1 Integrating process-based-related information into an ANN for

2 root-zone soil moisture prediction

- 3 Roiya Souissi¹, Mehrez Zribi¹, Chiara Corbari², Marco Mancini², Sekhar Muddu³, Sat Kumar
- 4 Tomer⁴, Deepti B Upadhyaya^{3,4}, Ahmad Al Bitar¹
- 5 ¹CESBIO—Centre d'Etudes Spatiales de la Biosphère, Université de Toulouse, CNES/CNRS/INRAE/IRD/UPS,

6 Toulouse, France

- ⁷ ²Department of Civil and Environmental Engineering (DICA), Polytechnic University of Milan, 20133 Milano, Italy
- 8 ³Department of Civil Engineering, Indian Institute of Science, Bangalore 560 012, India
- 9 ⁴Satyukt analytics Pvt Ltd, Sanjay Nagar Main Rd, MET Layout, Bengaluru, Karnataka 560094, India
- 10 Correspondence to: Roiya Souissi (roiya.souissi@cesbio.cnes.fr)

11 Abstract. Quantification of root-zone soil moisture (RZSM) is crucial for agricultural applications and soil sciences. 12 RZSM impacts processes such as vegetation transpiration and water percolation. Surface soil moisture (SSM) can be 13 assessed through active and passive microwave remote sensing methods, but no current sensor enables direct RZSM 14 retrieval. Spatial maps of RZSM can be retrieved via proxy observations (vegetation stress, water storage change, and 15 surface soil moisture) or via land surface model predictions. In this study, we investigated the combination of surface 16 soil moisture information with process-based-related inferred features involving artificial neural networks (ANNs). We 17 considered the infiltration process through the soil water index (SWI) computed with a recursive exponential filter and 18 the evaporation process through the evaporative evaporation efficiency computed based on a MODIS remote sensing 19 dataset and simplified analytical model, while vegetation growth was not modeled and only inferred expressed-through 20 normalized difference vegetation index (NDVI) time series. Several ANN models with different sets of features were 21 developed. Training was conducted considering in situ stations distributed several areas worldwide characterized by 22 different soil and climate patterns of the International Soil Moisture Network (ISMN), and testing was applied to 23 stations of the same data hosting facility. The results indicate that the integration of process-based-related features into 24 ANN models increased the overall performance over the reference model level in which only SSM features were 25 considered. In arid and semi-arid areas, for instance, performance enhancement was observed when the evaporative 26 evaporation efficiency was integrated into the ANN models. To assess the robustness of the approach, the trained 27 models were applied on observation sites in Tunisia, Italy and South-India that are not part of ISMN. The results reveal 28 that joint use of surface soil moisture, evaporative evaporation efficiency, NDVI and recursive exponential filter 29 represented the best alternative for more accurate predictions in the case of Tunisia, where the mean correlation of the 30 predicted RZSM based on SSM only sharply increased from 0.443 to 0.801 when process-based-related features were 31 integrated into the ANN models in addition to SSM. However, process-based-related features have no to little added 32 value in temperate to tropical conditions.

33 Keywords: root-zone soil moisture, artificial neural networks, evaporative evaporation efficiency, exponential filter.

34 1 Introduction

35 Soil moisture is a major land parameter integrated into several agricultural, hydrological and meteorological 36 applications (Koster et al., 2004; Anguela et al., 2008) This essential climate variable (ECV) consists of two 37 components, namely, surface soil moisture (SSM) (0-5 cm) and root-zone soil moisture (RZSM)-(30 cm to 1 m). 38 RZSM corresponds to the soil moisture in the region in which the main vegetation rooting network is developing. Its 39 definition varies depending on vegetation type and pedoclimatic conditions. The importance of RZSM is mainly 40 highlighted in agricultural applications through vegetation stress and water needs and in carbon and nitrogen cycles, as 41 RZSM influences biogeochemical activities in soil (Martínez-Espinosa et al., 2021). RZSM is nonlinearly related to 42 SSM through different hydrological processes, such as diffusion processes. RZSM may be extracted by evaporation at 43 the surface, through root extraction or by capillary rises (Calvet et Noilhan, 2000). SSM quantification is achieved 44 through three main sources: in situ measurements, model estimates and remote sensing-based products. Microwave 45 remote sensing technologies involving sensors such as the Soil Moisture and Ocean Salinity (SMOS) mission (Kerr et 46 al., 2010), Soil Moisture Active Passive (SMAP) mission (Entekhabi et al., 2010) Advanced Microwave Scanning 47 Radiometer (AMSR) (Owe et al., 2008) and Advanced Scatterometer (ASCAT) (Wagner et al., 2013) haves been 48 employed to retrieve SSM at coarse resolutions. Current satellite sensors can only provide surface soil moisture 49 information because of the shallow penetration depth of spaceborne data (on the order of a few centimeters) (Wagner et 50 al., 2007). Fine-spatial resolution synthetic aperture radar (SAR) data can also be applied in synergy with optical data 51 to retrieve soil moisture (Zribi et al., 2011; Hajj et al., 2014; Dorigo et al., 2011), but again for surface soil moisture. 52 The International Soil Moisture Network (ISMN) is an exhaustive data hosting facility focused on soil moisture data 53 and associated ancillary information. The ISMN provides in situ soil moisture measurements collected from operational 54 soil moisture networks worldwide (Dorigo et al., 2011). Various models can be adopted to estimate RZSM, such as 55 land surface models (Surfex (Masson et al., 2013), ISBA (Noilhan et al., 1996), CLM (Oleson et al., 2010), JULES 56 (Best et al., 2011), etc.) or dedicated crop models such as Aquacrop (Raes et al., 2009) or SAFYE (Battude et al., 57 2017). While these models provide the advantage of physical process-based estimates, these estimates depend on the 58 availability and accuracy of ancillary information. Model predictions are often enhanced by the implementation of data 59 assimilation techniques, such as the land data assimilation system (LDAS) (Sabater et al., 2007; Entekhabi et al., 2020).

60 Data-driven methods such as artificial neural networks (ANNs) have also been commonly applied in hydrology as 61 detailed for instance by the ASCE Task Committee on Application of Artificial Neural Networks in Hydrology (2020) 62 and in (Tanty el al., 2015). One of their advantages is that these models do not require an explicit model structure to 63 accurately represent the involved hydrological processes but instead construct a relationship between the given inputs 64 and the process of interest. Therefore, ANNs are regarded as dynamic input-output mapping models heavily relying on 65 the provided training data relevant to target values (Pan et al., 2017). Moreover, ANNs only require a one-time 66 calibration to provide soil moisture estimations once instrument data are loaded and thus generate relatively low 67 computational costs (Kolassa et al., 2018). These advantages explain the approach to estimate RZSM based on surface 68 information with ANNs in various methodologies (Pan et al., 2017; Grillakis et al., 2021; Souissi et al., 2020). In this 69 paper, we do not address ANN applications as a model twin where the ANN model is trained on the target for 70 mimicking purposes and subsequently generates predictions while requiring a short computation time or fewer input 71 simplifications. Here, we are instead interested in the adoption of ANNs as independent models trained on in situ 72 observations. Within this context, Pan et al. (2017) (Pan et al., 2017)-successfully applied an ANN as a model for 73 shallow 20-cm root zone soil moisture prediction with a global correlation coefficient of 0.7. Grillakis et al. (2021) 74 (Grillakis et al., 2021)-proposed employing an ANN as a means to calibrate and regionalize the time constant of a 75 recursive exponential filter, which was thereafter applied at the regional scale. A combined implementation of Bayesian

76 probabilistic approach and an ANN to infer RZSM at different depths from optical UAV acquisitions via local training 77 was also applied (Hassan-Esfahani et al., 2017). Multitemporal averaged features to predict RZSM based on only SSM 78 and investigated the transferability of a trained ANN across different climatic conditions globally were proposed in 79 (Souissi et al., 2020). Temporal information can be considered in ANNs through recurrent neural networks (RNNs), 80 long short-term memory (LSTM) architectures (Liu et al., 2021), 1D convolutional neural networks (CNNs), or 81 multitemporal averaging. In (Souissi et al., 2020), median, maximum, and minimum correlation values of 0.77, 0.96, 82 and 0.65 were respectively reported across training, validation and test datasets. The use of climatic variables such as 83 precipitation and surface temperature and intrinsic surface properties such as soil texture and land cover has also been 84 considered in ANNs (Liu et al., 2021). The choice of variables depends not only on the data availability but also on the objectives. Finally, ANN-based approaches pertain to the more general term of machine learning (ML)-approaches, and 85 86 within this framework, the random forest approach has been applied to root zone soil moisture prediction (Carranza et 87 al., 2021). The aforementioned studies have investigated the application of multiple information sources to predict root 88 zone soil moisture. The input features are commonly curated for quality, and correlation analysis is conducted to 89 determine the useful inputs, while physical processes are not considered. In this paper, we introduce process-based 90 related features based on simplified analytical models representing the major processes contributing to root zone soil 91 moisture dynamics. In this work, RZSM refers to a point observation of water content in a depth ranging between 30 92 and 55cm. We investigate the impact of the application of different process-based-related variables on the precision of 93 RZSM predictions as well as the robustness of our approach. (1) We start from a previously developed MLP-ANN 94 model (Souissi et al., 2020), and we extend the feature list to include NDVI time series, surface soil temperature and 95 process-based-related variables, namely, the soil water index given by a recursive exponential filter, and remote 96 sensing-based evaporative evaporation efficiency, and NDVI time series. (2) The robustness of the approach is assessed 97 through additional tests involving stations not included in the ISMN database in Tunisia, Italy, and South-India. (3) 98 Climatic analysis is conducted to infer the most indicative process-based-related features for each climate pattern. 99 2 **Materials and Methods** 100 The proposed methodology entails the construction of several ANN models with both direct (SSM, surface temperature, 101 and NDVI) and intermediate sets of features (soil water index and evaporative evaporation efficiency) computed based 102 on simplified analytical models. An overview of the processing configuration is shown in Figure 1. Standard scaling is

- applied to each dataset separately so that the different inputs fall into the same range of values, then the ANN outputs
- 104 <u>are descaled to make the comparison with actual values of RZSM possible.</u>



Figure 1. Overview of the processing configuration- showing the components of the model: the tested models are variations of this
 ANN with a different combination of inputs (see Table 1). The scaling and descaling are applied to each dataset separately.

