



**HAL**  
open science

## An operational approach to the forecasting of the dynamics of the Earth's radiation belts

Guillaume Bernoux, Antoine Brunet, Éric Buchlin, Miho Janvier, Angélica Sicard

► **To cite this version:**

Guillaume Bernoux, Antoine Brunet, Éric Buchlin, Miho Janvier, Angélica Sicard. An operational approach to the forecasting of the dynamics of the Earth's radiation belts. In press. hal-03242557v1

**HAL Id: hal-03242557**

**<https://hal.science/hal-03242557v1>**

Preprint submitted on 31 May 2021 (v1), last revised 11 Feb 2022 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# An operational approach to the forecasting of the dynamics of the Earth's radiation belts

Guillaume Bernoux<sup>1,\*</sup>, Antoine Brunet<sup>1</sup>, Éric Buchlin<sup>2</sup>, Miho Janvier<sup>2</sup>, and Angélica Sicard<sup>1</sup>

<sup>1</sup> ONERA / DPHY, Université de Toulouse, F-51005 Toulouse, France

<sup>2</sup> Université Paris-Saclay, CNRS, Institut d'Astrophysique Spatiale, Orsay, France

## ABSTRACT

The *Ca* index is a time-integrated geomagnetic index that correlates well with the dynamics of high-energy electron fluxes in the outer radiation belts. Therefore it can be used as an indicator for the state of filling of the radiation belts for those electrons, with the advantage of being a ground-based measurement with extensive historical records. In this work we propose a data-driven model to forecast *Ca* up to 24 hours in advance from near-Earth solar wind parameters. Our model relies mainly on a recurrent neural network called Long Short Term Memory network that has shown good performances in forecasting other geomagnetic indices in previous papers. Most implementation choices in this study have been arbitrated from the point of view of a space system operator, including the data selection and split, the definition of a binary classification threshold and the evaluation methodology. We evaluate our model (against a linear baseline) using both classic and novel (in the space weather field) measures. In particular, we use the Temporal Distortion Mix (TDM) to assess the propensity of two time series to exhibit time lags. We also evaluate the ability of our model to detect storm onsets during calm periods. It is shown that our model has a high overall accuracy, with evaluation measures deteriorating in an almost linear trend over time. However, using the TDM and the metrics for evaluating the binary classification forecast, it is shown that after a time horizon of approximately 6 hours the forecasts lose some of their usefulness in an operational context. This behaviour was not observable when evaluating the model only with metrics such as the root-mean-square error or the Pearson linear correlation. Considering the physics of the problem, this result is not surprising and suggests that the use of more spatially remote data (such as solar imaging) could improve space weather forecasts.

**Key words.** space weather – forecasting – radiation belts – machine learning – solar wind

## 1. Introduction

One of the current main topics of interest in the space weather field is the forecasting of geomagnetic indices based on machine learning methods. These methods have allowed for a great improvement

\* Corresponding author: [guillaume.bernoux@onera.fr](mailto:guillaume.bernoux@onera.fr)

33 in short-term forecasts of geomagnetic indices such as the global index  $Kp$  (Wintoft et al., 2017;  
34 Tan et al., 2018; Chakraborty and Morley, 2020) or  $Dst$  index (Gruet et al., 2018; Lethy et al.,  
35 2018). Space weather induced events can have heavy-to-extreme consequences on human-made  
36 infrastructures, as for instance space-borne hardware or even ground-based facilities (Riley et al.,  
37 2017). That is why the reliable forecast of geomagnetic indices and other space-weather relevant  
38 physical quantities (e.g. relativistic electron or proton fluxes in the radiation belts) is of paramount  
39 importance.

40 The extent of the effects of the space radiative environment on satellites ranges from single events  
41 caused by high energy charged particles from cosmic rays or solar energetic particles (SEP) to in-  
42 ternal charging, surface charging or total ionising dose (Horne et al., 2013). Therefore, being able  
43 to accurately and reliably forecast the fluxes of high-energy electrons (from dozens of kiloelectron-  
44 volts to a few megaelectronvolts) in the radiation belts would represent a great leap towards a better  
45 mitigation of the radiation-induced risks in space. Extensive efforts have already been conducted  
46 to forecast such electron fluxes. A considerable review of the methods used to forecast these elec-  
47 trons was recently proposed by Camporeale (2019), where it is detailed that feed-forward neural  
48 networks and recurrent neural networks (RNNs) are used to obtain forecasts up to a few hours or a  
49 few days ahead (see e.g. Ling et al. (2010); Wei et al. (2018)).

50 However Camporeale (2019) notes that although many approaches have been tested, it remains  
51 difficult to predict these fluxes due in particular to certain physical phenomena that are difficult to  
52 take into account for a "black-box" type model. Thus, many more recent models based on machine  
53 learning methods do not seem to perform better than older models. In addition, using data-driven  
54 approaches to predict radiation belt dynamics with in-situ data is challenging since it is important  
55 to have large databases that are properly calibrated (which is more complicated when using space-  
56 borne instruments rather than ground-based ones).

57 Recently, Bernoux and Maget (2020) have proposed a new time-integrated geomagnetic index  
58 which aims to be more representative of the state of filling of radiation belts. Thus we propose in  
59 this study to focus on the prediction of the dynamics of radiation belts through this so-called  $Ca$   
60 index. In order to do this we will use Deep Learning methods that have already been successfully  
61 tested with other geomagnetic indices. However, and in contrast to other studies, we will concentrate  
62 on evaluating our models by trying to take into account the point of view of a spacecraft operator  
63 and therefore using evaluation methods other than the classic metrics such as the root-mean-square  
64 error and the linear correlation, which can only account for global behaviour and are consequently  
65 largely insufficient to quantify other phenomena such as time shifts.

66 In this work we create a neural-network-based model to forecast the  $Ca$  index up to 24 hours in  
67 advance. We then evaluate the model using classical metrics and also using a method to detect the  
68 systematic existence of time shifts in our predictions. We then transform the regression problem  
69 into a binary classification problem aimed at predicting danger periods in terms of surface charging  
70 and we evaluate it accordingly. In Section 2 we present the data sets we used in our models and  
71 explain why they have been chosen and how they have been pre-processed. In Section 3 we present  
72 the regression and classification models and their dedicated evaluation methods. In Section 4 we  
73 present and discuss the results before concluding in Section 5.

## 74 2. Data analysis

75 In this section we describe and analyse the data sets used in this paper. Firstly we list the solar wind  
76 parameters and geomagnetic indices used here and explain where and how they can be obtained.  
77 Then we focus on the geomagnetic index  $Ca$  and explain its relevance and why it has been chosen  
78 for this study. Finally we explain how the time periods used for the training and the evaluation of  
79 the different models were selected.

### 80 2.1. Data sets

81 It has now been well known for decades that the geomagnetic indices representing the state of  
82 the magnetosphere are predominantly driven by the solar wind dynamics (Akasofu, 1981; Baker  
83 et al., 1981). That is why, as in many other studies (e.g. Lundstedt and Wintoft, 1994; Wu and  
84 Lundstedt, 1997; Wing et al., 2005; Chandorkar et al., 2017; Chakraborty and Morley, 2020),  
85 we have decided to use solar wind parameters available in the OMNIweb database (King and  
86 Papitashvili, 2005) as inputs to our geomagnetic index forecast models. The OMNIweb database  
87 (<https://omniweb.gsfc.nasa.gov/>) grants access to hourly spacecraft-interspersed near-Earth  
88 measurements of solar wind parameters. Earliest solar wind parameters are available since late  
89 1963. In particular we selected the plasma bulk velocity  $V_{sw}$ , the ion density  $\rho$ , the southward com-  
90 ponent of the interplanetary magnetic field (IMF)  $B_z$  and the plasma temperature  $T$  as our inputs to  
91 our models. It is now well known that these parameters correlate well with geomagnetic indices and  
92 with the dynamics of electron fluxes in the radiation belts (Burton et al., 1975; Wing et al., 2016).

93 The geomagnetic index studied here is the  $Ca$  index, that was first introduced by Bernoux and  
94 Maget (2020) based on a previous study by Rochel et al. (2016). Therefore the following paragraphs  
95 rephrase some information on the purpose and relevance of this index that was contained in these  
96 papers.

97 The  $Ca$  index is an index derived from the well-known  $aa$  index. The  $aa$  index is a 3-hr K-  
98 based index first introduced by Mayaud (1971) and computed from data provided by two subauroral  
99 opposite observatories.  $aa$  index is the geomagnetic index having the longest available track record  
100 with data available since 1868. This gives us more than 150 years of homogeneous (Mayaud, 1980)  
101 and exploitable geomagnetic data with a time cadence of 3 hours. This is particularly useful when  
102 dealing with topics such as statistical analysis which require a great amount of data. In particular,  $aa$   
103 index covers a time range equivalent to 14 solar cycles. Nowadays, the  $aa$  index is made available by  
104 the International Service of Geomagnetic Indices (IGSI) and can be downloaded from their website  
105 ([http://isgi.unistra.fr/data\\_download.php](http://isgi.unistra.fr/data_download.php)).

As stated in Bernoux and Maget (2020),  $Ca$  index has been designed to quantify the geoeffective-  
ness of solar wind structures impacting the magnetosphere from the radiation belts perspective. The  
relaxation characteristic time in the radiation belts for high-energy electrons after a strong magne-  
tospheric disturbance is of the order of 4 days (Meredith et al., 2006; Rochel et al., 2016). Therefore  
the  $Ca$  index is defined as follows:

$$Ca(t) = \frac{1}{\tau} \int_0^{\infty} aa(t-t') e^{-\frac{t'}{\tau}} dt' \quad (1)$$

106 with  $\tau = 4$  days being the relaxation characteristic time and  $aa$  representing the geomagnetic activ-  
107 ity. Being directly derived from  $aa$  index,  $Ca$  index shares the same above-mentioned qualities and

108 properties. Further details on the interest and relevance of using  $Ca$  index is provided in subsection  
109 [2.2](#).

