

## ©A Machine Learning Approach to Mitigate Ground Clutter Effects in the GPM Combined Radar–Radiometer Algorithm (CORRA) Precipitation Estimates

MIRCEA GRECU<sup>a,b</sup>, GERALD M. HEYMSFIELD,<sup>a</sup> STEPHEN NICHOLLS,<sup>a,c</sup> STEPHEN LANG,<sup>a,c</sup>  
AND WILLIAM S. OLSON<sup>a,d</sup>

<sup>a</sup> NASA Goddard Space Flight Center, Greenbelt, Maryland

<sup>b</sup> Morgan State University, Baltimore, Maryland

<sup>c</sup> Science Systems and Applications, Inc., Greenbelt, Maryland

<sup>d</sup> University of Maryland, Baltimore County, Baltimore, Maryland

(Manuscript received 25 April 2024, in final form 27 August 2024, accepted 4 November 2024)

**ABSTRACT:** In this study, a machine learning–based methodology is developed to mitigate the effects of ground clutter on precipitation estimates from the Global Precipitation Measurement Combined Radar–Radiometer Algorithm. Ground clutter can corrupt and obscure the precipitation echo in radar observations, leading to inaccuracies in precipitation estimates. To improve upon previous work, this study introduces a general machine learning (ML) approach that enables a systematic investigation and a better understanding of uncertainties in clutter mitigation. To allow for a less restrictive exploration of conditional relations between precipitation above the lowest clutter-free bin and surface precipitation, reflectivity observations above the clutter are included in a fixed-size set of predictors along with the precipitation type, surface type, and freezing level to estimate surface precipitation rates, and several ML-based estimation methods are investigated. A neural network (NN) model is ultimately identified as the best candidate for systematic evaluations, as it is computationally fast to apply while effective in applications. The NN provides unbiased estimates; however, it does not significantly outperform a simple bias correction approach in reducing random errors in the estimates. The similar performance of other ML approaches suggests that the NN’s limited improvement beyond bias removal is due to indeterminacies in the data rather than limitations in the ML approach itself.

**SIGNIFICANCE STATEMENT:** Ground clutter can obscure and corrupt the precipitation echo in the reflectivity observations by spaceborne radar, leading to inaccuracies and biases in the surface precipitation estimates. In this study, a machine learning approach is developed to mitigate the effects of ground clutter on precipitation estimates from the Global Precipitation Measurement (GPM) Combined Radar–Radiometer Algorithm (CORRA). The approach is shown to be effective in removing the biases associated with the simplest ground clutter mitigation approach and reducing the random errors associated with more complex climatologically based bias-removal approaches.

**KEYWORDS:** Algorithms; Radars/Radar observations; Machine learning

### 1. Introduction

In radar meteorology, the echo originating in power emitted by the radar and reflected by the ground is called ground clutter. Ground clutter has a negative impact on observations collected by dual-frequency precipitation radar (DPR) of the NASA Global Precipitation Measurement (GPM) mission (Skofronick-Jackson et al. 2017), as it may obscure or corrupt radar signal associated with precipitation. The extent of ground clutter in spaceborne radar observations increases with incidence angle (Kubota et al. 2016). Shown in Fig. 1 is a single scan representation of the Ku-band reflectivity observed by the GPM DPR for orbit 50853 on 9 February 2023. The enhanced reflectivity values at ranges close to (and larger than) 170 are contaminated by ground clutter. The lowest bins that are deemed clutter-free by the DPR algorithm

(Iguchi et al. 2021) are indicated by the black line in the plot. As apparent from the figure, the number of bins affected by clutter is quite significant for observations near the edge of the swath. Relative to the Earth ellipsoid, 20 range bins (bin width of 125 m) affected by clutter at the maximum scan incidence angle of 17° are equivalent to a clutter height of about 2.4 km at the edge of the DPR swath. While the assumption that the precipitation flux does not change significantly with height may be reasonable in some situations, it is likely to result in significant biases in the surface precipitation estimates in weather systems with freezing levels close to the ground. This is because ice processes such as riming and depositional growth can result in significant flux changes and reflected power changes.

To mitigate such biases, statistical correction methodologies, akin to those used to estimate the surface reflectivity from ground-based precipitation radar observations, may be used. Specifically, in ground-based radar, as the height of horizontally scanning radar beams increases with range, the lowest-elevation reflectivity observations may be significantly elevated above the ground at large ranges. For example, the beam center height is about 1500 m for an elevation angle of

<sup>a</sup> Denotes content that is immediately available upon publication as open access.

Corresponding author: Mircea Grecu, mircea.grecu-1@nasa.gov

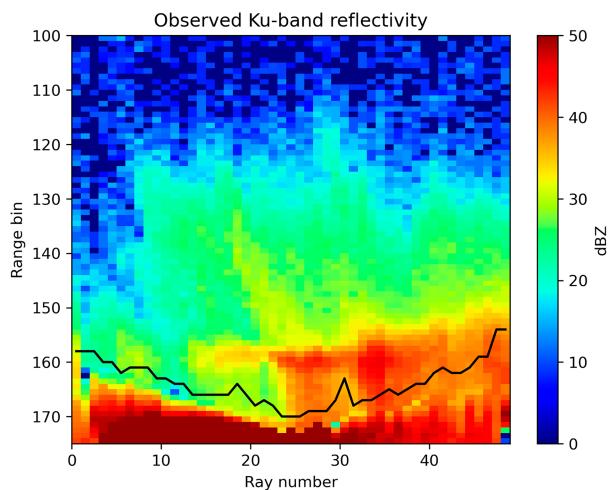


FIG. 1. Cross-track section of observed Ku-band reflectivity field for orbit (50853) on 9 Feb 2023. Range bin spacing is 125 m, and the black line indicates the LCFB. Bin 175 corresponds to the Earth ellipsoid.

0.5° and a range of 100.0 km (NWS 2023). Traditionally, to estimate the surface reflectivity from the lowest elevation angle reflectivity observations of ground radars, short-range reflectivity observations from multiple elevation angle scans are used to derive statistical relationships between surface reflectivities and reflectivities aloft (Koistinen 1991). A similar approach can be applied to mitigate ground clutter in spaceborne radar observations, with the difference being that the relationships between surface precipitation rates and precipitation rates aloft are derived from near-nadir spaceborne radar observations and associated precipitation estimates that are minimally impacted by ground clutter. This approach has already been applied by Hirose et al. (2021) to refine the GPM DPR surface precipitation estimates. In this study, we present a machine learning (ML)-based methodology to investigate and mitigate ground clutter effects on precipitation estimates from the GPM Combined Radar–Radiometer Algorithm (CORRA). To facilitate understanding, a list of acronyms used throughout this paper is provided in the appendix. While conceptually similar to the approach of Hirose et al. (2021), our methodology is different in several key aspects and provides additional insight into ground-clutter-related uncertainties in the surface precipitation estimates and the best strategies to mitigate them. Specifically, our study focuses on investigating the fundamental benefits and limitations of ground clutter mitigation methodologies using profile-level information. Moreover, unlike Hirose et al. (2021), we use reflectivity profile observations (rather than relying solely on profiles of estimated precipitation rates) in the derivation of relationships between the precipitation rate in the lowest clutter-free bin (LCFB) and the surface precipitation rate. The benefit of using reflectivity along with the LCFB precipitation rate estimate rather than exclusively precipitation profiles is that it facilitates the development of more physically consistent estimates. That is, radar profiling algorithms (Iguchi et al.

2021; Grecu et al. 2016) require assumptions regarding precipitation structure in the clutter to accurately incorporate estimates of the path-integrated attenuation (PIA) from the surface reference technique (SRT) to correct for attenuation down to the surface. However, if the clutter mitigation technique requires precipitation estimates, it can only be applied after the radar estimation process is complete. This may result in inconsistencies between the assumptions regarding the attenuation due to precipitation in the clutter and the actual precipitation estimates. While such inconsistencies may be addressed through iterative procedures, they result in a more computationally intensive retrieval process. In contrast, a clutter mitigation technique that uses reflectivity observations directly to derive relations between information above the clutter and precipitation in the clutter can be explicitly incorporated into the attenuation correction and precipitation estimation process, and this eliminates the need for iterative procedures to ensure the consistency of results. It should be mentioned, however, that the benefit (if any) of estimating the reflectivity in the clutter is limited in deep convection because there are large uncertainties in the attenuation correction process both above and in the clutter. In this case, additional uncertainties caused by physical inconsistencies may not matter. Another distinction relative to Hirose et al. (2021) is that our methodology is based on ML, which is beneficial from the feature engineering perspective (Zheng and Casari 2018). Specifically, machine learning models can effectively extract relevant information from the data without having to resort to the explicit identification of features (defined as numerical attributes uniquely derived through a computational procedure applied to input data), thereby reducing the need for manual feature engineering, which can be time consuming and error prone for human experts. For example, the precipitation gradient with respect to radar range is an intuitively derived feature in the surface precipitation estimation approach of Hirose et al. (2021). While features that make intuitive sense are valuable, questions regarding their optimality are difficult to objectively address without tedious investigation. From this perspective, ML procedures that do not require explicit features are worth considering. Additionally, it is important to note that ML has been successfully applied to a variety of problems involving the estimation of precipitation from spaceborne and ground-based radar and radiometer observations (Chase et al. 2021; King et al. 2022; Rahimi et al. 2024).

The paper is organized as follows. In section 2, we present the ML methodology used to estimate the surface precipitation rate from reflectivity observations not affected by clutter as well as additional information such as the precipitation and surface type and the 0°C isotherm height. In section 3, we present the results of the application of the ML methodology to the GPM CORRA precipitation estimates. In section 4, we offer some conclusions from the study.

