

DRAFT VERSION JULY 8, 2019

Typeset using L<sup>A</sup>T<sub>E</sub>X preprint style in AASTeX62

## A Comparison of Flare Forecasting Methods. II. Benchmarks, Metrics and Performance Results for Operational Solar Flare Forecasting Systems

K.D. LEKA,<sup>1,2</sup> SUNG-HONG PARK,<sup>1</sup> KANYA KUSANO,<sup>1</sup> JESSE ANDRIES,<sup>3</sup> GRAHAM BARNES,<sup>2</sup> SUZY BINGHAM,<sup>4</sup> D. SHAUN BLOOMFIELD,<sup>5</sup> AOIFE E. MCCLOSKEY,<sup>6</sup> VERONIQUE DELOUILLE,<sup>3</sup> DAVID FALCONER,<sup>7</sup> PETER T. GALLAGHER,<sup>8</sup> MANOLIS K. GEORGIOULIS,<sup>9,10</sup> YUKI KUBO,<sup>11</sup> KANGJIN LEE,<sup>12,13</sup> SANGWOO LEE,<sup>14</sup> VASILY LOBZIN,<sup>15</sup> JUNCHUL MUN,<sup>16</sup> SOPHIE A. MURRAY,<sup>6,8</sup> TAREK A.M. HAMAD NAGEEM,<sup>17</sup> RAMI QAHWAJI,<sup>18</sup> MICHAEL SHARPE,<sup>4</sup> ROB STEENBURGH,<sup>19</sup> GRAHAM STEWARD,<sup>15</sup> AND MICHAEL TERKILDSEN<sup>15</sup>

<sup>1</sup>*Institute for Space-Earth Environmental Research Nagoya University Furo-cho Chikusa-ku Nagoya, Aichi 464-8601 JAPAN*

<sup>2</sup>*NorthWest Research Associates 3380 Mitchell Lane Boulder, CO 80301 USA*

<sup>3</sup>*STCE - Royal Observatory of Belgium Avenue Circulaire, 3 B-1180 Brussels BELGIUM*

<sup>4</sup>*Met Office FitzRoy Road Exeter, Devon, EX1 3PB, UNITED KINGDOM*

<sup>5</sup>*Northumbria University Newcastle upon Tyne NE1 8ST, UNITED KINGDOM*

<sup>6</sup>*School of Physics, Trinity College Dublin, College Green, Dublin 2, IRELAND*

<sup>7</sup>*NASA/NSSTC Mail Code ST13 320 Sparkman Drive Huntsville, AL 35805, USA*

<sup>8</sup>*School of Cosmic Physics, Dublin Institute for Advanced Studies, 31 Fitzwilliam Place, Dublin, D02 XF86, IRELAND*

<sup>9</sup>*Department of Physics & Astronomy Georgia State University 1 Park Place, Rm #715, Atlanta, GA 30303, USA*

<sup>10</sup>*Academy of Athens 4 Soranou Efessiou Street, 11527 Athens, GREECE*

<sup>11</sup>*National Institute of Information and Communications Technology  
Space Environment Laboratory*

*4-2-1 NukuiKita Koganei Tokyo 184-8795 JAPAN*

<sup>12</sup>*Meteorological Satellite Ground Segment Development Center*

*Electronics and Telecommunications Research Institute, Daejeon*

*218 Gajeong-ro, Yuseong-gu, Daejeon, 34129, REPUBLIC OF KOREA*

<sup>13</sup>*Kyung Hee University*

*1732, Deogyong-daero, Giheung-gu, Yongin, 17104, REPUBLIC OF KOREA*

<sup>14</sup>*SELab, Inc. 150-8, Nonhyeon-ro, Gangnam-gu, Seoul, 06049, REPUBLIC OF KOREA*

<sup>15</sup>*Bureau of Meteorology Space Weather Services PO Box 1386 Haymarket NSW 1240 AUSTRALIA*

<sup>16</sup>*Korean Space Weather Center 198-6, Gwideok-ro, Hallim-eup, Jeju-si, 63025, REPUBLIC OF KOREA*

<sup>17</sup>*University of Bradford Bradford West Yorkshire BD7 1DP UK*

<sup>18</sup>*University of Bradford Bradford West Yorkshire BD7 1DP UNITED KINGDOM*

<sup>19</sup>*NOAA/National Weather Service National Centers for Environmental Prediction Space Weather Prediction Center, W/NP9 325 Broadway, Boulder CO 80305 USA*

(Received; Revised; Accepted)

Submitted to ApJSS

Corresponding author: K. D. Leka  
kdleka@isee.nagoya-u.ac.jp, leka@nwra.com

## ABSTRACT

Solar flares are extremely energetic phenomena in our Solar System. Their impulsive, often drastic radiative increases, in particular at short wavelengths, bring immediate impacts that motivate solar physics and space weather research to understand solar flares to the point of being able to forecast them. As data and algorithms improve dramatically, questions must be asked concerning how well the forecasting performs; crucially, we must ask *how* to rigorously measure performance in order to critically gauge any improvements. Building upon earlier-developed methodology (Barnes et al. 2016, Paper I), international representatives of regional warning centers and research facilities assembled in 2017 at the Institute for Space-Earth Environmental Research, Nagoya University, Japan to – for the first time – directly compare the performance of operational solar flare forecasting methods. Multiple quantitative evaluation metrics are employed, with focus and discussion on evaluation methodologies given the restrictions of operational forecasting. Numerous methods performed consistently above the “no skill” level, although which method scored top marks is decisively a function of flare event definition and the metric used; there was no single winner. Following in this paper series we ask why the performances differ by examining implementation details (Leka et al. 2019, Paper III), and then we present a novel analysis method to evaluate temporal patterns of forecasting errors in (Park et al. 2019, Paper IV). With these works, this team presents a well-defined and robust methodology for evaluating solar flare forecasting methods in both research and operational frameworks, and today’s performance benchmarks against which improvements and new methods may be compared.

*Keywords:* methods: statistical – Sun: flares – Sun: magnetic fields

## 1. INTRODUCTION

Solar flares can be considered the initiating event for many Space Weather phenomena and impacts. The impact of solar flare radiation is almost immediate in the case of sudden ionospheric disturbances, particularly with M- and X-class flares, which disrupt radar and terrestrial communications systems in the sunlit hemisphere. Solar flares are also intimately associated with other pertinent space weather phenomena such as energetic particle storms and coronal mass ejections whose impacts may be delayed relative to flare impacts, but can incur broader effects. Predicting solar flare likelihood has thus long been a defined and required operational product, now with several facilities world-wide providing operational forecasts to a variety of customers.

Predicting solar flares is also the ultimate test of understanding their cause, or causes. They have long been associated with certain morphological aspects of solar active regions such as complex structures, strong-gradient polarity inversion lines and indications of significant energy storage in the magnetic field itself (see *e.g.*, and references cited by Sawyer et al. 1986; Leka & Barnes 2003; Schrijver 2007). The only appropriate energy source is the stored free magnetic energy in solar active region magnetic fields, and the only appropriate release mechanism invokes magnetic reconnection and reconfiguration to release that free magnetic energy. Indeed, as discussed below and further in Leka et al. (2019, Paper III), quantitative “modern” forecasts incorporate this physical understand-

ing as they often characterize coronal magnetic energy storage by proxy, with the parametrizations of photospheric magnetograms. In these contexts, however, pinpointing a unique triggering mechanism has remained elusive. Alternatively, solar flares may inherently be stochastic in nature (see for example [Wheatland 2000](#); [Strugarek et al. 2014](#); [Aschwanden et al. 2016](#)), thus essentially unpredictable in a deterministic sense. The state of the research is presently at a point where it is still unknown in which regime the physics operates. While their heliospheric and societal impacts provides motivation for predicting these energetic events, success or failure at forecasting also provides a key indicator as to whether stochastic physics is or is not involved.

In 2009, the first in a series of workshops was held to compare and evaluate the newly-emerging plethora of methods aimed at distinguishing solar active regions with imminent flare threat. Data from the Solar and Heliospheric Observatory (SoHO; [Domingo et al. 1995](#)) and specifically the Michelson Doppler Imager (MDI; [Scherrer et al. 1995](#)) were provided to the methods for analysis. The performance results (see [Barnes et al. 2016](#)) are of secondary importance to the methodology that was established, identifying the importance of common definitions and standard metrics when determining what constitutes “good performance.”

During Solar Cycle 24, the availability of significantly improved data sources, such as the Helioseismic and Magnetic Imager (HMI) on the Solar Dynamics Observatory (SDO [Pesnell 2008](#); [Scherrer et al. 2012](#); [Schou et al. 2012](#); [Centeno et al. 2014](#); [Hoeksema et al. 2014](#); [Pesnell 2008](#)), has made possible a growing variety of flare forecasting systems that are running in an operational mode (some of which were in development phase in 2009). Consequently, an international collaboration effort was initiated through the Center for International Collaborative Research (CICR), at the Institute for Space-Earth Environmental Research (ISEE), Nagoya University, Japan, to bring together the operational forecasting teams from a variety of institutions (government, private, academic) to evaluate the performance of different techniques. The goals of that workshop and subsequent analysis are to (1) establish benchmarks and comparison methodologies for operational flare-forecasting facilities, and (2) begin to understand what particular forecasting methodologies enable the best forecasting performance.

The participating systems are listed in Section 2 with additional relevant (unpublished) details elaborated upon in Appendix A. Although additional research into improving forecasts is being published frequently of late ([Bobra & Couvidat 2015](#); [Nishizuka et al. 2017](#); [Florios et al. 2018](#)), for this research the comparisons were limited to those truly running in an operational manner, which the group describes as “providing a forecast on a routine, consistent basis using only data available prior to the issuance time.” Many methods, especially the long-standing government-institutional methods, rely on sunspot classification and historical flaring rates ([McIntosh 1990](#); [Sawyer et al. 1986](#)). A few, now, are employing more sophisticated analysis of the host sunspot groups and statistical classifiers or machine-learning algorithms. Forecasts were not required to be fully automatic – human intervention, a “Forecaster In The Loop” (FITL) was explicitly allowed. Providing a forecast on a daily basis was also not a requirement, although as an operational system, not doing so was effectively penalized by the evaluation metrics, as described in Section 2.2. No further restrictions were placed on the data employed or interval used for training, except that it could not overlap with the testing interval (see Section 2.1). The impacts of long- vs short- training intervals (*e.g.* whether more than one solar cycle was used for training the method) and other details are discussed further in Paper III.

