Global prediction of soil saturated hydraulic conductivity using random forest in a Covariate-based Geo Transfer Functions (CoGTF) framework

Surya Gupta¹, Tomislav Hengl², Peter Lehmann³, Sara Bonetti¹, Andreas Papritz¹, and Dani Or¹

¹Soil and Terrestrial Environmental Physics, Department of Environmental Systems Science, ETH Zurich, Switzerland ²Envirometrix Ltd., Wageningen, the Netherlands ³Soil and Terrestrial Environmental Physics, Department of Environmental Systems Science, ETH Zurich, Switzerland

November 22, 2022

Abstract

The saturated hydraulic conductivity (Ksat) is a key soil hydraulic parameter for representing infiltration and drainage in Earth system and land surface models. For large scale applications, Ksat is often estimated from pedotransfer functions (PTFs) based on easy-to-measure soil properties like soil texture and bulk density. The reliance of PTFs on data from uniform arable lands and omission of soil structure limits the applicability of texture-based predictions of Ksat in vegetated lands. A method to harness technological advances in machine learning and availability of remotely sensed surrogate information to derive a new global Ksat map at 1 km resolution using terrain, climate, vegetation, and soil covariates is proposed. For model training and testing, global compilation of 6,814 georeferenced Ksat measurements from the literature across the globe were used. The accuracy assessment results based on model cross-validations with re-fitting show a concordance correlation coefficient of 0.79 and root mean square error of 0.72 (in log10Ksat given in cm/day). The generated maps of Ksat represent spatial patterns of the vegetation-induced soil structure formation and clay mineralogy, more distinctly than previous global maps of Ksat such as computed with Rosetta 3 pedotransfer function. The validation of the model indicates that Ksat could be more accurately modeled using covariate-based Geo Transfer Functions (CoGTFs) that harness spatially distributed surface and climate attributes, compared to pedotransfer functions that rely only on soil information.

Global prediction of soil saturated hydraulic conductivity using 1 random forest in a Covariate-based Geo Transfer Functions 2 (CoGTF) framework 3

Surya Gupta¹, Tomislav Hengl^{2,3}, Peter Lehmann¹, Sara Bonetti^{1,4}, Andreas Papritz¹, Dani 4 **Or**^{1,5}

5

6	¹ Soil and Terrestrial Environmental Physics, Department of Environmental Systems Science, ETH Zürich, Zürich,
7	Switzerland
8	² Envirometrix Ltd., Wageningen, the Netherlands
9	³ OpenGeoHub, Wageningen, the Netherlands
10	⁴ Bartlett School of Environment, Energy and Resources, University College London, London, UK
11	⁵ Division of Hydrologic Sciences, Desert Research Institute, Reno, NV, USA

12	Key Points:
13	Climate, vegetation and terrain affect spatial patterns of saturated hydraulic conductiv-
14	ity (Ksat)
15	• The effect of these covariates on Ksat is quantified using remote sensing data and machine
16	learning
17	• We introduce geotransfer functions to improve Ksat predictions based on pedotransfer func-
18	tions

Corresponding author: Surya Gupta, surya.gupta@usys.ethz.ch

19 Abstract

The saturated hydraulic conductivity (Ksat) is a key soil hydraulic parameter for representing in-20 filtration and drainage in Earth system and land surface models. For large scale applications, Ksat 21 is often estimated from pedotransfer functions (PTFs) based on easy-to-measure soil properties 22 like soil texture and bulk density. The reliance of PTFs on data from uniform arable lands and 23 omission of soil structure limits the applicability of texture-based predictions of Ksat in vege-24 tated lands. A method to harness technological advances in machine learning and availability of 25 remotely sensed surrogate information to derive a new global Ksat map at 1-km resolution us-26 ing terrain, climate, vegetation, and soil covariates is proposed. For model training and testing, 27 global compilation of 6,814 georeferenced Ksat measurements from the literature across the globe 28 were used. The accuracy assessment results based on model cross-validations with re-fitting show 29 a concordance correlation coefficient of 0.79 and root mean square error of 0.72 (in \log_{10} Ksat 30 given in cm/day). The generated maps of Ksat represent spatial patterns of the vegetation-induced 31 soil structure formation and clay mineralogy, more distinctly than previous global maps of Ksat 32 such as computed with Rosetta 3 pedotransfer function. The validation of the model indicates 33 that Ksat could be more accurately modeled using covariate-based Geo Transfer Functions (CoGTFs) 34 that harness spatially distributed surface and climate attributes, compared to pedotransfer func-35 tions that rely only on soil information. 36

37 Plain Language Summary

The soil saturated hydraulic conductivity Ksat defines how fast water can infiltrate and per-38 colate through the soil. To model water flow at large scale, accurate maps of Ksat are needed. Usu-30 ally, Ksat is not measured directly but deduced from well known basic soil properties (soil tex-40 ture, packing density). But these estimates neglect the influence of vegetation and climate on for-41 mation of soil structures that control Ksat. To improve predictions of Ksat, we use a new spa-42 tially referenced Ksat data collection and apply Machine Learning to find correlations between 43 Ksat and other properties (soil information, terrain, climate and vegetation). These correlations 44 are then implemented at global scale using maps of all relevant properties (so called 'covariates' 45 that were measured by remote sensing). We called this new approach to predictive soil mapping 46 the "Covariate-based Geotransfer functions" (CoGTF) to highlight the difference to other maps 47 that neglect spatial correlation with soil formatting properties and are based only on soil infor-48 mation (so called "pedotransfer functions" or PTFs). We show that the new maps based on CoGTF 49 perform better than approaches based on PTFs. 50

51 1 Introduction

The description of water, energy, and carbon fluxes between the land surface and the atmosphere relies heavily on the availability of soil hydraulic data (Gutmann & Small, 2007; Fashi et al., 2016; Montzka et al., 2017). A prominent soil hydraulic property is the soil saturated hydraulic conductivity (Ksat) that affects the partitioning of rainfall between runoff and infiltration

⁵⁶ (Zimmermann et al., 2013), and plays a critical role in a variety of hydrological and climatolog-

ical applications (Gutmann & Small, 2007; Or, 2019; Fatichi et al., 2020). At global scale, maps

-2-

- of soil hydraulic properties at ever increasing resolution are required for building Land Surface
- ⁵⁹ Models (LSMs) (Montzka et al., 2017).

For large scale applications (regional and global), soil hydraulic parameters are often es-60 timated from easy-to-measure soil properties (e.g., texture, organic content, bulk density) by means 61 of pedotransfer functions (PTFs) (Bouma, 1989; Santra & Das, 2008). PTFs are usually devel-62 oped for specific geographic regions thus only representing local conditions of soil forming pro-63 cesses (e.g. Tomasella & Hodnett, 1998; Wösten et al., 1999; Nemes et al., 2005; Saxton & Rawls, 64 2006; Jorda et al., 2015; Khlosi et al., 2016). This hinders their transferability across large ge-65 ographical regions (Vereecken et al., 2016). In addition, PTFs generally ignore soil structure and 66 pedogenic information and rely heavily on soil textural information (Fatichi et al., 2020), lim-67 iting their applicability in soils characterized by aggregation and formation of biopores. More-68 over, PTFs are generally defined as a function of clay content, without consideration of the ef-69 fect of different clay minerals on soil hydraulic properties (Hodnett & Tomasella, 2002). Dai et 70 al. (2019) have recently produced 1-km resolution global maps of soil hydraulic properties (and 71 thermal soil conductivity) using the median values of multiple PTFs to estimate Ksat. Likewise, 72 Y. Zhang and Schaap (2017) have developed a global map of van Genuchten parameters and Ksat 73 based on the Rosetta 3 PTF (an extension of Schaap et al., 2001), making use of three data sets 74 from North America and Europe (i.e., Rawls et al., 1982; Ahuja et al., 1989) and UNSODA (Un-75 saturated Soil Hydraulic Database) as described in Leij et al. (1996) and Nemes et al. (2001) and 76 employing Artificial Neural Network and bootstrap sampling. 77