## 2. Methodology

### a. General considerations

The simplest method to estimate the precipitation rate at a given height above the Earth ellipsoid [and for a given

### Stratiform precipitation over ocean

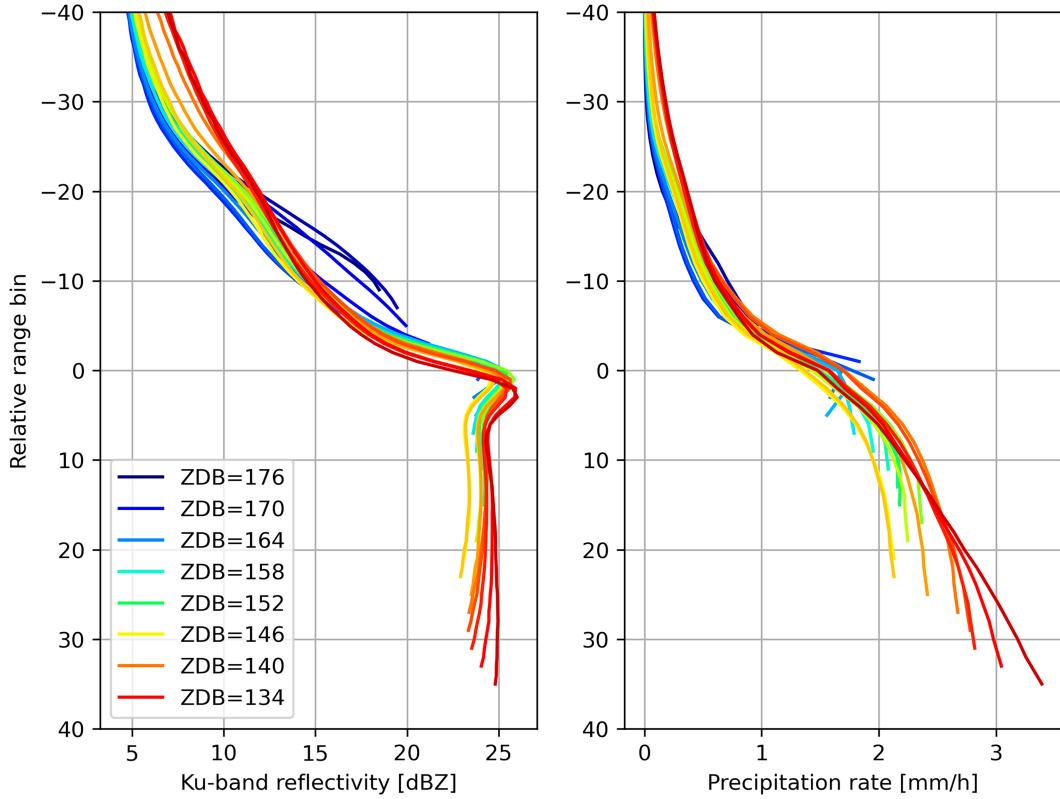


FIG. 2. Conditional mean reflectivity and precipitation rate profiles over global oceans for stratiform precipitation with various FLHs. The mean profiles were derived from 1 year (2018) of data characterized by fewer than eight bins affected by clutter and calculated conditionally on the location of the radar bin associated with the 0°C isotherm relative to the Earth ellipsoid.

precipitation type (PT), surface type (ST), and freezing level (FL)] from a precipitation rate at a higher level is to rescale the higher-level value by the ratio of the climatological mean precipitation rates at the two levels. Mathematically, this may be written as

$$P_{\text{rate}}(H_1, \text{PT}, \text{ST}, \text{FL}) = P_{\text{rate}}(H_2, \text{PT}, \text{ST}, \text{FL}) \\ \times \frac{\langle P_{\text{rate}}(H_1, \text{PT}, \text{ST}, \text{FL}) \rangle}{\langle P_{\text{rate}}(H_2, \text{PT}, \text{ST}, \text{FL}) \rangle}, \quad (1)$$

where  $P_{\text{rate}}$  is the precipitation rate;  $H_1$  is the height where the estimate is needed, but for which no direct measurement is available;  $H_2$  is LCFB ( $H_2 > H_1$ ) where a radar measurement is available; and operator  $\langle \rangle$  denotes the climatological mean over a large dataset characterized by the same freezing level, surface, and precipitation type.

While simple in form, the challenge in applying a clutter correction methodology based on Eq. (1) is the derivation of the correction factors  $\langle P_{\text{rate}}(H_1) \rangle / \langle P_{\text{rate}}(H_2) \rangle$  for all possible  $(H_1, H_2)$  pairs. Nevertheless, because the ground clutter depth is a function of the scanning incidence angle, estimates of the climatological correction factor derived from near-nadir

reflectivity observations and precipitation estimates may be used to mitigate the clutter near the edges of the swath. Shown in Fig. 2 are overoceans conditional mean reflectivity and precipitation rate profiles from the GPM CORRA algorithm (Grecu et al. 2016) for stratiform precipitation with various freezing-level heights (FLHs). The profiles are plotted relative to the 0°C bin to emphasize similarities rather than differences due to temperature-dependent processes. One year's worth (i.e., 2018) of DPR observations and associated GPM CORRA retrievals characterized by fewer than eight bins affected by clutter are selected and used in calculations of the mean profiles. The data are partitioned based on the FLH in 12 distinct subsets, with the FLHs of each subset within 125 m of  $1.875 + k \times 0.25$  km with  $k$  varying from 0 to 11, resulting in 12 conditional mean profiles. As shown in the figure, the mean reflectivity and the associated precipitation profiles tend to align with one another. This behavior may be used to mitigate the impact of clutter, even in near-nadir DPR observations that are affected by clutter at relatively low altitudes that make direct precipitation rate estimation at or near the surface impossible. Specifically, the data in Fig. 2 suggest that

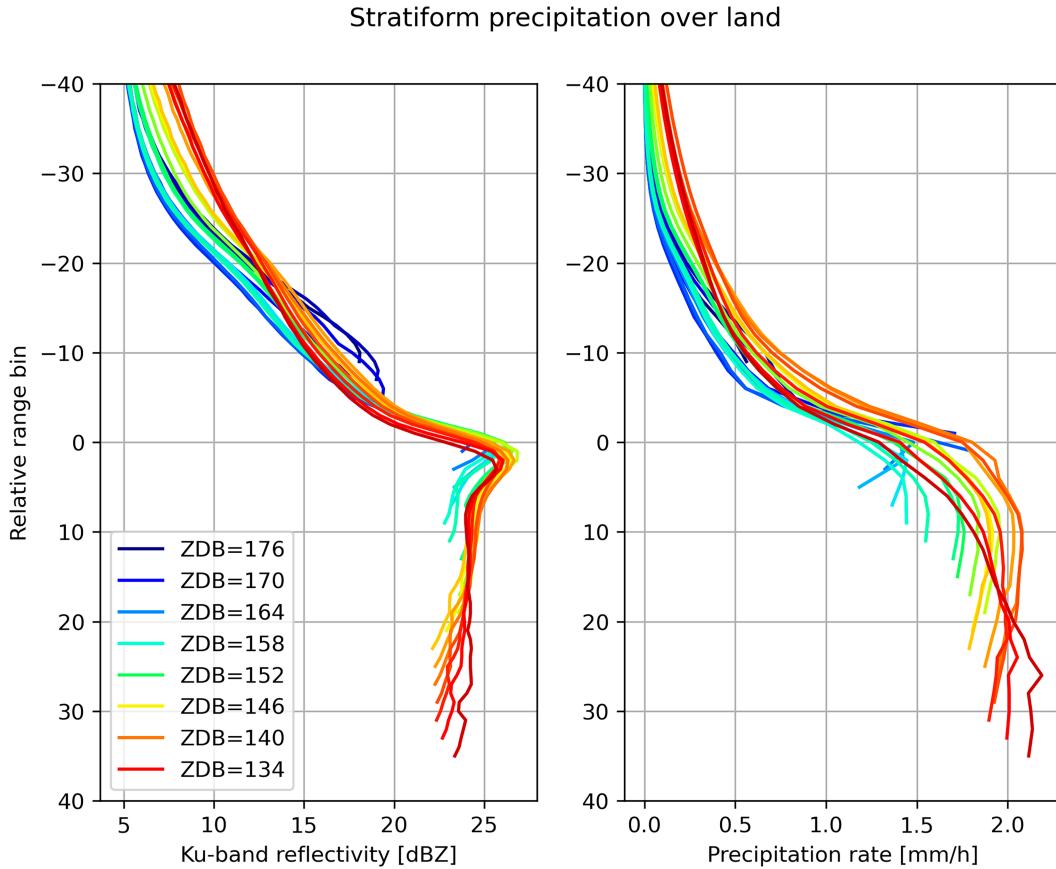


FIG. 3. As in Fig. 2, but over land.

$$\frac{\langle P_{\text{rate}}(H_1, \text{PT}, \text{ST}, \text{FL}) \rangle}{\langle P_{\text{rate}}(H_2, \text{PT}, \text{ST}, \text{FL}) \rangle} \approx \frac{\langle P_{\text{rate}}(H_1 + d\text{FL}, \text{PT}, \text{ST}, \text{FL} + d\text{FL}) \rangle}{\langle P_{\text{rate}}(H_2 + d\text{FL}, \text{PT}, \text{ST}, \text{FL} + d\text{FL}) \rangle}, \quad (2)$$

where  $d\text{FL}$  is the difference between two distinct FLHs. The validity of Eq. (2) is supported by the fact that in plots relative to the  $0^{\circ}\text{C}$  isotherm, the conditional mean precipitation profiles in Fig. 2 look very similar to profiles characterized by higher FLH and extending to greater depths below the  $0^{\circ}\text{C}$  level. Here, the conditional mean precipitation rate refers to the mean precipitation rate in situations where (nonzero) precipitation is occurring in the LCFB. The selection of profiles (with a maximum of eight radar bins impacted by ground clutter) results in a minimum value of  $H_1$  of 1000 m (for a climatology derived from profiles with at most six bins affected by clutter). However, one can use Eq. (2) to approximate  $\langle P_{\text{rate}}(0, \text{PT}, \text{ST}, \text{FLH}) \rangle / \langle P_{\text{rate}}(H_2, \text{PT}, \text{ST}, \text{FLH}) \rangle$  as  $\langle P_{\text{rate}}(1000\text{ m}, \text{PT}, \text{ST}, \text{FLH} + 1000\text{ m}) \rangle / \langle P_{\text{rate}}(H_2 + 1000\text{ m}, \text{PT}, \text{ST}, \text{FLH} + 1000\text{ m}) \rangle$ . Alternatively, microphysical models that incorporate information such as relative humidity from numerical weather prediction (NWP) products may be used to estimate the precipitation rate at the surface from the precipitation rate at 1000 m above the Earth ellipsoid. Microphysical models combined with hydrodynamic models can potentially quantify orographic effects that are not detectable in spaceborne

radar products with limited clutter extent, as such products are preponderantly derived from observations over relatively flat terrain. It is important to note, however, that the use of such models may introduce biases due to uncertainties in the microphysical parameterizations and the NWP products. The use of Eq. (2) or ML models is a more straightforward approach that does not require the explicit use of microphysical and hydrodynamic models.