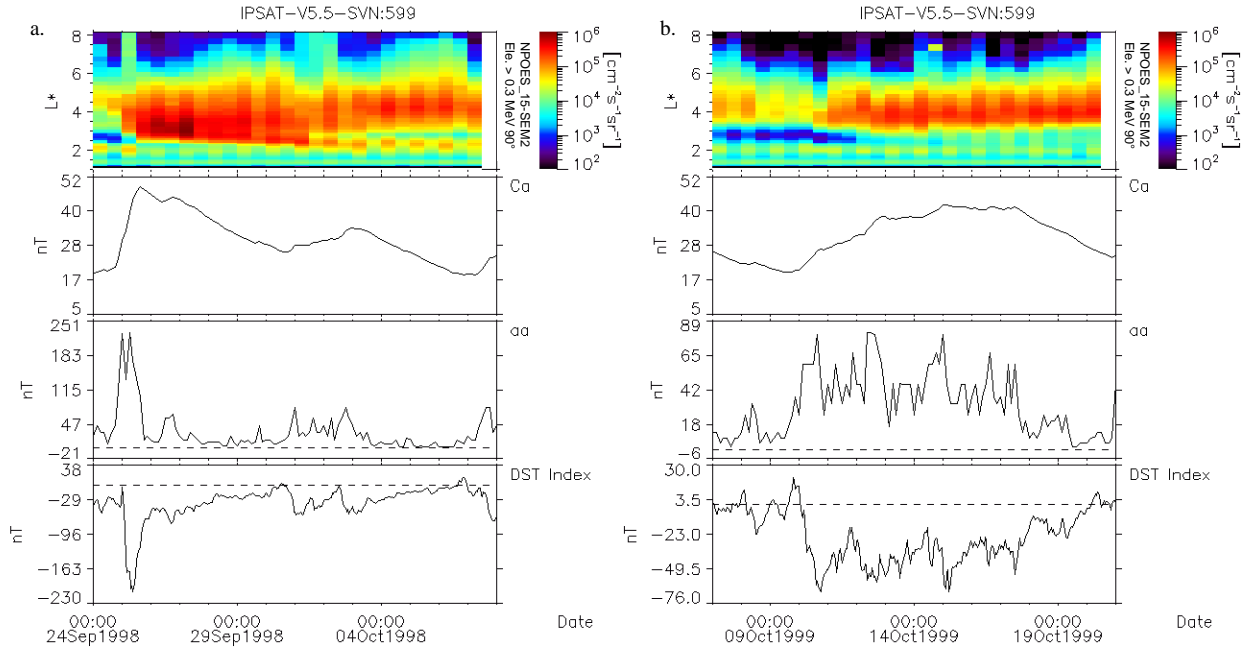
## 110 *2.2. Why study and forecast the $Ca$ index?*

111 Many studies have already been conducted on the topic of the nowcasting and forecasting of ge-  
112 omagnetic indices. Most of them focus on the  $Kp$  index or the  $Dst$  index, that are two very well  
113 known indices that have been thoroughly studied for decades. However it should be reminded that  
114 all geomagnetic indices are not interchangeable and that those indices have physical meanings. For  
115 instance, [Borovsky and Shprits \(2017\)](#) makes clear that the  $Dst$  index is unable to capture all types  
116 of geomagnetic storms behaviour and is in reality a very poor index when studying space-weather-  
117 relevant phenomena such as the dynamics of the electrons in the outer radiation belts induced by  
118 long-duration Corotating Interaction Regions (CIR)-driven storms. This is why it is important not to  
119 direct the research effort solely to the problem of forecasting the  $Kp$  and  $Dst$  indices, but to diversify  
120 the indices studied, in order to include a greater diversity of space-weather-relevant phenomena.

121 The  $Ca$  index was created to account for geomagnetic storms during which an intensification of  
122 relativistic electrons trapped in the radiation belts is observed. It was shown in [Rochel et al. \(2016\)](#)  
123 and [Bernoux and Maget \(2020\)](#) that this index correlates well with electron fluxes ( $E > 30$  keV) in  
124 the radiation belts and is able to take into account phenomena such as energy accumulation due to  
125 long duration Stream Interaction Region (SIR)-driven storms, but also due to multiple successive  
126 Interplanetary Coronal Mass Ejection (ICME)-driven events. Figure 1 display examples of the typ-  
127 ical behaviour of the  $Ca$  index during ICME- and SIR-driven storms. During ICME-driven storms,  
128 the  $aa$  index tends to reach higher values (in this example  $aa$  reaches 228 nT) quickly, but it also  
129 decreases rapidly, whereas during SIR-driven storms the disturbance lasts longer even though the  
130  $aa$  index usually does not reach such high values (in this example it only reaches 81 nT). Therefore  
131 the  $Ca$  index reaches its peak value much faster during the ICME-driven storm. However the value  
132 of the peak is similar during both these events as  $Ca$  accounts better for energy accumulation (48.6  
133 nT during the ICME-driven storm against 42.4 nT during the SIR-driven storm).

134 It was also stated in those papers that by changing the value of the parameter  $\tau$  it is possible to  
135 easily create an index that accounts better for a given specific orbit (but then less for the others). It  
136 is interesting to note that the  $Ca$  index is not the only attempt to create an index with such properties  
137 and another approach was proposed by [Borovsky and Yakymenko \(2017\)](#).

138 From an operational perspective, the prediction of the  $Ca$  index could serve as a basis for an alert  
139 service for the accumulation of high-energy electrons in the radiation belts. In such a context, the  $Ca$   
140 index would act as a proxy for relativistic electron fluxes, which is monitored from ground-based  
141 magnetometers. As stated in Section 1, using a data set that already has decades of cross-calibrated  
142 samples is also a great asset when dealing with data-driven approaches that require lots of data  
143 to be efficient. Besides, it may also be more reliable in terms of continuity of service to rely on  
144 ground-based instruments rather than on-board instruments that are subject to the risks associated  
145 with their being in space, at least as a back-up. Thus, the prediction of the  $Ca$  index is of immediate  
146 interest to the operators of space-borne systems.



**Fig. 1.** From bottom to top: evolution of the geomagnetic indices  $Dst$ ,  $aa$ , and  $C_a$ , and of the flux of electrons in the radiation belts for the  $E \geq 300$  keV energy range, measured by the SEM instrument aboard the POES-15 spacecraft a) from 24 September to 8 October 1998 during a period that displayed an ICME-induced disturbance starting on 25 September 1998, and b) from 7 October to 21 October 1999 during a period that displayed a SIR-induced disturbance starting on 9 October 1999.

### 147 2.3. Establishing the training, validation and test sets

#### 148 2.3.1. Splitting the data sets

149 In this subsection we briefly analyse the time series supplied by the OMNIweb database in order  
 150 to detect any important data gaps (that would be prejudicial for the training of a machine learning  
 151 algorithm) and to carefully chose the time periods used to train, validate and evaluate our models.  
 152 Dividing a data set into training, validation and test sets is very common practice in machine learn-  
 153 ing applications. If needed the reader is referred to [Carè and Camporeale \(2018\)](#) for more details.

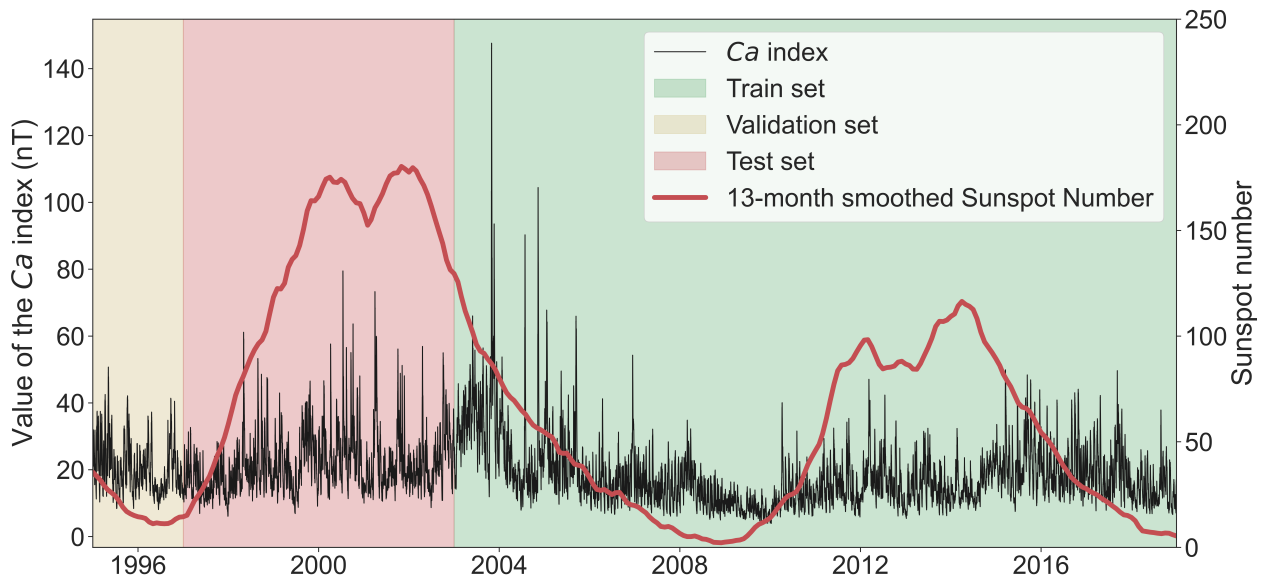
154 Before the availability of the Wind/Solar Wind Experiment (Wind/SWE) and the Advanced  
 155 Composition Explorer magnetometer and Solar Wind Electron, Proton, and Alpha Monitor  
 156 (ACE/MAG and ACE/SWEPAM) data starting in 1995 and 1998, the OMNIweb database has a  
 157 high percentage of missing data. Therefore in our study we only use data from 1995 onward. For  
 158 the 1995-2019 period there was in average 2.41% of missing data per year. The data gaps were  
 159 filled with a simple linear interpolation as the gaps were mostly very short ones.

160 The choice of the data that will be used to train, validate and test the neural network is of critical  
 161 importance. This includes the appropriate choice of how the data set will be temporally subdivided  
 162 into training, validation and test data sets ([Lazzús et al., 2017](#)). In order to correctly train a machine  
 163 learning algorithm, the training data set should be comprised of a representative period during which  
 164 all kinds of space weather phenomena, including extreme events were observed. The testing (and



165 the validation) period should also be comprised of both quiet and agitated periods. Eventually we  
 166 have chosen the following periods, highlighted in Figure 2:

- 167 – Training set: 2003-01-01 – 2018-12-31
- 168 – Validation set: 1995-01-01 – 1996-12-31
- 169 – Test set : 1997-01-01 – 2002-12-31



**Fig. 2.** Plot of the values taken by the  $Ca$  index between 1995 and 2018 included (black thin line). The training (green area), validation (yellow area) and test (red area) sets are highlighted. The 13-month smoothed Sunspot Number is also plotted as an indicator for the solar cycle (red thick line).

170 The train set is composed of 16 continuous years including the declining phase of one cycle and  
 171 a full second cycle. The train set includes several extreme and even most extreme events, including  
 172 the "Halloween storm" of November 2003 that reached a maximum value of  $Ca$  of 147.6 nT and  
 173 was found to be the only 1-in-100 year event (in terms of  $Ca$  index) witnessed since the beginning  
 174 of the Space Era (Bernoux and Maget, 2020). The validation set is composed of a 2-year long period  
 175 during a solar minimum. The test set is composed of 6 continuous years including the ascending  
 176 phase, the maximum and the beginning of the descending phase of a solar cycle. The test set in-  
 177 cludes intense and extreme storms ( $\geq 67$  nT), which is a good step towards a fair evaluation of our  
 178 model. The chosen split should ensure that our sets are representative enough of the space weather  
 179 phenomena that can be observed through  $Ca$ .