Maps produced by Dai et al. (2019) and Y. Zhang and Schaap (2017) are limited by the small 78 number and unevenly distributed Ksat measurements (N = 1306) used for model training and 79 large spatial gaps i.e. missing training points in tropics. Moreover, the training points used to pro-80 duce estimates of Ksat were usually dominated by particular land use and land cover, mainly col-81 lected in arable land. Furthermore, only a limited set of basic soil variables (i.e., bulk density and 82 texture) was employed in the derivation of the Rosetta 3 map (Y. Zhang & Schaap, 2019), while 83 several studies have shown that also other soil properties such as organic carbon, soil depth and 84 pH may increase accuracy of PTFs (Wösten et al., 1999; Mayr & Jarvis, 1999; Tóth et al., 2015). 85 The availability of highly resolved remote sensing (RS) and landscape covariates offer new op-86 portunities for injecting new and local information into the modeling of Ksat. Examples of the 87 potential usefulness of such covariates are reported by Obi et al. (2014) that developed a PTF us-88 ing terrain attributes for many soil hydraulic properties; Sharma et al. (2006) combined PTFs with 89 vegetation and topography indices; Jana and Mohanty (2011) showed that the introduction of to-90 pographic attributes (i.e., Digital Elevation Model, DEM) and information on vegetation (i.e., Leaf 91 Area Index, LAI) along with *in situ* soil basic properties could improve predictions of soil hy-92 draulic properties. 93

Many of the recent PTFs use Machine Learning (ML) algorithms to quantify the relations between hydraulic properties and various covariates (Schaap et al., 2001; Jana & Mohanty, 2011; Araya & Ghezzehei, 2019). In this paper, we hypothesize that Ksat predictions could be improved using a combination of soil variables and remote sensing covariate layers integrated by using machine learning (ML) framework. We profit from the advancement in remote sensing techniques (providing spatial information on different ecological parameters with unprecedented resolution) to improve the predictions for soil hydraulic parameters and bridge the gap between site-specific

manuscript submitted to Journal of Advances in Modeling Earth Systems (JAMES)

soil properties and landscape variability. We merge concepts of predictive soil mapping with a 101 large data set of Ksat measurements and local information (soil, vegetation, climate) into covariate-102 based "Geo Transfer Functions" (CoGTFs) to generate global estimates of Ksat values (to high-103 light the impact of Geo-referencing soil properties and RS-covariates we use the term GTF and 104 not PTF). We compare mapping accuracy using global and local/regional assessment including 105 visual interpretation of produced spatial predictions. We show how this method (providing novel 106 covariate-based maps of Ksat) could be used to overcome some of the limitations of traditional 107 PTFs.

108

109

Our specific objectives are:

- 1. to improve accuracy and spatial detail of global Ksat maps by harnessing the state-of-the-110 art global remote sensing data products at 1 km spatial resolution, 111
- 2. to generate global maps of Ksat at different soil depths (0, 30, 60 and 100 cm), 112
- 3. to identify the key environmental variables explaining the spatial distribution of Ksat. 113

We first describe the model training for Ksat mapping using a random forest ML algorithm, 114 and then compare the results against maps generated with Rosetta 3 and the map shown in Dai 115 et al. (2019). Note that for a detailed comparison of global maps, we focus on Rosetta 3 because 116 the map in Dai et al. (2019) is heavily influenced by the application of a different soil textural 117 map (see Supplementary Information file). Then, we validated the CoGTF map, Rosetta 3 map 118 and the map of Dai et al. (2019) with independent dataset. We finally show the importance of us-119 ing RS covariates to capture spatial patterns and improve the accuracy of soil hydraulic proper-120 ties. 121

122

2 MATERIALS AND METHODS

123

2.1 Covariate-based Geo Transfer Functions (CoGTF) framework

We propose here an integrated Predictive Soil Modeling (PSM) framework where soil vari-124 ables are combined with RS-based covariates using random forest method (Figure 1). We refer 125 to this approach as the "Covariate-based Geo Transfer Functions" (CoGTF) framework and en-126 visage it as a combination of traditional PTF approach and purely data science approach where 127 RS-based covariates are used to map patterns in soil properties. The CoGTF framework follows 128 six principal steps: 129

- 1. Prepare georeferenced dataset of response variable (Ksat), 130
- 2. Overlay training points and covariates (including predictions of basic soil properties), and 131 produce a regression matrix, 132
- 3. Optimize the hyper-parameters in the random forest approach (mtry), 133
- 4. Fit the random forest model, 134
- 5. Evaluate the performance of the Ksat model, 135
- 6. Produce spatial predictions of Ksat. 136

A central hypothesis in this study is that spatial and climatic covariates could be harnessed 137 to improve the global mapping of Ksat (Jana & Mohanty, 2011). The basis for such hypothesis 138

- is the dominant role of climate, topography, and vegetation in soil formation and thus in shap-
- ing local hydraulic transport properties. For each location with Ksat measurement, the values of
- the remote sensing covariates were extracted together with modeled soil information from Open-
- LandMap.org. We implement the spatial predictions and the creation of Ksat maps in the R en-
- vironment (R Core Team, 2013) for statistical computing and provide code examples via the https://
- github.com/ETHZ-repositories/Ksat_mapping_2020/.



Figure 1. Computational workflow used to generate the soil Ksat map. See text for more details about the specific steps.

145After extracting all covariates, a regression matrix was formed, and the best hyperparam-146eter (mtry) was computed by five-fold cross-validation, using the R packages 'caret' version 6.0-14785 (Kuhn, 2012) and 'ranger' version 0.12.1 (Wright & Ziegler, 2015). Then, log-transformed

- (log_{10}) Ksat was modeled as a function of depth using random forest (RF) algorithm.
- 149 2.2 Training point data

Our first task was to enlarge the Ksat measurement database beyond the ≈1,300 values used
 to train Rosetta 3 by compiling available and georeferenced Ksat values from the literature. The

- Ksat values were log-transformed (log_{10} Ksat) and cm/day was selected as a standardized unit.
- A detailed description of the data collection and processing is provided in Gupta et al. (2020).
- We managed to compile a total of 13,267 samples coming from 1,910 sites across the globe. Most
- training data are from the USA, followed by Europe, Asia, South America, Africa, and Australia
- as shown in Figure 2. The collected Ksat database (SoilKsatDB) includes both field (N = 4,460) and lab (N = 8,807) measurements.

- To limit the over-representation of Florida (mainly arable land not representative of soils with natural vegetation), we randomly selected approximately only 1% of the 6,532 Florida samples, so that a total of 6,814 Ksat values were finally used for Ksat mapping. This resulted in geographical balance between other national data sets (the effect of this selection of Florida data
- is discussed in Supplementary information file).



Figure 2. Spatial distribution of measured Ksat values (6,814 samples in total) used to produce the global Ksat map. Colors refer to laboratory (red) and field (blue) measurements. The map is presented in the Goode equal-area homolosine projection. For more details and access to the Ksat data see Gupta et al. (2020).

2.3 Soil and environmental covariates

As environmental and soil covariate layers for Ksat modeling at global scale, we used global maps of soil properties (sand, clay, and bulk density) and other 24 RS-based covariates available from https://openlandmap.org/. These were selected to represent ecological conditions essential in soil-forming processes according to Jenny (1994). The covariates can be divided into five groups:

- Climate-based covariates, including mean annual precipitation, temperature, temperature seasonality, maximum temperature of warmest month, minimum temperature of coldest month, precipitation of wettest month, precipitation of driest month (Chelsa products, Karger et al., 2017), cloud fraction (Wilson & Jetz, 2016), diffuse irradiation, direct irradiation, annual land surface temperature, monthly precipitation and its standard deviation (Brocca et al., 2019).
- Digital terrain model (DTM)-based covariates (Yamazaki et al., 2017), including land scape metrics (such as slope, aspect, topographic wetness index) derived from SAGA GIS
 (Conrad et al., 2015) and landform classification and lithological maps.
- Surface reflectance-based covariates, including surface reflectance from Landsat and MODIS
 dataset for different wavelength bands (Hansen et al., 2013), snow probability (Buchhorn
 et al., 2017) and regularly flooded wetlands (Tootchi et al., 2019).

4. Vegetation-based covariate, represented by the annual fraction of absorbed photosynthet-181 ically active radiation (FAPAR), averaged over the 2014-2019 period. 182 5. Basic soil properties, comprising sand, clay content and bulk density for different soil depths 183 (matching the sampling depth of Ksat), which were obtained from OpenLandMap (Hengl 184 et al., 2017). Soil depth is used as a covariate to model the change of Ksat with depth (the 185 methodology to use depth as a covariate is described in Hengl & MacMillan, 2019).

