# Refining Predictive Models for Sea Surface Currents:
# A Focus on Variable Configuration and Time Sequence Analysis

Ittaka Aldini[1,2], Adhistya Erna Permanasari[1], Risanuri Hidayat[1], Andri Ramdhani[2]

5 [1]Department of Electrical & Information Engineering, University of Gadjah Mada, Yogyakarta, 55281, Indonesia.
[2]Indonesian Agency for Meteorology, Climatology, and Geophysics, Jakarta, 10610, Indonesia.

*Correspondence to*: Adhistya Erna Permanasari (adhistya@ugm.ac.id)

**Abstract.** Accurate prediction of sea surface currents is crucial for understanding ocean dynamics, climate variability, and marine ecosystem health. Despite advancements in statistical modeling, challenges remain in terms of optimizing model parameters and variable configurations to enhance prediction accuracy. This study employed high-frequency (HF) radar data from the Bali Strait (2018-2021) to develop a statistical modeling approach for sea surface current prediction. We utilize random forest regression (RFR) as the primary machine learning technique. The data were subjected to a rigorous preprocessing pipeline to ensure robustness, including selection, cleaning, and imputation. We define 11 distinct model configurations with various input parameters, such as moving averages (avgh3, avgh6, or avgh12) and previous day values (h-24, h-48, and h-72). Our analysis focused on three prediction schemes: seasonal (P1) and monthly (P2 and P3), each with tailored training and testing data allocations. This study evaluates the models using root mean square error (RMSE) and Coefficient of Determination (R²). Results indicate that combining moving-average predictors significantly enhances the accuracy of long-term forecasts, whereas short-term predictions benefit from utilizing recent data. Our findings highlight specific variable configurations, particularly those incorporating moving averages, which lead to superior performance in sea surface current prediction. The results indicate that models employing configurations F1, F5, and F8 yield the best results, highlighting the importance of optimizing model variables to achieve high-accuracy predictions.

**Keywords**: Sea Surface Current Prediction, RF Regression, Variable Configuration, Model Optimization.

## 1. Introduction

25 The prediction of sea surface currents plays a crucial role in understanding ocean dynamics and their implications for various environmental and socioeconomic activities, such as navigation (Barrick et al., 2012), fisheries management (Ren et al., 2017), and coastal protection (Cosoli and de Vos, 2019). Accurate predictions of sea surface currents can significantly enhance the operational efficiency in maritime activities (Marmain et al., 2014) and improve safety measures in coastal regions (Saviano et al., 2019). High-Frequency (HF) Radar technology has emerged as a valuable tool for monitoring ocean surface currents,

30    providing high-resolution and real-time data essential for comprehending the complex patterns of sea surface currents (Kirincich et al., 2019). The U and V component radial current velocity data obtained from HF Radar systems offer insights into these dynamics, enabling researchers to develop more accurate predictive models (Pascual et al., 2015).

Recent advancements in ensemble modeling techniques have significant potential to enhance the capabilities of sea surface current predictions. Ensemble approaches have consistently demonstrated superior performance compared to single statistical

35    or machine learning approach, offering enhanced accuracy and reliability (Wu and Levinson, 2021), especially in complex and dynamic environments such as ocean coastal regions (Werner and Blanton, 2019). These techniques have proven particularly effective in addressing the multifaceted challenges associated with coastal zone management and prediction, including ocean waves (O'Donncha et al., 2019), significant wave heights (Ali and Prasad, 2019), sea level variations (Balogun and Adebisi, 2021), storm surge forecasting (Rezuanul Islam et al., 2023), and shoreline evolution modeling (Montaño et al., 2020). These

40    findings underscore the potential of ensemble modeling techniques to enhance the reliability and accuracy of predictions in oceanographic studies. However, despite these advancements, a notable gap remains in the literature regarding the systematic evaluation of variable configurations and temporal sequences in predictive models for sea surface currents.

The importance of accurately predicting sea surface currents cannot be overstated, particularly in the context of climate change and its impacts on oceanographic processes (Hardman and Wyatt, 2019). Current models often rely on historical data

45    that may not fully capture the dynamic nature of ocean currents, leading to potential inaccuracies in predictions (Zhao et al., 2020). Consider that while certain studies have utilized long-term datasets for training models, such as 365 days (1 year) (Thongniran et al., 2019a, b), 669 days (22 months) (Jitkajornwanich et al., 2017), 855 days (Jirakittayakorn et al., 2017), 1309 days (43 months) (Li et al., 2022),  and even up to 9862 days (324 months) (Xiao et al., 2019). Conversely, others have focused on shortened time scales, from 70 days (Kusnanti et al., 2022), 54 days (Sarkar et al., 2018), 48 days, 21 days (Pramesti et al.,

50    2022; Sarkar et al., 2018), 5 days (Putri et al., 2022), 3 days, and 1 day (Kusnanti et al., 2022; Zulfa et al., 2021). These differences result in varying degrees of predictive success. This disparity highlights the need for a more nuanced understanding of how different variable configurations and time sequences affect model performance.

Moreover, existing research has primarily focused on comparing different algorithms without adequately addressing the optimization of model parameters and input variables (Jirakittayakorn et al., 2017; Thongniran et al., 2019b). This gap

55    represents an opportunity to enhance the robustness of predictive models by systematically evaluating the impact of various configurations and temporal sequences on prediction outcomes.

Against this backdrop, the present research aims to refine predictive models for sea surface currents through a focused analysis of variable configuration and time sequence. Specifically, this study seeks to address the existing gap in the literature by optimizing model parameters and input variables to enhance the prediction accuracy over diverse time scales. By

60    systematically evaluating the impact of various configurations and temporal sequences on model performance, we aim to identify the most effective strategies for improving the robustness of sea surface current predictions.

The research questions guiding this study include how different variable configurations influence the accuracy of sea surface current predictions and the impact of temporal sequences on model performance. We hypothesized that optimizing the variable selection and temporal configurations would lead to improved prediction accuracy and reduced error rates for sea surface current forecasting.

In summary, this study builds upon the foundation of HF Radar technology and advanced statistical modeling to address the complexities of sea surface currents predictions. The subsequent sections of this paper will detail the methodological framework, present the findings of the analysis, and discuss the implications of the research.

## 2. Methodology

### 2.1. Data Collection and Postprocessing

The first step involve collecting HFR sea surface current data, including extracting the U and V component radial current velocity data. The values of u and v were gathered from two sites in the Bali Strait CODAR HF Radar every hour from December 1st, 2018 to 30 November 30th, 2021. The data is divided into training and test sets. We selected the all-year dataset to preserve the seasonal characteristics of the dynamic ocean surface currents. The dataset undergoes post-processing selection to ensure the quality and integrity of the data for predictive modeling. In this study, several approaches to data post-processing exist, namely data selection, imputation, and partitioning.

### 2.1.1. Data Selection

Data selection involves identifying and extracting relevant data from a large dataset for a specific analysis or modeling tasks. The process involves carefully curating and extracting data pertinent to the research objectives, ensuring that the selected data align with the particular requirements of the analysis or modeling exercise. In the context of sea surface current prediction using HFR data, data selection encompasses identifying and extracting U and V component radial current velocity data from the HFR sea surface current dataset. In this study, the original HFR data from the TUV format contained 17 columns of variables, including latitude, longitude, velocity, direction, and 13 other variables. Since this study only utilized U and V data, 15 other variables were dismissed and reduced, thus generating time-series data.

### 2.1.2. Data Imputation

Prior to doing data imputation, it is essential to identify the extent of missing data following the data gathering. Upon completing an investigation, it was discovered that there are 12,40%, 12,86%, and 11,39% missing data for 2019, 2020, and 2021, as shown in Table 1.

**Table 1. Percentage of HF Radar Data Availability, U and V**

| Data | Availability |
|------|--------------|
| DEC 2018 – NOV 2019 | 87.60% |
| DEC 2019 - NOV 2020 | 87.14% |
| DEC 2020 - NOV 2021 | 88.61% |

95

Thus, this study has 12,22% (6427) missing data from the 52.608 datasets. The causes of missing data can result from a poor radar signal and HFR stations processing errors. It was evident that grid-point data quality (regarding data availability) was a significant issue that could affect the prediction accuracy. In order to prevent this, this study employs the data-filling linear interpolation method. Linear interpolation is among the simplest forms of interpolation, and it connects two data points

100 by a straight line. The equation of the linear interpolation function is given by Eq. (1):

$$f_1(x) = b_0 + b_1(x - x_0) \tag{1}$$

where $f_1(x)$ is the estimated value at point $x$, $b_0$ is the value of the function at the known point $x_0$, defined as $b_0 = f(x_0)$; b1 represents the slope of the line connecting the two known points, which is calculated as follows:

$$b_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0} \tag{2}$$

105 Here, $f_1(x_1)$ is the value of the function at another known point $x_1$.

### 2.1.3. Data Partitioning

Data partitioning was conducted to facilitate three distinct prediction schemes: Prediction Scheme 1 (long-term seasonal), Prediction Scheme 2 (short-term monthly), and Prediction Scheme 3 (short-term monthly comprehensive). Each scheme

110 employs a tailored data division approach to optimize the training and testing processes to realize accurate predictions.

