

Predicting and forecasting root zone soil moisture with Random Forests

Coleen Carranza¹ and Martine van der Ploeg¹

¹Wageningen University

November 25, 2022

Abstract

The importance of soil moisture is recognized globally because it controls hydrological processes that are relevant to agriculture and climate studies. Currently, estimation of root zone soil moisture is largely accomplished using physical models, which are based on flow and transport equations. However, with the complexity of the processes operating in the vadose zone as well as their interactions with each other, parameterizing all the relevant processes is quite a challenge. This complexity is further enhanced by spatio-temporal variability in soil and vegetation properties which demand model parameters to be dynamic. Alternatively, purely data-based methods for root zone soil moisture estimation are still limited despite the growing availability of datasets from networks established within the last decade. Currently, these datasets are used largely for calibration and validation of physical models and retrieval methods from satellites. In this study, we explored the utility of Random Forest (RF) as an approach for predicting and forecasting daily root zone soil moisture from selected stations in the Raam and Twente network. We trained a single RF using meteorological datasets, soil type, land cover type, and LAI as predictor variables. The model was also tuned in order to obtain the optimal hyperparameters (mtry and ntree) and number of training samples. A comparison with model simulation results using Hydrus-1D was also performed. Our results show that RF can accurately predict and forecast root zone soil moisture at the study sites based on RMSE of 0.02 – 0.12 m³m⁻³, in comparison with Hydrus-1D simulations having RMSE of 0.05-0.22 m³m⁻³. However, poor results were obtained for saturated water conditions. In addition, 5-95% RF prediction intervals become wider at saturated water conditions for some sites, which indicates higher prediction and forecast uncertainties. RF can be used for root zone soil moisture estimation, especially at data poor regions where information on soil hydraulic parameters are sparse or lacking. It can also be used for estimating missing values at gaps in time series datasets.

Predicting and forecasting root zone soil moisture with Random Forests

Coleen Carranza¹ and Martine van der Ploeg²

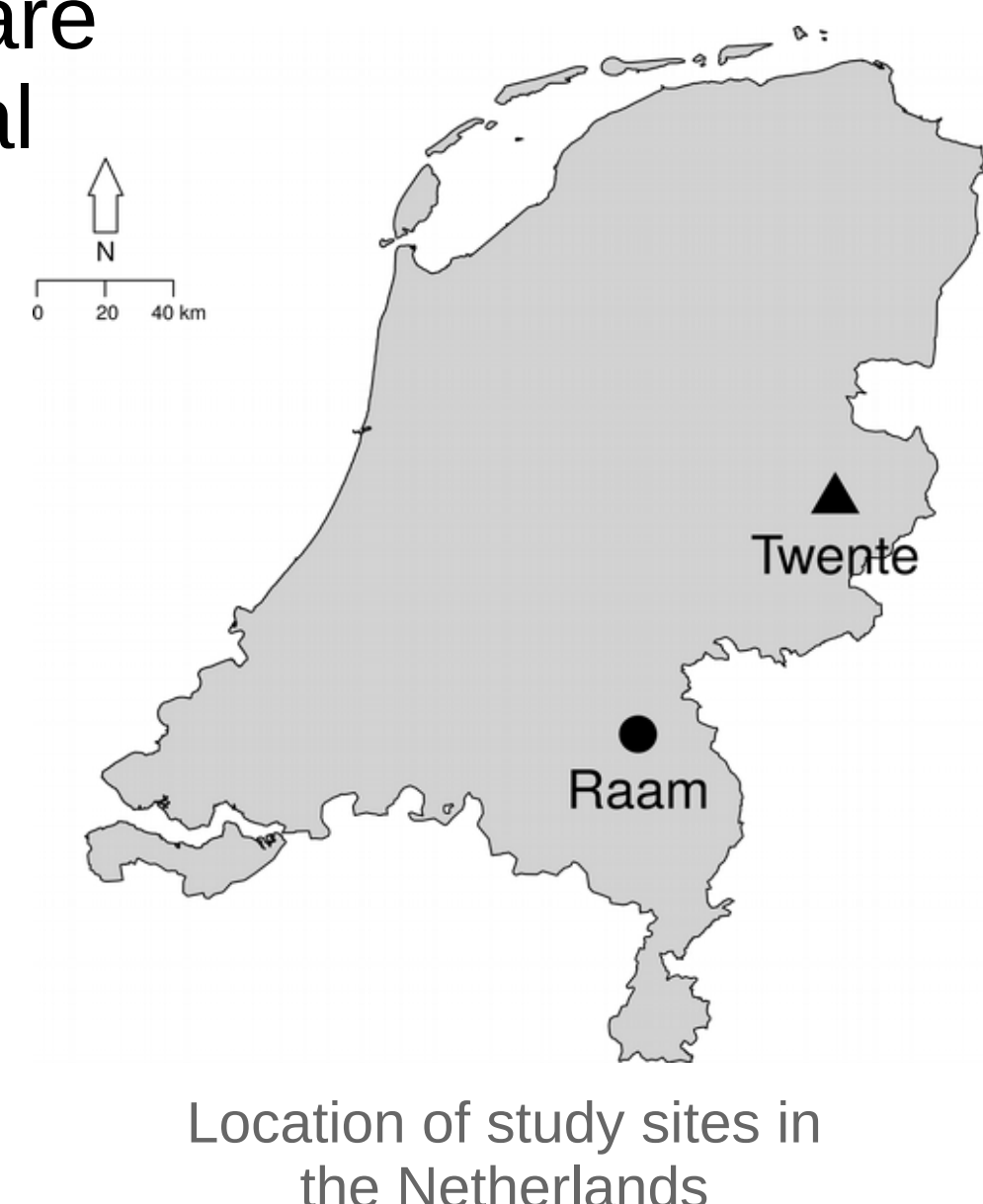
¹Soil Physics and Land Management Group, Wageningen University, Netherlands

