



How well can we forecast local magnetic ground perturbations with existing space weather monitoring resources?

Stephen Omondi¹, Spencer Mark Hatch², Andreas Kvammen³, Magnar Gullikstad Johnsen³, Mathew J. Owens⁴, Kristian Solheim Thinn⁵, and Rodrigo López⁶

¹Geophysical Institute, University of Alaska Fairbanks, Fairbanks, AK, USA

²Department of Physics and Technology, University of Bergen, Bergen, Norway

³Tromsø Geophysical Observatory (TGO), University of Tromsø, Norway

⁴Department of Meteorology, University of Reading, UK

⁵Electric Power Engineering, SINTEF Energi AS, Trondheim, Norway

⁶Statnett, Oslo, Norway

Correspondence: Spencer Mark Hatch (spencer.hatch@uib.no)

Abstract. In this study we examine how a deep-learning based forecast of local, ground-based geomagnetic field variations trained on solar wind parameters available in real time might be improved by including information contained in an accurate forecast of solar wind conditions. This is accomplished using a long short-term memory (LSTM) model together with magnetic field measurements made at the Rørvik magnetometer station in Mid-Norway. We use Advanced Composition Explorer (ACE) satellite measurements of solar wind and interplanetary magnetic field (IMF) conditions at the first Sun-Earth Lagrange point, and historical lists of coronal mass ejection (CME) impacts at Earth to train and validate the LSTM model. We find that accurate information about the IMF B_z component and solar wind speed are important for obtaining a reasonably accurate ($r^2 \geq 0.5$) forecast of local geomagnetic activity over forecasting horizons beyond ~ 3 h. Information about CME arrival time is only important when simultaneously accompanied by accurate, relatively high-resolution information about IMF B_z . In the absence of the latter, CME arrival time information does not contribute to model performance. This empirical result amounts to a quantitative demonstration of the widely recognized impact of IMF orientation on CME geoeffectiveness. This result also highlights that new innovations, probably in the form of new prediction capabilities of conditions in interplanetary space, will be required to produce accurate forecasts of local geomagnetic disturbances beyond a forecast horizon of 1 h.

1 Introduction

Geomagnetic field disturbances at Earth's surface are driven by interactions between the solar wind and Earth's coupled magnetosphere-ionosphere system. The intensity of these disturbances is strongly tied to variations in the solar wind speed and the interplanetary magnetic field (IMF) that is borne by the solar wind. The most intense disturbances are associated with transient solar phenomena such as coronal mass ejections (CMEs), with the degree of influence that a particular CME has on space weather conditions in near-Earth geospace—its so-called "geoeffectiveness"—critically dependent on IMF orientation upon arrival at Earth. Periods of strongly negative (southward) IMF B_z are invariably connected with the most intense dis-



turbances geomagnetic storm activity (e.g., Gonzalez et al., 2011; Lakhina and Tsurutani, 2016; Sierra-Porta et al., 2024, and references therein).

At the surface of Earth and in its interior, geomagnetic disturbances give rise to geomagnetically induced electric fields (GEFs) that themselves give rise to geomagnetically induced currents (GICs). GICs can have deleterious effects on critical ground infrastructure including power grids, gas pipelines, and some navigation and communication systems (Ngwira and Pulkkinen, 2019; Press, 2021; CIGRE, 2019; EPRI, 2020; Patterson et al., 2023). Operational forecasts of geomagnetic disturbances, GEFs, and GICs are therefore naturally a central objective in space weather research (see Hapgood, 2011; Fry, 2012; Pulkkinen et al., 2017; Abda et al., 2020; Press, 2021, and references therein).

Forecasting of GEFs and GICs, as well as proxies for them, has seen an explosion of interest during the past decade or so (Hapgood, 2011; Tóth et al., 2014; Pulkkinen, 2015; Keesee et al., 2020; Siddique and Mahmud, 2022; Conde et al., 2023; Wang et al., 2025), likely as a result of increasing awareness of the threat posed by GICs and space weather generally (Fry, 2012; Gannon et al., 2019; Abda et al., 2020) as well as increasing awareness and acceptance of machine learning-based solutions (Camporeale et al., 2018; Camporeale, 2019).

With few exceptions, forecasting methods rely on measurements of solar wind speed and density as well as the strength and orientation of the IMF made by solar wind monitors such as the Advanced Composition Explorer (ACE) and Deep Space Climate Observatory (DSCOVR) satellites located at L1, the first Sun-Earth Lagrange point. As L1 is $\sim 1.5 \times 10^6$ km from Earth, such measurements provide lead times of ~ 15 –60 min depending on solar wind speed. The accuracy of forecasts that rely on measurements at L1 therefore degrades rapidly with increasing forecasting horizon beyond 1 h.

One source of information that, at present, is little used in quantitative geomagnetic forecasting models described in the literature are outputs from solar wind models such as the Wang-Sheeley-Argge (WSA)-Enlil model (Argge and Pizzo, 2000; Argge et al., 2004; Odstroil et al., 2004) and the heliospheric upwind extrapolation with time-dependence (HUXt) model (Barnard and Owens, 2022). The former provides a forecast of solar wind density and radial velocity. The latter produces a probabilistic forecast of solar wind speed, and a probability distribution of the time of coronal mass ejection (CME) impact at Earth for any currently active CMEs.

In this study we seek to answer the following question: How accurately could one predict local geomagnetic activity if high-quality forecasts of solar wind properties such as speed, density, and IMF strength and orientation, as well as precise information about CME arrival time, were available 24 h in advance?

We answer this question via a deep learning-based forecasting model of the spectral power of dH/dt (the time rate of change of horizontal magnetic field disturbances, a well known proxy for ground induced currents) that is provided with four main sources of data: (a) ambient parameters such as time of day and the tilt of Earth's magnetic dipole; (b) real-time measurements of local magnetic disturbances measured by a ground-based magnetometer station in Rørvik, Norway; (c) real-time solar wind and interplanetary magnetic field properties measured by the ACE satellite at L1; (d) "perfect" forecasts of solar wind and IMF properties as well as CME arrival times. The latter are "perfect" in the sense that they are based on actual solar wind and IMF observations and CME arrival times at Earth.



55 This paper is organized as follows: In Section 2 we present the datasets used to train the model as well as the methodologies used for producing a training dataset. In Section 3 we present the model architecture. In Section 4 we present a comparison of model predictions with observations for three storms in 2021, 2023, and 2024. In Section 5 we discuss our results and present an evaluation of the model performance against simpler statistical models. In Section 6 we conclude our study.

2 Data

60 There are three primary sets of measurements used in this study: magnetometer measurements made at the Rørvik magnetometer station (64.95° N, 10.99° E geographic; 62.46°, 91.21° geomagnetic) maintained by the Tromsø Geophysical Observatory (Johnsen, 2025); solar wind and IMF measurements made by the ACE satellite (Advanced Composition Explorer (ACE), 2025); and the Richardson and Cane (2024) CME list, which is extended from the original work of Cane and Richardson (2003). All of these datasets extend over the 15-year time period (2010–2024) that is considered in this study.

65 Our choice to focus on Rørvik magnetometer measurements is based on information that ground-based conducting infrastructure close to the area was responsive to elevated levels of geomagnetic activity in 2021 (Statnett, private communication, 2025). We therefore specifically examine dH/dt , the time rate of change of the horizontal magnetic field and a commonly used proxy for GICs Viljanen et al. (2001), as estimated from Rørvik magnetometer measurements.

2.1 GIC proxy from Rørvik magnetometer data

70 The starting point for producing model training data is the set of Rørvik magnetometer measurements with a temporal resolution of 10 s. An example time series of the horizontal component H and the time rate of change dH/dt estimated using first-order backward differencing during 3–4 November, 2021, are shown in Figures 1a and 1b, respectively. An interplanetary CME impacted Earth at 19:42 UT on 3 November, 2021, as indicated by the red dash-dotted vertical line in all panels.

Figure 1c shows the power spectral density of dH/dt up to the Nyquist frequency 0.05 Hz calculated using a Python 75 implementation (Prieto, 2022) of the multitaper technique (Thomson, 1982; Hatch and LaBelle, 2018). We use the multitaper technique because it yields a less biased estimate of the power spectral density than is obtained, for example, via Fourier transformation of the underlying time series. From Figure 1c it is clear that the power is concentrated at the lowest frequencies both before and after CME impact.

Figure 1d shows the integrated power spectral density, or spectral power, in various frequency bands. The majority of the 80 power is generally concentrated at the lowest frequencies (here represented by 0–0.01 Hz, the green solid line), denoted by P_0 in Table 1, although during periods of intense geomagnetic activity the power in other in contributions from the other frequency bands can reach or exceed the spectral power in the lowest frequency band. This occurs, for example, near 09:00 UT on 4 November, 2021.

In this study we choose to focus on the spectral power contained in the lowest frequency band based on the findings of 85 Oyedokun et al. (2020) that the most dominant frequencies in actual GIC measurements tend to be below 50 mHz.

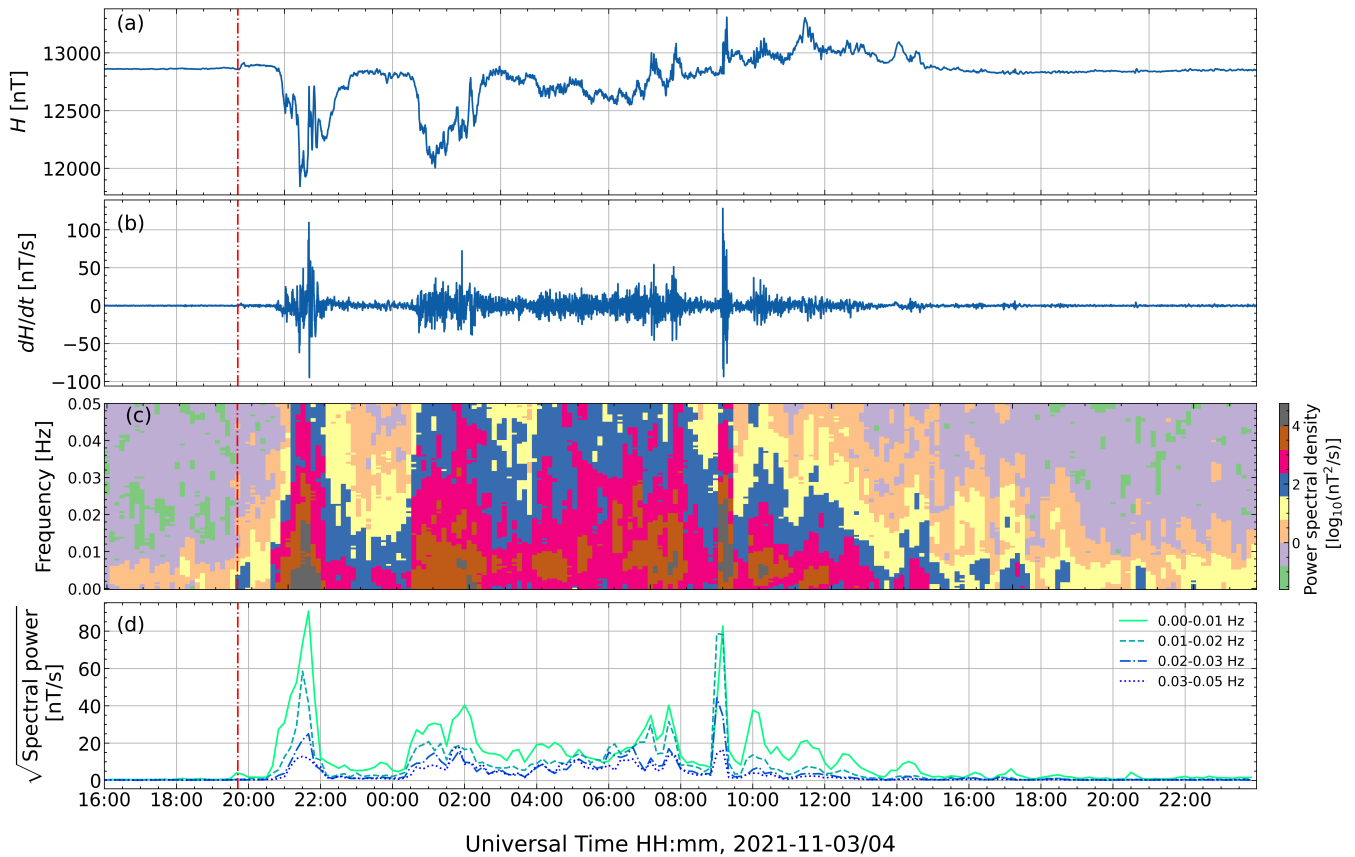


Figure 1. Illustration of procedure for preprocessing Rørvik magnetometer data for measurements made on 3–4 November, 2021. (a) Horizontal magnetic field component H . (b) Time rate of change of H (dH/dt) estimated using backward differencing. (c) Spectrogram of logarithmic dH/dt power spectral density. (d) Square root of integrated dH/dt power spectral density over different frequency bands. The focus of this study is the power in the lowest frequency band (0–0.01 Hz, green solid line). The red dash-dotted vertical line at 19:42 UT on November 3, 2021 indicates the time of impact of an interplanetary CME as given in the Richardson and Cane (2024) CME list.