A detailed list and description of all the covariates is provided in Table S1 in the Supple-187 mentary Information (SI). All covariate maps were resampled to the standard grid at a spatial res-188 olution of 1 km covering latitudes between -62.0 and 87.37. We did not map Antarctica as this 189 continent is dominantly covered with permanent ice and lacks training points. 190

191

186

2.4 Evaluating the performance of Ksat predictive models

The model-fitting results were evaluated using out of bag (OOB) error reported by the ranger 192 package by default. A bootstrap sampling is used to construct each tree in the random forest and 193 different bootstrap samples are used for each tree containing approximately 2/3 of the total ob-194 servations. The samples not used in the bootstrap samples are called out-of-bag (OOB) samples 195 (sub-dataset) (Peters et al., 2007; Rad et al., 2014). The relative importance of the covariates was 196 assessed by the increase in node purity. It is calculated using gini criterion from all the splits (in 197 our case 200 splits) in the forest based on a particular variable (Breiman, 2001; Rodrigues & de la 198 Riva, 2014). 199

The performance of the Ksat model was evaluated using 5-fold cross-validation. This means 200 that models were refitted 5 times using 80% of the data and the predictions for remaining 20% 201 estimated using these models were compared with observations. The process was repeated three 202 times to produce stable results. The final results are shown using hexbin plot with the LOWESS 203 (Locally Weighted Scatterplot Smoothing) line to present the conditional bias of the Ksat val-204 ues. The accuracy of the cross-validation predictions was evaluated using bias (mean error), root 205 mean square error (RMSE), coefficient of determination (R^2) and concordance correlation co-206 efficient (CCC). 207

Bias and RMSE are defined by: 208

$$bias = \sum_{i=1}^{n} \frac{(y_i - \hat{y}_i)}{n} \tag{1}$$

$$RMSE = \sqrt{\sum_{i=1}^{n} \frac{(\hat{y}_i - y_i)^2}{n}}$$
(2)

where y and \hat{y} are observed and predicted values and n is the total number of cross-validation points. 209 R^2 is defined as: 210

$$R^2 = \left[1 - \frac{SSE}{SST}\right]\%\tag{3}$$

where SSE is the sum of squared errors between the cross-validation predictions \hat{y} and the measurements y, and SST is the total sum of squares (proportional to variance of measurements). A coefficient of determination equal to 1 indicates that variance of the prediction errors is equal to zero but the bias may differ from zero.

In addition, Concordance Correlation Coefficient (CCC) (as measure of the agreement be tween observed and predicted Ksat values) of cross validation (CV) (Lawrence & Lin, 1989) is
 given by:

$$CCC = \frac{2 \cdot \rho \cdot \sigma_{\hat{y}} \cdot \sigma_{y}}{\sigma_{\hat{y}}^{2} + \sigma_{y}^{2} + (\mu_{\hat{y}} - \mu_{y})^{2}}$$
(4)

where $\mu_{\hat{y}}$ and μ_y are predicted and observed means, $\sigma_{\hat{y}}$ and σ_y are are predicted and observed variances and ρ is the Pearson correlation coefficient between predicted and observed values. CCC is equal to 1 for perfect model.

221 222

2.5 Comparision of accuracy of Ksat maps: CoGTF, Rosetta 3 and the map of Dai et al. (2019)

The accuracy of the predictions of Ksat by the three approaches was evaluated with a sub-223 set of the Ksat database that was selected in the following way: First, the surface of the Earth was 224 divided into blocks of 5 degrees as shown in Figure S1 in the SI. For fair comparison, Ksat mea-225 surements in blocks in North America or Europe were dropped because Rosetta 3 was mostly 226 calibrated with data from these regions (2525 Ksat values were outside of these regions). Then 227 we randomly selected blocks until about 20% of the remaining Ksat measurements had been cho-228 sen. These 508 Ksat measurements formed the test set for which predictions were extracted from 229 the Rosetta 3 and the Dai et al. (2019) maps. CoGTF predictions of Ksat were computed for these 508 test observations. The accuracy of the predictions by the three approaches was then evalu-231 ated with the same criteria as used for cross-validation. 232

233 3 Results

234

3.1 Model fitting and accuracy of modeled Ksat values

The CoGTF model fitted the logarithms of the Ksat measurements reasonably well (outof-bag RMSE = 0.73 (log_{10} Ksat in cm/day) and $R^2 = 0.66$). Figure 3 shows the list of most important covariates for Ksat modelling. The *x*-axis displays the average increase in node purity. The higher the value, the more important is a covariate. Figure 3 shows that sand content was found the most important covariate followed by elevation (important for soil formation and water flow), clay content, and bulk density. Climate covariates are dominating after the fifth covariate.

The results of the 5-fold cross-validation are presented in Figure 4a using hexbin density plots. For predictions of Ksat greater than equal to 10 cm/day the line of LOWESS falls onto the 1:1-line, hence the predictions were conditionally unbiased here. A slight positive conditional bias is visible for predictions less than 10 cm/day where the LOWESS line is below the 1:1 line. CoGTF tended to overestimate small Ksat values, but this bias remains small. Hence, RF predictions were both marginally and conditionally approximately unbiased. Cross-validation re-



Figure 3. Importance of the covariates for modeling Ksat by a random forest model. The x-axis displays the average increase in node purity (the larger the value, the more important is a covariate). The 14 most important covariates are shown here: sand content, bulk density (BD), and clay content belong to soil covariates. Elevation and topographic wetness index (TWI) are topography covariates. Temperature seasonality (TS), precipitation of driest month (PDM), cloud fraction (CF), minimum temperature of coldest month (MTCM), annual average land surface temperature (LST), maximum temperature of warmest month (MTWM), mean annual temperature (AMT), and mean annual rainfall (AMR) belong to climate category. Shortwave infrared (SWIR) Landsat-7 band is from the surface reflectance group.

sults show a reasonable overall model accuracy, with R^2 , CCC, and RMSE and bias equal to 0.66,

0.79, 0.72, and 0.0039 (log_{10} of Ksat in cm/day for RMSE and bias), respectively. The obser-

vations were also correlated with Rosetta 3 Ksat map (for this comparison, a total of 5,255 sam-

ples from shallow soil depth were selected out of 6,814 to compare with Rosetta 3 map for top

15 cm as shown in Figure 4b. RMSE and CCC was observed 1.23 and 0.12 (log_{10} of Ksat in cm/day

²⁵² for RMSE), respectively.

3.2 Global map of Ksat

Global Ksat maps were produced for four soil depths (0, 30, 60, and 100 cm). Figure 5a 254 shows the CoGTF map of Ksat at 0 cm soil depth, while results for other soil depths are provided 255 in Figure S2 (SI). Ksat values in the top layer (0 cm depth) vary between 0.05 to 31,600 cm/day. 256 High Ksat values were predicted for the equatorial belt and for parts of Russia and Canada, while 257 low Ksat values were produced in East America, Europe and parts of Asia (mainly India and North-258 East part of China). In general, Ksat value decreased with depth, with the most significant reduc-259 tion observed in North America, South America, China, India, and Russia (see Figures S2-S3 in 260 the SI). Figure 6 compares the probability distribution of the global Ksat map values with the dis-261 tribution of measured and fitted Ksat values for the 6,814 Ksat samples. Results show a more peaked 262 distribution of global Ksat map compared to the measured and fitted Ksat at the sampling loca-263 tions. Both measured (red) and fitted log_{10} Ksat showed the same mean values of 1.64 with stan-264 dard deviations 1.25 and 1.01, respectively, whereas the mean and standard deviation of global 265 map were observed 1.99 and 0.30 respectively. 266



Figure 4. Accuracy plots based on cross-validation: (a) correlation between observations and cross-validation predictions of log_{10} Ksat based on CoGTF model, (b) correlation between observations (0-30 cm soil depth) and Rosetta 3 predicted values from 0-15 cm map. The color codes the number of observations in each hexagonal pixel. The solid black line is 1:1 line and the blue dashed line is LOWESS curve (locally weighted scatterplot smoothing). The model accuracy of CoGTF was assessed using CCC (0.79) and RMSE (0.72). The RMSE and CCC between observations and Rosetta 3 predicted Ksat values were observed 1.23 and 0.12, respectively. The unit of RMSE is log_{10} of Ksat in cm/day.