²Hydrology and Quantitative Water Management Group, Wageningen University, Netherlands

1. Introduction

- Root zone soil moisture (θ_{rz}) is important in the hydrologic cycle and has been used for drought monitoring, water storage estimation and carbon cycle monitoring;
- Using physical hydrological models are the most common approach for estimation of θ_{rz} ;
- Soil moisture datasets from different monitoring sites are continuously increasing and becoming more available (e.g. International Soil Moisture Network (ISMN));
- There is an opportunity to test data-driven methods such as Random Forests to estimate real-world θ_{rz} conditions which are currently far less common than physical models.

Objective: Perform Random Forest (RF) to predict and forecast RZSM and compare results with simulations from a physical hydrological model (Hydrus-1D).

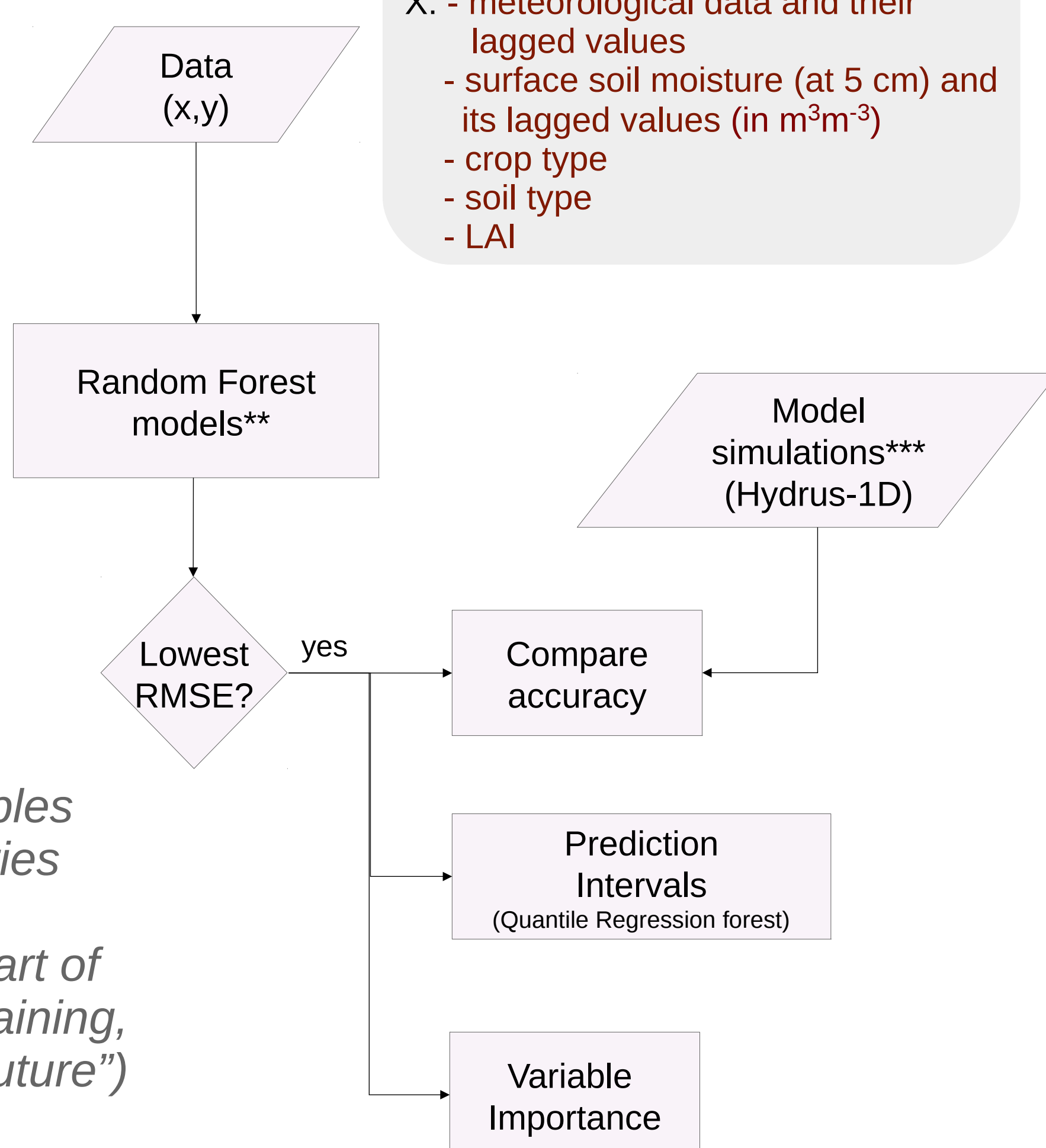


2. Materials and Methods

Study area: 4 stations within existing soil moisture networks (1 in Raam and 3 in Twente) which are closest to meteorological stations of KNMI (Royal Dutch Meteorological Institute).

Y: Depth-average θ_{rz} up to 40 cm (in m^3m^{-3})
X: - meteorological data and their lagged values
- surface soil moisture (at 5 cm) and its lagged values (in m^3m^{-3})
- crop type
- soil type
- LAI

Combinations of parameters:
mtry: 3,4,5,6,7
ntree: 500,600,700,800, 900,1000
training prop*: 0.5, 0.6, 0.7, 0.8