We view the spectral power of dH/dt as a natural choice for modeling power dissipation due to GICs in power lines and transformer windings, since the spectral power of dH/dt as a function of frequency f , $|dH/dt(f)|^2$ is proportional to the spectral power of the GIC current $|I(f)|^2$ through a transformer, which is in turn directly proportional to the power dissipation $P = I^2 R$. Here R represents the winding resistance of copper windings in a transformer. (The other form of impedance in a transformer, steel core resistance, is the main heat source. It unfortunately cannot be approximated via $P = I^2 R$ as it is highly frequency dependent.)

2.2 Exogenous data inputs

Table 1 lists all inputs used for training the deep learning models we present in the next section. A selection of these inputs are shown in Figure 2 for a four-day period during November 2–5, 2021. These are, in display order from top to bottom, dipole tilt ψ , IMF B_y and B_z , solar wind density n , solar wind speed v , and the “predicted” CME arrival time distribution. Quantities in panels b–e are generated from ACE measurements, while the arrival time distribution in panel f is produced by generating a Gaussian distribution with a standard deviation of ~ 5 –6 h around the time of CME impact as recorded in the Richardson and Cane (2024) CME list. The total area under the curve is 1, and the arrival time distribution is shifted 24 h backward in time so that the model has information about incoming CMEs 24 h ahead of time.

In addition to observed solar wind and IMF properties, panels b–e of Figure 2 also show two sets of artificially generated space weather forecast quantities labeled as “Set A” and “Set B”, respectively denoted by blue lines and orange lines. Both “Set A” and “Set B” quantities consist of ACE measurements shifted 24 h backwards in time, but “Set B” quantities are additionally smoothed using a second-order Savitzky-Golay filter with a 36-h window. Each set of artificial space weather forecast products is used to train a separate model, as described in Section 3. Quantities in panels d–f are among the outputs of the WSA-Enlil and HUXt space weather models. (Note that although the heliospheric magnetic field polarity is routinely forecasted, no currently existing space weather forecasting model is capable of predicting IMF B_y and B_z .)

3 Modeling

In this section we briefly describe the design of the long short-term memory (LSTM) networks we employ, deep learning, and a persistence-based benchmark model for evaluating the performance of the LSTM networks. A summary of the four models used in this study are shown in Table 2.

3.1 Long Short Term Memory

LSTM networks are a variant of recurrent neural networks that handles short-term memory problems by introducing a *cell state* (Hochreiter, 1997). The cell state keeps past information for a long time, solving the vanishing gradient problem and allowing the network to “remember” past information in the sequence (Wei et al., 2018).

As illustrated in Figure 3a, the recurrent structure of the LSTM unit is made up of three gates: forget, input, and output (Siciliano et al., 2021). The forget gate keeps the relevant information in the network by filtering both the current and past

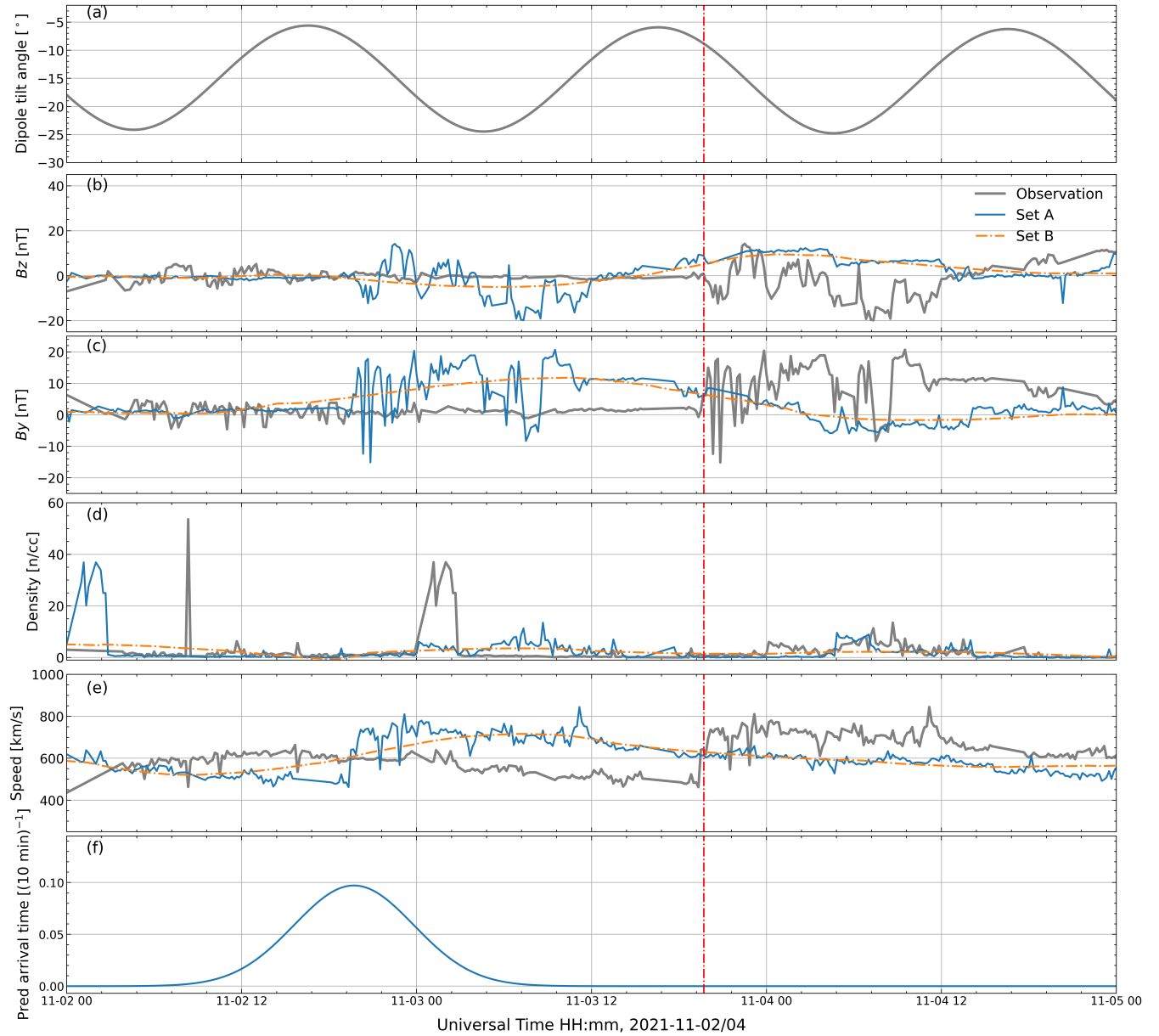


Figure 2. Illustration of exogenous model inputs during November 2–5, 2021. The quantities labeled “Set A” and “Set B” in panels b–e (indicated by blue solid lines and orange dash-dotted lines, respectively) are the artificial solar wind forecasts generated as described in Section 2.2. Both “Set A” and “Set B” are ACE measurements at L1 shifted 24 h backwards in time, but “Set B” is additionally smoothed using a second-order Savitzky-Golay filter with a 36-h window. The red dash-dotted vertical line at 19:42 UT on November 3, 2021 indicates the time of impact of an interplanetary CME according to the Richardson and Cane (2024) CME list. (a) Dipole tilt angle Ψ . (b) Observed and forecasted IMF B_y . (c) Observed and forecasted IMF B_z . (d) Observed and forecasted solar wind density n . (e) Observed and forecasted solar wind speed v . (f) CME arrival time distribution “forecasted” 24 h in advance. Note that the peak of the CME arrival time distribution is deliberately shifted 24 h prior to the actual arrival of the CME indicated by the red dash-dotted line.

**Table 1.** List of input variables used in training of deep learning models

Variable	Symbol
Log dH/dt spectral power over 0–0.01 Hz	$\log_{10}(P_0)$
Cosine UT	$\cos(UT)$
Sine UT	$\sin(UT)$
Dipole tilt angle	ψ
IMF B_y	B_y
IMF B_z	B_z
Solar wind density	n
Solar wind velocity	v
CME arrival time distribution	pred_t
Predicted IMF B_y *	B_y^P
Predicted IMF B_z *	B_z^P
Predicted solar wind density*	n^P
Predicted solar wind velocity*	v^P
Standard normal random number	$\mathcal{N}(0, 1)$

*For Model A, predicted quantities are generated by shifting the original time series

24 h backward. For Model B, predicted quantities are first smoothed using a second-order

Savitzky-Golay filter with a 36-h window and then shifted 24 h backward.

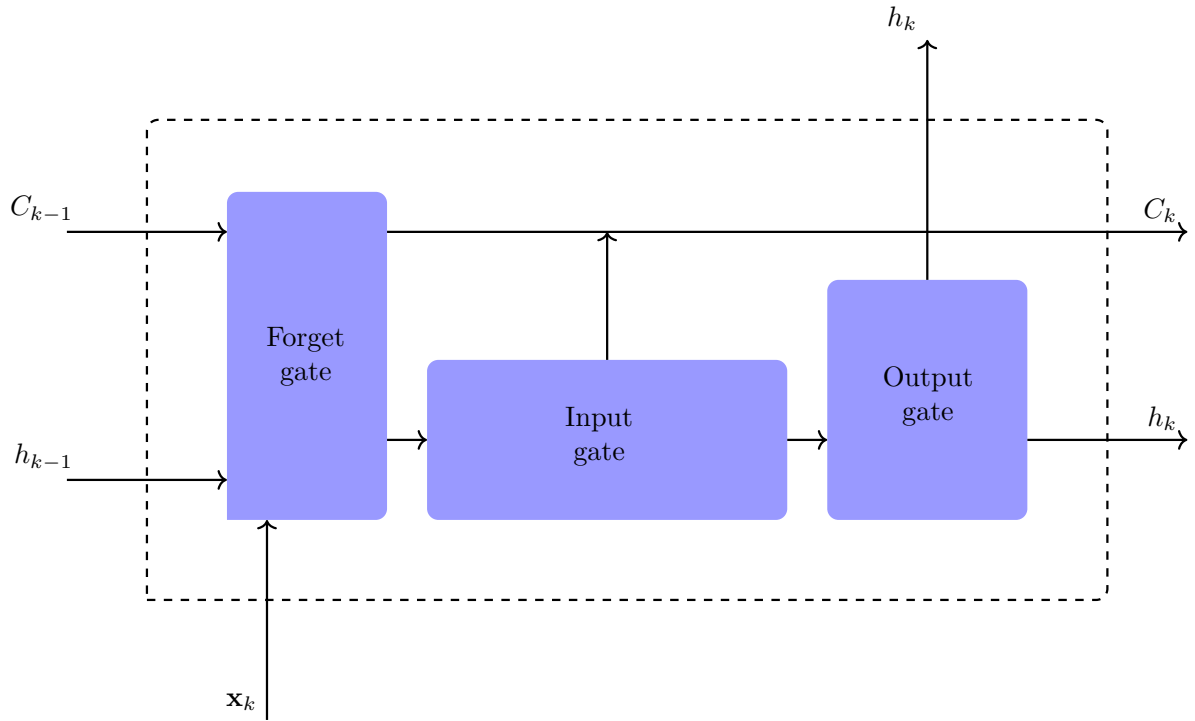
Table 2. Summary of models used in this study.