Prediction Scheme 1 utilizes seasonal data categorized according to the four monsoon seasons: December-January-February (DJF), March-April-May (MAM), June-July-August (JJA), and September-October-November (SON). For this scheme, the training dataset comprised data from 2019 and 2020, and the testing dataset comprised of data from 2021. Predictions are generated separately for each season, ensuring that the model is trained on relevant seasonal patterns. Prediction

115 Scheme 2 focuses exclusively on 2021, encompassing data from December 2020 to November 2021. The dataset is divided monthly, with 80% allocated to training and 20% to testing. This structure allows for a detailed analysis of monthly variations in sea surface currents, enabling the model to adapt to short-term fluctuations. Prediction Scheme 3 combines data from three years—2019, 2020, and 2021—providing a more comprehensive dataset for analysis. Similar to Scheme 2, the data are first split by months. However, in this scheme, the training dataset includes all relevant month data from 2019 and 2020 along with

120 the same testing dataset used in Prediction Scheme 2. This approach enhances the model's ability to capture the temporal trends and characteristics of sea surface currents by leveraging a broader historical context.

## 2.2. Model Development and Prediction Modeling

### 2.2.1. Random Forest Regression Implementation

In this study, we primarily used random forest regression (RFR), a machine learning algorithm, to predict sea surface currents.
125    Specifically, RFR is employed to forecast the U and V radial velocity components derived from high-frequency (HF) radar observations. RFR is an ensemble method that combines the predictions of multiple decision trees to produce more accurate and stable results (Ambhika et al., 2024; Kunapuli, 2023). The RFR model is generate several decision trees, each trained on a different subset of the training data. This process involves two key randomization steps: selecting a random subset of the training data and choosing a random subset of features for each tree. The final prediction was derived from the weighted
130    average of the predictions made by all individual trees in the forest (Bateman et al., 2020; Schonlau and Zou, 2020).

When applying RFR to our dataset, we used an ensemble of decision trees specifically tailored to predict continuous target variables associated with sea surface currents. The model was trained on a carefully partitioned dataset, as outlined in Section 2.1.3, to ensure adequate representation of both temporal and seasonal characteristics. The training process involves bagging, a technique in which subsets of the training data are repeatedly sampled to create smaller decision trees. These smaller trees
135    are then combined to form the overall Random Forest model. One of the notable advantages of RFR is its robustness against missing values in the dataset, which allows it to maintain performance even when faced with incomplete information (Tang and Ishwaran, 2017). This characteristic is particularly beneficial for oceanographic studies, where data gaps can occur due to environmental factors or instrument limitations.

Overall, the implementation of Random Forest Regression in this study leverages its ensemble capabilities, including
140    bagging, to provide accurate and stable predictions of sea surface currents. The subsequent sections will detail the results obtained from this model, highlighting its effectiveness in capturing the dynamics of U and V components across different prediction schemes.

### 2.2.2. Variable Configurations of the Model

Following the data preparation, prediction models were created. There are eleven (11) prediction models made for each U and
145    V component of the radial velocity data of sea surface current. Each model targeted either the U or V components of the radial velocity data. A detailed description of these 11 models is given in Table 2.

When we observe Table 2, each of the 11 models has a different combination of variables. The variables of (h-1) to (h-3) represent the last three hours of data. The variable avgh3, avgh6, avgh12, and avgh24 are the average measurement values for the last three, six, twelve, and twenty-four hours, respectively. Meanwhile, the variables of (h-24), (h-48), and (h-72)
150    respectively, represent data for the previous 24 hours, the previous 48 hours, and the previous 72 hours.

These combinations of models were carefully selected to capture diverse timescales and characteristics of the current dynamics of the sea surface, allowing for a comprehensive assessment of the predictive capabilities of the method.

155

**Table 2. Eleven prediction models**

| Model | Variable |
|---|---|
| F1 | (h-1) to (h-3), avgh3 |
| F2 | (h-1) to (h-3), avgh3, (h-24) |
| F3 | (h-1) to (h-3), avgh3, (h-24), (h-48) |
| F4 | (h-1) to (h-3), avgh3, (h-24), (h-48), (h-72) |
| F5 | (h-1) to (h-3), avgh3, avgh6 |
| F6 | (h-1) to (h-3), avgh3, avgh6, (h-24) |
| F7 | (h-1) to (h-3), avgh3, avgh6, (h-24), (h-48) |
| F8 | (h-1) to (h-3), avgh3, avgh6, avgh12 |
| F9 | (h-1) to (h-3), avgh3, avgh6, avgh12, (h-24) |
| F10 | (h-1) to (h-3), avgh3, avgh6, avgh12, (h-24), (h-48) |
| F11 | (h-1) to (h-3), avgh3, avgh6, avgh12, avgh24 |

After the creation of the model, the dataset with the defined model set is applied using the random forest regression algorithm. Random Forest Regression is a machine learning algorithm used for regression tasks. The proposed ensemble method combines the predictions of multiple decision trees to produce more accurate and stable predictions. The algorithm works by

160 constructing multiple decision trees from different subsets of the training data. Each tree was trained on a random subset of models and samples. The final prediction is the average of the predictions of all trees. One of the key advantages of random forest regression is its ability to handle nonlinear relationships between models and target variables. It can also handle missing data values and is overfitting resistant.

165 **2.3. Performance Matrix and Evaluation**

The prediction results obtained using these 11 models were compared to determine the optimal model for predicting sea surface currents. The comparative analysis aims to identify the most effective set of models to enhance the accuracy and robustness of sea surface current predictions. To evaluate the performance of the statistical models, we used the following performance metrics: RMSE and R2.

170 The root mean square error (RMSE) of each step is given by,

$$RMSE = \sqrt{\sum_{i=1}^{n} \frac{(\widehat{y_i} - y_i)^2}{n}}, \tag{3}$$

Where $y_i$ are the observed values; $\widehat{y_i}$ are the predicted values; $n$ is the number of data points. The RMSE quantifies the average magnitude of the residuals (prediction errors). It is sensitive to outliers because squared differences give more weight to significant errors. A lower RMSE value indicates better model performance, where a perfect fit has an RMSE value of 0.

6

Another important metric for assessing model performance is R-squared ($R^2$). This represents the proportion of variance in the dependent variable that is predictable from the independent variables. R-squared ranges from 0 to 1, with higher values indicating better model fit. The formula for R-squared is:

$$R^2 = 1 - \frac{SSR}{SST} = 1 - \frac{\sum_i (y_i - \widehat{y_i})^2}{\sum_i (y_i - \overline{y_i})^2},$$ (4)

Where $y_i$ are the observed values; $\widehat{y_i}$ are the predicted values; $\overline{y_i}$ is the mean of observed values; and $n$ is the number of data points. The RMSE and R-Square values are in the same unit as the data.



**Figure 1: The workflow of Sea Surface Current Prediction Scheme.**

We compared the performance of the models using these metrics and identified the most accurate statistical model for sea surface current prediction across diverse timescales and models. The overall sea surface current prediction process involves

7

several key steps, as illustrated in Figure 1. First, sea surface current data is collected from high-frequency (HF) radar systems over a multi-year period (2018-2021). Subsequently, the U (zonal) and V (meridional) components of the current time series are then selected for further processing. An imputation process is applied to the time series data to handle missing data. Next, 11 prediction models are created, each with a different combination of input variables. The data is then divided into training and testing sets based on three prediction schemes: P1 for seasonal prediction and P2/P3 for monthly predictions. The specific grouping and allocation of data for training and testing varied between these schemes. A Random Forest Regression algorithm based on decision trees is then applied to train the models and make predictions. Finally, the predictions were evaluated using the R-squared ($R^2$) and Root Mean Squared Error (RMSE) metrics to assess the accuracy of each model.

## 3. Results and Discussion

The experiment was conducted on a portable computer equipped with an Intel i7-1065G7 Octa-core @ 130 Ghz, 16GB DDR4, and Intel Iris Plus graphics specification to model random forest regression machine learning model.

### 3.1. Data Postprocessing Outcomes

#### 3.1.1. Data Selection Results

The data collection yielded high-frequency (HF) radar data spanning December 2018 to November 2021, with a availability rate of 86.26%, as illustrated in Figure 2, which serves as a representative depiction of the HF radar dataset.



**Figure 2: HF radar data over the three years (2018-2021).**

### 3.1.2. Data Imputation Results

The results of data imputation using linear interpolation for sea surface current data reveal notable implications for predictive methods. The linear interpolation method was employed to fill the missing data in the U and V components of the radial current velocity. A graph showing the outcomes of linear interpolation is shown in Figure 3. Subsequently, the interpolated data are employed in the Random Forest Regression model to forecast the U and V components. The imputed data were then used in Random Forest Regression models to predict the U and V components across seasons.

8

**Figure 3: Example graph of linear interpolation results for sea surface current data,**
**The U component (left) and V component (right), before and after interpolation.**

### 3.1.3. Data Partitioning Results

#### 3.1.3.1. Data Partitioning for Long-term Prediction (P1)

In this study, we employed data partitioning for each season to facilitate long-term (seasonal) sea surface current prediction. The U and V components of the ocean surface current data are characterized by various statistical parameters, including the count, minimum, standard deviation, and maximum values, which were categorized based on four distinct monsoon seasons: DJF (December to February), MAM (March to May), JJA (June to August), and SON (September to November). These seasonal statistics are presented in Table A1.