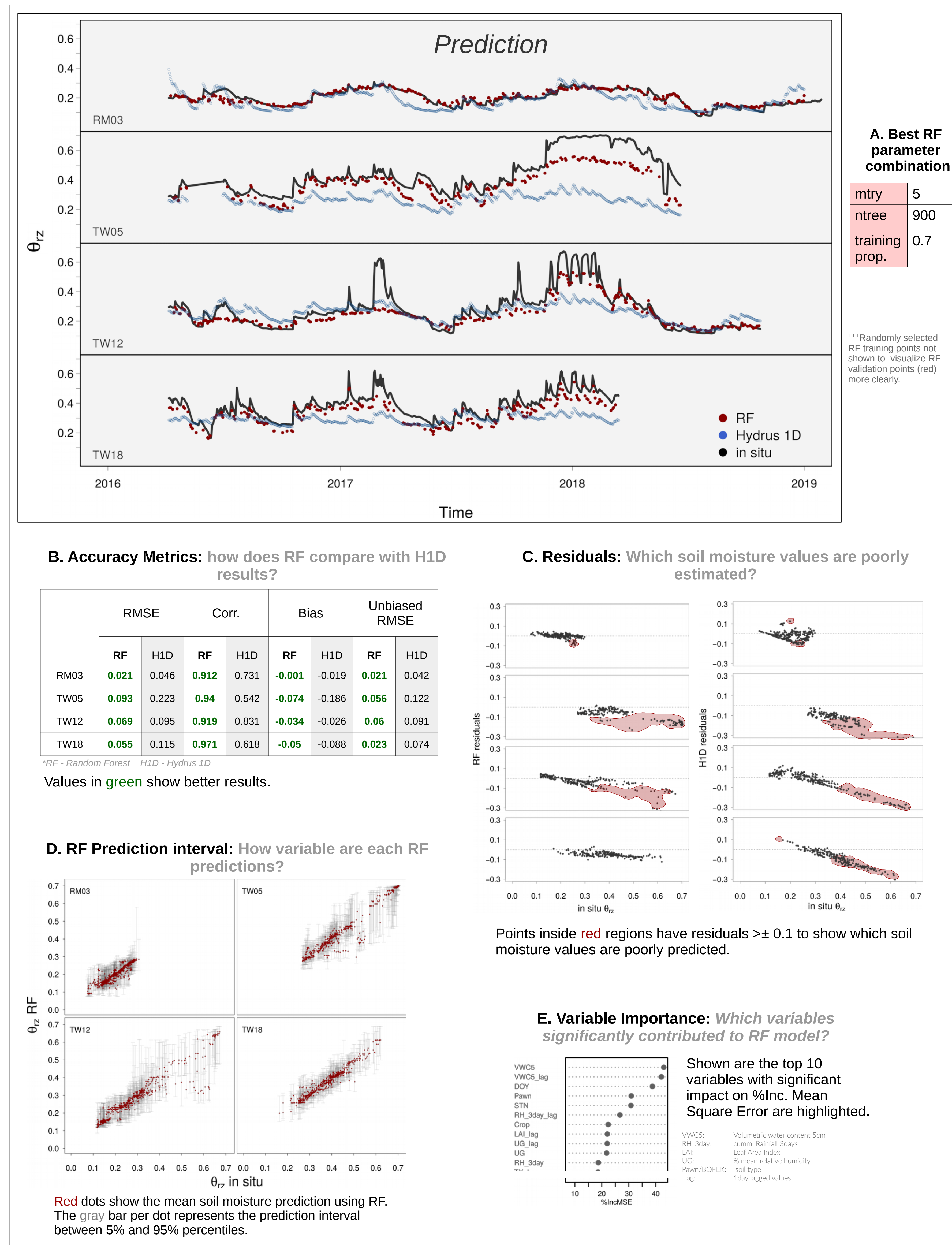


*prediction: random samples from time series data
forecasting: subset first part of dataset for training, remaining ("future") for validation

** 120 sets of parameter combinations (based on list above) were used for training RF models. Samples from each site were combined train a single RF model