Model	Description	Section
A	Based on conditions at L1 up to 24 h in advance	3.2
B	Based on smoothed conditions at L1 up to 24 h in advance	3.2
C	Same as Model A but without CME arrival time information	3.2
Benchmark	Weighted combination of persistence and climatology models	3.3

values in the sequence. The input gate enables the network to determine which information to retain from the current values of the sequence. Thus, the input gate is where the current information is added to the network's long-term memory. The previous cell state, C_{k-1} , is then updated to the final cell state, C_k , which is finally uploaded to the output gate. The output gate is the final destination of the information after transitioning from the forget and input gates. The output gate then determines what information is allowed to the next element. In each LSTM unit, there is an element of a sequence being processed, X_k , and the hidden states are also calculated as h_k . In this case, both the hidden state h_k and the cell state C_k are passed on to the preceding element of the sequence, to ensure the past information is used as input in the subsequent element in the sequence being processed. This unique architecture makes the LSTM models better at handling time series problems. In a skillful model, several LSTM cell units per layer are looped together, forming a block of layers such as the one illustrated in Figure 3b.



(a)



(b)

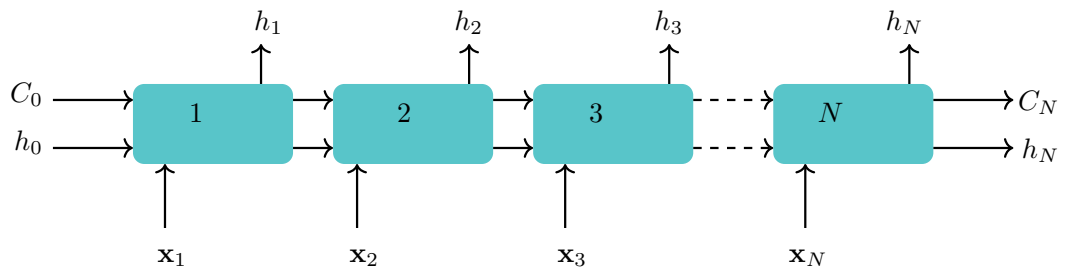


Figure 3. Simplified schematic of an LSTM network at two levels of detail: (a) An LSTM cell; (b) block stacking of LSTM cell units. The symbols X_t , h_t , and C_t are respectively input sequence data, the hidden state, and the cell state at time t .

3.2 Deep Learning

In big data analytics, deep learning models attempt to learn a functional consistency between input features and future values of the label feature y (In this study $y = \log_{10}(P_0)$). The resulting model provides forecasts for the target feature at future time steps. Given a time series vector $\mathbf{x}(t)$ of length n (corresponding to n input features) measured at time t , developing a model to
 130 forecast a target variable $y(t+h)$ at a future time $t+h$ requires historical data $x(t-1)$ at time $t-1$. The model thus requires



uniform input sequences of a fixed length using a sliding lookback window $\Delta t_b = N_b \Delta t$, where N_b is the number of past samples that comprise a sequence and Δt is the cadence at which input features are measured or sampled. In our study these are $\Delta t = 10$ min and $N_b = 204$ so that $\Delta t_b = 34$ h.

Mathematically, we may express the functional relationship between the input features and the target feature that the model has learned as

$$\hat{y}_t^h = g_h(\mathbf{x}_{t-\Delta t_b}, \dots, \mathbf{x}_{t-\Delta t}, y_{t-\Delta t_b}, \dots, y_{t-\Delta t}), \quad (1)$$

where \hat{y}_t^h is the expected forecast for the time t and h indicates the forecast horizon, while $y_{t-\Delta t_b}$ and $x_{t-\Delta t_b}$ are the input target and covariate input variables observed from the time $t - \Delta t_b$ to $t - \Delta t$. The forecast horizon h is the length of time to be forecasted into the future. For the models we present in the following section, the forecast horizon h takes on the values 10 min, 1 h, 3 h, 6 h, and 12 h.

The deep learning model we use consists of an input layer, LSTM units in the hidden layers, and a fully connected or dense layer in the output layer. The hidden layers comprised three layers, each with sequential LSTM units of 327, 56, and 249, and a tanh activation function. The model was optimised using the Optuna framework to obtain the best-performing model hyperparameters Akiba et al. (2019). The tuned parameters had a batch size of 70, a look-back window $\Delta t_b = 34$ h (or $N_b = 204$ as previously mentioned), an optimizer with a learning rate of $5.4\text{e-}05$, and early stopping at a patience value of 8.

This multivariate modeling framework, with a single output for each of the five forecast horizons, results in five distinct models. The choice of this architecture over the multivariate input multi-output allows for the significant contribution of every input variable to the model output in each forecast horizon.

Two LSTM networks were trained with the 14 input variables listed in Table 1 together with a gradient descent approach with the ADAM optimizer Kingma (2014). The first network, hereafter “Model A”, was trained using “Set A” quantities (see Table 1) and thus effectively had access to a perfect forecast of solar wind conditions at L1, as measured by ACE, up to 24 h in advance. The second network, “Model B”, was trained using lower-resolution “Set B” quantities. We therefore expected the performance of Model A to be relatively much higher than that of Model B. Last, to quantify the value of CME arrival time information, we trained a third neural network (“Model C”) using all Set A quantities except for the predicted arrival time “pred_t”. We hypothesized that Model A would perform better than Model C.

Each network was trained using an early stopping mechanism and a maximum of 100 epochs.

We used the Mean Squared Error

$$\text{MSE} = \frac{1}{m} \sum_t (y_t - \hat{y}_t^h)^2 \quad (2)$$

as the loss function in training each LSTM network, where y and \hat{y}^h are as defined in Equation 1 and m is the total number of observations in validation data set (20%, as we describe below). This loss term is used because it increases the sensitivity of the resulting model to outliers relative to the root mean square error (RMSE) Girosi et al. (1995)—or in the case of this study, to the relatively infrequent but extreme values of P_0 that occur during periods of elevated geomagnetic activity. The model that generates the minimum MSE on the validation set is saved, and its performance is evaluated using a test data set with the root



mean square error (the square root of Equation 2). We address the possibility of overfitting using the dropout method. This method consists of ignoring randomly selected neurons and their connections during training.

To calculate the prediction accuracy with the LSTM for forecasting the power of the rate of change of magnetic field, we used the coefficient of determination

$$R^2 = 1 - \frac{\sum_t (y_t - \hat{y}_t^h)^2}{\sum_t (y_t - \bar{y})^2}, \quad (3)$$

where \bar{y} is the mean of the observed values. This metric indicates how much of the variance of the actual measurement is explained by the model (e.g., Liemohn et al., 2018).

Regarding the choice to train the models using the logarithm of the spectral power with an MSE loss function, we also trained sets of models directly using the spectral power (not the logarithm) as well as the RMSE. We found that the specific combination of the logarithm of the spectral power and an MSE loss function yielded the most predictive power. Chu et al. (2025) also found that training using the logarithm of their heavily tailed training data set yielded the most predictive power.

3.3 Benchmark model

We use a simple combination of a persistence model and a climatology model as a performance benchmark for the deep learning model.

A persistence model is defined as a model that renders a naive forecast of a given time series problem (Bailey et al., 2022). The basis of persistence modeling is to use the zero-rule algorithm to predict the future measurement. The zero-rule algorithm is built on the principle that the future forecast value is a continuation of the past measurement as described by Equation $f(t) = f(t-1)$ (Hu et al., 2024). For a regression problem, the mean of the real values is preferred to render a good forecast. Persistence models work well in near-term forecasting and degrade as the forecast horizon increases.

To make a meaningful forecast with a longer forecast horizon, we begin with the simple persistence model:

$$f(t) = \sum \left(\frac{f(t-h-n)}{w} \right), \quad (4)$$

where $f(t)$ is the forecast value and $f(t-h)$ is the previous value. Given that $h = w$ & $n = 0, 1, 2, \dots, w$ and h is the forecast horizon, Equation 4 represents a rolling mean.

The climatology model is obtained by first grouping the entire ~ 15 -year time series of dH/dt spectral power over 0–0.01 Hz (denoted by P_0) by month, and then for each month subgroup calculating the average value of P_0 within a one-hour window for each hour in UT. This procedure yields, for example, the average P_0 during 00–01 UT for the month of June. With 12 months per year and 24 h per day, we end up with 288 average values of P_0 . To obtain a climatology model prediction for a specific date, we use the month of the date to first select a group of 24 coefficients, and then convert the timestamp of the observation to fractional hours and perform a linear interpolation between the two nearest hourly averages.

Finally, we weight the persistence and climatology model outputs equally:

$$P_0^{\text{hybrid}} = \frac{1}{2} P_0^{\text{clim}} + \frac{1}{2} P_0^{\text{pers}}. \quad (5)$$



195 As we show in the following sections, the resulting benchmark model (i.e., combined persistence and climatology model) performs well for near-term forecasting horizons and degrades with increasing forecast horizon.

3.4 Data Preprocessing

We normalized the data using robust scaling and windowing techniques to create suitable sequences for time series forecasting. First, we scaled the target P_0 by calculating the logarithm of the square root of P_0 . The logarithm reduces data skewness and
 200 effectively leads to faster convergence and more predictive power for events of interest in space weather forecasting, such as CMEs, which occur sporadically but are generally accompanied by extreme target values (see, e.g., Chu et al., 2025). Secondly, all input and output data were scaled using the robust scaler method. This method scales data by subtracting the median from each point and scaling it according to the quantile range. Data scaled in this way are robust to outliers (Pedregosa et al., 2011).

Conventionally, the data is split into training, validation, and testing sets for time series forecasting, based on statistical ratios;
 205 for example 70:20:10, respectively. This is desirable for static learning. Here we are however interested in dynamic learning, whereby only the best-performing data in the entire training dataset are used to train the model. In specific, we employ the time series cross-validation method using the TimeSeriesSplit software module (Pedregosa et al., 2011). In this scenario, the training data extends from January 2010 up to and including September 2021, and the test set extends from October 2021 up to and including September 2024. The training set was subjected to time series cross-validation, which splits the training set
 210 further into 80:20 proportions per fold number for training and validation sets, respectively. In this work, we optimized fold number with Optuna, and out of 4 (0, 1, 2, 3), the optimal fold number used for training was 2. The training and validation sets are used to train and find the best model. Eventually, the test set evaluates the model's performance. Finally, all the transformed data were sequenced using the sliding window method (Figure 4). The sliding window was created by taking a time series data sample windowed with a fixed value and rolling it over to the next value per forecast steps. Illustration, taking a subsample of
 215 length n , the forecasted length m , and the rolling window $n - m$, to create different sequences fixed rolling window length is fixed and shifted to the next value as the forecast step m shifts to the next value per sequence. This procedure is summarized in Figure 4.

