a particular product,
1308 TCA should be applied to all possible triplets and redundant estimates should be averaged
1309 (*Gruber et al.*, 2016b). The spread between redundant estimates should be used as a diagnostic
1310 to verify if orthogonality and zero error correlation assumptions are met (*Dorigo et al.*, 2010;
1311 *Draper et al.*, 2013; *Chen et al.*, 2017).

1312 **4.3.3 Confidence intervals**

1313 For each metric, 80-95% confidence intervals should be calculated using their analytical esti-
1314 mators (Eqs. (14)-(17)) or, if not available, block-bootstrapping. The latter should be based
1315 on at least 1000 bootstrap samples (*Efron and Tibshirani*, 1986) or possibly less if tested for
1316 convergence, and all confidence intervals should be corrected for sample auto-correlation.

1317 **4.4 Presentation**

1318 Validation metrics together with sample sizes and confidence intervals (and/or their upper and
1319 lower confidence limits) should be presented for each location where they are calculated, either
1320 by means of spatial maps or, if not meaningful (for example for core validation sites), in tabular
1321 form. Additionally, summary statistics (representing average conditions and spatial variability)
1322 of both validation metrics and their confidence intervals (and/or limits) should be provided,
1323 e.g., in the form of boxplots (i.e. median, inter-quartile-range and 5th/95th percentiles). The
1324 presentation can be further customized, for example by stratifying the summary statistics for
1325 climatological or land surface conditions.

1326 Ratio-based metrics (i.e. Pearson and TCA-based correlation coefficients as well as SNRs)
1327 must not be averaged. Differences between these metrics must always be related to their absolute
1328 values and be interpreted with care (see Sec. 3.7). SNR-related properties of different products
1329 may be compared in terms of SNR ratios or SNR differences in decibel space (Eq. (11)).

1330 Examples of how validation metrics and associated confidence intervals can be presented are
1331 provided in Appendix A.

1332 5 Final remarks: towards best practices

1333 In this paper we have reviewed state-of-the-art validation methods, including reference data
1334 sources and data pre-processing procedures, and provided good practice guidelines for the vali-
1335 dation of satellite soil moisture products. Moreover, we have identified several weak links that
1336 require careful attention to increase the reliability of soil moisture data quality assessments.
1337 Specifically, the following research gaps should be addressed in the near future:

- 1338 • On assumptions: the majority of studies assume that estimated biases and uncertainties
1339 are stationary (i.e. constant over time) or at least that they represent the average data
1340 quality of a product. However, given the strong link between soil moisture data quality
1341 and vegetation (*van der Schalie et al., 2018; Zwieback et al., 2018; Gruber et al., 2019a*),
1342 retrieval accuracy can be expected to vary strongly between seasons and many applications
1343 could greatly benefit from temporally varying quality information. Given the rapidly
1344 growing temporal coverage of soil moisture products, efforts should be made to provide
1345 bias and uncertainty estimates at different time scales, which also requires the use of
1346 seasonally varying bias correction (i.e. rescaling) parameters.
- 1347 • On pre-processing: very little is known about how spatial and temporal collocation mis-
1348 matches contribute to bias and uncertainty estimates. Using simple NN or IDW approaches
1349 to find match-ups between measurements and/or estimates that sample (represent) very
1350 different soil volumes or were taken at different times will give rise to representativeness
1351 errors that may considerably affect the overall picture of the quality of a product. More
1352 research is needed to quantify these representativeness errors and to develop resampling
1353 methods that more rigorously take actual measurement or model resolution into account.
- 1354 • On metric calculation: most current studies neglect the impact of second-order biases on
1355 various validation metrics such as the temporal mean difference or the ubRMSD. Several

1356 attempts are made to mitigate their impact using rescaling methods that match the sta-
1357 tistical moments of the data sets, yet most of these methods do not account for random
1358 errors and therefore match the moments in an insufficient manner. More research is needed
1359 to quantify the impact of suboptimal rescaling on second-order biases, on the impact of
1360 uncorrected second-order biases on validation metrics, and on how such uncorrected biases
1361 can be accounted for.

1362 • On reference data: validation targets are typically defined against an unknown truth.
1363 Comparing metrics against error-prone estimates of this truth (i.e. reference data) will
1364 be inflated by some unknown amount. Efforts should be made to obtain proper bias and
1365 uncertainty estimates for reference data sets, which should be further used to correct over-
1366 or underestimated validation metrics (*Miralles et al.*, 2010; *Chen et al.*, 2017).

1367 • On statistical uncertainty: most validation studies do not report confidence intervals,
1368 even though they are critical for a reliable interpretation of validation results. Although
1369 an accurate analytical calculation of confidence intervals for large-scale validation is not
1370 trivial for all metrics, bootstrapping provides an easy and robust alternative. However,
1371 care must be taken to properly account for spatial and temporal auto-correlation in the
1372 data.

1373 • On data merging: in recent years, several data merging algorithms have been developed
1374 that aim at providing consistent long-term soil moisture data records, whose temporal
1375 coverage extends beyond the lifetime of single satellite missions (*van der Schalie et al.*,
1376 2018; *Gruber et al.*, 2019a). Such merging procedures give rise to unique error charac-
1377 teristics such as highly non-stationary errors due to the intermittent and weighted use of
1378 retrievals from different sensors (*Gruber et al.*, 2017) or inhomogeneities between sensor
1379 transition periods (*Su et al.*, 2016). More research is needed to understand the impact of
1380 different transformation steps in data merging algorithms (e.g., data harmonization using
1381 cdf-matching) on final product quality, and good-practice validation guidelines need to be
1382 developed to comprehensively characterize such products.

1383 • On continuity: given the perpetual changes in the land surface character and climate as
1384 well as progressively increasing data record lengths, sensor drifts, changing reference data
1385 availability, and improving soil moisture retrieval algorithms, validation should be a con-
1386 tinuous process and validation reports frequently (at least annually) updated throughout

1387 and beyond the lifetime of the various satellite missions.

1388 • On accuracy requirements: the well-known soil moisture mission target accuracy require-
1389 ment of $0.04 \text{ m}^3\text{m}^{-3}$ (as specified by the Global Climate Observing System as well as for
1390 individual products and missions), against which soil moisture products are typically eval-
1391 uated, does not relate to the fitness-for-purpose for a specific application and no product
1392 can be declared “valid” if no meaningful validity requirements are available. We there-
1393 fore strongly encourage a closer collaboration between satellite data providers and the soil
1394 moisture user community to determine application specific accuracy requirements that
1395 provide deeper insight into what constitutes “good” or “bad” soil moisture data quality,
1396 thereby fostering the development of improved satellite products. To that end, we stress
1397 that only definitions of *relative* accuracy targets are meaningful as no reference for absolute
1398 soil moisture levels at a satellite scale is available (nor is it likely to be in the near future).

1399 Finally, many of the discussed principles and methods are not exclusively restricted to soil
1400 moisture. By setting this example, we hope to also nurture the development and evolution of
1401 validation good practice guidelines in other Earth observation communities.

1402 6 Acknowledgements

1403 We acknowledge the support from the International Space Science Institute (ISSI; <http://www.issibern.ch/>;
1404 last access: 1 July 2019). This publication is an outcome of the ISSI’s Team
1405 on “Adding value to soil moisture information for climate studies” and has received funding
1406 from the earthH2Observe project (European Union’s Seventh Framework Programme, Grant
1407 Agreement No. 603608), from the KU Leuven C1 internal fund C14/16/045, from the Research
1408 Foundation Flanders (FWO-1224320N and FWO-1530019N), and from the European Space
1409 Agency’s Climate Change Initiative for Soil Moisture.

1410 Appendix

1411 A Validation example

1412 Sec. 4 compiles the validation good practice guidelines provided in this paper into a recom-
1413 mended validation protocol. In this appendix, we provide an example that follows this protocol,
1414 not to actually assess the quality of certain products, but to provide an illustration that can be
1415 easily extrapolated to more specific validation tasks that readers may face. This includes a com-
1416 prehensive description of the validation setup, demonstrative examples of how validation results
1417 may be presented, and a discussion on where the currently available satellite soil moisture vali-
1418 dation literature often fails to comply with the good practice recommendations presented here.
1419 Results shown in this appendix have been generated using the python programming language.
1420 All source code is available at https://github.com/alexgruber/validation_good_practice/
1421 (last access: 1 July 2019). Metric calculation routines have been additionally translated into
1422 MATLAB.

1423 A.1 Data sets and study area

1424 Select validation examples are shown for soil moisture retrievals from the Advanced SCATterom-
1425 eter (ASCAT; *Naeimi et al.*, 2009), the Soil Moisture and Ocean Salinity (SMOS) mission (*Kerr*
1426 *et al.*, 2010), and the Soil Moisture Active Passive (SMAP) mission (*Entekhabi et al.*, 2010a).
1427 Reference data used are coarse-resolution model estimates from the Modern-Era Retrospective
1428 analysis for Research and Applications, Version 2 (MERRA-2; *Gelaro et al.*, 2017). This analy-
1429 sis is performed over the Contiguous United States (CONUS) using data from the beginning of
1430 2015 through the end of 2018.

1431 ASCAT data used are the EUMETSAT H SAF H113 data record and its extension H114,
1432 which are Level 2 (L2) soil moisture products that have been retrieved from inter-calibrated
1433 backscatter measurements from identical ASCAT instruments onboard the MetOp-A and MetOp-
1434 B satellites using the TU Wien Water Retrieval Package (WARP) algorithm (*Wagner et al.*,
1435 1999; *Naeimi et al.*, 2009). ASCAT is an active C-band radar with a spatial resolution of 25 km.
1436 Soil moisture is retrieved as the degree of saturation and sampled onto a 12.5 km discrete global
1437 grid. Data can be obtained upon registration from [http://hsaf.meteoam.it/soil-moisture.](http://hsaf.meteoam.it/soil-moisture.php)
1438 [php](http://hsaf.meteoam.it/soil-moisture.php) (last access: 1 July 2019).

1439 SMOS data are the reprocessed L2 soil moisture retrievals version V650, which can be ob-

1440 tained upon registration from <https://smos-diss.eo.esa.int/> (last access: 1 July 2019; *Kerr*
1441 *et al.*, 2012). SMOS is a passive L-band interferometric radiometer with an average spatial res-
1442 olution of 43 km. Soil moisture is retrieved in volumetric units and sampled on a 15 km discrete
1443 global grid.

1444 SMAP data used are the 36 km L2 radiometer-only soil moisture retrievals (SPL2SMP), al-
1445 gorithm version 5 (R16010) (*O’Neill et al.*, 2018, DOI: 10.5067/SODMLCE6LGLL). The passive
1446 SMAP radiometer operates at L-band at a spatial resolution of 40 km. Soil moisture is retrieved
1447 in volumetric units and sampled on the 36 km EASE grid version 2 (*Brodzik et al.*, 2012).

1448 MERRA-2 (*Gelaro et al.*, 2017) is the latest atmospheric reanalysis produced by NASA’s
1449 Global Modelling and Assimilation Office. Soil moisture is estimated on a $0.5^\circ \times 0.625^\circ$ grid in
1450 volumetric units as internal state variable of its land surface component, the Catchment Land
1451 Surface Model (*Koster et al.*, 2000). Here we use soil moisture estimates of the surface layer,
1452 which refers to the top 5 cm of the soil (*GMAO*, 2015). MERRA-2 data can be downloaded
1453 from https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/data_access/ (last access: 1 July
1454 2019).

1455 A.2 Pre-processing

1456 Unreliable soil moisture retrievals of the individual satellite products are masked out following
1457 the recommendations of the data providers. ASCAT soil moisture retrievals are masked out if
1458 the correction flag has a value other than 0 or 4, if the confidence flag and the processing flag
1459 have values other than 0, or if the surface state flag (*Naeimi et al.*, 2012) has a value other than
1460 1. SMOS retrievals are masked out if the RFI probability exceeds 0.1 or if the Chi-2 probability
1461 drops below 0.05. SMAP data are masked out if the retrieval quality flag has a value other than
1462 0 or 8. In addition, soil moisture retrievals of all satellite products are masked out at time steps
1463 where MERRA-2 estimates a soil temperature below 4°C or non-zero snow mass.

1464 ASCAT, SMOS and MERRA-2 are resampled to the 36 km EASE v2 grid that is used
1465 for SMAP retrievals using a nearest-neighbor approach. Note that ASCAT data is, although
1466 sampled on a 12.5 km grid, not aggregated as the actual measurement resolution (25 km) is
1467 already close to the EASE v2 grid resolution. Data sets are collocated in time by resampling them
1468 to fixed reference time steps with 24 hour intervals using a nearest-neighbor search. Reference
1469 time steps are selected for each grid cell separately such that they maximize the number of
1470 collocated time steps where all data sets provide valid soil moisture estimates. Note that the

1471 choice of this reference time step can increase or decrease the sample size - depending on the
1472 spatial location of the grid cell - by up to a factor of two.

1473 After spatial and temporal collocation, short-term anomalies are calculated for each data set
1474 using a 35-day moving average window. Long-term anomalies are not considered here because
1475 the study period of four years (2015-2018) is too short to calculate reliable long-term clima-
1476 tologies. The term “raw time series” is used to refer to the non-decomposed data, i.e. before
1477 anomalies have been calculated. For the estimation of unbiased RMSDs, data sets (both raw
1478 and anomaly time series) are rescaled by matching their temporal mean and standard deviation
1479 using MERRA-2 as scaling reference for comparability.

1480 **A.3 Skill metrics and presentation**

1481 **A.3.1 Sample size**

1482 All metrics are calculated from the same collocated data points, i.e. days where all four data
1483 sets provide valid soil moisture estimates. The number of temporal matches at each grid cell
1484 within our study domain is shown in Figure A.1. As discussed in Sec. 3, sample size directly
1485 translates into statistical power, i.e. reliability (in terms of confidence intervals) of the calculated
1486 skill metrics. Sample sizes obtained here, which range from 150 in the more mountainous areas
1487 to up to about 300-500 in the rest of the CONUS, are typically considered high and associated
1488 with reasonably low confidence intervals for validation purposes.

1489 However, as discussed in Sec. 3.6, confidence intervals are affected by temporal auto-
1490 correlation. “Effective” sample sizes, corrected for auto-correlation using Eq. (18), are ad-
1491 ditionally shown in Figure A.1 considering all data sets (for TCA metrics), and in Figure A.2
1492 for raw soil moisture time series and Figure A.3 for soil moisture anomalies considering different
1493 data set pairs. Effective sample sizes are considerably smaller than actual sample sizes, especially
1494 for raw time series due to the strong auto-correlation of the seasonal soil moisture cycle. Since
1495 auto-correlation levels vary between data sets, effective sample sizes vary when calculated for
1496 different data set pairs (albeit only slightly), which in turn leads to differences in the confidence
1497 intervals of relative skill metrics that are calculated between these data pairs.

1498 In the following, all analytical confidence intervals (Eqs. (14), (15), and (17)) are calculated
1499 using these auto-correlation corrected effective sample sizes. For bootstrapped confidence in-
1500 tervals, temporal auto-correlation is accounted for using block-bootstrapping (see Sec. 3.6.2)
1501 where block-lengths are estimated from the same auto-correlation levels that are underlying the

1502 calculation of effective sample sizes (see Eq. (21)).

1503 **A.3.2 Relative metrics**

1504 Figures A.4, A.5 and A.6 show spatial plots of relative (mean) bias, ubRMSD and R^2 (coefficient
1505 of determination or squared Pearson correlation) estimates for raw soil moisture values, respec-
1506 tively, and Figures A.7 and A.8 show ubRMSD and R^2 estimates for soil moisture anomalies,
1507 respectively.

1508 Biases are only calculated for raw soil moisture time series and between soil moisture esti-
1509 mates that are expressed in the same unit, i.e. for SMOS, SMAP, and MERRA-2 which provide
1510 estimates of volumetric soil moisture. ASCAT estimates of the degree of saturation could be
1511 converted into volumetric units using porosity information, but since the quality of soil texture
1512 maps on these scales is questionable, this is not recommended for bias estimation purposes. Note
1513 also, that the biases between the remaining three data sets also include collocation and (vertical
1514 and horizontal) scale mismatches and should therefore be interpreted with care.

1515 Along with the skill estimates, maps of confidence intervals are shown as the difference
1516 between the upper and lower confidence limits, chosen to be the 90th and the 10th percentile of
1517 the sampling distribution, respectively. Important to note is that confidence intervals for R^2 and
1518 ubRMSD estimates depend on the magnitude of the respective skill estimate, and are for R^2 not
1519 centered around the skill estimate. Misinterpretations may be avoided by directly presenting
1520 the actual confidence limits (see Sec. 3.7).

1521 We choose a confidence level of 80% because confidence intervals at the more common (yet
1522 completely arbitrary) 95% confidence level typically become excessively large for the sample
1523 sizes available from collocated satellite products (*Gruber et al.*, 2019a), especially when taking
1524 temporal auto-correlation into account.

1525 Figure A.9 shows spatial summary statistics of the relative skill metrics as well as of their
1526 upper and lower confidence limits. Hardly any skill differences would be considered significant
1527 when tested in the common way of checking for overlap between upper and lower confidence
1528 limits, even though Figures A.4 - A.8 show clear differences in spatial patterns.

1529 **A.3.3 Triple collocation metrics**

1530 As discussed in Sec. 3, TCA requires three data sets with independent random errors. Since
1531 errors of SMAP and SMOS are expected to be correlated (see Sec. 4.3), two independent data

1532 set triplets can be formed, i.e. ASCAT - SMOS - MERRA-2 and ASCAT - SMAP - MERRA-2.
1533 This results in unambiguous skill estimates for SMAP and SMOS, and in two skill estimates for
1534 ASCAT, which are averaged for increased precision.