Shown in Fig. 3 are conditional mean reflectivity and precipitation rate profiles from CORRA for stratiform precipitation with various freezing-level heights over land. The conditional mean precipitation profiles over land exhibit more variability than over oceans. However, this may be a consequence of precipitation retrieval artifacts rather than differences in temperature-dependent physical processes. Specifically, given that the SRT PIA estimates are noisier and less reliable over land, their impact on precipitation estimates may be less systematic, which could result in a larger spread of conditional mean estimates. Nevertheless, Eq. (2) is still a reasonable assumption.

The mean reflectivity profiles shown in Figs. 2 and 3 are stratified by precipitation type (stratiform), freezing level, and surface type only, but it is conceivable that features that further separate the relationships between the reflectivity observations and the final precipitation estimates exist. As previously mentioned, Hirose et al. (2021) use the precipitation slope to stratify

Stratiform precipitation over oceans  
Zero-degree bin=144

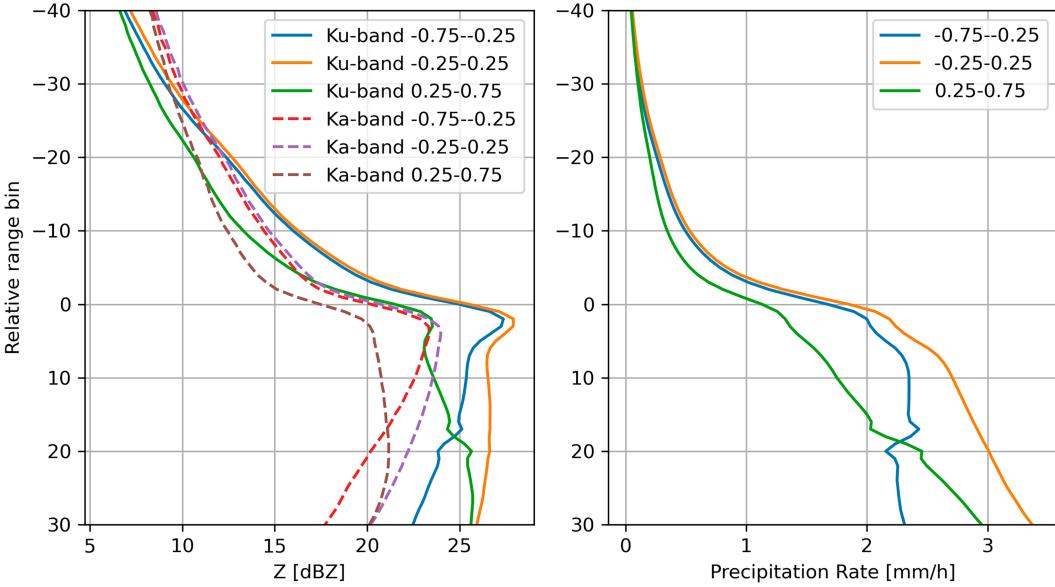


FIG. 4. As in Fig. 2, but for a ZDB of 144 and stratified by reflectivity slopes. The dashed lines in the left-hand side panel indicate the conditional mean reflectivity profiles at the Ka-band associated with the three classes.

the database of near-nadir precipitation supporting their precipitation refinement process. In the current study, we also investigate the slope of the reflectivity profile below the freezing level as a feature potentially useful for predicting the surface precipitation rate from the LCFB precipitation rate. Specifically, the slope of the Ku-band reflectivity observations in the first six radar bins below the brightband bottom (defined as in Iguchi et al. 2021) is used to stratify the observations into five categories. The resulting mean reflectivity profiles and the associated mean precipitation profiles are shown in Fig. 4 for three of these categories. The other two (i.e., associated with slopes with absolute values larger than 0.75 dB per bin) account for less than 10% of the total number of profiles. As seen in the figure, distinct mean reflectivity profiles result in distinct mean precipitation profiles. This behavior may be used to derive more accurate surface precipitation estimates than those derived from Eq. (1). However, to make effective use of the reflectivity slope and other such features, questions regarding the optimal strategy to calculate the slopes and partition them by value, especially when the ground clutter extends close to or above the freezing level, need to be addressed. As there is no obvious strategy to address such questions, we resort to an ML approach. The approach is applied to derive surface precipitation estimates for both stratiform and convective precipitation, applicable over land and oceans.

As previously mentioned, ML approaches do not require the explicit use of Eq. (1) and manually designed and optimized features. Instead, they require the organization of the dataset into a design matrix and a response matrix (Bishop 2006). In machine learning, the concepts of design and response matrices are borrowed from regression analysis, with

the design matrix representing the array of predictor variables while the response matrix representing the array of predicted variables. Each row of the design matrix corresponds to a single observation or data point, while each column represents a different predictor variable or feature.

In our study, the design matrix (i.e., input features) is an array of reflectivity observations and associated information, with each row containing the reflectivity values from a fixed-size portion of an observed profile. In addition to the reflectivity information, the zero-degree bin (ZDB), the position of the LCFB bin relative to the ZDB, the position of the surface relative to the ZDB, and the LCFB precipitation rate are included in the design matrix. To make the ML models computationally efficient, the number of reflectivity observations above the LCFB is set to 30. Larger numbers of reflectivity observations above the LCFB were tested but did not result in improved results. The response matrix is one-dimensional, i.e., a vector, and it contains the precipitation rates associated with the lowest bin in the training/evaluation dataset, which, as explained above, is eight bins above the Earth ellipsoid.

The structured organization of the dataset facilitates the exploration of multiple ML models with minimal effort, enabling the selection of the optimal one. While ML models are generally physics-agnostic in the sense that they do not explicitly make use of physical laws, they can exploit physical causality embedded in the dataset. If the slopes of the reflectivity profiles above the clutter are reliable predictors of the precipitation rate at the surface relative to the LCFB rate (as suggested by Fig. 4), then a machine learning model will be able to exploit this relationship. This is because similar reflectivity profiles in the design matrix are associated with similar slopes,

enabling the model to capture the underlying pattern effectively. However, some models may be potentially more accurate or computationally more efficient than others, and consequently, we consider multiple models from the scikit-learn library (Pedregosa et al. 2011). Details on the models considered and the strategy used to identify the best option are provided in the next subsection.

### b. Implementation details

A particular characteristic of our problem is that it involves ordinal variables—variables that are discrete in nature but have a natural order. Specifically, the ZDB (Iguchi et al. 2021) and the LCFB are ordinal variables, while precipitation and surface type are categorical nominal variables. One approach to handling these variables is to convert the nominal ones into multivariate binary variables through one-hot encoding (Hancock and Khoshgoftaar 2020) while treating the ordinal variables as continuous. Alternatively, an independent submodel can be developed for each combination of ordinal and nominal variables. Although the latter approach may seem computationally more expensive, it is straightforward and more interpretable, allowing performance to be analyzed conditionally based on precipitation and surface type, the ZDB, and clutter extent (expressed as the number of bins above the Earth ellipsoid). In this study, we follow the latter approach. To make the results more robust and limit the number of conditional models, we sort profiles into categories characterized by two consecutive ZDB values, e.g., 130 and 131, and 132 and 133. The ordinal values describing the clutter extent are preserved in their original resolution.

To determine the range of LCFB that needs to be considered in the model development, we analyze the cumulative distribution of a number of bins affected by clutter for rays in the DPR's outer swath shown in Fig. 5. As apparent in the figure, more than eight bins are affected by clutter for the vast majority of the DPR outer swath profiles over land, with about half of profiles characterized by more than 15 bins affected by clutter and slightly more than 10% of profiles characterized by more than 25 clutter-affected bins. To simulate clutter effects,  $n_c$  bins are assumed affected by clutter, where  $n_c$  varies between 8 and 25. Given the number of profiles with  $n_c > 25$  is relatively small, deriving ML models for  $n_c \leq 25$  and truncating values greater than 25 to 25 is not necessarily a poor choice.