The dataset has been carefully divided into training and testing subsets; data from 2019 and 2020 allocated for training purposes, while data from 2021 is reserved for testing. This deliberate separation of the dataset based on seasons and the allocation of data for training and testing is crucial for ensuring the robustness and reliability of the statistical models developed for predicting ocean surface currents. By using data from different years for training and testing, the models can be rigorously evaluated for their predictive performance across seasonal variations, thereby enhancing their applicability and generalizability.

Analysis of the statistical characteristics presented in Table A1 revealed substantial variability in the sea surface current data across various seasons. The range of the U and V components, as indicated by the maximum and minimum values for each season, provides insights into the extent of the current data. Furthermore, the mean and standard deviation offer a deeper understanding of the central tendency and distribution of the sea surface currents. Examination of this data range suggests that ocean surface currents exhibit significant fluctuations throughout the monsoon seasons, with each season displaying unique characteristics and dynamics that can potentially improve prediction accuracy when incorporated into the model.

9

This variability underscores the importance of considering seasonal influences when developing statistical models for sea

270  surface current prediction, as the range of the data reflects the diverse and dynamic nature of sea currents across different

temporal contexts.

### 3.1.3.2. Data Partitioning for Short-Term Prediction (P2 & P3)

In contrast to the seasonal data partitioning employed for long-term prediction, the data distribution for prediction scheme 2

275  (P2) focuses on the short term and utilizes data from December 2020 to November 2021. The training and testing data were

divided on a monthly basis, with the process starting in December and ending in November.

This partitioning allocates 80% of the data to training and 20% to testing. After data division and segregation, the training

and testing data descriptions for each month are presented in Tables A2 and A3, respectively.

The data distribution for prediction scheme 3 (P3) closely resembles that of prediction scheme 2 (P2). However, in

280  prediction scheme 3, the training data include not only the data used in prediction scheme 2, but also 100% of the data for the

relevant month in 2019 and 100% of the data for the relevant month in 2020. The testing data were unchanged from those in

prediction scheme 2. The data descriptions for prediction scheme 3 can be found in Tables A3 and A4. The implementation of

these short-term data partitioning strategies enabled the models developed for prediction schemes 2 and 3 to effectively capture

the temporal dynamics and patterns present in the ocean surface current data on a monthly basis. This approach facilitates

285  more targeted and accurate predictions for short-term forecasts.

### 3.2. Prediction Model Creation Results

The sea surface current prediction models generates additional models from the U and V component. These models

consist of several variables, as explained in Table 2. These variables include (h-1) to (h-3) in the last 3 hours of data.

The variable avgh3, avgh6, avgh12, and avgh24 are the average measurement values for the last three, six, twelve,

290  and twenty-four hours, respectively. Meanwhile, variables (h-24), (h-48), and (h-72) respectively represent data for the

previous 24 hours, the previous 48 hours, and the previous 72 hours. The correlation matrix can be observed in Figure 4a. The

prediction model plot is also shown in Figure 4b.

This creation of this model is crucial for enhancing the model's predictive capability by incorporating historical data

and creating new variables that capture the temporal evolution of sea surface currents. By integrating these models into the

295  prediction method, it becomes possible to account for the influence of past observations and trends in sea surface currents,

thereby improving the prediction accuracy and reliability.

300

**Figure 4: (a) Heatmap correlation between variables in prediction model.**
**(b) Plot of original data and prediction models, first 80 Data.**

## 3.3. Model Performance Analysis

### 3.3.1. Performance of Prediction Scheme 1 (P1)

320　Prediction Scheme 1 employs seasonal data for training and testing, utilizing data from 2019 and 2020 for training purposes, and data from 2021 serves as the testing dataset for each season. The results of this scheme are detailed in Tables B1 to B4 under column P1. The analysis focuses on the U and V components across four seasons: December-January-February (DJF), March-April-May (MAM), June-July-August (JJA), and September-October-November (SON), using 11 distinct model configurations (F1 to F11), as outlined in Subchapter 2.2.2 and Table 2. To facilitate comparison with other prediction schemes,

325　we present the results for Prediction Scheme 1 for the first month of each season: December for DJF, March for MAM, June for JJA, and September for SON.

　　The Root Mean Square Error (RMSE) values for Prediction Scheme 1, as shown in Tables B2 and B4, under column P1, indicate the magnitude of the prediction errors. The RMSE values ranged from 8.8 to 13.87 for the U component and 11.24 to 17.60 for the V component, with the lowest RMSE observed in F1 during the DJF season for both components.

330　The correlation coefficient $R^2$ values for Prediction Scheme 1 are presented in Tables B1 and B3 under column P1. These values generally ranged from 0.94 to 0.96 for the U component and from 0.95 to 0.98 for the V component. The highest accuracy was achieved with models F1, F2, F5, F6, F8, and F11 during DJF, as well as with models F1 and F5 during SON, both yielding an $R^2$ value of 0.96 for the U component; model F1 during JJA achieved an $R^2$ value of 0.98 for the V component. The graphical representation of these optimal prediction results is illustrated in Figure 5a.

335    Across all seasons, models F1, F5, and F8 consistently demonstrated superior RMSE and R2 values compared to the other models. The elevated RMSE values and reduced R2 values observed in models other than F1, F5, and F8 suggest that certain variables significantly enhance prediction accuracy—specifically, avgh3 in F1; avgh3 and avgh6 in F5; and avgh3, avgh6, and avgh12 in F8. Conversely, variables such as (h-24) and (h-48) in models F3 and F10, as well as (h-24), (h-48), and (h-72) in model F4, did not contribute positively to prediction accuracy. This variance underscores the importance of identifying

340    effective variable combinations while avoiding those that introduce noise into the model.



**Figure 5: The Highest prediction results for (a) U component (DJF, F1) and V component (JJA, F1) in Prediction Scheme 1, (b) U component (NOV, F1) and V component (JUL, F1) in Prediction Scheme 2, (c) U component (AUG, F1) and V component (NOV, F1) in Prediction Scheme 3**

365    **3.3.2. Performance of Prediction Scheme 2 (P2)**

Prediction Scheme 2 utilizes monthly data for training and testing, with 80% of the 2021 dataset allocated for training and 20% reserved for testing each month. The results of this scheme are summarized in Tables B1 to B4 under column P2. The

analysis focuses on the U and V components across each month from December to November (12th to 11th months) using 11 model configurations (F1 to F11), as detailed in Subchapter 2.2.2 and Table 2.

370     The root mean square error (RMSE) values for Prediction Scheme 2, presented in Tables B2 and B4 under column P2, indicate the magnitude of the prediction errors, ranging from 12.30 to 29.32 for the U component and 18.13 to 44.00 for the V component. The lowest RMSE values were recorded in model F1 for November for both components.

The correlation coefficient $R^2$ values for Prediction Scheme 2 are shown in Tables B1 and B3 under column P2, with $R^2$ values ranging from 0.90 to 0.57 for the U component and 0.94 to 0.45 for the V component. The highest $R^2$ value of 0.90 for
375   the U component was observed for F1 in November, while F1, F5, and F8 achieved an $R^2$ of 0.94 for the V component in July. Figure 5b shows the highest prediction results obtained by this scheme.

    Across all months, models F1, F5, and F8 consistently demonstrated superior RMSE and $R^2$ values compared to the other models. The elevated RMSE values and lower $R^2$ values in models other than F1, F5, and F8 suggest that certain variables enhance prediction accuracy; specifically, avgh3 in F1; avgh3 and avgh6 in F5; and avgh3, avgh6, and avgh12 in F8.
380   Conversely, variables such as (h-24) and (h-48) found in models F3, F7, and F10 did not positively contribute to prediction accuracy. This variance underscores the necessity of identifying effective variable combinations while minimizing noise within the model.

### 3.3.3. Performance of Prediction Scheme 3 (P3)

Prediction Scheme 3 employs a comprehensive dataset where all data from 2019 and 2020, along with 80% of the data from
385   2021, are used for training, while the remaining 20% of the data from 2021 serves as the testing dataset for each month. The results of this scheme are summarized in Tables B1 to B4 under column P3, analyzing both U and V components across each month from December to November (12th to 11th months) using eleven model configurations (F1 to F11), as described in Subchapter 2.2.2 and Table 2.

    The RMSE values for Prediction Scheme 3 are presented in Tables B2 and B4 under column P3, which indicates error
390   magnitudes that range from 7.68 to 26.84 for the U component and from 9.52 to 28.85 for the V component, with the lowest RMSE values recorded in model F1 during November for both components. Figure 5c shows the highest prediction results obtained by this scheme.

The correlation coefficient $R^2$ values are detailed in Tables B1 and B3 under column P3, with values ranging from 0.96 to 0.73 for the U component and 0.98 to 0.81 for the V component. The highest $R^2$ value of 0.96 was achieved by models F1 and
395   F5 in November as well as by model F1 in February for the U component; similarly, models F1, F5, F8, and F11 achieved this value in August for the V component. Throughout all the analyzed months, models F1, F5, and F8 consistently outperformed the other models in terms of RMSE and $R^2$ values.