The participants provided forecasts for an agreed-upon interval with agreed-upon event definitions as described in Section 2.1. Representatives from most participating groups attended (in person or remotely) a 3-day workshop during which the approaches and initial results were discussed in depth. The results of those days, plus further discussions and analysis which occurred in the subsequent months, are now presented here and in Papers III, IV.

## 2. COMPARISON METHODOLOGY

The participating facilities and methods (with their monikers and published references, as available) are listed in Table 1, and specific details which are not available by published literature (or modifications that have been made since the relevant publications) are briefly described in Appendix A. Some methods have multiple options for producing forecasts, and those are also delineated both in Table 1 and Appendix A. In Paper III we distinguish the methods according to broad categorizations of their implementations, such as data sources, training intervals, imposed limits, forecast approach (*e.g.*, statistical, FITL) *etc.*, and hence we leave such level of detail to that paper.

**Table 1.** Participating Operational Forecasting Methods (Alphabetical by Label Used)

Institution	Name of Method/Code <sup>†</sup>	Label	Symbol	Reference(s)
ESA/SSA A-EFFORT Service	Athens Effective Solar Flare Forecasting	A-EFFORT		Georgoulis & Rust (2007)
Korean Meteorological Administration & Kyung Hee University	Automatic McIntosh-based Occurrence probability of Solar activity	AMOS		Lee et al. (2012)
University of Bradford (UK)	Automated Solar Activity Prediction	ASAP		Colak & Qahwaji (2008, 2009)
Korean Space Weather Center (by SELab, Inc)	Automatic Solar Synoptic Analyzer	ASSA		Hong et al. (2014), Lee et al. (2013)
Bureau of Meteorology (Australia)	FlarecastII	BOM		Steward et al. (2011, 2017)
120-day No-Skill Forecast	Constructed from NOAA event lists	CLIM120		Sharpe & Murray (2017)
NorthWest Research Associates (US)	Discriminant Analysis Flare Forecasting System	DAFFS		Leka et al. (2018)
” ”	GONG+GOES only	DAFFS-G		” ”
NASA/Marshall Space Flight Center (US)	MAG4 (+according to	MAG4W		Falconer et al. (2011);
” ”	magnetogram source	MAG4WF		also see Appendix A
” ”	and flare-history	MAG4VW		
” ”	inclusion)	MAG4VWF		
Trinity College Dublin (Ireland)	SolarMonitor.org Flare Prediction System (FPS)	MCSTAT		Gallagher et al. (2002); Bloomfield et al. (2012)
” ”	FPS with evolutionary history	MCEVOL		McCloskey et al. (2018)
MetOffice (UK)	Met Office Space Weather Operational Center human-edited forecasts	MOSWOC		Murray et al. (2017)
National Institute of Information and Communications Technology (Japan)	NICT-human	NICT		Kubo et al. (2017)
New Jersey Institute of Technology (UK)	NJIT-helicity	NJIT		Park et al. (2010)
NOAA/Space Weather Prediction Center (US)		NOAA		Crown (2012)
Royal Observatory Belgium Regional Warning Center	Solar Influences Data Analysis Center human-generated	SIDC		Berghmans et al. (2005); Devos et al. (2014)

†: if applicable

**Table 2.** 24 hr Event Rates for 2016.01.01 – 2017.12.31

Class	# Quiet Days	# Event Days	Climatology (Event Day Rate)
C1.0+	543	188	0.257
M1.0+	705	26	0.036
X1.0+	728	3	0.004

### 2.1. Event Definitions and Testing Interval

The participants agreed on a testing interval of 01 January 2016 – 31 December 2017 for evaluating forecasts. This is arguably a very short testing interval; in the present situation, it was chosen to balance both training and testing data for those methods relying upon data from *SDO/HMI*, since the near-real-time data from HMI are only available from late 2012. The resulting activity levels are summarized in Table 2. Evaluation was performed on full-disk forecasts only, to avoid the requirement of standardizing the different active-region identification methods in use (combining region-based forecasts to full disk is described in Appendix B.1).

Event definition choices were dictated by the need for common definitions across methods and the fact that these are operational methods, hence most already produce forecasts that match the NOAA/SWPC-established event definition and timings.

As such, the group agreed upon event thresholds as “lower-limits plus exceedance” following the NOAA/SWPC definition, based on the NOAA *Geostationary Observing Earth Satellite* (GOES) X-Ray Sensor (XRS) 1–8 Å bands: C1.0+ and M1.0+ corresponding to lower limits of  $1.0 \times 10^{-6}$  and  $1.0 \times 10^{-5} \text{ Wm}^{-2}$ , respectively, with no upper limit (*i.e.*, “exceedance” forecasts). All forecasts were put onto an exceedance basis; calculating exceedance forecasts from category-limited forecasts (*i.e.* including an upper limit), as were provided by some methods, is discussed in Appendix B.2. No background or pre-flare subtraction was performed for the evaluation data, which is consistent with none generally being performed by any operational method during either training or event prediction (see also Wheatland 2005, for a discussion on the impact of background subtraction.). The event definitions include 24 hr validity periods and effectively 0 hr latencies (the time periods between forecast issuance and the start of the validity period) for the initial comparisons (*i.e.* only “one-day” forecasts, not longer-range forecasts). Longer effective latencies may be implied due to data acquisition times, but these are ignored here for delays  $< 1$  hr. Additionally, it is noted that a number of centers produce additional forecasts (with variations in frequency of forecast, event thresholds, latencies, or validity periods); for this comparison, we chose the event definitions to assure the most overlap between methods. We refer now to these two event definitions using the shorthand “C1.0+/0/24” and “M1.0+/0/24”, noting that the nomenclature includes all three aspects of the event definition (thresholds, latency in hours, and validity period in hours).

The C1.0+/0/24 exceedance definition provided 188 event-days, and the M1.0+/0/24 exceedance definition provided 26 event-days over the 731 days of the testing interval (2016 was a leap-year; see Table 2). Not all methods produce C1.0+/0/24 forecasts. While most methods produce a forecast for X1.0+ ( $1.0 \times 10^{-4} \text{ Wm}^{-2}$  and larger), in practice the short testing interval produced too few of these largest events to provide meaningful evaluations.

Most methods issue a forecast in the neighborhood of midnight Universal Time. Within approximately one hour, any discrepancy from midnight was ignored. Beyond that, the discrepancies in event lists would become problematic. For methods which issue forecasts significantly different from midnight (SIDC at 12:30 UT, NICT at 06:00 UT), custom event lists were constructed based on that issuance time. Although these custom lists do change the number of events slightly (C1.0+/0/24 becomes 183 and 185 event-days for NICT and SIDC respectively; M1.0+/0/24 becomes 27 event-days for both), they provide the most appropriate approach to enable cross-comparisons. Almost all methods issue multiple forecasts throughout the day; in the course of these comparisons the forecast issued closest to midnight Universal Time (UT) was used and others were ignored.

## 2.2. Standard Metrics and Evaluation Tools

Different performance metrics inform on different performance aspects. This is discussed in Jolliffe & Stephenson (2012) and other references specifically with regards to flare forecasting in Bloomfield et al. (2012); Barnes et al. (2016); Kubo et al. (2017); Steward et al. (2017); Murray et al. (2018). Hence, we present a number of metrics and evaluation tools, but for brevity we refer to any of the above references for the definitions of specific metrics<sup>1</sup>.

Graphical representations of performance are used due to the wealth of information available in a compact form. Reliability Plots (also known as Attribute Diagrams) plot bins of the predicted probability against the observed number of instances in that event frequency bin. A perfect reliability displays points along the  $x = y$  line. A perfect forecast is one in which an event is only and always predicted with a probability of 100%; such a service will only have points in the first and last probability bins. Also included in these plots are the climatological rate (event rate) for the testing period (a  $y = \text{constant}$  line at the event rate for that testing period) and the “no skill” line which is defined as the bisector between the testing-interval climatology and the “perfectly reliable”  $x = y$  line. Additionally, we indicate the relative population of the full sample proportion of forecasts within each bin.

Relative (Receiver) Operating Characteristic (Curve) or “ROC” diagrams are constructed by plotting the Probability of Detection (POD) *vs.* the Probability of False Detection (POFD) as a threshold is varied by which a forecast outcome becomes a “yes” forecast. This threshold is commonly referred to as the Probability Threshold  $P_{\text{th}}$  as it is applied to forecast probabilities, but is applied here even though some methods may not strictly produce probabilities. ROC diagrams measure resolution but not reliability. ROC diagrams include the  $x = y$  line to indicate “no skill”; on a ROC plot, perfect forecasts trace the path from (0, 0) to (0, 1) to (1, 1).

Supplementing the graphical evaluation tools are quantitative metrics. Skill score metrics in particular compare performance to that of a reference forecast. These are normalized such that perfect forecasts result in a metric of 1.0, and “no skill” as compared to the reference results in 0.0. The reference forecast may take various forms; commonly used is the climatology of the testing period or a random forecast (Jolliffe & Stephenson 2012), but it may be any other valid forecast method.