3.5 Model Hyperparameter Tuning

In deep learning problems, hyperparameter search is a prerequisite step for successful optimal modeling. The growing popular-
 220 ity of deep learning and its complexity require an efficient automatic hyperparameter tuning method. There are well-developed hyperparameter optimization software packages such as Spearmint (Snoek et al., 2012), SMAC (Hutter et al., 2011), Hyperopt (Bergstra et al., 2015), Google Vizier (Golovin et al., 2017), and Autotune (Koch et al., 2018) to address this need. To accelerate the optimization process, distributed computing is required to enable parallel processing of multiple trials. However, the need for a next-generation optimization framework that can dynamically construct the search space, is easy to set up, and
 225 provide efficient sampling and pruning algorithms, Optuna is recommended. Optuna is the next-generation open-source optimization framework that addresses those desired needs (Akiba et al., 2019). The Optuna framework was applied to search for ideal hyperparameters for the LSTM model. Recently, this tuning strategy has been reported to be robust in optimizing model

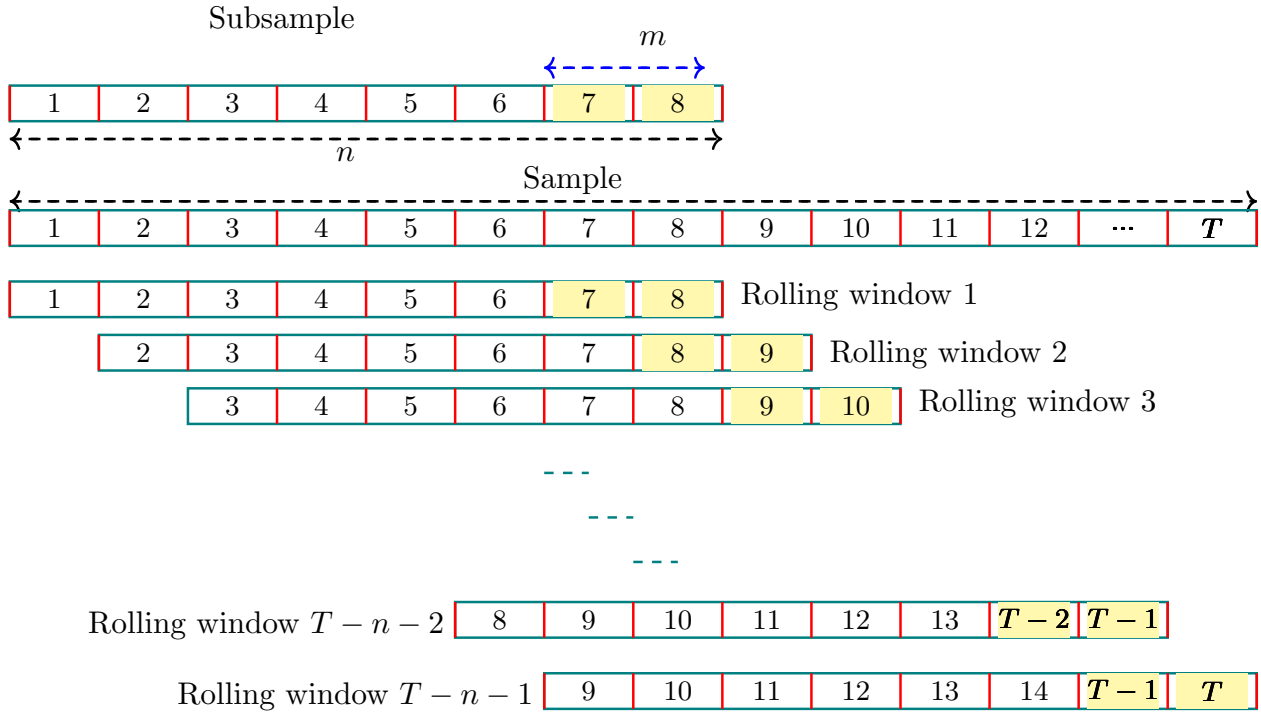


Figure 4. Data windowing using rolling strategies

hyperparameters (Conde et al., 2023). In the present work, we examined sequence window size, number of layers, LSTM hidden units, batch size, patience, learning rate, and fold number. The ADAM optimizer used in the model was preselected based on its successful performance in other modeling works (Kingma, 2014). To achieve the best optimizations with Optuna, various approaches are employed for each tuning. Traditional grid and random-search strategies have proven to be inefficient when the hyperparameter search space is large. The Hyperband approach, which is a bandit-based algorithm (Li et al., 2018), combined with the tree-structured Parzen estimator (TPE) Sampler, is used to approximate the objective function (Bergstra et al., 2011). Given n trials in the optimizations using the Optuna tuner, many training cases may yield the worst outcome, which eventually will be ignored in selecting the optimal model. To save on computational cost, we pruned the training cases that would not yield good models using TFpruningCallback and HyperbandPruner with the minimum resource of three epochs and reduction factor of 3 to the maximum resource of a fixed training epoch number. The final useful trained models listed are again stored and ranked to produce the best optimal model using their loss ranking basis.



3.6 Feature importance

Neural networks are often regarded as black-box algorithms, though some useful external inference methods are used to explain deep learning models. Computing the feature importance of deep learning models, such as LSTM, provides a vital aspect of model understanding and interpretation. The degree of feature or input variable contribution to the output of a model is measured by feature importance. Feature importance is thus calculated based on Gini importance, such that more important features have higher values and less important features have lower values. Understanding the feature importance of a model can help improve model performance by training the model with the most relevant features. It also provides useful insights into the underlying dynamics, relationships, and patterns available in the examined time-series data.

In this study, we calculate the feature importance via the gradient tape technique as implemented in TensorFlow (Abadi et al., 2015). Gradient tape is a gradient-based, post-hoc method that works on an already trained model, and that keeps track of automatic differentiation operations of tensor variables. This watching capability makes it a useful technical tool in evaluating the contributions of trainable input variables in a model. We use the gradient tape technique instead of other widely used feature importance techniques, such as Shapley Additive exPlanations (Lundberg and Lee, 2017; Long et al., 2022) and input permutation (Breiman, 2001; Fisher et al., 2019), that are not necessarily optimized for evaluating deep learning-based models. We present feature importance analysis of Models A–C in section 5.

4 Results

In this section, we present case study results for three space weather events. The first occurred at 19:42 UT on November 3, 2021, and involved a CME impact at Earth. The second involved consecutive CME impacts at 17:05 UT on May 10 and 09:17 UT on May 12 in 2024. We also present results from a high-speed stream/stream interaction region (HSS/SIR) storm that occurred during May 5–7, 2023.

Figure 5 shows the minor geomagnetic storm of November 3–4. It includes measured values of P_0 at Rørvik (gray line) as well as the forecasted power at forecasting horizons of 10 min, 1 h, 3 h, 6 h, and 12 h (panels a–e, respectively). Results from Models A–C (blue, orange, and green lines, respectively) and the benchmark model (dashed black line) for each forecasting horizon are also shown. In 10-minute and 1-hour forecasts, Models A–C and the benchmark model yielded similar forecast output as demonstrated in Figures 5a–b.

For Models A–C, the trend in the forecasted time series is mostly consistent for forecasting horizons extending from $h = 10$ min up to $h = 3$ h. The forecast output amplitudes decrease with increasing forecast horizon. The deep learning models perform better than the benchmark model, which exhibits worsening performance with increasing forecast horizon. In panels d and e, 6 h and 12 h forecast, the model demonstrated very low forecast power, yielding output with diminished amplitudes. During this period the solar wind speed reached a maximum value of 850 km/s, and IMF B_z reached a minimum value of -20 nT, as shown in Figure 2.

Figure 6 illustrates the forecast output of the May 9–12, 2024 geomagnetic storm. Models A–C yielded a generally good correlation with observed P_0 values (thick gray line) up to 1 hour in advance, as shown in panels a–b. However, the enhanced



P_0 values that accompanied the last phase of this storm were only captured by the 10-minute forecast of each model, while the enhancements associated with the last storm phase are mostly absent in the models with a 1-h or greater forecasting horizon. We return to this point in the discussion.

275 Panels c–d show results for forecasting horizons of 3–6 h; the storm initial and main phases were well captured in the forecast, while the last phase was completely missed. The model forecast results for $h = 12$ h, shown in panel e, shows some resemblance to the observed time series, but does not evince any of the rapid, larger amplitude variations that appear in the observed time series. In particular the forecasted values of P_0 are everywhere less than 30 nT/s for all models, whereas the observed time series frequently exceeds 60 nT/s.

280 Figure 7 illustrates observations and model results for an HSS/SIR-driven storm that occurred during May 5–7, 2023. During this period the maximum solar wind speed was 580 km/s, and IMF B_z reached a minimum value of -18 nT (not shown). For forecasting horizons beyond 10 min, none of the deep learning models manage to predict the enhanced levels of P_0 that were observed at Rørvik.

Figure 8 summarizes performance metrics of each model for each of the five forecasting horizons. Results for Models A and
285 B are respectively indicated by solid blue and orange dash-dotted lines, while metrics for the benchmark model are indicated by black dotted lines. (Results for Model C are essentially identical to those of Model A and are omitted.) The coefficients of determination R^2 for Models A and B decrease uniformly from 88% for forecast horizon $h = 10$ min to $\sim 50\%$ for $h = 3$ h. The performance of Model A recorded performance metrics with a slight increasing trend over $h = 3$ –12 h with values of 50%, 52%, and 56%. On the other hand, Model B for $h = 3$ –12 h continues to decrease gradually with values 52%, 42%, and
290 42% respectively. The R^2 value for the benchmark model monotonically decreases with increasing forecast horizon h and is everywhere lower than R^2 values for Models A and B.

Examining the RMSE of every forecast horizon, there are similar trends as a function of forecasting horizon h for all models, as shown in Figure 8b. Models A and B had the same in lower forecast horizons up to 3 hours in advance, with RMSE values of 1.4–2 nT/s, respectively. For longer forecasting horizons ($h = 3$ –12 h), for Model A there is a slight decrease from 2.06
295 to 1.97 nT/s, while for Model B, RMSE increases from 2.0 to 2.2 nT/s. The RMSE values are calculated from the target and forecast variables $y(= \log_{10} P_0)$ and \hat{y} , and are then quoted here after raising them to the power of 10.

In the benchmark model, there is a steady decrease in model performance with increasing forecast horizon h , with $h = 12$ h showing the worst RMSE value of 2.75 nT/s, and $h = 10$ min slightly higher than those of intelligent models, with a value of 1.56 nT/s.

300 In the current study, feature importance was calculated using the gradient-based method described in Section 3.6. Figure 9 shows contributions of every input feature to the model output for all forecast horizons, demonstrated in panels a–e. For $h = 10$ min (Figure 9a) the input feature P_0 contributes the most, with a gradient slightly higher than 0.005. The other significant input features (gradients above 0.005) are IMF B_z , v , $\sin(\text{UT})$, and $\cos(\text{UT})$. All other features have little influence on the performance of the models, as indicated by their having gradients similar to that of the normally distributed random input
305 denoted by $\mathcal{N}(0, 1)$ in Table 1.



For $h = 1$ h (Figure 9b) the importance of IMF B_z is greatest, followed by P_0 and v with gradient values above 0.0015. Forecast input features, such as predicted IMF B_z and v , also make a non-negligible contribution to the 1 h forecast.

For $h = 3$ h (Figure 9c) Models A and B respond differently to the various input features. The importance of IMF B_z , v , $\sin UT$, and pred_arrival time dominate for Model A. (Predicted arrival time information is not provided to Model C and therefore has zero importance.) On the other hand, predicted arrival time and predicted IMF B_z are the two most important features for Model B, with multiple other features making lesser but nevertheless non-negligible contributions.

For $h = 6$ h (Figure 9d), P_0 has little to no importance. Predicted IMF B_z and solar wind speed v are important for all three models, and predicted arrival time is also relatively important for Models A and B.

The trends for $h = 12$ h (Figure 9e) are generally the same as those for $h = 6$ h. The most important features for all models are externally driven, with the gradients having higher values for Models A and C than those for Model B.