    The elevated RMSE values alongside lower $R^2$ values noted in models outside of F1, F5, and F8 indicate that specific variables can enhance prediction accuracy—namely, avgh3 in F1; avgh3 and avgh6 in F5; as well as avgh3, avgh6, and avgh12
400   in F8—while certain variables such as (h-24), (h-48), and (h-72) present in models F4 did not contribute positively to prediction

accuracy along with (h-24) and (h-48) found in model F10. These differences highlight the importance of optimizing variable selection in prediction models to improve overall effectiveness while minimizing extraneous variables that may introduce noise.

### 3.3.4. Comprehensive Prediction Schemes Analysis

405    In this section, we present a comparative analysis of the $R^2$ correlation coefficients and RMSE for the U and V components for the three prediction schemes. Figure 6 presents a comparative analysis of the $R^2$ correlation coefficients for the U and V components for the three prediction schemes. The arrangement of graphs from left to right includes two graphs for Prediction Scheme 1, two for Prediction Scheme 2, and two for Prediction Scheme 3. The data reveals that the F1 prediction model consistently achieves the highest accuracy for both components (U and V), followed closely by models F5 and F8. Notably,
410    Prediction Scheme 1 exhibits the smallest range of $R^2$ values, indicating more consistent performance across seasons. In contrast, Prediction Scheme 2 displays the most extensive range, reflecting greater variability in prediction accuracy. Both Prediction Schemes 1 and 3 yield high R2 correlation results, reaching up to 0.96 for the U component and 0.98 for the V component. These findings underscore the significant impact of data selection on predictive quality.

     Prediction Scheme 1 benefits from using seasonal data from 2019 and 2020, which allows it to maintain a narrow range of
415    correlation results. This consistency is due to its comprehensive dataset, which includes three months of data for each season. Interestingly, although Prediction Scheme 3 utilizes less data than Scheme 1, it incorporates data from both 2019 and 2020 for the same month, effectively enhancing the model's ability to capture the characteristics of sea surface currents. This result demonstrates that carefully selecting training data can lead to high correlation results similar to those achieved in Prediction Scheme 1 even with limited data.

420    In addition to identifying factors that enhance prediction accuracy, it is crucial to recognize elements that diminish correlation values. The middle graphs in Figure 6 illustrate the results obtained by Prediction Scheme 2, which relies solely on data from 2021 for both training and testing. The R2 correlation values for May were shallow, with values of 0.57 for the U component and 0.45 for the V component. This decline can be attributed to insufficient data availability in May 2021, where only 177 data points were recorded (as can be referred to in Table A2).

425    This result highlights the necessity for future studies to ensure adequate data availability to improve prediction correlations. Almost similar to Figure 6, Figure 7 complements this analysis by comparing thw RMSE values across the three prediction schemes, while maintaining the same arrangement as in Figure 6. Similar to the correlation results, Prediction Scheme 1 exhibits the smallest error range and Prediction Scheme 2 shows the most extensive error range for both components. However, it is noteworthy that Prediction Scheme 3 achieved superior performance with RMSE values as low as 8.8 for the U component
430    and 7.68 for the V component. Despite utilizing more training data in Prediction Scheme 1, it is evident that Prediction Scheme 3 can match its high correlation while outperforming it in terms of error rates. This finding emphasized that while increasing data volume can enhance prediction outcomes, selecting appropriate data types is equally essential for reducing error rates. In addition, attention must be given to factors that can increase error rates in predictions.

**Figure 6: Comparison of R2 Correlations for Prediction Schemes 1, 2, and 3**



**Figure 7: Comparison of RMSE for Prediction Schemes 1, 2, and 3**

The results from Prediction Scheme 1 indicate that significant errors occur during specific seasons; notably, MAM for the U component and JJA for the V component exhibit the highest errors among all seasons. An examination of Table A1 reveals

470 that MAM has a standard deviation of 59.83 for component U. In contrast, JJA has an even higher standard deviation of 93.70 for component V, indicating substantial variability within these datasets. In summary, this comprehensive analysis illustrates that careful consideration of data selection and model configuration is essential for optimizing sea surface current prediction accuracy. The results affirm that high standard deviations can adversely affect prediction outcomes; therefore, researchers must thoroughly understand their datasets to mitigate potential errors effectively

475 **4. Conclusion**

This study significantly advances the field of sea surface current prediction by utilizing high-frequency (HF) radar data and developing innovative prediction schemes. By addressing a critical gap in the literature regarding seasonal and monthly data segmentation, we demonstrate how these factors influence prediction accuracy. The findings revealed that the selection of training data and prediction models, particularly those incorporating relevant variables, such as moving average variables

480 in the F1, F5, and F8 models, greatly affected correlation values and prediction errors.

The analysis of three distinct prediction schemes—seasonal (P1), monthly (P2), and comprehensive monthly (P3)— demonstrates that larger datasets yield higher correlation coefficients, whereas tailored models effectively reduce prediction errors. Notably, Prediction Scheme 1 achieved robust correlation coefficients ($R^2$) ranging from 0.94 to 0.98 for both U and V components, while Prediction Scheme 3 exhibited the lowest root mean square error (RMSE) values, showcasing the

485 effectiveness of integrating extensive datasets.

The implications of this research extend beyond theoretical advancements, offering practical applications in marine navigation, environmental monitoring, and disaster response. Accurate predictions are crucial for optimizing shipping routes and enhancing search and rescue operations. Furthermore, this study advocates integrating HF radar data into operational oceanographic systems to improve real-time monitoring capabilities.

490 In conclusion, the present research underscores the importance of data selection and model configuration for predicting sea surface currents. The insights gained contribute valuable knowledge on oceanography and highlight the potential of HF radar technology as a powerful tool for real-time ocean monitoring. As challenges from climate change and human activities escalate, accurate predictive models are increasingly critical for sustainable marine practices. Future research should explore integrating additional data sources and advanced modeling techniques to further enhance the predictive accuracy and understanding of

495 ocean dynamics.

500 article. This work is a culmination of extensive effort, and their invaluable assistance has been instrumental in bringing it to fruition. Their essential support is greatly appreciated.

**Author Contribution**

**Ittaka Aldini**: Writing – original draft, Methodology, Data curation, Conceptualization. **Adhistya Erna Permanasari**:
505 Writing – review & editing, Supervision, Methodology, Formal analysis, Conceptualization. **Risanuri Hidayat**: Writing – review & editing, Methodology, Formal analysis, Conceptualization. **Andri Ramdhani**: Writing – review & editing, Validation, Formal analysis, Conceptualization. All authors have read and approved the final version of the manuscript and agree with its submission to Ocean Sciences Journal.

510 **Competing Interests**

The authors declare no competing interests related to this work.

**Data availability**

The data used is referenced or stated in the paper.

515 **References**

Ali, M. and Prasad, R.: Significant wave height forecasting via an extreme learning machine model integrated with improved complete ensemble empirical mode decomposition, Renewable and Sustainable Energy Reviews, 104, 281–295, https://doi.org/https://doi.org/10.1016/j.rser.2019.01.014, 2019.

Ambhika, C., Anish, T. P., Dhinakaran, D., Elavarasan, E., Harish, K., and Kalyan, N. U. P.: Precisely predicting heart disease
520 by examining the data analysis, in: Advances in Networks, Intelligence and Computing, CRC Press, 72–81, 2024.

Balogun, A.-L. and Adebisi, N.: Sea level prediction using ARIMA, SVR and LSTM neural network: assessing the impact of ensemble Ocean-Atmospheric processes on models' accuracy, Geomatics, Natural Hazards and Risk, 12, 653–674, 2021.

Barrick, D., Fernandez, V., Ferrer, M. I., Whelan, C., and Breivik, Ø.: A short-term predictive system for surface currents from a rapidly deployed coastal HF radar network, Ocean Dyn, 62, 725–740, 2012.

525 Bateman, B., Jha, A. R., Johnston, B., and Mathur, I.: The The Supervised Learning Workshop: A New, Interactive Approach to Understanding Supervised Learning Algorithms, Packt Publishing Ltd, 2020.

Cosoli, S. and de Vos, S.: Interoperability of direction-finding and beam-forming high-frequency radar systems: An example from the Australian high-frequency ocean radar network, Remote Sens (Basel), 11, 291, 2019.

Hardman, R. L. and Wyatt, L. R.: Inversion of HF radar Doppler spectra using a neural network, J Mar Sci Eng, 7, 255, 2019.

530     Jirakittayakorn, A., Kormongkolkul, T., Vateekul, P., Jitkajornwanich, K., and Lawawirojwong, S.: Temporal kNN for short-term ocean current prediction based on HF radar observations, in: 2017 14th International Joint Conference on Computer Science and Software Engineering (JCSSE), 1–6, https://doi.org/10.1109/JCSSE.2017.8025921, 2017.

Jitkajornwanich, K., Vateekul, P., Gupta, U., Kormongkolkul, T., Jirakittayakorn, A., Lawawirojwong, S., and Srisonphan, S.: Ocean surface current prediction based on HF radar observations using trajectory-oriented association rule mining, in: 2017

535     IEEE International Conference on Big Data (Big Data), 4293–4300, 2017.

Kirincich, A., Emery, B., Washburn, L., and Flament, P.: Improving Surface Current Resolution Using Direction Finding Algorithms for Multiantenna High-Frequency Radars, J Atmos Ocean Technol, 36, 1997–2014, https://doi.org/https://doi.org/10.1175/JTECH-D-19-0029.1, 2019.