1535 Figures A.10 and A.11 show spatial plots of TCA-based ubRMSE and R^2 (coefficient of
1536 determination w.r.t. the unknown truth) estimates, respectively, and Figures A.12 and A.13
1537 show ubRMSE and R^2 estimates for short-term soil moisture anomalies, respectively. The skill
1538 estimates represent the median of the bootstrapped sampling distribution, which are more robust
1539 than the direct estimates, and 80 % confidence intervals (i.e. the range between the 90th and
1540 the 10th percentile of the bootstrapped sampling distribution) are provided. Spatial summary
1541 statistics of the TCA estimates (sampling distribution median) as well as of the upper and lower
1542 confidence limits are shown in Figure A.14.

1543 The two degrees of freedom in TCA-based ASCAT skill estimates can not only be used for
1544 increasing the precision of the estimates by averaging them, but also to verify if TCA assumptions
1545 (i.e. zero error cross-correlation and error orthogonality) are met because if so, skill estimates
1546 should be identical. To this end, Figure A.15 shows the differences between R^2 and ubRMSE
1547 estimates for ASCAT when calculated once using SMOS as third data set and once using SMAP
1548 as third data set.

1549 On average, differences are close to zero and especially R^2 estimates do not exhibit spatial
1550 patterns of notable magnitude, which suggests that differences are mainly caused by sampling
1551 errors and hence that the TCA assumptions are generally respected. Some positive skill biases
1552 for raw soil moisture estimation for ASCAT are apparent in some northern and western parts
1553 of the CONUS, with skill estimates being slightly higher when using SMOS rather than SMAP
1554 in the triplet. These areas strongly coincide with regions of generally poor ASCAT performance
1555 (see Figure A.11), which is more pronounced in the ubRMSD because SNR biases of a given
1556 magnitude are associated with larger biases in error variance at low SNR levels than at high
1557 SNR levels. (see Sec. 3.7). Poor ASCAT performance in the northern CONUS is associated
1558 with issues in the vegetation correction of the WARP retrieval algorithm (see Sec. A.1). These
1559 uncorrected vegetation signals are removed when using soil moisture anomalies, which results in
1560 a considerable increase in skill metrics (see Figure A.13) and also removes the non-zero difference
1561 in ASCAT skill estimates when using SMOS versus SMAP for TCA, i.e. spurious error cross-
1562 correlations (see Figure A.15).

1563 A.4 Final remarks

1564 In this appendix, we provide an illustrative validation example that follows the good practice
1565 guidelines presented in this paper. For brevity, we omit the presentation of ground data compar-
1566 isons, which can be calculated and presented in the exact same way as the area-wide coarse-scale
1567 comparisons shown above. For simplicity, results are presented in spatial maps and boxplots
1568 that cover all of CONUS without further stratification. For summary information or if metrics
1569 are only computed at a few locations using ground reference data, results could be further pre-
1570 sented in tabular format. Some examples of comprehensive ground reference data comparison
1571 including both sparse networks and core validation sites can be found in *Dorigo et al. (2015)*;
1572 *Chen et al. (2017)*; *Colliander et al. (2017a)*.

1573 References

- 1574 Aitkin, A. (1936), On least squares and linear combination of observations, *Proceedings of the*
1575 *Royal Society of Edinburgh*, **55**, p. 42–48, doi:10.1017/S0370164600014346.
- 1576 Al-Yaari, A., J.-P. Wigneron, A. Ducharne, Y. Kerr, W. Wagner, G. De Lannoy, R. Reichle,
1577 A. Al Bitar, W. Dorigo, P. Richaume, et al. (2014), Global-scale comparison of passive (SMOS)
1578 and active (ASCAT) satellite based microwave soil moisture retrievals with soil moisture
1579 simulations (MERRA-Land), *Remote Sensing of Environment*, **152**, p. 614–626, doi:10.1016/
1580 j.rse.2014.07.013.
- 1581 Albergel, C., C. Ruediger, T. Pellarin, J. Calvet, N. Fritz, F. Froissard, D. Suquia, A. Pe-
1582 titpa, B. Piguet, and E. Martin (2008), From near-surface to root-zone soil moisture us-
1583 ing an exponential filter: an assessment of the method based on in-situ observations
1584 and model simulations., *Hydrology and earth system sciences.*, **12**(6), p. 1323–1337, doi:
1585 10.5194/hess-12-1323-2008.
- 1586 Albergel, C., E. Zakharova, J.-C. Calvet, M. Zribi, M. Pardé, J.-P. Wigneron, N. Novello,
1587 Y. Kerr, A. Mialon, and N. ed Dine Fritz (2011), A first assessment of the smos data in
1588 southwestern france using in situ and airborne soil moisture estimates: The carols airborne
1589 campaign, *Remote Sensing of Environment*, **115**(10), p. 2718 – 2728, doi:10.1016/j.rse.2011.
1590 06.012.
- 1591 Albergel, C., P. de Rosnay, C. Gruhier, J. Munoz-Sabater, S. Hasenauer, L. Isaksen, Y. Kerr, and

1592 W. Wagner (2012), Evaluation of remotely sensed and modelled soil moisture products using
1593 global ground-based in situ observations, *Remote Sensing of Environment*, **118**, p. 215–226,
1594 doi:10.1016/j.rse.2011.11.017.

1595 Albergel, C., W. Dorigo, R. Reichle, G. Balsamo, P. De Rosnay, J. Muñoz-Sabater, L. Isaksen,
1596 R. De Jeu, and W. Wagner (2013), Skill and global trend analysis of soil moisture from
1597 reanalyses and microwave remote sensing, *Journal of Hydrometeorology*, **14**(4), p. 1259–1277,
1598 doi:10.1175/JHM-D-12-0161.1.

1599 Babaeian, E., M. Sadeghi, S. B. Jones, C. Montzka, H. Vereecken, and M. Tuller (2019), Ground,
1600 proximal, and satellite remote sensing of soil moisture, *Reviews of Geophysics*, **57**, doi:10.1029/
1601 2018RG000618.

1602 Baker, M. (2016), 1,500 scientists lift the lid on reproducibility, *Nature News*, **533**(7604), p. 452,
1603 doi:10.1038/533452a.

1604 Balsamo, G., C. Albergel, A. Beljaars, S. Boussetta, E. Brun, H. Cloke, D. Dee, E. Dutra,
1605 J. Muñoz-Sabater, F. Pappenberger, et al. (2015), ERA-Interim/Land: a global land surface
1606 reanalysis data set, *Hydrology and Earth System Sciences*, **19**(1), p. 389–407, doi:10.5194/
1607 hess-19-389-2015.

1608 Bartalis, Z., R. Kidd, and K. Scipal (2006), Development and implementation of a discrete
1609 global grid system for soil moisture retrieval using the MetOp ASCAT scatterometer, in *1st*
1610 *EPS/MetOp RAO Workshop*, vol. ESA SP-618, ESRIN, Frascati, Italy.

1611 Bauer-Marschallinger, B., D. Sabel, and W. Wagner (2014), Optimisation of global grids for
1612 high-resolution remote sensing data, *Computers & Geosciences*, **72**, p. 84–93, doi:10.1016/j.
1613 cageo.2014.07.005.

1614 Bauer-Marschallinger, B., C. Paulik, S. Hochstöger, T. Mistelbauer, S. Modanesi, L. Ciabatta,
1615 C. Massari, L. Brocca, and W. Wagner (2018), Soil moisture from fusion of scatterometer
1616 and sar: Closing the scale gap with temporal filtering, *Remote Sensing*, **10**(7), p. 1030, doi:
1617 10.3390/rs10071030.

1618 Bindlish, R., T. J. Jackson, A. J. Gasiewski, M. Klein, and E. G. Njoku (2006), Soil moisture
1619 mapping and AMSR-E validation using the PSR in SMEX02, *Remote Sensing of Environment*,
1620 **103**(2), p. 127–139, doi:10.1016/j.rse.2005.02.003.

- 1621 Bindlish, R., T. Jackson, A. Gasiewski, B. Stankov, M. Klein, M. Cosh, I. Mladenova, C. Watts,
1622 E. Vivoni, V. Lakshmi, et al. (2008), Aircraft based soil moisture retrievals under mixed
1623 vegetation and topographic conditions, *Remote Sensing of Environment*, **112**(2), p. 375–390,
1624 doi:10.1016/j.rse.2007.01.024.
- 1625 Bircher, S., N. Skou, K. H. Jensen, J. Walker, and L. Rasmussen (2012), A soil moisture and
1626 temperature network for SMOS validation in western denmark, *Hydrology and Earth System
1627 Sciences*, **16**(5), p. 1445–1463.
- 1628 Blyth, C. R. (1972), On Simpson’s paradox and the sure-thing principle, *Journal of the American
1629 Statistical Association*, **67**(338), p. 364–366, doi:10.1080/01621459.1972.10482387.
- 1630 Bogena, H., C. Montzka, J. Huisman, A. Graf, M. Schmidt, M. Stockinger, C. von Hebel,
1631 H. Hendricks-Franssen, J. van der Kruk, W. Tappe, et al. (2018), The TERENO-Rur hydro-
1632 logical observatory: A multiscale multi-compartment research platform for the advancement
1633 of hydrological science, *Vadose Zone Journal*, **17**(1), doi:10.2136/vzj2018.03.0055.
- 1634 Bogena, H. R., J. A. Huisman, B. Schilling, A. Weuthen, and H. Vereecken (2017), Effective cal-
1635 ibration of low-cost soil water content sensors, *Sensors*, **17**(1), p. 208, doi:10.3390/s17010208.
- 1636 Bolten, J. D., W. T. Crow, X. Zhan, T. J. Jackson, and C. A. Reynolds (2010), Evaluating the
1637 utility of remotely sensed soil moisture retrievals for operational agricultural drought moni-
1638 toring, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*,
1639 **3**(1), p. 57–66, doi:10.1109/JSTARS.2009.2037163.
- 1640 Bonett, D. G., and T. A. Wright (2000), Sample size requirements for estimating pearson, kendall
1641 and spearman correlations, *Psychometrika*, **65**(1), p. 23–28, doi:10.1007/BF02294183.
- 1642 Brocca, L., F. Melone, T. Moramarco, and R. Morbidelli (2010a), Spatial-temporal variability of
1643 soil moisture and its estimation across scales, *Water Resources Research*, **46**(2), doi:10.1029/
1644 2009WR008016.
- 1645 Brocca, L., F. Melone, T. Moramarco, W. Wagner, and S. Hasenauer (2010b), ASCAT soil
1646 wetness index validation through in situ and modeled soil moisture data in central italy,
1647 *Remote Sensing of Environment*, **114**(11), p. 2745–2755, doi:10.1016/j.rse.2010.06.009.
- 1648 Brocca, L., S. Hasenauer, T. Lacava, F. Melone, T. Moramarco, W. Wagner, W. Dorigo, P. Mat-
1649 gen, J. Martinez-Fernandez, P. Llorens, J. Latron, C. Martin, and M. Bittelli (2011), Soil

1650 moisture estimation through ASCAT and AMSR-E sensors: An intercomparison and vali-
1651 dation study across europe, *Remote Sensing of Environment*, **115**(12), p. 3390–3408, doi:
1652 10.1016/j.rse.2011.08.003.

1653 Brocca, L., T. Tullo, F. Melone, T. Moramarco, and R. Morbidelli (2012), Catchment scale soil
1654 moisture spatial-temporal variability, *Journal of Hydrology*, **422-423**, p. 63–75, doi:10.1016/
1655 j.jhydrol.2011.12.039.

1656 Brodzik, M. J., B. Billingsley, T. Haran, B. Raup, and M. H. Savoie (2012), EASE-Grid 2.0:
1657 Incremental but significant improvements for earth-gridded data sets, *ISPRS International*
1658 *Journal of Geo-Information*, **1**(1), p. 32–45, doi:10.3390/ijgi1010032.

1659 Burgin, M. S., A. Colliander, E. G. Njoku, S. K. Chan, F. Cabot, Y. H. Kerr, R. Bindlish, T. J.
1660 Jackson, D. Entekhabi, and S. H. Yueh (2017), A comparative study of the SMAP passive
1661 soil moisture product with existing satellite-based soil moisture products, *IEEE Transactions*
1662 *on Geoscience and Remote Sensing*, **55**(5), p. 2959–2971, doi:10.1109/TGRS.2017.2656859.

1663 Caires, S., and A. Sterl (2003), Validation of ocean wind and wave data using triple collocation,
1664 *Journal of Geophysical Research: Oceans*, **108**(C3), doi:10.1029/2002JC001491.

1665 Caldwell, T. G., T. Bongiovanni, M. H. Cosh, C. Halley, and M. H. Young (2018), Field and
1666 laboratory evaluation of the cs655 soil water content sensor, *Vadose Zone Journal*, **17**(1),
1667 doi:10.2136/vzj2017.12.0214.

1668 Caldwell, T. G., T. Bongiovanni, M. H. Cosh, T. J. Jackson, A. Colliander, C. J. Abolt, R. Cas-
1669 teel, B. R. Scanlon, and M. H. Young (2019), The texas soil observation network: A compre-
1670 hensive soil moisture dataset for remote sensing and land surface model validation, *Vadose*
1671 *Zone Journal*, doi:10.2136/vzj2019.04.0034.

1672 Chauhan, N. S., S. Miller, and P. Ardanuy (2003), Spaceborne soil moisture estimation at
1673 high resolution: a microwave-optical/ir synergistic approach, *International Journal of Remote*
1674 *Sensing*, **24**(22), p. 4599–4622, doi:10.1080/0143116031000156837.

1675 Chen, F., W. T. Crow, A. Colliander, M. H. Cosh, T. J. Jackson, R. Bindlish, R. H. Reichle,
1676 S. K. Chan, D. D. Bosch, P. J. Starks, et al. (2017), Application of triple collocation in ground-
1677 based validation of soil moisture active/passive (SMAP) level 2 data products, *IEEE Journal*

1678 *of Selected Topics in Applied Earth Observations and Remote Sensing*, **10**(2), p. 489–502,
1679 doi:10.1109/JSTARS.2016.2569998.

1680 Chen, F., W. T. Crow, R. Bindlish, A. Colliander, M. S. Burgin, J. Asanuma, and K. Aida (2018),
1681 Global-scale evaluation of SMAP, SMOS and ASCAT soil moisture products using triple
1682 collocation, *Remote Sensing of Environment*, **214**, p. 1–13, doi:10.1016/j.rse.2018.05.008.

1683 Chen, F., W. T. Crow, M. H. Cosh, A. Colliander, J. Asanuma, A. Berg, D. D. Bosch, T. G.
1684 Caldwell, C. H. Collins, K. H. Jensen, J. Martínez-Fernández, H. McNairn, P. J. Starks,
1685 Z. Su, and J. P. Walker (2019), Uncertainty of reference pixel soil moisture averages sampled
1686 at smap core validation sites, *Journal of Hydrometeorology*, **20**(8), p. 1553–1569, doi:10.1175/
1687 JHM-D-19-0049.1.

1688 Colliander, A., T. J. Jackson, R. Bindlish, S. Chan, N. Das, S. Kim, M. Cosh, R. Dunbar,
1689 L. Dang, L. Pashaian, et al. (2017a), Validation of SMAP surface soil moisture products with
1690 core validation sites, *Remote sensing of environment*, **191**, p. 215–231, doi:10.1016/j.rse.2017.
1691 01.021.

1692 Colliander, A., M. H. Cosh, S. Misra, T. J. Jackson, W. T. Crow, S. Chan, R. Bindlish, C. Chae,
1693 C. H. Collins, and S. H. Yueh (2017b), Validation and scaling of soil moisture in a semi-arid en-
1694 vironment: Smap validation experiment 2015 (smavex15), *Remote Sensing of Environment*,
1695 **196**, p. 101 – 112, doi:10.1016/j.rse.2017.04.022.

1696 Colliander, A., M. H. Cosh, S. Misra, T. J. Jackson, W. T. Crow, J. Powers, H. McNairn,
1697 P. Bullock, A. Berg, R. Magagi, Y. Gao, R. Bindlish, R. Williamson, I. Ramos, B. Latham,
1698 P. O’Neill, and S. Yueh (2019), Comparison of high-resolution airborne soil moisture retrievals
1699 to smap soil moisture during the smap validation experiment 2016 (smavex16), *Remote*
1700 *Sensing of Environment*, **227**, p. 137 – 150, doi:10.1016/j.rse.2019.04.004.

1701 Corey, D. M., W. P. Dunlap, and M. J. Burke (1998), Averaging correlations: Expected val-
1702 ues and bias in combined pearson rs and fisher’s z transformations, *The Journal of general*
1703 *psychology*, **125**(3), p. 245–261, doi:10.1080/00221309809595548.

1704 Cosh, M., T. J. Jackson, R. Bindlish, J. S. Famiglietti, and D. Ryu (2005), Calibration of an
1705 impedance probe for estimation of surface soil water content over large areas, *Journal of*
1706 *Hydrology*, **311**, p. 49–58, doi:10.1016/j.jhydrol.2005.01.003.

