**Reply on RC1**

General comments: This study proposes a comprehensive data-driven framework for selecting the optimal observing operations (data-worth analysis) and updating the predictions for soil moisture dynamics. The fully data-driven approach provides a complement to physics-based models, especially for complex real-world scenarios. While the quality of the manuscript is good, there are still some issues that require clarification.

Specific comments:

1. A major concern is the conclusions drawn from applying the Gaussian processes and EnKF assimilation techniques. While efficient and simple to implement, these methods have inherent limitations such as excessively smooth predictions (GP) and optimality only for Gaussian linear problems (EnKF). As the soil moisture dynamics are not fully met by these assumptions, the proposed method may experience difficulties, such as the mentioned localized surges. Therefore, some conclusions "high-quality and small data may be better than unfiltered big data" and "the soil water content in the middle layer exhibits remarkable superiority in comparison to the surface with its highest-level variability" may be case-specific rather than generalizable. It is important to consider other data-driven and assimilation methods, such as deep neural networks, particle filtering, and MCMC, leading to potentially different outcomes. I would like to see some clarifications regarding this issue.

**Answer:**

Thank you for your constructive comments. We have accepted your suggestions and evaluated a new NP-DWA framework where EnKF is replaced by particle filtering (PF). **Fig. S1** depicts the expected data-worth of potential observations of $\theta_S$, $\theta_M$, and $\theta_D$ regarding the retrieval of $\theta_{0.30}^{ave}$, $\theta_{0.60}^{ave}$, and $\theta_{1.00}^{ave}$, respectively. A comparison of **Fig. S1** and **Fig. 4** reveals that the spatio-temporal changes of expected data-worth under these two assimilation methods are remarkably similar. This demonstrates the generalizability of our proposed framework and related conclusions under different data assimilation schemes. To avoid duplication of research, we are sorry that we finally decided not to add the results of PF in the main text, but rather to include them as supplementary material in the revised manuscript (please see Lines 165-170 and Supplement).

In addition, we also tested two other NP-DWA frameworks where GP was replaced by support vector machines (SVM) and random forests (RF), respectively. The temporal changes of expected data-worth metrics are depicted in **Fig. S2**. Only the results at DAHRA are presented here. A comparison of **Fig. S2** and **Fig. 4** indicates that although the magnitude and trends of data-worth vary slightly across different machine learning methods, the selection of the optimal monitoring depths for specific targets is quite consistent. For example, the optimal observation depth shifted as the prediction target varied, and soil water content in the middle layer robustly exhibited remarkable superiority in the construction of model-free soil moisture models. Moreover, the performance comparison of various machine learning algorithms in reproducing soil moisture dynamics has been widely discussed in previous studies (Dubois et al., 2021; Liu et al., 2020; Prakash et al., 2018). In particular, the ability of GP to reproduce the nonlinearity of soil water problems has also been demonstrated in (He et al., 2023; Ju et al., 2018; Wang et al., 2021). Therefore, we finally decided only to include these results as supplementary material as well in the revised manuscript.
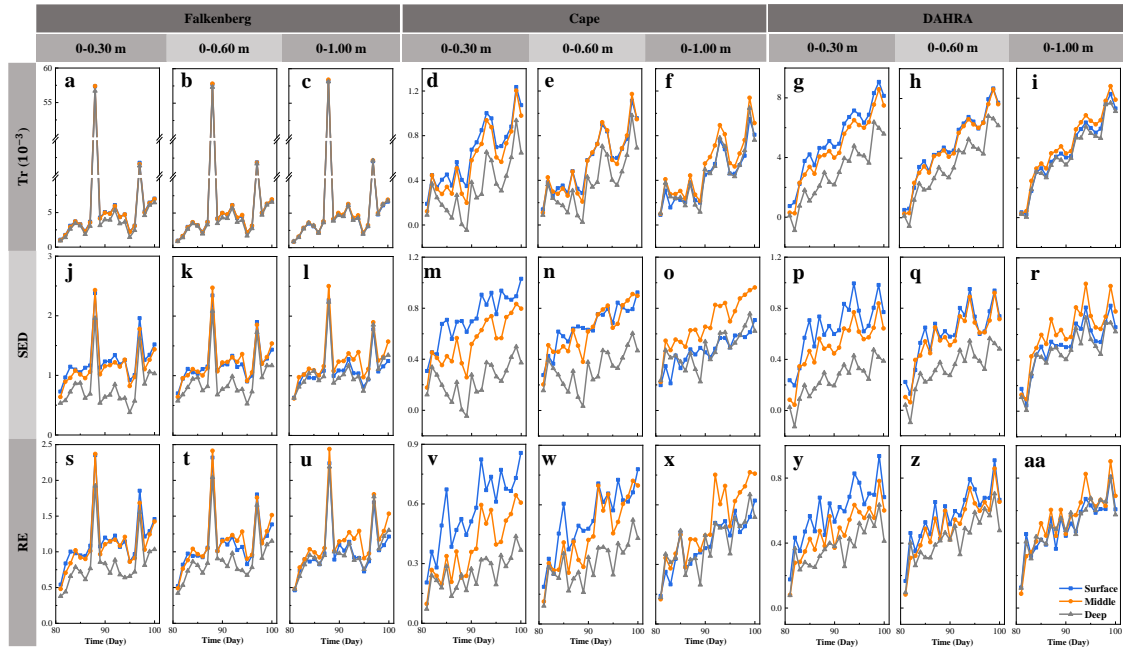
**Figure S1.** The expected data-worth of potential soil moisture observations in the surface, middle, and deep layers in the form of trace ($T_r$), Shannon entropy difference (*SED*), and relative entropy (*RE*), respectively, regarding the retrieval of average soil moisture in the top 0.30 m, 0.60 m, and 1.00 m at three sites, when EnKF is replaced by particle filtering (PF) in the proposed NP-DWA framework
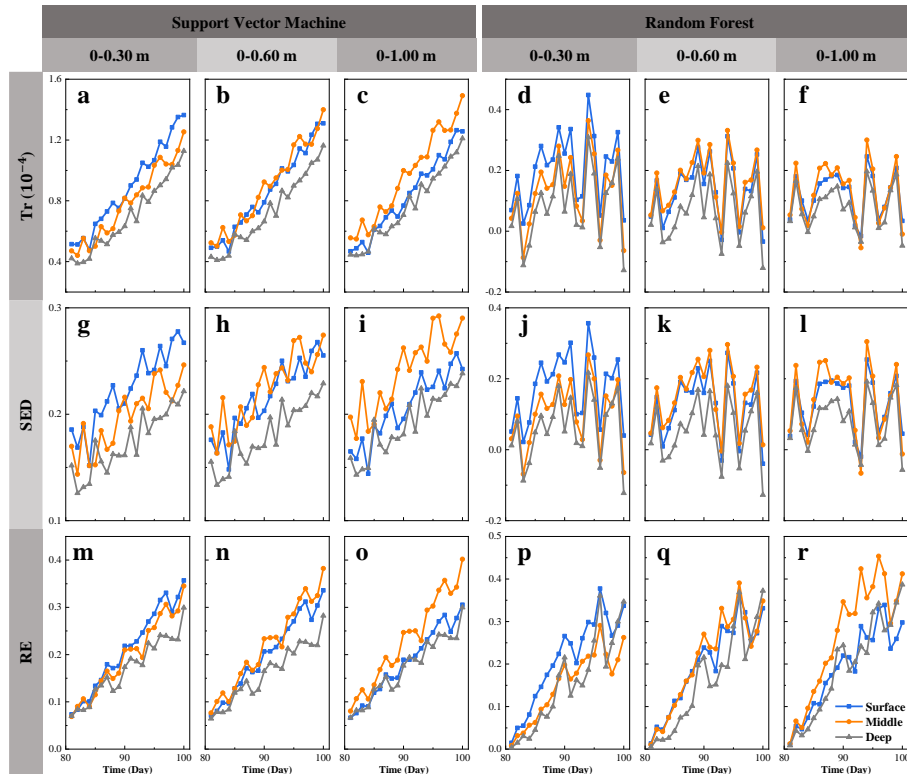