3.3 Comparison with Rosetta 3 global Ksat map

267

The CoGTF Ksat map is compared with the Rosetta 3 map (Y. Zhang & Schaap, 2019) in 268 Figure 5. Note that there are different models of Rosetta 3 according to the soil information used 269 to build the neural network: H1w (information on soil textural class), H2w (sand, silt, and clay 270 percentage), H3w (sand, silt, and clay percentage plus bulk density), H4w (same information as 271 H3w plus water content at 330 cm suction), and H5w (same as H3w plus water content at 330 272 cm and at 15,000 cm) (X. Zhang et al., 2019). As standard model H3w is often chosen (see map 273 in Y. Zhang & Schaap, 2019) because information on water content at 330 cm and 15,000 cm is 274 sparse at global scale compared to bulk density and soil texture information. For comparison with 275 CoGTF, we chose H3w model as well. 276

The main differences between the CoGTF map and Rosetta 3 are the low Ksat values pre-277 dicted by Rosetta 3 for tropical regions and the abrupt change in Rosetta 3 predictions in high 278 latitude regions of Canada and Russia as a consequence of the strong sensitivity of Rosetta 3 pre-279 dictions on bulk density. In general, lower Ksat values were observed in the Rosetta 3 map com-280 pared to the CoGTF map for most regions worldwide except the northern regions (Canada and 281 Russia), while regions with coarser soils such as Sahara and middle East showed higher Ksat val-282 ues in Rosetta 3. The lower values of Ksat in Rosetta 3 than the in CoGTF map is evident in Fig-283 ure 6a. Medians of the common logarithm of Ksat (unit cm/day) were equal to 1.62 and 2.00, re-284 spectively (Figure 6b). 285



a) CoGTF map (0 cm)

Figure 5. Visual comparison between (a) CoGTF Ksat map, and (b) map based on Rosetta 3 PTF. Ksat values predicted by Rosetta 3 were higher for sandy soils (Sahara) and in northern regions with smaller bulk density. The scale of the maps was truncated at minimum and maximum values of 10 and 1000 cm/day to show the significant variations in the maps

3.4 Validation of global Ksat maps

286

Table 1 shows the results of the comparison of the accuracy of Ksat predictions for the CoGTF, 287 Rosetta 3 and Dai et al. (2019) maps (see Figures S7 and S8 for the map of Dai et al., 2019, with 288 CoGTF map). A total of 372 Ksat samples out of the validation dataset with 508 samples (we 289 selected samples with soil depth 0-30 cm) were compared with measured Ksat values and RMSE 290 values of 1.02, 1.29, and 1.15 were computed (log10 of Ksat in cm/day) for the CoGTF map, Rosetta 291 3, and Dai et al. (2019) map, respectively. The RMSE illustrates that CoGTF map showed bet-292 ter performance than the other maps. However, RMSE of 1 also shows that the precision is lim-293 ited for CoGTF as well. The better performance of CoGTF is manifested in the much lower bias 294 compared to the two other models. 295



Figure 6. Difference in probability density functions (PDF): (a) between global CoGTF map (yellow) and Rosetta 3 (black) Ksat values at 0 cm depth, measured (red) and fitted (blue) Ksat values at the sampling sites, (b) cumulative distribution functions for Rosetta 3 map (black) and CoGTF map (yellow) for soil depth 0 cm.

Table 1. Root mean square error (RMSE) and bias of predictions of log_{10} (Ksat) (units cm/day) for test data. A total of 372 Ksat sample points were selected to investigate the accuracy of Ksat predictions (0-30 cm soil depth were used). The negative signs in bias demonstrate that all three models underestimated Ksat values. The range shows the minimum and maximum values of 372 samples.

Models	Samples used	RMSE	bias	Extracted points range
CoGTF (0 -15 cm)	372	1.02	-0.19	0.85-2.60
Rosetta 3 (0 -15 cm)	372	1.29	-0.75	0.83-2.64
Dai et al. (2019) (0 -15 cm)	372	1.15	-0.51	0.68-2.30

4 Discussion

4.1 Characteristics of the CoGTF global Ksat maps

In this paper we have produced global estimates of Ksat by linking terrain, climate, veg-298 etation and soil spatial covariates to measured Ksat values, thus injecting local information usu-299 ally ignored by traditional PTFs. We refer to this approach as the Covariate-based Geo Trans-300 fer Functions (CoGTF) framework. The newly developed global CoGTF map of Ksat (Figure 5) 301 shows high values in the Northern part of South America, the central part of Africa and South-302 east Asia (mainly Indonesia, Malaysia, Myanmar (Burma), Philippines, Singapore, and Thailand), 303 most likely due to high rainfall, temperature, and vegetation. Our results shows (Figure 3) that 304 rainfall, temperature and their variation are the most important climate covariates for the Ksat 305 mapping (Shoji et al., 2006). This indicates that these climatic factors not only act as catalyst in 306 soil chemical reactions but also determine the type and biomass of vegetation that is important 307 for soil structure formation. This impact of vegetation on soil Ksat is in line with the research 308

²⁹⁷

by Niemeyer et al. (2014) who compared the leaf area index with Ksat and observed that high leaf area index increases the Ksat (with R-square = 0.33).

The central part of India, eastern part of Australia, and parts of China showed low Ksat val-311 ues due to the presence of high clay content that reduces the soil permeability (see as well dis-312 cussion on role of clay mineral type in section 4.2). The west part of North America, middle east 313 countries (Tibet, Iran, Turkey), and northern parts of Algeria have low Ksat values that may be 314 related to high elevation, low rainfall, less vegetation and thus less structure formation processes. 315 Many studies have recognized the indirect influence of elevation on soil proprieties (Leij et al., 316 2004; Carter & Ciolkosz, 1991). Similarly, different land-use (forest or pasture) directly impact 317 Ksat. Chandler et al. (2018) showed that forests had larger soil hydraulic conductivity than pas-318 tures. 319

Likewise, high values of Ksat up to around 100 to 300 cm/day are observed in desert re-320 gions such as Thar desert in India, northern and southern Africa, and central Australia, where dom-321 inating fractions of sand cause high water permeability. Similarly, Colombia and Peru showed 322 high Ksat values due to high organic carbon content (Allison, 1973). Furthermore, high Ksat val-323 ues were observed in parts of Brazil that strongly decreased with depth. Similar results were re-324 ported by Belk et al. (2007). They conducted a study in the tropical forest of Brazil and measured 325 the Ksat at various depths for different sites. The authors found that Ksat values at surface were 326 mainly between 100 to 1000 cm/day and decreased with depth. 327

328

4.2 Effect of clay type — active and inactive clay minerals

Pedotransfer functions like Rosetta 3 and the ensemble of PTFs used in (Dai et al., 2019) 329 to estimate soil hydraulic properties based on clay fraction and do not take into account the large 330 differences in microstructure and hydration of different clay minerals. The remarkable spatial 331 segregation in climatic regions of different clay minerals (see Ito & Wagai, 2017) and the differ-332 ent hydraulic properties of the clay minerals indicate that PTFs built for temperate regions with 333 swelling clays cannot be applied for tropical regions with non-swelling clays (see Ottoni et al., 334 2018). In tropical soils, dense vegetation, and non-swelling ('inactive') kaolinite clay minerals 335 result in higher conductivities (Hodnett & Tomasella, 2002) in contrast to PTFs that are trained 336 with data from temperate soils with swelling (more 'active') clays. This is further discussed for 337 estimates relevant to Brazil shown in Figure 7. 338