- 1707 Cosh, M. H., T. J. Jackson, R. Bindlish, and J. H. Prueger (2004), Watershed scale temporal and
1708 spatial stability of soil moisture and its role in validating satellite estimates, *Remote sensing*
1709 *of Environment*, **92**(4), p. 427–435, doi:10.1016/j.rse.2004.02.016.
- 1710 Cosh, M. H., T. J. Jackson, P. Starks, and G. Heathman (2006), Temporal stability of surface
1711 soil moisture in the little washita river watershed and its applications in satellite soil moisture
1712 product validation, *Journal of Hydrology*, **323**(1–4), p. 168–177, doi:10.1016/j.jhydrol.2005.
1713 08.020.
- 1714 Cosh, M. H., T. J. Jackson, S. Moran, and R. Bindlish (2008), Temporal persistence and stability
1715 of surface soil moisture in a semi-arid watershed, *Remote Sensing of Environment*, **112**(2),
1716 p. 304 – 313, doi:10.1016/j.rse.2007.07.001, soil Moisture Experiments 2004 (SMEX04) Special
1717 Issue.
- 1718 Crow, W. T., A. A. Berg, M. H. Cosh, A. Loew, B. P. Mohanty, R. Panciera, P. de Rosnay,
1719 D. Ryu, and J. P. Walker (2012), Upscaling sparse ground-based soil moisture observations
1720 for the validation of coarse-resolution satellite soil moisture products, *Rev. Geophys.*, **50**(2),
1721 p. RG2002, doi:10.1029/2011RG000372.
- 1722 Cuenca, R. H., D. E. Stangel, and S. F. Kelly (1997), Soil water balance in a boreal forest, *Journal*
1723 *of Geophysical Research-Atmospheres*, **102**(D 24), p. 29,355–29,365, doi:10.1029/97JD02312.
- 1724 Das, N. N., D. Entekhabi, S. Kim, T. Jagdhuber, S. Dunbar, S. Yueh, and A. Colliander (2017),
1725 High-resolution enhanced product based on smap active-passive approach using sentinel 1a
1726 and 1b sar data, in *2017 IEEE International Geoscience and Remote Sensing Symposium*
1727 *(IGARSS)*, p. 2543–2545, IEEE, doi:10.1109/IGARSS.2017.8127513.
- 1728 Dawdy, D., and N. Matalas (1964), *Statistical and probability analysis of hydrologic data, part*
1729 *III: Analysis of variance, covariance and time series*, McGraw-Hill.
- 1730 De Lannoy, G. J., and R. H. Reichle (2016), Assimilation of SMOS brightness temperatures
1731 or soil moisture retrievals into a land surface model, *Hydrology and Earth System Sciences*,
1732 **20**(12), p. 4895–4911, doi:10.5194/hess-20-4895-2016.
- 1733 de Nijs, A. H., R. M. Parinussa, R. A. de Jeu, J. Schellekens, and T. R. Holmes (2015), A
1734 methodology to determine radio-frequency interference in AMSR2 observations, *Geoscience*

- 1735 *and Remote Sensing, IEEE Transactions on*, **53**(9), p. 5148–5159, doi:10.1109/TGRS.2015.
1736 2417653.
- 1737 Dee, D. P. (2005), Bias and data assimilation, *Quarterly Journal of the Royal Meteorological*
1738 *Society*, **131**(613), p. 3323–3343, doi:10.1256/qj.05.137.
- 1739 Djamai, N., R. Magagi, K. Goïta, M. Hosseini, M. H. Cosh, A. Berg, and B. Toth (2015),
1740 Evaluation of SMOS soil moisture products over the CanEx-SM10 area, *Journal of hydrology*,
1741 **520**, p. 254–267, doi:10.1016/j.jhydrol.2014.11.026.
- 1742 Dorigo, W., P. van Oevelen, W. Wagner, M. Drusch, S. Mecklenburg, A. Robock, and T. Jackson
1743 (2011a), A new international network for in situ soil moisture data, *Eos Transactions AGU*,
1744 **92**(17), p. 141–142, doi:10.1029/2011EO170001.
- 1745 Dorigo, W., R. de Jeu, D. Chung, R. Parinussa, Y. Liu, W. Wagner, and D. Fernández-Prieto
1746 (2012), Evaluating global trends (1988–2010) in harmonized multi-satellite surface soil mois-
1747 ture, *Geophysical Research Letters*, **39**(18), doi:10.1029/2012GL052988.
- 1748 Dorigo, W., A. Xaver, M. Vreugdenhil, A. Gruber, H. A. A. Sanchis-Dufau, D. Zamojski,
1749 C. Cordes, W. Wagner, and M. Drusch (2013), Global automated quality control of in situ
1750 soil moisture data from the international soil moisture network, *Vadose Zone Journal*, **12**(3),
1751 doi:10.2136/vzj2012.0097.
- 1752 Dorigo, W., A. Gruber, R. De Jeu, W. Wagner, T. Stacke, A. Loew, C. Albergel, L. Brocca,
1753 D. Chung, R. Parinussa, et al. (2015), Evaluation of the ESA CCI soil moisture product using
1754 ground-based observations, *Remote Sensing of Environment*, **162**, p. 380–395, doi:10.1016/j.
1755 rse.2014.07.023.
- 1756 Dorigo, W., W. Wagner, C. Albergel, F. Albrecht, G. Balsamo, L. Brocca, D. Chung, M. Ertl,
1757 M. Forkel, A. Gruber, et al. (2017), ESA CCI soil moisture for improved earth system un-
1758 derstanding: state-of-the art and future directions, *Remote Sensing of Environment*, **203**,
1759 p. 185–215, doi:10.1016/j.rse.2017.07.001.
- 1760 Dorigo, W. A., K. Scipal, R. M. Parinussa, Y. Y. Liu, W. Wagner, R. A. M. de Jeu, and
1761 V. Naeimi (2010), Error characterisation of global active and passive microwave soil moisture
1762 datasets, *Hydrol. Earth Syst. Sci.*, **14**(12), p. 2605–2616, doi:10.5194/hessd-7-5621-2010.

- 1763 Dorigo, W. A., W. Wagner, R. Hohensinn, S. Hahn, C. Paulik, A. Xaver, A. Gruber, M. Drusch,
1764 S. Mecklenburg, P. van Oevelen, A. Robock, and T. Jackson (2011b), The international soil
1765 moisture network: a data hosting facility for global in situ soil moisture measurements, *Hydrol.*
1766 *Earth Syst. Sci.*, **15**(5), p. 1675–1698, doi:10.5194/hess-15-1675-2011.
- 1767 Draper, C., and R. Reichle (2015), The impact of near-surface soil moisture assimilation at
1768 subseasonal, seasonal, and inter-annual timescales, *Hydrology and Earth System Sciences*,
1769 **19**(12), p. 4831, doi:10.5194/hess-19-4831-2015.
- 1770 Draper, C., R. Reichle, G. De Lannoy, and Q. Liu (2012), Assimilation of passive and ac-
1771 tive microwave soil moisture retrievals, *Geophysical Research Letters*, **39**(4), doi:10.1029/
1772 2011GL050655.
- 1773 Draper, C., R. Reichle, R. de Jeu, V. Naeimi, R. Parinussa, and W. Wagner (2013), Estimating
1774 root mean square errors in remotely sensed soil moisture over continental scale domains,
1775 *Remote Sensing of Environment*, **137**, p. 288–298, doi:10.1016/j.rse.2013.06.013.
- 1776 Efron, B., and R. Tibshirani (1986), Bootstrap methods for standard errors, confidence intervals,
1777 and other measures of statistical accuracy, *Statistical science*, **1**(1), p. 54–75, doi:10.1214/ss/
1778 1177013815.
- 1779 Entekhabi, D., E. Njoku, P. O’Neill, K. Kellogg, W. Crow, W. Edelstein, J. Entin, S. Good-
1780 man, T. Jackson, J. Johnson, J. Kimball, J. Piepmeier, R. Koster, N. Martin, K. McDonald,
1781 M. Moghaddam, S. Moran, R. Reichle, J. Shi, M. Spencer, S. Thurman, L. Tsang, and
1782 J. Van Zyl (2010a), The soil moisture active passive (SMAP) mission, *Proceedings of the*
1783 *IEEE*, **98**(5), p. 704–716, doi:10.1109/JPROC.2010.2043918.
- 1784 Entekhabi, D., R. H. Reichle, R. D. Koster, and W. T. Crow (2010b), Performance metrics
1785 for soil moisture retrievals and application requirements, *J. Hydrometeor*, **11**(3), p. 832–840,
1786 doi:10.1175/2010JHM1223.1.
- 1787 Famiglietti, J., J. Devereaux, C. Laymon, T. Tsegaye, P. Houser, T. Jackson, S. Graham,
1788 M. Rodell, and P. v. Oevelen (1999), Ground-based investigation of soil moisture variability
1789 within remote sensing footprints during the southern great plains 1997 (SGP97) hydrology
1790 experiment, *Water Resources Management (1999)*, **35**(6), p. 1839–1851.

1791 Famiglietti, J. S., D. Ryu, A. A. Berg, M. Rodell, and T. J. Jackson (2008), Field observations
1792 of soil moisture variability across scales, *Water Resour. Res.*, **44**(1), p. W01,423, doi:10.1029/
1793 2006WR005804.

1794 Figa-Saldaña, J., J. J. Wilson, E. Attema, R. Gelsthorpe, M. Drinkwater, and A. Stoffelen
1795 (2002), The advanced scatterometer (ASCAT) on the meteorological operational (MetOp)
1796 platform: A follow on for european wind scatterometers, *Canadian Journal of Remote Sensing*,
1797 **28**(3), p. 404–412, doi:10.5589/m02-035.

1798 Fox, N. (2010), A guide to “reference standards” in support of quality assur-
1799 ance requirements of GEO, *Tech. Rep. QA4EO-QAEO-GEN-DQK-003, v4.0*, QA4EO,
1800 http://qa4eo.org/docs/QA4EO-QAEO-GEN-DQK-003_v4.0.pdf, last access: 1 July 2019.

1801 Gelaro, R., W. McCarty, M. J. Suárez, R. Todling, A. Molod, L. Takacs, C. A. Randles, A. Dar-
1802 menov, M. G. Bosilovich, R. Reichle, et al. (2017), The modern-era retrospective analysis for
1803 research and applications, version 2 (MERRA-2), *Journal of Climate*, **30**(14), p. 5419–5454,
1804 doi:10.1175/JCLI-D-16-0758.1.

1805 Gelman, A., and H. Stern (2006), The difference between “significant” and “not significant” is
1806 not itself statistically significant, *The American Statistician*, **60**(4), p. 328–331, doi:10.1198/
1807 000313006X152649.

1808 Gilleland, E. (2010), Confidence intervals for forecast verification, *NCAR Technical Note*, **TN-**
1809 **479**, doi:10.5065/D6WD3XJM.

1810 GMAO (2015), Global Modeling and Assimilation Office (GMAO), MERRA-2 tavg1_2d_lnd_Nx:
1811 2d,1-Hourly,Time-Averaged,Single-Level,Assimilation,Land Surface Diagnostics V5.12.4,
1812 Greenbelt, MD, USA, Goddard Earth Sciences Data and Information Services Center (GES
1813 DISC), Accessed: 1 Nov 2018, doi:10.5067/RKPHT8KC1Y1T.

1814 Greenland, S., S. J. Senn, K. J. Rothman, J. B. Carlin, C. Poole, S. N. Goodman, and D. G. Alt-
1815 man (2016), Statistical tests, p values, confidence intervals, and power: a guide to misinterpre-
1816 tations, *European journal of epidemiology*, **31**(4), p. 337–350, doi:10.1007/s10654-016-0149-3.

1817 Gruber, A., W. Dorigo, S. Zwieback, A. Xaver, and W. Wagner (2013a), Characterizing coarse-
1818 scale representativeness of in situ soil moisture measurements from the international soil mois-
1819 ture network, *Vadose Zone Journal*, **12**(2), doi:10.2136/vzj2012.0170.

1820 Gruber, A., W. Wagner, A. Hegyiova, F. Greifeneder, and S. Schlaffer (2013b), Potential of
1821 sentinel-1 for high-resolution soil moisture monitoring, in *Geoscience and Remote Sensing*
1822 *Symposium (IGARSS), 2013 IEEE International*, p. 4030–4033, IEEE, doi:10.1109/IGARSS.
1823 2017.8127513.

1824 Gruber, A., W. Crow, W. Dorigo, and W. Wagner (2015), The potential of 2D Kalman filtering
1825 for soil moisture data assimilation, *Remote Sensing of Environment*, **171**, p. 137–148, doi:
1826 10.1016/j.rse.2015.10.019.

1827 Gruber, A., C.-H. Su, S. Zwieback, W. Crow, W. Dorigo, and W. Wagner (2016a), Recent
1828 advances in (soil moisture) triple collocation analysis, *International Journal of Applied Earth*
1829 *Observation and Geoinformation*, **45**, p. 200–211, doi:10.1016/j.jag.2015.09.002.

1830 Gruber, A., C.-H. Su, W. Crow, S. Zwieback, W. Dorigo, and W. Wagner (2016b), Estimating
1831 error cross-correlations in soil moisture data sets using extended collocation analysis, *Journal*
1832 *of Geophysical Research: Atmospheres*, **121(3)**, p. 1208–1219, doi:10.1002/2015JD024027.

1833 Gruber, A., W. A. Dorigo, W. Crow, and W. Wagner (2017), Triple collocation-based merging
1834 of satellite soil moisture retrievals, *IEEE Transactions on Geoscience and Remote Sensing*,
1835 **55(12)**, p. 6780–6792, doi:10.1109/TGRS.2017.2734070.

1836 Gruber, A., W. Crow, and W. Dorigo (2018), Assimilation of spatially sparse in situ soil moisture
1837 networks into a continuous model domain, *Water Resources Research*, **54(2)**, p. 1353–1367,
1838 doi:10.1002/2017WR021277.

1839 Gruber, A., T. Scanlon, R. van der Schalie, W. Wagner, and W. Dorigo (2019a), Evolution
1840 of the esa cci soil moisture climate data records and their underlying merging methodology,
1841 *Earth System Science Data*, **11(2)**, p. 717–739, doi:10.5194/essd-11-717-2019.

1842 Gruber, A., G. D. Lannoy, and W. Crow (2019b), A monte carlo based adaptive kalman filtering
1843 framework for soil moisture data assimilation, *Remote Sensing of Environment*, **228**, p. 105
1844 – 114, doi:10.1016/j.rse.2019.04.003.

1845 Gupta, H. V., H. Kling, K. K. Yilmaz, and G. F. Martinez (2009), Decomposition of the mean
1846 squared error and NSE performance criteria: Implications for improving hydrological mod-
1847 elling, *Journal of Hydrology*, **377(1)**, p. 80–91, doi:10.1016/j.jhydrol.2009.08.003.

1848 H-SAF (2017), Product validation report (PVR) h111 metop ASCAT soil mois-
1849 ture, *Tech. Rep. SAF/HSAF/CDOP3/PVR/H111, v0.3*, EUMETSAT H SAF reports,
1850 http://hsaf.meteoam.it/documents/PVR/H111_ASCAT_SSM_CDR_PVR_v0.3.pdf (last ac-
1851 cess: 1 July 2019).

1852 H-SAF (2018), Algorithm theoretical baseline document (ATBD) soil mois-
1853 ture data records, metop ASCAT soil moisture time series, *Tech.*
1854 *Rep. SAF/HSAF/CDOP3/ATBD, v0.7*, EUMETSAT H SAF reports,
1855 http://hsaf.meteoam.it/documents/ATDD/ASCAT_SSM_CDR_ATBD_v0.7.pdf (last ac-
1856 cess: 1 July 2019).

1857 Jackson, T., M. Cosh, R. Bindlish, P. Starks, D. Bosch, M. Seyfried, D. Goodrich, M. Moran, and
1858 J. Du (2010), Validation of advanced microwave scanning radiometer soil moisture products,
1859 *Geoscience and Remote Sensing, IEEE Transactions on*, **48**(12), p. 4256–4272, doi:10.1109/
1860 TGRS.2010.2051035.

1861 Jackson, T., A. Colliander, J. Kimball, R. Reichle, W. Crow, D. Entekhabi, and P. Neill (2012),
1862 Science data calibration and validation plan, *SMAP Mission, NASA Jet Propuls. Lab.*

1863 Jackson, T. J., D. M. Le Vine, C. T. Swift, T. J. Schmugge, and F. R. Schiebe (1995), Large
1864 area mapping of soil moisture using the ESTAR passive microwave radiometer in washita’92,
1865 *Remote sensing of Environment*, **54**(1), p. 27–37, doi:10.1016/0034-4257(95)00084-E.

1866 Jackson, T. J., D. M. Le Vine, A. Y. Hsu, A. Oldak, P. J. Starks, C. T. Swift, J. D. Isham,
1867 and M. Haken (1999), Soil moisture mapping at regional scales using microwave radiometry:
1868 The southern great plains hydrology experiment, *IEEE transactions on geoscience and remote*
1869 *sensing*, **37**(5), p. 2136–2151, doi:10.1109/36.789610.

1870 Jackson, T. J., R. Bindlish, A. J. Gasiewski, B. Stankov, M. Klein, E. G. Njoku, D. Bosch,
1871 T. L. Coleman, C. A. Laymon, and P. Starks (2005), Polarimetric scanning radiometer C-
1872 and X-band microwave observations during SMEX03, *IEEE Transactions on Geoscience and*
1873 *Remote Sensing*, **43**(11), p. 2418–2430, doi:10.1109/TGRS.2005.857625.

1874 JCGM (2008), Evaluation of measurement data—guide to the expression of uncer-
1875 tainty in measurement (GUM), *Tech. Rep. JCGM 100:2008*, Bureau International des
1876 Poids et Mesures (BIPM), Joint Committee for Guides in Metrology (JCGM), URL:
1877 <https://www.bipm.org/en/publications/guides/gum.html>, last access: 1 July 2019.

1878 JCGM (2012), International vocabulary of metrology—basic and general concepts and asso-
1879 ciated terms (VIM 3rd edition), *Tech. Rep. JCGM 200:2012*, Bureau International des
1880 Poids et Mesures (BIPM), Joint Committee for Guides in Metrology (JCGM), URL:
1881 <https://www.bipm.org/en/publications/guides/vim.html>, last access: 1 July 2019.