Kunapuli, G.: Ensemble methods for machine learning, Simon and Schuster, 2023.

540     Kusnanti, E. A., Rini Novitasari, D. C., Setiawan, F., Fanani, A., Hafiyusholeh, M., and Permata Sari, G. I.: Predicting Velocity and Direction of Ocean Surface Currents using Elman Recurrent Neural Network Method., Journal of Information Systems Engineering & Business Intelligence, 8, 2022.

Li, X., Cao, J., Guo, J., Liu, C., Wang, W., Jia, Z., and Su, T.: Multi-step forecasting of ocean wave height using gate recurrent unit networks with multivariate time series, Ocean Engineering, 248, 110689,
545     https://doi.org/https://doi.org/10.1016/j.oceaneng.2022.110689, 2022.

Marmain, J., Molcard, A., Forget, P., Barth, A., and Ourmières, Y.: Assimilation of HF radar surface currents to optimize forcing in the northwestern Mediterranean Sea, Nonlin. Processes Geophys., 21, 659–675, https://doi.org/10.5194/npg-21-659-2014, 2014.

Montaño, J., Coco, G., Antolínez, J. A. A., Beuzen, T., Bryan, K. R., Cagigal, L., Castelle, B., Davidson, M. A., Goldstein, E.
550     B., Ibaceta, R., Idier, D., Ludka, B. C., Masoud-Ansari, S., Méndez, F. J., Murray, A. B., Plant, N. G., Ratliff, K. M., Robinet, A., Rueda, A., Sénéchal, N., Simmons, J. A., Splinter, K. D., Stephens, S., Townend, I., Vitousek, S., and Vos, K.: Blind testing of shoreline evolution models, Sci Rep, 10, 2137, https://doi.org/10.1038/s41598-020-59018-y, 2020.

O'Donncha, F., Zhang, Y., Chen, B., and James, S. C.: Ensemble model aggregation using a computationally lightweight machine-learning model to forecast ocean waves, Journal of Marine Systems, 199, 103206, 2019.

555     Pascual, A., Lana, A., Troupin, C., Ruiz, S., Faugère, Y., Escudier, R., and Tintoré, J.: Assessing SARAL/AltiKa Data in the Coastal Zone: Comparisons with HF Radar Observations, Marine Geodesy, 38, 260–276, https://doi.org/10.1080/01490419.2015.1019656, 2015.

Pramesti, D. D., Novitasari, D. C. R., Setiawan, F., and Khaulasari, H.: Long-Short Term Memory (Lstm) for Predicting Velocity and Direction Sea Surface Current on Bali Strait, BAREKENG: Jurnal Ilmu Matematika dan Terapan, 16, 451–
560     462, 2022.

Putri, E. R. S., Novitasari, D. C. R., Setiawan, F., Hamid, A., Susanto, D., and Fahmi, M.: Prediction of Sea Surface Current Velocity and Direction using Gated Recurrent Unit (GRU), in: 2022 IEEE International Conference on Industry 4.0, Artificial Intelligence, and Communications Technology (IAICT), 144–149, 2022.

Ren, Y., Li, X.-M., Gao, G., and Busche, T. E.: Derivation of Sea Surface Tidal Current From Spaceborne SAR Constellation Data, IEEE Transactions on Geoscience and Remote Sensing, 55, 3236–3247, https://doi.org/10.1109/TGRS.2017.2666086, 2017.

Rezuanul Islam, Md., Duc, L., and Sawada, Y.: Assessing Storm Surge Multiscenarios Based on Ensemble Tropical Cyclone Forecasting, Journal of Geophysical Research: Atmospheres, 128, e2023JD038903, https://doi.org/https://doi.org/10.1029/2023JD038903, 2023.

Sarkar, D., Osborne, M. A., and Adcock, T. A. A.: Prediction of tidal currents using Bayesian machine learning, Ocean Engineering, 158, 221–231, https://doi.org/https://doi.org/10.1016/j.oceaneng.2018.03.007, 2018.

Saviano, S., Kalampokis, A., Zambianchi, E., and Uttieri, M.: A year-long assessment of wave measurements retrieved from an HF radar network in the Gulf of Naples (Tyrrhenian Sea, Western Mediterranean Sea), Journal of Operational Oceanography, 12, 1–15, https://doi.org/10.1080/1755876X.2019.1565853, 2019.

Schonlau, M. and Zou, R. Y.: The random forest algorithm for statistical learning, Stata J, 20, 3–29, https://doi.org/10.1177/1536867X20909688, 2020.

Tang, F. and Ishwaran, H.: Random forest missing data algorithms, Statistical Analysis and Data Mining: The ASA Data Science Journal, 10, 363–377, https://doi.org/https://doi.org/10.1002/sam.11348, 2017.

Thongniran, N., Jitkajornwanich, K., Lawawirojwong, S., Srestasathiern, P., and Vateekul, P.: Combining attentional CNN and GRU networks for ocean current prediction based on HF radar observations, in: Proceedings of the 2019 8th international conference on computing and pattern recognition, 440–446, 2019a.

Thongniran, N., Vateekul, P., Jitkajornwanich, K., Lawawirojwong, S., and Srestasathiern, P.: Spatio-temporal deep learning for ocean current prediction based on HF radar data, in: 2019 16th International Joint Conference on Computer Science and Software Engineering (JCSSE), 254–259, 2019b.

Werner, F. and Blanton, J. O.: Coastal Circulation Models☆, in: Encyclopedia of Ocean Sciences (Third Edition), edited by: Cochran, J. K., Bokuniewicz, H. J., and Yager, P. L., Academic Press, Oxford, 467–474, https://doi.org/https://doi.org/10.1016/B978-0-12-409548-9.11412-5, 2019.

Wu, H. and Levinson, D.: The ensemble approach to forecasting: A review and synthesis, Transp Res Part C Emerg Technol, 132, 103357, https://doi.org/https://doi.org/10.1016/j.trc.2021.103357, 2021.

Xiao, C., Chen, N., Hu, C., Wang, K., Gong, J., and Chen, Z.: Short and mid-term sea surface temperature prediction using time-series satellite data and LSTM-AdaBoost combination approach, Remote Sens Environ, 233, 111358, 2019.

Zhao, C., Chen, Z., Li, J., Ding, F., Huang, W., and Fan, L.: Validation and evaluation of a ship echo-based array phase manifold calibration method for HF surface wave radar DOA estimation and current measurement, Remote Sens (Basel), 12, 2761, 2020.

Zulfa, I. I., Candra, D., Novitasari, R., Setiawan, F., Fanani, A., and Hafiyusholeh, M.: Prediction of sea surface current velocity and direction using LSTM, Indones. J. Electron. Instrum. Syst, 11, 93–102, 2021.

# Appendix A

600

**Table A1. Description of Train and Test Data for Prediction Scheme 1**

| PARAMETER | TRAIN | | | | | | | | TEST | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DJF | | MAM | | JJA | | SON | | DJF | | MAM | | JJA | | SON | |
| | U | V | U | V | U | V | U | V | U | V | U | V | U | V | U | V |
| COUNT | 3436 | 3436 | 3587 | 3587 | 3070 | 3070 | 2968 | 2968 | 1586 | 1586 | 1457 | 1457 | 1496 | 1496 | 1699 | 1699 |
| MEAN | -29.96 | 41.85 | -17.08 | 23.76 | -4.27 | -9.77 | -19.33 | 17.28 | -24.08 | 29.29 | -22.67 | 30.44 | -10.01 | -5.65 | -20.55 | 19.95 |
| STD | 53.66 | 62.06 | 59.83 | 72.11 | 52.01 | 93.70 | 52.24 | 77.27 | 45.27 | 59.57 | 54.49 | 70.69 | 45.71 | 96.18 | 45.59 | 77.16 |
| MIN | -263.74 | -206.37 | -262.34 | -267.36 | -170.48 | -233.29 | -271.66 | -247.81 | -188.25 | -183.76 | -259.71 | -234.99 | -162.48 | -248.33 | -199.31 | -214.43 |
| MAX | 219.48 | 211.53 | 299.73 | 229.30 | 191.63 | 203.18 | 232.80 | 229.89 | 147.63 | 230.39 | 201.97 | 210.94 | 200.25 | 228.88 | 164.38 | 214.91 |