In addition to the foregoing case study and feature importance results, we have also carried out a binary classification analysis of events using the class-wise performance measures

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (6)$$

and

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}. \quad (7)$$

In these expressions “true positives” (TP) are those for which the predicted values \hat{y}_{t+h} of the target variable $y = \log_{10} P_0$ for forecast horizon h both reach or exceed an arbitrarily selected percentile threshold of 0.985 in their respective data sets. “True negatives” (TN) are those for which both forecasted and observed values are less than the threshold. For wrongly forecasted false negatives (FN) the forecasted value does not reach the threshold percentile while the observed value exceeds it, and vice versa for false positives (FP) (Welling et al., 2018).

We also use True Skill Statistics

$$TSS = \text{Sensitivity} + \text{Specificity} - 1, \quad (8)$$

where

$$\text{Specificity} = \frac{TN}{TN + FP}, \quad (9)$$

and the Matthew Correlation Coefficient (MCC) as measures of model performance, accuracy, and unbiased balance between predicted classes. The hit rate (elsewhere known as the “probability of detection” or POD) was calculated using sensitivity, and the false alarm rate (elsewhere the “probability of false detection” or POFD) by specificity. The metrics represented by Equations 6–9 are used in a number of previous investigations of space weather forecasting ?e.g.,>[]Hu2024,Baily2021SW.

Regarding our choice of 0.985 as a percentile threshold for calculating the confusion matrix, this threshold corresponds to ~ 10 nT/s in the target variable P_0 , which serves as a practical threshold for assessing the potential impact of geomagnetic activity on nearby deployed infrastructure. We also found that a threshold of 0.985 gives the best comparison with Table 3 in



Hu et al. (2024) for a 1-h forecasting horizon. We also tested increasing the percentile threshold to 0.998 (corresponding to $P_0 \approx 20$ nT/s) and found that this only slightly decreased the TSS and MCC. We therefore deem the results with a threshold of 0.985 to be representative.

Metrics	Model	Forecast Horizons				
		10 min	1 h	3 h	6 h	12 h
TP	Model A*	1612	867	329	524	753
	Model B	1594	854	457	340	295
	Benchmark	602	266	136	106	64
TN	Model A	150951	150739	151229	150809	150574
	Model B	150979	150749	151090	151285	151375
	Benchmark	151757	151663	151672	151690	1512733
FP	Model A	629	841	351	771	1006
	Model B	601	831	490	295	205
	Benchmark	27	121	112	94	51
FN	Model A	800	1545	2083	1888	1659
	Model B	818	1558	1955	2072	2117
	Benchmark	1811	2147	2277	2308	2349
MCC	Model A	0.69	0.42	0.25	0.29	0.36
	Model B	0.69	0.42	0.30	0.27	0.26
	Benchmark	0.49	0.27	0.17	0.15	0.12
TSS	Model A	0.66	0.35	0.13	0.21	0.31
	Model B	0.66	0.35	0.19	0.14	0.12
	Benchmark	0.25	0.11	0.06	0.04	0.03

*Model C values are very similar to those for Model A and are therefore not shown.

Table 3. Comparison of MCC and TSS score between LSTM and benchmark models with forecast horizons between 10 min and 12 h.

340 Table 3 presents the confusion matrix and several other statistical classification metrics of model performance for Models A and B as well as the benchmark model. (Metrics for Model C are almost identical to those of Model A and are not shown.) All of the models perform well for $h = 10$ min, with Models A–C having the highest values of TP, MCC, and TSS.

Compared with Table 3 in Hu et al. (2024), only TSS for both models in a 10-min forecasting horizon matches the value TSS=0.66 that they obtained, while for a forecasting horizon $h = 1$ h we find our models have slightly higher TSS values.
 345 Here it is important to observe that their study examined the geoelectric field at mid-latitudes, whereas this study is focused on a GIC proxy at high latitudes, so this indirect comparison is only suggestive.

We also examined the performance of a simple recurrence model (not shown) in which the current value was predicted to be the same as the value 24 h earlier; this model performed exceptionally poorly, with a coefficient of determination $R^2 = -0.44$.



5 Discussion

Results presented in the previous section illustrate that the deep learning models we employ exhibit modest learning levels and predictive power. The forecasting power of the three deep learning models initially decreases for forecasting horizons up to about 3 h, and then show a trend in Figure 8 toward constant or slightly increasing forecasting power for horizons above 3 h, with a similar degree of performance ($R^2 = 50\%$). A similar finding was reported in the study of Zewdie et al. (2021).

Despite the rather modest levels of predictive power exhibited by the deep learning models, they all outperform the benchmark model for all horizons. Previous studies have shown that deep learning models generally perform better than conventional machine learning models (e.g., neural networks and nonlinear models), which are only capable of predicting one step at a time (Keese et al., 2020; Zewdie et al., 2021; Long et al., 2022; Hu et al., 2023).

The feature importance shown in Figure 9 was calculated for each model and h value to explain the deep learning model. This was done using the post hoc gradient based method, which we find to be most appropriate for evaluation of deep learning models of the sort we have presented in this study. Results in Figure 9 shows that the importance of the various input features varies significantly with forecasting horizon h . For instance, for $h = 10$ min and $h = 1$ h, only IMF B_z and solar wind velocity v significantly contribute beyond that of the target variable P_0 itself. We also see that as forecasting horizon h increases, the models rely increasingly on accurate information about predicted IMF B_z , and predicted solar wind speed v .

Figure 9 indicates that CME arrival time is an important input feature, but Figure 8 clearly demonstrates that Model C performs just as well as Model A without CME arrival time information. We therefore conclude that high-resolution solar wind and IMF forecasts are necessary to achieve the model performances we have reported, whereas CME arrival time information is unimportant when high-resolution solar wind and IMF forecasts are available.

In contrast, it is somewhat unclear what role CME arrival time information plays when only low-resolution (but nevertheless accurate) solar wind and IMF forecasts are provided to the model: According to Figures 9c–e, CME arrival time is the single most important piece of information for Model B’s predictions. To test this we trained another version of Model B (low-resolution space weather forecast inputs) in which information about IMF B_z and IMF B_y was excluded but CME arrival time was included. In this model (not shown), CME arrival time had no importance for any of the five forecast horizons; this model also performed generally much worse than Model B for forecasting horizons $h = 6$ h and $h = 12$ h. We therefore conclude that when only low-resolution solar wind and IMF forecasts are available, CME arrival time could contribute to model performance. This conclusion incidentally highlights that a study of model performance as a function of temporal resolution and accuracy of solar wind forecast parameters would likely provide clarity around this point.

One unresolved aspect of the model performance is that the model inputs do not always seem to provide the information necessary for the model to predict enhanced geomagnetic activity. This is illustrated in Figure 7, where we observe that the deep learning models performed poorly in longer forecast horizons beyond 1 h. It is interesting to note that this particular event is associated with HSS/SIR. By contrast, the event in Figure 5 is a purely CME-driven storm. This leads us to speculate that perhaps the nature of the solar wind and IMF signatures of some types of disturbances, such as HSS/SIR events, are not clearly identifiable from the inputs provided to the model for longer forecasting horizons. This same seeming insensitivity to elevated



solar and IMF conditions is apparent in model outputs for the latter half of the time window shown in Figure 6 for the May 2024 storm: a large enhancement in geomagnetic activity was observed at Rørvik, several hours before the arrival of a second CME which was not geoeffective at 09:17 UT on May 12, 2024. The second enhancement was not forecasted by the models for forecasting horizons beyond 1 h. We reserve determining whether model predictive power in fact varies based on the type of disturbance carried by the solar wind as the topic of a potential future study.

Our finding about the limited utility of existing space weather forecast products likely applies to other domains concerned with space weather, such as predictions of total electron content (TEC), GPS scintillation, and satellite drag. A possible counterexample to this suggestion is the work of Adolfs et al. (2024), who find that the performance of their TEC forecasting model did not improve significantly with the addition of external inputs.

The findings on the class-wise classification of the models' outputs showed MCC and TSS metrics were consistent with previous studies, for instance the work of Hu et al. (2024). The MCC and TSS summarize the significance of TP, TN, FP, and FN classes in model performance. Thus, the models' imperfection in space weather forecasting for longer forecast horizons was still reliable in forecasting events. In addition to RMSE and R^2 , we observe a similar trend in model performance with the MCC and TSS metrics, indicating reliable intelligence in the model benchmark by the persistence model.

6 Conclusions

A central goal of this study was to determine whether we could produce an accurate forecast of local geomagnetic activity at the Rørvik magnetometer station up to 12 h in advance using artificial forecasts of space weather conditions and information about CME arrival. We find that high-resolution (time resolution of 10 min) space weather forecasts of IMF B_z and solar wind speed v 24 h in advance would enable a reasonably accurate ($R^2 = 55\%$) forecast of local geomagnetic activity 12 h in advance. In this case the inclusion of CME arrival time provides no useful information to the model. When only low-resolution (time resolution of 10-min but smoothed with a 36-h window) space weather forecasts are available, CME arrival time information may improve model performance. In the absence of any information about IMF B_z , CME arrival time is apparently unimportant.

Increasing the predictive power of forecasts with horizons beyond 1 h therefore seems to require either heliosphere / solar wind models that are able to accurately predict the IMF — a capability that does not exist at present — or in situ solar wind monitors that, for example, have a different viewing angle, such as the European Space Agency's upcoming Vigil mission which will be located at L5.

This study is, to our knowledge, the first to demonstrate via artificial intelligence-based models what is otherwise well established in the space weather community: That information about the IMF B_z component is crucial to predicting the geoeffectiveness of an inbound CME. This study therefore suggests how one might incorporate other potentially relevant measures of space weather, such as solar energetic proton intensity or X-ray flux. We speculate that these quantities might allow a deep learning-based model to further modulate the intensity of predicted local geomagnetic activity.



415 *Data availability.* ACE satellite measurements of IMF and solar wind speed are available at <https://sohoftp.nascom.nasa.gov/sdb/goes/ace/daily/>. Rørvik magnetometer data are made available by the Tromsø Geophysical Observatory at <https://flux.phys.uit.no/ascii/>. The latest version of the Richardson and Cane near-Earth ICME list (Richardson and Cane, 2024) is available at <https://izw1.caltech.edu/ACE/ASC/DATA/level3/icmetable2.htm>.

420 *Author contributions.* SO: model design and creation (lead), formal analysis, writing (equal), and visualization (lead). SMH: project administration, funding acquisition, supervision (lead), training dataset production, study conceptualization (lead), writing (equal). AK: model design, conceptualization, validation, and writing. MGJ: supervision, conceptualization, validation, writing, and curation of the Rørvik magnetometer dataset. MO: conceptualization including expertise on solar wind forecasting methods, writing – review and editing, and validation. KST: supervision, contribution of expertise on induction in power systems, writing – review and editing. RL: validation, contribution of expertise on GIC events, writing – review and editing.

425 *Competing interests.* The authors declare no competing interests.

Acknowledgements. The work of SO and SMH was supported by the University of Bergen via the UiB Idé program. SMH was additionally supported by NFR grant no. 355484.