The model definition and training procedure are summarized in Algorithm 1. Specifically, for each combination of surface type, precipitation type, ZDB, and assumed LCFB (which is  $n_c$  above the lowest bin in the training dataset or  $n_c + 8$  above the Earth ellipsoid), all corresponding reflectivity profiles in the training dataset are selected. A total of  $n_z = 30$  reflectivity observations above the clutter are used, along with the precipitation rate at the assumed LCFB, to define the model input. To enhance model interpretability, the reflectivity observations are first normalized using standard normalization. Reflectivities with values smaller than 12.0 dBZ (which is the noise level in the GPM DPR data) are set to 0.0 prior to the normalization. Principal component analysis

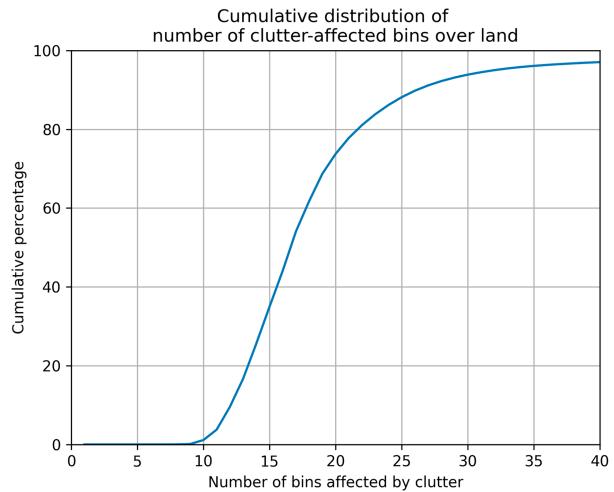


FIG. 5. The cumulative distribution of a number of bins affected by clutter for rays in the DPR's outer swath (defined as the portion of the swath within 12 rays from the edges) over land. The distribution is derived from 1 year (2018) of global DPR observations over land.

(PCA) is then applied to reduce the dimensionality of the input data. Five principal components explain more than 97.5% of the variance in the reflectivity observations, and therefore, we set the number of principal components to five. The PCA-transformed data are then concatenated with the normalized LCFB precipitation rate to form the input data for the ML model. The target data are the normalized precipitation rate at the lowest bin in the training dataset (i.e., eight bins above the Earth ellipsoid).

At the most granular level, the conditional model can utilize any regressor capable of predicting the target variable from the input data. In this study, we evaluated multiple regressors available in the scikit-learn library and ultimately selected a feedforward neural network (NN) with one hidden layer as the best candidate. The NN model is trained using the Adam optimizer (Kingma and Ba 2014) and the mean-square error (MSE) loss function. The hidden layer contains 32 neurons, utilizing the rectified linear unit (ReLU) activation function (Nair and Hinton 2010). The output layer is a linear layer. For hyperparameter optimization, we use a randomized search (Bergstra and Bengio 2012) with 100 iterations. It should be mentioned though that the performance of the NN model does not change significantly with the number of hidden neurons or the number of hidden layers.

Although more advanced models are accessible through the Skorch library (Mishra 2023), we found them unnecessary for our purposes. Conversely, we also considered simpler linear models, such as the ridge regression (Pedregosa et al. 2011). The overall performance differences between ridge regression and the NN were not significant. This suggests that the relationship between the precipitation rate at the LCFB and eight bins above the Earth ellipsoid is conditionally linear, with reflectivity information contributing only marginally to the prediction of the surface precipitation rate.

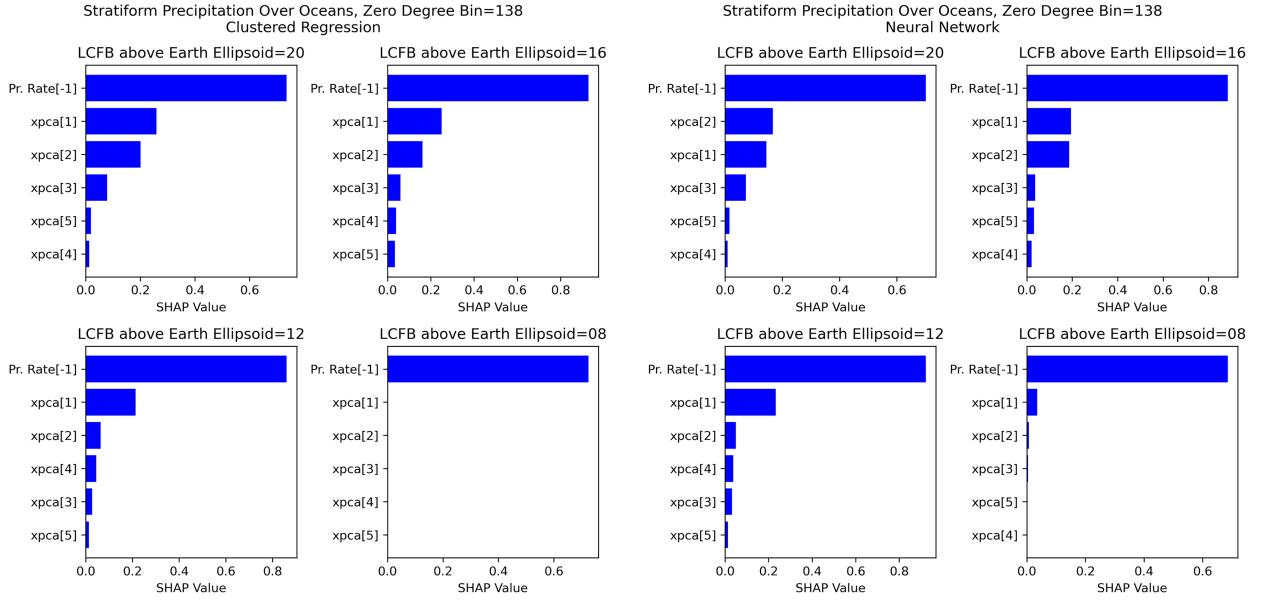


FIG. 6. SHAP analysis results showing feature importance for ridge regression and NN models.

**ALGORITHM 1: MODEL DEFINITION AND TRAINING**

**Input:** Reflectivity and precipitation dataset; initialized but not trained ML model  
**Output:** Trained ML model  
**for** surfacetype ∈ {land, ocean} **do**  
**for** precipstype ∈ {convective, stratiform} **do**  
**for** zdb ∈ {zdb categories} **do**  
**for** nc ← 1 **to** 18 **do**  
 $Z_{\text{train}} = Z[\text{zdb}][\text{nbins} - \text{nc} - \text{nz}: \text{nbins} - \text{nc}]$ ;  
 $P_{\text{rate},\text{train}} = P_{\text{rate},\text{ate}}[\text{zdb}][\text{nbins}]$ ;  
 $P_{\text{rate},\text{target},\text{train}} = \text{Precip}[\text{zdb}][\text{nbins}]$ ;  
Normalize  $Z_{\text{train}}$ ;  
 $X_{\text{PCA}} = \text{PCA}(Z_{\text{train},\text{norm}})$ ;  
Normalize  $X_{\text{PCA}}$ ;  
Normalize  $P_{\text{rate},\text{train}}$ ;  
Normalize  $P_{\text{rate},\text{target},\text{train}}$ ;  
 $X_{\text{train}} = \text{concentrate}(X_{\text{PCA},\text{norm}}, P_{\text{rate},\text{train},\text{norm}})$ ;  
 $\text{ML}[\text{zdb}][\text{nc}] = \text{train}(X_{\text{train}}, P_{\text{rate},\text{target},\text{train}})$ ;  
**return** ML;

To investigate this hypothesis, we performed a Shapley additive explanation (SHAP) analysis (Lundberg and Lee 2017) on both the ridge regression and the NN model. Results are shown in Fig. 6. The SHAP analysis reveals that the LCFB precipitation rate is the most important feature, with the contribution of the reflectivity information increasing with LCFB. The ridge regression and the NN model show similar feature importance, which is consistent with the similar performance of the two models. It should be mentioned that an evaluation using the original reflectivity observations, rather than the PCA-transformed data, also showed qualitative agreement between the two models. However, the models did not fully agree on ranking the importance of the reflectivity observations. Since this discrepancy was an artifact of the high correlation

among reflectivity observations, the PCA-transformed data were used in the final model.

To evaluate the performance of the ML models, we use a holdout validation approach. Specifically, the DPR dataset is split into a training and a testing dataset with 70% of profiles in the training dataset and the remaining 30% in the testing dataset. To ensure the independence of the training and testing datasets, a systematic splitting strategy is used. Specifically, blocks of 10 days are processed sequentially, and the first 7 days of data are assigned to the training dataset, while the remaining 3 days are assigned to the testing dataset. The process is repeated until all days in the dataset are assigned to either the training or testing dataset. As the number of profiles needed to train the ML models is large, we used 2 years' worth of data (i.e., 2018 and 2019) over oceans and 5 years of data over land (i.e., 2018–2022). The number of profiles in the training dataset as a function of the freezing level, precipitation, and surface type is shown in Fig. 7. The training dataset is used to optimize the ML models, while the testing dataset is used to evaluate them. The evaluation is based on calculations of the correlation coefficient and bias between the predicted and observed surface precipitation rates. The training dataset is used to optimize the ML models, while the testing dataset is used to evaluate them. The evaluation is based on calculations of the correlation coefficient and bias between the predicted and observed surface precipitation rates. Results are presented in the next section.

### 3. Results

The reason for considering several ML model architectures is to ensure that there is no latent information in the input data that are not properly captured. The inclusion of multiple ML models reduces the likelihood of such a possibility, as the models are based on different statistical modeling paradigms.

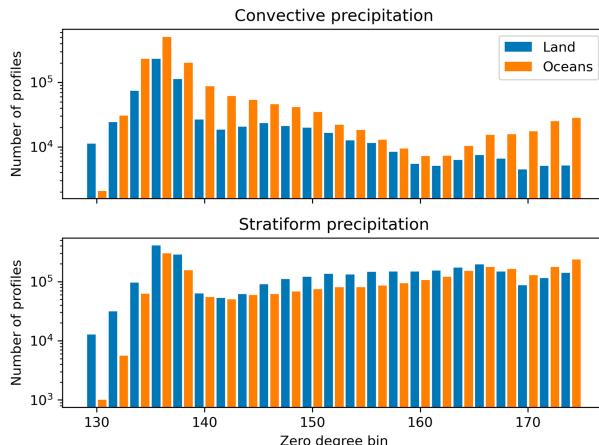


FIG. 7. Number of profiles in the training dataset as a function of FL, PT, and ST.