1882 Justice, C., A. Belward, J. Morissette, P. Lewis, J. Privette, and F. Baret (2000), Developments
1883 in the ‘validation’ of satellite sensor products for the study of the land surface, *International*
1884 *Journal of Remote Sensing*, **21**(17), p. 3383–3390, doi:10.1080/014311600750020000.

1885 Kerr, Y., P. Waldteufel, J.-P. Wigneron, S. Delwart, F. Cabot, J. Boutin, M. Escorihuela,
1886 J. Font, N. Reul, C. Gruhier, S. Juglea, M. Drinkwater, A. Hahne, M. Martin-Neira, and
1887 S. Mecklenburg (2010), The SMOS mission: New tool for monitoring key elements of the global
1888 water cycle, *Proceedings of the IEEE*, **98**(5), p. 666–687, doi:10.1109/JPROC.2010.2043032.

1889 Kerr, Y. H., P. Waldteufel, J.-P. Wigneron, J. Martinuzzi, J. Font, and M. Berger (2001), Soil
1890 moisture retrieval from space: The soil moisture and ocean salinity (SMOS) mission, *IEEE*
1891 *transactions on Geoscience and remote sensing*, **39**(8), p. 1729–1735, doi:10.1109/36.942551.

1892 Kerr, Y. H., P. Waldteufel, P. Richaume, J. P. Wigneron, P. Ferrazzoli, A. Mahmoodi,
1893 A. Al Bitar, F. Cabot, C. Gruhier, S. E. Juglea, et al. (2012), The SMOS soil moisture
1894 retrieval algorithm, *IEEE Transactions on Geoscience and Remote Sensing*, **50**(5), p. 1384–
1895 1403, doi:10.1109/TGRS.2012.2184548.

1896 Kerr, Y. H., A. Al-Yaari, N. Rodriguez-Fernandez, M. Parrens, B. Molero, D. Leroux, S. Bircher,
1897 A. Mahmoodi, A. Mialon, P. Richaume, et al. (2016), Overview of SMOS performance in terms
1898 of global soil moisture monitoring after six years in operation, *Remote Sensing of Environment*,
1899 **180**, p. 40–63, doi:10.1016/j.rse.2016.02.042.

1900 Kolassa, J., P. Gentine, C. Prigent, F. Aires, and S. Alemohammad (2017), Soil moisture retrieval
1901 from amsr-e and ascats microwave observation synergy. part 2: Product evaluation, *Remote*
1902 *Sensing of Environment*, **195**, p. 202 – 217, doi:<https://doi.org/10.1016/j.rse.2017.04.020>.

1903 Koster, R. D., M. J. Suarez, A. Ducharne, M. Stieglitz, and P. Kumar (2000), A catchment-
1904 based approach to modeling land surface processes in a general circulation model: 1. model
1905 structure, *Journal of Geophysical Research: Atmospheres*, **105**(D20), p. 24,809–24,822, doi:
1906 10.1029/2000JD900327.

- 1907 Koster, R. D., Z. Guo, R. Yang, P. A. Dirmeyer, K. Mitchell, and M. J. Puma (2009), On
1908 the nature of soil moisture in land surface models, *Journal of Climate*, **22**(16), p. 4322–4335,
1909 doi:10.1175/2009JCLI2832.1.
- 1910 Kumar, S. V., R. H. Reichle, K. W. Harrison, C. D. Peters-Lidard, S. Yatheendradas, and J. A.
1911 Santanello (2012), A comparison of methods for a priori bias correction in soil moisture data
1912 assimilation, *Water Resour. Res.*, **48**(3), p. W03,515, doi:10.1029/2010WR010261.
- 1913 Lahoz, W., and G. De Lannoy (2014), Closing the gaps in our knowledge of the hydrological
1914 cycle over land: Conceptual problems, *Surveys in Geophysics*, **35**(3), p. 623–660, doi:10.1007/
1915 s10712-013-9221-7.
- 1916 Loew, A., W. Bell, L. Brocca, C. E. Bulgin, J. Burdanowitz, X. Calbet, R. V. Donner,
1917 D. Ghent, A. Gruber, T. Kaminski, et al. (2017), Validation practices for satellite-based
1918 earth observation data across communities, *Reviews of Geophysics*, **55**(3), p. 779–817, doi:
1919 10.1002/2017RG000562.
- 1920 Macelloni, G., M. Brogioni, P. Pampaloni, A. Cagnati, and M. R. Drinkwater (2006), DOMEX
1921 2004: An experimental campaign at Dome-C antarctica for the calibration of spaceborne
1922 low-frequency microwave radiometers, *IEEE transactions on geoscience and remote sensing*,
1923 **44**(10), p. 2642–2653, doi:10.1109/TGRS.2006.882801.
- 1924 Magagi, R., A. A. Berg, K. Goïta, S. Bélair, T. J. Jackson, B. Toth, A. Walker, H. McNairn,
1925 P. E. O’Neill, M. Moghaddam, et al. (2013), Canadian experiment for soil moisture in 2010
1926 (CanEx-SM10): Overview and preliminary results, *IEEE Transactions on Geoscience and*
1927 *Remote Sensing*, **51**(1), p. 347–363, doi:10.1109/TGRS.2012.2198920.
- 1928 Martínez-Fernández, J., and A. Ceballos (2005), Mean soil moisture estimation using temporal
1929 stability analysis, *Journal of Hydrology*, **312**(1), p. 28 – 38, doi:10.1016/j.jhydrol.2005.02.007.
- 1930 McColl, K. A., J. Vogelzang, A. G. Konings, D. Entekhabi, M. Piles, and A. Stoffelen
1931 (2014), Extended triple collocation: Estimating errors and correlation coefficients with
1932 respect to an unknown target, *Geophysical Research Letters*, **41**(17), p. 6229–6236, doi:
1933 10.1002/2014GL061322.
- 1934 McColl, K. A., A. Roy, C. Derksen, A. G. Konings, S. H. Alemohammed, and D. En-
1935 tekhabi (2016), Triple collocation for binary and categorical variables: Application to val-

- 1936 idating landscape freeze/thaw retrievals, *Remote Sensing of Environment*, **176**, p. 31–42,
1937 doi:10.1016/j.rse.2016.01.010.
- 1938 McNairn, H., T. J. Jackson, G. Wiseman, S. Belair, A. Berg, P. Bullock, A. Colliander, M. H.
1939 Cosh, S.-B. Kim, R. Magagi, et al. (2015), The soil moisture active passive validation experi-
1940 ment 2012 (SMAPVEX12): Prelaunch calibration and validation of the SMAP soil moisture
1941 algorithms, *IEEE Transactions on Geoscience and Remote Sensing*, **53**(5), p. 2784–2801, doi:
1942 10.1109/TGRS.2014.2364913.
- 1943 Merchant, C. J., F. Paul, T. Popp, M. Ablain, S. Bontemps, P. Defourny, R. Hollmann,
1944 T. Lavergne, A. Laeng, G. d. Leeuw, et al. (2017), Uncertainty information in climate
1945 data records from earth observation, *Earth System Science Data*, **9**(2), p. 511–527, doi:
1946 10.5194/essd-9-511-2017.
- 1947 Miralles, D. G., W. T. Crow, and M. H. Cosh (2010), Estimating spatial sampling errors in
1948 coarse-scale soil moisture estimates derived from point-scale observations, *J. Hydrometeorol.*,
1949 **11**(6), p. 1423–1429, doi:10.1175/2010JHM1285.1.
- 1950 Miyaoka, K., A. Gruber, F. Ticconi, S. Hahn, W. Wagner, J. Figa-Saldana, and C. Anderson
1951 (2017), Triple collocation analysis of soil moisture from Metop-A ASCAT and SMOS against
1952 JRA-55 and ERA-Interim, *IEEE Journal of Selected Topics in Applied Earth Observations*
1953 *and Remote Sensing*, **10**(5), p. 2274–2284, doi:10.1109/JSTARS.2016.2632306.
- 1954 Moghaddam, M., D. Entekhabi, Y. Goykhman, K. Li, M. Liu, A. Mahajan, A. Nayyar,
1955 D. Shuman, and D. Teneketzis (2010), A wireless soil moisture smart sensor web using
1956 physics-based optimal control: Concept and initial demonstrations, *IEEE Journal of Se-*
1957 *lected Topics in Applied Earth Observations and Remote Sensing*, **3**(4), p. 522–535, doi:
1958 10.1109/JSTARS.2010.2052918.
- 1959 Molero, B., D. Leroux, P. Richaume, Y. Kerr, O. Merlin, M. Cosh, and R. Bindlish (2018),
1960 Multi-timescale analysis of the spatial representativeness of in situ soil moisture data within
1961 satellite footprints, *Journal of Geophysical Research: Atmospheres*, **123**(1), p. 3–21, doi:10.
1962 1002/2017JD027478.
- 1963 Naeimi, V., K. Scipal, Z. Bartalis, S. Hasenauer, and W. Wagner (2009), An improved soil
1964 moisture retrieval algorithm for ERS and METOP scatterometer observations, *Geoscience*

- 1965 *and Remote Sensing, IEEE Transactions on*, **47**(7), p. 1999–2013, doi:10.1109/TGRS.2008.
1966 2011617.
- 1967 Naeimi, V., C. Paulik, A. Bartsch, W. Wagner, R. Kidd, S.-E. Park, K. Elger, and J. Boike
1968 (2012), ASCAT surface state flag (SSF): Extracting information on surface freeze/thaw con-
1969 ditions from backscatter data using an empirical threshold-analysis algorithm, *Geoscience*
1970 *and Remote Sensing, IEEE Transactions on*, **50**(7), p. 2566–2582, doi:10.1109/TGRS.2011.
1971 2177667.
- 1972 Narapusetty, B., T. DelSole, and M. K. Tippett (2009), Optimal estimation of the climatological
1973 mean, *Journal of Climate*, **22**(18), p. 4845–4859, doi:10.1175/2009JCLI2944.1.
- 1974 Neyman, J. (1937), X—outline of a theory of statistical estimation based on the classical theory
1975 of probability, *Philosophical Transactions of the Royal Society of London. Series A, Mathe-*
1976 *matical and Physical Sciences*, **236**(767), p. 333–380, doi:10.1098/rsta.1937.0005.
- 1977 Nicolai-Shaw, N., M. Hirschi, H. Mittelbach, and S. I. Seneviratne (2015), Spatial representa-
1978 tiveness of soil moisture using in situ, remote sensing, and land reanalysis data, *Journal of*
1979 *Geophysical Research: Atmospheres*, **120**(19), p. 9955–9964, doi:10.1002/2015JD023305.
- 1980 Noilhan, J., P. Lacarrère, and P. Bougeault (1991), An experiment with an advanced surface pa-
1981 rameterization in a mesobeta-scale model. part III: Comparison with the HAPEX-MOBILHY
1982 dataset, *Monthly weather review*, **119**(10), p. 2393–2413, doi:10.1175/1520-0493(1991)
1983 119(2393:AEWAAS)2.0.CO;2.
- 1984 Ochsner, T. E., M. H. Cosh, R. H. Cuenca, W. A. Dorigo, C. S. Draper, Y. Hagimoto, Y. H.
1985 Kerr, E. G. Njoku, E. E. Small, M. Zreda, et al. (2013), State of the art in large-scale
1986 soil moisture monitoring, *Soil Science Society of America Journal*, **77**(6), p. 1888–1919, doi:
1987 10.2136/sssaj2013.03.0093.
- 1988 Ólafsdóttir, K., and M. Mudelsee (2014), More accurate, calibrated bootstrap confidence inter-
1989 vals for estimating the correlation between two time series, *Mathematical Geosciences*, **46**(4),
1990 p. 411–427, doi:10.1007/s11004-014-9523-4.
- 1991 O’Neill, P., S. Chan, E. Njoku, T. Jackson, and R. Bindlish (2012), SMAP level 2 & 3 soil
1992 moisture (passive) algorithm theoretical basis document (ATBD), *Initial Release, version*, **1**.

- 1993 O'Neill, P., S. Chan, E. Njoku, T. Jackson, and R. Blindish (2018), SMAP L2 radiometer half-
1994 orbit 36 km EASE-grid soil moisture, version 5, *Boulder, Colorado USA. NASA National*
1995 *Snow and ice Data Center Distributed Active Archive Center*, doi:[https://doi.org/10.5067/](https://doi.org/10.5067/SODMLCE6LGLL)
1996 SODMLCE6LGLL.
- 1997 Pan, M., C. K. Fisher, N. W. Chaney, W. Zhan, W. T. Crow, F. Aires, D. Entekhabi, and E. F.
1998 Wood (2015), Triple collocation: Beyond three estimates and separation of structural/non-
1999 structural errors, *Remote Sensing of Environment*, **171**, p. 299–310, doi:[doi.org/10.1016/j.rse.](https://doi.org/10.1016/j.rse.2015.10.028)
2000 2015.10.028.
- 2001 Panciera, R., J. P. Walker, J. D. Kalma, E. J. Kim, J. M. Hacker, O. Merlin, M. Berger, and
2002 N. Skou (2008), The NAFE'05/CoSMOS data set: Toward SMOS soil moisture retrieval,
2003 downscaling, and assimilation, *IEEE Transactions on Geoscience and Remote Sensing*, **46**(3),
2004 p. 736–745, doi:[10.1109/TGRS.2007.915403](https://doi.org/10.1109/TGRS.2007.915403).
- 2005 Parinussa, R. M., A. G. Meesters, Y. Y. Liu, W. Dorigo, W. Wagner, and R. A. De Jeu (2011),
2006 Error estimates for near-real-time satellite soil moisture as derived from the land parameter
2007 retrieval model, *Geoscience and Remote Sensing Letters, IEEE*, **8**(4), p. 779–783, doi:[10.1109/](https://doi.org/10.1109/LGRS.2011.2114872)
2008 LGRS.2011.2114872.
- 2009 Parinussa, R. M., T. R. Holmes, N. Wanders, W. A. Dorigo, and R. A. de Jeu (2015), A
2010 preliminary study toward consistent soil moisture from AMSR2, *Journal of Hydrometeorology*,
2011 **16**(2), p. 932–947, doi:[10.1175/JHM-D-13-0200.1](https://doi.org/10.1175/JHM-D-13-0200.1).
- 2012 Pathe, C., W. Wagner, D. Sabel, M. Doubkova, and J. B. Basara (2009), Using envisat asar
2013 global mode data for surface soil moisture retrieval over oklahoma, usa, *IEEE Transactions*
2014 *on Geoscience and Remote Sensing*, **47**(2), p. 468–480, doi:[10.1109/TGRS.2008.2004711](https://doi.org/10.1109/TGRS.2008.2004711).
- 2015 Peischl, S., J. P. Walker, C. Rüdiger, N. Ye, Y. H. Kerr, E. Kim, R. Bandara, and M. Al-
2016 lahmoradi (2012), The AACES field experiments: SMOS calibration and validation across
2017 the murrumbidgee river catchment., *Hydrology & Earth System Sciences Discussions*, **9**(3),
2018 doi:[10.5194/hessd-9-2763-2012](https://doi.org/10.5194/hessd-9-2763-2012).
- 2019 Peng, J., A. Loew, S. Zhang, J. Wang, and J. Niesel (2015), Spatial downscaling of satellite
2020 soil moisture data using a vegetation temperature condition index, *IEEE Transactions on*
2021 *Geoscience and Remote Sensing*, **54**(1), p. 558–566, doi:[10.1109/TGRS.2015.2462074](https://doi.org/10.1109/TGRS.2015.2462074).

- 2022 Peng, J., A. Loew, O. Merlin, and N. E. Verhoest (2017), A review of spatial downscaling
2023 of satellite remotely sensed soil moisture, *Reviews of Geophysics*, **55**(2), p. 341–366, doi:
2024 10.1002/2016RG000543.
- 2025 Pierdicca, N., F. Fascetti, L. Pulvirenti, and R. Crapolicchio (2017), Error characterization of
2026 soil moisture satellite products: Retrieving error cross-correlation through extended quadru-
2027 ple collocation, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote*
2028 *Sensing*, **10**(10), p. 4522–4530, doi:10.1109/JSTARS.2017.2714025.
- 2029 QA4EO (2010), *A Quality Assurance Framework for Earth Observation: Principles*, version 4.0
2030 ed.
- 2031 Quast, R., and W. Wagner (2016), Analytical solution for first-order scattering in bistatic ra-
2032 diative transfer interaction problems of layered media, *Applied optics*, **55**(20), p. 5379–5386,
2033 doi:10.1364/AO.55.005379.
- 2034 Reichle, R., G. De Lannoy, Q. Liu, R. Koster, J. Kimball, W. Crow, J. Ardizzone,
2035 P. Chakraborty, D. Collins, L. Conaty, et al. (2017a), Global assessment of the SMAP level-4
2036 surface and root-zone soil moisture product using assimilation diagnostics, *Journal of Hy-*
2037 *drometeorology*, **18**(12), p. 3217–3237, doi:10.1175/JHM-D-17-0130.1.
- 2038 Reichle, R., G. De Lannoy, Q. Liu, J. Ardizzone, A. Colliander, A. Conaty, W. Crow, T. Jackson,
2039 L. Jones, J. Kimball, et al. (2017b), Assessment of the SMAP level-4 surface and root-zone soil
2040 moisture product using in situ measurements, *Journal of hydrometeorology*, **18**(10), p. 2621–
2041 2645, doi:10.1175/JHM-D-17-0063.1.
- 2042 Reichle, R. H., and R. D. Koster (2004), Bias reduction in short records of satellite soil moisture,
2043 *Geophys. Res. Lett.*, **31**(19), p. L19,501, doi:10.1029/2004GL020938.
- 2044 Reichle, R. H., R. D. Koster, G. J. De Lannoy, B. A. Forman, Q. Liu, S. P. Mahanama, and
2045 A. Touré (2011), Assessment and enhancement of MERRA land surface hydrology estimates,
2046 *Journal of climate*, **24**(24), p. 6322–6338, doi:10.1175/JCLI-D-10-05033.1.
- 2047 Reichle, R. H., C. S. Draper, Q. Liu, M. Girotto, S. P. Mahanama, R. D. Koster, and G. J.
2048 De Lannoy (2017c), Assessment of MERRA-2 land surface hydrology estimates, *Journal of*
2049 *Climate*, **30**(8), p. 2937–2960, doi:10.1175/JCLI-D-16-0720.1.