180 To evaluate our model in an even more detailed way, we divide the test set into subparts corre-  
 181 sponding to periods of disturbances induced on the one hand by ICMEs and on the other hand by  
 182 Stream Interaction Regions (SIRs), including CIRs. For this purpose we use the ICME database  
 183 provided by Chi et al. (2016) and the SIR database provided by Chi et al. (2018). These databases  
 184 include the time of beginning and time of ending for several ICME and SIR induced geomagnetic  
 185 disturbances between 1995 and 2015 (2016 for SIRs). According to these databases, 212 SIRs and

186 204 ICMEs were observed in the near-Earth environment between 1997 and 2002 included. In our  
 187 study we define an ICME- (respectively SIR-) induced disturbance period as the time period dur-  
 188 ing which an ICME- (respectively SIR-) induced geomagnetic disturbance has an influence on the  
 189 dynamics of the  $Ca$  index. The beginning of the disturbance period is given by the beginning of the  
 190 storm as indicated in the database. The ending of the disturbance period is given by adding  $\tau = 4$   
 191 days to the ending of the storm as indicated in the database. We can hence evaluate our models using  
 192 only the ICME- or SIR-induced disturbance periods and be able to better understand the accuracy  
 193 of our forecasts.

194 Let us note that the train set is composed of 139512 samples, the validation set of 16801 sam-  
 195 ples and the full test set of 51841 samples. The SIR-induced disturbance period includes 212  
 196 recorded SIRs, which makes 27776 samples, and the ICME-induced disturbance period includes  
 197 204 recorded ICMEs, which makes 24058 samples. 5407 samples belong both to the SIR-induced  
 198 period and to the ICME induced period.

### 199 2.3.2. Preprocessing the data

200 Before being fed into the neural-network based model, the data are processed as follows:

- 201 – We interpolate the values of the  $Ca$  index in order to have hourly values instead a value every  
 202 3 hours (this is meaningful since  $Ca$  is a very smooth time-integrated index and thus doing this  
 203 interpolation does not change neither the physics nor the statistics of the problem).
- 204 – Missing values in the other data sets are filled with linear interpolation.
- 205 – Inputs and outputs are rescaled so that their mean is 0 and their standard deviation is 1. The  
 206 weights for performing the transformation are calculated only from the training set data in order  
 207 not to include bias for validation and testing. This procedure is standard when working with  
 208 recurrent networks.

## 209 3. Models and evaluation methods

210 In this section we present the models used to predict the  $Ca$  index as well as the machine learning  
 211 algorithms used in these models. We also describe the methods and measures used to evaluate the  
 212 model.

### 213 3.1. Model description

214 The model developed in this study receives as input the past values of four solar wind parameters  
 215 listed in Subsection 2.1, namely the plasma bulk velocity ( $V_{sw}$ ), the ion density ( $\rho$ ), the southward  
 216 component of the interplanetary magnetic field (IMF)  $B_z$  and the plasma temperature ( $T$ ). Unlike  
 217 other studies we decided not to include the past values of the geomagnetic index as an input to  
 218 the models because we position ourselves in an operational-like context. Indeed, even though the  
 219 ISGI provides quick-look  $aa$  index values, reliance on two different data sources always presents a  
 220 higher risk of data unavailability from one source, which is prejudicial when establishing a near-real  
 221 time forecasting service. Ideally for such a service one would have both models (with and without



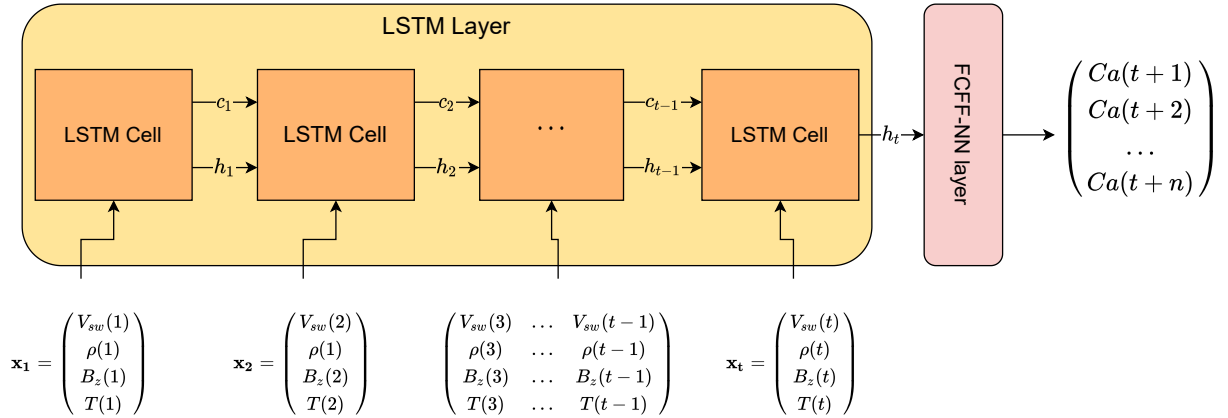
222 historical geomagnetic indices as inputs), but this is out of the scope of this study and for clarity we  
 223 only study one model in this paper. Here we have used the 30 last days for each input (*i.e.* the 720  
 224 last hourly values). The inputs/outputs link can be summarised as follows:

$$\begin{pmatrix} V_{sw}(t-719) & \dots & V_{sw}(t-1) & V_{sw}(t) \\ \rho(t-719) & \dots & \rho(t-1) & \rho(t) \\ B_z(t-719) & \dots & B_z(t-1) & B_z(t) \\ T(t-719) & \dots & T(t-1) & T(t) \end{pmatrix} \longrightarrow \begin{pmatrix} Ca(t+1) \\ Ca(t+2) \\ \dots \\ Ca(t+n) \end{pmatrix}, \text{ where } n \text{ is the forecast horizon.}$$

225 In section 4 we will analyse the results for a model trained and tested with a forecast horizon  
 226  $n = 24$  hours.

227 Our main model is a neural network-based model. It consists of a single layer Long-Short Term  
 228 Memory network (LSTM) combined with a linear fully-connected feed-forward (FCFF-NN) layer.  
 229 LSTMs are a type of recurrent neural networks first introduced in [Hochreiter and Schmidhuber](#)  
 230 (1997). LSTMs were created to address problems involving sequentially-structured data such as  
 231 time-series or natural language. In particular, LSTMs possess two internal memory states that are  
 232 designed to help addressing the gradient vanishing issue that occurs when handling long sequences  
 233 ([Hochreiter, 1998](#)). For an in-depth understanding of deep learning methods, including recurrent  
 234 and LSTM networks, the reader is referred to the above-mentioned papers as well as to reference  
 235 textbooks such as [Goodfellow et al. \(2016\)](#).

236 Our model is summarised in Figure 3.



**Fig. 3.** Simple scheme representing the LSTM-based model to forecast the values of the  $Ca$  index up to  $n$  hours in advance. The mechanism inside the LSTM cell was voluntarily not detailed.

237 Let us summarise the functioning of the LSTM network here. For each sample corresponding to  
 238 a time step  $t - p$ , the LSTM cell is fed with our solar wind parameters  $\mathbf{x}_{t-p}$  and the two memory  
 239 states computed at the previous time step: the hidden state  $h_{t-p}$  and the cell state  $c_{t-p}$ . The LSTM  
 240 cell processes and transforms the input and updates its hidden state and cell state (now  $h_{t-p+1}$  and  
 241  $c_{t-p+1}$ ) using three "gates": the input gate, the output gate, and the forget gate. To put it in simple  
 242 words, the LSTM cell decides which information from the past is "worth" being kept, forgotten,

243 or updated according to the last input. The latest memory states are again fed to the LSTM cell  
 244 along with the solar wind parameters at the next time step  $\mathbf{x}_{t-p+1}$ . After all time steps have been  
 245 given to the network, the LSTM layer outputs the final hidden state  $h_t$  that serves as the input to  
 246 the FCFF-NN layer, which itself outputs the  $t + 1$  to  $t + n$  next values of  $Ca$ ,  $n$  being the forecast  
 247 horizon.

248 Let us note that LSTMs have already demonstrated a good efficiency on geomagnetic index  
 249 prediction problems (see e.g. [Gruet et al., 2018](#); [Chakraborty and Morley, 2020](#); [Laperre et al.,](#)  
 250 [2020](#)).

251 Since this is the first study that focuses on the forecast of the  $Ca$  index, there is no immediate  
 252 baseline for us to compare our model to. The usual baseline used in such situation is the "persistence  
 253 model" (also known as the "naive model"), which simply consists in assuming that the predicted  
 254 value is the same as the last observed value. However that baseline cannot be pertinently used here as  
 255 we do not include the past values of  $Ca$  index among the inputs to our model. That is why we have  
 256 also trained a simple linear regression model to forecast the  $Ca$  index from the same solar wind  
 257 parameters as with the neural network-based model, with the notable exception that the baseline  
 258 linear model only uses the last value for each solar wind parameter as input (and not several past  
 259 values as with the neural network-based model).

### 260 *3.2. Training and parameters of the model*

261 Our model was trained using the classical backpropagation method ([Rumelhart et al., 1986](#)). The  
 262 optimisation method used is the Adam algorithm ([Kingma and Ba, 2017](#)). We have used a learning  
 263 rate  $lr = 5 \times 10^{-4}$  that is halved every 10 epochs. The loss function is the mean-square error (MSE).  
 264 The parameters of the model were hand-picked using cross-validation and iteration. We list below  
 265 the main parameters of our model and some implementation choices, so that the replicability of our  
 266 results is made easier. Let the reader be advised that even after changing some of these parameters  
 267 (e.g. in order to reduce the computational cost) it is possible to obtain very similar results.

- 268 – The LSTM cell state has dimension 256.
- 269 – The LSTM layer is mono-directional.
- 270 – We use L2-regularisation with weight  $5 \times 10^{-3}$ . L2-regularisation consists in adding the squared-  
 271 sum of the network's weights (with a multiplicative constant) to the loss function in order to  
 272 avoid overfitting.
- 273 – Size of each mini-batch: 256.
- 274 – The training is done with 30 epochs and with early stopping. Early stopping consists in stopping  
 275 the training of the network as soon as clear signs of overfitting are observed.

276 The model was developed using the PyTorch (v1.6) library for Python ([Paszke et al., 2019](#)).

### 277 *3.3. Detection of events*

278 Our models as described above offer predictions in the form of a regression problem. However,  
 279 it is more often useful for an end-user in a decision-making context to benefit from a predictive

280 alert system. The problem to be solved is then no longer a regression problem, but a classification  
 281 problem. In this study we will transform our prediction model into a simple binary classification  
 282 model (*i.e.* with only two classes) based on threshold detection: if we predict that  $Ca$  will exceed a  
 283 given threshold value during the next  $t$  hours then we issue an alert (class 1), if we predict that we  
 284 stay below this threshold then we issue no alert (class 0). The only difficulty lies in the choice of a  
 285 suitable threshold.

286 In our example we will choose a threshold based as much as possible on operational criteria.  
 287 The threshold must be meaningful to the end user, *i.e.* the triggering of an alert must correspond  
 288 to a situation for which the operator is expected to make a decision or take an action. As the  $Ca$   
 289 index represents the filling state of radiation belts with high energy electrons, we will choose a  $Ca$   
 290 threshold associated with a non-negligible risk of damage due to surface charging.