In Figure 7, the CoGTF and Rosetta 3 Ksat maps are shown together with six covariates 339 and clay mineral map. The Ksat values predicted with CoGTF are one order of magnitude higher 340 than based on Rosetta 3. The difference stems from the dominant role of soil texture for Rosetta 341 3 as illustrated with a black polygon in Figure 7: the polygon marks a region of high sand con-342 tent and low clay content that is manifested in relatively high values of Ksat for Rosetta 3, with 343 values typical for temperate regions. For CoGTF, the conductivity in this 'sand band' is relatively 344 low because other covariates and processes are more important. These lower values coincidence 345 with low elevation. The important role of elevation in CoGTF is also manifested in the high Ksat 346 values in the mountainous region in the south and the low Ksat values in the Amazon region. An-347 other reason for the lack of correlation between Ksat and texture for CoGTF in Brazil is the in-348 active clay mineral type (kaolinite) that does not limit Ksat the same way as in case of more ac-349



Figure 7. Predicted Ksat values for Brazil (a), spatial patterns of the Rosetta 3 Ksat map in \log_{10} cm/day (b) and of the first four most important covariates (c-f, see Figure 3): sand fraction (%), elevation (meters above sea level), bulk density (g/cm³) and clay fraction (%). Other covariates that are related to soil formation to link with Ksat are shown as well (g-i): mean annual rainfall (mm), Copernicus fraction of absorbed photosynthetically active radiation (FAPAR, values in %) and kaolinite (in %) clay mineral. The region with black polygon marks a region with high sand and low clay content that is expressed in Rosetta 3 as band of relatively high Ksat values. In contrast to Rosetta 3, CoGTF is not dominated by soil texture but takes into account covariates that are important for soil formation (here mainly the elevation).

- tive clay in temperate regions. Precipitation and temperature are the main reasons for the strong
- ³⁵¹ chemical weathering of the rock and the formation of the non-swelling, kaolinite clay minerals
- (Montes et al., 2002). It is evident in Figure 7(g to i) that in the region with low rainfall and veg-
- etation kaolinite percentage is lower than other regions with high rainfall and vegetation.

In contrast to Brazil, India is a region with active (swelling) clay minerals. In contrast to 354 the inactive kaolinite in Brazil, for the active clays in India, low Ksat values can be expected. Fig-355 ure S4a and S4i show the correlation between low Ksat values and high contents of smectite clay 356 mineral. The low values of Ksat in central India directly relate to high clay content, low vege-357 tation biomass and low mean annual rainfall (see Figure S4 for covariates in SI). Figure S4b il-358 lustrates the Ksat values from Rosetta 3 for India. The patterns of high active clay fraction in In-359 dia is not captured by Rosetta 3 model. This might be the effect of considering only soil basic 360 properties or using samples from only temperate region. 361

362

4.3 Effects of information clustering — The Florida database example

Out of 13,267 Ksat values, only 6,814 values were used for the Ksat mapping to avoid a 363 distortion of the Ksat predictions by the many data from Florida. The dataset contained 6,532 364 Ksat values from Florida but we used only 1% of these points for mapping. Figure S5a and S5b 365 compares the map computed with all 13,267 Ksat measurements with the map trained on 6,814 366 measurements. The difference between these maps (Figure S5c) showed a large impact on the 367 sandy regions such as Sahara and center part of Africa and middle east with significantly higher 368 Ksat values when all Florida points are included in the fitting. A similar effect was observed in 369 parts of South America and Australia. On the other hand, the south of Africa and the higher north-370 ern latitudes showed higher Ksat values when only 1% of the Florida data was used. 371

372

4.4 Improved model performance using remote sensing covariates

As we described above, the RS (vegetation, topography, climatic) covariates could be used to harness the heterogeneity produced by these environmental variables as these factors shape clay activity and soil-forming processes that control saturated hydraulic conductivity (Ottoni et al., 2018; Hao et al., 2019). To investigate this effect of RS covariates in the predictions, we fitted the RF model only with soil properties or remote sensing covariates. The maps are shown in Figure S6a and S6b in the SI.

Table 2. Root mean square error (RMSE) and coefficient of determination (R^2) for different models.

Models	RMSE	R ²	Total covariates	Best mtry
CoGTF	0.72	0.66	28	6
Only soil covariates	0.75	0.63	4	2
Only RS covariates	0.73	0.65	24	16

Table 2 shows the RMSE and R^2 using different models where we used only soil covariates, only remote sensing covariates, and the CoGTF model. Remote sensing (RMSE = 0.73; R^2 = 0.65) predicted the Ksat better than only soil covariates (RMSE = 0.75; R^2 = 0.63). Similarly, the CoGTF model showed lower RMSE (0.72) and higher R^2 (0.66) than only RS covariate. Hence, consideration of RS covariates in predicting hydraulic properties could increase the accuracy of the predictions of soil hydraulic properties compared to a model that is based only on soil information.

4.5 Usage of the global CoGTF Ksat maps and future developments

We observed that RMSE in the model validation for the CoGTF map was better than the other maps. However, RMSE with 1 (log_{10} Ksat in units of cm/day) also shows that the precision of even the CoGTF map is not so good. On the other hand, the bias for CoGTF map was much better than for the other maps. Although, the predictions are not so accurate, it shows the one step ahead in terms of improvement in the predictions using distributed Ksat dataset and consideration of RS covariates.

The global CoGTF maps can be used to extract the information of Ksat at different depths 393 for local, regional, and global scale studies. On the local scale, these maps can be helpful in agro-394 nomic processes such as soil interpretation, water-plant relationships, and assessing soil suitabil-395 ity for agriculture. For regional and global scale, the maps could provide unique values to each 396 pixel in watershed scale and Earth surface models and would enhance the heterogeneity and ac-397 curacy in the area. The maps could also be useful for the soil water management policies as guide-398 line to show where soil reclamation is required to reduce and enhance the hydraulic conductiv-399 ity. 400

The actual CoGTF map has a resolution of 1 km. This resolution may be improved in the near future considering various initiatives to estimate soil and RS information with higher resolution. But independent of improved resolutions, subgrid information on Ksat may be required for a catchment when specific information on soil texture or vegetation type is available. For such applications, we are actually developing a parametric model of CoGTF so that Ksat can be estimated as a linear combination of most important covariates.

407 5 Conclusions

Soil saturated hydraulic conductivity is an important soil property for the parameterization 408 of Earth system and land surface models. The major limitations of currently available maps are 409 that (1) they are developed using a limited number of Ksat measurements mainly from temper-410 ate regions, (2) they are derived only from basic soil properties thus ignoring the effect of biologically-411 induced soil structure as well as clay mineralogy, and (3) they are not benefiting from the wealth 412 of local remote sensing (RS) covariates. Therefore, we proposed a new global map of Ksat ob-413 tained by linking the measured Ksat values (6,814 samples) with 24 remote sensing covariates 414 and 3 soil properties (sand content, clay content and bulk density) to add local information on 415 vegetation, climate, and topography. The new map combines georeferenced information of soil 416 properties and remote sensing covariates and is called covariate-based Geo Transfer Functions 417 (CoGTF) map. We used the random forest machine-learning algorithm to fit the Ksat models and 418 the performance was assessed using CCC and RMSE which was computed using 5 fold cross-419 validation. The CCC and RMSE (in log₁₀ Ksat given in cm/day) were observed 0.79 and 0.72, 420 respectively. The CoGTF global Ksat map was compared with the map calculated with the well 421 known Rosetta 3 PTF and major differences between the two maps were found. Firstly, Ksat val-422 ues in Rosetta 3 were much lower for tropical regions compared to the CoGTF map. The trop-423 ical regions are expected to have rather high Ksat values due to intense soil formation processes 424 and presence of more conductive clay minerals (kaolinite). The effects of active and inactive clay 425 minerals on Ksat are captured in CoGTF map as formation of clay minerals are linked to precip-426

- itation, temperature and dense vegetation. Secondly, in CoGTF there is no abrupt change in Ksat
- as shown in Rosetta 3 map for the higher latitude regions such as Canada and Russia. This large
- 429 contrast is related to a change in bulk density that is dominant in Rosetta 3. In CoGTF, RS co-
- variates pattern cover this contrast. Furthermore, the CoGTF map, Rosetta 3 map, and the map
- of Dai et al. (2019) were validated using test data that were not used to calibrate the models, and
- the result showed that the CoGTF map performed better than the other models. Consequently,
- 433 we propose to transition from PTFs based only on soil texture and bulk density to spatial-association
- of climate and vegetation covariates ("GTFs") to estimate Ksat. The study provides a blueprint
- for how georeferenced covariates could be used within the machine learning framework to im-
- 436 prove Ksat predictive mapping. Moreover, the resulting CoGTF global maps are readily updat-
- able as more information becomes available (covariates of measured Ksat).