The Reliability Plots can be summarized by the Brier Skill Score (BSS), a metric based on the probability forecasts, and for which the reference is specifically the no-skill climatological forecast of

<sup>1</sup> See also [http://www.cawcr.gov.au/projects/verification/#What\\_makes\\_a\\_forecast\\_good](http://www.cawcr.gov.au/projects/verification/#What_makes_a_forecast_good) and <https://www.nssl.noaa.gov/users/brooks/public.html/feda/note/reliroc.html> for broad discussion and numerous definitions

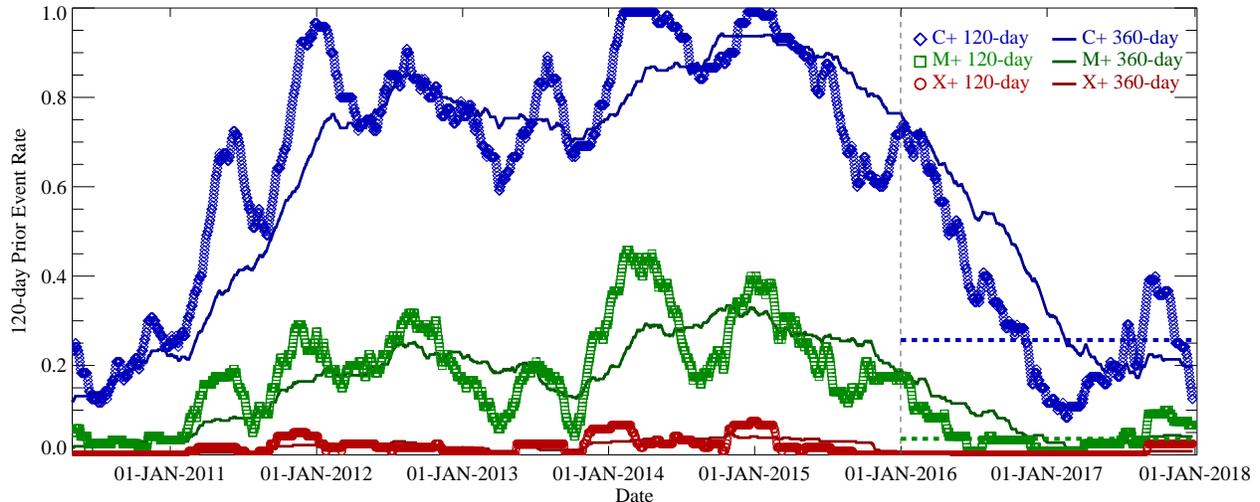
the testing period (see Table 2). This metric answers the question, “how well did this method do compared to the underlying climatology?”.

The ROC curves are summarized here by the ROC Skill Score (ROCSS) also known as the Gini Coefficient, both of which are related to the Area Under the Curve (AUC) but provide more discrimination (Jolliffe & Stephenson 2012; Leka et al. 2018). The ROCSS and Gini coefficient are normalized such that no skill provides a score of 0.0, and perfect forecasts provide a score of 1.0.

Deterministic (or categorical) forecasts can be valuable when preparing forecasts for a particular customer who may require a specified acceptable rate of false alarms, for example, rather than simply a probabilistic forecast. Four additional metrics based on dichotomous (yes/no) forecasts are included: the Appleman Skill Score (ApSS) uses the testing interval to construct an “across the board” climatology reference forecast (a single reference forecast according to the event day rate in the testing interval), the Equitable Threat Score (ETS) invokes a random forecast, and the Hanssen & Kuiper Skill Score / Peirce Skill Score / True Skill Statistic (here just PSS/TSS) is the difference between the POD and the POFD (see definitions and discussions in Woodcock (1976); Murphy (1996); Barnes & Leka (2008); Bloomfield et al. (2012); Barnes et al. (2016); Murray et al. (2017); Kubo et al. (2017)). These metrics are all based on permutations of the “truth table” entries that compare Predicted *vs.* Observed outcomes, and are discussed at length in the references cited above. Additional numeric metrics such as the Proportion Correct (PC, also called Rate Correct or Accuracy) and the Frequency Bias (FB) (Jolliffe & Stephenson 2012) do not compare to reference forecasts per se, and may or may not have a similar normalization as required for a true skill score. The PC metric is common (but can be misleadingly high even for unskilled forecasts in highly unbalanced samples) and the FB indicates systematic over- or under-forecasting, a necessary complement to the TSS metric.

A deterministic forecast is produced by imposing a threshold  $P_{th}$  for assigning the probabilities or forecast outcomes to yes/no forecasts. This threshold reflects a probability level for an event at which a “real-world” action/no-action decision has to be taken based on, for example, economic losses incurred from one or the other type of error. This threshold is then also used for the dichotomous-based metrics (PC, ApSS, ETS, PSS/TSS, FB) by which that method is evaluated. The performance of a method according to a dichotomous-based metric may vary as a function of  $P_{th}$  – this is demonstrated in ROC curves where the vertical distance of each point of the curve from the no-skill  $x = y$  line reflects the PSS/TSS and thus the method’s discrimination between events and non-events as  $P_{th}$  is varied (see the discussion in Barnes et al. 2016). Generally speaking, the methods here are either not explicitly optimized for a particular  $P_{th}$  during their training or the training method implicitly maximizes a particular metric that effectively optimizes the system at  $P_{th} = 0.5$ . All but one method produced probabilistic forecasts; for the one that did not, outputs of 0.0 and 1.0 were assigned “no” and “yes” forecasts, respectively.

Hence, we adopt  $P_{th} = 0.5$  to compute dichotomous-based metrics for all methods. A few methods provide custom forecasts to customers with different  $P_{th}$ , or routinely provide their alerts above a particular  $P_{th}$ , and those were invited for evaluation with a custom  $P_{th}$  (none were submitted). Unless specified otherwise, selecting  $P_{th} = 0.5$  for categorical-based metrics is an allowable choice



**Figure 1.** The “120-day prior climatology” and “360-day prior climatology are plotted for the C1.0+/0/24 and M1.0+/0/24 event definitions, plus the same for an X1.0+ threshold for completeness, from the start of the SDO mission (2010.05.01) through the testing interval, whose start is indicated by a vertical dashed line. The climatological event rate of the testing interval is indicated by horizontal dashed lines over that time period. Each symbol-point (as indicated) represents the daily full-disk event rate for the prior 120 days (up until but not including the date on which the point falls), similarly for the curves indicating the 360-day prior climatology. The 120-day prior climatology is used as the unskilled reference forecast in the “MSESS\_clim” and “ApSS\_clim” metrics in Figure 4.

for all methods. All probabilities for all forecast methods accompany this publication<sup>2</sup> and are thus available for readers to calculate additional metrics, for example with  $P_{th} \neq 0.5$ .

For all methods, missing forecasts were assigned a probability  $p = 0.0$  for that day. This is appropriate for operational forecasts, where missed or skipped forecasts should be penalized. Most operational methods have built in backup sources of data, forecasts, or the ability to forecast prior climatology in the event of, for example, data interruption (see additional details in Paper III).

We do not present the popular “maximum TSS” ( $TSS_{max}$ ) for two reasons. First, an “optimal  $P_{th}$ ” with which  $TSS_{max}$  is calculated should be established based on information obtainable only from the training interval, rather than the testing interval itself, as is common practice. No method supplied such a customized  $P_{th}$  to use. Determining an “optimal  $P_{th}$ ” from which to achieve a maximum TSS score based on testing-period information is not consistent with a purely operational approach. The optimal  $P_{th}$  can have a correspondence to the underlying event rate (Bloomfield et al. 2012; Barnes et al. 2016), which varies according to the solar cycle and from one cycle to the next as discussed below.<sup>3</sup> Hence, there is limited “actionable information” in determining the optimal  $P_{th}$  from a training period for future forecasting. Second, the  $P_{th}$  for each method used to achieve  $TSS_{max}$  will differ from each other and will depend on the event definition, so interpreting these results is challenging (see discussion in Barnes et al. 2016). That being said, one can roughly estimate  $TSS_{max}$

<sup>2</sup> Leka and Park 2019, Harvard Dataverse, doi:10.7910/DVN/HYP740

<sup>3</sup> Some methods (*e.g.* A-EFFORT) do establish optimal  $P_{th}$  levels during training and apply them in order to issue alerts. They elected to not invoke these  $P_{th}$  for the evaluations here.

for each method from the shape of its ROC plot (*i.e.* the point of maximum vertical departure from the no-skill  $x = y$  line).

### 2.3. Highlighted Metrics: Comparison against No-Skill Operational Forecasts

All metrics discussed thus far explicitly evaluate the performance of forecasts against the outcome of the testing interval. In true operational settings, however, an appropriate reference forecast against which to judge performance is more appropriately the best “unskilled” forecast available (Sharpe & Murray 2017; Murray et al. 2017). In other words, for operational forecasting it is appropriate to separately and specifically ask “to what extent is the method in question an improvement beyond what would be otherwise available by simply using an unskilled forecast?” If a forecasting method cannot perform better than this unskilled forecast, then it does not add any skill or value beyond that unskilled forecast.