108 This approach results in a combination of ANN models (Table 1). Each model has one <u>or more process-related</u>

109 physical process based or a geophysical_features in addition to the three SSM features which correspond to backward

110 <u>rolling averages of in-situ SSM computed over 10,30 and 90 days</u>. All the ANN model hyperparameters remain the

same except the number of input features, as described at the end of this section.

105

Table 1. ANN model configurations with the respective input variables : (*: rolling averages of SSM over 10 days; **: rolling

averages of SSM over 30 days; ***: rolling averages of SSM over 90 days; ****: number of parameters of the ANN model.--

Model	SSM 104 PAV*	SSM 204 DAV**	SSM OOD DAV***	сст	NDVI	SWI	EVAD	NIb****
Woder	55W_104_KAV	55W_50U_KAV	55WI_700_KAV	166	ND VI	5 11	LVAI	110
Features								
ANN_SSM	Х	Х	Х					<u>101</u>
ANN_SSM_TEMP	Х	Х	Х	Х				121
ANN_SSM_NDVI	Х	Х	Х		Х			<u>121</u>
ANN_SSM_EXP-	Х	Х	Х			Х		<u>121</u>
FILT-T5								

ANN_SSM_EVAP-	Х	Х	Х		Х	<u>121</u>
EFF_B60						
ANN_SSM_NDVI_E	Х	Х	Х	Х	X X	<u>161</u>
VAP-EFF-B60_EXP-						
FILT-T5						

The model with the simplest starting point is ANN_SSM based on (Souissi et al., 2020). The most complex model includes the full set of inputs. Intercomparison of the model performance provides information on the added value of each input. All input features are scaled, and training is performed on each of these features based on scaled in situ RZSM data retrieved from the ISMN. The RZSM model predictions are validated against an independent set of observations.

120 **2.1 Datasets**

121 2.1.1 ISMN soil moisture data

122 The first training and test operations were conducted on eight ISMN networks previously considered in (Souissi et al., 123 2020). Figure 2 shows the distribution of the considered soil moisture networks with different soil textures and climatic 124 parameters. The selected stations exhibit a root zone depth varying between 30 and 60 cm (Table 2). (cf. appendix B). 125 For each station, the RZSM observation point is located between 30 and 55cm (Table 2). For each soil moisture hourly 126 acquisition, ISMN provides quality flags. Quality flags can be marked as 'C' (exceeding plausible geophysical range),' 127 D' (questionable/dubious), 'M' (missing), or 'G' (good) (Dorigo et al.,2011). Category 'D' has subset flags namely 128 'D01' for which in situ soil temperature $\leq 0^{\circ}$ C, 'D02' that flags points at which in situ air temperature $\leq 0^{\circ}$ C as well as 129 'D03' that also flags areas where GLDAS soil temperature $< 0^{\circ}$ C. In our study, only soil moisture data which quality 130 flag is marked 'G' were retained.



- 132 Figure 2. International Soil Moisture Network (ISMN) network distribution (adapted from the ISMN web data portal
- 133 (https://www.geo.tuwien.ac.at/insitu/data_viewer/); scale: 1 cm=1000 km).

	0	Number of Selected	Selected RZSM Depth	SM
Inetwork	Country	Stations	(cm)	Sensors
AMMA-CATCH	Benin, Niger	5 (3 in Benin and 2 in Niger)	40	CS616
BIEBRZA-S-1	Poland	3	50	GS-3
CTP-SMTMN	China	54	40	EC-TM/5TM
HOBE	Denmark	29	55	Decagon-5TE
FR-Aqui	France	5	30, 34, 50	ThetaProbe ML2X
OZNET	Australia	19	30	Hydra Probe-CS616
SCAN	USA	209	50	Hydraprobe-Sdi-12/Ana
SMOSMANIA	France	22	30	ThetaProbe ML2X

134 **Table 2.** Overview of the considered ISMN and external networks.

135

136 **2.1.2 External soil moisture data**

137 The external networks only considered to assess the transferability and robustness of the approach were employed for 138 validation. The trained models are run for predictions only over these sites. They have been selected to cover semi-arid, 139 moderate and tropical semi-arid climates.

- <u>Tunisian site</u>: The Merguellil site is located in central Tunisia (9°54 E; 35°35 N). This site is characterized by a semiarid climate with highly variable rainfall patterns <u>(average equal to 300mm/year)</u>, very dry summer seasons, and wet winters. The Merguellil site represents an agricultural region where croplands, namely, olive groves and cereal fields, prevail (Zribi et al., 2021). At this study site, a network of continuous thetaprobe stations installed at bare soil locations provided moisture measurements at depths of 5 and 40 cm. All measurements were calibrated against gravimetric estimations. Data were obtained from the Système d'Information Environmental (SIE) web application catalog.
- 147 Italian site: The Landriano site is located in northern Italy (Pavia province, Lombardia region). This station is • 148 located in a maize field, which was monitored in 2006 and from 2010 to 2011 (Masseroni et al., 2014). The soil 149 texture is sandy loam, The average rainfall in Pavia province is of 650-700 mm, the climate is classified as 150 'Cfa' (cf. appendix A) and the field is irrigated by the border method with an average irrigation amount of 151 approximately 100 to 200 mm per application with one to two applications per season due to the presence of a 152 shallow groundwater table. Soil moisture measurements were performed with time domain reflectometer 153 (TDR) soil moisture sensors. Five TDR soil moisture sensors were installed along a profile at depths of 5, 20, 154 35, 50, and 70 cm.

- 155 -Indian site: The Berambadi watershed is located in Gundalpet taluk, Chamarajanagara district, in the southern 156 part of Karnataka state in India and covers an area of approximately 84 km². The aridity index (P/PET) is 0.7, 157 with an average rainfall of is equal to -800 mm/year and a PET value of 1100 mm/year. The and the climate is 158 classified as Aw (cf. appendix A)., and the major soil types in the region vary between sandy loam, sandy clay 159 loam and sandy clay. Hydrological variables have been intensively monitored since 2009 in the Berambadi 160 watershed by the Environmental Research Observatory ORE BVET and AMBHAS Observatory. The soil 161 moisture levels at the surface (5 cm) and root zone (50 cm) are monitored with a HydraProbe sensor at different 162 agricultural sites across the watershed, and in the current study, 4 stations were chosen. The 3-major cropping 163 seasons include kharif (June to September), during which the first crop is grown, which is usually rainfed 164 during the rabi season (October to January), and summer (February to May), during which the second and third 165 crops are grown, which are usually irrigated. The major crops grown in the region include turmeric, maize, 166 sunflower, marigold and vegetables.
- 167

168 **2.1.3 Surface soil temperature**

In addition to in situ soil moisture, the ISMN optionally includes meteorological and soil variables that are available over specific time periods. Values of the situ surface soil temperature among these variables can be employed as a useful indicator of the soil moisture data quality. The soil temperature was provided in Celsius, and the plausible values range from -60 to 60 °C. Regarding soil moisture data, surface soil temperature data were also provided with quality flags (Dorigo et al., 2011). However, the drawback is that this variable is not available in all networks, which is the case with the AMMA-CATCH network.

175 **2.1.4 Normalized difference vegetation index**

176 We considered the remote sensing-based normalized difference vegetation index (NDVI) to quantify-infer vegetation 177 dynamics. We extracted this index from the Moderate Resolution Imaging Spectroradiometer (MODIS) Vegetation 178 Indices product (MOD13Q1 version 6). MODIS Vegetation Indices (MOD13Q1) version 6 data are generated at 16-179 day intervals and a 250-m spatial resolution as a Level 3 product. This product provides two primary vegetation layers. 180 The first vegetation layer is the NDVI, which is referred to as the continuity index of the existing National Oceanic and 181 Atmospheric Administration-Advanced Very High Resolution Radiometer (NOAA-AVHRR)-derived NDVI. The 182 algorithm chooses the best available pixel value from all the acquisitions over the 16-day period. The criteria 183 considered are low cloud coverage, low view angle, and highest NDVI value (Huete et al., 1999 Didan, 2015). To 184 obtain daily NDVI values, we conducted linear interpolation of the 16-day product.

185 **2.1.5 Potential evapotranspiration**

Similarly, we assessed the impact of considering a remote sensing-based <u>evaporative-evaporation</u> efficiency, <u>which is</u> <u>initially defined as the ratio of actual to potential soil evaporation</u>, on RZSM prediction. The computation details of this variable will be detailed later (cf. Section <u>32.2.2</u>). We employed the remote sensing-based potential evapotranspiration (PET) to compute the <u>evaporative evaporation</u> efficiency. We extracted the PET from the MOD16A2 Evapotranspiration/Latent Heat Flux version 6 product, which is an 8-day composite dataset produced at a 500-m pixel resolution. The algorithm used for MOD16 data product collection is based on the logic of the Penman–Monteith equation, which employs inputs of daily meteorological reanalysis data along with MODIS remote sensing data 193 products such as vegetation property dynamics, albedo, and land cover. The MOD16A2 product provides layers for the

194 composite evapotranspiration (ET), latent heat flux (LE), potential ET (PET) and potential LE (PLE). The pixel values

195 for the PET layer include the sum of all eight days within the composite period (Running et al., 2017). To obtain daily

196 PET values, we performed <u>a linear interpolation of over the 8-day product and then we divided by eight the interpolated</u>

197 <u>value</u>.