References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Good-
430 fellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S.,
Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F.,
Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X.: TensorFlow: Large-Scale Machine Learning on Heterogeneous
Systems, <https://www.tensorflow.org/>, software available from tensorflow.org, 2015.
- Abda, Z. M. K., Aziz, N. F. A., Kadir, M. Z. A. A., and Rhazali, Z. A.: A Review of Geomagnetically Induced Current Effects on Electrical
435 Power System: Principles and Theory, *IEEE Access*, 8, 200 237–200 258, <https://doi.org/10.1109/ACCESS.2020.3034347>, 2020.
- Adolfs, M., Hoque, M. M., and Shprits, Y. Y.: Forecasting 24-Hr Total Electron Content With Long Short-Term Memory Neural Network,
Journal of Geophysical Research: Machine Learning and Computation, 1, <https://doi.org/10.1029/2024JH000123>, 2024.
- Advanced Composition Explorer (ACE): High-resolution ACE solar wind measurements [dataset], retrieved from <https://sohftp.nascom.nasa.gov/sdb/goes/ace/daily/>, 2025.
- 440 Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M.: Optuna: A next-generation hyperparameter optimization framework, in: *Proceed-*
ings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining, pp. 2623–2631, 2019.
- Arge, C. N. and Pizzo, V. J.: Improvement in the prediction of solar wind conditions using near-real time solar magnetic field updates, *Journal*
of Geophysical Research: Space Physics, 105, 10 465–10 479, <https://doi.org/https://doi.org/10.1029/1999JA000262>, 2000.
- Arge, C. N., Luhmann, J., Odstrcil, D., Schrijver, C., and Li, Y.: Stream structure and coronal sources of the so-
445 lar wind during the May 12th, 1997 CME, *Journal of Atmospheric and Solar-Terrestrial Physics*, 66, 1295–1309,
<https://doi.org/https://doi.org/10.1016/j.jastp.2004.03.018>, towards an Integrated Model of the Space Weather System, 2004.
- Bailey, R. L., Leonhardt, R., Möstl, C., Beggan, C., Reiss, M. A., Bhaskar, A., and Weiss, A. J.: Forecasting GICs and
Geoelectric Fields From Solar Wind Data Using LSTMs: Application in Austria, *Space Weather*, 20, e2021SW002907,
<https://doi.org/https://doi.org/10.1029/2021SW002907>, e2021SW002907 2021SW002907, 2022.
- 450 Barnard, L. and Owens, M.: HUXt—An open source, computationally efficient reduced-physics solar wind model, written in Python, *Front-*
iers in Physics, 10, <https://doi.org/10.3389/fphy.2022.1005621>, 2022.
- Bergstra, J., Bardenet, R., Bengio, Y., and Kégl, B.: Algorithms for hyper-parameter optimization, *Advances in neural information processing*
systems, 24, 2011.
- Bergstra, J., Komer, B., Eliasmith, C., Yamins, D., and Cox, D. D.: Hyperopt: a python library for model selection and hyperparameter
455 optimization, *Computational Science & Discovery*, 8, 014 008, 2015.
- Breiman, L.: Random forests, *Machine learning*, 45, 5–32, 2001.
- Camporeale, E.: The Challenge of Machine Learning in Space Weather: Nowcasting and Forecasting, *Space Weather*, 17, 1166–1207,
<https://doi.org/10.1029/2018SW002061>, 2019.
- Camporeale, E., Wing, S., Johnson, J., Jackman, C. M., and McGranaghan, R.: Space Weather in the Machine Learning Era: A Multidisci-
460 plinary Approach, *Space Weather*, 16, 2–4, <https://doi.org/10.1002/2017SW001775>, 2018.
- Cane, H. V. and Richardson, I. G.: Interplanetary coronal mass ejections in the near-Earth solar wind during 1996–2002, *Journal of Geo-*
physical Research, 108, 1156, <https://doi.org/10.1029/2002JA009817>, 2003.
- Chu, X., Jia, L., McPherron, R. L., Li, X., and Bortnik, J.: Imbalanced Regression Artificial Neural Network Model for Auroral Electrojet
Indices (IRANNA): Can We Predict Strong Events?, *Space Weather*, 23, <https://doi.org/10.1029/2024SW004236>, 2025.



- 465 CIGRE: TB 780 - Understanding of geomagnetic storm environment for high voltage power grids, Tech. rep., <https://www.e-cigre.org/publications/detail/780-understanding-of-geomagnetic-storm-environment-for-high-voltage-power-grids.html>, 2019.
- Conde, D., Castillo, F. L., Escobar, C., García, C., García, J. E., Sanz, V., Zaldívar, B., Curto, J. J., Marsal, S., and Torta, J. M.: Forecasting Geomagnetic Storm Disturbances and Their Uncertainties Using Deep Learning, *Space Weather*, 21, e2023SW003474, <https://doi.org/https://doi.org/10.1029/2023SW003474>, 2023.
- 470 EPRI: Research Findings for Geomagnetic Disturbance Research Work Plan: Summary Report, Tech. rep., <https://www.epri.com/research/products/000000003002019720>, 2020.
- Fisher, A., Rudin, C., and Dominici, F.: All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously, *Journal of Machine Learning Research*, 20, 1–81, 2019.
- Fry, E. K.: The risks and impacts of space weather: Policy recommendations and initiatives, *Space Policy*, 28, 180–184, <https://doi.org/10.1016/j.spacepol.2012.06.005>, 2012.
- 475 Gannon, J. L., Swidinsky, A., and Xu, Z., eds.: Geomagnetically Induced Currents from the Sun to the Power Grid, Wiley, ISBN 9781119434344, <https://doi.org/10.1002/9781119434412>, 2019.
- Girosi, F., Jones, M., and Poggio, T.: Regularization theory and neural networks architectures, *Neural computation*, 7, 219–269, 1995.
- Golovin, D., Solnik, B., Moitra, S., Kochanski, G., Karro, J., and Sculley, D.: Google vizier: A service for black-box optimization, in: Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining, pp. 1487–1495, 2017.
- 480 Gonzalez, W. D., Echer, E., Tsurutani, B. T., Gonzalez, A. L. C., and Lago, A. D.: Interplanetary Origin of Intense, Superintense and Extreme Geomagnetic Storms, *Space Science Reviews*, 158, 69–89, <https://doi.org/10.1007/s11214-010-9715-2>, 2011.
- Hapgood, M.: Towards a scientific understanding of the risk from extreme space weather, *Advances in Space Research*, 47, 2059–2072, <https://doi.org/10.1016/j.asr.2010.02.007>, 2011.
- 485 Hatch, S. M. and LaBelle, J.: Application of a new method for calculation of low-frequency wave vectors, in: PLANETARY RADIO EMISSIONS VIII, edited by Fischer, G., Mann, G., Panchenko, M., and Zarka, P., pp. 247–260, Austrian Academy of Sciences Press, ISBN 978-3-7001-8263-4, <https://doi.org/10.1553/PRE8s247>, 2018.
- Hochreiter, S.: Long Short-term Memory, Neural Computation MIT-Press, 1997.
- Hu, A., Camporeale, E., and Swiger, B.: Multi-Hour-Ahead Dst Index Prediction Using Multi-Fidelity Boosted Neural Networks, *Space*
- 490 *Weather*, 21, <https://doi.org/10.1029/2022SW003286>, 2023.
- Hu, A., Camporeale, E., Lucas, G., and Berger, T.: LiveWire: Horizontal Geoelectric Field Prediction With 1-hr Lead-Time Using Multi-Fidelity Boosted Neural Networks, *Journal of Geophysical Research: Machine Learning and Computation*, 1, <https://doi.org/10.1029/2024JH000151>, 2024.
- Hutter, F., Hoos, H. H., and Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration, in: Learning and intelligent optimization: 5th international conference, LION 5, rome, Italy, January 17–21, 2011. selected papers 5, pp. 507–523, Springer, 2011.
- 495 Iong, D., Chen, Y., Toth, G., Zou, S., Pulkkinen, T., Ren, J., Camporeale, E., and Gombosi, T.: New Findings From Explainable SYM-H Forecasting Using Gradient Boosting Machines, *Space Weather*, 20, e2021SW002928, <https://doi.org/https://doi.org/10.1029/2021SW002928>, 2022.
- 500 Johnsen, M. G.: Tromsø Geophysical Observatory geomagnetic data, magnetometer measurements retrieved from <https://flux.phys.uit.no/ascii/>, 2025.



- Keesee, A. M., Pinto, V., Coughlan, M., Lennox, C., Mahmud, M. S., and Connor, H. K.: Comparison of Deep Learning Techniques to Model Connections Between Solar Wind and Ground Magnetic Perturbations, *Frontiers in Astronomy and Space Sciences*, 7, <https://doi.org/10.3389/fspas.2020.550874>, 2020.
- 505 Kingma, D. P.: Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980, 2014.
- Koch, P., Golovidov, O., Gardner, S., Wujek, B., Griffin, J., and Xu, Y.: Autotune: A derivative-free optimization framework for hyperparameter tuning, in: *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 443–452, 2018.
- Lakhina, G. S. and Tsurutani, B. T.: Geomagnetic storms: historical perspective to modern view, *Geoscience Letters*, 3, 5, <https://doi.org/10.1186/s40562-016-0037-4>, 2016.
- 510 Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and Talwalkar, A.: Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization, *Journal of Machine Learning Research*, 18, 1–52, <http://jmlr.org/papers/v18/16-558.html>, 2018.
- Liemohn, M. W., McCollough, J. P., Jordanova, V. K., Ngwira, C. M., Morley, S. K., Cid, C., Tobiska, W. K., Wintoft, P., Ganushkina, N. Y., Welling, D. T., Bingham, S., Balikhin, M. A., Opgenoorth, H. J., Engel, M. A., Weigel, R. S., Singer, H. J., Buresova, D., Bruinsma, S., Zhelavskaya, I. S., Shprits, Y. Y., and Vasile, R.: Model Evaluation Guidelines for Geomagnetic Index Predictions, *Space Weather*, 16, 2079–2102, <https://doi.org/https://doi.org/10.1029/2018SW002067>, 2018.
- 515 Lundberg, S. M. and Lee, S.-I.: A unified approach to interpreting model predictions, *Advances in neural information processing systems*, 30, 2017.
- Ngwira, C. M. and Pulkkinen, A. A.: An Introduction to Geomagnetically Induced Currents, pp. 1–13, <https://doi.org/10.1002/9781119434412.ch1>, 2019.
- 520 Odstřil, D., Riley, P., and Zhao, X. P.: Numerical simulation of the 12 May 1997 interplanetary CME event, *Journal of Geophysical Research: Space Physics*, 109, <https://doi.org/https://doi.org/10.1029/2003JA010135>, 2004.
- Oyedokun, D., Heyns, M., Cilliers, P., and Gaunt, C.: Frequency Components of Geomagnetically Induced Currents for Power System Modelling, in: *2020 International SAUPEC/RobMech/PRASA Conference*, pp. 1–6, IEEE, ISBN 978-1-7281-4162-6, <https://doi.org/10.1109/SAUPEC/RobMech/PRASA48453.2020.9041021>, 2020.
- 525 Patterson, C. J., Wild, J. A., and Boteler, D. H.: Modeling “Wrong Side” Failures Caused by Geomagnetically Induced Currents in Electrified Railway Signaling Systems in the UK, *Space Weather*, 21, <https://doi.org/10.1029/2023SW003625>, 2023.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E.: Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research*, 12, 2825–2830, 2011.
- 530 Press, N. A.: *Planning the Future Space Weather Operations and Research Infrastructure*, National Academies Press, ISBN 978-0-309-45433-9, <https://doi.org/10.17226/26128>, 2021.
- Prieto, G. A.: The Multitaper Spectrum Analysis Package in Python, *Seismological Research Letters*, 93, 1922–1929, <https://doi.org/10.1785/0220210332>, 2022.
- 535 Pulkkinen, A.: Geomagnetically Induced Currents Modeling and Forecasting, *Space Weather*, 13, 734–736, <https://doi.org/10.1002/2015SW001316>, 2015.
- Pulkkinen, A., Bernabeu, E., Thomson, A., Viljanen, A., Pirjola, R., Boteler, D., Eichner, J., Cilliers, P. J., Welling, D., Savani, N. P., Weigel, R. S., Love, J. J., Balch, C., Ngwira, C. M., Crowley, G., Schultz, A., Kataoka, R., Anderson, B., Fugate, D., Simpson, J. J.,