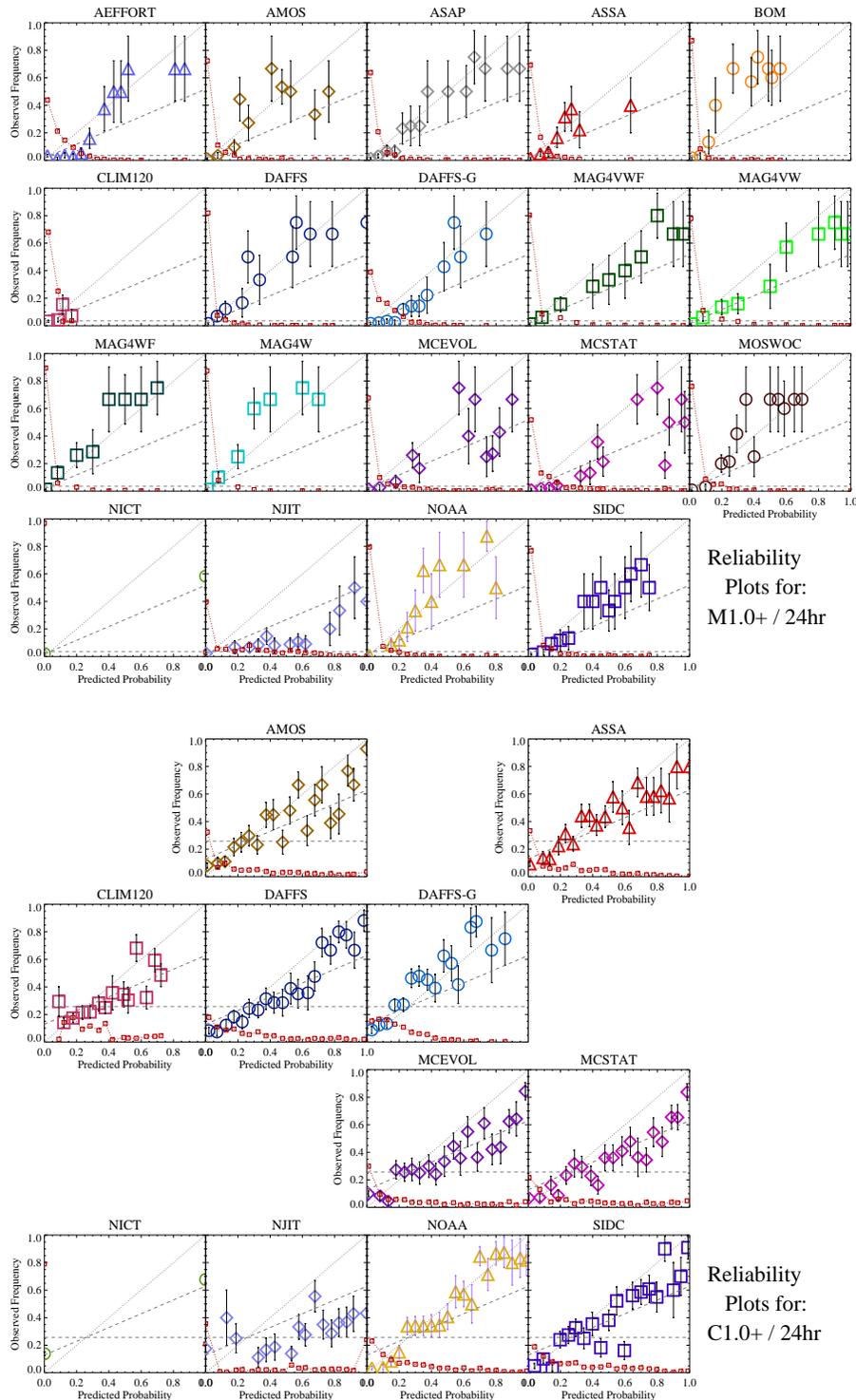
To construct a “no-skill” forecast for day  $t$  for the event definition in question, we use an event rate determined over the prior  $N$  days up to and including  $t - 1$ . The resulting event rate is then used as the reference forecast’s predicted probability for that date  $t$ . We choose  $N = 120$  days as suggested by Sharpe & Murray (2017). This unskilled reference forecast does vary, as shown in Figure 1 – in particular decreasing from  $> 0.5$  to  $< 0.5$  for C1.0+/0/24 within the testing interval. Its abrupt variation on short timescales (*e.g.* around September 2017, see also Sharpe & Murray (2017) Figure 5) likely reflects active-region recurrence patterns and space weather effects rather than reflecting longer-range climatology (see discussion on climatology variations in McCloskey et al. 2018, and the 360-day prior climatology curves also shown here in Figure 1). However, a 120-day prior climatology forecast avoids significant lag against the fairly rapid event-rate changes that occur at the beginning and end of the solar magnetic cycles evident in the 360-day prior climatology curves. Either provides a valid unskilled forecast and a valid reference forecast for associated metrics, with expected performance differences and resulting scores – as would a “no-skill” forecast using yet another value for  $N$ . The 120-day prior climatology forecast (“CLIM120”) is included for evaluation along with all other methods as a “sanity check” on the performance of this reference forecast.

Two metrics are constructed using this unskilled forecast as the reference. A metric “MSESS\_clim” is analogous to the Brier Skill score as based on the mean square error (MSE) of the forecast probabilities. However, instead of the testing-period climatology as defined for the BSS, the MSESS\_clim uses the prior 120-day event rate (“120-day prior climatology”) as the reference forecast. Analogously, we compute an Appleman skill score for which the “across-the-board” forecast for any given day is dictated by this reference; the resulting accuracy is computed and used as the reference forecast in the “ApSS\_clim” score.

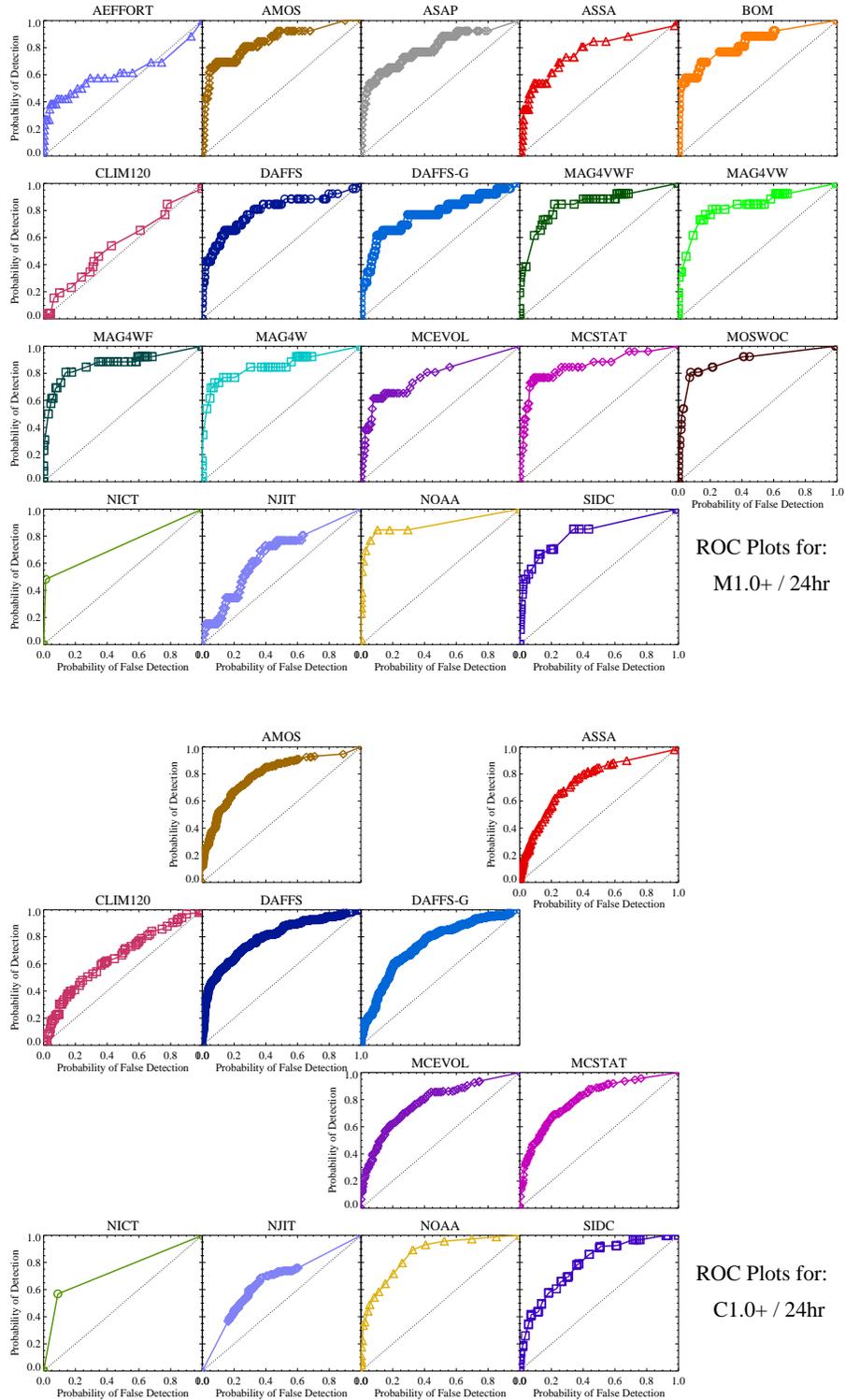
## 3. THE METHOD PERFORMANCES

Results are shown here for the metrics and evaluation methodology described in the previous section. Of note, if a particular method is highlighted in the text as an example of a particular trend it will rarely be the only example, and such a call-out does not mean other methods are exempt from said trend. Such call-outs refer to M1.0+/0/24 results unless otherwise noted.