198 **2.2 Methods**

199 2.2.1 Recursive exponential filter

Two ANN models presented in Table 1 contained extra knowledge on infiltration process information based on the
 outputs of the recursive exponential filter (Stroud, 1999) as a feature. The recursive exponential filter was first
 introduced by (Wagner et al., 1999) Wagner et al. (1999) -to estimate the soil water index (SWI) from surface soil
 moisture. The equation for the recursive formulation can be written as follows SWI is computed as follows:

204	SIMI	-SIMI	$\perp K$	(mc(t))	-SM/I	(1)
204	m(n)	$-5m_{m(n-1)}$	n _n	$(ms(r_{H})$	5111	m(n=1)

205

 $SWI_{t_n} = SWI_{t_{n-1}} + K_n(ms(t_n) - SWI_{t_{n-1}})$ <u>(1)</u>

206 where:

207	- $\frac{SWI_{m(n)}}{SWI_{m(n)}}$ SWIt _n is the soil water index at time t _n .
208	$- ms(t_n)$ is the estimated scaled-surface soil moisture at time t _n (scaled between maximum and
200	minimum voluos)
209	<u>imminum values)</u> ,
210	- K_n is the gain at time t_n , which occurs in [0,1] and is given by:
211	$K_n = \frac{K_{n-1}}{K_{n-1} + e^{-\frac{(t_n - t_{n-1})}{T}}}$ (2) and
212	T is a time constant and is the only required tuning parameter to compute the recursive
213	exponential filter.
214 215	- For the initialisation of the filter, gain $K_1 = 1$ and $SWI_{(t1)}^* = ms(t_1)$ Regarding T values, we considered an empirical list ([1,3,5,7,10,13,15,20,40,60]), which was partly inspired by (Paulik
216	et al., 2014) (T \in [1,5,10,15,20,40,60,100]). Given the list of T values, recursive exponential filter outputs were
217	computed for all of the stations (346 stations) given each T value. Based on the correlation values between the in situ
218	RZSM values and the recursive exponential filter-based RZSM pre-estimates, we established the optimal time variable
219	T, hereafter referred to as T _{best} , for each station.
220	A large proportion of the stations attained an optimal time constant (Tbest) value equal to 60 days which suggests an
221	abnormally long infiltration time. These stations belong to the SCAN network and exhibit an RZSM acquisition depth
222	of 50 cm, in contrast other networks such as SMOSMANIA, for instance, where RZSM is retrieved at 30 cm. The high
223	values correspond to correlation with seasonal dynamics rather than infiltration processes. This depth could explain the
224	anomalously long infiltration time. This has been demonstrated in (Paulik et al., 2014), who demonstrated that the
225	average T value with the highest correlation (T _{best}) increased with increasing depth of the in situ observations.
226	For comparison purposes, (Paulik et al., 2014) found that 23.98% of the stations achieved T _{best} =5 days, while 21.58% of
227	the stations achieved $T_{best} \ge 60$ days (60 or 100 days).

234 2.2.2 Evaporative efficiency

240 241

242

243

245

248

249

An ANN model with <u>evaporative evaporation</u> efficiency input was also developed. This variable, which is defined as the ratio of the actual to potential soil evaporation, was first introduced in (Noilhan, J. and Planton, 1989; Jacquemin et al., 1990; Lee et al., 1992) and thereafter readapted in (Merlin et al., <u>20112010</u>) to include the soil thickness<u>- and is</u> expressed as follows: In our work, we use a modified evaporation efficiency formulation, based on the third model developed in (Merlin et al., 2010), which can be expressed as follows (cf. appendix C):

$\beta_3 = [\frac{1}{2} -$	$\frac{1}{2}\cos(\pi\theta_L/\theta_{max})]^p$	$for \theta_{L} \leq \theta_{max} (3)$
	$\beta_{3} = 1 for \theta_{L}$	> θ_{max}

$$\beta = \left[\frac{1}{2} - \frac{1}{2}\cos(\pi\theta/\theta_{max})\right]^{P*}$$
(3)

244 where: $-\beta$ is evaporation efficiency

 $-\theta_{\rm L}$ $\theta_{\rm L}$ is the water content in the soil layer of thickness L.

 $-\theta_{max}$ is the maximum soil moisture at each station.

 $- P^{\pm}$ is a parameter computed as follows:

$$P = \frac{1}{2} + A_3 \frac{L - L_{\pm}}{L_{\pm}} \frac{L - L_{\pm}}{B_3} \frac{L - E_{\pm}}{B_3} (4)$$
$$P^* = \frac{PET}{2B} (4)$$

250 $-\theta_{max}$ is the soil moisture at field capacity, as reported in (Noilhan, J. and Planton, 1989; Jacquemin et al.,2511990; Lee et al., 1992), or the soil moisture at saturation, as considered in (Merlin et al., 2011). In our case,252this variable denotes the maximum soil moisture at each station.

253 $-LE_p$ is the potential evaporation. In our case, we replaced this variable with the potential evapotranspiration254(PET) extracted from the MODIS 500-m 8-day product (MOD16A2). P was then replaced by proxy P*. As the255ANN model performed its own calibrations on the set of features, this adaptation of the P term did not impact256the process.

257 $-L_1$ is the thinnest represented soil layer, and A_3 (unitless) and B_3 (W/m²) are the two best fit parameters a 258 priori depending on the soil texture and structure, respectively. As we were interested in the evaporative 259 efficiency at the surface, L=L1=5 cm, P^{*} is thus expressed as:

$$P^* = \frac{PET}{2B_3}(5)$$

- P*, a proxy of parameter P (cf. appendix C), represents an equilibrium state controlled by retention forces in
 the soil, which increase with the thickness L of considered soil and by evaporative demands at the soil surface.
 -PET is the potential evapotranspiration (PET) extracted from the MODIS 500-m 8-day product (MOD16A2).
- 264 The soil evaporation efficiency computed by model 3, developed in (Merlin et al., 2010), decreases when PET
- 265 increases. Retention force and evaporative demand make the term P increase (replaced by P*), as if an increase of
- 266 <u>potential evaporation LE_p (here replaced by PET) at the soil surface would make the retention force in the soil greater.</u>
- 267 Merlin et al. (2010) tested this approach at two sites in southwestern France using in situ measurements of actual
- 268 evaporation, potential evaporation, and soil moisture at five different depths collected in summer. Model 3 was able to
- 269 represent the soil evaporation process with a similar accuracy as the classical resistance-based approach for various soil
- 270 <u>thicknesses up to 100 cm. Merlin et al. (2010) affirm the parameterization of P as function of LE_p (here PET) indicates</u>
- 271 that β cannot be considered as a function of soil moisture alone since it also depends on potential evaporation.
- 272 Moreover, the effect of potential evaporation on β appears to be equivalent to that of soil thickness on β . This
- 273 equivalence is physically interpreted as an increase of retention forces in the soil in reaction to an increase in potential
- 274 <u>evaporation.</u>

275 **2.2.3** Artificial neural network implementation

276 The multilayer perceptron (MLP), which is a multilayer feed-forward ANN, is one of the most widely applied ANNs, 277 mainly in the field of water resources (Abrahart and See, 2007) The multilayer perceptron contains one or more hidden 278 layers between its input and output layers. Neurons are organized in layers such that the neurons of the same layer are 279 not interconnected and that any connections are directed from lower to upper layers (Ramchoun et al., 2016). Each 280 neuron returns an output based on the weighted sum of all inputs and according to a nonlinear function referred to as 281 the transfer or activation function (Oyebode and Stretch, 2019). The input layer, consisting of SSM values and/or other 282 process-based related variables, is connected to the hidden layer(s), which comprises hidden neurons. The final ANN-283 derived estimates of the ANN are given by an activation function associated with the final layer denoted as the output 284 layer, based on the sum of the weighted outputs of the hidden neurons.

285 We started with the ANN model developed in (Souissi et al., 2020), whose architecture consists of one hidden layer of 286 20 hidden neurons, a tangent sigmoid function as the activation function of the hidden layer, a quadratic cost function 287 as the loss function and the stochastic gradient descent (SGD) technique as the optimization algorithm. This model was 288 developed to estimate RZSM based on only in situ SSM information. SSM was not applied as a feature of hourly values 289 but was employed in the form of three features, namely, SSM rolling averages over 10, 30 and 90 days. Two additional 290 Additional ANN models were developed to study, through each model, the impact of the application of the NDVI, 291 which describes vegetation dynamics and the surface soil temperature as features SWI, evaporation efficiency and the 292 surface soil temperature as features. A model combining surface soil moisture, NDVI, evaporative evaporation 293 efficiency and recursive exponential filter was further considered. These ANN models were trained and validated on the 294 122 ISMN stations among the 346 stations of the ISMN based on (Souissi et al., 2020). considered of good quality after 295 a data filtering step as detailed in (Souissi et al., 2020).-Training of the above ANN models was conducted considering 296 70% of these 122 stations. Thirty percent was reserved for validation, and testing was conducted at all-the rest of 297 stations. So in summary, 122 stations were considered for the training/validation of the ANN models and 224 stations, 298 if all input data are available, were used for testing. In a second step, tests were conducted on data external to the ISMN

database namely on sites of Tunisia, Italy and India. The trained models over ISMN are used only in prediction mode
 over these sites. The data for SSM in addition to the other features are used as inputs and RZSM is predicted in outputs.

301 <u>3</u> Results

302 <u>3.1 Exponential filter characteristic time length</u>

A large proportion of the stations attained an optimal time constant (T_{best}) value equal to 60 days which suggests an abnormally long infiltration time. These stations belong to the SCAN network and exhibit an RZSM acquisition depth of 50 cm, in contrast other networks such as SMOSMANIA, for instance, where RZSM is retrieved at 30 cm. The high values correspond to correlation with seasonal dynamics rather than infiltration processes. This depth could explain the anomalously long infiltration time. This is consistent with (Paulik et al., 2014) in which the average T value with the highest correlation (T_{best}) increased with increasing depth of the in situ observations.

309For comparison purposes, Paulik et al. (2014) found that 23.98% of the stations achieved T_{best} =5 days, while 21.58% of310the stations achieved $T_{best} \ge 60$ days (60 or 100 days).

Albergel et al. (2008) considered an average T_{best} value of 6 days for the SMOSMANIA network. This value represented the average T_{best} value for all stations belonging to the SMOSMANIA network. In our case, the average T_{best} value for all stations of the SMOSMANIA network reached 9 days. In this study, an average T_{best} value could be established for each station or each network. However, this is not relevant to our work because we aim to evaluate maps of remote sensing data in next steps, and thus, we did not compute T_{best} at each location. We fixed the value of T to 5 days as a median infiltration time.

317 **3.1–2** Intercomparison of the ANN models

The generated correlation histograms distribution histograms for training, validation and test stations (Fig. 3) and performance metrics presented in Table 3 demonstrateshow that the integration of the considered process-based-related features improved the prediction accuracy in certain cases compared to the reference. <u>Time series of good and less good</u> quality of fit were provided in appendix E for training, validation and test stations using reference model ANN_SSM and the most complex ANN model.