291 Figure 4 in [Bernoux and Maget \(2020\)](#) shows that  $Ca$  index has a quite high correlation coefficient  
 292 ( $R \approx 0.83$ ) with the dynamics of the integrated  $E \geq 30$  keV electron flux at  $L^* \approx 6$ . Moreover,  
 293 [Matéo-Vélez et al. \(2018\)](#) shows that the risk of damage due to surface charging for a spacecraft in  
 294 geostationary orbit (*i.e.* at  $L^* \approx 6$ ) is well correlated with the  $10 \leq E \leq 50$  keV electron flux when  
 295 the latter is greater than  $1 \times 10^8 \text{ cm}^{-2}\text{s}^{-1}\text{sr}^{-1}$ . A day during which the  $10 \leq E \leq 50$  keV electron  
 296 flux always stayed above this value has a minimum daily fluence of  $8.64 \times 10^{12} \text{ cm}^{-2}\text{sr}^{-1}$ . From this  
 297 value we define a fluence threshold equals  $8 \times 10^{12} \text{ cm}^{-2}\text{sr}^{-1}$ .

298 We then tried and find a  $Ca$  threshold that gives the highest correlation between the monthly  
 299 exceedances of the electron fluence and the monthly exceedances of the  $Ca$  threshold (using the  
 300 daily  $Ca$  maximum). For the daily fluences we have taken data provided by the Magnetospheric  
 301 Plasma Analyzer (MPA) instrument onboard the Geosynchronous Equatorial Orbit (GEO) LANL  
 302 1991-80 spacecraft between 1997 and 2006 for the energy range 35-46 keV ([McComas et al., 1993](#)).  
 303 It was found that the number of monthly fluence exceedances is best correlated with the monthly  
 304  $Ca$  exceedances when the  $Ca$  threshold is  $Ca_{\text{threshold}} = 38$  nT. This is also illustrated in Figure 4.

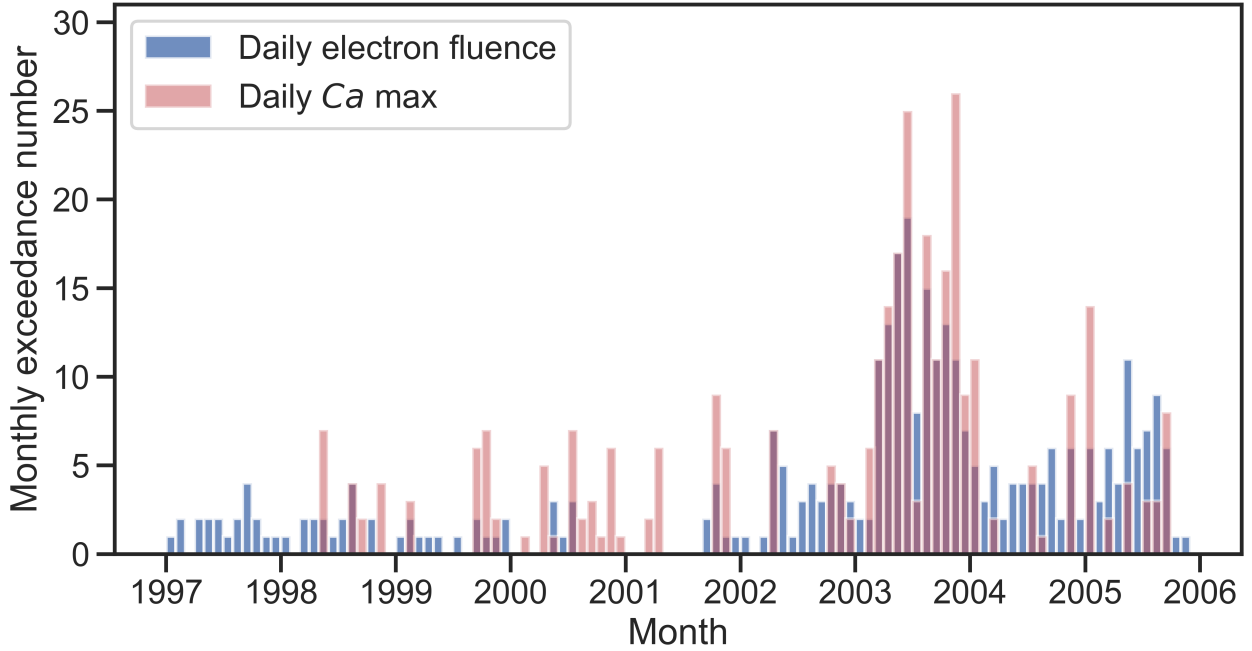
305 It should be noted in hindsight that the  $Ca$  value of 38 nT corresponds approximately to the  
 306 0.95 percentile of all  $Ca$  values, which seems statistically satisfactory. Indeed, it is a value that is  
 307 therefore rare enough to make a credible and useful alert threshold (an operator would probably not  
 308 want to receive an alert when the  $Ca$  value only exceeds the median, for example). But it is also a  
 309 value that is not too high, which allows better learning for the neural network (indeed, the higher the  
 310 threshold, the fewer samples we would have to train and evaluate the model). Let us also insist on  
 311 the fact that this threshold value used to define our binary classes in our study is only an example,  
 312 and that depending on the effect considered (internal charging, surface charging, singular events,  
 313 etc), the orbit considered, or even the satellite considered (and thus its structure) it would be more  
 314 interesting to use other thresholds, and probably to increase the number of classes.

### 315 3.4. Model evaluation

316 In this subsection we describe the measures used to evaluate the forecast performance of our models.

#### 317 3.4.1. Regression metrics

318 Since our problem is designed as a regression problem we first evaluate our model using two very  
 319 common regression metrics: the root-mean-square error (RMSE) and the Pearson (linear) correla-



**Fig. 4.** Count of days per month for which LANL 1991-80/MPA instrument measured a daily electron fluence above  $8 \times 10^{12} \text{ cm}^{-2} \text{ sr}^{-1}$  along with the count of days per month for which the daily  $Ca$  max was above 38 nT.

320 tion coefficient ( $R$ ). Let us define  $y_i$  the real observed values and  $\bar{y}_i$  the values forecast by a model  
 321 for  $i \in 1, \dots, N$ ,  $N$  being the number of samples.

- The RMSE is a measure of the global accuracy of the model, with more emphasis put on higher values (e.g. here the emphasis is on periods of more intense geomagnetic activity). A lower RMSE means a more accurate forecast. The RMSE is given by:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (\bar{y}_i - y_i)^2} \quad (2)$$

- The Pearson correlation indicates if the forecast values globally follow the same trends as the real values. The Pearson correlation ranges between 0 and 1 (higher is better). It is given by:

$$R = \frac{\text{Cov}(\bar{y}_i, y_i)}{\sqrt{\text{Var}(\bar{y}_i) \times \text{Var}(y_i)}} \quad (3)$$

322 Both metrics are widely used in the geomagnetic indices forecasting literature (e.g. in [Lazzús](#)  
 323 [et al. \(2017\)](#); [Tan et al. \(2018\)](#); [Gruet et al. \(2018\)](#); [Sexton et al. \(2019\)](#)). However, these metrics do  
 324 not capture the full performance of a model in all situations. Indeed, these metrics indicate overall  
 325 trends. Most of the time geomagnetic activity is fairly quiet, so quiet periods will weigh much more  
 326 heavily on the evaluation metrics than periods of high activity, thus creating a bias. While it is  
 327 very interesting for a satellite operator to be able to accurately predict quiet periods, it is also very

328 important to be able to accurately predict periods of geomagnetic disturbance. This type of bias  
 329 can be partially counterbalanced by taking adapted test sets, as we have done in Section 2.3. In the  
 330 following subsections we describe two other methods for evaluating our predictions that allow us  
 331 to better capture other types of behaviours.

### 332 3.4.2. Measuring time lags

333 Some studies such as [Wintoft and Wik \(2018\)](#) and [Laperre et al. \(2020\)](#) highlight the fact that  
 334 some forecasting models that display a great RMSE or Pearson correlation actually fail to reliably  
 335 forecast high disturbance periods in advance. [Laperre et al. \(2020\)](#) shows that some prediction  
 336 models exhibit systematic time lags between the observed time series and the predicted time series.  
 337 This systematic time lag would most often be of the order of magnitude of the model’s prediction  
 338 horizon. This would indicate that the model in reality would fail to predict a disturbance before it  
 339 has actually been observed, which is of very limited interest to an operator.

340 To quantify this behaviour [Laperre et al. \(2020\)](#) use the Dynamic Time Warping (DTW) algo-  
 341 rithm, which measures the time difference between two time series ([Berndt and Clifford, 1994](#)).  
 342 By applying this algorithm to the observed series and the predicted series shifted successively  
 343 by several consecutive time steps the authors are able determine the extent of the systematic lag.  
 344 Nonetheless in our study we do not use the exact same approach but a very similar one. Indeed,  
 345 the main drawback of the DTW method is that for a given prediction horizon  $n$ , it requires circa  $n^2$   
 346 iterations of the DTW algorithm with different time shifts to accurately assess the systematic time  
 347 lag. Besides, the computational complexity of the DTW algorithm is high even with modern now  
 348 methods to fasten the computation of the DTW measure (e.g. [Gold and Sharir, 2018](#)). This is why  
 349 we use instead the Temporal Distortion Mix (TDM).

350 The Temporal Distortion Mix is a metric proposed in [Vallance et al. \(2017\)](#) to characterise the  
 351 propensity of a time series to be late or early relative to a reference series. This metric is also based  
 352 on the DTW algorithm. Based on this algorithm, [Frías-Paredes et al. \(2016\)](#) proposes the Temporal  
 353 Distortion Index (TDI), which indicates to what extent the two time series are systematically (or not)  
 354 late (or early). Unlike the approach proposed by [Laperre et al. \(2020\)](#), the TDI does not indicate  
 355 the value of a possible systematic time lag, but whether the two time series exhibit this type of  
 356 behaviour and to which extent. In return, there is no need for several computations of the DTW  
 357 measure as only one (per forecast horizon) is sufficient to get the TDI. [Guen and Thome \(2019\)](#)  
 358 has even suggested that the TDI could be used as a part of the loss function when training a neural  
 359 network but this is out of scope of our paper.

360 To obtain the TDM, the TDI is decomposed into two components, which characterise the lateness  
 361 and the advance, so that  $\text{TDI} = \text{TDI}_{adv} + \text{TDI}_{late}$ . The TDM is then given by:

$$\text{TDM} = 1 - 2 \times \frac{\text{TDI}_{adv}}{\text{TDI}} \quad (4)$$

362 The TDM is hence a normalised version of the TDI. It ranges between -1 and 1. Let  $\mathbf{s}_1$  and  $\mathbf{s}_2$  be  
 363 two time series.

- 364 – if  $\text{TDM}(\mathbf{s}_1, \mathbf{s}_2) = -1$  then  $\mathbf{s}_1$  is systematically in advance compared to  $\mathbf{s}_2$
- 365 – if  $\text{TDM}(\mathbf{s}_1, \mathbf{s}_2) = 1$  then  $\mathbf{s}_1$  is systematically late compared to  $\mathbf{s}_2$

366 – if  $TDM(\mathbf{s}_1, \mathbf{s}_2) = 0$  then both time series are temporally aligned

367 For instance, the TDM between a given time series and its corresponding naive forecast is always  
 368 1. A good forecast is hence a forecast that has a TDM close to 0. The TDM is a very interesting  
 369 evaluation measure since it only requires one run of the DTW algorithm and it is possible to compare  
 370 the TDM between several forecasts (e.g. several forecast horizons). The TDM was first introduced in  
 371 a study dealing with the topic of solar irradiance forecasting, which is also a time-series forecasting  
 372 problem that shares structural similarities with ours.