However, in our initial model testing, no particular ML model emerged as significantly better than the others. This outcome, which is not totally surprising, may be an indication that the relations between the surface precipitation rate and the precipitation rate at a given height above the surface depend on a multitude of factors that cannot be directly observed or do not have a clear signature in the reflectivity observations. Although incorporating more comprehensive microphysical and dynamic data derived from NWP models could potentially lead to more accurate ML models, doing so might introduce biases in the precipitation estimates due to the potentially biased representation of microphysical processes in the NWP models (Morrison et al. 2020). Additionally, developing more complex but robust ML models would require a large number

of simulations covering multiple regimes over both land and oceans. Given these considerations, we limit this investigation to models derived exclusively from GPM observations. Despite the similar performance across models, some are preferable to others. Therefore, based on this initial testing, we choose the NN as the best option, and instead of exploring additional methodologies or carrying out further tuning, we focus on characterizing its performance, especially in relation to a simple estimation methodology.

#### a. Stratiform precipitation over land

Before describing the performance of the different ML estimation methods, we will first examine the simplest solution as a benchmark. In this solution, the precipitation rate at the LCFB is assumed to be the same as the target precipitation rate (defined as the precipitation rate at 1.0 km above the Earth ellipsoid). As previously discussed, when the surface is below 1.0 km from the Earth ellipsoid, estimating the surface precipitation rate from its 1.0-km-level value using the relation in Eq. (2) introduces additional uncertainties that are not quantified in the study due to the lack of evaluation data.

Shown in the left-hand side panel of Fig. 8 is the correlation coefficient between the precipitation rate at the lowest level in the evaluation dataset and the LCFB precipitation rate up to 25 bins above the Earth ellipsoid for stratiform precipitation events over land. The vertical axis is the LCFB relative to the Earth ellipsoid, and the horizontal axis represents the ZDB. As seen in the panel, the correlation decreases with the position of the LCFB above the surface. The bins marking a more significant correlation decrease generally occur in the mixed and ice phases. The biases associated with the LCFB-derived precipitation rate relative to the target precipitation rate are shown in the right-hand side panel of Fig. 8. This type

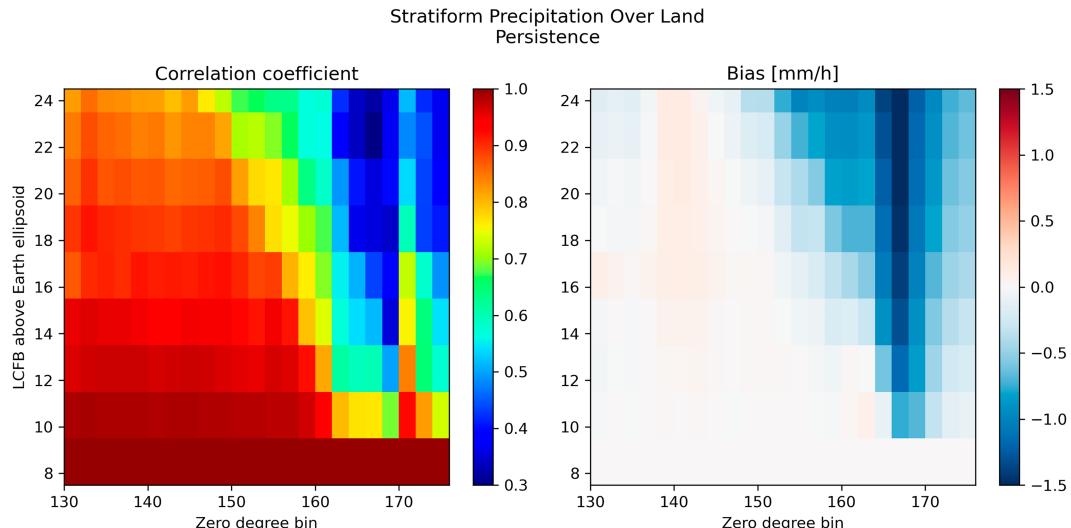


FIG. 8. Performance of a persistence-based clutter mitigation method for stratiform precipitation over land. (left) The correlation coefficient of the target precipitation rates (1.0 km above the Earth ellipsoid) with the precipitation rates in the LCFB (which serves as the target estimate in the persistence-based scheme). (right) Mean differences between the LCFB precipitation rates and the target precipitate rates. Values are plotted for different LCFBs (vertical axis) and for different ZDBs (horizontal axis).

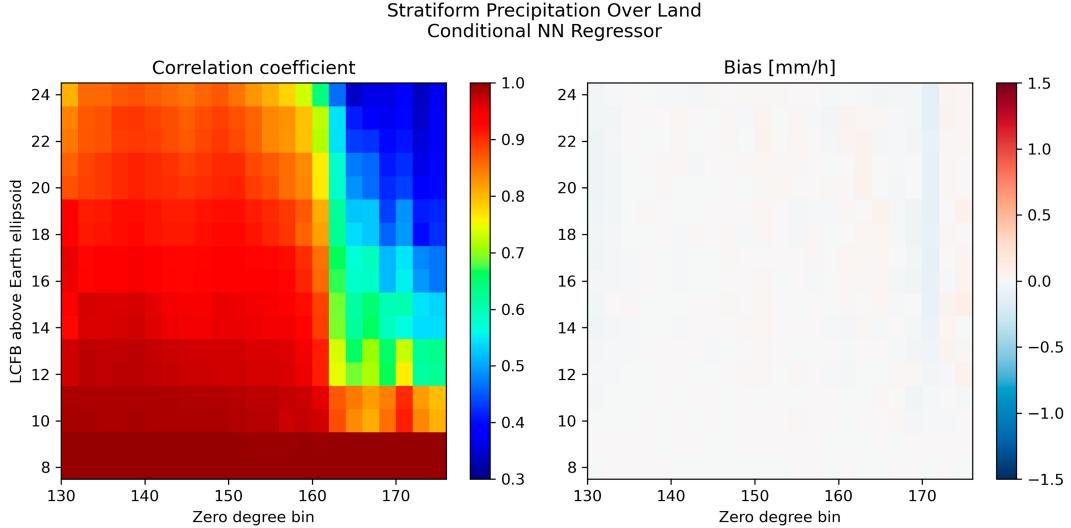


FIG. 9. Performance of the NN clutter mitigation method for stratiform precipitation over land. That is, as in Fig. 8, but for NN surface precipitation rate estimates instead of the persistence-based estimates.

of estimation is referred to as persistence in the figure and henceforth. Similar to the correlation coefficient, the largest biases occur when the LCFB is in the ice phase. Both the correlation coefficients and biases exhibit a discontinuous distribution for profiles with a ZDB around 170. This behavior is likely a consequence of the fact that precipitation estimates in the mixed layer may be biased and noisy. At the same time, the DPR detection capabilities deteriorate for profiles with only snow above the clutter or if the melting layer is close to the clutter.

Unlike persistence-based estimates, surface precipitation estimates based on Eq. (1) would be bias-free (assuming that the precipitation climatology is bias-free in the training dataset). However, the distribution of correlation coefficients between the estimates and the true surface values would not be different from that shown in Fig. 8. In other words, systematic errors are zero in estimates based on Eq. (1), but the random differences remain largely the same. An estimation superior to bias removal would also show an improvement in the distribution of the correlation coefficients and an overall reduction in the root-mean-square error (RMSE). Shown in Fig. 9 are the results for the NN-based method. As seen in the figure, the correlation coefficients increase slightly relative to those shown in Fig. 8, while biases are almost zero. In particular, the biases in the ice phase associated with the persistence-based estimates are largely removed. However, the marginal improvement in the correlations between the estimated surface precipitation rates and those in the databases suggests that there is significant variability of precipitation profiles in the clutter that cannot be reliably predicted from observations in the clutter-free portion of the reflectivity profile. A quantitative comparison of the performance of the persistence-based and NN-based methods is provided in a subsequent subsection.

#### b. Convective precipitation over land

Shown in Fig. 10 are the distributions of correlation coefficients and biases of the persistence-based estimator of surface

convective precipitation over land. Results are qualitatively similar to those obtained for stratiform precipitation over land but with larger biases when the LCFB is in the ice phase. Some positive biases for bins in the mixed phase are also obvious. These biases are most likely the consequence of artifacts in the precipitation estimates across the melting layer due to the use of different reflectivity/precipitation lookup tables. The distributions of correlation coefficients and biases associated with the NN model for convective precipitation over land are shown in Fig. 11. As seen in the figure, both the correlation coefficient and the bias improve relatively to results in Fig. 10. However, the bias distribution exhibits more variability around zero than the bias associated with stratiform precipitation over land. Moreover, coherent bias patterns are apparent for precipitation profiles with ZDB greater than 170. This is most likely a consequence of convective precipitation exhibiting more vertical variability while, as apparent in Fig. 7, being about five times less frequent than stratiform precipitation in the warm season and significantly less frequent in the cold season. This makes the statistics of convective precipitation profiles in the training dataset more uncertain than those of stratiform precipitation. Uncertainty can be mitigated by extending the dataset through the inclusion of DPR observations and CORRA estimates from other periods, but it is likely that part of it is caused by artifacts due to multiple scattering and nonuniform beam filling in the precipitation estimation procedure. From this perspective, it is beneficial that convective dataset extension be considered at the same time with or after a refinement of the convective precipitation estimation methods in CORRA.