2050 Rodell, M., P. Houser, U. e. a. Jambor, J. Gottschalck, K. Mitchell, C. Meng, K. Arse-
2051 nault, B. Cosgrove, J. Radakovich, M. Bosilovich, et al. (2004), The global land data as-
2052 simulation system, *Bulletin of the American Meteorological Society*, **85**(3), p. 381–394, doi:
2053 10.1175/BAMS-85-3-381.

2054 Rüdiger, C., A. W. Western, J. P. Walker, A. B. Smith, J. D. Kalma, and G. R. Willgoose (2010),
2055 Towards a general equation for frequency domain reflectometers, *Journal of hydrology*, **383**(3-
2056 4), p. 319–329, doi:10.1016/j.jhydrol.2009.12.046.

2057 Rykiel Jr, E. J. (1996), Testing ecological models: the meaning of validation, *Ecological mod-
2058 elling*, **90**(3), p. 229–244, doi:10.1016/0304-3800(95)00152-2.

2059 Sabaghy, S., J. Walker, L. Renzullo, R. Akbar, S. Chan, J. Chaubell, N. Das, R. Dunbar,
2060 D. Entekhabi, A. Gevaert, T. Jackson, A. Loew, O. Merlin, M. Moghaddam, J. Peng, J. Peng,
2061 J. Piepmeier, C. Rüdiger, V. Stefan, X. Wu, N. Ye, and S. Yueh (in review), Comprehensive
2062 analysis of alternative downscaled soil moisture products, *Remote Sensing of Environment*.

2063 Sahoo, A. K., G. J. D. Lannoy, R. H. Reichle, and P. R. Houser (2013), Assimilation and
2064 downscaling of satellite observed soil moisture over the little river experimental watershed in
2065 georgia, usa, *Advances in Water Resources*, **52**, p. 19 – 33, doi:10.1016/j.advwatres.2012.08.
2066 007.

2067 Scanlon, T., J. Nightingale, F. Boersma, J.-P. Muller, C. Farquhar, S. Compernelle, and J.-
2068 C. Lambert (2017), Outline of QA4ECV quality assurance service (version 2.0), *Tech. rep.*,
2069 QA4ECV, <http://www.qa4ecv.eu/qa-system>, last access: 1 July 2019.

2070 Scipal, K., M. Drusch, and W. Wagner (2008a), Assimilation of a ERS scatterometer derived
2071 soil moisture index in the ECMWF numerical weather prediction system, *Advances in water
2072 resources*, **31**(8), p. 1101–1112, doi:10.1016/j.advwatres.2008.04.013.

2073 Scipal, K., T. Holmes, R. de Jeu, V. Naeimi, and W. Wagner (2008b), A possible solution for
2074 the problem of estimating the error structure of global soil moisture data sets, *Geophys. Res.
2075 Lett.*, **35**(24), p. L24,403, doi:10.1029/2008GL035599.

2076 Seyfried, M., L. Grant, E. Du, and K. Humes (2005), Dielectric loss and calibration of the hydra
2077 probe soil water sensor, *Vadose Zone Journal*, **4**(4), p. 1070–1079, doi:10.2136/vzj2004.0148.

2078 Smith, A., J. Walker, A. Western, R. Young, K. Ellett, R. Pipunic, R. Grayson, L. Siriwardena,
2079 F. Chiew, and H. Richter (2012), The murrumbidgee soil moisture monitoring network data
2080 set, *Water Resources Research*, **48**(7).

2081 Starks, P. J., G. C. Heathman, T. J. Jackson, and M. H. Cosh (2006), Temporal stability of soil
2082 moisture profile, *Journal of Hydrology*, **324**, p. 400–411, doi:10.1016/j.jhydrol.2005.09.024.

2083 Stoffelen, A. (1998), Toward the true near-surface wind speed: Error modeling and calibration
2084 using triple collocation, *J. Geophys. Res.*, **103**(C4), p. 7755–7766, doi:10.1029/97JC03180.

2085 Su, C.-H., and D. Ryu (2015), Multi-scale analysis of bias correction of soil moisture, *Hydrology
2086 and Earth System Sciences*, **19**(1), p. 17–31, doi:10.5194/hess-19-17-2015.

2087 Su, C.-H., D. Ryu, W. T. Crow, and A. W. Western (2014), Beyond triple collocation: Appli-
2088 cations to soil moisture monitoring, *Journal of Geophysical Research: Atmospheres*, **119**(11),
2089 p. 6419–6439, doi:10.1002/2013JD021043.

2090 Su, C.-H., D. Ryu, W. Dorigo, S. Zwieback, A. Gruber, C. Albergel, R. H. Reichle, and W. Wag-
2091 ner (2016), Homogeneity of a global multisatellite soil moisture climate data record, *Geophys-
2092 ical Research Letters*, **43**(21), p. 11–245, doi:10.1002/2016GL070458.

2093 Su, Z., W. Timmermans, Y. Zeng, J. Schulz, V. O. John, R. A. Roebeling, P. Poli, D. Tan,
2094 F. Kaspar, A. K. Kaiser-Weiss, E. Swinnen, C. Toté, H. Gregow, T. Manninen, A. Riihelä,
2095 J.-C. Calvet, Y. Ma, and J. Wen (2018), An overview of european efforts in generating climate
2096 data records, *Bulletin of the American Meteorological Society*, **99**(2), p. 349–359, doi:10.1175/
2097 BAMS-D-16-0074.1.

2098 Tong, C. (2019), Statistical inference enables bad science; statistical thinking enables good sci-
2099 ence, *The American Statistician*, **73**(sup1), p. 246–261, doi:10.1080/00031305.2018.1518264.

2100 Ulaby, F. T., D. G. Long, W. J. Blackwell, C. Elachi, A. K. Fung, C. Ruf, K. Sarabandi,
2101 H. A. Zebker, and J. Van Zyl (2014), *Microwave radar and radiometric remote sensing*, vol. 4,
2102 University of Michigan Press Ann Arbor.

2103 Vachaud, G., A. Passerat De Silans, P. Balabanis, and M. Vauclin (1985), Temporal stability
2104 of spatially measured soil water probability density function, *Soil Sci. Soc. Am. J.*, **49**(4),
2105 p. 822–828, doi:10.2136/sssaj1985.03615995004900040006x.

2106 van der Schalie, R., R. de Jeu, R. Parinussa, N. Rodríguez-Fernández, Y. Kerr, A. Al-Yaari,
2107 J.-P. Wigneron, and M. Drusch (2018), The effect of three different data fusion approaches
2108 on the quality of soil moisture retrievals from multiple passive microwave sensors, *Remote*
2109 *Sensing*, **10**(1), p. 107, doi:10.3390/rs10010107.

2110 Van Leeuwen, P. J. (2015), Representation errors and retrievals in linear and nonlinear data
2111 assimilation, *Quarterly Journal of the Royal Meteorological Society*, **141**(690), p. 1612–1623.

2112 Vogelzang, J., and A. Stoffelen (2012), Triple collocation, *EUMETSAT Report. Available at*
2113 *http://research.metoffice.gov.uk/research/interproj/nwpsaf/scatterometer/TripleCollocation_NWPSAF_TR.L*
2114 *last access: 1 July 2019.*

2115 Wagner, W., G. Lemoine, and H. Rott (1999), A method for estimating soil moisture from
2116 ERS scatterometer and soil data, *Remote Sensing of Environment*, **70**(2), p. 191–207, doi:
2117 10.1016/S0034-4257(99)00036-X.

2118 Wagner, W., L. Brocca, V. Naeimi, R. Reichle, C. Draper, R. de Jeu, D. Ryu, C.-H. Su, A. West-
2119 ern, J.-C. Calvet, et al. (2014), Clarifications on the “comparison between SMOS, VUA, AS-
2120 CAT, and ECMWF soil moisture products over four watersheds in US”, *IEEE Transactions*
2121 *on Geoscience and Remote Sensing*, **52**(3), p. 1901–1906, doi:10.1109/TGRS.2013.2282172.

2122 Walker, J. P., G. R. Willgoose, and J. D. Kalma (2004), In situ measurement of soil moisture:
2123 a comparison of techniques, *Journal of Hydrology*, **293**, p. 85–99, doi:10.1016/j.jhydrol.2004.
2124 01.008.

2125 Wang, G., D. Garcia, Y. Liu, R. De Jeu, and A. J. Dolman (2012), A three-dimensional gap
2126 filling method for large geophysical datasets: Application to global satellite soil moisture
2127 observations, *Environmental Modelling & Software*, **30**, p. 139–142, doi:10.1016/j.envsoft.
2128 2011.10.015.

2129 Wasserstein, R. L., and N. A. Lazar (2016), The ASA’s statement on p-values: context, pro-
2130 cess, and purpose, *The American Statistician*, **70**(2), p. 129–133, doi:10.1080/00031305.2016.
2131 1154108.

2132 Wasserstein, R. L., A. L. Schirm, and N. A. Lazar (2019), Moving to a world beyond “p<0.05”,
2133 *The American Statistician*, **73**(sup1), p. 1–19, doi:10.1080/00031305.2019.1583913.

2134 Wigneron, J.-P., T. Jackson, P. O’neill, G. De Lannoy, P. De Rosnay, J. Walker, P. Ferrazzoli,
2135 V. Mironov, S. Bircher, J. Grant, et al. (2017), Modelling the passive microwave signature
2136 from land surfaces: A review of recent results and application to the l-band smos & smap
2137 soil moisture retrieval algorithms, *Remote Sensing of Environment*, **192**, p. 238–262, doi:
2138 10.1016/j.rse.2017.01.024.

2139 Wilks, D. S. (2011), *Statistical Methods in the Atmospheric Sciences*, vol. 100, 3rd ed., Academic
2140 Press.

2141 WMO (2016), The global observing system for climate: Implementation needs, *Implementation*
2142 *Plan GCOS-200*, World Meteorological Organization.

2143 Yee, M. S., J. P. Walker, A. Monerris, C. Rüdiger, and T. J. Jackson (2016), On the identification
2144 of representative in situ soil moisture monitoring stations for the validation of smap soil
2145 moisture products in australia, *Journal of Hydrology*, **537**, p. 367 – 381, doi:10.1016/j.jhydrol.
2146 2016.03.060.

2147 Yilmaz, M. T., and W. T. Crow (2013), The optimality of potential rescaling approaches in land
2148 data assimilation., *Journal of Hydrometeorology*, **14**(2), doi:10.1175/JHM-D-12-052.1.

2149 Yilmaz, M. T., and W. T. Crow (2014), Evaluation of assumptions in soil moisture triple collocation
2150 analysis, *Journal of Hydrometeorology*, **15**(3), p. 1293–1302, doi:10.1175/JHM-D-13-0158.
2151 1.

2152 Zeng, Y., Z. Su, J.-C. Calvet, T. Manninen, E. Swinnen, J. Schulz, R. Roebeling, P. Poli, D. Tan,
2153 A. Riihelä, C.-M. Tanis, A.-N. Arslan, A. Obregon, A. Kaiser-Weiss, V. John, W. Timmer-
2154 mans, J. Timmermans, F. Kaspar, H. Gregow, A.-L. Barbu, D. Fairbairn, E. Gelati, and
2155 C. Meurey (2015), Analysis of current validation practices in europe for space-based climate
2156 data records of essential climate variables, *International Journal of Applied Earth Observation*
2157 *and Geoinformation*, **42**, p. 150 – 161, doi:https://doi.org/10.1016/j.jag.2015.06.006.

2158 Zribi, M., M. Pardé, J. Boutin, P. Fanise, D. Hauser, M. Dechambre, Y. Kerr, M. Leduc-
2159 Leballeur, G. Reverdin, N. Skou, S. Søbjaerg, C. Albergel, J. C. Calvet, J. P. Wigneron,
2160 E. Lopez-Baeza, A. Rius, and J. Tenerelli (2011), Carols: A new airborne l-band radiometer
2161 for ocean surface and land observations, *Sensors*, **11**(1), p. 719–742, doi:10.3390/s110100719.

2162 Zwieback, S., K. Scipal, W. Dorigo, and W. Wagner (2012), Structural and statistical properties
2163 of the collocation technique for error characterization, *Nonlin. Processes Geophys.*, **19**(1),
2164 p. 69–80, doi:10.5194/npg-19-69-2012.

2165 Zwieback, S., A. Colliander, M. H. Cosh, J. Martínez-Fernández, H. McNairn, P. J. Starks,
2166 M. Thibeault, and A. Berg (2018), Estimating time-dependent vegetation biases in the SMAP
2167 soil moisture product, *Hydrology and Earth System Sciences*, **22**(8), p. 4473–4489, doi:10.
2168 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	http://ismn.geo.tuwien.ac.at/ (<i>Dorigo et al., 2011a,b</i>)
CVS	Openly available Core Validation Site (CVS) data that have been specifically processed for SMAP validation.	https://nsidc.org/data/nsidc-0712
GLDAS	NASA’s global modelling and data assimilation system	https://ldas.gsfc.nasa.gov/gldas/
MERRA	NASA’s global reanalysis data sets	https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/
ERA	ECMWF’s global reanalysis data sets	https://www.ecmwf.int/en/forecasts/datasets/browse-reanalysis-datasets/

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	https://github.com/alexgruber/validation_good_practice/
pytesmo	Geospatial time series validation toolbox	python	https://doi.org/10.5281/zenodo.1215760/
poets	Geospatial image resampling toolbox	python	https://pypi.org/project/poets/

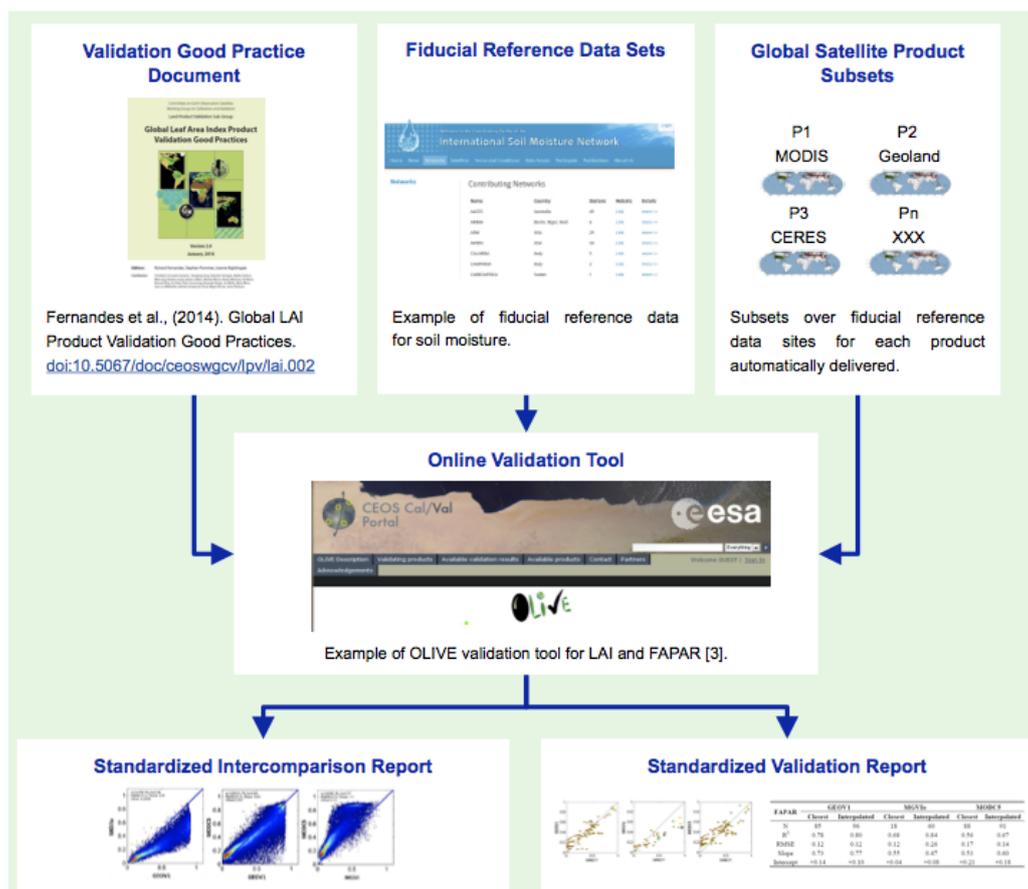


Figure 1: Validation framework as defined by CEOS (from <https://lpvs.gsfc.nasa.gov/>; last access: 1 July 2019).



Figure 2: Currently available stations from sparse networks hosted by the ISMN (from https://www.geo.tuwien.ac.at/insitu/data_viewer/, last access: 1 July 2019). Colors represent different station hosting networks.