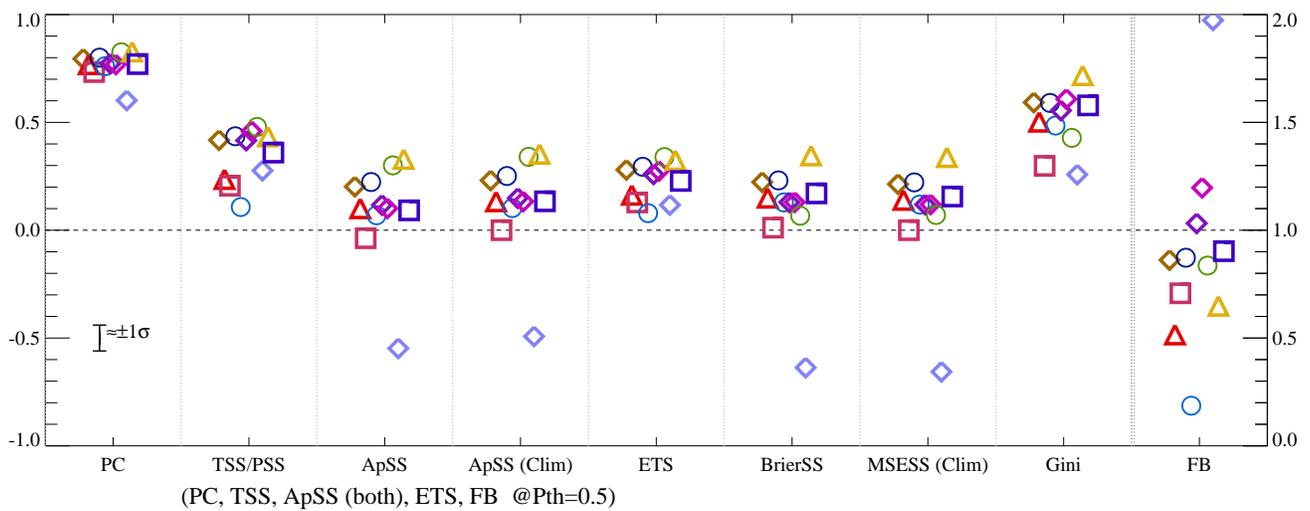
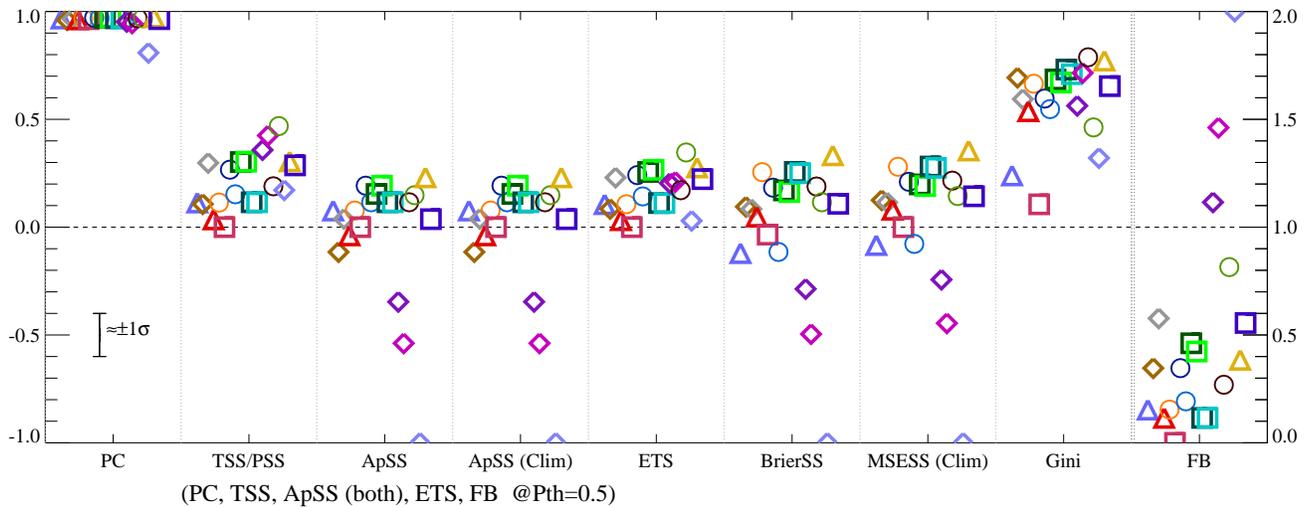
First, in Figure 2 the Reliability Diagrams (Attribute Diagrams) are shown, comparing predicted probabilities to the observed frequencies across 20 probability bins. The predicted probabilities are indicated on the  $x$ -axis by the average of the probabilities in that bin. Points in each bin are accumulated, and thus accurately reflect the distribution whether from continuous probabilities or



**Figure 2.** Reliability plots (Attribute Diagrams) for each method, indicating the performance of the probabilistic forecasts as named, the  $x = y$  “perfect reliability” line, the (horizontal) climatology level (dashed line), and the “no skill” line (dotted line) that lies between the two. Additionally (red dotted line and small square) is the fraction of the total sample for which a forecast exists for each bin. Each method has an assigned color / symbol combination (Table 1), where related methods (*e.g.* from the same institution) have the same symbols and are plotted with colors in the same family (“nearby” in hue). Results are shown for M1.0+/0/24 (top) and C1.0+/0/24 (bottom); fewer methods predict the latter than the former. Results were not calculated for X1.0+ due to extremely small number of events in the testing interval.



**Figure 3.** Relative Operating Characteristic plots with the  $x = y$  “no skill” line, following the color/symbol scheme of Figure 2.



**Figure 4.** Results from the direct comparison of flare forecasting methods for a variety of performance metrics (Left to Right): the Proportion Correct, the True Skill Statistic / Peirce Skill Score, Appleman Skill Score (testing period), the Appleman Skill Score with the 120-day prior climatology reference forecast, the Equitable Threat Score, Brier Skill Score, a Mean Square Error Skill Score with the 120-day prior climatology as a reference forecast, and the Gini Coefficient. A lower limit of  $-1.0$  was imposed for the plotting. The final metric is the Frequency Bias, whose displayed range is indicated on the right-hand axis; a  $+2.0$  limit was placed on this plot. Metrics based on “truth tables” are calculated using  $P_{th} = 0.5$ ; the Brier Skill Score and Mean Square Error Skill Score (clim), and Gini Coefficient are independent of  $P_{th}$ . The symbols follow the scheme in Figures 2, 3 and are offset slightly in the  $x$ -dir for clarity in the same order as they appear in Figures 2, 3, and Table 1). Results are shown for M1.0+/ $0/24$  (top) and C1.0+/ $0/24$  (bottom); fewer methods predict the latter than the former. Results were not calculated for X1.0+ due to the extremely small number of events in the testing interval.

discrete forecast probabilities. This figure also displays the symbol and color schemes devised to both compare methods and inter-compare between related methods (*e.g.* variations from the same institution, *c.f.* Table 1). Most methods provided a form of M1.0+/0/24 forecasts (natively, or computed as per Appendix B.2). A subset of methods also produce forecasts for C1.0+/0/24 and those are displayed as well. The decision regarding whether to produce forecasts for these smaller flares rests on the facility or agency according to resources, customer needs, and perceived threat; if publicly available, these forecasts were included. Most methods do provide a forecast for X1.0+, however the number of events was so small during the testing period as to be uninformative (see Table 2). The error bars are determined by the number of points and events in each bin (Wheatland 2005; Barnes et al. 2016): for a reliability value  $R$  in a particular bin,  $\sigma_R = (R((1 - R)/N_{\text{bin}} + 3))^{1/2}$  with  $N_{\text{bin}}$  being the number of points in that bin.

The Reliability Diagrams graphically display trends of over-forecasting (*cf.* MCSTAT) or under-forecasting (*cf.* MAG4W) the issued probabilities. Some methods more systematically perform errors of one type (*e.g.* BOM) while others display a mix according to the probability bin (*e.g.* AMOS) but not an obvious dominance of one error or the other. The reliability plots also highlight that some probabilistic methods provide predictions covering the full range of probabilities (*e.g.* MAG4VWF) while others do not provide predictions at the highest probabilities (*e.g.* ASSA). The case of NICT, as the sole fully deterministic forecast, appears different due to the assignment of probabilities (see Kubo et al. (2017) for more on evaluation methods for fully deterministic forecasting). This lack of high probability forecasting is more pronounced for larger event-magnitude thresholds (*e.g.* more prevalent here for M1.0+/0/24 forecasts as compared to C1.0+/0/24 forecasts), a trend noted in Barnes et al. (2016). Most of the methods here are probabilistic with the exception of the NICT facility which produces deterministic forecasts. Larger flares are less frequent, and probability-based forecasts will train to reflect that fact, which reduces the presence of high-probability forecast values.

The ROC curves for all methods are presented in Figure 3, using the same color and symbol scheme. The  $x = y$  line indicates no ability to discriminate between the two forecast outcomes (forecast for or against an event in the present case). The points on the ROC curve are computed for each distinct probability presented by a method. Hence, methods which provide forecasts in discrete probability bins present with fewer points than those which provide continuous-probability forecasts (*cf.* NICT *vs.* DAFFS). We see a slight increase in the ability of the models that provided forecasts for both event definitions to discriminate for the M1.0+/0/24 results as compared to C1.0+/0/24. This is a generally observed trend (Leka et al. 2018; Murray et al. 2017).

Comparing the Reliability *vs.* the ROC plots for a particular method highlights the different information presented by each. As an example, the MAG4 results using line-of-sight magnetograms (MAG4W, MAG4WF) *vs.* those using vector magnetograms (MAG4VW, MAG4VWF) appear to show very similar ROC plots while displaying systematically different behavior in the Reliability plots (even with different training particulars with regards to longitudinal limitations). Also of interest are the comparative performances of methods which are ostensibly based on the same basic approaches (*e.g.* Poisson statistics applied to historical region flaring rates *e.g.*, MCSTAT *vs.* ASSA – or those with human forecasters involved – *e.g.*, NOAA *vs.* MOSWOC).

Figure 4 shows the variety of skill scores and quantitative metrics described in Section 2.2, with approximate  $1\sigma$  error bars also indicated. There is no straightforward way to estimate uncertainties on the metrics, given the operational approach (*e.g.*, data for a bootstrap evaluation are not gener-

ally available). However, we estimate the uncertainties in two ways. First, there are other studies which have employed bootstrap or similar methods to calculate the uncertainties in skill scores (e.g. Bobra & Couvidat 2015; Leka et al. 2018), although the underlying event populations are somewhat different. Adjusting for the smaller sample sizes here, one can estimate a general level of uncertainty in the skill metrics of  $\approx 0.06$  for C1.0+/0/24, and  $\approx 0.10$  for M1.0+/0/24. To supplement this estimate, the DAFFS facility (specifically the magnetic field parameter component) was re-run for the testing interval (2016.01.01 – 2017.12.31) using 100-draw (with replacement) bootstrap analysis. Across numerous metrics and variables available in DAFFS, we find the uncertainties range over 0.04 – 0.09 for C1.0+/0/24 and over 0.05 – 0.17 for M1.0+/0/24, with the ranges due to whether 1- or 2-variables were tested and the particular metrics used. These estimates are only guidance and do not necessarily reflect the full uncertainty situation. These uncertainties are also likely to be underestimates, because they only account for the random error and no separate bias is calculated for the error estimate itself. For example when using the full-disk bootstrap, individual days are drawn rather than full disk-passages of individual active regions. Additionally, given the change in event rate between training and testing intervals, there is likely to be a significant bias present for most methods.