In terms of the NDVI, 55.56% of the stations attained better correlation values with ANN_SSM_NDVI than those
 obtained with ANN_SSM. Additionally, 44.44% of the stations achieved a correlation value higher than 0.7 with model

ANN_SSM_NDVI, versus 38.41% of the total stations achieving a similar correlation value with model ANN_SSM.





In terms of the NDVI, 65.82%, 45.71% and 55.22% stations attained better correlation values with ANN_SSM_NDVI

than those obtained with ANN_SSM for the training, validation and test stations, respectively. RMSE decreased for

44.3%, 40.0% and 40.3% of the stations with ANN SSM NDVI compared to model ANN SSM for training.

336 <u>validation and test stations, respectively (Table 3).</u>

In regard to the ANN_SSM_TEMP model that integrates the soil surface temperature, 49.4%, 55.56% and 59.35% of

the training, validation and test stations exhibited higher correlation values than those obtained with the ANN_SSM

339 model, respectively. RMSE decreased with ANN_SSM_TEMP compared to model ANN_SSM for 25.3%, 38.89% and

340 <u>42.99% of the training, validation and test stations, respectively.</u>

In addition, model ANN_SSM_EXP-FILT-T5 that integrates the simplified infiltration based features yielded slightly
 better correlations, and 64.56%, 60.61% and 63.68% 62.62% of the training, validation and test stations attained better
 correlations than those obtained with model ANN_SSM, respectively. Besides, RMSE decreased for 36.71%, 42.42%
 and 50.25% of the training, validation and test stations with ANN_SSM_EXP-FILT-T5 compared to model
 ANN_SSM, respectively.

Regarding the evaporation efficiency, we considered different values of fitting parameter B (Eq. 4) such that B

- 347 remained within the [50,60] interval. This parameter can be fitted using different variables, such as the wind speed or
- 348 relative humidity. Comparisons based on the correlation values provided by the different models for each B value

- 349 indicated that the performance was insensitive to the B value. Thus, we fixed the B value to 60 W m-². Comparison of
- 350 models ANN_SSM and ANN_SSM_EVAP-EFF-B60 revealed that 54.55%, 52.94% and 52.33% of the training,
- 351 <u>validation and test stations attained higher correlation values with the latter model, respectively. RMSE was reduced for</u>
- 352 <u>28.57%, 41.18% and 48.19% of the training, validation and test stations with ANN_SSM_EVAP-EFF-B60 compared</u>
- 353 to model ANN SSM, respectively.
- 554 <u>Finally, we investigated the impact of the joint application of the NDVI, recursive exponential filter (T= 5 days)</u>
- and evaporation efficiency (B=60 W m⁻²) in the ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT-T5 model. The
- 356 surface soil temperature was not included, as its effect is included in the evaporation process. At 84.06%, 61.29% and
- 357 <u>62.07% of the training, validation and test stations, the correlation value obtained with this model was higher than that</u>
- obtained with the ANN_SSM model, respectively. In addition, RMSE was minimized for 62.32%, 54.84% and 54.02%
- 359 of the training, validation and test stations with ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT-T5 compared to
- 360 <u>model ANN_SSM</u>, respectively.
- <u>Considering model ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT-T5, only one training station had a decrease in</u>
 correlation by more than 0.1 namely station 'Lind#1' (network 'SCAN') compared to reference model ANN_SSM. All
 inputs were not available at the same dates which implied a significant reduction in data points (cf. appendix F). The
 decrease in correlation and increase in RMSE didn't exceed 0.1 and 0.01 m³/m³, respectively, for the rest of stations of
 lower performance metrics with the most complex ANN.
- Similarly for validation stations, only one station had a decrease in correlation above 0.1, namely station 'PineNut'
 (network 'SCAN'), with model ANN_SSM_NDVI EVAP-EFF-B60 EXP-FILT-T5. This decrease can be also
 explained because of data shortage (cf. appendix F). The decrease in correlation and increase in RMSE didn't exceed
 0.1 and 0.01 m³/m³, respectively, for the rest of stations of lower performance metrics with the most complex ANN.
- Regarding test stations, correlation decrease by more than 0.1 and RMSE increase by more than 0.01 m³/m³with model
 ANN SSM NDVI EVAP-EFF-B60 EXP-FILT-T5 compared to model ANN SSM was detected for only 2 stations.
 Both stations, namely station 'S-Coleambally' and 'Widgiewa' which belong to network 'OZNET', significantly lose in
 data volume when process-related variables are integrated in ANN and more precisely because of NDVI data
 availability (cf. appendix F). For the rest of test stations, correlation decreased and RMSE increased simultaneously by
 less than 0.1 and 0.01 m³/m³, respectively, whith model ANN SSM NDVI EVAP-EFF-B60 EXP-FILT-T5.
- Table 3. Proportion of the stations which performance enhances using the ANN models enriched with process-related features
 compared to model ANN SSM (*: % of stations at which the correlation improves over the model ANN SSM level; **: % of stations
 at which RMSE improves over the model ANN SSM level)

Model	Training stations		Validation stations		Test stations	
	% of stations	% of stations	% of stations	% of stations	% of stations	% of stations
	<u>(corr ↑)*</u>	(RMSE ↓)**	<u>(corr ↑)*</u>	(RMSE ↓)**	<u>(corr ↑)*</u>	<u>(RMSE ↓)**</u>
<u>ANN SSM NDVI</u>	<u>65.82</u>	<u>44.3</u>	<u>45.71</u>	<u>40.0</u>	55.22	<u>40.3</u>
ANN SSM TEMP	<u>49.4</u>	<u>25.3</u>	<u>55.56</u>	<u>38.89</u>	<u>59.35</u>	<u>42.99</u>

ANN SSM EXP-FILT-T5	<u>64.56</u>	<u>36.71</u>	<u>60.61</u>	42.42	<u>63.68</u>	<u>50.25</u>
ANN_SSM_EVAP-EFF-B60	<u>54.55</u>	<u>28.57</u>	<u>52.94</u>	<u>41.18</u>	<u>52.33</u>	<u>48.19</u>
ANN SSM NDVI EVAP- EFF-B60 EXP-FILT-T5	<u>84.06</u>	<u>62.32</u>	<u>61.29</u>	<u>54.84</u>	<u>62.07</u>	<u>54.02</u>

Table 4. Proportion of the stations which correlation decreases using the ANN models enriched with process-related features

281 <u>compared to model ANN_SSM (Δ_{corr}=corr_{ANN_SSM} – corr_{ANN_SSM x}, X denotes a or a combination of process-related variables)</u>

Model	<u>Training</u>	stations	Validation stations		<u>Test stations</u>	
	% of stations	% of stations	% of stations	% of stations	% of stations	<u>% of</u>
	<u>corr ↓ and</u>	<u>corr↓and</u>	<u>corr ↓ and</u>	<u>corr↓ and</u>	$\underline{\operatorname{corr}} \downarrow \operatorname{and}$	stations
	$\underline{0.05} < \!\! \Delta_{corr} \! < \!\! 0.1^*$	$\Delta_{\rm corr} > 0.1^*$	$\underline{0.05} < \Delta_{corr} \leq 0.1^*$	$\Delta corr > 0.1^*$	$\underline{0.05 \leq \Delta_{corr} \leq 0.1}$	<u>corr ↓ and</u>
					*	$\Delta_{\rm corr} \ge 0.1^*$
	2.0	0	2.96	0	0.05	5.07
<u>ANN SSM NDVI</u>	<u>3.8</u>	<u>U</u>	2.80	<u>U</u>	<u>9.95</u>	<u>5.97</u>
ANN SSM TEMP	<u>0</u>	<u>1.2</u>	<u>0</u>	<u>2.78</u>	<u>4.67</u>	<u>3.27</u>
ANN SSM EXP-FILT-T5	<u>6.33</u>	<u>1.27</u>	<u>3.03</u>	<u>9.09</u>	<u>6.97</u>	<u>3.48</u>
ANN SSM EVAP-EFF-B60	<u>10.39</u>	<u>1.3</u>	<u>0</u>	<u>2.94</u>	<u>6.74</u>	<u>5.7</u>
ANN SSM NDVI EVAP- EFF-B60 EXP-FILT-T5	<u>4.35</u>	<u>1.45</u>	<u>6.45</u>	3.23	<u>9.2</u>	<u>6.9</u>

382

383 Always in terms of the general performance of model ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT-T5, about 384 75% of the stations have an RMSE less than $0.05 \text{ m}^3/\text{m}^3$ and around half of the stations have an RMSE less than 0.04385 m³/m³. This accuracy is consistent, for instance, with the target value in SMAP (Entekhabi et al., 2010) and SMOS 386 (Kerr et al., 2010) missions which is equal to 0.04 m³/m³ and also to the average sensor accuracy adopted by Dorigo et 387 al. (2013) which is equal to 0.05 m³/m³. Overall, the most complex model ANN_SSM_NDVI_EVAP-EFF-B60_EXP-388 FILT-T5 can successfully characterize the soil moisture dynamics in the root zone since half of the stations have a 389 correlation value greater than 0.7. Pan et al. (2017) developed different ANN models to estimate RZSM at depth of 390 20cm and 50cm over the continental United States using surface information. They found that half of the stations have 391 RMSE less than 0.06 m³/m³ and more than 70% of stations have correlation above 0.7 when predicting RZSM at 20cm. 392 However, the developed ANN was less effective in RZSM prediction at 50cm which is also in accordance with 393 (Kornelsen and Coulibaly, 2014). In our study, the densest soil moisture network is 'SCAN', located in the USA. Soil 394 moisture was predicted at a depth of 50cm over this network. Around half of the stations have a correlation value of 395 above 0.6 and RMSE less than 0.04 m³/m³ after the integration of process-related inputs. Pan et al., (2017) suggests that 396 the use of only time-dependent variables may not be sufficient for the ANN models to accurately predict RZSM and 397 suggests adding soil texture data.