### 373 3.4.3. Evaluation of the classification-based alert system

374 As we have already established, in an operational context in space weather it is important not only  
 375 to have regression type predictions but also to have warning systems based on class predictions.  
 376 In Section 3.3 we discussed how to transform our regression problem into a binary classification  
 377 problem (with a threshold of  $Ca_{threshold} = 38$  nT). In order to evaluate this derived alert system we  
 378 use several following metrics and measures. TP, FP, FN and TN are the true positive, false positive,  
 379 false negative and true negative counts.

– the precision: it is the ratio of issued alerts that match a true threshold excess. It gives an indication of how relevant the issued alerts are. It ranges between 0 and 1. Higher is better. It is given by:

$$\text{precision} = \frac{TP}{TP + FN} \quad (5)$$

– the recall: it is the ratio of true threshold exceedances that match an issued alert. It gives an indication of our ability to issue relevant alerts. It ranges between 0 and 1. Higher is better. It is given by:

$$\text{recall} = \frac{TP}{TP + FP} \quad (6)$$

– the  $F_{score}$ : it is the harmonic mean of precision and recall. It ranges between 0 and 1. Higher is better. It is given by:

$$F_{score} = \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (7)$$

– the threat score (TS): it gives an indication of how well true threshold exceedances were forecast, penalising both false alarms and false negatives. It ranges between 0 and 1. Higher is better. It is given by:

$$TS = \frac{TP}{TP + FN + FP} \quad (8)$$

– the Heidke skill score (HSS): it could be seen as a generalised skill score, giving the overall accuracy of the model against that of a random model. It ranges between -1 and 1. Higher is better, 0 denotes no skill. It is given by:

$$HSS = \frac{2 \times (TP \times TN - FP \times FN)}{(TP + FN)(FN + TN) + (TP + FP)(FP + TN)} \quad (9)$$



380 – The percentage of threshold-exceedance periods for which the model actually issues an alert  
 381 before the threshold was exceeded (*i.e.* the amount of active periods that were forecast before  
 382 they started and not only forecast after the threshold was exceeded for the first time). This is not  
 383 a classical metric, but perhaps one of the most useful ones here, since this gives an indication of  
 384 how well the model is able to forecast disturbance periods before they happened, not including  
 385 the performance of the model once the disturbance period has already started. Let us note that  
 386 there are 42 disturbance onsets (above the threshold  $Ca = 38$  nT) in the test set.

## 387 4. Results and discussion

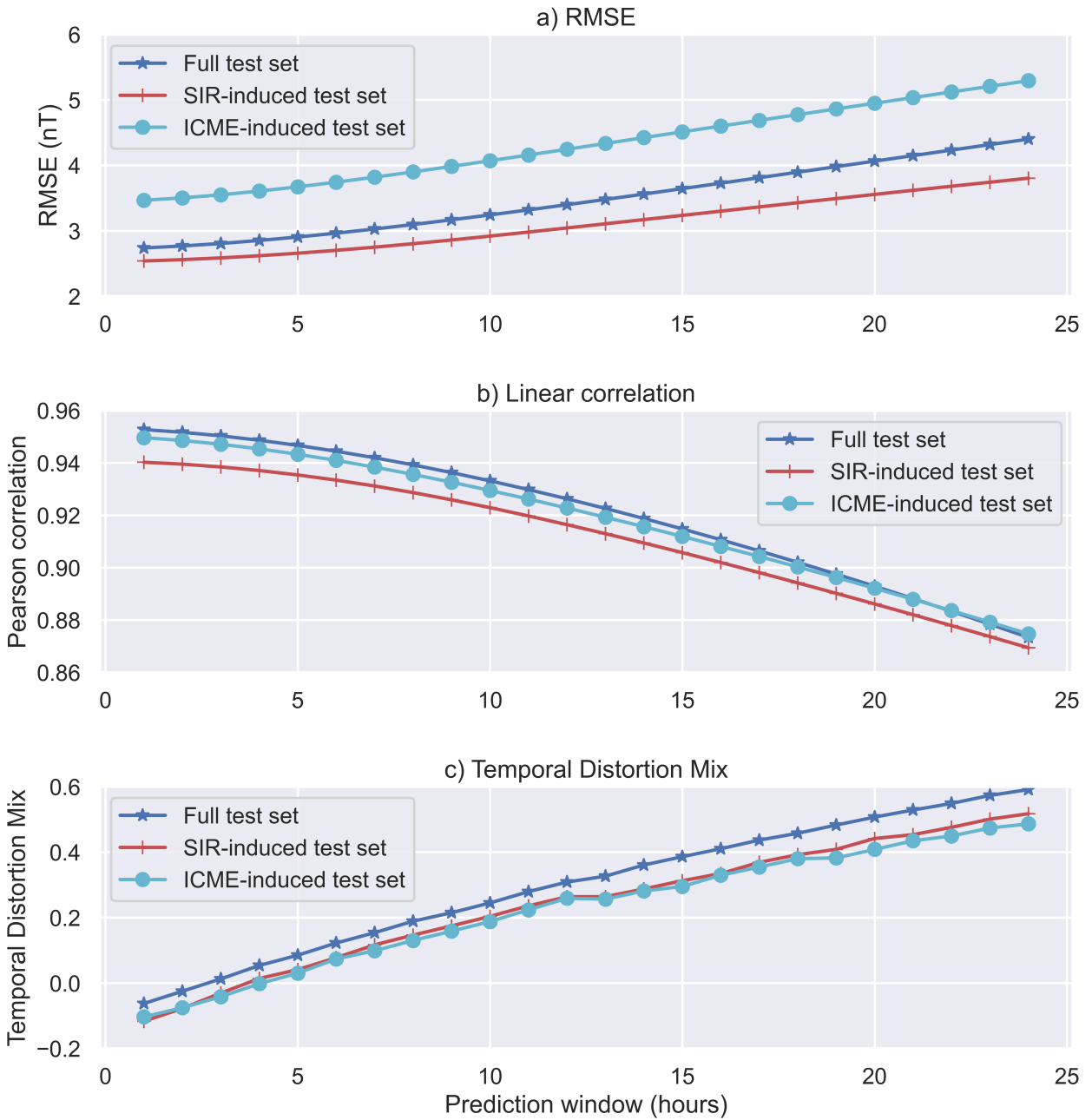
### 388 4.1. Regression results

389 The regression results obtained with the baseline model and the LSTM-NN model are presented  
 390 in Table 1. Firstly, we can see that the classical metrics (RMSE, R) give much better values with  
 391 the LSTM-NN model than with the linear baseline. For a time horizon of 3 hours, the RMSE with  
 392 the LSTM-NN is about 2.9 times lower than with the baseline (2.81 instead of 8.13), and for a  
 393 time horizon of 24 hours this ratio is 1.9 (4.40 instead of 8.16). This is an additional indication  
 394 to the fact that LSTM-NN networks are efficient for understanding the solar wind-magnetosphere  
 395 coupling. The RMSE values should be put into perspective with the statistical distribution of the  
 396  $Ca$  index, which over the test period has a variance of 8.9 nT and an interquartile range of 10.5 nT.  
 397 This comparison allows us to state that the RMSE values are satisfactory, especially for a model  
 398 that does not include the  $Ca$  index among its inputs. We also find that the Pearson correlation values  
 399 are quite high ( $\geq 0.9$  for all test sets up to a time horizon of 15 hours, instead of  $\leq 0.65$  with the  
 400 baseline), which is very satisfactory.

401 The TDM gives values close to 0 for a time horizon of 3 hours and up to 6 hours, for test sets  
 402 based on periods of disturbance. This indicates that up to about 6 hours, our forecasts are well  
 403 aligned in time with the target values. Beyond that, the TDM value increases up to 0.59 for a 24  
 404 hour time horizon with the full test set, indicating that there is an almost systematic delay between  
 405 the predicted values and the target values.

406 Unsurprisingly, the values of the conventional metrics all degrade as the time horizon increases.  
 407 This degradation (increase for RMSE and TDM, decrease for the Pearson correlation) appears to  
 408 be slow and smooth, as shown in Figure 5. However, for this reason it becomes difficult to tell from  
 409 these metrics alone from which time horizon the model is no longer operationally valid.

410 We also observe that, in general, the LSTM-NN model performs better during periods of SIR-  
 411 induced disturbances than during periods of ICME-induced disturbances. For a time horizon of 3  
 412 hours, the RMSE is 1.37 times higher for the ICME-induced period than for the SIR-induced period,  
 413 which is far from negligible. Figure 6 shows several examples of forecasts for two geomagnetic  
 414 storms: one induced by an ICME and the other by a SIR, the same storms already shown in Figure  
 415 1. This figure shows the forecast values for 4 different time horizons (3, 6, 12 and 24 hours) made  
 416 with both the LSTM-NN model and the linear baseline model. It is clear from this figure that the  
 417 neural network-based model outperforms the linear model, as already indicated by the evaluation  
 418 measures for the regression problem. In these examples the dynamics of the storm appear to be well  
 419 captured, and the forecast values are indeed close to the observed values, as indicated by the RMSE.  
 420 Furthermore it becomes apparent that the negative TDM values measured with the linear model are



**Fig. 5.** Evaluation of the LSTM-NN model with three measures (RMSE, R and TDM) for values of time horizon ranging from 1 hour to 24 hours. Three evaluation sets (full test set, SIR-induced set and ICME-induced set) were used.

421 due to the fact that the model has difficulty in correctly modeling the decay phase of a storm, which  
 422 decreases too fast and hence appears "ahead" in comparison to the true series.

423 Besides, the fact that the predicted (with the LSTM-NN model) and observed time series show  
 424 a time delay as the time horizon increases is evident in these examples. It would appear that this  
 425 time shift is more pronounced during the beginning of the disturbance period than during the decay

426 phase of the storm, which in the SIR-induced storm example remains well predicted even 24 hours  
 427 in advance. We should be able to better quantify this behaviour using the measures for the evaluation  
 428 of the classification problem.