The answer to the question of which methods perform “best” depends on event definition and the metric under consideration. The rank order of performance changes between metrics and between event definitions. This is demonstrated poignantly by MCSTAT/MCEVOL which score near bottom rank for ApSS, but near top rank for TSS/PSS for the M1.0+/0/24 tests.

Some metrics can differentiate performance better than others in these applications. The Proportion Correct metric for M1.0+/0/24 is uninformative in trying to differentiate between methods due to the large percentage of correct negatives, however it provides some information for the C1.0+/0/24 analysis. Because the climatology rate does not vary across the 0.5 threshold for M1.0+/0/24, the two Appleman scores (ApSS and ApSS\_clim) are identical in this case. In the case of the C1.0+/0/24 event definition, the climatology rate does vary across the 0.5 threshold, and the results for the two scores are slightly different.

That being said, the majority of methods perform similarly to each other – that is, their scores are consistent with each other across metrics. This is particularly the case for the M1.0+/0/24 tests given the estimated uncertainties, although there are arguably performance differences beyond the uncertainties for the C1.0+/0/24 test.

Comparing the Reliability plots (based on probabilities) to the Frequency Bias (which is a dichotomous-based metric employing a single  $P_{th} = 0.5$ ) it appears that the vast majority of methods tend toward underforecasting for larger-flare M1.0+/0/24 tests by varying degrees ( $FB < 1.0$ ), with a less pronounced deviation from  $FB = 1.0$  for most methods that underwent the C1.0+/0/24 tests. As mentioned above, the FB score ‘checks’ the TSS, in that for low event rates such as typical for solar flares, an over-forecasting system can attain a high TSS while an under-forecasting system is less likely to – so comparing TSS scores should only be performed in the context of an accompanying FB score. As such, for example, confidence in the TSS scores for MCSTAT for the M1.0+/0/24 test should be tempered somewhat, while the NICT TSS result is more robust.

Different implementations of otherwise the same method can be differentiated and the hoped-for “improvements” confirmed (or not). The implementations using vector magnetic field data do perform better (albeit only slightly by most metrics) than implementations using  $B_{los}$  data within

the same general method (*e.g.* MAG4W\* *vs.* MAG4V\*, DAFFS *vs.* DAFFS-G). By most metrics, MCEVOL’s addition of an evolutionary component to MCSTAT does improve performance, although notably not in the Gini (as visible by the shape of the ROC curve). However, the inclusion of prior flaring history makes almost no difference in performance across the MAG\* method (*e.g.* MAG4W *vs.* MAG4WF, MAG4VW *vs.* MAG4VWF).

None of the operational methods are exceptionally good (*i.e.*, close to 1.0 on any metric, except Gini and Proportion Correct), although the majority consistently score above “no skill” for the metrics considered here. Three methods demonstrate arguably poor performance specifically for the metrics that refer to climatology; these three also show  $FB > 1.0$  (over-forecasting). The case of NJIT is fairly well understood and discussed below, while the others will be discussed further in Paper III (Leka et al. 2019).

#### 4. DISCUSSION

In this study we demonstrate two things: first, a methodology to provide meaningful head-to-head comparisons, and second, the present state of operational flare forecasting. With this first direct comparison of forecast methods, benchmarks of performance by a variety of measures are now provided against which future developments can be tested – an important element of measuring progress in space weather prediction capability.

Regarding the methodology, all forecasting facilities are placed on a level evaluation platform with respect to the full event definition (including thresholds, validity periods and latencies). Those whose forecasting time differed significantly were afforded custom event lists for evaluation, and those producing both upper- and lower- threshold-limited forecasts were converted to exceedance forecasts to match other methods. Full-disk forecasts ensured that differences in defining “solar active regions” would not impede the comparisons. The time period chosen was not ideal – too short with arguably a very small event list – but in the face of new data sources and a very quiet solar cycle, it was an acceptable and necessary compromise. Most important was how the time period was chosen – a period that was common to all methods which also afforded those methods relying on SDO/HMI data an adequate training interval.

The second component of the methodology is the choice of evaluation metrics, and this is arguably a challenge in the context of a direct comparison because it is crucial to ensure that the metrics are all fair (or equally unfair) to all methods. For the presentation here, we select a representative array of dichotomous-based and probability-based metrics, with accompanying graphical evaluation tools, to try and provide as complete a picture as possible. As discussed in Barnes et al. (2016) and elsewhere, applying dichotomous-based metrics to probabilistic-based forecasts require thresholds to be set which may or may not be ideal for a particular method, resulting in unfair penalties. In operational practice, it is challenging to choose the threshold that would ensure optimum performance (by measure of various dichotomous-based metrics) at the time of forecast issuance. As discussed in Bloomfield et al. (2012); Barnes et al. (2016), an optimum threshold for TSS/PSS is usually close to the climatological event rate – which is itself found only after long-term averages are taken in the testing period. Such information is not available at the time of forecast issuance, and may not be optimal for a different metric. For the evaluations here we encouraged methods to submit deterministic forecasts or submit probabilistic forecasts and specify thresholds that may have been used to produce customized deterministic forecasts for particular customers or needs (such as an acceptable error rate of one type or the other). None chose to provide other thresholds and thus

$P_{\text{th}} = 0.5$  was applied to all. As such, we examine how well the methods perform in a deterministic sense if action is only taken when an event is forecast with a probability 50% or higher.

We make note of metrics which are appropriate specifically for evaluating operational systems, since they specifically query what value the system brings above an available unskilled forecast. The Appleman and Brier skill scores by definition employ reference forecasts based on the climatology of the testing period but, as discussed, this information is not actionable for improved future performance. We promote evaluations against an unskilled forecast. Here we provide analogous Mean-Square-Error Skill Score and an Appleman Skill Score that employ a 120-day prior climatology as the reference unskilled forecast (as described in Sharpe & Murray 2017), although others may obviously be used. For the testing period herein, the results did not differ substantially from the original version of the metrics. However, the question asked differs in a distinct way and these metrics are highlighted as part of this work’s focus on methodology.

There was not universal agreement in this group regarding evaluation philosophy, specifically with regards to utilizing dichotomous metrics for probabilistic forecasts. The discussion centers on performance variation as a function of assigned  $P_{\text{th}}$  in the context of an operational system. While a system may be trained to optimize a particular metric and  $P_{\text{th}}$ , there is no guarantee the performance will be the same with that  $P_{\text{th}}$  during the testing interval; evaluating a method using a new optimal  $P_{\text{th}}$  from the testing interval mis-represents the performance when the information needed to assign an optimal  $P_{\text{th}}$  is unknown at the time of the forecast. One approach for evaluating probabilistic forecasts is to only employ graphical methods such as the Reliability Plots and ROC curves and apply metrics such as the Brier Skill Score and ROCSS (Gini score) for which no  $P_{\text{th}}$  is required; this approach is fair (except to the inherently deterministic method(s)) but dismisses some metrics that the community find informative and popular. A second approach is to present all dichotomous metrics in a manner similar to ROC curves, displaying their outcomes as  $P_{\text{th}}$  is varied and reporting the maximum attained score (with its associated  $P_{\text{th}}$ ); but this approach can imply performance better than is attainable in an operational setting and is unlikely to provide guidance for improvement. Hence, the group recognizes that the primary reason for setting a particular  $P_{\text{th}}$  to apply to probabilistic forecasts is to define a threshold upon which action should be considered according to a particular customer’s cost/benefit analysis and resilience against forecasting errors. The full forecast data and evaluation tools used in the present analysis accompany this publication<sup>4</sup> so that additional metrics using, for example, a different  $P_{\text{th}}$ , may be calculated by the interested reader.

Regarding the results, generally speaking no method is working extraordinarily well; although we demonstrate that a fair number of methods consistently perform better than various “no skill” measures, meaning that they do show definitive skill across more than one metric. No method scores above 0.5 (*i.e.* halfway between “no skill” and “perfect”) across all evaluation metrics, and for a number of metrics *no* method provides results above 0.5. The specific ordering of performance varies according to metric and event definition: *there is no single “best” method*, especially given the estimated uncertainties in the metrics. Amongst methods which provide different versions, the versions generally behave similarly in some of the gross characteristics (*e.g.*, shapes and sampling for the ROC curves) with subtle offsets reflecting the refinements made between each.

<sup>4</sup> Leka and Park 2019, Harvard Dataverse, doi:10.7910/DVN/HYP740

Three particular impacts on forecast method success are worth noting. First, the underlying event rate obviously varies within the solar cycle (Fig. 1), and possibly across solar cycles (McCloskey et al. 2018). This will impact the forecasting methods, although the degree of impact will vary depending on training methodology. One example would be that if a method is trained to have high reliability during a time of high solar activity, it may then systematically over-forecast during times of declining or lower solar activity. Alternatively a method may not in fact be particularly reliable during training, but when faced with a particular epoch of the solar cycle (*e.g.* such as the declining phase with more isolated sunspot groups), it may perform better.

Second, there are always flares which occur that are not assigned to any particular active region, or occur behind the visible limb and may be assigned to a region *post facto*. During the testing period, there were 41 unassigned C1.0--C9.9-flares and 3 unassigned M1.0--M9.9-flares, in some cases such unassigned flares were the sole cause of an “event day” (this is discussed further in Paper IV). Unassigned regions have consequences for training operational systems as well as for evaluating and testing them. The vast majority of methods train on individual regions, and in doing so, they will then underforecast systematically for full-disk forecasts. All region-based forecasting methods will miss days where events are produced by no assigned or detected region.

Third, we can highlight here a distinct case of the impact arising from the lack of a full transition to operational functionality. The NJIT method arguably employs one of the more sophisticated analyses of magnetic field data and shows distinct skill in the TSS and Gini metrics. However, it arguably performs the worst according to other metrics. Of all the methods, the NJIT system most reflects the “research” stage of flare forecasting. It was implemented without calibration across a change in instrumentation between training and testing intervals, which in this case (given the analysis method) could easily cause the systematic over-forecasting as evidenced by the metrics. This is an issue faced by many methods in light of aging or changing data sources and the assumed advantage of longer training sets (see Paper III for additional discussion on that point). Additionally, no provisions were made for issuing forecasts in the event of missing or delayed data, and this severely impacted the metrics in a negative manner. Research methods often report encouraging results, but these must be interpreted in the appropriate context. In parallel, the challenge and effort required to bring research into a fully operational mode to the point that it is ready to undergo evaluation in an operational context must not be underestimated.