398	In regard to the ANN_SSM_1EMP model that integrates the soil surface temperature, 54.35% of the stations (except
399	the stations of the AMMA CATCH network, as no surface temperature data were available) exhibited higher
400	correlation values than those obtained with the ANN_SSM model. Additionally, 40.24% of the stations achieved a
401	correlation value higher than 0.7 with model ANN_SSM_TEMP versus 36.94% of the stations with model ANN_SSM.
402	In addition, model ANN_SSM_EXP-FILT-T5 that integrates the simplified infiltration based features yielded slightly
403	better correlations, and 62.62% of the total stations attained better correlations than those obtained with model
404	ANN_SSM. A total of 45.37% of the stations achieved a correlation value higher than 0.7 with model
405	ANN_SSM_EXP FILT T5, in contrast to 38.98% of the stations achieving a similar correlation value with model
406	ANN_SSM.
407	Regarding the evaporative efficiency, we considered different values of fitting parameter B ₃ (Eq. 4) such that B ₃
408	remained within the [50,60] interval. This parameter can be fitted different variables, such as the wind speed or relative
409	humidity. Comparisons based on the correlation values provided by the different models for each B3 value indicated
410	that the performance was insensitive to the B ₃ value. Thus, we fixed the B ₃ value to 60 W m ⁻² . Comparison of models
411	ANN_SSM and ANN_SSM_EVAP EFF B60 revealed that 57.89% of the stations attained higher correlation values
412	with the latter model. A total of 41.12% of the stations exhibited a correlation value higher than 0.7 with model
413	ANN_SSM_EVAP-B60 versus 38.48% of the stations with model ANN_SSM.
414	Finally, we investigated the impact of the joint application of the NDVI, recursive exponential SWI (T= 5 days) and
415	evaporative efficiency (B ₃ =60 W m ⁻²) in the ANN_SSM_NDVI_EVAP EFF B60_EXP FILT_T5 model. The surface
416	soil temperature was not included, as its effect is included in the evaporation process. At 64.6% of the stations, the
417	correlation value obtained with this model was higher than that obtained with the ANN_SSM model. In addition, 51.1%

- 418 of the stations achieved a correlation value higher than 0.7 with model ANN_SSM_NDVI_EVAP EFF B60_EXP
- 419 FILT_T5, in contrast to 39.42% of the stations with model ANN_SSM.
- Table 3. Proportion of the stations exhibiting performance enhancement using the ANN models with the process based features
 model compared to model ANN_SSM.

improves over the model ANN_SSM level	the model ANN_SSM level
55.56	48.25
54.35	<u>46.25</u>
<u>62.62</u>	51.44
57.89	48.68
64.6	57.3
-	55.56 54.35 62.62 57.89 64.6

hoo

423 **3.2-3** Robustness of the approach

- 424 To <u>further</u> assess the robustness of our approach, which involves RZSM prediction using the <u>various different</u> ANN
- 425 models with different features, we predicted RZSM at 40 cm at sites not previously considered in previous parts of the
- 426 study. The selected stations are located in: the Kairouan Plain, a semiarid region in central Tunisia, Landriano site
- 427 located in the North of Italy, and the Berambadi watershed located in Gundalpet taluk, South-India. In the case of the
- 428 Kairouan Tunisia, model ANN_SSM yielded moderate- to low-precision predictions, as highlighted by the performance
- 429 metrics listed in Table 4<u>5</u>. The time series (cf. appendix G)-indicated show that the RZSM predictions followed the
- 430 SSM seasonality, which was reflected by the false peaks generated in the RZSM predictions whenever a sharp increase
- 431 or decrease occurred in the SSM values. This observation was already demonstrated by also found in (Souissi et al.,
- 432 2020). Actually, the Kairouan Plain is characterized by a semiarid environment where rainfall events infrequently occur
- and the level of evaporation is high. <u>The reference model ANN_SSM shows its limitations to accurately predict RZSM</u>
- 434 <u>in areas with no alternate wet and dry cycles.</u>
- 435 <u>However, the consideration of additional features, namely, the NDVI, evaporation efficiency and SWI in the ANN</u>
- 436 <u>models resulted in a good agreement between the in situ and predicted RZSM values (Fig. 4). The correlation values</u>
- 437 were improved by 60.04%, 169.5%, 112.02%, 80.23% and 53.7% at stations Barrouta-160, Hmidate 163,
- 438 Barrage 162, Bouhajla 164 and P12, respectively, with the ANN SSM NDVI EVAP-EFF-B60 EXP-FILT-T5 model
- 439 over ANN SSM model values. Similarly, RMSE values were reduced (Table 5). As shown in figure 4, the most
- 440 <u>complex ANN model is able to capture the variations of RZSM. This finding highlights the added value of our hybrid</u>
- 441 <u>approach based on an association of a machine learning method with process-related variables. Instead of injecting</u>
- 442 <u>uncertain information in physical models, such as soil properties, we used a nonparametric method related to physical</u>
- 443 processes without using forcing data that may be subject to errors and potentially lead to inaccurate tracking of the
- 444 <u>long-term evolution of soil moisture.</u>





445 Date 446 **Figure 4.** In situ SSM, in situ RZSM, and predicted RZSM series at the stations in the Kairouan Plain (Tunisia) with model 447 ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT__T5 (cf. appendix G for larger figure format).

448 However, the consideration of additional features, namely, the NDVI, evaporative efficiency and recursive exponential

- 449 filter SWI, in the ANN models resulted in a good agreement between the in situ and predicted RZSM values (Fig. 4).
- 450 The correlation values were improved by 60.04%, 169.5%, 112.02%, 80.23% and 53.7% at stations Barrouta 160,
- 451 Hmidate_163, Barrage_162, Bouhajla_164 and P12, respectively, with the ANN_SSM_NDVI_EVAP EFF B60_EXP
- 452 FILT_T5 model over ANN_SSM model values. Similarly, RMSE values were reduced (Table 4).
- 453 <u>A second comparison can be conducted between the quality of fit of these independent datasets and training datasets.</u>
- 454 <u>Actually, the climate class of the Tunisian stations is 'Bsh' (cf. appendix A). At the training stage, no station falls into</u>
- 455 the climate class 'Bsh' (cf. appendix A). However, some training stations fall under a similar climate class which is
- 456 <u>'Bsk' (cf. appendix B). Table 5 presents correlation and RMSE values for these training stations and Tunisian sites with</u>
- 457 both models ANN SSM and ANN SSM NDVI EVAP-EFF-B60 EXP-FILT-T5. For all training stations,
- 458 <u>performance metrics are slightly enhanced with the most complex ANN model compared to reference model</u>
- 459 <u>ANN SSM, except for stations GrouseCreek, Harmsway and Lind#1 which performance decreases. Overall, the range</u>
- 460 of correlation values is similar for training and external validation stations with model ANN SSM NDVI EVAP-EFF-
- 461 <u>B60 EXP-FILT-T5 and RMSE is well reduced for Tunisian stations compared to training stations. Given the results on</u>
- 462 <u>unseen datasets, namely on Tunisia, the performance of the most complex ANN model is good as it is able to generalize</u>
- 463 <u>the patterns present in the training dataset.</u>
- 464 <u>Table 5. Performance metrics of models ANN_SSM and ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT-T5 at training stations of</u>
 465 <u>climate "Bsk" and Tunisian stations of climate "Bsh".</u>

Model	ANN SSM	ANN SSM NDVI EVAP-EFF-B60 EXP-FILT-T5				
		<u>Trainir</u>	ng stations (climate class 'Bsh'	<u>)</u>		
Station	Correlation	<u>RMSE</u>	Correlation	RMSE		
<u>Banandra</u> (OZNET)	0.701	0.05	<u>0.764</u>	<u>0.046</u>		
DRY-LAKE (OZNET)	<u>0.674</u>	<u>0.031</u>	<u>0.692</u>	<u>0.03</u>		
CPER (SCAN)	<u>0.691</u>	0.032	0.695	0.032		
<u>EPHRAIM</u> (SCAN)	<u>0.758</u>	<u>0.051</u>	<u>0.791</u>	<u>0.046</u>		
<u>GrouseGreek</u> (SCAN)	<u>0.818</u>	<u>0.033</u>	<u>0.802</u>	0.035		
<u>HarmsWay</u> (SCAN)	<u>0.705</u>	<u>0.034</u>	<u>0.622</u>	<u>0.038</u>		
Lind#1 (SCAN)	0.605	0.055	<u>0.483</u>	0.022		
		Ext	ternal test stations (Tunisia)			
Station	Correlation	<u>RMSE</u>	Correlation	RMSE		
Barrouta 160	<u>0.463</u>	0.021	<u>0.714</u>	<u>0.016</u>		
Hmidate 163	<u>0.318</u>	<u>0.019</u>	<u>0.834</u>	0.011		
Barrage_162	<u>0.416</u>	0.035	<u>0.864</u>	<u>0.019</u>		
Bouhajla_164	0.435	0.016	0.733	<u>0.01</u>		
P12	<u>0.581</u>	0.047	<u>0.861</u>	<u>0.029</u>		

467 At the South-Indian stations, the ANN_SSM model yielded a good agreement even without the integration of process-468 based-related features (Table 6). The NDVI added little to nonsignificant improvement at station Bheemanbidu. The 469 same observation was made at the Italian site. The application of multiple features performed the best under arid 470 conditions, e.g., in Tunisia. In the tropical and temperate climate regions, this was not the case. The presence of clouds 471 in the MODIS NDVI and potential evapotranspiration products could explain this observation at sites of South-India 472 and North-Italy. In South-India, for instance, the maximum variability in soil moisture occurred during the monsoon 473 season, which is characterized by a large amount of clouds. Moreover, the coarse resolution of MODIS NDVI product 474 makes it sometimes not adapted to the considered site. (Chen et al., 2016) investigated the impact of sample impurity 475 and landscape heterogeneity on crop classification using coarse spatial resolution MODIS imagery. They showed that 476 the sample impurity such as mixed crop types in a specific sample, compositional landscape heterogeneity that is the 477 richness and evenness of land cover types in a landscape, and configurational heterogeneity that is the complexity of 478 spatial structure of land cover types in a specific landscape are sources of uncertainty affecting crop area mapping when using coarse spatial resolution imagery. High resolution NDVI from sensors like Sentinel-2 could have been used in
 this exercise to mitigate the spatial resolution issue, however, MODIS data were privileged in order to provide NDVI
 and PET from the same sensor.

Table 46. Performance metrics of models ANN_SSM, ANN_SSM_NDVI and ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT__T5
 at the sites in <u>central Tunisia</u>, <u>South-India and northern-Northern</u> Italy-and South-India.