#### c. Precipitation over oceans

The statistics for precipitation over oceans are qualitatively similar to those over land; see Figs. 12 and 13. The most significant difference is that, as suggested by Figs. 3 and 4, the mean precipitation profiles have different shapes, with the

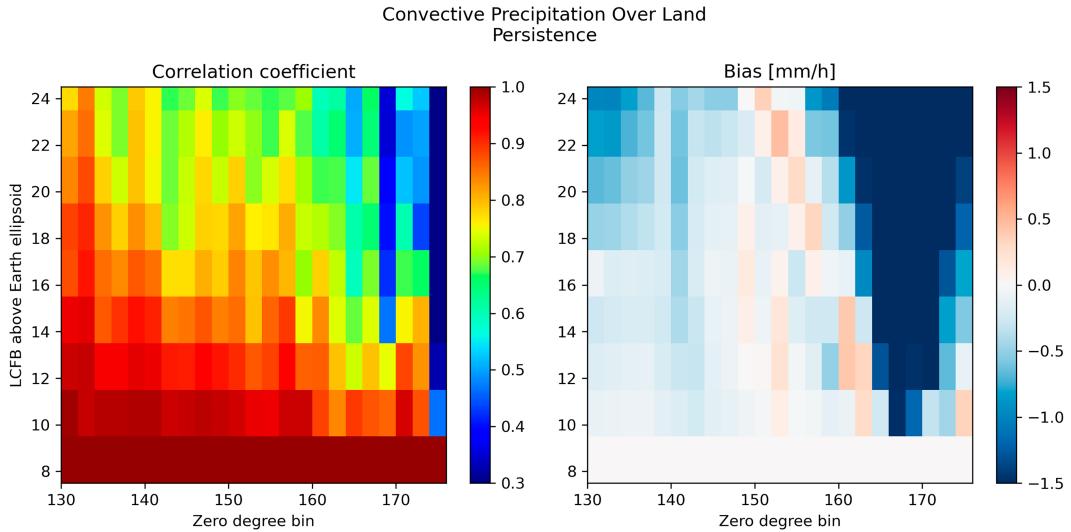


FIG. 10. As in Fig. 8, but for convective precipitation over land.

oceanic precipitation generally exhibiting more systematic increases with a range below the freezing level than precipitation over land. However, the NN clutter correction schemes exhibit behaviors similar to those over land for both stratiform and convective precipitation types. This is shown in Fig. 12 for stratiform precipitation and in Fig. 13 for convective precipitation. A notable difference compared to the estimates over land is that the bias patterns in profiles with a ZDB around 170 are less pronounced. This is likely due to the larger number of profiles and the better-defined climatology in the ocean dataset.

#### *d. Evaluation of the random errors in the correction*

In the previous section, the NN method was shown to be effective in removing the biases associated with the persistence-

based estimates. However, the random errors in the estimates were not necessarily reduced. Specifically, the correlation coefficients between NN surface precipitation estimates and actual surface precipitation estimates did not appear to be significantly improved relative to those associated with the persistence-based estimates.

To investigate this quantitatively, we calculate the relative RMSE associated with both a climatological scaling correction based on Eq. (1) and the NN estimates. The relative RMSE involves normalization by the standard deviation of the conditional surface precipitation rates. Results are given in Table 1. Here, the NN estimates exhibit RMSEs that are about 10% smaller than those of the climatological scaling estimates. This suggests that a simple bias-removal methodology based on Eq. (1) in section 2a is likely to be satisfactory in many respects. Nevertheless, the application of the NN

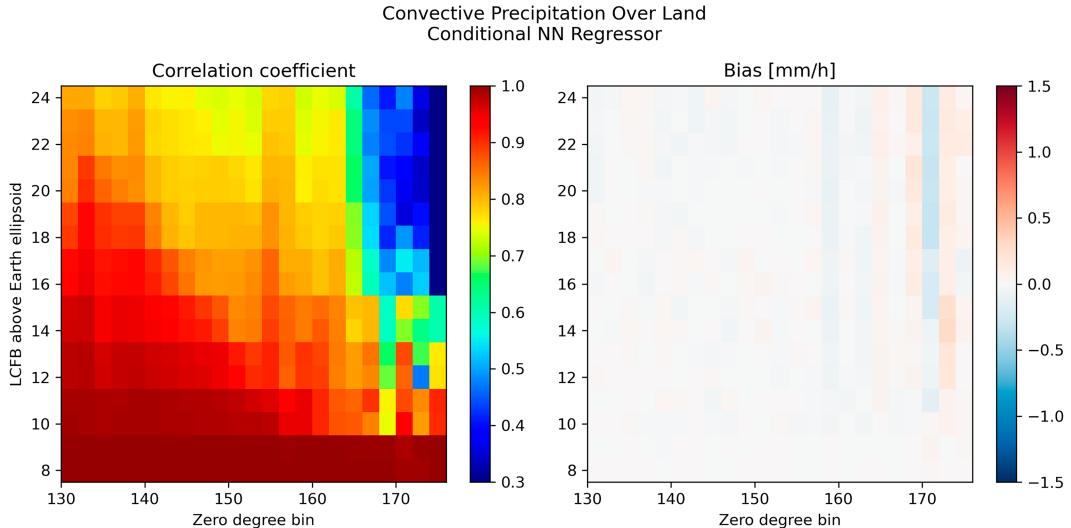


FIG. 11. As in Fig. 9, but for convective precipitation over land.

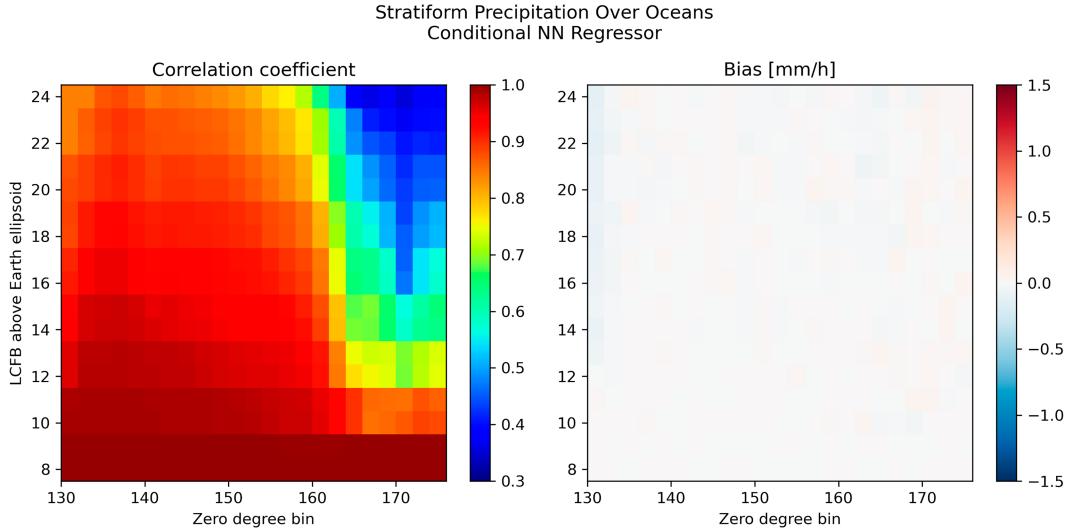


FIG. 12. Performance of the NN clutter mitigation method for stratiform precipitation over oceans.

method results in some RMSE reduction. As expected, the relative RMSEs are greater in convective than in stratiform precipitation. The fact that the NN method (which is representative of a broader class of one-dimensional clutter mitigation ML techniques) does not result in significant improvements relative to the simple bias correction provided by Eq. (1) is not necessarily an indication that ML techniques offer no benefit to the clutter mitigation problem. One potential advantage of the ML techniques is that they can incorporate radiometer observations, which may yield a significant benefit in the estimation of light precipitation over oceans. Also, the correction methods explored in this study as well as in the previous work of Hirose et al. (2021) make exclusive use of profile-level information. However, modern deep learning architectures such as U-Nets (Siddique et al. 2021) can readily process 3D information that may be useful for identifying the impacts of phenomena such as

the wind shear on reflectivity observations and use this kind of information to more accurately predict the distribution of precipitation in the clutter. Particularly, the U-Net formulation of King et al. (2024) appears promising and will be explored in future studies.

#### 4. Application to GPM CORRA precipitation estimates over the contiguous United States in the cold season

To investigate the impacts of clutter mitigation on the estimation of precipitation over the contiguous United States (CONUS) in the cold season, we apply the NN method to all GPM CORRA retrievals over CONUS from 1 December 2021 to 28 February 2022. While the same type of analysis can be applied to the entire GPM domain over all seasons, given that the focus of this study is on the fundamental benefits and

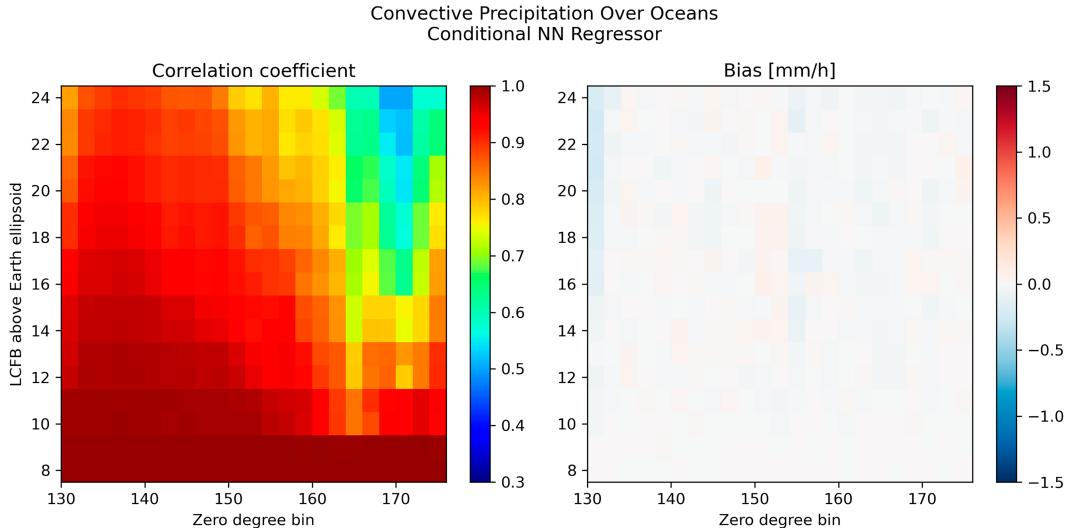


FIG. 13. As in Fig. 12, but for convective precipitation.