From this presentation it is not possible to further determine why performances differ. Established methods on which national warning centers rely (*e.g.*, NICT *vs.* NOAA) display very different characteristics in the Reliability and ROC plots, but track fairly well amongst the evaluation metrics. Newer methods show both improvements and degradation against established ones (*e.g.*, MCEVOL and DAFFS *vs.* MOSWOC and SIDC). However, these differences are fairly subtle (that is, within uncertainties) when examined across all evaluation metrics.

We delve further into the “why” question of performance differences in Paper III (Leka et al. 2019) by examining the impact of six distinct categories of implementation differences, finding performance advantages to including prior flare information and a human forecaster, and performance disadvantages to restricting forecast-relevant data to disk-center observations. We use a novel analysis method to evaluate temporal patterns of forecasting errors of both types (*i.e.*, misses and false alarms) for Paper IV (Park et al. 2019), finding weak support for a hypothesis that including temporal infor-

mation such as active region evolution improves a method’s ability to successfully forecast, *e.g.*, a region’s first flare.

The obvious conclusions from this work are actually broad challenges: new forecasting methods, whether empirical or physics based, need to be evaluated *against these established benchmarks* with the goals of improved characteristics in Reliability and ROC plots, and metrics (specifically TSS, ApSS, ETS and BrierSS) all consistently measuring above 0.5 across the full range of event definitions.

We wish to acknowledge funding from the Institute for Space-Earth Environmental Research, Nagoya University for supporting the workshop and its participants. We would also like to acknowledge the “big picture” perspective brought by Dr. M. Leila Mays during her participation in the workshop. S.-H.P. gratefully acknowledges Dr. Ju Jing for maintaining the NJIT flare forecasting system and providing the archive forecasts. KDL and GB acknowledge that the DAFFS and DAFFS-G tools were developed under NOAA SBIR contracts WC-133R-13-CN-0079 (Phase-I) and WC-133R-14-CN-0103 (Phase-II) with additional support from Lockheed-Martin Space Systems contract #4103056734 for Solar-B FPP Phase E support. A.E.McC. was supported by an Irish Research Council Government of Ireland Postgraduate Scholarship. D.S.B. and M.K.G. were supported by the European Union Horizon 2020 research and innovation programme under grant agreement No. 640216 (FLARECAST project; <http://flarecast.eu>). MKG also acknowledges research performed under the A-EFFort project and subsequent service implementation, supported under ESA Contract number 4000111994/14/D/ MPR. S. A. M. is supported by the Irish Research Council Postdoctoral Fellowship Programme and the US Air Force Office of Scientific Research award FA9550-17-1-039. The operational Space Weather services of ROB/SIDC are partially funded through the STCE, a collaborative framework funded by the Belgian Science Policy Office. The authors thank the referees for their constructive comments.

*Facilities:* SDO(HMI), GONG, GOES(XRS)

## APPENDIX

### A. OPERATIONAL FORECASTING METHODS: ADDITIONAL DETAILS

Here we list the methods involved in the comparisons. Pertinent details are provided beyond the descriptions provided in the references listed in Table 1; all times here are quoted in UT. For additional details we also suggest referring to [Leka et al. \(2019, Paper III\)](#), where performance is compared according to specific distinctions.

#### A.1. A-EFFORT (*Academy of Athens*)

A-EFFORT is a Space Situational Awareness (SSA) service of the European Space Agency (ESA), available at <http://a-effort.academyofathens.gr> (with registration). Forecasts are issued at about 00:00 UT and refresh every three hours. Four exceedance thresholds are used: M1.0+, M5.0+, X1.0+ and X5.0+, with a fixed forecast window of 24 hr and 0 hr latency.

There is a single parameter computed from magnetic field data, namely the effective connected magnetic field strength (“Beff” [Georgoulis & Rust 2007](#)) whose values are translated into probabilities using elements of Bayesian analysis and Laplace’s rule of succession. Beff is calculated directly up

to central meridian distances of  $\pm 50^\circ$ ; from this limit to  $\pm 70^\circ$  a magnetic flux-based proxy of  $B_{\text{eff}}$  is calculated to avoid the impact of severe projection effects.

Each of the four forecasts is computed for each of the active regions present within a solar meridional zone of  $\pm 70^\circ$ , identified using a custom active region identification algorithm (see [LaBonte et al. 2007](#)); full-disk probabilities are computed as per Eqn. B1.

#### A.2. *AMOS (Korean Meteorological Administration and Kyung Hee University)*

The Automatic McIntosh-based Occurrence probability of Solar activity (AMOS) model provides daily occurrence probabilities separately for C, M, and X-class flares for each NOAA active region and full disk using McIntosh sunspot group classes and the daily change in area for the sunspot groups. The details are well described in [Lee et al. \(2012\)](#).

#### A.3. *ASAP (U. Bradford, UK)*

Described in [Colak & Qahwaji \(2008, 2009\)](#), ASAP also participated in the ‘‘All Clear’’ workshop in 2009 ([Barnes et al. 2016](#)).

#### A.4. *ASSA (Korean Space Weather Center)*

The Automatic Solar Synoptic Analyzer (ASSA) system at the Korean Space Weather Center identifies and predicts for a variety of solar activity, including sunspot groups and associated flaring. Flare forecast results are issued hourly at :00, with a McIntosh-class-based forecast extending for 24h (used here, initiated in late 2013) and a new ‘‘parameter-based’’ forecast using six major parameters extending for 12h. The McIntosh-class-based forecast uses an independent ASSA algorithm (not NOAA determinations) to identify sunspot groups and determines their McIntosh class by estimating their morphological characteristics, and produces an independent flaring probability according to the ASSA sunspot-flare archive (not based on otherwise published rates). The ASSA sunspot-flare archive was produced based on statistical matching between ASSA’s sunspot group catalog and NOAA’s GOES Soft X-ray events catalog during 1996–2013. A parameter-based method was initiated in late 2016, and provides flare forecasts based on multi-component linear regression using parameters such as the number of sunspots in a sunspot group, the total area of sunspots in a group, and the group’s longitudinal extent. Unfortunately, forecasts from this second method were not submitted. ASSA forecasts rely on SDO/HMI continuum and line-of-sight magnetogram images with no correction for limb-ward effects. Additional details may be found in the user manual ([Lee et al. 2013](#)).

#### A.5. *BOM (Flarecast, Bureau of Meteorology, Australia)*

The details of the probabilistic model are well described in [Steward et al. \(2011, 2017\)](#). Flarecast II (not yet published but results are submitted here) uses the SDO HMI magnetogram imagery analysis capability developed for the original Flarecast model ([Steward et al. 2017](#)) plus prior flaring history, and adds a machine learning technique (logistic regression) to generate a probabilistic forecast. Variables that describe HMI  $B_{\text{los}}$  magnetograms are selected to minimize Aikake’s Information Criteria (AIC), and logistic regression is used to estimate the coefficients of the model, and then used to generate M+, X+, region and full-disk, probabilistic and categorical deterministic forecasts output for flaring activity over the next 24 hours. In operational mode the predictions are updated at 00, 06, 12, 18 UT.

### A.6. *DAFFS, DAFFS-G (Discriminant Analysis Flare Forecasting System, NorthWest Research Associates (NWRA))*

DAFFS is well described in [Leka et al. \(2018\)](#), but it should be noted that it is a fairly young, recently-released system. Of note, being the only method to primarily rely on quantitative analysis of vector magnetic field data from a non-operational data source (*SDO/HMI*), this method suffered from data problems arising from the data-acquisition mode change that incurred a temporary data mis-alignment<sup>5</sup> (MAG4V\* methods use *SDO/HMI* data in a more limited fashion, see below). The impacted data spanned April 2016 – September 2017, and was most damaging for data away from disk center. (The “definitive” data have subsequently been re-processed; the “near real time” data will not be). We noted that it most dramatically impacted some parameters in top-performing combinations, but not others. For the results here, we modified DAFFS to run using parameter combinations that performed essentially identically (within the metric error bars) in the training phase but were not as susceptible to the HMI vector data problem: specifically for the C1.0+/0/24 event definition the parameter combination was changed to  $[E_e, \log(\mathcal{R}_{\text{nwra}})]$  and the M1.0+/0/24 event definition parameter pair was changed to  $[\text{FL}_{24}, \log(\mathcal{R}_{\text{nwra}})]$  from what is described in [Leka et al. \(2018\)](#).

The DAFFS-G ( tool runs simultaneously and is based primarily on GONG  $B_{\text{los}}$  data and persistence (NOAA near real time (NRT) event reports). DAFFS-G is a very “young” release, and has not yet been fully optimized for performance. For the forecasts submitted here, the parameter combinations were  $[\nabla(B_z^{\text{pot}}), \Phi_{\text{tot}}^{\text{pot}}]$  for C1.0+/0/24, and the parameters for M1.0+/0/24 were  $[\sigma(\nabla(B_h^{\text{pot}})), \Phi_{\text{tot}}^{\text{pot}}]$ , where the “pot” moniker refers to the potential field calculated from the  $B_{\text{los}}$  data ([Leka et al. 2017](#)).

### A.7. *MAG4\* (NASA/Marshall Space Flight Center)*

MAG4 is described in ([Falconer et al. 2011, 2014](#)). This study included four versions:

- MAG4W: Free-energy Proxy Only using Line-of-Sight Magnetogram
- MAG4WF: Free-energy Proxy and Previous Flare History using Line-of-Sight Magnetograms
- MAG4VW: Free-energy Proxy Only using Deprojected HMI Vector Magnetogram
- MAG4VWF Free-energy Proxy and Previous Flare History using Deprojected HMI Vector Magnetograms

MAG4W[F] uses the HMI NRT  $B_{\text{los}}$  data with no further correction. The MAG4VW and MAG4VWF, like DAFFS, use *SDO/HMI* vector magnetic field data, however only to 30° from disk center, which were minimally impacted by the data misalignment. In MAG4\*F, previous flare information is used, although a region is assumed to be non-flaring if that information is not available.