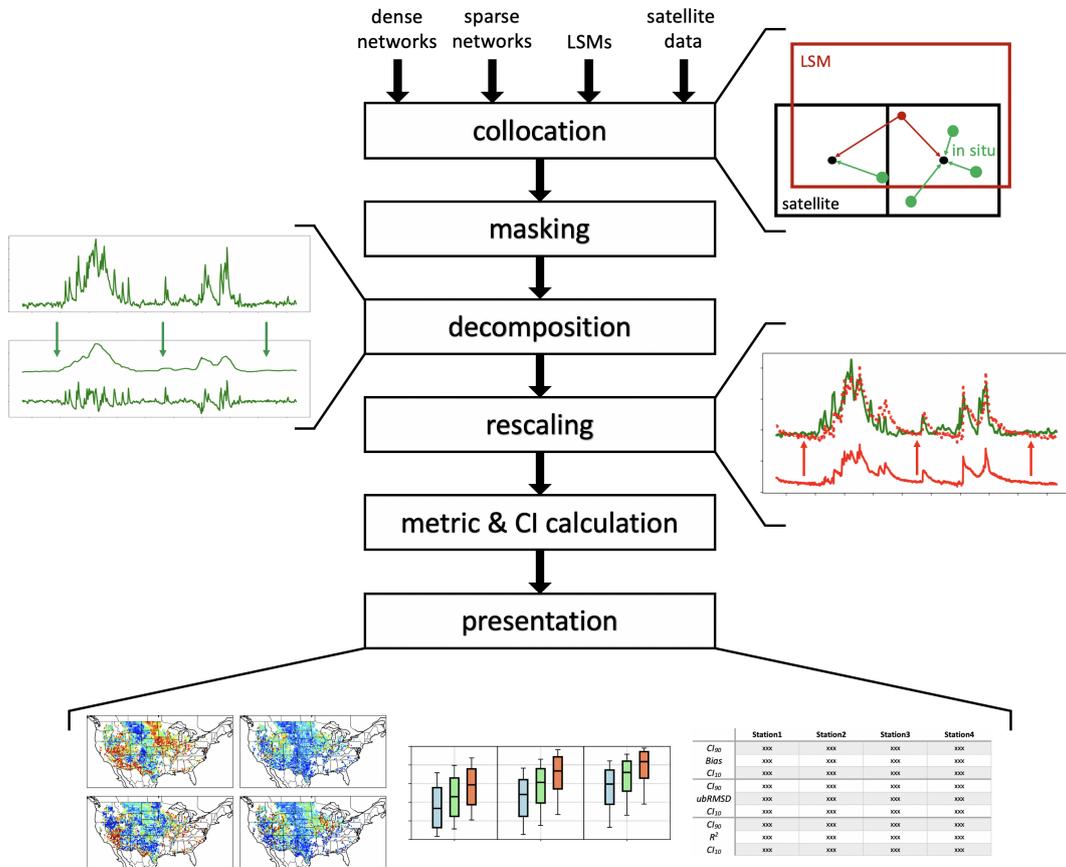


Figure 3: Validation good practice protocol illustration.

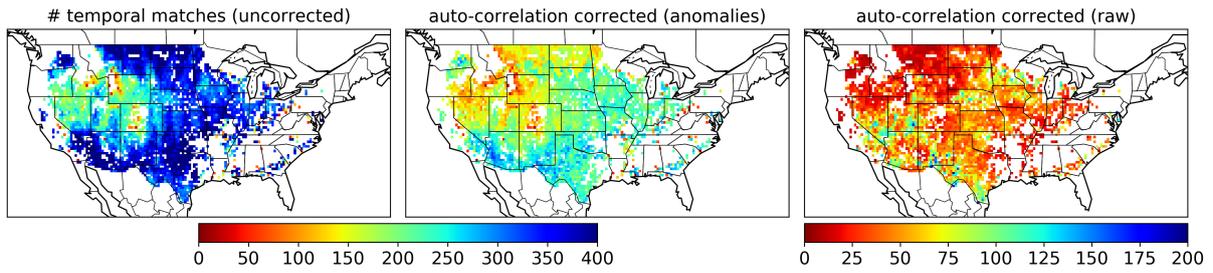


Figure A.1: Sample size for temporal matches between ASCAT, SMOS, SMAP and MERRA-2 between 2015 and 2018 (left), effective sample size when correcting for anomaly auto-correlation (middle), and effective sample size when correcting for auto-correlation in the raw time series (right).

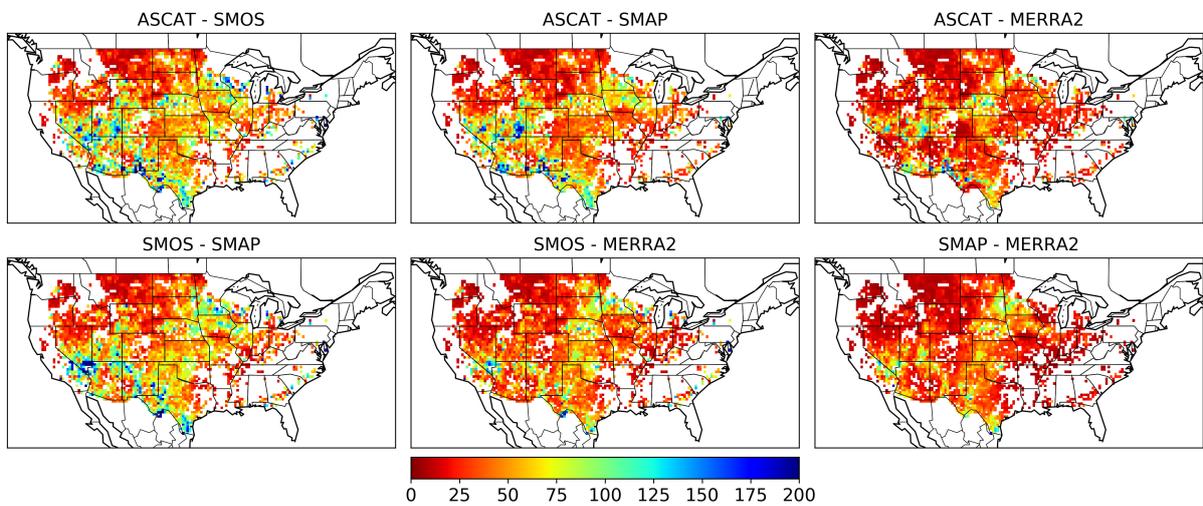


Figure A.2: Effective raw time series sample size, corrected for auto-correlation, for different data set combinations.

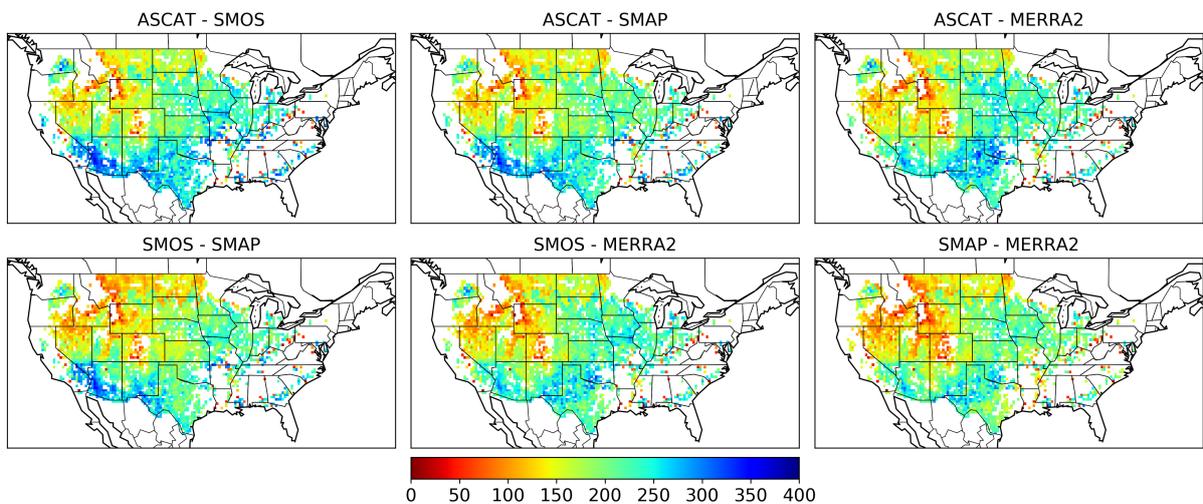


Figure A.3: Effective anomaly sample size, corrected for auto-correlation, for different data set combinations.

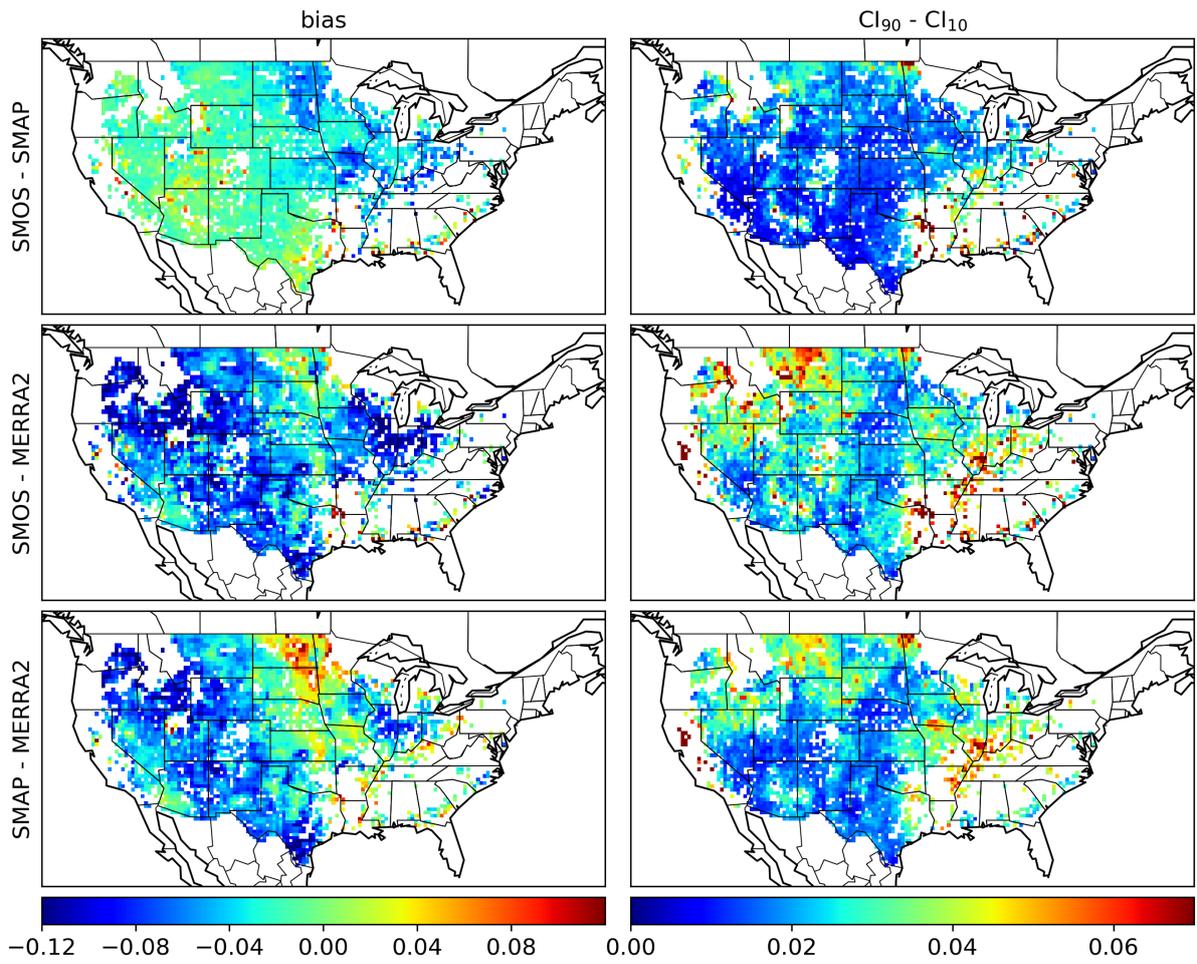


Figure A.4: Temporal mean biases [m^3m^{-3}] (left) and associated 80% confidence intervals (right) between raw soil moisture estimates of SMOS, SMAP and MERRA-2.

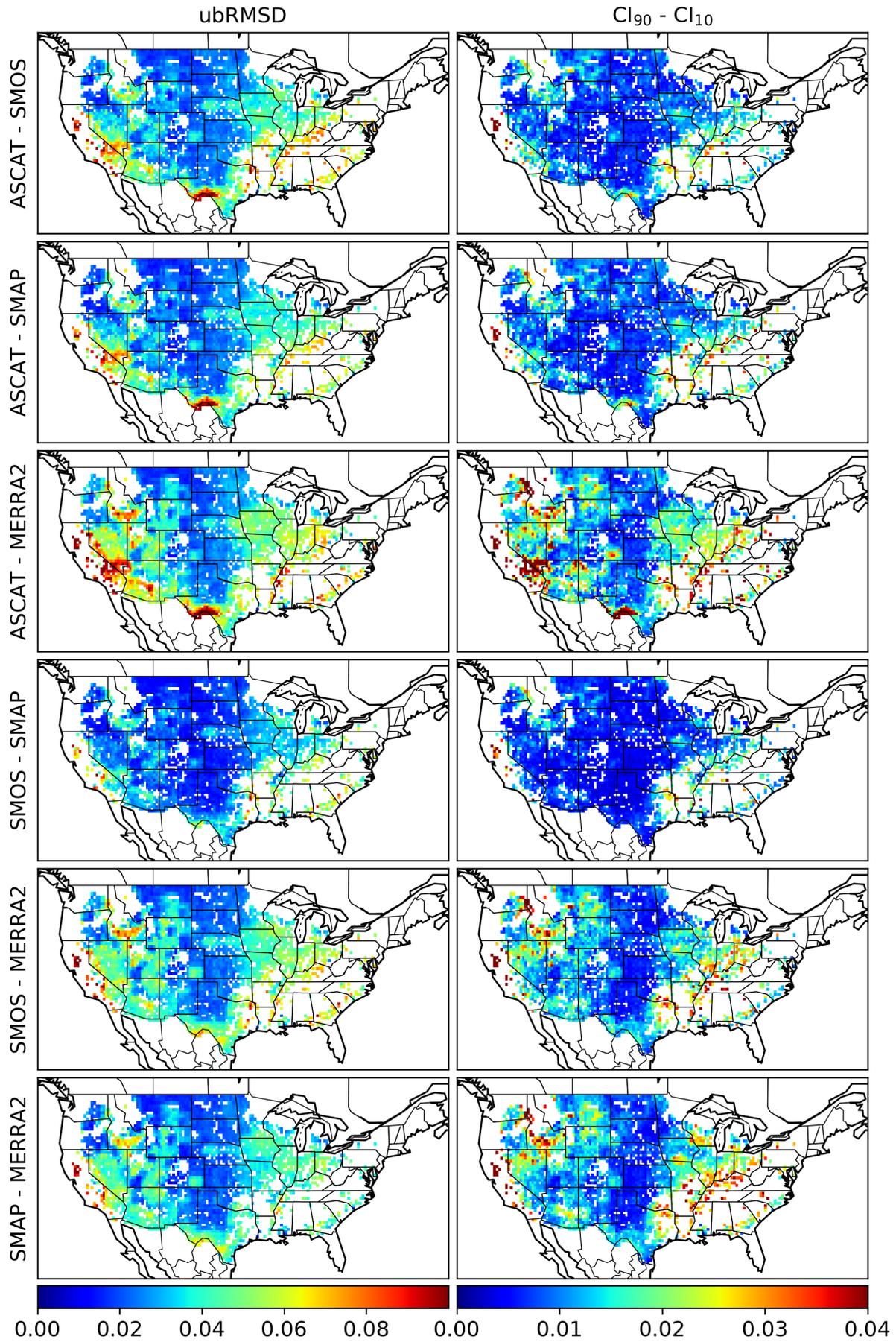


Figure A.5: Unbiased (in mean and standard deviation) root-mean-square-differences [m^3m^{-3}] (left) and associated 80% confidence intervals (right) between raw soil moisture estimates of ASCAT, SMOS, SMAP and MERRA-2. 82

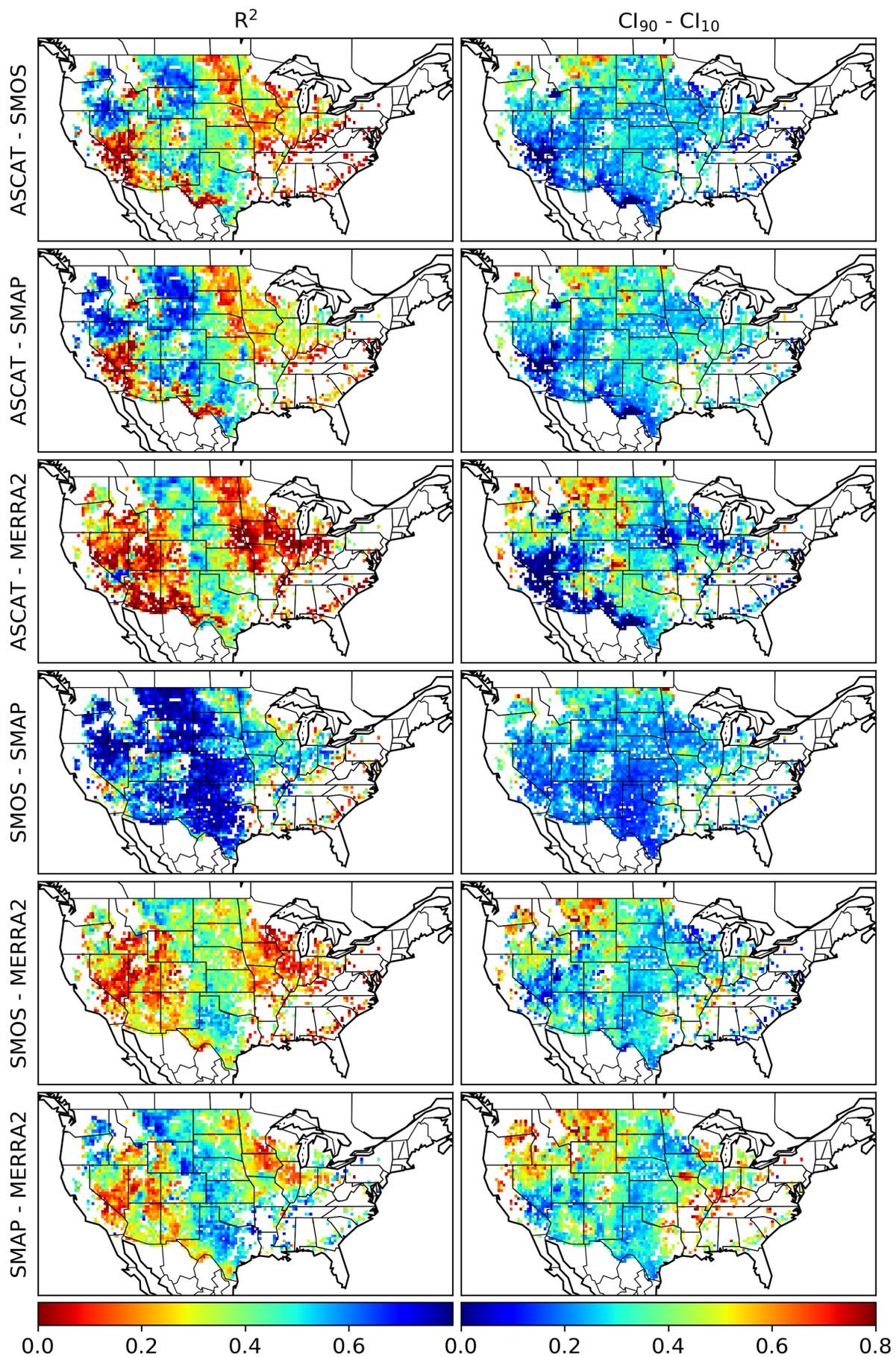


Figure A.6: Coefficients of determination [-] (left) and associated 80% confidence intervals (right) between raw soil moisture estimates of ASCAT, SMOS, SMAP and MERRA-2.