**Figure S2.** The expected data-worth of potential soil moisture observations in the surface, middle, and deep layers in the form of $T_r$, *SED*, and *RE* regarding the retrieval of average soil moisture in the top 0.30 m, 0.60 m, and 1.00 m at DAHRA site, when GP is replaced by support vector machine (SVM) and random forest (RF) in the proposed NP-DWA framework, respectively

**References:**

Dubois, A., Teytaud, F. and Verel, S., 2021. Short term soil moisture forecasts for potato crop farming: A machine learning approach. Computers and Electronics in Agriculture, 180: 105902.

He, L. et al., 2023. Physics-constrained Gaussian process regression for soil moisture dynamics. Journal of Hydrology, 616: 128779.

Ju, L., Zhang, J., Meng, L., Wu, L. and Zeng, L., 2018. An adaptive Gaussian process-based iterative ensemble smoother for data assimilation. Advances in water resources, 115: 125-135.

Liu, Y., Jing, W., Wang, Q. and Xia, X., 2020. Generating high-resolution daily soil moisture by using spatial downscaling techniques: A comparison of six machine learning algorithms. Advances in Water Resources, 141: 103601.

Prakash, S., Sharma, A. and Sahu, S.S., 2018. Soil moisture prediction using machine earning. IEEE, pp. 1-6.

Wang, Y. et al., 2021. A nonparametric sequential data assimilation scheme for soil moisture flow. Journal of Hydrology, 593: 125865.

2. It is recommended that the methodology section of this paper be better presented. Specifically, the problem setup for moisture prediction, an explicit list of the contents of vectors X and y should be provided prior to section 2.1. This will enable the reader to better understand the proposed data-driven framework.

**Answer:**

Thank you for your valuable suggestions. We have revised the methodology section. A clearer description of vectors $X$ and $y$ has also been added in section 2.1 of the revised manuscript (please see Lines155-160 and 170-185).

3. Some techniques have been proposed for better performance in nonlinear problems, e.g., restart, iterations. How will these techniques perform in NP-DWA?

**Answer:**

We thank the reviewer for the constructive comments. In fact, the procedure of constructing GP models in a sequential manner in our NP-DWA framework resembles a restart operation. At any time step $t=k$, the construction of the GP model does not solely rely on the information from the previous time step, instead, its training data includes all available soil moisture data from $t=1$ to $t=(k-1)$. This restart-like operation ensures that the training database is sequentially augmented to include more diverse training scenarios, so that actual observations can be accurately "captured" by the generated potential observation samples. Ultimately, the accuracy (or reliability) of our NP-DWA framework for data-worth assessment can be guaranteed. Related descriptions have been added in the revised manuscript (please see Lines 170-185). The performance improvements of these techniques such as restart and iterations for our NP-DWA will be explored in our future study.

4. L31:"An alternative monitoring strategy with a larger data-worth was prone to a higher DW assessment accuracy within the proposed NP-DWA framework" This sentence is meaningless and should be removed.

**Answer:**

We have accepted the reviewer's suggestion and deleted this sentence (please see Lines 35 and 665).

5. Please provide the dimensionality for all the involved vectors and matrices.

**Answer:**

We have added the dimensionality for all the involved vectors and matrices in Section 2 of the revised manuscript.