Model	ANN_SSM	ANN_SSM_NDVI			ANN_SSM_NDVI_EVAP-	
					EFF_B60_EXP-FILTT5	
			INDIA			
Station	Correlation	RMSE	Correlation	RMSE	Correlation	RMSE
Madyanahundi	0.813	0.04	0.78	0.042	0.744	0.044
Bheemanbidu	0.76	0.046	0.784	0.044	0.763	0.046
Beechanalli2	0.825	0.038	0.787	0.04	0.743	0.044
Beechanalli1	0.713	0.024	0.713	0.024	0.633	0.025
			Italy			
Station	Correlation	RMSE	Correlation	RMSE	Correlation	RMSE
Landriano	0.861	0.038	0.827	0.041	0.841	0.038

484

485 **<u>34</u> Discussion**

486 Climate analysis of the results yielded by the different models indicated that among all models, the climate class with 487 the highest mean correlation change rate (Fig. 5) was class BWk (cf. appendix A), which regroups desert areas where 488 the link between SSM and RZSM is weak due to high evaporative rates. Class Dfa (cf. appendix A), which includes 489 areas experiencing harsh and cold winters, also yielded a high mean correlation change rate (>100%). Similarly, at 490 stations of this climate type, the link between the surface and root zone is poor. In regard to class Cfa (cf. appendix A), 491 in which 88.6% more than 80% of the total stations belongs to SCAN network, the high mean correlation change rate 492 could be explained by the surface-subsurface decoupling phenomena detected within this network, as previously 493 reported in (Souissi et al., 2020). The model with the largest number of stations with improved predictions over the 494 ANN_SSM model predictions was ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT_T5. Actually, the coupled use of 495 process-based-related features in the ANN models exerted a greater impact on the prediction accuracy than that exerted 496 by the one-at-a-time application of these features. In model ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT__T5, the 497 three process-based features jointly employed seemed to counterbalance the weight of these three SSM features. In this 498 model, the process-based-related features were equally represented versus the SSM information depicted by these three 499 features. The redundancy of the considered SSM information could explain the limited impact of the one-at-a-time 500 addition of process-based-related features the joint addition of the three process-based-related features.

In addition, (Karthikeyan and Mishra, 2021) Karthikeyan and Mishra (2021) demonstrated that at root depths beyond 20 cm, the importance of SSM was notably lower than that at the 20-cm depth, signifying decorrelation between surface and deeper SM values, which is in accordance with the findings in (Souissi et al., 2020), and it was further revealed that vegetation exhibits a higher importance than that of meteorological predictors LST and precipitation. (Kornelsen and Coulibaly, 2014) Kornelsen and Coulibaly (2014) indicated that evapotranspiration is the most important meteorological input for the prediction of soil moisture in the root zone with the MLP, which reflects the importance of the water vapor flux in soil moisture state determination.

508



509

*Mean correlation change rate per climate class

510 **Figure 5.** Climate classification of the stations performing better with models (a) ANN_SSM_NDVI (b) ANN_SSM_EXP-T5 (c) 511 ANN_SSM_EVAP-60 (d) ANN_SSM_TEMP and (e) ANN_SSM_NDVI_EVAP-EFF-B60_EXP-FILT__T5 compared to model 512 ANN_SSM_(Dark green corresponds to stations which correlation improved with complexified models, light green corresponds to 513 total stations, rate in blue correspond to mean correlation change rate per climate class).-

514
$$* corr_change_rate = mean(\frac{corr_{ANN_SSM_X} - corr_{ANN_SSM}}{corr_{ANN_SSM}} * 100) (65)$$

515 where X denotes a process-based-related variable (X \in ['NDVI', 'EXP-FILT-T5', 'EVAP-<u>EFF-B60'</u>, 'TEMP']

516 The <u>world</u> map illustrated in Fig. 6, shows the best-performing ANN models based on the mean correlation change rate 517 (Eq. 65). We assumed that the results in a given area of a specific climate class could be extended to other areas of the 518 same climate class even if we did not consider the data for these areas. The climate classes without at least one station

519 were marked in black and labeled with 'NO DATA'.



- Figure 6. Best World map of best-performing ANN models per climate class based on the mean correlation change rate; colors
 correspond to climate classes (cf. appendix A), hatches correspond to the most contributive input to the predictions namely: EVAP
 (evaporation efficiency), EXP (exponential filter SWI), NDVI, TEMP (surface soil temperature).
- In arid areas such as the eastern and western sides of the USA with high evaporation rates, ANN_SSM_EVAP-EFF-60 was the best performing model. Similarly, in bare areas of Africa, the Middle East and Australia where the Bwh climate class prevailed (arid desert hot climate; cf. appendix A), the evaporative evaporation efficiency was the best informative variable.
- In the internal part of continental Europe and near the Mediterranean Basin, the NDVI was the most relevant indicator for RZSM estimation, where agricultural fields dominated. Similarly, the Great Plains region in the USA was deeply affected by the NDVI, as this region is a cultivated area. The same result could be obtained for regions belonging to climate class Bsh (arid steppe hot; cf. appendix A) and mainly covered by grassland and shrubland areas according to ESA CCI land cover maps.
- 533 In Nordic areas characterized by the ET climate class, the soil temperature was the most important root zone soil 534 moisture indicator mainly because of the freeze-thaw events encountered in these regions. In tropical savannah wet 535 areas (class Aw; cf. appendix A), the ANN_SSM_TEMP model was the best-performing model.
- This classification definitely suffered limitations mainly provoked by the generalization of the climatic analysis results to areas not considered in this study. For instance, in regions of climate class Dfc (cold dry without a dry season, cold summer climate: cf. appendix A), we expect the temperature to serve as the most relevant indicator instead of the evaporative-evaporation efficiency.

540 45 Conclusion

In this study, we developed several ANN models to estimate RZSM based either on solely in situ SSM information or on a group of process-<u>based-related</u> features, in addition to SSM, namely, the soil water index computed with a recursive exponential filter, <u>evaporative-evaporation</u> efficiency, NDVI and surface soil temperature. Different regions across the globe with distinct land cover and climate patterns were considered. The main conclusion of this study was that the consideration of more features in addition to SSM information could enhance the accuracy of RZSM predictions mainly in regions where the link between SSM and RZSM is weak.

In arid areas with high evaporation rates, the most informative feature was the <u>evaporative evaporation</u> efficiency. In regions with agricultural fields, the NDVI was, for example, the most relevant indicator to predict RZSM. Overall, the best performing model included the surface soil moisture, NDVI, <u>recursive exponential filterSWI</u> and <u>evaporative</u> <u>evaporation</u> efficiency as features. Approximately 61.68% of the tested stations experienced correlation enhancement due to the joint consideration of process based features over RZSM model predictions based on only surface soil <u>moisture information</u>.

553 The robustness of the approach was further assessed through additional tests considering external sites in central 554 Tunisia, India and Italy. Similarly, the process-based-related features exerted a positive impact on the prediction 555 accuracy when combined with surface soil moisture in the case of Tunisia. The mean correlation across the five 556 Tunisian stations sharply increased from 0.44 when only SSM was considered to 0.8 when all process-based-related 557 features were combined with SSM. In India and Italy, the correlations were already high with the reference model 558 ANN_SSM, and the addition of process based features, namely, NDVI, did not improve the performance potentially 559 because of the cloudy conditions in India and noisy MODIS products.. The change in correlation after the addition of 560 process-related features, namely NDVI, is about -0.04 which is nonsignificant, and is potentially because of the cloudy 561 conditions in India and noisy MODIS products. Also the crop heterogeneity and sample impurity makes MODIS NDVI 562 products not adapted to all sites.

<u>As a research perspective, datasets can be separated in clusters corresponding to major climate classes and/or soil types.</u>
 <u>More analysis can be conducted in this direction to eventually make connections between the different inputs and</u>
 <u>climate/soil configurations.</u>

Future work will examine the ability of the developed model to estimate RZSM across larger areas based on remote sensing global soil moisture products. The use of remote sensing derived soil moisture products may yield lower correlations with the reference model ANN_SSM which potentially implies further improvement when process-based related features are added.

570 Acknowledgments

- 571 The PhD thesis of R. Souissi was financed by the ERANET RET-SIF project, and complementary financing was
- 572 provided by the PRIMA Programme SMARTIES project. The authors thank the International Soil Moisture Network
- 573 (ISMN) and supporting networks for providing the soil moisture data.

574 References

- Abrahart, R. J. and See, L. M.: Neural network modelling of non-linear hydrological relationships, 11, 1563–1579,
 https://doi.org/10.5194/hess-11-1563-2007, 2007.
- 577 Albergel, C., Rüdiger, C., Pellarin, T., Calvet, J.-C., Fritz, N., Froissard, F., Suquia, D., Petitpa, A., Piguet, B., and
- 578 Martin, E.: From near-surface to root-zone soil moisture using an exponential filter: an assessment of the method based
- 579 on in-situ observations and model simulations, 12, 1323–1337, https://doi.org/https://doi.org/10.5194/hess-12-1323-
- 580 2008, 2008.
- 581 ASCE Task Committee on Application of Artificial Neural Networks in Hydrology, Artificial Neural Networks in
- 582 Hydrology. II: Hydrologic Applications, 5, 124–137, https://doi.org/10.1061/(ASCE)1084-0699(2000)5:2(124), 2000.
- 583 Battude, M., Al Bitar, A., Brut, A., Tallec, T., Huc, M., Cros, J., Weber, J.-J., Lhuissier, L., Simonneaux, V., and
- 584 Demarez, V.: Modeling water needs and total irrigation depths of maize crop in the south west of France using high
- 585 spatial and temporal resolution satellite imagery, Agricultural Water Management, 189, 123–136,
- 586 https://doi.org/10.1016/j.agwat.2017.04.018, 2017.
- 587 Best, M. J., Pryor, M., Clark, D. B., Rooney, G. G., Essery, R. L. H., Ménard, C. B., Edwards, J. M., Hendry, M. A.,
- Porson, A., Gedney, N., Mercado, L. M., Sitch, S., Blyth, E., Boucher, O., Cox, P. M., Grimmond, C. S. B., and
 Harding, R. J.: The Joint UK Land Environment Simulator (JULES), model description Part 1: Energy and water
 fluxes, 4, 677–699, https://doi.org/10.5194/gmd-4-677-2011, 2011.
- 591 <u>Calvet, J.-C. and Noilhan, J.: From Near-Surface to Root-Zone Soil Moisture Using Year-Round Data, 1, 393–411,</u>
 592 https://doi.org/10.1175/1525-7541(2000)001<0393:FNSTRZ>2.0.CO;2, 2000.
- 593

Carranza, C., Nolet, C., Pezij, M., and van der Ploeg, M.: Root zone soil moisture estimation with Random Forest,
Journal of Hydrology, 593, 125840, https://doi.org/10.1016/j.jhydrol.2020.125840, 2021.