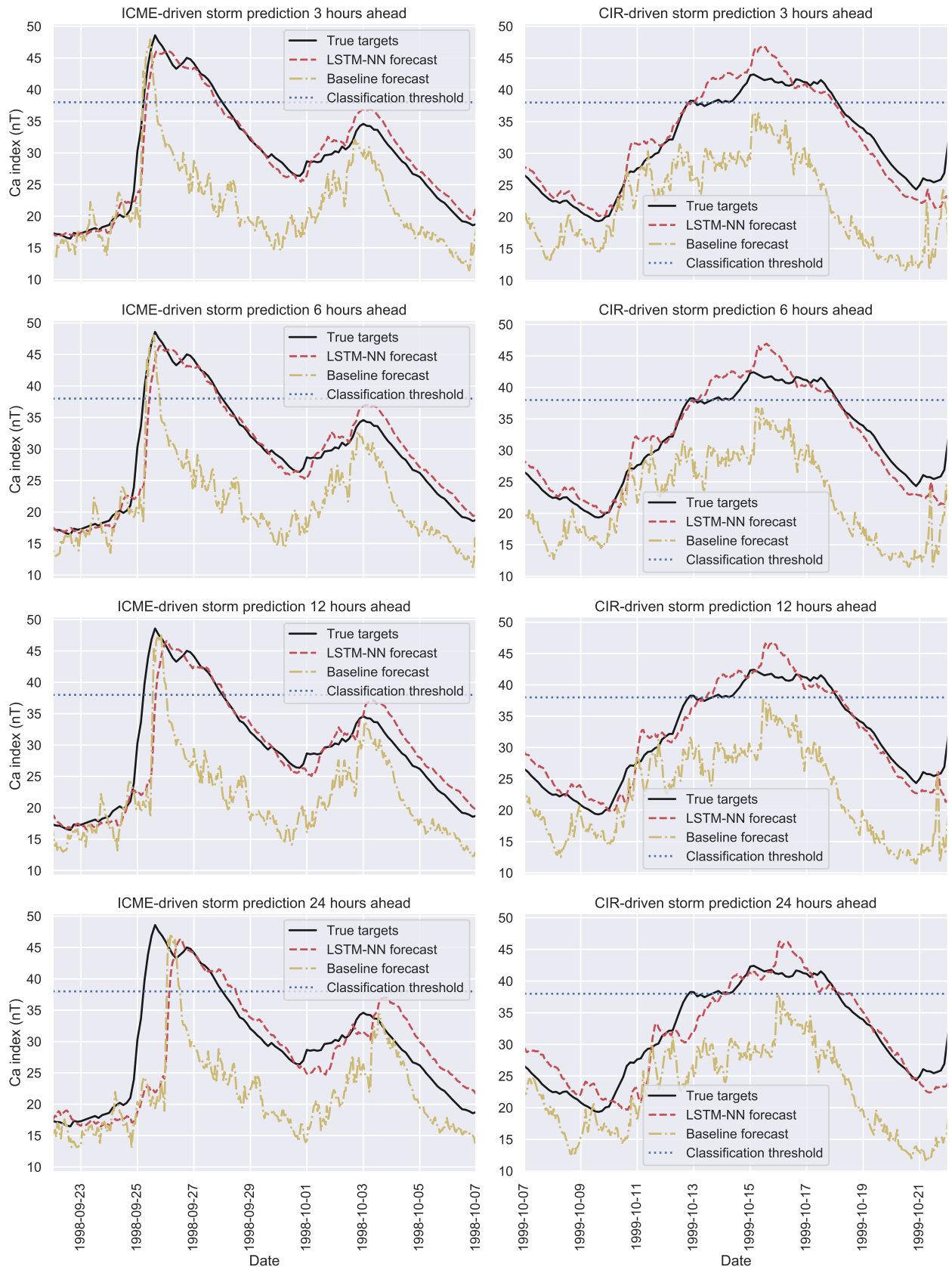
429 The difference of performance between ICME-induced and SIR-induced storms could hence at  
 430 least partly be explained by the fact that  $Ca$  increases more rapidly during ICME-induced distur-  
 431 bances. As indicated by the TDM values (and as we will see below with the classification measures),  
 432 the LSTM-NN model seems to be under-performing during the initial phase of a disturbance. Since  
 433 during SIR-induced disturbances the initial increase is slower than during ICME-induced distur-  
 434 bances, the RMSE during the beginning of the disturbance period should be lower in the first case,  
 435 which contributes to the overall RMSE being lower for the SIR-induced test set than for the ICME-  
 436 induced test set.

#### 437 4.2. Classification results

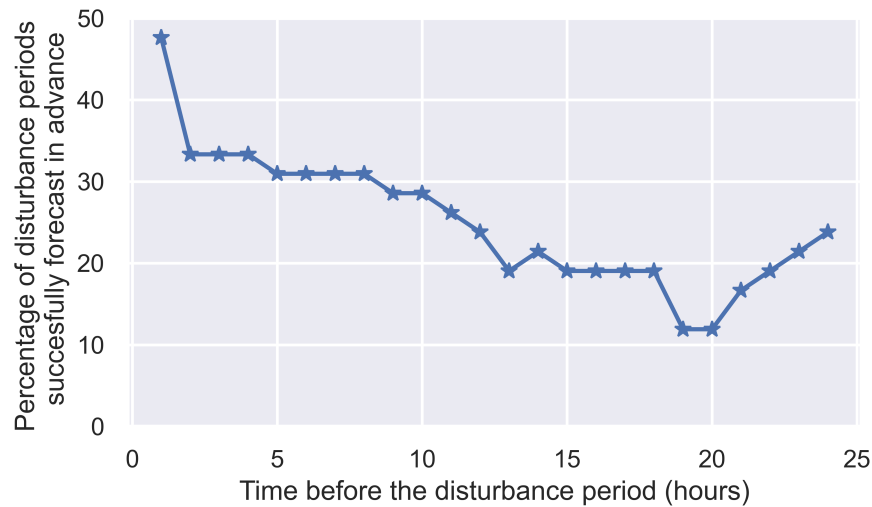
438 The classification results are given in Table 2 and Figure 7. For a time horizon of 3 hours, nearly  
 439 89% of the alerts issued were true positives, while 74% of the threshold exceedances were detected.  
 440 For a time horizon of 24 hours, these numbers rise and fall respectively to 93% and 63%. The fact  
 441 that the precision increases with the time horizon is due to the definition of our binary classes.  
 442 Indeed, we are trying to forecast if the threshold will be exceeded at any given time in the next  $t$   
 443 hours (and not at a precise given time). In our case, as the threshold increases, the model forecasts  
 444 less often true and false positives and more false negatives. That is why the precision increases  
 445 somewhat counter-intuitively. This highlights the need for several evaluation methods in order to  
 446 obtain a more exhaustive idea of the true performance of the model. Let us note that the  $F_{score}$ ,  
 447 which is the harmonic mean of precision and recall, decreases from 0.81 (for a time horizon of 3  
 448 hours) to 0.75 (for a time horizon of 3 hours), further indicating that the model performs better for  
 449 shorter time horizons.

450 It is difficult to argue at what percentage of precision and recall the model becomes satisfactory.  
 451 In absolute terms, correctly predicting more than two out of three periods of disturbance while  
 452 making only  $\approx 25\%$  false positives might seem to be a satisfactory target. However, depending  
 453 on the economic constraints due to spacecraft operation this could be largely insufficient. Here  
 454 we cannot definitively conclude about the absolute quality of our model but only about criteria that  
 455 would be defined by an operator and that depend on each space mission or on the targeted objective.  
 456 It should be noted, however, that the score values are also quite high, especially for the HSS. In  
 457 absolute terms, these values are rather difficult to interpret and should serve above all as a point of  
 458 comparison for possible future studies focusing on the forecast of similar physical quantities.

459 A result that is easier to interpret and that gives a user-friendly information is the percentage of  
 460 disturbance periods forecast in advance, given in Figure 7. To obtain this figure we calculated the  
 461 percentage of times and how long before the model was able to correctly answer the question: "will  
 462 the threshold be exceeded during the next 24 hours?" Therefore here we are only interested in the  
 463 model's ability to predict the beginning of a period of disturbance (without taking into account the  
 464 continuation of such a period). It appears that the model is able to answer this question correctly  
 465 less than 50% of the time 1 hour before the threshold is exceeded and less than 25% of the time 12  
 466 hours and longer before. This shows that even though 65% of the total exceedances were detected  
 467 somewhere between 1 hour and 24 hours before they happened, only less than one out of two



**Fig. 6.** Example of forecasts obtained with the LSTM-NN model and the linear model during two geomagnetic storms, the first one (left-hand side) being an ICME-driven storm and the second one being a CIR-driven storm (right-hand side). 4 different forecast horizons were used (3, 6, 12 and 24 hours). The value of  $Ca$  used for the binary classification is given in blue dotted line as a landmark.



**Fig. 7.** Percentage of times the 24h-binary classification problem was correctly forecast during calm periods previous to a threshold exceedance depending on how much time (from 1 hour to 24 hours) there was left before the exceedance.

468 disturbance periods were detected 1 hour before they happened and less than one out of four were  
 469 detected 24 hours before they happened. This is a much more significant measure of the operational  
 470 nature of our model and confirms the point we made earlier about the difficulty of predicting the  
 471 onset of a geomagnetic storm.

#### 472 4.3. Further discussion

473 In fact, the above-mentioned results are not very surprising since we rely on solar wind parameters  
 474 measured close to the Earth and thus the temporal hindsight to predict the dynamics of radiation  
 475 belts is small. This is reflected in the TDM measurements which indicate that the forecasts are  
 476 very well temporally aligned with the observations for forecast horizon values shorter than 6 hours,  
 477 which corresponds approximately to the reaction time of the geomagnetosphere interacting with a  
 478 disturbance arriving near Earth. We can therefore deduce on the one hand that our model seems to  
 479 be in agreement with the physics of the problem, but also that this same physics stops us, if we do  
 480 not change inputs, from having good operational performances for prediction horizons greater than  
 481 6 hours. The limit of 6 hours was also found in other papers dealing with the forecasting of the  $Dst$   
 482 index (Lazzús et al., 2017; Gruet et al., 2018). It would also be interesting to evaluate with these  
 483 methods (TDM and evaluating only the ability to predict the onset of a storm) the models presented  
 484 e.g. in Tan et al. (2018); Sexton et al. (2019) that aim at forecasting the  $Kp$  index up to 24 hours  
 485 in advance, in order to have a more comprehensive understanding of their actual effectiveness for  
 486 prediction horizons between 6 and 24 hours. Let us insist, however, on the fact that the difficulty  
 487 for these prediction horizons lies in the beginning of the storm and not in its continuity, because the  
 488 accumulation of energy makes it possible to find a link between the solar wind parameters and the  
 489 geomagnetic indices even after 6 hours of course. This is particularly the case with a time-integrated  
 490 index such as  $Ca$ , which allows for good overall forecast performances up to 24 hours in advance.

491 It might be tempting to compare our results to the results presented in e.g. [Forsyth et al. \(2020\)](#)  
492 where the authors present a deterministic model to forecast the GOES  $15 \geq 2$  MeV electron fluxes  
493 from solar wind data and also evaluate their model with classification measures. For instance, one  
494 of their models (when maximising the average Receiver Operating Characteristic score) for a time  
495 horizon of 6 hours gives a hit rate (or precision) of 0.75 whereas for the same time horizon ours  
496 give a higher hit rate of 0.87. However this comparison does not stand because we are not focusing  
497 on the same energy range and our model does not use the same classification thresholds and crite-  
498 ria. Indeed, here we answer the question: will the threshold be exceeded somewhere in the next  $t$   
499 hours? In [Forsyth et al. \(2020\)](#) the question is: will the threshold be exceeded in exactly  $t$  hours?  
500 We have chosen to approach the problem in this way because we believe that a warning system  
501 defined in this way is more useful, especially if we ask this question for several time horizons  $t$ .  
502 However this is an arbitrary choice and it could be argued otherwise. We wanted to stress here that,  
503 as highlighted in [Camporeale \(2019\)](#), comparing the performance of one model relative to another  
504 is not straightforward, and one should be cautious when doing it.