- and MacAlester, M.: Geomagnetically induced currents: Science, engineering, and applications readiness, *Space Weather*, 15, 828–856, <https://doi.org/10.1002/2016SW001501>, 2017.
- Richardson, I. and Cane, H.: Near-Earth Interplanetary Coronal Mass Ejections Since January 1996, <https://doi.org/10.7910/DVN/C2MHTH>, 2024.
- Siciliano, F., Consolini, G., Tozzi, R., Gentili, M., Giannattasio, F., and De Michelis, P.: Forecasting SYM-H Index: A Comparison Between Long Short-Term Memory and Convolutional Neural Networks, *Space Weather*, 19, e2020SW002589, <https://doi.org/https://doi.org/10.1029/2020SW002589>, e2020SW002589 10.1029/2020SW002589, 2021.
- Siddique, T. and Mahmud, M. S.: Ensemble deep learning models for prediction and uncertainty quantification of ground magnetic perturbation, *Frontiers in Astronomy and Space Sciences*, 9, <https://doi.org/10.3389/fspas.2022.1031407>, 2022.
- Sierra-Porta, D., Petro-Ramos, J., Ruiz-Morales, D., Herrera-Acevedo, D., García-Teheran, A., and Alvarado, M. T.: Machine learning models for predicting geomagnetic storms across five solar cycles using Dst index and heliospheric variables, *Advances in Space Research*, 74, 3483–3495, <https://doi.org/10.1016/j.asr.2024.08.031>, 2024.
- Snoek, J., Larochelle, H., and Adams, R. P.: Practical bayesian optimization of machine learning algorithms, *Advances in neural information processing systems*, 25, 2012.
- Thomson, D.: Spectrum estimation and harmonic analysis, *Proceedings of the IEEE*, 70, 1055–1096, <https://doi.org/10.1109/PROC.1982.12433>, 1982.
- Tóth, G., Meng, X., Gombosi, T. I., and Rastätter, L.: Predicting the time derivative of local magnetic perturbations, *Journal of Geophysical Research: Space Physics*, 119, 310–321, <https://doi.org/10.1002/2013JA019456>, 2014.
- Viljanen, A., Nevanlinna, H., Pajunpää, K., and Pulkkinen, A.: Time derivative of the horizontal geomagnetic field as an activity indicator, *Annales Geophysicae*, 19, 1107–1118, <https://doi.org/10.5194/angeo-19-1107-2001>, 2001.
- Wang, T., Luo, B., Wang, J., Ao, X., Shi, L., Zhong, Q., and Liu, S.: Forecasting of the Geomagnetic Activity for the Next 3 Days Utilizing Neural Networks Based on Parameters Related to Large-Scale Structures of the Solar Corona, *Space Weather*, 23, <https://doi.org/10.1029/2024SW004090>, 2025.
- Wei, L., Zhong, Q., Lin, R., Wang, J., Liu, S., and Cao, Y.: Quantitative Prediction of High-Energy Electron Integral Flux at Geostationary Orbit Based on Deep Learning, *Space Weather*, 16, 903–916, <https://doi.org/https://doi.org/10.1029/2018SW001829>, 2018.
- Welling, D. T., Ngwira, C. M., Opgenoorth, H., Haiducek, J. D., Savani, N. P., Morley, S. K., Cid, C., Weigel, R., Weygand, J. M., Woodroffe, J. R., Singer, H., Rosenqvist, L., and Liemohn, M.: Recommendations for Next-Generation Ground Magnetic Perturbation Validation, *Space Weather*, 16, 1912–1920, <https://doi.org/https://doi.org/10.1029/2018SW002064>, 2018.
- Zewdie, G. K., Valladares, C., Cohen, M. B., Lary, D. J., Ramani, D., and Tsidu, G. M.: Data-Driven Forecasting of Low-Latitude Ionospheric Total Electron Content Using the Random Forest and LSTM Machine Learning Methods, *Space Weather*, 19, e2020SW002639, <https://doi.org/https://doi.org/10.1029/2020SW002639>, 2021.

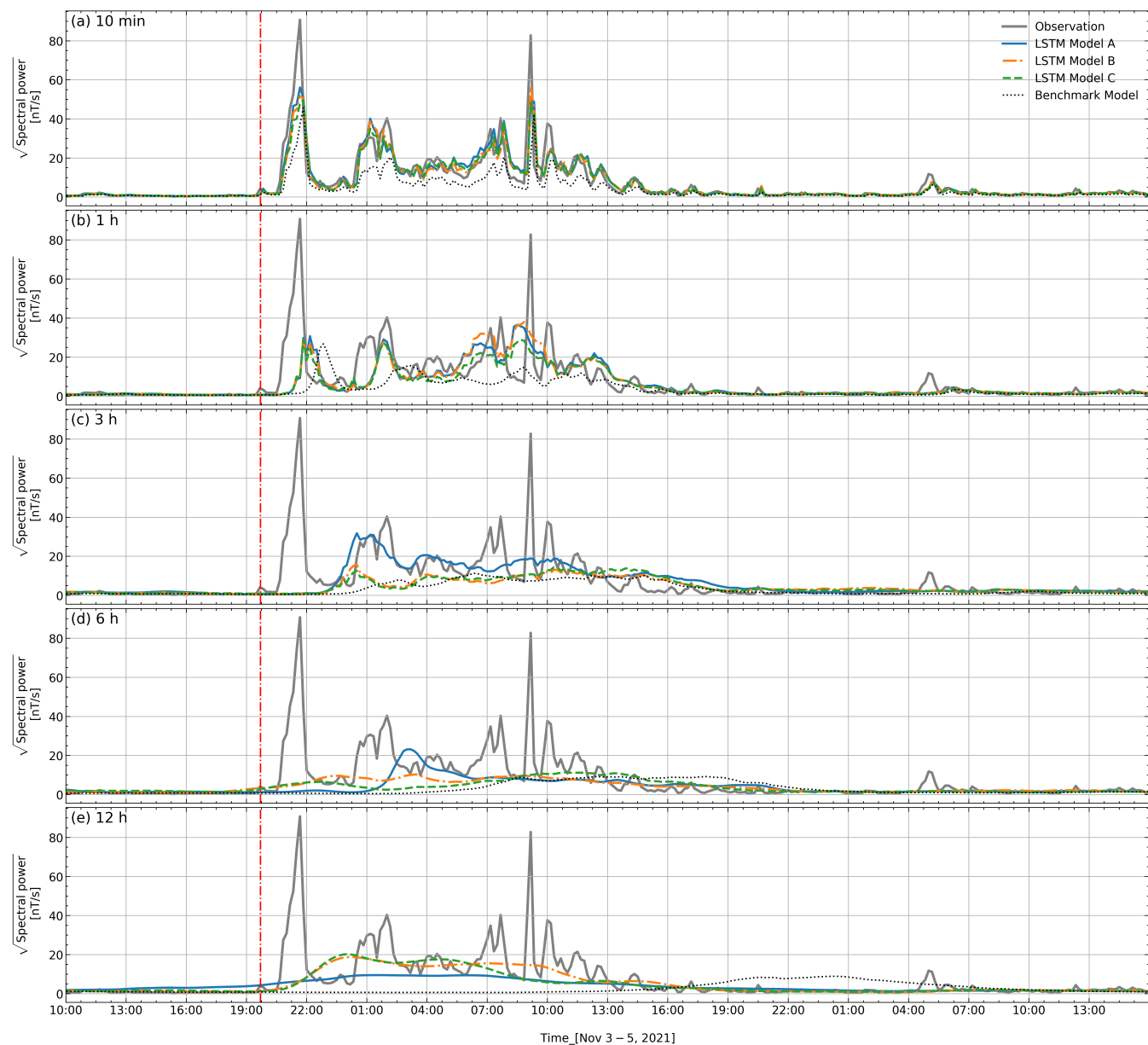


Figure 5. Comparison of observed and forecasted values of the GIC proxy P_0 (measured dH/dt power over 0–0.01 Hz) derived from Rørvik magnetometer measurements during November 3–5, 2021. Forecasts with time horizons of 10 min, 1 h, 3 h, 6 h, and 12 h are shown in panels a–e as indicated in the panel caption. Forecast outputs from Model A ("high-resolution" artificial space weather forecast inputs), Model B (low-resolution inputs), and Model C (high-resolution inputs except for predicted arrival time) are respectively indicated by solid blue lines, orange dash-dotted lines, and green dashed lines. In each panel The observed P_0 time series and the benchmark model predictions are respectively indicated by a thick, gray line and a black dotted line. CME arrival at 19:42 UT on November 3, 2021, is indicated by the vertical dash-dotted red line.

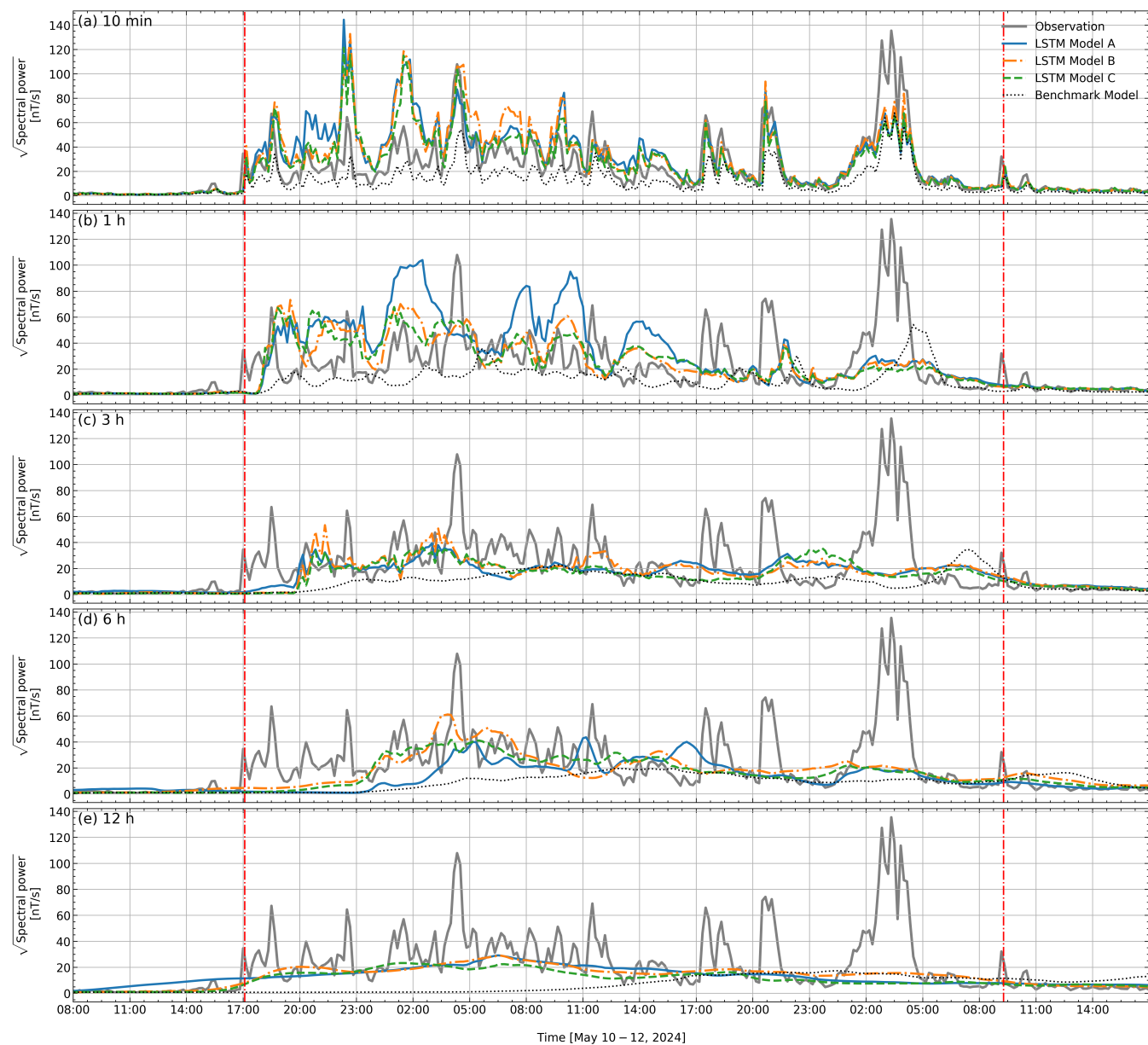


Figure 6. Measured dH/dt power over 0–0.01 Hz measured at Rørvik (thick, gray line) as well as model forecast and benchmark model forecast for May 10–12, 2024, in the same format as Figure 5. CME arrivals at 17:05 UT on May 10 and 09:17 UT on May 12 are indicated by vertical dash-dotted red lines.