**Table A2. Description of data training for Prediction Scheme 2**

| TRAIN | DEC | | JAN | | FEB | | MAR | |
|---|---|---|---|---|---|---|---|---|
| | U | V | U | V | U | V | U | V |
| COUNT | 320 | 320 | 492 | 492 | 455 | 455 | 490 | 490 |
| MEAN | -27.15 | 35.66 | 26.27 | 29.80 | 22.57 | 28.11 | -23.58 | 38.86 |
| STD | 40.81 | 54.20 | 42.59 | 58.18 | 51.93 | 59.91 | 57.59 | 59.56 |
| MIN | -137.27 | -133.06 | -187.23 | -152.87 | -188.25 | -122.82 | -242.93 | -153.16 |
| MAX | 118.57 | 160.54 | 104.49 | 187.25 | 145.98 | 199.142 | 139.43 | 210.94 |
| TRAIN | APR | | MAY | | JUN | | JUL | |
| | U | V | U | V | U | V | U | V |
| COUNT | 497 | 497 | 177 | 177 | 293 | 293 | 595 | 595 |
| MEAN | -24.19 | 30.76 | -3.03 | -12.40 | -26.81 | 36.84 | -9.55 | -10.07 |
| STD | 50.87 | 66.92 | 60.33 | 89.43 | 52.21 | 76.27 | 42.12 | 98.31 |
| MIN | -259.71 | -163.82 | -125.75 | -234.99 | -162.48 | -165.43 | -130.41 | -248.33 |
| MAX | 201.97 | 205.03 | 194.82 | 168.15 | 200.25 | 193.88 | 167.21 | 205.06 |
| TRAIN | AUG | | SEP | | OCT | | NOV | |
| | U | V | U | V | U | V | U | V |
| COUNT | 308 | 308 | 415 | 415 | 439 | 439 | 504 | 504 |
| MEAN | 2.30 | -37.23 | 12.31 | 6.28 | -23.07 | 16.26 | -27.09 | 36.38 |
| STD | 412.89 | 90.83 | 51.42 | 91.31 | 41.97 | 84.50 | 43.94 | 56.14 |
| MIN | -112.51 | -219.15 | -153.55 | -216.43 | -154.79 | -209.75 | -199.31 | -122.19 |
| MAX | 179.3 | 176.09 | 149.1 | 214.91 | 164.38 | 210.52 | 90.143 | 165.64 |

605

610

**Table A3. Description of data testing for Prediction Schemes 2 & 3**

| TEST | DEC | | JAN | | FEB | | MAR | |
|---|---|---|---|---|---|---|---|---|
| | U | V | U | V | U | V | U | V |
| COUNT | 81 | 81 | 124 | 124 | 114 | 114 | 123 | 123 |
| MEAN | -23.75 | 24.25 | -18.64 | 29.15 | -18.29 | 17.61 | -29.86 | 43.57 |
| STD | 46.98 | 59.65 | 48.59 | 60.39 | 48.54 | 74.80 | 51.04 | 67.77 |
| MIN | -133.05 | -108.38 | -142.58 | -99.63 | -151.00 | -183.76 | -155.32 | -172.55 |
| MAX | 89.95 | 141.51 | 92.14 | 183.07 | 147.63 | 230.39 | 105.13 | 200.72 |
| TEST | APR | | MAY | | JUN | | JUL | |
| | U | V | U | V | U | V | U | V |
| COUNT | 125 | 125 | 45 | 45 | 74 | 74 | 149 | 149 |
| MEAN | -32.86 | 50.33 | -25.44 | 12.61 | -9.01 | -0.96 | -7.93 | -0.16 |
| STD | 50.08 | 81.33 | 36.68 | 44.87 | 52.19 | 96.00 | 44.34 | 107.89 |
| MIN | -198.36 | -191.59 | -91.29 | -55.81 | -133.55 | -186.76 | -99.08 | -228.24 |
| MAX | 174.00 | 205.35 | 47.99 | 114.79 | 166.97 | 188.48 | 173.915 | 228.88 |
| TEST | AUG | | SEP | | OCT | | NOV | |
| | U | V | U | V | U | V | U | V |
| COUNT | 77 | 77 | 104 | 104 | 110 | 110 | 127 | 127 |
| MEAN | -3.85 | -21.92 | -13.36 | -13.81 | -17.08 | 26.41 | -21.68 | 34.21 |
| STD | 38.35 | 87.87 | 43.66 | 79.38 | 48.07 | 72.06 | 38.32 | 53.49 |
| MIN | -89.08 | -176.97 | -126.28 | -215.69 | -103.96 | -161.76 | -113.28 | -96.90 |
| MAX | 84.47 | 148.44 | 97.73 | 142.36 | 133.05 | 149.95 | 95.05 | 140.00 |

**Table A4. Description of data training for Prediction Scheme 3**

| TRAIN | DEC | | JAN | | FEB | | MAR | |
|---|---|---|---|---|---|---|---|---|
| | U | V | U | V | U | V | U | V |
| COUNT | 1263 | 1263 | 1734 | 1734 | 1706 | 1706 | 1519 | 1519 |
| MEAN | -28.43 | 34.54 | -30.26 | 41.65 | -27.24 | 39.15 | -23.69 | 36.28 |
| STD | 47.02 | 61.57 | 51.84 | 61.89 | 54.74 | 59.95 | 60.63 | 62.49 |
| MIN | -210.30 | -197.15 | -227.20 | -152.87 | -263.74 | -206.37 | -262.34 | -243.86 |
| MAX | 181.35 | 187.27 | 167.50 | 211.53 | 219.48 | 199.14 | 299.73 | 213.72 |
| TRAIN | APR | | MAY | | JUN | | JUL | |
| | U | V | U | V | U | V | U | V |
| COUNT | 1937 | 1937 | 1295 | 1295 | 810 | 810 | 1660 | 1660 |
| MEAN | -16.88 | 29.32 | -12.91 | 4.23 | -15.03 | 12.60 | -0.65 | -21.61 |
| STD | 60.86 | 68.02 | 52.91 | 82.00 | 53.08 | 91.73 | 47.02 | 95.21 |
| MIN | -259.71 | -238.6 | -199.96 | -267.36 | -162.48 | -213.61 | -151.94 | -248.33 |
| MAX | 292.17 | 207.94 | 241.98 | 229.3 | 200.25 | 200.65 | 191.63 | 205.06 |
| TRAIN | AUG | | SEP | | OCT | | NOV | |
| | U | V | U | V | U | V | U | V |
| COUNT | 1796 | 1796 | 1448 | 1448 | 1216 | 1216 | 1662 | 1662 |
| MEAN | -7.06 | -6.13 | -16.97 | 5.81 | -22.47 | 20.50 | -20.67 | 27.71 |
| STD | 51.67 | 92.27 | 48.33 | 84.19 | 45.38 | 79.37 | 55.33 | 68.62 |
| MIN | -170.48 | -233.29 | -164.55 | -232.9 | -203.68 | -209.75 | -271.66 | -247.81 |
| MAX | 183.28 | 203.18 | 170.68 | 214.91 | 164.38 | 210.52 | 232.80 | 229.89 |

# Appendix B

## Table B1. R² values for the U component across Prediction Schemes 1, 2, and 3

| U | 12 (D) | | | 1 (J) | | | 2 (F) | | | 3 (M) | | | 4 (A) | | | 5 (M) | | | 6 (J) | | | 7 (J) | | | 8 (A) | | | 9 (S) | | | 10 (O) | | | 11 (N) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 |
| F1 | 0.96 | 0.8 | 0.93 | | 0.85 | 0.95 | | 0.89 | 0.96 | 0.95 | 0.85 | 0.94 | | 0.75 | 0.89 | | 0.76 | 0.93 | 0.95 | 0.74 | 0.83 | | 0.82 | 0.94 | | 0.82 | 0.93 | 0.96 | 0.83 | 0.92 | | 0.84 | 0.88 | | 0.9 | 0.96 |
| F2 | 0.96 | 0.77 | 0.92 | | 0.84 | 0.94 | | 0.88 | 0.95 | 0.94 | 0.84 | 0.93 | | 0.74 | 0.88 | | 0.63 | 0.89 | 0.95 | 0.71 | 0.78 | | 0.8 | 0.92 | | 0.81 | 0.93 | 0.95 | 0.81 | 0.91 | | 0.83 | 0.87 | | 0.88 | 0.95 |
| F3 | 0.95 | 0.77 | 0.92 | | 0.82 | 0.93 | | 0.86 | 0.95 | 0.94 | 0.83 | 0.93 | | 0.72 | 0.87 | | 0.57 | 0.89 | 0.94 | 0.7 | 0.75 | | 0.79 | 0.92 | | 0.8 | 0.92 | 0.95 | 0.8 | 0.91 | | 0.81 | 0.86 | | 0.87 | 0.94 |
| F4 | 0.95 | 0.76 | 0.91 | | 0.8 | 0.93 | | 0.86 | 0.94 | 0.94 | 0.82 | 0.92 | | 0.72 | 0.86 | | 0.58 | 0.88 | 0.94 | 0.69 | 0.73 | | 0.78 | 0.92 | | 0.79 | 0.93 | 0.94 | 0.8 | 0.9 | | 0.81 | 0.86 | | 0.86 | 0.94 |
| F5 | 0.96 | 0.81 | 0.92 | | 0.83 | 0.94 | | 0.88 | 0.95 | 0.94 | 0.85 | 0.93 | | 0.73 | 0.88 | | 0.74 | 0.92 | 0.95 | 0.71 | 0.82 | | 0.81 | 0.93 | | 0.81 | 0.93 | 0.96 | 0.81 | 0.91 | | 0.84 | 0.88 | | 0.89 | 0.96 |
| F6 | 0.96 | 0.78 | 0.91 | | 0.82 | 0.94 | | 0.87 | 0.95 | 0.94 | 0.84 | 0.93 | | 0.72 | 0.87 | | 0.64 | 0.89 | 0.94 | 0.69 | 0.77 | | 0.79 | 0.91 | | 0.8 | 0.93 | 0.95 | 0.8 | 0.91 | | 0.82 | 0.87 | | 0.87 | 0.95 |
| F7 | 0.95 | 0.77 | 0.9 | | 0.81 | 0.93 | | 0.86 | 0.94 | 0.94 | 0.83 | 0.92 | | 0.71 | 0.86 | | 0.58 | 0.88 | 0.94 | 0.69 | 0.74 | | 0.79 | 0.91 | | 0.79 | 0.92 | 0.94 | 0.79 | 0.9 | | 0.81 | 0.86 | | 0.86 | 0.94 |
| F8 | 0.96 | 0.79 | 0.91 | | 0.82 | 0.94 | | 0.88 | 0.95 | 0.94 | 0.83 | 0.92 | | 0.71 | 0.88 | | 0.73 | 0.92 | 0.95 | 0.71 | 0.81 | | 0.81 | 0.93 | | 0.81 | 0.93 | 0.95 | 0.8 | 0.91 | | 0.82 | 0.88 | | 0.88 | 0.95 |
| F9 | 0.95 | 0.76 | 0.91 | | 0.81 | 0.93 | | 0.86 | 0.94 | 0.94 | 0.83 | 0.92 | | 0.7 | 0.87 | | 0.62 | 0.89 | 0.94 | 0.69 | 0.76 | | 0.79 | 0.91 | | 0.81 | 0.93 | 0.95 | 0.79 | 0.91 | | 0.81 | 0.86 | | 0.87 | 0.94 |
| F10 | 0.95 | 0.76 | 0.9 | | 0.8 | 0.93 | | 0.86 | 0.94 | 0.94 | 0.82 | 0.92 | | 0.69 | 0.85 | | 0.57 | 0.88 | 0.94 | 0.68 | 0.74 | | 0.78 | 0.91 | | 0.79 | 0.91 | 0.94 | 0.78 | 0.9 | | 0.8 | 0.85 | | 0.86 | 0.94 |
| F11 | 0.96 | 0.78 | 0.91 | | 0.8 | 0.94 | | 0.87 | 0.95 | 0.94 | 0.82 | 0.92 | | 0.7 | 0.87 | | 0.72 | 0.91 | 0.94 | 0.69 | 0.81 | | 0.8 | 0.93 | | 0.8 | 0.92 | 0.95 | 0.8 | 0.91 | | 0.82 | 0.88 | | 0.88 | 0.95 |