TABLE 1. Relative RMSE and bias of the NN and climatological scaling methods for stratiform and convective precipitation over land and oceans.

	Conditional NN	Climatological scaling	Conditional NN	Climatological scaling
Surface		Land		Oceans
Precip type			Stratiform	
RMS	62.1%	68.9%	62.8%	68.6%
Bias	-0.65%	-0.13%	-0.74%	0.8%
Precip type			Convective	
RMS	76.8%	85.1%	86.6%	95.8%
Bias	0.88%	-0.41%	-0.23%	-1.14%

limitations of profile-level corrections rather than their climatological impact, we limit the investigation to a single region and season and defer more extensive analyses to future studies. Only profiles with freezing levels below 1250 m are considered in the analysis because they are conducive to the largest corrections (and errors in the absence of any correction), as the LCFB may be associated with temperatures below freezing, while the surface precipitation may be rain.

Shown in Fig. 14 is the mean Ku-band reflectivity conditioned on the observed profiles being classified as precipitating. The means are conditioned on the associated profiles being classified as precipitating. As seen in the figure, the region contaminated by clutter (characterized by large reflectivity values) increases in height with the incidence angle. Some artifacts related to the processing of the received power to mitigate sidelobe clutter (Kubota et al. 2016) are also apparent in the figure. Specifically, while some enhanced echo is visible above 4.0 km, a slight reduction in the reflectivities is apparent near the center of the swath (roughly from ray 20 to ray 30). The reduction is more significant below the average height of the LCFB (blue line in the figure), but that reduction

does not directly impact the precipitation estimation, as the pixels associated with it are classified as clutter.

Shown in the top panel of Fig. 15 are the conditional near-surface precipitation estimated by CORRA and the surface precipitation predicted by the NN method. As seen in this figure, the clutter mitigation methodology has a significant impact on the precipitation estimates, with the impact increasing from the center toward the edges of the swath. This behavior is, most likely, a consequence of the fact that the DPR's detection capabilities deteriorate near the edges of the swath for precipitation systems with low FLH. The detected profiles are fewer but are characterized by more intense (and deeper) precipitation that results in reflectivity observations that can be reliably distinguished from clutter. This hypothesis is consistent with the distribution of the number of precipitation profiles as a function of rays, shown in the bottom panel of Fig. 15. However, the clutter correction technique does not result in artificial increases of intensity with distance from the center of the swath in the overall (unconditional) precipitation

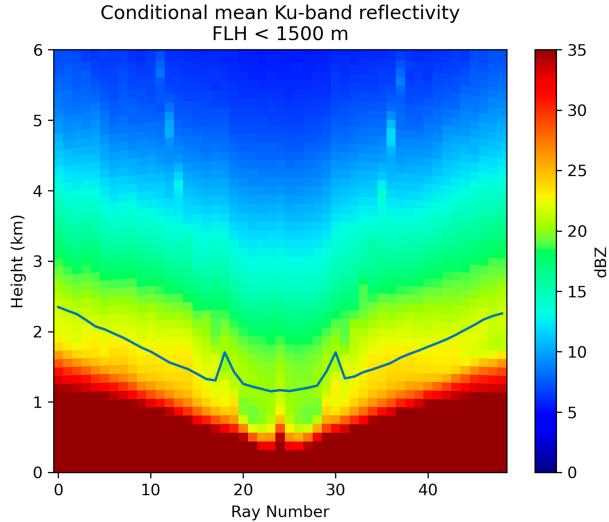


FIG. 14. Mean Ku-band reflectivity conditioned on the observed profile being classified as precipitating for all observations with  $\text{FLH} < 1500 \text{ m}$  over CONUS from 1 Dec 2021 to 28 Feb 2022. The blue line indicates the average height of the LCFB.

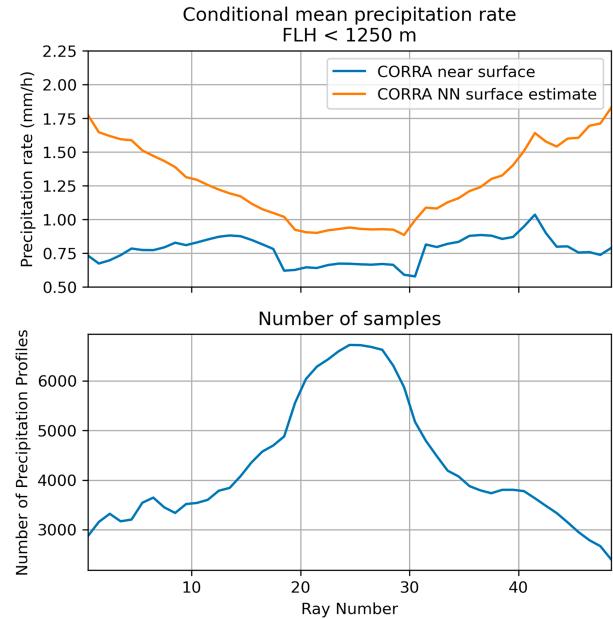


FIG. 15. (top) Conditional near-surface mean precipitation rate from the CORRA and the surface mean precipitation rate predicted by the NN method. (bottom) Number of detected precipitation profiles as a function of ray index.

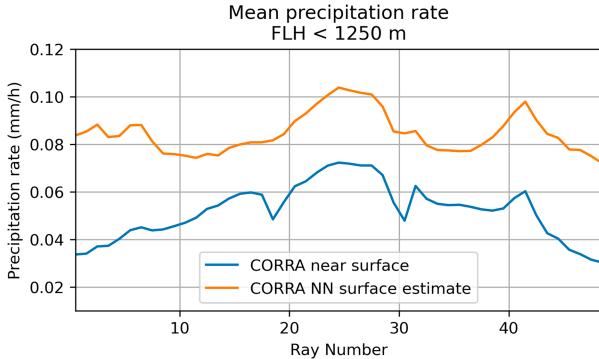


FIG. 16. Near-surface mean precipitation rate from CORRA and the surface mean precipitation rate predicted by the NN method.

rate. This point is illustrated in Fig. 16. Instead, the opposite effect, i.e., a reduction of the unconditional precipitation rate estimates with distance from the swath center (consistent with the DPR precipitation detection capabilities near the edges of the swath), is apparent in the figure. The overall impact of the NN surface precipitation rate estimation procedure is significant for precipitation systems with freezing-level heights below 1250 m over CONUS. This is an indication that significant precipitation growth processes such as water vapor deposition and riming occur in the clutter region.

Shown in Fig. 17 is an example of the extension of the CORRA precipitation estimates into the clutter zone for orbit 50632 over CONUS on 26 January 2023. The top panel shows the nominal precipitation rate from CORRA consisting of estimates derived directly from reflectivity observations not affected by clutter, while the bottom panel shows the precipitation rates extended into the clutter zone using the NN method. The figure illustrates the significant impact of the clutter mitigation methodology on the precipitation estimates, with the impact increasing from the center toward the edges of the swath. For the extension of the NN estimates to surfaces closer than 1.00 km from the Earth ellipsoid, we use Eq. (2).

## 5. Summary and conclusions

In this study, a new method for mitigating ground clutter effects in precipitation estimates derived from the GPM mission's CORRA algorithm is developed. CORRA combines data from the DPR and GMI on the GPM core satellite to estimate precipitation rate, and ground clutter is a significant problem for spaceborne radar observations, as it can obscure or corrupt the signal associated with precipitation. An approach to mitigate ground clutter using statistical relationships based on precipitation estimates from near-nadir scans has already been developed (Hirose et al. 2021) and applied to precipitation estimates from the DPR algorithm (Iguchi et al. 2021). However, the study of Hirose et al. (2021) did not fully explore the benefits and limitations of statistical methods to mitigate clutter in the DPR reflectivity observations.

To build upon the previous work, ML approaches are investigated to gain further insight into the uncertainties of surface precipitate rates derived from information in the portion

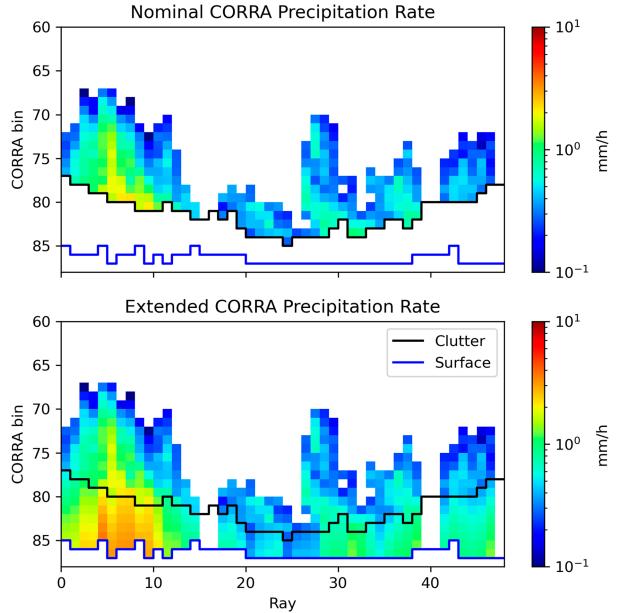


FIG. 17. Example of extension of CORRA precipitation estimates into the clutter zone for orbit 50632 over CONUS on 26 Jan 2023. CORRA bins are characterized by a coarser resolution than the DPR bins, with every two consecutive DPR bins being aggregated into a CORRA bin.

of the reflectivity profile not affected by clutter. The ML model uses reflectivity observations, along with additional information such as precipitation type, surface type, and freezing level, to estimate the surface precipitation rate. The benefits of this approach include the use of ML models efficient at leveraging existing features and capturing complex relationships within the data without relying on explicit feature engineering. Specifically, various machine learning architectures are investigated to automatically extract information from the data without resorting to subjective efforts. A preliminary evaluation suggests that no architecture offers significantly better performance than the others, and so we select the NN as the best candidate for further systematic evaluations since it is computationally fast to train and deploy while effective in application. Nevertheless, the ML approach results in a reduction of about 10% in the relative RMSE of the precipitation rate estimates relative to a simple climatological scaling method.