- 596 Chen, Y., Song, X., Wang, S., Huang, J., and Mansaray, L. R.: Impacts of spatial heterogeneity on crop area mapping in
 597 Canada using MODIS data, ISPRS Journal of Photogrammetry and Remote Sensing, 119, 451–461,
 598 https://doi.org/10.1016/j.isprsjprs.2016.07.007, 2016.
- 599

Didan, K., MOD13Q1 MODIS/Terra Vegetation Indices 16 Day L3 Global 250m SIN Grid V006 [Data set], NASA
 EOSDIS Land Processes DAAC, 2015. Available: https://doi.org/10.5067/MODIS/MOD13Q1.006, last access: 2
 december 2021.

Dorigo, W. A., Wagner, W., Hohensinn, R., Hahn, S., Paulik, C., Drusch, M., Mecklenburg, S., van Oevelen, P.,
Robock, A., and Jackson, T.: The International Soil Moisture Network: a data hosting facility for global in situ soil
moisture measurements, Hydrol. Earth Syst. Sci. Discuss., 8, 1609–1663, https://doi.org/10.5194/hessd-8-1609-2011,
2011.

- 607 Dorigo, W. A., Xaver, A., Vreugdenhil, M., Gruber, A., Hegyiová, A., Sanchis-Dufau, A. D., Zamojski, D., Cordes, C.,
- 608 Wagner, W., and Drusch, M.: Global Automated Quality Control of In Situ Soil Moisture Data from the International
- 609 Soil Moisture Network, Vadose Zone Journal, 12, vzj2012.0097, https://doi.org/10.2136/vzj2012.0097, 2013.
- 610
- 611 Entekhabi, D., Nakamura, H., and Njoku, E. G.: Retrieval of soil moisture profile by combined remote sensing and
- modeling, in: Retrieval of soil moisture profile by combined remote sensing and modeling, De Gruyter, 485–498, 2020.
- 613 Entekhabi, D., Njoku, E. G., O'Neill, P. E., Kellogg, K. H., Crow, W. T., Edelstein, W. N., Entin, J. K., Goodman, S.
- 614 D., Jackson, T. J., Johnson, J., Kimball, J., Piepmeier, J. R., Koster, R. D., Martin, N., McDonald, K. C., Moghaddam,
- 615 M., Moran, S., Reichle, R., Shi, J. C., Spencer, M. W., Thurman, S. W., Tsang, L., and Van Zyl, J.: The Soil Moisture
- 616 Active Passive (SMAP) Mission, Proc. IEEE, 98, 704–716, https://doi.org/10.1109/JPROC.2010.2043918, 2010.
- 617 Fieuzal, R., Baup, F., and Marais-Sicre, C.: Monitoring Wheat and Rapeseed by Using Synchronous Optical and Radar
- 618 Satellite Data—From Temporal Signatures to Crop Parameters Estimation, ARS, 02, 162–180,
- 619 https://doi.org/10.4236/ars.2013.22020, 2013.
- 620 Grillakis, M. G., Koutroulis, A. G., Alexakis, D. D., Polykretis, C., and Daliakopoulos, I. N.: Regionalizing Root-Zone
- 621 Soil Moisture Estimates From ESA CCI Soil Water Index Using Machine Learning and Information on Soil,
- 622 Vegetation, and Climate, 57, e2020WR029249, https://doi.org/10.1029/2020WR029249, 2021.
- Hajj, M., Baghdadi, N., Belaud, G., Zribi, M., Cheviron, B., Courault, D., Hagolle, O., and Charron, F.: Irrigated
- 624 Grassland Monitoring Using a Time Series of TerraSAR-X and COSMO-SkyMed X-Band SAR Data, Remote Sensing,
- 625 6, 10002–10032, https://doi.org/10.3390/rs61010002, 2014.
- 626 Han, H., Choi, C., Kim, J., Morrison, R. R., Jung, J., and Kim, H. S.: Multiple-Depth Soil Moisture Estimates Using
- 627 Artificial Neural Network and Long Short-Term Memory Models, Water, 13, 2584,
 - 628 https://doi.org/10.3390/w13182584, 2021.
 - Hassan-Esfahani, L., Torres-Rua, A., Jensen, A., and Mckee, M.: Spatial Root Zone Soil Water Content Estimation in
 - Agricultural Lands Using Bayesian-Based Artificial Neural Networks and High- Resolution Visual, NIR, and Thermal
 Imagery, 66, 273–288, https://doi.org/10.1002/ird.2098, 2017.
 - Huete, A., Didan, K., Leeuwen, W., Jacobson, A., Solanos, R., and Laing, T.: MODIS VEGETATION INDEX (MOD
 ALGORITHM THEORETICAL BASIS DOCUMENT Version 3. 1 Principal Investigators, 1999.
 - G34 Jacquemin, B. and Noilhan, J.: Sensitivity study and validation of a land surface parameterization using the HAPEX-
 - 635 MOBILHY data set, Boundary-Layer Meteorol, 52, 93–134, https://doi.org/10.1007/BF00123180, 1990.
 - 636 Karthikeyan, L. and Mishra, A. K.: Multi-layer high-resolution soil moisture estimation using machine learning over
 - the United States, Remote Sensing of Environment, 266, 112706, https://doi.org/10.1016/j.rse.2021.112706, 2021.
 - 638 Kerr, Y. H., Waldteufel, P., Wigneron, J.-P., Delwart, S., Cabot, F., Boutin, J., Escorihuela, M.-J., Font, J., Reul, N.,
 - 639 Gruhier, C., Juglea, S. E., Drinkwater, M. R., Hahne, A., Martín-Neira, M., and Mecklenburg, S.: The SMOS Mission:
 - 640 New Tool for Monitoring Key Elements of the Global Water Cycle, 98, 666–687,
 - 641 https://doi.org/10.1109/JPROC.2010.2043032, 2010.

- 642 Kolassa, J., Reichle, R. H., Liu, Q., Alemohammad, S. H., Gentine, P., Aida, K., Asanuma, J., Bircher, S., Caldwell, T.,
- 643 Colliander, A., Cosh, M., Holifield Collins, C., Jackson, T. J., Martínez-Fernández, J., McNairn, H., Pacheco, A.,
- 644 Thibeault, M., and Walker, J. P.: Estimating surface soil moisture from SMAP observations using a Neural Network
- 645 technique, Remote Sensing of Environment, 204, 43–59, https://doi.org/10.1016/j.rse.2017.10.045, 2018.
- Kornelsen, K. C. and Coulibaly, P.: Root-zone soil moisture estimation using data-driven methods, Water Resour. Res.,
 50, 2946–2962, https://doi.org/10.1002/2013WR014127, 2014.
- Koster, R. D., Dirmeyer, P. A., Guo, Z., Bonan, G., Chan, E., Cox, P., Gordon, C. T., Kanae, S., Kowalczyk, E.,
- 649 Lawrence, D., Liu, P., Lu, C.-H., Malyshev, S., McAvaney, B., Mitchell, K., Mocko, D., Oki, T., Oleson, K., Pitman,
- A., Sud, Y. C., Taylor, C. M., Verseghy, D., Vasic, R., Xue, Y., and Yamada, T.: Regions of Strong Coupling Between
- 651 Soil Moisture and Precipitation, 305, 1138–1140, https://doi.org/10.1126/science.1100217, 2004.
- Lee, T. J. and Pielke, R. A.: Estimating the Soil Surface Specific Humidity, 31, 480–484, https://doi.org/10.1175/15200450(1992)031<0480:ETSSSH>2.0.CO;2, 1992.
- Liu, Y., Chen, D., Mouatadid, S., Lu, X., Chen, M., Cheng, Y., Xie, Z., Jia, B., Wu, H., and Gentine, P.: Development
- of a Daily Multilayer Cropland Soil Moisture Dataset for China Using Machine Learning and Application to Cropping
- 656 Patterns, 22, 445–461, https://doi.org/10.1175/JHM-D-19-0301.1, 2021.
- 657 Martínez-Espinosa, C., Sauvage, S., Al Bitar, A., Green, P. A., Vörösmarty, C. J., and Sánchez-Pérez, J. M.:
- Denitrification in wetlands: A review towards a quantification at global scale, Science of The Total Environment, 754,
- 659 142398, https://doi.org/10.1016/j.scitotenv.2020.142398, 2021.
- 660 Masseroni, D., Corbari, C., and Mancini, M.: Validation of theoretical footprint models using experimental
- measurements of turbulent fluxes over maize fields in Po Valley, Environ Earth Sci, 72, 1213–1225,
 https://doi.org/10.1007/s12665-013-3040-5, 2014.
- 663 Masson, V., Le Moigne, P., Martin, E., Faroux, S., Alias, A., Alkama, R., Belamari, S., Barbu, A., Boone, A., Bouyssel,
- 664 F., Brousseau, P., Brun, E., Calvet, J.-C., Carrer, D., Decharme, B., Delire, C., Donier, S., Essaouini, K., Gibelin, A.-L.,
- 665 Giordani, H., Habets, F., Jidane, M., Kerdraon, G., Kourzeneva, E., Lafaysse, M., Lafont, S., Lebeaupin Brossier, C.,
- 666 Lemonsu, A., Mahfouf, J.-F., Marguinaud, P., Mokhtari, M., Morin, S., Pigeon, G., Salgado, R., Seity, Y., Taillefer, F.,
- Tanguy, G., Tulet, P., Vincendon, B., Vionnet, V., and Voldoire, A.: The SURFEXv7.2 land and ocean surface
- platform for coupled or offline simulation of earth surface variables and fluxes, Geosci. Model Dev., 6, 929–960,
 https://doi.org/10.5194/gmd-6-929-2013, 2013.
- 670 Merlin, O., Bitar, A. A., Rivalland, V., Béziat, P., Ceschia, E., and Dedieu, G.: An Analytical Model of Evaporation
- 671 Efficiency for Unsaturated Soil Surfaces with an Arbitrary Thickness, 50, 457–471,
- 672 https://doi.org/10.1175/2010JAMC2418.1, 20112010.
- Noilhan, J. and Mahfouf, J.-F.: The ISBA land surface parameterisation scheme, Global and Planetary Change, 13,
 145–159, https://doi.org/10.1016/0921-8181(95)00043-7, 1996.
- Noilhan, J. and Planton, S.: A Simple Parameterization of Land Surface Processes for Meteorological Models, 117,
- 676 536–549, https://doi.org/10.1175/1520-0493(1989)117<0536:ASPOLS>2.0.CO;2, 1989.