**Table 1.** Evaluation of the NN-based and the baseline models in the context of the regression problem. The model was evaluated with the full test set and also with the SIR-induced test set and the ICME-induced test set.

Time horizon (hours)	RMSE (nT)			R			TDM		
	Full	SIR	ICME	Full	SIR	ICME	Full	SIR	ICME
3	<b>2.81</b> (8.13)	<b>2.59</b> (6.37)	<b>3.55</b> (11.18)	<b>0.95</b> (0.63)	<b>0.94</b> (0.64)	<b>0.95</b> (0.64)	<b>0.01</b> (-0.42)	<b>-0.03</b> (-0.30)	<b>-0.04</b> (-0.55)
6	<b>2.96</b> (8.08)	<b>2.70</b> (6.37)	<b>3.74</b> (11.11)	<b>0.94</b> (0.64)	<b>0.93</b> (0.64)	<b>0.94</b> (0.64)	<b>0.12</b> (-0.37)	<b>0.08</b> (-0.26)	<b>0.07</b> (-0.51)
9	<b>3.17</b> (8.05)	<b>2.86</b> (6.39)	<b>3.98</b> (11.06)	<b>0.94</b> (0.64)	<b>0.93</b> (0.64)	<b>0.93</b> (0.65)	<b>0.21</b> (-0.33)	<b>0.17</b> (-0.22)	<b>0.16</b> (-0.50)
12	<b>3.40</b> (8.05)	<b>3.04</b> (6.43)	<b>4.25</b> (11.02)	<b>0.93</b> (0.64)	<b>0.92</b> (0.64)	<b>0.92</b> (0.64)	<b>0.31</b> (-0.28)	<b>0.26</b> (-0.17)	<b>0.26</b> (-0.46)
15	<b>3.64</b> (8.06)	<b>3.24</b> (6.48)	<b>4.51</b> (10.99)	<b>0.91</b> (0.64)	<b>0.91</b> (0.63)	<b>0.91</b> (0.64)	<b>0.39</b> (-0.23)	<b>0.31</b> (-0.12)	<b>0.30</b> (-0.42)
18	<b>3.90</b> (8.08)	<b>3.43</b> (6.54)	<b>4.77</b> (10.96)	<b>0.90</b> (0.63)	<b>0.89</b> (0.62)	<b>0.90</b> (0.64)	<b>0.46</b> (-0.18)	<b>0.39</b> (-0.06)	<b>0.38</b> (-0.35)
21	<b>4.15</b> (8.12)	<b>3.62</b> (6.59)	<b>5.04</b> (10.94)	<b>0.89</b> (0.62)	<b>0.88</b> (0.62)	<b>0.88</b> (0.63)	<b>0.53</b> (-0.14)	<b>0.45</b> (-0.03)	<b>0.43</b> (-0.31)
24	<b>4.40</b> (8.16)	<b>3.80</b> (6.65)	<b>5.29</b> (10.92)	<b>0.87</b> (0.61)	<b>0.87</b> (0.61)	<b>0.87</b> (0.62)	<b>0.59</b> (-0.11)	<b>0.52</b> (0.01)	<b>0.49</b> (-0.28)

**Notes.** The results obtained with the NN-based model are given in bold. The results obtained with the baseline are given in brackets.

**Table 2.** Evaluation of the NN-based and the baseline models in the context of the classification problem.

Time horizon (hours)	Precision	Recall	$F_{score}$	Threat score	Heidke Skill Score
6	<b>0.90</b> (0.66)	<b>0.72</b> (0.06)	<b>0.80</b> (0.12)	<b>0.67</b> (0.06)	<b>0.79</b> (0.11)
9	<b>0.91</b> (0.68)	<b>0.71</b> (0.06)	<b>0.80</b> (0.12)	<b>0.66</b> (0.06)	<b>0.78</b> (0.11)
12	<b>0.92</b> (0.69)	<b>0.69</b> (0.06)	<b>0.79</b> (0.11)	<b>0.65</b> (0.06)	<b>0.78</b> (0.10)
15	<b>0.92</b> (0.70)	<b>0.68</b> (0.06)	<b>0.78</b> (0.11)	<b>0.64</b> (0.06)	<b>0.77</b> (0.10)
18	<b>0.92</b> (0.70)	<b>0.66</b> (0.06)	<b>0.77</b> (0.11)	<b>0.63</b> (0.06)	<b>0.76</b> (0.10)
21	<b>0.93</b> (0.71)	<b>0.65</b> (0.06)	<b>0.76</b> (0.11)	<b>0.62</b> (0.06)	<b>0.75</b> (0.10)
24	<b>0.93</b> (0.71)	<b>0.63</b> (0.06)	<b>0.75</b> (0.11)	<b>0.60</b> (0.06)	<b>0.74</b> (0.10)

**Notes.** The results obtained with the NN-based model are given in bold. The results obtained with the baseline are given in brackets.

## 5. Conclusion

In this study we propose a recurrent-network based approach to forecast the fairly new geomagnetic index  $Ca$ . The main reason for focusing on this index is that this index is well correlated with the high-energy electron fluxes in the radiation belts and could hence be used as an indicator for their state of filling, without the drawbacks inherent to measuring *in-situ* fluxes with spacecrafts.

The implementation choices made in this paper were made by keeping in mind an operational context. These choices include the geomagnetic index to be forecast, the inputs used in our models and the whole evaluation methodology. To this end we have highlighted the importance of choosing statistically and physically representative train and test sets. We have also stressed the need to use adequate measures to evaluate the model, since classical metrics such as the RMSE or the Pearson correlation are not able to give an exhaustive report on the performance of the model, in particular during disturbance periods. That is why we use the Temporal Distortion Mix to measure the tendency for a forecast to be late or in advance in regards to the true observations.

We also transform the forecast problem from a regression problem to a binary classification one. The choice of the threshold used to define the binary classes was made taking into account risk for GEO spacecrafts to suffer damage from surface charging effect. The evaluation of the binary classification forecasts shows that even though the regression measures seemed great, the network does not show outstanding performance when it comes to forecasting the onset of a disturbance period. This is most certainly due to the spatial (and hence temporal) proximity between the solar wind parameters used as inputs and the geomagnetosphere. In order to improve the forecast results for time horizons of 12 hours, 24 hours and beyond it could be interesting to go back to the Sun and use data originating from solar imaging as inputs to a model. This topic will be the main focus of future studies. For now even though the measures are good and much better than the linear baseline, it would be difficult to claim that this model is fully adequate for use in an operational situation. It represents however a first and great step towards this purpose.

Other possibilities that remained out of scope of this study are the use of probabilistic forecasts (as done with other indices e.g. in Chandorkar et al., 2017; Chakraborty and Morley, 2020) or grey-box models. This paper being the first one dealing with the topic of forecasting the  $Ca$  index we voluntarily kept those possibilities aside for the sake of clarity and so as not to dilute the purpose of this study. However, we acknowledge that these are important avenues to explore, which will be done in future studies.

*Acknowledgements.* The authors are thankful to the NOAA-POES for online data access available on the CDAweb (at <http://cdaweb.gsfc.nasa.gov/>). The results presented in this paper rely on geomagnetic indices calculated and made available by ISGI Collaborating Institutes from data collected at magnetic observatories. We thank the involved national institutes, the INTERMAGNET network and ISGI ([isgi.unistra.fr](http://isgi.unistra.fr)). The OMNI data were obtained from the GSFC/SPDF OMNIWeb interface (at <https://omniweb.gsfc.nasa.gov>). Sunspot data from the World Data Center SILSO, Royal Observatory of Belgium, Brussels.

G. Bernoux is thankful for funding from Région Occitanie and ONERA, under Grant Agreements 19008721/ALDOCT and 30196.

545 **References**

- 546 Akasofu, S.-I., 1981. Prediction of Development of Geomagnetic Storms Using the Solar Wind-  
547 Magnetosphere Energy Coupling Function  $\epsilon$ . *Planetary and Space Science*, **29**(11), 1151–1158.  
548 10.1016/0032-0633(81)90121-5. [2.1](#)
- 549 Baker, D. N., E. W. Hones, J. B. Payne, and W. C. Feldman, 1981. A High Time Resolution Study  
550 of Interplanetary Parameter Correlations with AE. *Geophysical Research Letters*, **8**(2), 179–182.  
551 10.1029/GL008i002p00179. [2.1](#)
- 552 Berndt, D. J., and J. Clifford, 1994. Using Dynamic Time Warping to Find Patterns in Time Series. In  
553 Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining, AAAIWS'94,  
554 359–370. AAAI Press, Seattle, WA. [3.4.2](#)
- 555 Bernoux, G., and V. Maget, 2020. Characterizing Extreme Geomagnetic Storms Using Extreme Value  
556 Analysis: A Discussion on the Representativeness of Short Data Sets. *Space Weather*, **18**(6),  
557 e2020SW002,450. 10.1029/2020SW002450. [1](#), [2.1](#), [2.2](#), [2.3.1](#), [3.3](#)
- 558 Borovsky, J. E., and Y. Y. Shprits, 2017. Is the Dst Index Sufficient to Define All Geospace Storms? *Journal*  
559 *of Geophysical Research: Space Physics*, **122**(11), 11,543–11,547. 10.1002/2017JA024679. [2.2](#)
- 560 Borovsky, J. E., and K. Yakymenko, 2017. Systems Science of the Magnetosphere: Creating Indices of  
561 Substorm Activity, of the Substorm-Injected Electron Population, and of the Electron Radiation Belt.  
562 *Journal of Geophysical Research: Space Physics*, **122**(10), 10,012–10,035. 10.1002/2017JA024250. [2.2](#)
- 563 Burton, R. K., R. L. McPherron, and C. T. Russell, 1975. An Empirical Relationship between  
564 Interplanetary Conditions and Dst. *Journal of Geophysical Research (1896-1977)*, **80**(31), 4204–4214.  
565 10.1029/JA080i031p04204. [2.1](#)
- 566 Camporeale, E., 2019. The Challenge of Machine Learning in Space Weather: Nowcasting and Forecasting.  
567 *Space Weather*, **17**(8), 1166–1207. 10.1029/2018SW002061. [1](#), [4.3](#)
- 568 Carè, A., and E. Camporeale, 2018. Chapter 4 - Regression. In E. Camporeale, S. Wing, and J. R. Johnson,  
569 eds., *Machine Learning Techniques for Space Weather*, 71–112. Elsevier. ISBN 978-0-12-811788-0.  
570 10.1016/B978-0-12-811788-0.00004-4. [2.3.1](#)
- 571 Chakraborty, S., and S. K. Morley, 2020. Probabilistic Prediction of Geomagnetic Storms and the Kp Index.  
572 *Journal of Space Weather and Space Climate*, **10**, 36. 10.1051/swsc/2020037. [1](#), [2.1](#), [3.1](#), [5](#)
- 573 Chandorkar, M., E. Camporeale, and S. Wing, 2017. Probabilistic Forecasting of the Disturbance Storm  
574 Time Index: An Autoregressive Gaussian Process Approach. *Space Weather*, **15**(8), 1004–1019.  
575 10.1002/2017SW001627. [2.1](#), [5](#)
- 576 Chi, Y., C. Shen, B. Luo, Y. Wang, and M. Xu, 2018. Geoeffectiveness of Stream Interaction Regions From  
577 1995 to 2016. *Space Weather*, **16**(12), 1960–1971. 10.1029/2018SW001894. [2.3.1](#)
- 578 Chi, Y., C. Shen, Y. Wang, M. Xu, P. Ye, and S. Wang, 2016. Statistical Study of the Interplanetary Coronal  
579 Mass Ejections from 1995 to 2015. *Solar Physics*, **291**(8), 2419–2439. 10.1007/s11207-016-0971-5. [2.3.1](#)