438 Acknowledgments

- 439 The study was supported by ETH Zurich (Grant ETH-18 18-1). We would like to thank Zhong-
- wang Wei, Samuel Bickel and Simone Fatichi (ETH Zurich) for insightful discussions. The data
- sets produced in this study are available at https://doi.org/10.5281/zenodo.3934853.

442 References

- Ahuja, L., Cassel, D., Bruce, R., & Barnes, B. (1989). Evaluation of spatial distribution of
 hydraulic conductivity using effective porosity data. *Soil Science*, *148*(6), 404–411.
- Allison, F. E. (1973). Soil organic matter and its role in crop production. Elsevier.
- Araya, S. N., & Ghezzehei, T. A. (2019). Using machine learning for prediction of saturated hydraulic conductivity and its sensitivity to soil structural perturbations. *Water Resources Research*, 55(7), 5715–5737.
- Belk, E. L., Markewitz, D., Rasmussen, T. C., Carvalho, E. J. M., Nepstad, D. C., & David-
- 450 son, E. A. (2007). Modeling the effects of throughfall reduction on soil water content 451 in a brazilian oxisol under a moist tropical forest. *Water Resources Research*, 43(8).
- Bouma, J. (1989). Using soil survey data for quantitative land evaluation. In *Advances in soil science* (pp. 177–213). Springer.
- ⁴⁵⁴ Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5–32.
- Brocca, L., Filippucci, P., Hahn, S., Ciabatta, L., Massari, C., Camici, S., ... Wagner, W.
 (2019). SM2RAIN-ASCAT (2007–2018): Global daily satellite rainfall from ASCAT soil moisture. *Earth Syst. Sci. Data Discuss*, 1–31.
- Buchhorn, M., Bertels, L., Smets, B., Lesiv, M., & Wur, N. (2017). Copernicus Global Land
 Operations "Vegetation and Energy". *Copernicus Global Land Operations "Vegetation and Energy*.
- Carter, B. J., & Ciolkosz, E. J. (1991). Slope gradient and aspect effects on soils developed
 from sandstone in pennsylvania. *Geoderma*, 49(3-4), 199–213.
- Chandler, K., Stevens, C., Binley, A., & Keith, A. (2018). Influence of tree species and forest
 land use on soil hydraulic conductivity and implications for surface runoff generation.
 Geoderma, *310*, 120–127.
- Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., ... Böhner, J.
- 467 (2015). System for automated geoscientific analyses (saga) v. 2.1.4. *Geoscientific*

468	Model Development, 8(7), 1991-2007. doi: 10.5194/gmd-8-1991-2015
469	Dai, Y., Xin, Q., Wei, N., Zhang, Y., Shangguan, W., Yuan, H., Lu, X. (2019). A global
470	high-resolution dataset of soil hydraulic and thermal properties for land surface model-
471	ing. Journal of Advances in Modeling Earth Systems.
472	Fashi, F. H., Gorji, M., & Shorafa, M. (2016). Estimation of soil hydraulic parameters for
473	different land-uses. Modeling Earth Systems and Environment, 2(4), 1-7.
474	Fatichi, S., Or, D., Walko, R., Vereecken, H., Young, M. H., Ghezzehei, T. A., Avissar,
475	R. (2020). Soil structure is an important omission in earth system models. Nature
476	Communications, 11.
477	Gupta, S., Hengl, T., Lehmann, P., Bonetti, S., & Or, D. (2020, April). SoilKsatDB: a global
478	compilation of soil saturated hydraulic conductivity measurements.
479	doi: 10.5281/zenodo.3752721
480	Gutmann, E. D., & Small, E. E. (2007). A comparison of land surface model soil hydraulic
481	properties estimated by inverse modeling and pedotransfer functions. Water Resources
482	<i>Research</i> , 43(5).
483	Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A.,
484	Townshend, J. R. G. (2013). High-resolution global maps of 21st-century forest
485	cover change. Science, 342(6160), 850-853. doi: 10.1126/science.1244693
486	Hao, M., Zhang, J., Meng, M., Chen, H. Y., Guo, X., Liu, S., & Ye, L. (2019). Impacts of
487	changes in vegetation on saturated hydraulic conductivity of soil in subtropical forests.
488	Scientific reports, 9(1), 1–9.
489	Hengl, T., de Jesus, J. M., Heuvelink, G. B., Gonzalez, M. R., Kilibarda, M., Blagotić, A.,
490	others (2017). Soilgrids250m: Global gridded soil information based on machine
491	learning. <i>PLoS one</i> , <i>12</i> (2), e0169748.
492	Hengl, T., & MacMillan, R. A. (2019). Predictive Soil Mapping with R. Lulu. com.
493	Hodnett, M., & Tomasella, J. (2002). Marked differences between van genuchten soil
494	water-retention parameters for temperate and tropical soils: a new water-retention
495	pedo-transfer functions developed for tropical soils. Geoderma, 108(3-4), 155–180.
496	Ito, A., & Wagai, R. (2017). Global distribution of clay-size minerals on land surface for bio-
497	geochemical and climatological studies. Scientific data, 4, 170103.
498	Jana, R. B., & Mohanty, B. P. (2011). Enhancing ptfs with remotely sensed data for multi-
499	scale soil water retention estimation. <i>Journal of hydrology</i> , 399(3-4), 201–211.
500	Jenny, H. (1994). Factors of soil formation: a system of quantitative pedology. Courier Cor-
501	poration.
502	Jorda, H., Bechtold, M., Jarvis, N., & Koestel, J. (2015). Using boosted regression trees to
503	explore key factors controlling saturated and near-saturated hydraulic conductivity.
504	European Journal of Soil Science, 66(4), 744–756.
505	Karger, D. N., Conrad, O., Böhner, J., Kawohl, T., Kreft, H., Soria-Auza, R. W., Kessler,
506	M. (2017). Climatologies at high resolution for the earth's land surface areas. <i>Scien</i> -
507	<i>tific data</i> , <i>4</i> , 170122.
508	Khlosi, M., Alhamdoosh, M., Douaik, A., Gabriels, D., & Cornelis, W. (2016). Enhanced
509	pedotransfer functions with support vector machines to predict water retention of
510	calcareous soil. European Journal of Soil Science, 67(3), 276–284.