## Table B2. RMSE values for the U component in Prediction Schemes 1, 2, and 3

| U | 12 (D) | | | 1 (J) | | | 2 (F) | | | 3 (M) | | | 4 (A) | | | 5 (M) | | | 6 (J) | | | 7 (J) | | | 8 (A) | | | 9 (S) | | | 10 (O) | | | 11 (N) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 |
| F1 | 8.8 | 20.7 | 12.1 | | 18.8 | 11.1 | | 16.1 | 9.6 | 12.5 | 19.5 | 12.1 | | 24.8 | 16.5 | | 17.6 | 9.9 | 10.0 | 26.5 | 21.6 | | 18.9 | 11.1 | | 16.3 | 9.8 | 9.2 | 18.0 | 12.3 | | 18.9 | 16.4 | | 12.3 | 7.7 |
| F2 | 9.5 | 22.4 | 13.0 | | 19.5 | 12.1 | | 17.0 | 10.6 | 12.8 | 20.1 | 13.1 | | 25.5 | 17.5 | | 22.1 | 11.9 | 10.6 | 28.1 | 24.5 | | 19.9 | 12.5 | | 16.7 | 10.0 | 10.1 | 19.1 | 12.8 | | 19.8 | 17.3 | | 13.0 | 8.5 |
| F3 | 9.9 | 22.6 | 13.6 | | 20.5 | 12.6 | | 17.8 | 10.9 | 13.5 | 20.9 | 13.9 | | 26.3 | 18.3 | | 23.8 | 12.3 | 11.0 | 28.6 | 26.1 | | 20.1 | 12.5 | | 17.1 | 10.5 | 10.7 | 19.5 | 13.1 | | 20.6 | 17.7 | | 13.7 | 9.0 |
| F4 | 10.2 | 22.7 | 14.2 | | 21.4 | 13.0 | | 18.1 | 11.5 | 13.7 | 21.3 | 14.2 | | 26.5 | 18.7 | | 23.6 | 12.3 | 11.2 | 28.8 | 26.8 | | 20.6 | 12.5 | | 17.3 | 10.3 | 10.8 | 19.6 | 13.5 | | 20.9 | 17.8 | | 14.2 | 9.3 |
| F5 | 9.0 | 20.6 | 13.1 | | 20.0 | 11.6 | | 16.8 | 10.3 | 12.8 | 20.0 | 13.2 | | 25.7 | 16.9 | | 18.4 | 10.0 | 10.3 | 27.8 | 22.2 | | 19.3 | 11.5 | | 16.6 | 9.9 | 9.6 | 18.9 | 12.8 | | 19.3 | 16.5 | | 12.6 | 8.1 |
| F6 | 9.7 | 22.1 | 13.8 | | 20.4 | 12.2 | | 17.5 | 11.3 | 13.1 | 20.6 | 13.8 | | 26.4 | 17.7 | | 21.8 | 12.1 | 10.8 | 29.1 | 24.9 | | 20.0 | 13.0 | | 16.9 | 10.2 | 10.4 | 19.7 | 13.2 | | 20.4 | 17.5 | | 13.5 | 8.8 |
| F7 | 10.1 | 22.5 | 14.5 | | 21.1 | 12.8 | | 18.3 | 11.4 | 13.6 | 21.3 | 14.4 | | 26.9 | 18.5 | | 23.6 | 12.4 | 11.1 | 28.9 | 26.2 | | 20.4 | 13.0 | | 17.4 | 10.7 | 10.8 | 20.1 | 13.5 | | 21.0 | 17.9 | | 14.2 | 9.1 |
| F8 | 9.4 | 21.5 | 13.7 | | 20.8 | 11.8 | | 17.0 | 10.8 | 13.1 | 21.0 | 14.0 | | 26.9 | 17.3 | | 18.8 | 10.3 | 10.6 | 28.0 | 22.6 | | 19.3 | 11.7 | | 16.8 | 10.2 | 9.8 | 19.2 | 12.9 | | 20.3 | 16.7 | | 13.1 | 8.2 |
| F9 | 9.9 | 22.9 | 14.0 | | 21.0 | 12.6 | | 17.8 | 11.3 | 13.4 | 21.3 | 14.4 | | 27.3 | 18.1 | | 22.4 | 12.3 | 11.1 | 28.8 | 25.3 | | 20.5 | 13.4 | | 16.8 | 10.4 | 10.5 | 19.8 | 13.3 | | 21.0 | 17.8 | | 13.9 | 9.0 |
| F10 | 10.3 | 22.7 | 14.9 | | 21.6 | 12.8 | | 18.4 | 11.8 | 13.9 | 21.9 | 14.8 | | 27.9 | 19.0 | | 23.8 | 12.5 | 11.4 | 29.3 | 26.6 | | 20.6 | 13.3 | | 17.4 | 11.2 | 11.1 | 20.2 | 13.9 | | 21.3 | 18.3 | | 14.5 | 9.4 |
| F11 | 9.5 | 21.7 | 14.1 | | 21.6 | 12.2 | | 17.7 | 10.6 | 13.5 | 21.3 | 14.1 | | 27.2 | 17.8 | | 19.4 | 10.7 | 10.8 | 28.9 | 22.8 | | 19.8 | 11.7 | | 17.2 | 10.5 | 10.1 | 19.5 | 12.9 | | 20.6 | 16.7 | | 13.3 | 9.0 |