- 677 Oleson, W., Lawrence, M., Bonan, B., Flanner, G., Kluzek, E., Lawrence, J., Levis, S., Swenson, C., Thornton, E., Dai,
- 678 A., Decker, M., Dickinson, R., Feddema, J., Heald, L., Hoffman, F., Lamarque, J.-F., Mahowald, N., Niu, G.-Y., Qian,
- T., Randerson, J., Running, S., Sakaguchi, K., Slater, A., Stockli, R., Wang, A., Yang, Z.-L., Zeng, X., and Zeng, X.:
- 680 Technical Description of version 4.0 of the Community Land Model (CLM), https://doi.org/10.5065/D6FB50WZ,
- 681 2010.
- 682 Owe, M., de Jeu, R., and Holmes, T.: Multisensor historical climatology of satellite-derived global land surface
- 683 moisture, J. Geophys. Res., 113, F01002, https://doi.org/10.1029/2007JF000769, 2008.
- 684 Oyebode, O. and Stretch, D.: Neural network modeling of hydrological systems: A review of implementation
- 685 techniques, 32, e12189, https://doi.org/10.1111/nrm.12189, 2019.
- 686 Pan, X., Kornelsen, K. C., and Coulibaly, P.: Estimating Root Zone Soil Moisture at Continental Scale Using Neural
- 687 Networks, 53, 220–237, https://doi.org/10.1111/1752-1688.12491, 2017.
- 688 Paris Anguela, T., Zribi, M., Hasenauer, S., Habets, F., and Loumagne, C.: Analysis of surface and root-zone soil
- 689 moisture dynamics with ERS scatterometer and the hydrometeorological model SAFRAN-ISBA-MODCOU at Grand
- 690 Morin watershed (France), 12, 1415–1424, https://doi.org/10.5194/hess-12-1415-2008, 2008.
- 691 Paulik, C., Dorigo, W., Wagner, W., and Kidd, R.: Validation of the ASCAT Soil Water Index using in situ data from
- the International Soil Moisture Network, International Journal of Applied Earth Observation and Geoinformation, 30,
- 693 1–8, https://doi.org/10.1016/j.jag.2014.01.007, 2014.
- Raes, D., Steduto, P., Hsiao, T. C., and Fereres, E.: AquaCrop—The FAO Crop Model to Simulate Yield Response to
- Water: II. Main Algorithms and Software Description, 101, 438–447, https://doi.org/10.2134/agronj2008.0140s, 2009.
- Ramchoun, H., Amine, M., Idrissi, J., Ghanou, Y., and Ettaouil, M.: Multilayer Perceptron: Architecture Optimization
 and Training, IJIMAI, 4, 26, https://doi.org/10.9781/ijimai.2016.415, 2016.
- Running, S., Q. Mu, M. Zhao. MOD16A2 MODIS/Terra Net Evapotranspiration 8-Day L4 Global 500m SIN Grid
- V006. 2017, distributed by NASA EOSDIS Land Processes DAAC, https://doi.org/10.5067/MODIS/MOD16A2.006,
 last access: 2 december 2021.
- Sabater, J. M., Jarlan, L., Calvet, J.-C., Bouyssel, F., and De Rosnay, P.: From Near-Surface to Root-Zone Soil
- 702 Moisture Using Different Assimilation Techniques, J. Hydrometeor., 8, 194–206, https://doi.org/10.1175/JHM571.1,
- 703 2007.
- SIE, Available: https://osr-cesbio.ups-tlse.fr/, last access: 8 december 2021.
- 705 Souissi, R., Al Bitar, A., and Zribi, M.: Accuracy and Transferability of Artificial Neural Networks in Predicting in Situ
- Root-Zone Soil Moisture for Various Regions across the Globe, 12, 3109, https://doi.org/10.3390/w12113109, 2020.
- 707 Stroud, P. D.: A Recursive Exponential Filter For Time-Sensitive Data, 1999.
- 708 Tanty, R., Desmukh, T. S., and MANIT BHOPAL: Application of Artificial Neural Network in Hydrology- A Review,
- 709 IJERT, V4, IJERTV4IS060247, https://doi.org/10.17577/IJERTV4IS060247, 2015.

- 710 Wagner, W., Blöschl, G., Pampaloni, P., Calvet, J.-C., Bizzarri, B., Wigneron, J.-P., and Kerr, Y.: Operational
- readiness of microwave remote sensing of soil moisture for hydrologic applications, Hydrology Research, 38, 1–20,
- 712 https://doi.org/10.2166/nh.2007.029, 2007.
- 713 Wagner, W., Hahn, S., Kidd, R., Melzer, T., Bartalis, Z., Hasenauer, S., Figa-Saldaña, J., de Rosnay, P., Jann, A.,
- 714 Schneider, S., Komma, J., Kubu, G., Brugger, K., Aubrecht, C., Züger, J., Gangkofner, U., Kienberger, S., Brocca, L.,
- 715 Wang, Y., Blöschl, G., Eitzinger, J., and Steinnocher, K.: The ASCAT Soil Moisture Product: A Review of its
- 716 Specifications, Validation Results, and Emerging Applications, 5–33, https://doi.org/10.1127/0941-2948/2013/0399,
- 717 2013.
- 718 Wagner, W., Lemoine, G., and Rott, H.: A Method for Estimating Soil Moisture from ERS Scatterometer and Soil
- 719 Data, Remote Sensing of Environment, 70, 191–207, https://doi.org/10.1016/S0034-4257(99)00036-X, 1999.
- 720 Zribi, M., Chahbi, A., Shabou, M., Lili-Chabaane, Z., Duchemin, B., Baghdadi, N., Amri, R., and Chehbouni, A.: Soil
- surface moisture estimation over a semi-arid region using ENVISAT ASAR radar data for soil evaporation evaluation,
- 722 15, 345–358, https://doi.org/10.5194/hess-15-345-2011, 2011.
- 723 Zribi, M., Foucras, M., Baghdadi, N., Demarty, J., and Muddu, S.: A New Reflectivity Index for the Retrieval of
- Surface Soil Moisture From Radar Data, IEEE J. Sel. Top. Appl. Earth Observations Remote Sensing, 14, 818–826,
 https://doi.org/10.1109/JSTARS.2020.3033132, 2021.
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739

741 APPENDIX A 742 Climate classes (Köppen classification): 743 Af: Tropical Rainforest • 744 • Am: Tropical Monsoon 745 • As: Tropical Savanna Dry 746 Aw: Tropical Savanna Wet • 747 BWk: Arid Desert Cold • 748 BWh: Arid Desert Hot • 749 • BWn: Arid Desert with Frequent Fog 750 BSk: Arid Steppe Cold • 751 • BSh: Arid Steppe Hot 752 • BSn: Arid Steppe with Frequent Fog 753 Csa: Temperate Dry Hot Summer • 754 Csb: Temperate Dry Warm Summer • 755 Csc: Temperate Dry Cold Summer • Cwa: Temperate Dry Winter, Hot Summer 756 • 757 • Cwb: Temperate Dry Winter, Warm Summer 758 Cwc: Temperate Dry Winter, Cold Summer • 759 Cfa: Temperate without a Dry Season, Hot Summer • 760 Cfb: Temperate without a Dry Season, Warm Summer • 761 Cfc: Temperate without a Dry Season, Cold Summer • 762 Dsa: Cold Dry Summer, Hot Summer • 763 Dsb: Cold Dry Summer, Warm Summer • 764 Dsc: Cold Dry Summer, Cold Summer • 765 Dsd: Cold Dry Summer, Very Cold Winter • Dwa: Cold Dry Winter, Hot Summer 766 • 767 Dwb: Cold Dry Winter, Warm Summer • 768 Dwc: Cold Dry Winter, Cold Summer • 769 Dwd: Cold Dry Winter, Very Cold Winter • 770 Dfa: Cold Dry without a Dry Season, Hot Summer • 771 Dfb: Cold Dry without a Dry Season, Warm Summer • 772 Dfc: Cold Dry without a Dry Season, Cold Summer • 773 Dfd: Cold Dry without a Dry Season, Very Cold Winter • 774 • ET: Polar Tundra 775 EF: Polar Eternal Winter • 776 W: Water •

APPENDIX B



786	
787	<u>APPENDIX C</u>
788	Evaporation efficiency (section 2.2.2):
789	The standard equations to compute evaporation efficiency (β_3) in (Merlin et al., 2010) are as follows:
790 791	$\beta_{r} = \left[\frac{1}{r} - \frac{1}{r}\cos(\pi\theta_{r}/\theta_{r})\right]^{p} for \theta_{r} \leq \theta (C3)$
792	$\beta_3 = \frac{1}{2} \sum_{2} \cos(n\theta_L / \theta_{max}) + \frac{1}{2} \int \partial r \theta_L \le \theta_{max} + \frac{1}{2} \int \partial r \theta_L \le \theta_{max}$
793	where: $-\theta_L$ is the water content in the soil layer of thickness L.
794	- P is a parameter computed as follows:
795	$P = (\frac{1}{2} + A_3 \frac{L - L_1}{L_1}) \frac{L E_p}{B_3} (C4)$
796	$-\theta_{max}$ is the soil moisture at saturation.
797	$-LE_p$ is the potential evaporation.
798	- L_1 is the thinnest represented soil layer, and A_3 (unitless) and B_3 (W/m ²) are the two best-fit parameters a
799	priori depending on the soil texture and structure, respectively.
800	
001	
802	
803	
804	
805	
806	
807	
808	
809	
810	
811	
812	
813	







APPENDIX E



836 <u>Training stations:</u>











APPENDIX G