511	Kuhn, M. (2012). The caret package. R Foundation for Statistical Computing, Vienna, Aus-
512	tria. URL https://cran. r-project. org/package= caret.
513	Lawrence, I., & Lin, K. (1989). A concordance correlation coefficient to evaluate repro-
514	ducibility. Biometrics, 255-268.
515	Leij, F., Alves, W., Van Genuchten, M. T., & Williams, J. (1996). The unsoda unsaturated
516	soil hydraulic database; user's manual, version 1.0. Rep. EPA/600/R-96, 95, 103.
517	Leij, F., Romano, N., Palladino, M., Schaap, M. G., & Coppola, A. (2004). Topographical at-
518	tributes to predict soil hydraulic properties along a hillslope transect. Water Resources
519	<i>Research</i> , 40(2).
520	Mayr, T., & Jarvis, N. (1999). Pedotransfer functions to estimate soil water retention param-
521	eters for a modified brooks-corey type model. Geoderma, 91(1-2), 1-9.
522	Montes, C. R., Melfi, A. J., Carvalho, A., Vieira-Coelho, A. C., & Formoso, M. L. (2002).
523	Genesis, mineralogy and geochemistry of kaolin deposits of the jari river, amapá state,
524	brazil. Clays and Clay Minerals, 50(4), 494-503.
525	Montzka, C., Herbst, M., Weihermüller, L., Verhoef, A., & Vereecken, H. (2017). A global
526	data set of soil hydraulic properties and sub-grid variability of soil water retention and
527	hydraulic conductivity curves. Earth System Science Data, 9(2), 529–543.
528	Nemes, A., Rawls, W. J., & Pachepsky, Y. A. (2005). Influence of organic matter on the es-
529	timation of saturated hydraulic conductivity. Soil Science Society of America Journal,
530	69(4), 1330–1337.
531	Nemes, A., Schaap, M., Leij, F., & Wösten, J. (2001). Description of the unsaturated soil hy-
532	draulic database unsoda version 2.0. <i>Journal of Hydrology</i> , 251(3-4), 151–162.
533	Niemeyer, R., Fremier, A. K., Heinse, R., Chávez, W., & DeClerck, F. A. (2014). Woody
534	vegetation increases saturated hydraulic conductivity in dry tropical Nicaragua. <i>Vadose</i>
535	Zone Journal, 13(1).
536	Obi, J., Ogban, P., Ituen, U., & Udoh, B. (2014). Development of pedotransfer functions for
537	coastal plain soils using terrain attributes. <i>Catena</i> , <i>123</i> , 252–262.
538	Or, D. (2019). The tyranny of small scales—on representing soil processes in global land sur-
539	face models. <i>Water Resources Research</i> .
540	Ottoni, M. V., Ottoni Filho, I. B., Schaap, M. G., Lopes-Assad, M. L. R., & Rotunno Filho,
541	U. C. (2018). Hydrophysical database for Brazilian soils (HYBRAS) and pedotransier functions for water retention. <i>Values Zone Journal</i> , 17(1)
542	Initiality of water retention. <i>Values Zone Journal</i> , 17(1).
543	Peters, J., De Baets, B., vernoest, N. E., Samson, R., Degroeve, S., De Becker, P., & Huy-
544	alling coological modelling 207(2,4) 204, 218
545	P. Core Team. (2012). D: A Language and Environment for Statistical Computing [Computer
546	software manuall Vienna Austria Patriaved from http://www.R. project.org/
547	Pad M P P Toomanian N Khormali E Brungard C W Komaki C P & Boggart P
548	(2014) Undating soil survey maps using random forest and conditioned latin hyper
549	cube sampling in the loss derived soils of northern iran. <i>Geoderma</i> 232 97–106
550	Rawls W I Brakensiek D I & Saxtonn K (1982) Fetimation of soil water properties
553	Transactions of the ASAE 25(5) 1316–1320
552	Rodrigues M & de la Riva I (2014) An insight into machine-learning algorithms to
555	model human-caused wildfire occurrence Environmental Modelling & Software 57

555	192–201.
556	Santra, P., & Das, B. S. (2008). Pedotransfer functions for soil hydraulic properties devel-
557	oped from a hilly watershed of eastern india. Geoderma, 146(3-4), 439-448.
558	Saxton, K. E., & Rawls, W. J. (2006). Soil water characteristic estimates by texture and or-
559	ganic matter for hydrologic solutions. Soil science society of America Journal, 70(5),
560	1569–1578.
561	Schaap, M. G., Leij, F. J., & Van Genuchten, M. T. (2001). Rosetta: A computer program for
562	estimating soil hydraulic parameters with hierarchical pedotransfer functions. Journal
563	<i>of hydrology</i> , <i>251</i> (3-4), 163–176.
564	Sharma, S. K., Mohanty, B. P., & Zhu, J. (2006). Including topography and vegetation at-
565	tributes for developing pedotransfer functions. Soil Science Society of America Jour-
566	nal, 70(5), 1430-1440.
567	Shoji, S., Nanzyo, M., & Takahashi, T. (2006). Factors of soil formation: climate. as exem-
568	University Press Cambridge UK 131 149
569	Tomospile L & Hodpett M C (1008). Estimating coil water retention characteristics from
570	limited data in brazilian amazonia Soil science 163(3) 190–202
571	Tootchi A Jost A & Ducharne A (2019) Multi-source global wetland maps combin-
572	ing surface water imagery and groundwater constraints <i>Earth System Science Data</i>
574	11, 189–220.
575	Tóth, B., Weynants, M., Nemes, A., Makó, A., Bilas, G., & Tóth, G. (2015). New genera-
576	tion of hydraulic pedotransfer functions for europe. European journal of soil science,
577	66(1), 226–238.
578	Vereecken, H., Schnepf, A., Hopmans, J. W., Javaux, M., Or, D., Roose, T., others
579	(2016). Modeling soil processes: Review, key challenges, and new perspectives.
580	Vadose zone journal, 15(5).
581	Wilson, A. M., & Jetz, W. (2016, 03). Remotely sensed high-resolution global cloud dy-
582	namics for predicting ecosystem and biodiversity distributions. PLOS Biology, 14(3),
583	1-20. Retrieved from http://dx.doi.org/10.1371%2Fjournal.pbio.1002415
584	doi: 10.1371/journal.pbio.1002415
585	Wösten, J., Lilly, A., Nemes, A., & Le Bas, C. (1999). Development and use of a database of
586	hydraulic properties of european soils. <i>Geoderma</i> , 90(3-4), 169–185.
587	Wright, M. N., & Ziegler, A. (2015). ranger: A fast implementation of random forests for
588	high dimensional data in c++ and r. arxiv preprint arXiv:1508.04409.
589	Yamazaki, D., Ikeshima, D., Tawatari, R., Yamaguchi, T., O'Loughlin, F., Neal, J. C.,
590	Bates, P. D. (2017). A mgn-accuracy map of global terrain elevations. <i>Geophysical</i> <i>Pasagraph Letters</i> 44(11), 5844, 5853
591	Zhang X. Zhu I. Wandroth O. Matacha C. & Edwards D. (2010) Effort of marra
592	porosity on pedotransfer function estimates at the field scale Vadose Zong Journal
594	18(1).
595	Zhang, Y., & Schaap, M. G. (2017). Weighted recalibration of the Rosetta pedotransfer
596	model with improved estimates of hydraulic parameter distributions and summary

statistics (Rosetta3). Journal of Hydrology, 547, 39-53.

597

- Zhang, Y., & Schaap, M. G. (2019). Estimation of saturated hydraulic conductivity with pe dotransfer functions: A review. *Journal of Hydrology*, 575, 1011–1030.
- Zimmermann, A., Schinn, D. S., Francke, T., Elsenbeer, H., & Zimmermann, B. (2013).
 Uncovering patterns of near-surface saturated hydraulic conductivity in an overland
- flow-controlled landscape. *Geoderma*, 195, 1–11.



Journal of Advances in Modelling Earth Systems

Supporting Information for

Global prediction of soil saturated hydraulic conductivity using random forest in a Covariate-based Geo Transfer Functions (CoGTF) framework

Surya Gupta^{1,*}, Tomislav Hengl^{2, 3}, Peter Lehmann¹, Sara Bonetti^{1, 4}, Andreas Papritz¹, Dani Or^{1,5}

¹ Soil and Terrestrial Environmental Physics, Department of Environmental Systems Science, ETH Zurich, Zurich, Switzerland,

² Envirometrix Ltd., Wageningen, the Netherlands

³OpenGeoHub, Wageningen, the Netherlands

⁴ Bartlett School of Environment, Energy and Resources, University College London, London, UK

⁵ Division of Hydrologic Sciences, Desert Research Institute, Reno, NV, USA

* Correspondence to:

Surya Gupta (surya.gupta@usys.ethz.ch)

Contents of this file

Figures S1 to S8

Table S1

Introduction

This supplementary information provides the figures of block-wise cross-validation (Figure S1), Ksat values at different depth computed with CoGTF (Figures S2 and S3), spatial distributions of Ksat and different environmental covariates in India (Figure S4), effect of clustering of Ksat samples on global map (Figure S5), global Ksat maps predicted with remote sensing or soil covariates (Figure S6) and a comparison between CoGTF and Dai et al., 2019 (Figures S7 and S8). At the end, a table lists the environmental covariates used in this study.



Figure S1 Regionalization of global map for validation of CoGTF model. a) 5 degrees by 5 degrees grids plotted with positions of Ksat dataset (6,814 samples). A total of 168 grid cells contains the data points. b) 30 blocks of data were removed randomly (i.e 20% of 2,525 Ksat dataset) for validation. The 2,525 Ksat samples are a subset of the totally 6,814 samples because samples from Europe and North America were excluded (they were used to train Rosetta 3 model and could not be considered for model validation).



Figure S2 Ksat maps at different depths a) 0 cm b) 30 cm, b) 60 cm, c) 100 cm computed with CoGTF



Figure S3 Cumulative distribution function CDF for global maps of Ksat at different depths (0, 30, 60 and 100 cm) computed with CoGTF.