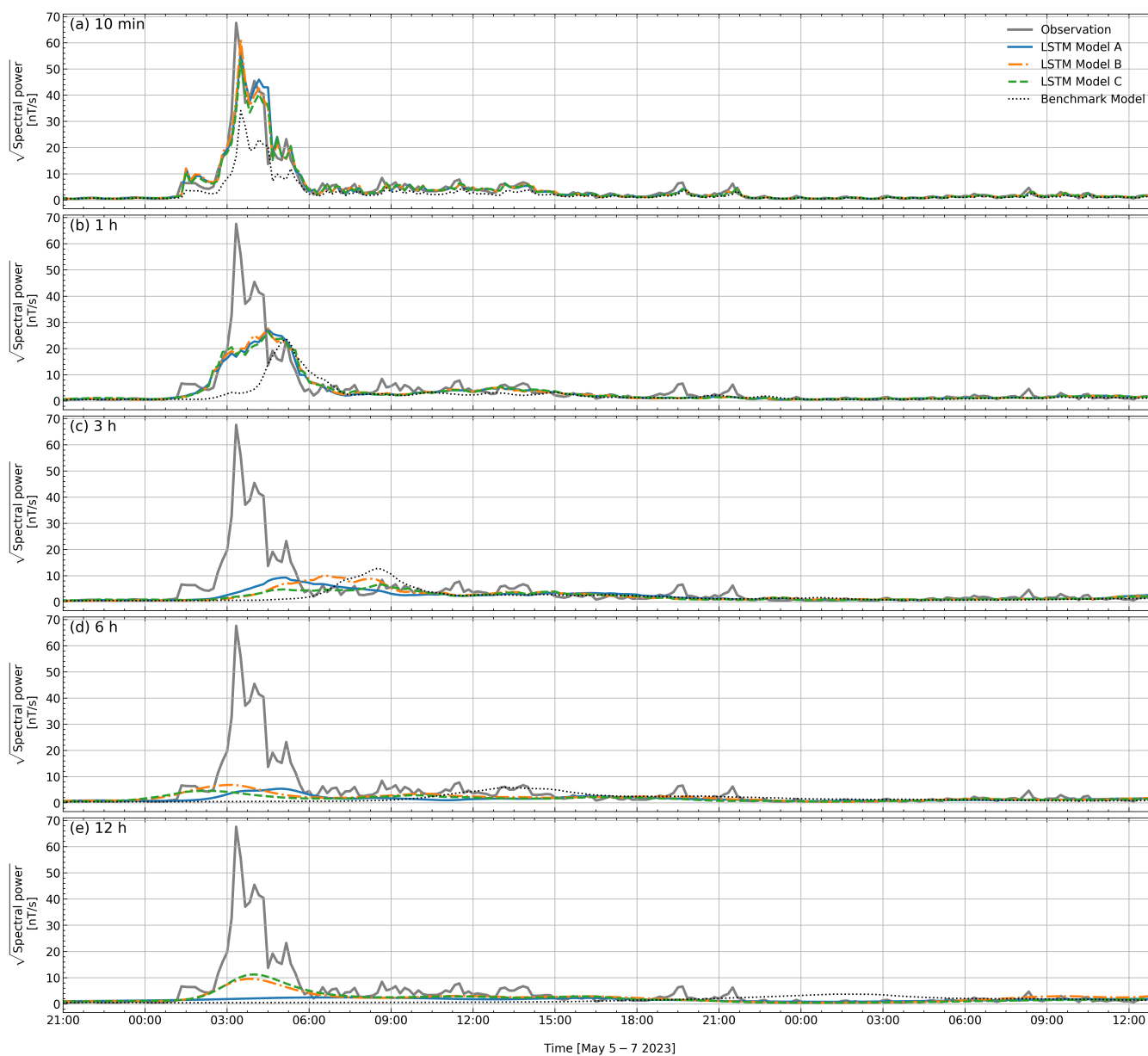


Figure 7. Measured dH/dt power over 0–0.01 Hz measured at Rørvik (thick, gray line) as well as model forecast and benchmark model forecast for May 5–7, 2023, in the same format as Figure 5.

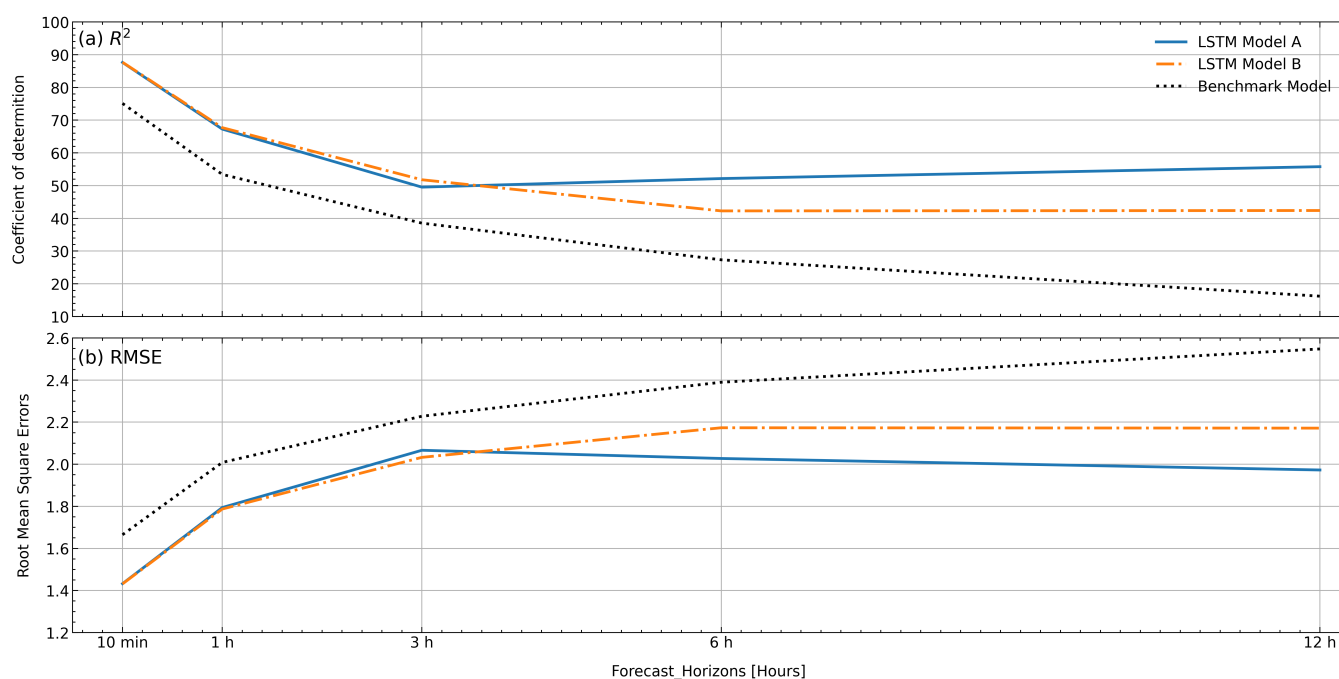


Figure 8. Statistical evaluation of model performance for each forecasting horizon h . (a) Coefficient of determination R^2 . (b) RMSE metrics for each model. Results for Model C are indistinguishable from those of Model A and are therefore omitted.

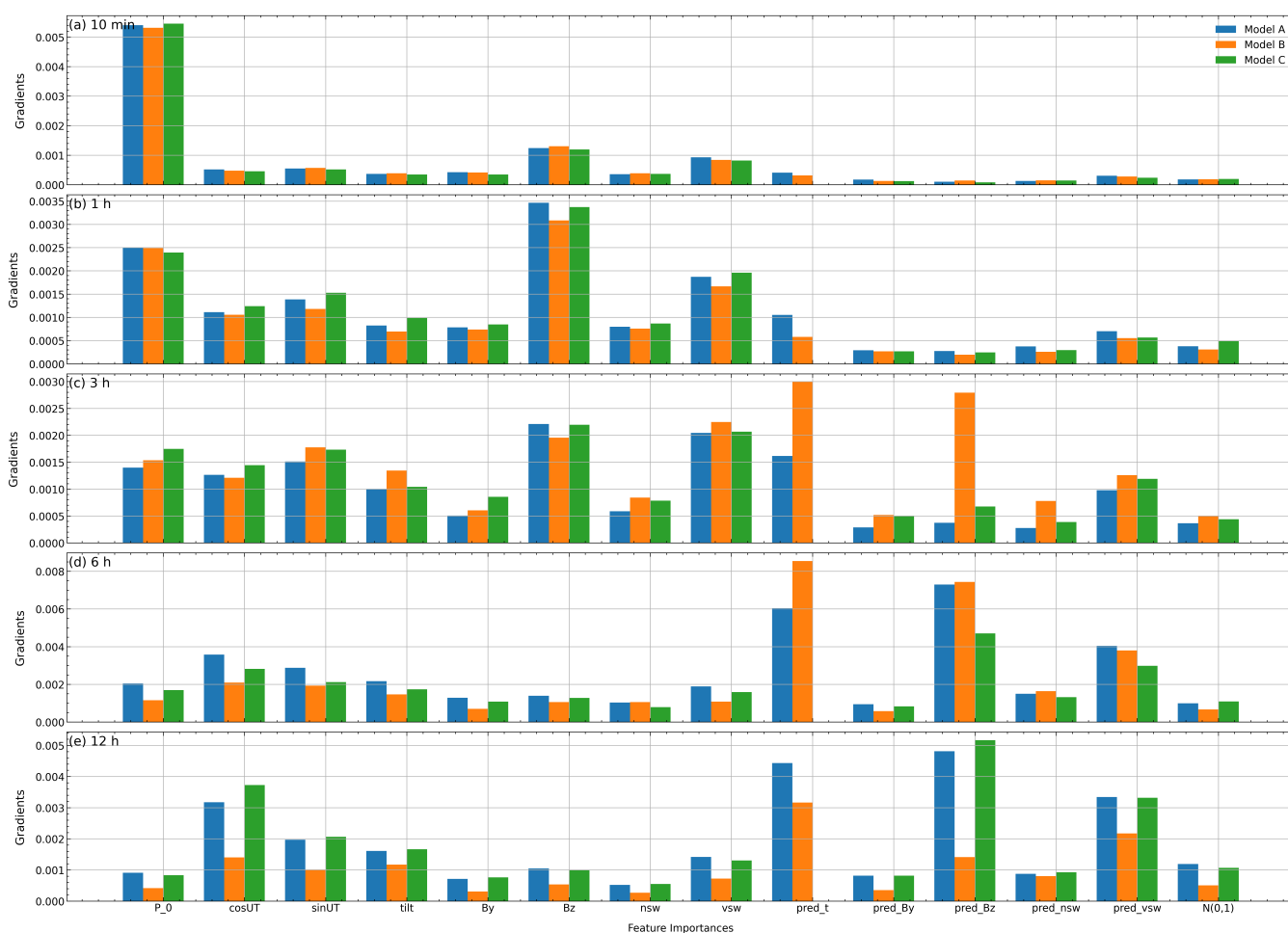


Figure 9. Feature importance for each model forecasting horizon defined using the gradient tape method described in Section 3.6. The y axis indicates the relative importance of each feature.