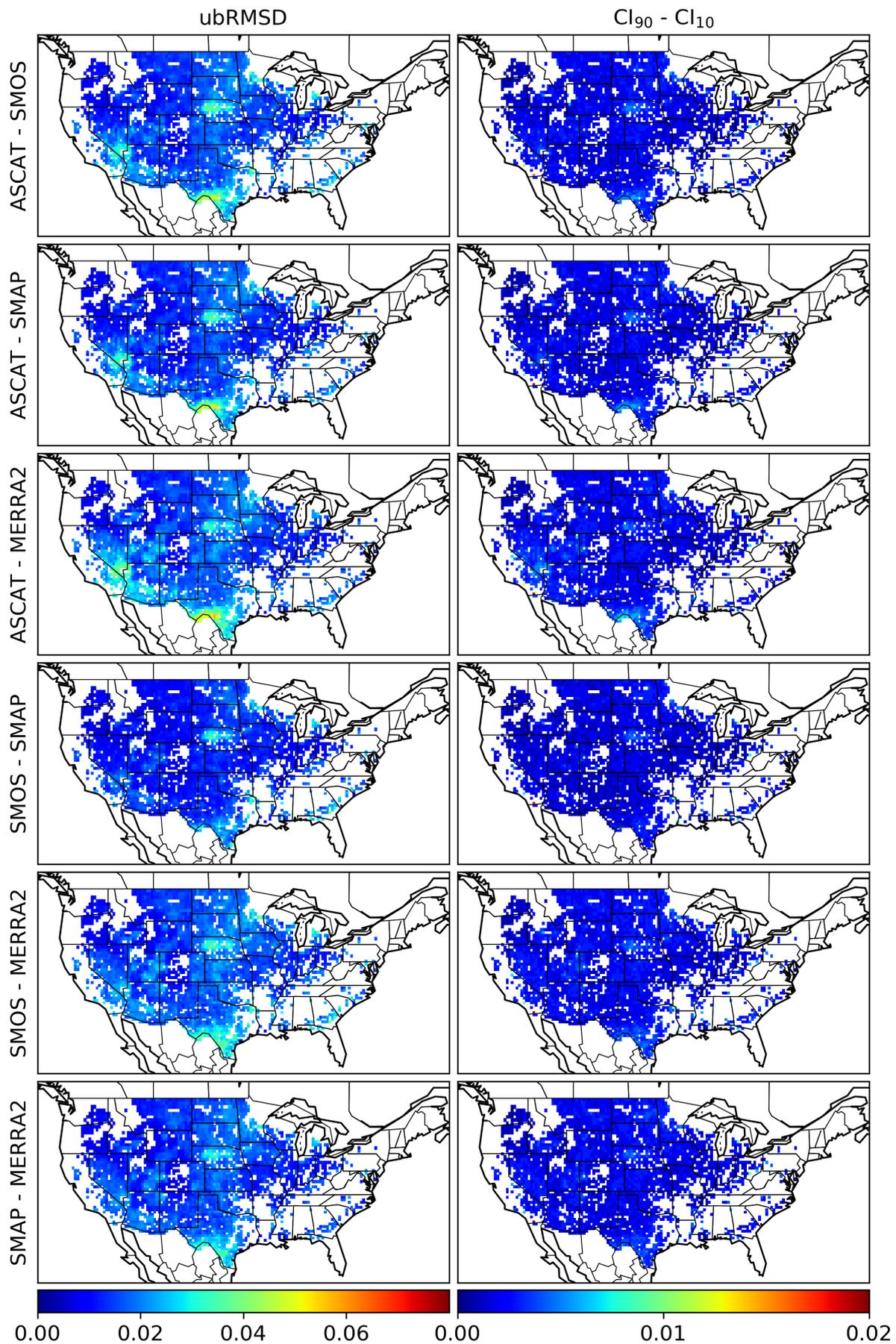


Figure A.7: Unbiased (in mean and standard deviation) $[m^3m^{-3}]$ root-mean-square-differences (left) and associated 80% confidence intervals (right) between soil moisture anomaly estimates of ASCAT, SMOS, SMAP and MERRA-2. 84

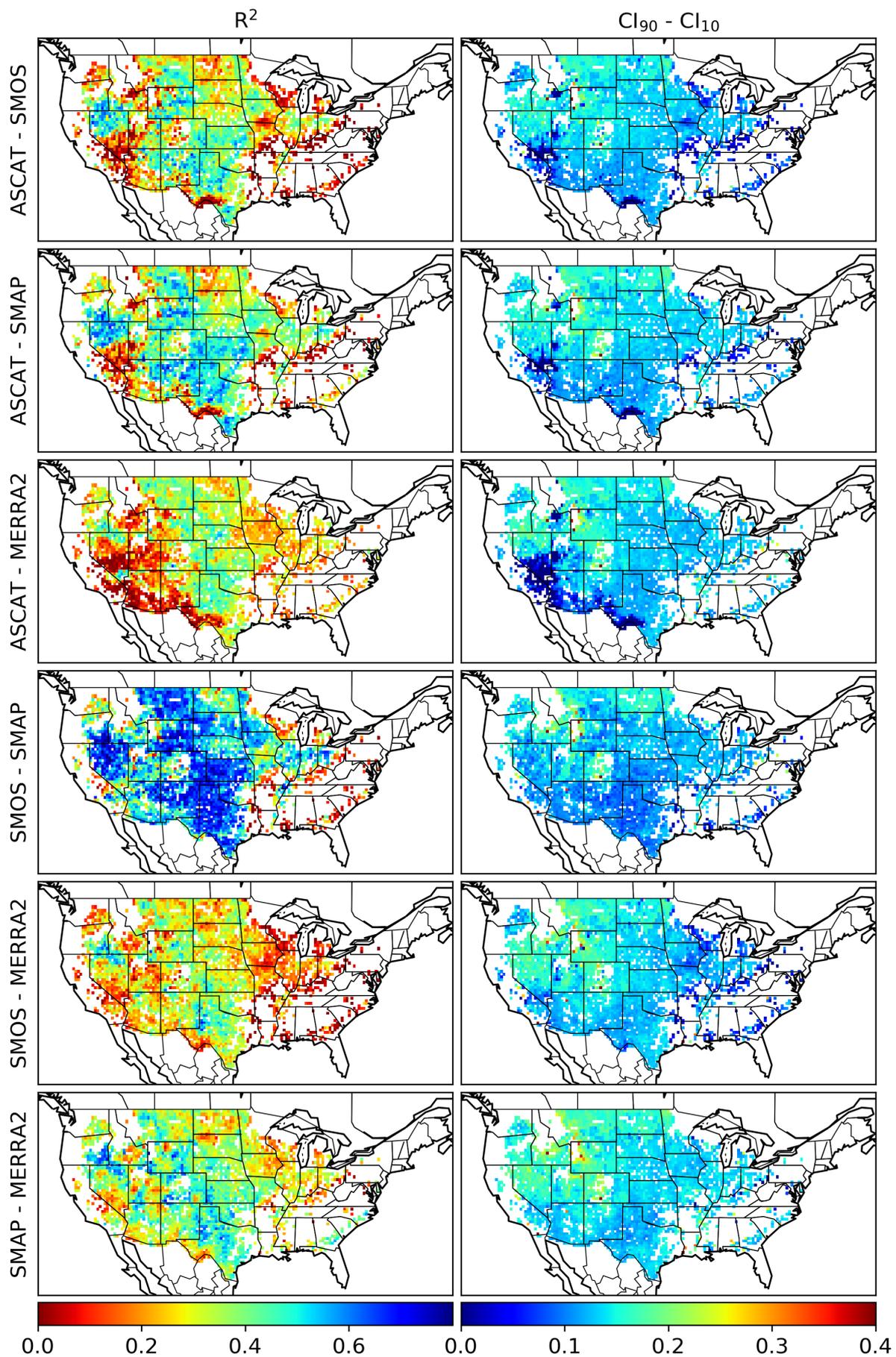


Figure A.8: Coefficients of determination [-] (left) and associated 80% confidence intervals (right) between soil moisture anomaly estimates of ASCAT, SMOS, SMAP and MERRA-2.

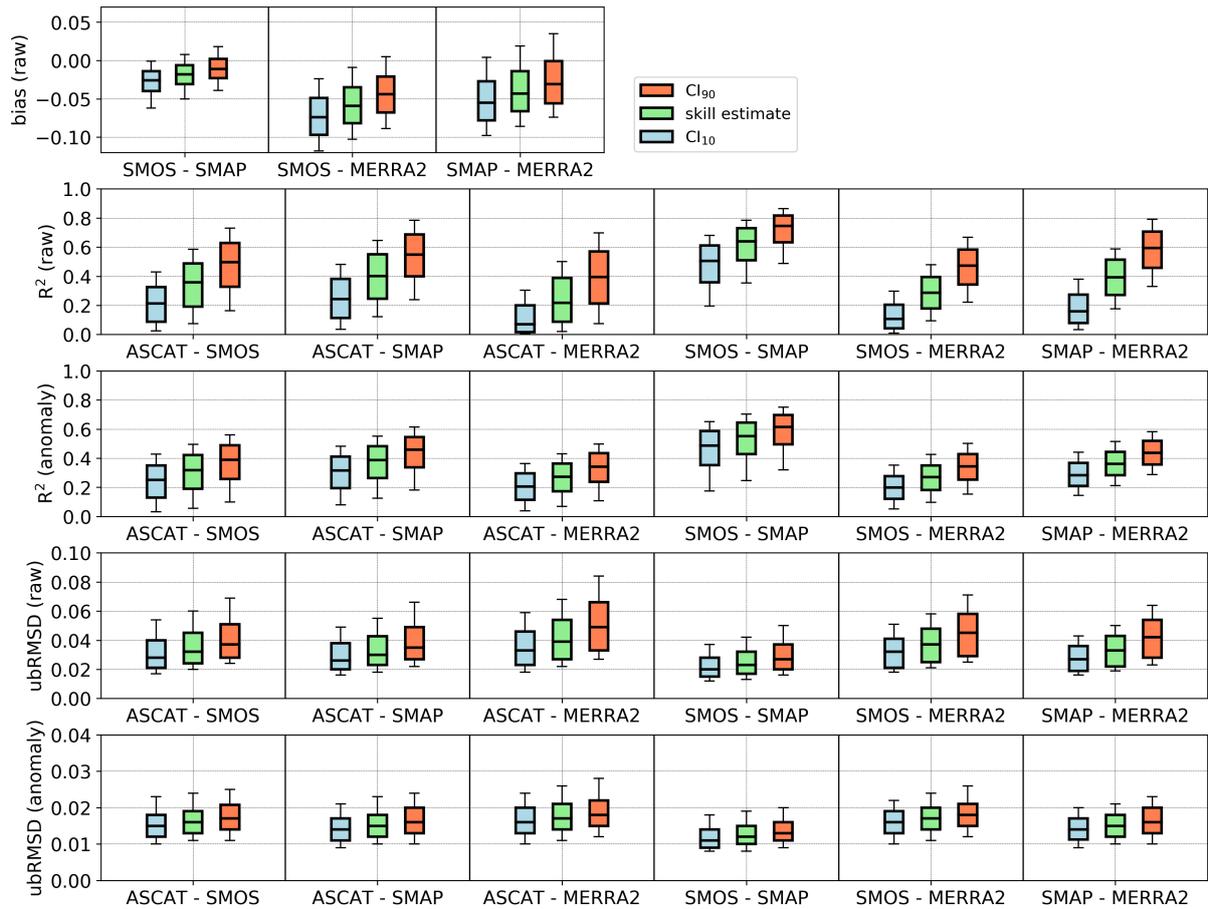


Figure A.9: Spatial summary statistics of biases [m^3m^{-3}], ubRMSDs [m^3m^{-3}], and coefficients of determination [-] and their 10% and 90% confidence limits, respectively, for raw soil moisture estimates and soil moisture anomalies of ASCAT, SMOS, SMAP and MERRA-2. Boxes represent the (spatial) median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

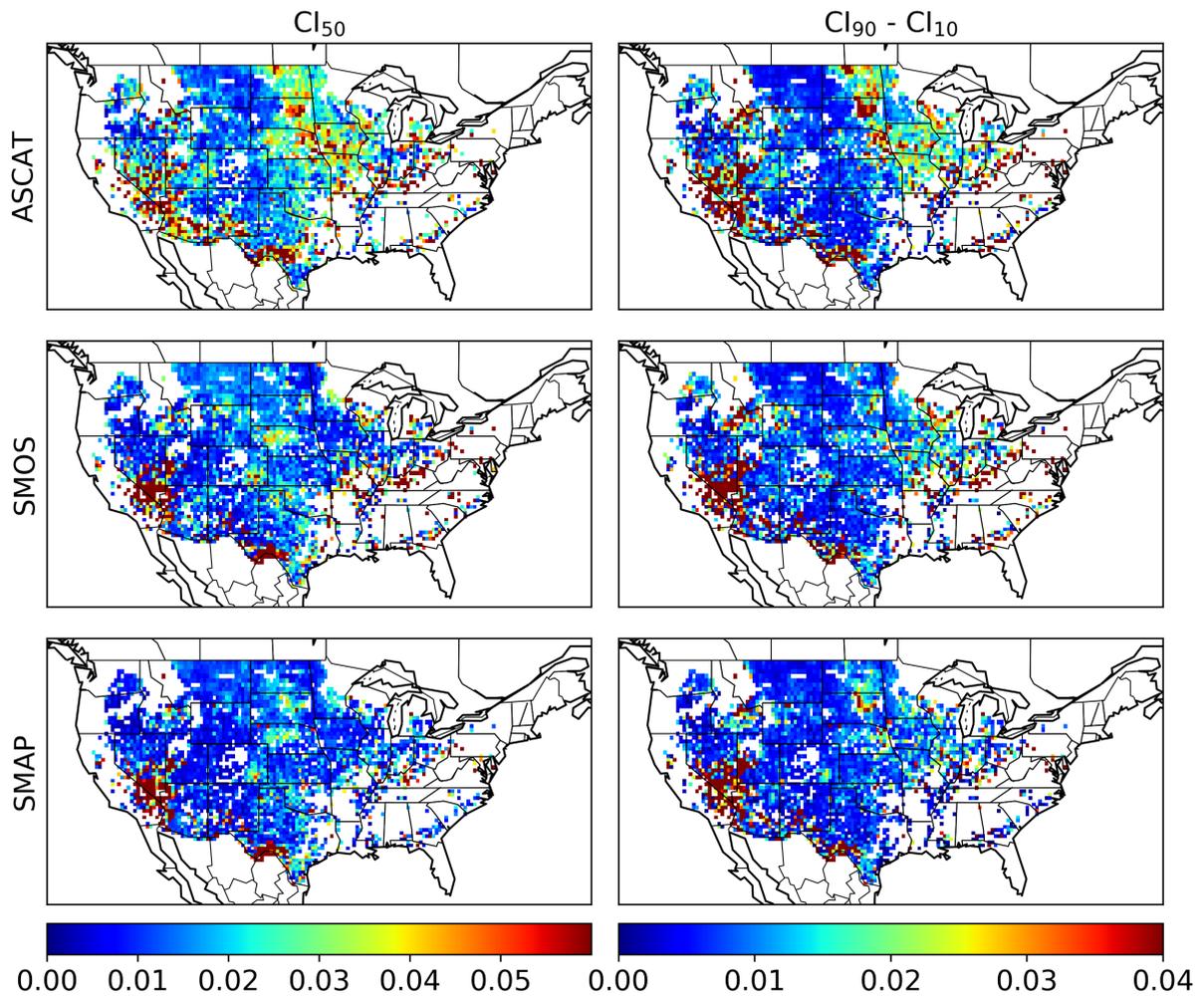


Figure A.10: Median of the bootstrapped TCA-based ubRMSEs [m^3m^{-3}] (left) and associated 80% confidence intervals (right) of raw soil moisture estimates of ASCAT, SMOS, and SMAP.