The database used for training and evaluating the ML models in this study is derived from 2 years of DPR near-nadir observed reflectivity profiles over oceans and 5 years of data over land, all of which are minimally affected by clutter. A minor limitation of this database is that, despite the small number of clutter-affected bins, statistical models and assumptions are still necessary to derive surface precipitation estimates. To extend the NN precipitation estimates from 1.0 km above the Earth ellipsoid to the surface, we hypothesized that precipitation profiles with freezing-level heights (FLHs) differing by no more than 1.0 km are highly similar, and thus, climatological scaling can be applied.

However, this hypothesis requires further evaluation. Over land, high-quality ground radar precipitation estimates, adjusted by rain gauges, such as those provided by the MRMS product (Zhang et al. 2016), may be used for this evaluation. Over oceans, the evaluation is likely to be more challenging due to the limited availability of data useful for direct validation of the estimates.

Estimates using the LCFB as a proxy for surface precipitation rate are systematically different from the actual surface precipitation rates in the training data, necessitating an evaluation of the NN model in this context. Specifically, relying on the LCFB's precipitation rate to estimate surface precipitation leads to biased results, prompting a closer examination of the NN model's ability to produce unbiased estimates. The NN model demonstrates effectiveness in providing unbiased estimates, but it only slightly outperforms a basic climatological scaling method in reducing random errors. Additionally, the performance of simpler machine learning techniques, such as regularized multivariate regressions, mirrors that of the NN model in initial assessments. This suggests that the NN model's limited improvement beyond bias removal is due to the inherent challenges of the problem rather than the limitations of the model itself.

**Acknowledgments.** This work was supported by the NASA Global Precipitation Measurement (GPM) mission project. The authors thank Dr. Tsengdar Lee and Dr. Will McCarty (NASA Headquarters) for their support of this effort.

**Data availability statement.** Version 7 of GPM DPR and CORRA data can be accessed online (<https://arthurhouhttps://pps.eosdis.nasa.gov/gpmdata/>). Code to train the ML models used in this study and training and evaluation data may be accessed online (<https://tinyurl.com/corraBlindZone>).

## APPENDIX

### List of Acronyms and Abbreviation

CORRA	Combined Radar–Radiometer Algorithm
DPR	Dual-frequency precipitation radar
GPM	Global Precipitation Measurement
LCFB	Lowest clutter-free bin
FLH	Freezing-level height
ML	Machine learning
NN	Neural network
NWP	Numerical weather prediction
PCA	Principal component analysis
PIA	Path-integrated attenuation
RMSE	Root-mean-square error
SHAP	Shapley additive explanation
SRT	Surface reference technique
ZDB	Zero-degree bin

## REFERENCES

- Bergstra, J., and Y. Bengio, 2012: Random search for hyperparameter optimization. *J. Mach. Learn. Res.*, **13**, 281–305.
- Bishop, C. M., 2006: *Pattern Recognition and Machine Learning*. Springer, 738 pp.
- Chase, R. J., S. W. Nesbitt, and G. M. McFarquhar, 2021: A dual-frequency radar retrieval of two parameters of the snowfall particle size distribution using a neural network. *J. Appl. Meteor. Climatol.*, **60**, 341–359, <https://doi.org/10.1175/JAMC-D-20-0177.1>.
- Grecu, M., W. S. Olson, S. J. Munchak, S. Ringerud, L. Liao, Z. Haddad, B. L. Kelley, and S. F. McLaughlin, 2016: The GPM combined algorithm. *J. Atmos. Oceanic Technol.*, **33**, 2225–2245, <https://doi.org/10.1175/JTECH-D-16-0019.1>.
- Hancock, J. T., and T. M. Khoshgoftaar, 2020: Survey on categorical data for neural networks. *J. Big Data*, **7**, 28, <https://doi.org/10.1186/s40537-020-00305-w>.
- Hirose, M., S. Shige, T. Kubota, F. A. Furuzawa, H. Minda, and H. Masunaga, 2021: Refinement of surface precipitation estimates for the Dual-Frequency Precipitation Radar on the GPM core observatory using near-nadir measurements. *J. Meteor. Soc. Japan*, **99**, 1231–1252, <https://doi.org/10.2151/jmsj.2021-060>.
- Iguchi, T., and Coauthors, 2021: GPM/DPR level-2 algorithm theoretical basis document. NASA Goddard Space Flight Center Tech Doc., 127 pp., [https://gpm.nasa.gov/sites/default/files/document\\_files/ATBD\\_DPR\\_201811\\_with\\_Appendix3b\\_0.pdf](https://gpm.nasa.gov/sites/default/files/document_files/ATBD_DPR_201811_with_Appendix3b_0.pdf).
- King, F., G. Duffy, L. Milani, C. G. Fletcher, C. Pettersen, and K. Ebell, 2022: DeepPrecip: A deep neural network for precipitation retrievals. *Atmos. Meas. Tech.*, **15**, 6035–6050, <https://doi.org/10.5194/amt-15-6035-2022>.
- , C. Pettersen, C. G. Fletcher, and A. Geiss, 2024: Development of a full-scale connected U-Net for reflectivity inpainting in spaceborne radar blind zones. *Artif. Intell. Earth Syst.*, **3**, e230063, <https://doi.org/10.1175/AIES-D-23-0063.1>.
- Kingma, D. P., and J. Ba, 2014: Adam: A method for stochastic optimization. arXiv, 1412.6980v9, <https://doi.org/10.48550/arXiv.1412.6980>.
- Koistinen, J., 1991: Operational correction of radar rainfall errors due to the vertical reflectivity profile. Preprints, *25th Int. Conf. on Radar Meteorology*, Paris, France, Amer. Meteor. Soc, 91–94.
- Kubota, T., T. Iguchi, M. Kojima, L. Liao, T. Masaki, H. Hanado, R. Meneghini, and R. Oki, 2016: A statistical method for reducing sidelobe clutter for the Ku-band precipitation radar on board the GPM core observatory. *J. Atmos. Oceanic Technol.*, **33**, 1413–1428, <https://doi.org/10.1175/JTECH-D-15-0202.1>.
- Lundberg, S. M., and S.-I. Lee, 2017: A unified approach to interpreting model predictions. *Proc. 31st Int. Conf. on Neural Information Processing Systems*, Long Beach, CA, Curran Associates Inc., 4768–4777, <https://dl.acm.org/doi/10.5555/3295222.3295230>.
- Mishra, P., 2023: Pytorch model interpretability and interface to sklearn. *PyTorch Recipes: A Problem-Solution Approach to Build, Train and Deploy Neural Network Models*, Springer, 237–260, [https://doi.org/10.1007/978-1-4842-8925-9\\_10](https://doi.org/10.1007/978-1-4842-8925-9_10).
- Morrison, H., and Coauthors, 2020: Confronting the challenge of modeling cloud and precipitation microphysics. *J. Adv. Model. Earth Syst.*, **12**, e2019MS001689, <https://doi.org/10.1029/2019MS001689>.
- Nair, V., and G. E. Hinton, 2010: Rectified linear units improve restricted Boltzmann machines. *Proc. 27th Int. Conf. On Machine Learning (ICML-10)*, Haifa, Israel, Omnipress, 807–814, <https://dl.acm.org/doi/10.5555/3104322.3104425>.
- NWS, 2023: National weather service beam property calculator. <https://training.weather.gov/wdtd/tools/beamwidth/>.

- Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.*, **12**, 2825–2830, <https://doi.org/10.5555/1953048.2078195>.
- Rahimi, R., P. Ravirathinam, A. Ebtehaj, A. Behrangi, J. Tan, and V. Kumar, 2024: Global precipitation nowcasting of Integrated Multi-satellitE Retrievals for GPM: A U-Net convolutional LSTM architecture. *J. Hydrometeor.*, **25**, 947–963, <https://doi.org/10.1175/JHM-D-23-0119.1>.
- Siddique, N., S. Paheding, C. P. Elkin, and V. Devabhaktuni, 2021: U-net and its variants for medical image segmentation: A review of theory and applications. *IEEE Access*, **9**, 82031–82057, <https://doi.org/10.1109/ACCESS.2021.3086020>.
- Skofronick-Jackson, G., and Coauthors, 2017: The Global Precipitation Measurement (GPM) mission for science and society. *Bull. Amer. Meteor. Soc.*, **98**, 1679–1695, <https://doi.org/10.1175/BAMS-D-15-00306.1>.
- Zhang, J., and Coauthors, 2016: Multi-Radar Multi-Sensor (MRMS) quantitative precipitation estimation: Initial operating capabilities. *Bull. Amer. Meteor. Soc.*, **97**, 621–638, <https://doi.org/10.1175/BAMS-D-14-00174.1>.
- Zheng, A., and A. Casari, 2018: *Feature Engineering for Machine Learning: Principles and Techniques for Data Scientists*. O'Reilly Media, 218 pp.