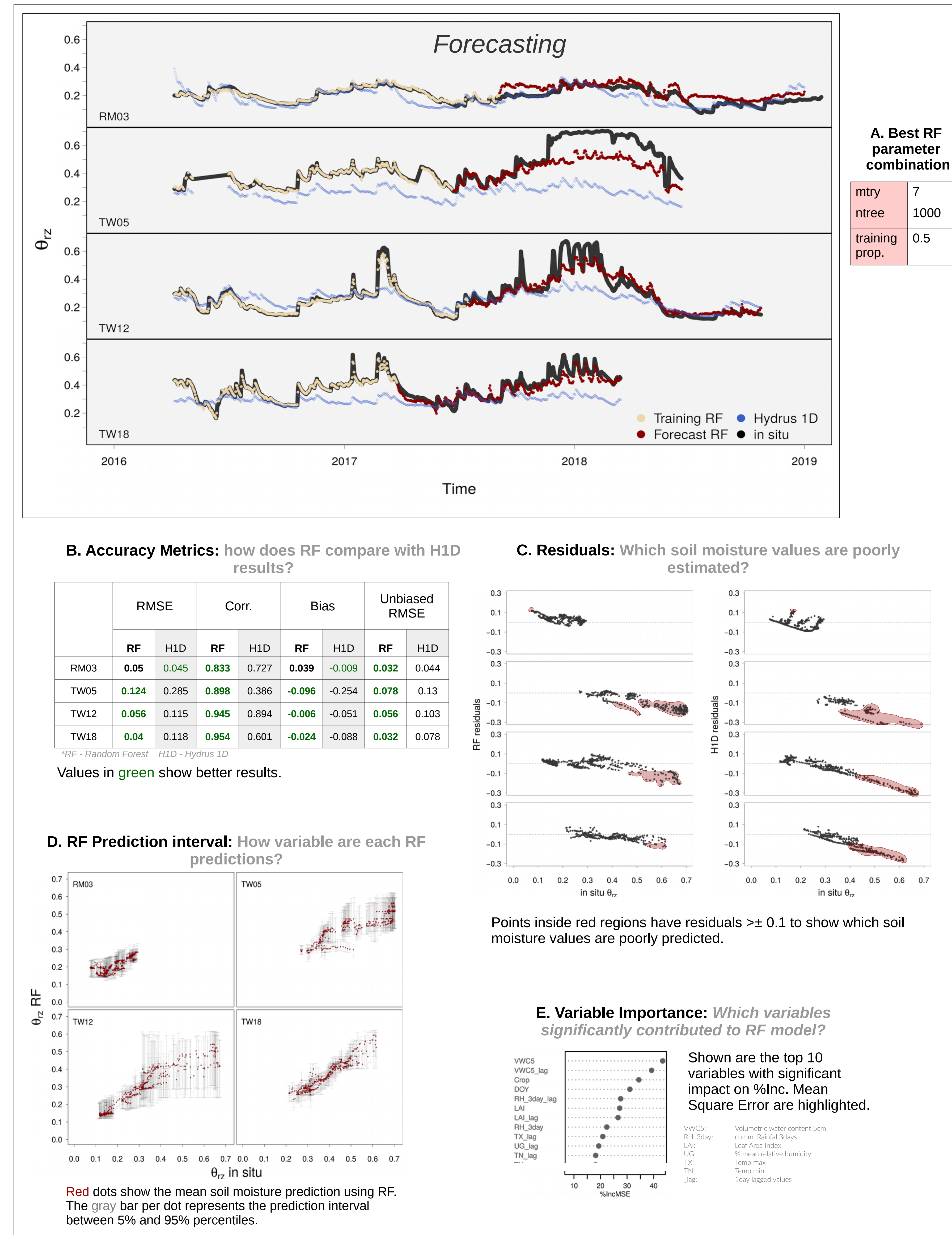
*** pore flow simulations with Hydrus-1D utilized the same independent variables (X) as RF. Additional inputs to perform the simulations were: soil hydraulic parameters, root water uptake function, site characteristics (e.g. Lat, Long) of meteorological stations

3. Results



4. Discussion

- In general, RF predictions showed better results than H1D simulations (B), for both prediction and forecasting.
- Both H1D and RF underestimate very high soil moisture values (C). Overestimation of very low values is also observed, but is not as common
- For RF, variables which represent processes that promote occurrence of very high soil moisture should be included (e.g. likelihood of macropore occurrence, bulk density).
- Range of prediction intervals appear to be site-specific (D). However, for a couple of sites (TW05 and TW12), larger intervals are observed for high soil moisture content.
- Important variables (E) show combination of meteorological conditions, vegetation and, soil properties are necessary for accurate prediction. Lagged values appear to be as important as current values.



Why choose RF (or another data-driven method) over a physical model?

- Results from the study sites show capability of RF to accurately estimate θ_{rz} . It even surpassed the accuracy of Hydrus-1D estimates using a pore flow model.
- Soil hydraulic parameters not required to run RF compared to physical models, which means they can be especially applicable in areas when these are not available. Prediction method also can be used to fill data gaps in soil moisture time series
- When the main objective is to estimate soil moisture states, RF can do the task. However, if processes that control soil moisture states are also sought, physical models should be applied. Although for RF, a glimpse into these processes are given by the list of important variables.

Acknowledgements

This study is part of the project entitled Operational Water Management using Sentinel-1 Satellites (OWAS15). The project is funded by Toegepaste en Technische Wetenschappen (TTW) which is part of the Netherlands Organization for Scientific Research (NWO).