615

**Table B3. R² values for the V component across Prediction Schemes 1, 2, and 3**

| V | 12 (D) | | | 1 (J) | | | 2 (F) | | | 3 (M) | | | 4 (A) | | | 5 (M) | | | 6 (J) | | | 7 (J) | | | 8 (A) | | | 9 (S) | | | 10 (O) | | | 11 (N) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 |
| F1 | 0.96 | 0.9 | 0.96 | | 0.88 | 0.96 | | 0.79 | 0.87 | 0.96 | 0.84 | 0.92 | | 0.8 | 0.91 | | 0.6 | 0.88 | 0.98 | 0.86 | 0.93 | | 0.94 | 0.97 | | 0.89 | 0.98 | 0.97 | 0.87 | 0.94 | | 0.9 | 0.95 | | 0.88 | 0.97 |
| F2 | 0.96 | 0.86 | 0.95 | | 0.86 | 0.95 | | 0.78 | 0.87 | 0.96 | 0.83 | 0.92 | | 0.79 | 0.9 | | 0.53 | 0.86 | 0.97 | 0.81 | 0.92 | | 0.93 | 0.95 | | 0.85 | 0.97 | 0.96 | 0.83 | 0.92 | | 0.82 | 0.92 | | 0.86 | 0.96 |
| F3 | 0.96 | 0.85 | 0.94 | | 0.85 | 0.95 | | 0.76 | 0.87 | 0.96 | 0.82 | 0.91 | | 0.79 | 0.9 | | 0.48 | 0.83 | 0.97 | 0.8 | 0.91 | | 0.92 | 0.95 | | 0.86 | 0.97 | 0.96 | 0.82 | 0.91 | | 0.81 | 0.92 | | 0.85 | 0.95 |
| F4 | 0.96 | 0.85 | 0.94 | | 0.84 | 0.94 | | 0.75 | 0.86 | 0.95 | 0.81 | 0.9 | | 0.78 | 0.9 | | 0.49 | 0.83 | 0.97 | 0.79 | 0.91 | | 0.92 | 0.94 | | 0.85 | 0.97 | 0.96 | 0.81 | 0.9 | | 0.81 | 0.92 | | 0.84 | 0.95 |
| F5 | 0.96 | 0.89 | 0.96 | | 0.88 | 0.96 | | 0.78 | 0.87 | 0.96 | 0.82 | 0.92 | | 0.79 | 0.91 | | 0.59 | 0.88 | 0.97 | 0.84 | 0.93 | | 0.94 | 0.96 | | 0.88 | 0.98 | 0.96 | 0.85 | 0.93 | | 0.88 | 0.95 | | 0.87 | 0.97 |
| F6 | 0.96 | 0.85 | 0.94 | | 0.86 | 0.95 | | 0.77 | 0.87 | 0.96 | 0.81 | 0.91 | | 0.79 | 0.9 | | 0.51 | 0.84 | 0.97 | 0.8 | 0.92 | | 0.92 | 0.95 | | 0.85 | 0.97 | 0.96 | 0.82 | 0.91 | | 0.81 | 0.92 | | 0.86 | 0.96 |
| F7 | 0.96 | 0.85 | 0.94 | | 0.85 | 0.95 | | 0.76 | 0.86 | 0.96 | 0.8 | 0.91 | | 0.78 | 0.9 | | 0.45 | 0.82 | 0.97 | 0.8 | 0.91 | | 0.92 | 0.94 | | 0.85 | 0.97 | 0.96 | 0.8 | 0.91 | | 0.81 | 0.91 | | 0.84 | 0.95 |
| F8 | 0.96 | 0.86 | 0.95 | | 0.87 | 0.95 | | 0.77 | 0.88 | 0.96 | 0.81 | 0.91 | | 0.78 | 0.91 | | 0.58 | 0.87 | 0.97 | 0.83 | 0.93 | | 0.94 | 0.96 | | 0.87 | 0.98 | 0.96 | 0.85 | 0.93 | | 0.88 | 0.95 | | 0.86 | 0.96 |
| F9 | 0.96 | 0.84 | 0.94 | | 0.85 | 0.95 | | 0.76 | 0.87 | 0.96 | 0.8 | 0.91 | | 0.78 | 0.9 | | 0.49 | 0.83 | 0.97 | 0.79 | 0.91 | | 0.92 | 0.94 | | 0.84 | 0.97 | 0.96 | 0.81 | 0.91 | | 0.81 | 0.91 | | 0.85 | 0.95 |
| F10 | 0.95 | 0.84 | 0.94 | | 0.84 | 0.94 | | 0.75 | 0.86 | 0.95 | 0.8 | 0.9 | | 0.78 | 0.9 | | 0.46 | 0.81 | 0.97 | 0.79 | 0.91 | | 0.91 | 0.94 | | 0.85 | 0.97 | 0.95 | 0.8 | 0.9 | | 0.8 | 0.91 | | 0.84 | 0.95 |
| F11 | 0.96 | 0.85 | 0.95 | | 0.86 | 0.95 | | 0.76 | 0.87 | 0.96 | 0.81 | 0.91 | | 0.78 | 0.9 | | 0.58 | 0.86 | 0.97 | 0.83 | 0.93 | | 0.93 | 0.96 | | 0.87 | 0.98 | 0.96 | 0.83 | 0.93 | | 0.88 | 0.94 | | 0.85 | 0.96 |

**Table B4. RMSE values for the V component in Prediction Schemes 1, 2, and 3**

| V | 12 (D) | | | 1 (J) | | | 2 (F) | | | 3 (M) | | | 4 (A) | | | 5 (M) | | | 6 (J) | | | 7 (J) | | | 8 (A) | | | 9 (S) | | | 10 (O) | | | 11 (N) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 | P1 | P2 | P3 |
| F1 | **11.2** | 18.7 | 11.6 | | 20.6 | 12.3 | | 34.2 | 26.5 | 13.5 | 27.1 | 18.8 | | 35.9 | 23.9 | | 28.1 | 15.1 | 15 | 36.2 | 24.8 | | 25.3 | 20 | | 29.4 | 12.6 | 14 | 28.4 | 19.4 | | 23.2 | 16.5 | | **18.1** | **9.52** |
| F2 | 11.8 | 22.3 | 13.9 | | 22.3 | 13.1 | | 34.9 | 26.8 | 14.3 | 28 | 19.4 | | 36.4 | 25.2 | | 30.5 | 16.8 | 16.2 | 42 | 27 | | 29 | 24.5 | | 33.3 | 15 | 15 | 32.9 | 22.9 | | 30.2 | 20 | | 19.7 | 10.8 |
| F3 | 12.2 | 22.7 | 14.4 | | 23.4 | 13.8 | | 36.4 | 27.3 | 14.8 | 28.9 | 20.2 | | 37.2 | 25.4 | | 32 | 18.1 | 16.9 | 42.3 | 28.2 | | 30.6 | 24.9 | | 33.1 | 14.9 | 15.7 | 33.6 | 23.7 | | 31 | 20.4 | | 20.8 | 11.6 |
| F4 | 12.4 | 23.3 | 14.7 | | 24.1 | 14.1 | | 37.2 | 27.5 | 15.1 | 29.7 | 21.2 | | 37.9 | 25.6 | | 31.8 | 18.4 | 17.6 | 43.5 | 28.7 | | 31.3 | 25.7 | | 33.5 | 15.1 | 16.1 | 34.7 | 24.4 | | 31.6 | 20.9 | | 21.3 | 11.9 |
| F5 | 11.6 | 19.9 | 12.4 | | 20.9 | 12.6 | | 35.1 | 26.6 | 13.9 | 28.6 | 19.2 | | 37.2 | 24.5 | | 28.6 | 15.6 | 15.6 | 37.7 | 24.7 | | 26.1 | 20.7 | | 29.8 | 12.9 | 14.5 | 30.2 | 20.2 | | 24.4 | 16.8 | | 19.1 | 9.8 |
| F6 | 12.1 | 23 | 14.1 | | 22.7 | 13.3 | | 35.8 | 27 | 14.6 | 29.1 | 19.9 | | 37.1 | 25.3 | | 30.9 | 17.7 | 16.7 | 42.6 | 27.3 | | 30.4 | 24.9 | | 33.8 | 15.4 | 15.5 | 33.9 | 23.6 | | 30.9 | 20.6 | | 20.3 | 11.3 |
| F7 | 12.4 | 23.2 | 14.7 | | 23.6 | 13.9 | | 36.5 | 27.4 | 15 | 30.1 | 20.5 | | 37.9 | 25.6 | | 32.8 | 18.9 | 17.3 | 42.9 | 28.5 | | 31.3 | 25.7 | | 33.5 | 15.4 | 16.2 | 34.9 | 24.3 | | 31.6 | 21 | | 21.3 | 12 |
| F8 | 11.9 | 22 | 12.6 | | 21.8 | 13 | | 36.1 | 26.3 | 14.3 | 29.1 | 20 | | 38.1 | 24.9 | | 28.8 | 16.3 | 16.2 | 39.3 | 25.5 | | 27.2 | 21.5 | | 31 | 13.4 | 15 | 30.8 | 20.9 | | 24.8 | 16.7 | | 19.8 | 10.4 |
| F9 | 12.3 | 24.1 | 14.8 | | 23 | 13.7 | | 36.7 | 27 | 14.8 | 30 | 20.4 | | 37.8 | 26 | | 31.5 | 18.4 | 17.2 | 43.7 | 27.9 | | 31.3 | 25.3 | | 34.7 | 15.6 | 15.8 | 34.9 | 24 | | 31.5 | 21.1 | | 20.8 | 11.6 |
| F10 | 12.7 | 24 | 15.1 | | 24.3 | 14.2 | | 37.5 | 27.5 | 15.3 | 30.2 | 21.3 | | 38.1 | 25.9 | | 32.8 | 19.3 | 17.6 | 44 | 28.9 | | 31.9 | 25.7 | | 34.3 | 15.5 | 16.4 | 35.7 | 24.6 | | 32 | 21.4 | | 21.6 | 12.1 |
| F11 | 12.1 | 22.6 | 12.7 | | 22.4 | 13.5 | | 36.2 | 27 | 14.7 | 29.5 | 20.5 | | 38.2 | 25.2 | | 28.7 | 16.7 | 16.6 | 39.4 | 26.1 | | 28.1 | 22.1 | | 31.2 | 13.5 | 15.2 | 32.4 | 21.3 | | 24.9 | 17.1 | | 20.3 | 10.8 |

a. **Color density label: Bold:** Best value at each prediction scheme; Light color: Higher quality / lower error ; Darker: Lower quality / higher error

620