Figure S4 Ksat values for India predicted with CoGTF (a), spatial patterns of the Rosetta 3 Ksat map b), and the first four most important covariates (c-f, see figure3 in the main text): sand fraction (%), elevation (meters above sea level), bulk density (g/cm³) and clay fraction (%). Other covariates important for soil formation liked with Ksat are shown as well (g-i): mean annual rainfall (mm), fraction of absorbed photosynthetically active radiation (FAPAR, values in %) and kaolinite (in %) clay mineral.



Figure S5 The difference between Ksat map including all Florida samples (a) and using only 1% of these Florida Ksat points (b) to build the CoGTF model. In the maps of differences (c), blue color represents higher values when all Florida points are included, yellow represents approximately the same value in both maps, and red shows locations with higher Ksat when only 1% of Florida samples are included. The Florida cluster showed a large impact on the sandy regions such as Sahara and center part of Africa and middle east as it significantly increased the Ksat values. A similar effect was observed in parts of South America and Australia. On the other hand, south of Africa and higher Nothern latitude showed high Ksat values for map that includes 1% of Florida samples.



Figure S6 Ksat maps computed with Random Forest approach for soil depth of 0 cm. a) Only 24 remote sensing covariates were used to build model and to compute the map. b) Only soil properties were used (sand content, clay content and bulk density). Note that high contrast in northern latitudes in Eurasia are controlled by changes in bulk density (a dominant pattern in Rosetta 3 map).



Figure S7 a) Ksat map at 0-5 cm depth from Dai et al. (2019) computed from an ensemble of 16 pedotransfer functions. The map used soil information from Global Soil Dataset for Earth System Models (GSDE; Shangguan et al., 2017) and SoilGrids (Hengl et al., 2017). b) CoGTF Ksat map at 0 cm.



Figure S8 Difference in probability density functions (PDF) of Ksat values: (a) between global CoGTF map (yellow) and Dai et al. (2019) (black dotted line), measured (red) and fitted (blue) Ksat values at the sampling sites; (b) cumulative distribution functions for Dai et al. (2019) map (black dotted line) and CoGTF map (yellow) for soil depth 0 cm.

S.	List of Covariates	Source
no		
	Climate	
1	clm_annual mean	http://chelsa-
	temperaturebio1_m_1km_s00cm_1979-	climate.org/bioclim/
-	2013_V1.0	$(K_{\text{reserved}}, 1, 2017)$
2	cim_temperature seasonalitybio4_ m_1km_	(Karger et al., 2017)
2	s00cm_1979-2015_v1.0	-
5	monthbio5 m 1km s0 0cm 1979-2013 v1 0	
Δ	clm_min_temperature_of_coldest_monthbio6	
-	m 1km s00cm 1979-2013 v1.0	
5	clm_annual precipitationbio12	-
-	m_1km_1979_2013_v1.0	
6	clm_precipitation of wettest monthbio13_	
	m_1km_1979_2013	
7	clm_precipitation of driest monthbio14_	
	m_1km_1979_2013	
8	clm_cloud.fraction_earthenv.modis.annual_m_1	http://www.earthenv.org/cloud
	km_s00cm_20002015_v1.0	
-		(Wilson & Jetz, 2016).
9	clm_diffuse.irradiation_solar.atlas.kwhm2.100_	
10	m_1km_s00cm_2016_v1	<u>https://globalsolaratias.info/do</u>
10	1km s0 0cm 2016 v1	wilload/world
11	IKII_S00CII_2010_V1	https://lpdaac.usgs.gov/product
11	2000 2017 v1 0	s/mod11a2v006/
12	clm lst mod11a2 annual day sd 1km s0 0cm	<u>5/110011027000/</u>
	2000.2017 v1.0	
13	clm_precipitation_sm2rain.annual_m_1km_s0	https://zenodo.org/record/3405
	0cm_20072018_v0.2	563#.XlgdNTFKhaQ
		(Brocca et al., 2019)
	Digital terrain model	
14	dtm_twi_merit.dem_m_1km_s00cm_2017_v1.	https://zenodo.org/record/1447
	0	210#.XllTejFKhaQ
15	dtm_slope_merit.dem_m_1km_s00cm_2017_v	
	1.0	(Yamazaki et al., 2017)
16	dtm_aspect.cosine_merit.dem_m_1km_s00cm	
17	_2018_v1.0	
17	dtm_elevation_merit.dem_m_1km_s00cm_20	
10	17_v1.0	
18	dtm_lithology_usgs.ecotapestry.acid.plutonics_	
	p 1km s00cm 2014 v1.0	

 $\label{eq:stable} \textbf{Table S1} \ \textbf{List} \ \textbf{of covariates used for creating the Ksat map}$

	Surface reflectance	
19	lcv_landsat.nir_wri.forestwatch_m_1km_s00c	Hansen et al. (2013)
	m_2018_v1.2	
20	lcv_landsat.red_wri.forestwatch_m_1km_s00c	
	m_2018_v1.2	
21	lcv_landsat.swir2_wri.forestwatch_m_1km_s0	
	0cm_2018_v1.2	
22	lcv_snow_probav.lc100_p_1km_s00cm_2017_	Tsendbazar et al. (2017)
	v1.0	
23	lcv_wetlands.regularly.flooded_upmc.wtd_p_1k	https://doi.pangaea.de/10.1594/
	m_b0200cm_20102015_v1.0	PANGAEA.892657
		(Tootchi et al., 2019)
	Vegetation covariates	
24	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20	https://land.copernicus.eu/glob
24	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0	https://land.copernicus.eu/glob al/products/fapar
24	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0 Predicted soil properties	https://land.copernicus.eu/glob al/products/fapar
24	Vegetation covariatesveg_fapar_proba.v.annnual_d_1km_s00cm_20142019_v1.0Predicted soil propertiessol_clay.wfraction_usda.3a1a1a_m_1km_b0_10	https://land.copernicus.eu/glob al/products/fapar
24	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0 Predicted soil properties sol_clay.wfraction_usda.3a1a1a_m_1km_b0_10 _30_60_100_200cm_	https://land.copernicus.eu/glob al/products/fapar
24	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0 Predicted soil properties sol_clay.wfraction_usda.3a1a1a_m_1km_b0_10 _30_60_100_200cm_ 19502017_v0.2	https://land.copernicus.eu/glob al/products/fapar
24 25 26	Vegetation covariatesveg_fapar_proba.v.annnual_d_1km_s00cm_20142019_v1.0Predicted soil propertiessol_clay.wfraction_usda.3a1a1a_m_1km_b0_10_30_60_100_200cm_19502017_v0.2sol_sand.wfraction_usda.3a1a1a_m_1km_	https://land.copernicus.eu/glob al/products/fapar https://www.openlandmap.org/
24 25 26	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0 Predicted soil properties sol_clay.wfraction_usda.3a1a1a_m_1km_b0_10 _30_60_100_200cm_ 19502017_v0.2 sol_sand.wfraction_usda.3a1a1a_m_1km_ b0_10_30_60_100_200cm_19502017_v0.2	https://land.copernicus.eu/glob al/products/fapar https://www.openlandmap.org/
24 25 26 27	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0 Predicted soil properties sol_clay.wfraction_usda.3a1a1a_m_1km_b0_10 _30_60_100_200cm_ 19502017_v0.2 sol_sand.wfraction_usda.3a1a1a_m_1km_ b0_10_30_60_100_200cm_19502017_v0.2 sol_bulk_density.wfraction_usda.3a1a1a_m_1k	https://land.copernicus.eu/glob al/products/fapar https://www.openlandmap.org/ Hengl et al. (2017)
24 25 26 27	Vegetation covariates veg_fapar_proba.v.annnual_d_1km_s00cm_20 142019_v1.0 Predicted soil properties sol_clay.wfraction_usda.3a1a1a_m_1km_b0_10 _30_60_100_200cm_ 19502017_v0.2 sol_sand.wfraction_usda.3a1a1a_m_1km_ b0_10_30_60_100_200cm_19502017_v0.2 sol_bulk_density.wfraction_usda.3a1a1a_m_1km_ b0_10_30_60_100_200cm_19502017_v0.2	https://land.copernicus.eu/glob al/products/fapar https://www.openlandmap.org/ Hengl et al. (2017)