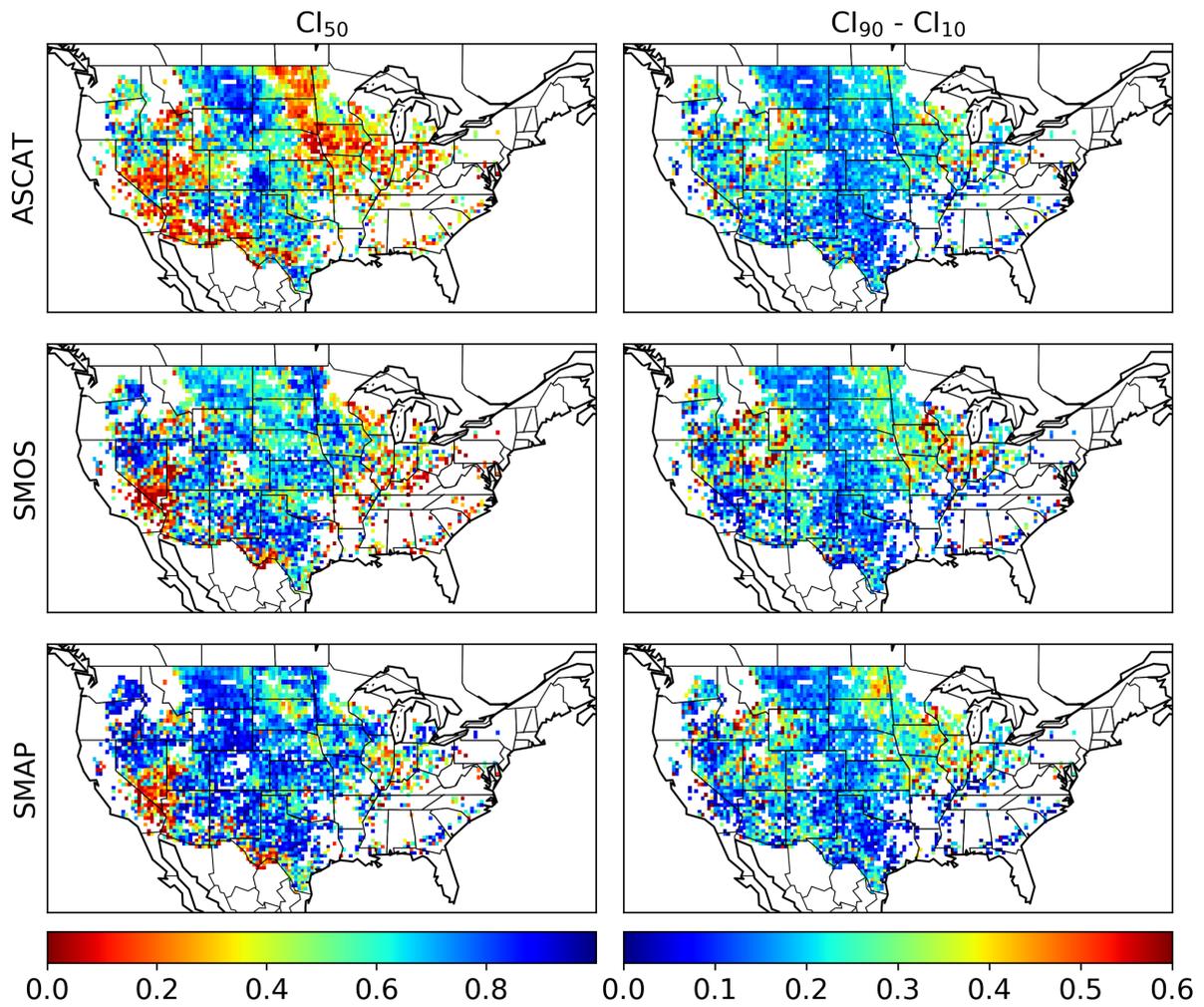


Figure A.11: Median of the bootstrapped TCA-based R^2 estimates [-] (left) and associated 80% confidence intervals (right) of raw soil moisture estimates of ASCAT, SMOS, and SMAP.

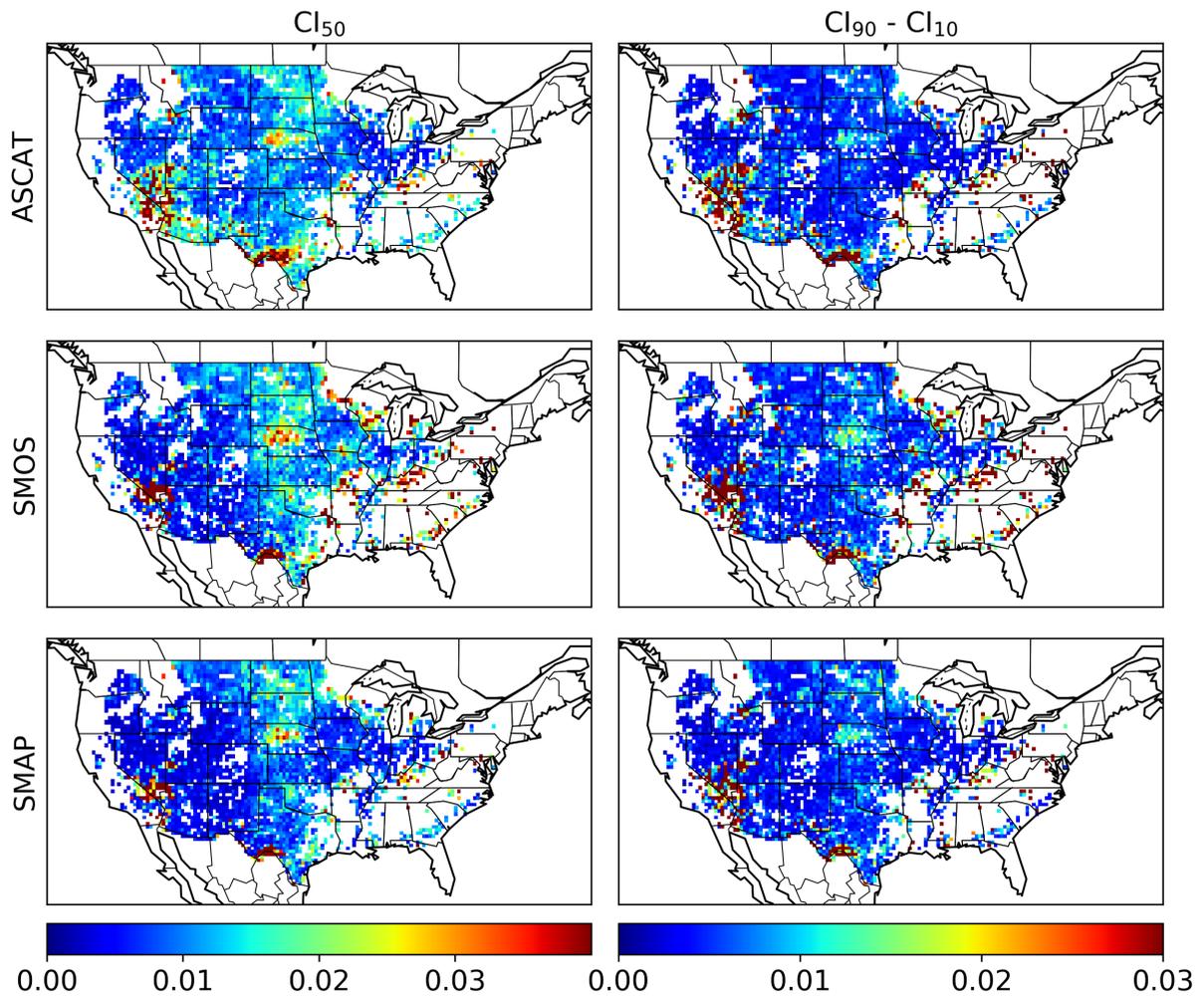


Figure A.12: Median of the bootstrapped TCA-based ubRMSEs [m^3m^{-3}] (left) and associated 80% confidence intervals (right) of soil moisture anomaly estimates of ASCAT, SMOS, and SMAP.

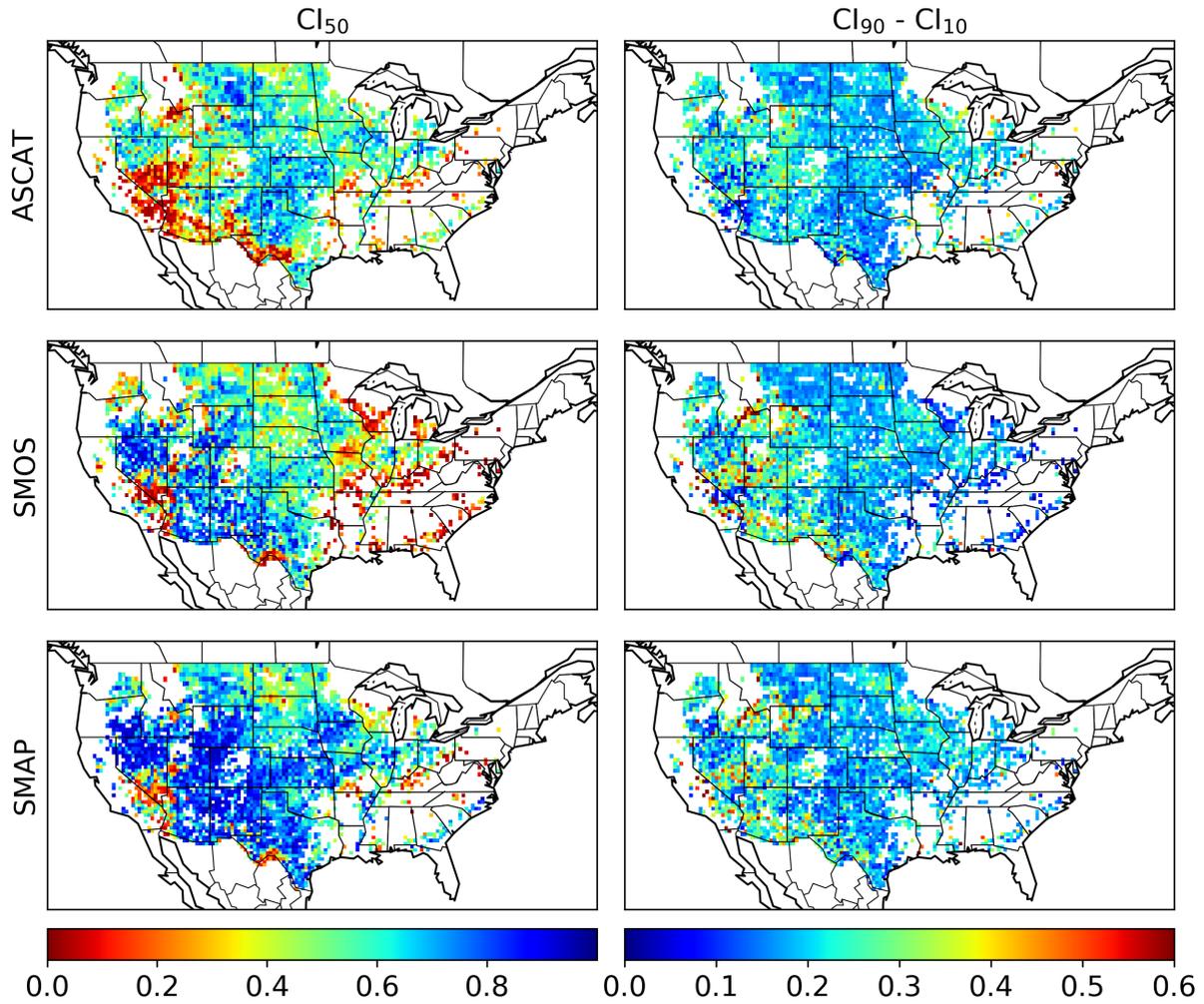


Figure A.13: Median of the bootstrapped TCA-based R^2 estimates [-] (left) and associated 80% confidence intervals (right) of soil moisture anomaly estimates of ASCAT, SMOS, and SMAP.

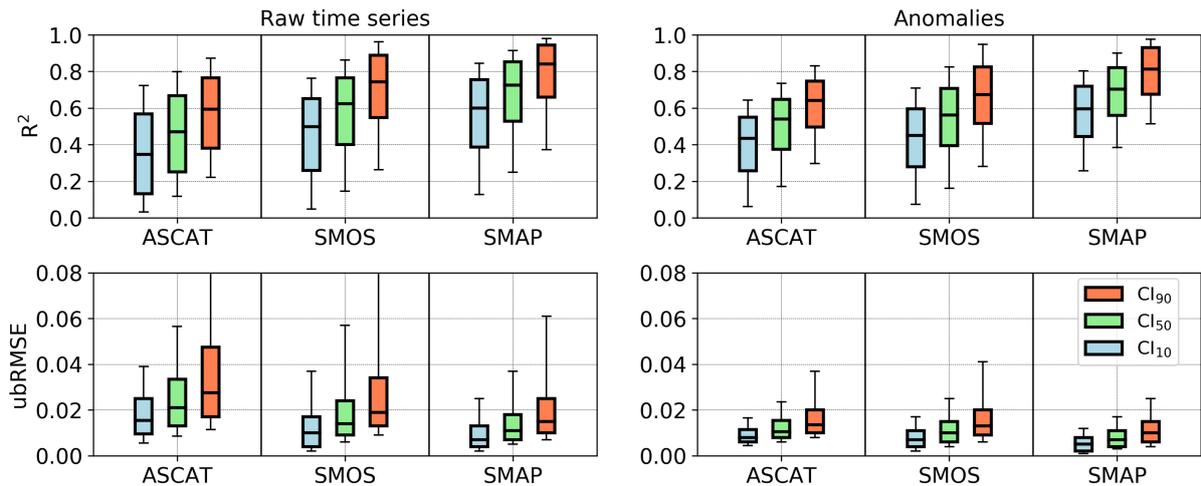


Figure A.14: Spatial summary statistics of the median of the bootstrapped TCA-based ubRMSEs [m^3m^{-3}], and R^2 estimates [-] and their 10% and 90% confidence limits, respectively, for raw soil moisture estimates and soil moisture anomalies of ASCAT, SMOS, and SMAP. Boxes represent the (spatial) median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

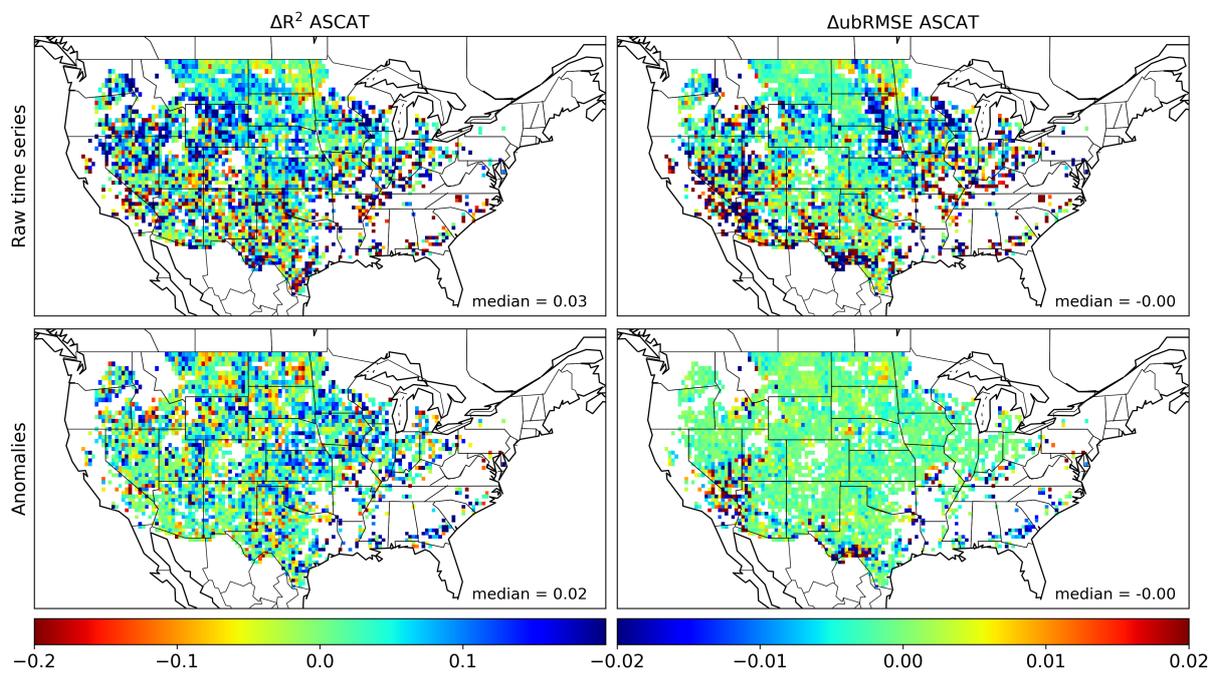


Figure A.15: Difference in TCA-based ubRMSE [m^3m^{-3}] and R^2 estimates [-] for raw soil moisture estimates (top) and soil moisture anomaly estimates (bottom) of ASCAT when using SMOS as third data set minus when using SMAP as third data set in the triplet.