- 580 Forsyth, C., C. E. J. Watt, M. K. Mooney, I. J. Rae, S. D. Walton, and R. B. Horne, 2020. Forecasting GOES  
581 15 >2 MeV Electron Fluxes From Solar Wind Data and Geomagnetic Indices. *Space Weather*, **18**(8),  
582 e2019SW002,416. 10.1029/2019SW002416. 4.3
- 583 Frías-Paredes, L., F. Mallor, T. León, and M. Gastón-Romeo, 2016. Introducing the Temporal Distortion  
584 Index to Perform a Bidimensional Analysis of Renewable Energy Forecast. *Energy*, **94**, 180–194.  
585 10.1016/j.energy.2015.10.093. 3.4.2
- 586 Gold, O., and M. Sharir, 2018. Dynamic Time Warping and Geometric Edit Distance: Breaking the Quadratic  
587 Barrier. *ACM Transactions on Algorithms*, **14**(4), 50:1–50:17. 10.1145/3230734. 3.4.2
- 588 Goodfellow, I., Y. Bengio, and A. Courville, 2016. *Deep Learning*. The MIT Press, Cambridge,  
589 Massachusetts, illustrated edition edn. ISBN 978-0-262-03561-3. 3.1
- 590 Gruet, M., M. Chandorkar, A. Sicard, and E. Camporeale, 2018. Multiple-Hour-Ahead Forecast of the Dst  
591 Index Using a Combination of Long Short-Term Memory Neural Network and Gaussian Process. *Space*  
592 *Weather*, **16**(11), 1882–1896. 10.1029/2018SW001898. 1, 3.1, 3.4.1, 4.3
- 593 Guen, V. L., and N. Thome, 2019. Shape and Time Distortion Loss for Training Deep Time Series Forecasting  
594 Models. *arXiv:1909.09020 [cs, stat]*. 1909.09020, URL <http://arxiv.org/abs/1909.09020>. 3.4.2
- 595 Hochreiter, S., 1998. The Vanishing Gradient Problem During Learning Recurrent Neural Nets and Problem  
596 Solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, **06**(02), 107–  
597 116. 10.1142/S0218488598000094. 3.1
- 598 Hochreiter, S., and J. Schmidhuber, 1997. Long Short-Term Memory. *Neural Computation*, **9**(8), 1735–1780.  
599 10.1162/neco.1997.9.8.1735. 3.1
- 600 Horne, R. B., S. A. Glauert, N. P. Meredith, D. Boscher, V. Maget, D. Heynderickx, and D. Pitchford,  
601 2013. Space Weather Impacts on Satellites and Forecasting the Earth’s Electron Radiation Belts with  
602 SPACECAST. *Space Weather*, **11**(4), 169–186. 10.1002/swe.20023. 1
- 603 King, J. H., and N. E. Papitashvili, 2005. Solar Wind Spatial Scales in and Comparisons of Hourly Wind  
604 and ACE Plasma and Magnetic Field Data. *Journal of Geophysical Research: Space Physics*, **110**(A2).  
605 10.1029/2004JA010649. 2.1
- 606 Kingma, D. P., and J. Ba, 2017. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*.  
607 1412.6980, URL <http://arxiv.org/abs/1412.6980>. 3.2
- 608 Laperre, B., J. Amaya, and G. Lapenta, 2020. Dynamic Time Warping as a New Evaluation for Dst Forecast  
609 With Machine Learning. *Frontiers in Astronomy and Space Sciences*, **7**. 10.3389/fspas.2020.00039, 2006.  
610 04667. 3.1, 3.4.2
- 611 Lazzús, J. A., P. Vega, P. Rojas, and I. Salfate, 2017. Forecasting the Dst Index Using a Swarm-Optimized  
612 Neural Network. *Space Weather*, **15**(8), 1068–1089. 10.1002/2017SW001608. 2.3.1, 3.4.1, 4.3
- 613 Lethy, A., M. A. El-Eraki, A. Samy, and H. A. Deebes, 2018. Prediction of the Dst Index and Analysis of  
614 Its Dependence on Solar Wind Parameters Using Neural Network. *Space Weather*, **16**(9), 1277–1290.  
615 10.1029/2018SW001863. 1

- 616 Ling, A. G., G. P. Ginet, R. V. Hilmer, and K. L. Perry, 2010. A Neural Network–Based Geosynchronous  
617 Relativistic Electron Flux Forecasting Model. *Space Weather*, **8**(9). 10.1029/2010SW000576. [1](#)
- 618 Lundstedt, H., and P. Wintoft, 1994. Prediction of Geomagnetic Storms from Solar Wind Data with the Use  
619 of a Neural Network. *Annales Geophysicae*, **12**(1), 19–24. 10.1007/s00585-994-0019-2. [2.1](#)
- 620 Matéo-Vélez, J.-C., A. Sicard, D. Payan, N. Ganushkina, N. P. Meredith, and I. Sillanpää, 2018. Spacecraft  
621 Surface Charging Induced by Severe Environments at Geosynchronous Orbit. *Space Weather*, **16**(1), 89–  
622 106. 10.1002/2017SW001689. [3.3](#)
- 623 Mayaud, P.-N., 1971. Une Mesure Planétaire d’activité Magnétique Basée Sur Deux Observatoires  
624 Antipodaux. *Annales Geophysicae*, **27**, 67–70. [2.1](#)
- 625 Mayaud, P.-N., 1980. Derivation, Meaning, and Use of Geomagnetic Indices. Geophysical Monograph ; 22.  
626 American Geophysical Union, Washington. ISBN 978-0-87590-022-3. [2.1](#)
- 627 McComas, D. J., S. J. Bame, B. L. Barraclough, J. R. Donart, R. C. Elphic, J. T. Gosling, M. B. Moldwin,  
628 K. R. Moore, and M. F. Thomsen, 1993. Magnetospheric Plasma Analyzer: Initial Three-Spacecraft  
629 Observations from Geosynchronous Orbit. *Journal of Geophysical Research: Space Physics*, **98**(A8),  
630 13,453–13,465. 10.1029/93JA00726. [3.3](#)
- 631 Meredith, N. P., R. B. Horne, S. A. Glauert, R. M. Thorne, D. Summers, J. M. Albert, and R. R. Anderson,  
632 2006. Energetic Outer Zone Electron Loss Timescales during Low Geomagnetic Activity. *Journal of*  
633 *Geophysical Research: Space Physics*, **111**(A5). 10.1029/2005JA011516. [2.1](#)
- 634 Paszke, A., S. Gross, F. Massa, A. Lerer, J. Bradbury, et al., 2019. PyTorch: An Imperative Style, High-  
635 Performance Deep Learning Library. *arXiv:1912.01703 [cs, stat]*. [1912.01703](#), URL <http://arxiv.org/abs/1912.01703>. [3.2](#)
- 637 Riley, P., D. Baker, Y. D. Liu, P. Verronen, H. Singer, and M. Güdel, 2017. Extreme Space Weather Events:  
638 From Cradle to Grave. *Space Science Reviews*, **214**(1), 21. 10.1007/s11214-017-0456-3. [1](#)
- 639 Rochel, S., D. Boscher, R. Benacquista, and J. F. Roussel, 2016. A Radiation Belt Disturbance Study from  
640 the Space Weather Point of View. *Acta Astronautica*, **128**, 650–656. 10.1016/j.actaastro.2016.07.012. [2.1](#),  
641 [2.2](#)
- 642 Rumelhart, D. E., G. E. Hinton, and R. J. Williams, 1986. Learning Representations by Back-Propagating  
643 Errors. *Nature*, **323**(6088), 533–536. 10.1038/323533a0. [3.2](#)
- 644 Sexton, E. S., K. Nykyri, and X. Ma, 2019. Kp Forecasting with a Recurrent Neural Network. *Journal of*  
645 *Space Weather and Space Climate*, **9**, A19. 10.1051/swsc/2019020. [3.4.1](#), [4.3](#)
- 646 Tan, Y., Q. Hu, Z. Wang, and Q. Zhong, 2018. Geomagnetic Index Kp Forecasting With LSTM. *Space*  
647 *Weather*, **16**(4), 406–416. 10.1002/2017SW001764. [1](#), [3.4.1](#), [4.3](#)
- 648 Vallance, L., B. Charbonnier, N. Paul, S. Dubost, and P. Blanc, 2017. Towards a Standardized Procedure to  
649 Assess Solar Forecast Accuracy: A New Ramp and Time Alignment Metric. *Solar Energy*, **150**, 408–422.  
650 10.1016/j.solener.2017.04.064. [3.4.2](#)
- 651 Wei, L., Q. Zhong, R. Lin, J. Wang, S. Liu, and Y. Cao, 2018. Quantitative Prediction of High-Energy  
652 Electron Integral Flux at Geostationary Orbit Based on Deep Learning. *Space Weather*, **16**(7), 903–916.  
653 10.1029/2018SW001829. [1](#)

- 654 Wing, S., J. R. Johnson, E. Camporeale, and G. D. Reeves, 2016. Information Theoretical Approach to  
655 Discovering Solar Wind Drivers of the Outer Radiation Belt. *Journal of Geophysical Research: Space*  
656 *Physics*, **121**(10), 9378–9399. 10.1002/2016JA022711. [2.1](#)
- 657 Wing, S., J. R. Johnson, J. Jen, C.-I. Meng, D. G. Sibeck, K. Bechtold, J. Freeman, K. Costello, M. Balikhin,  
658 and K. Takahashi, 2005. Kp Forecast Models. *Journal of Geophysical Research: Space Physics*, **110**(A4).  
659 10.1029/2004JA010500. [2.1](#)
- 660 Wintoft, P., and M. Wik, 2018. Evaluation of Kp and Dst Predictions Using ACE and DSCOVR Solar Wind  
661 Data. *Space Weather*, **16**(12), 1972–1983. 10.1029/2018SW001994. [3.4.2](#)
- 662 Wintoft, P., M. Wik, J. Matzka, and Y. Shprits, 2017. Forecasting Kp from Solar Wind Data: Input Parameter  
663 Study Using 3-Hour Averages and 3-Hour Range Values. *Journal of Space Weather and Space Climate*,  
664 **7**, A29. 10.1051/swsc/2017027. [1](#)
- 665 Wu, J.-G., and H. Lundstedt, 1997. Geomagnetic Storm Predictions from Solar Wind Data with the Use of  
666 Dynamic Neural Networks. *Journal of Geophysical Research: Space Physics*, **102**(A7), 14,255–14,268.  
667 10.1029/97JA00975. [2.1](#)