### A.8. *MCSTAT, MCEVOL (MaxMillenium Flare Prediction System)*

The MCSTAT approach is well described in [Gallagher et al. \(2002\)](#); [Bloomfield et al. \(2012\)](#) while the MCEVOL approach is well described in [McCloskey et al. \(2018\)](#).

<sup>5</sup> see <http://hmi.stanford.edu/hminuggets/?p=1596> and the *SolarNews* note of 01 September 2017 at <https://solarnews.nso.edu/2017.html#20170901>.

### A.9. *MetOffice (UK) MOSWOC*

The details are well described by [Murray et al. \(2017\)](#). Of note, the forecast closest to 00:00 was used, but is not necessarily the official forecast for that day from MOSWOC, as updates are applied through the (local) night.

### A.10. *NICT (National Institute of Information and Communications Technology, Japan)*

The details of this long-running system are well described in [Kubo et al. \(2017\)](#). Unique to the methods, the NICT-human approach provides four categorical deterministic forecasts of maximum flare size: “Quiet” (max: A/B-class), “Eruptive” (max: C-class), “Active” (max: M-class) or “Major Flare” (max: X-class). These were converted to probabilities of [0.0, 1.0] for the probabilistic-based analysis and converted to exceedance forecasts.

### A.11. *NJIT (New Jersey Institute of Technology)*

The basic methodology is described in [Park et al. \(2010\)](#). The NJIT method is operational in the sense it produces forecasts automatically, but has not been developed further since 2010. It provides probabilistic forecasts of at least one C-, M- and X-class flare occurrence only for a given NOAA-numbered active region within  $\pm 60^\circ$  of disk center; these were converted to exceedance forecasts. The method was trained on 300 primarily flare-productive active regions using SOHO/MDI line-of-sight active-region magnetic field data in solar cycle 23. However, the forecasts now use HMI line-of-sight data without any cross-calibration between the two data sources.

### A.12. *NOAA (Space Weather Prediction Center, US National Oceanic and Atmospheric Administration NOAA)*

The forecasts by NOAA/SWPC have long been considered a standard ([Crown 2012](#)) and have set the benchmarks against which methods are measured using the NOAA/SWPC event definitions (see commentary on this in [Leka & Barnes \(2017\)](#)). SWPC forecasters begin with a climatological basis according to an active region’s classification (SWPC’s assignment of active region class is also considered “The Standard”), according to the historical flaring rates of different sunspot region classes ([McIntosh 1990](#)). From this, a forecaster may modify a region’s probability according to region evolution, flaring trends, and forecaster experience and expertise. These region probability forecasts are combined for a full-disk forecast which itself may be modified based on flaring history of recently-rotated-off regions, or indications of a highly active region about to return. Forecasters may also incorporate other model data when available. Initial forecasts are issued at 22:00 (the “Geophysical Activity Report and Forecast” or RSGA) valid beginning at 00:00 the next day. These are incorporated into the “3-day Forecast” issued at 00:30, with a minimal but not zero probability of a forecast update in the intervening 2.5 hr. Forecasts can, but are not likely to, be updated again before the next 3-day forecast is issued at 12:30. The data used in this comparison arise from the 3-day forecasts but include the C1.0+/0/24 forecasts that are not generally published.

### A.13. *SIDC (Solar Influence Data Analysis Centre of the Royal Observatory of Belgium)*

The forecaster on duty at the SIDC produces each day (nominal issue time 12:30UT) a probabilistic forecast for the occurrence of X-ray flares over the next 24h. Probabilities are provided for flare classes C-, M- and X- separately. A full disk as well as an active region specific forecast is provided. The forecasters use various data sources, the main one being the flaring probability from active regions

with the same McIntosh classification. Such probability is then modulated using for example: the specific flare histories for the regions to be forecasted, *SDO/HMI* magnetogram movies, *SDO/AIA* movies, and *STEREO/EUVI* movies *e.g.* to assess the flaring activity of active regions rotating onto or off the solar disk. Details on flare forecasting at ROB/SIDC and its validation procedures are provided in Berghmans et al. (2005); Devos et al. (2014).

## B. STEPS TO PRODUCE FULL-DISK EXCEEDANCE FORECASTS.

### B.1. Full-Disk Forecasts from Region-Based Forecasts

The forecasts considered here are “full-disk” forecasts, meaning essentially treating the Sun as a star. In practice, only one method did not produce full-disk forecasts, meaning that they only provided forecasts for active regions individually. In that case, the region probabilities were combined according to,

$$P_{\text{FD}} = 1.0 - \prod_{\text{AR}}(1.0 - P_{\text{AR}}) \quad (\text{B1})$$

where  $P_{\text{AR}}$  is the probability of an event for each active region, and the product is performed over all active regions for which such a probability is provided. This equation is effectively how all region-forecasting methods produce their baseline full-disk forecasts.

### B.2. Class-Specific vs. Exceedance Forecasts

The results from methods producing class-specific forecasts (*e.g.* M1.0 - M9.9) were converted to exceedance forecasts (*e.g.* M1.0+ with no upper limit) using conditional probabilities over that method’s training interval, by the following methodology. Suppose one has the probabilities of occurrence of at least one C-, M- and X- class flares respectively, for a given forecast time window  $\tau$  denoted by  $P(\text{C})$  for C1.0--C9.9,  $P(\text{M})$  for M1.0--M9.9 and  $P(\geq \text{X1})=P(\text{X})$  for X1.0+. Then, the lower-bound only probabilities of  $P(\geq \text{C})$  and  $P(\geq \text{M})$  can be determined by combining the probabilities of  $P(\text{C})$ ,  $P(\text{M})$  and  $P(\geq \text{X1})$  with their associated conditional probabilities.

The probability of occurrence of at least one flare at the level greater than or equal to M1.0 during  $\tau$ , i.e.,  $P(\geq \text{M1})$ , can be derived as follows,

$$\begin{aligned} P(\geq \text{M1}) &= P(\text{M}) + P(\text{X}) - P(\text{M and X}) \\ &= P(\text{M}) + P(\text{X}) - P(\text{M}) \times P(\text{X}|\text{M}), \end{aligned} \quad (\text{B2})$$

where  $P(\text{M and X})$  is the probability that both M- and X-class flares will occur at least once during  $\tau$ , and  $P(\text{X}|\text{M})$  is the conditional probability of at least one X-class flare occurring given at least one M-class flare occurred during  $\tau$ .

Similarly,  $P(\geq \text{C1})$  can be determined as follows:

$$\begin{aligned} P(\geq \text{C1}) &= P(\text{C}) + P(\text{M}) + P(\text{X}) - P(\text{C and M}) - P(\text{C and X}) \\ &\quad - P(\text{M and X}) + P(\text{C and M and X}) \\ &= P(\text{C}) + P(\text{M}) + P(\text{X}) - P(\text{C}) \times P(\text{M}|\text{C}) - P(\text{C}) \times P(\text{X}|\text{C}) \\ &\quad - P(\text{M}) \times P(\text{X}|\text{M}) + P(\text{C}) \times P(\text{M}|\text{C}) \times P(\text{X}|\text{C and M}), \end{aligned} \quad (\text{B3})$$

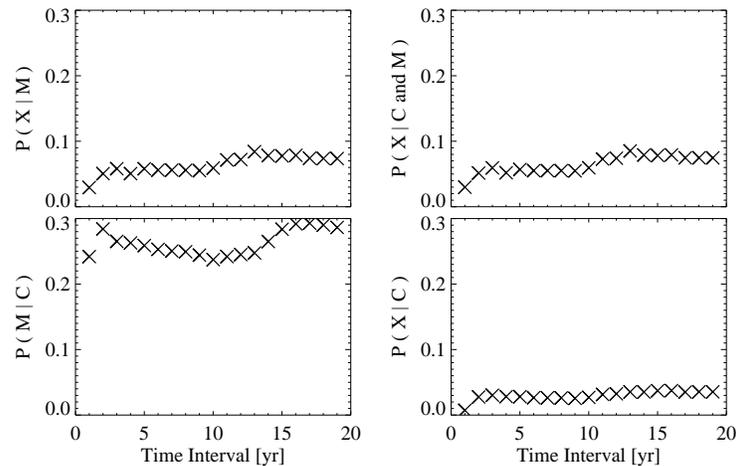
where  $P(\text{X}|\text{C and M})$  is the conditional probability of at least one X-class flare occurring given both C- and M-class flares occurred at least once during  $\tau$ .

The conditional probabilities are calculated using the NOAA/SWPC historical flare event list data and  $\tau$  as the prescribed validity interval (*e.g.*, 24 hr) starting from 00:00 UT of a given date. In this case, for example,  $P(X|M)$  can be determined as follows:

1. During the training interval for a given forecast method, we find the dates  $D(M)$  on which at least one M-class flare occurred.
2. From the dates  $D(M)$ , we determine the subset  $D(X|M)$  of dates on which at least one flare at the level greater than or equal to X1.0 occurred.
3. The conditional probability  $P(X|M)$  is then the total number of elements in  $D(X|M)$  divided by the total number of  $D(M)$ .

The other conditional probabilities can be calculated in the same way as  $P(X|M)$  explained above. Figure 5 shows the conditional probabilities for different time intervals used for their calculations. Note that the end date of all of the time intervals is fixed at 23:59 UT on 2017-Dec-31. The conditional probabilities do not significantly change as a function of the time interval. Because our goal is to calculate  $P(\geq C1)$  and  $P(\geq M1)$  from the probabilities of  $P(C)$ ,  $P(M)$  and  $P(\geq X1)$  that a given forecast method provides, the proper time interval to use for calculating the conditional probabilities is the training interval for that specific forecast method.

Forecasts for flare-class specific probabilities are converted to exceedance forecasts for the following methods: AMOS, ASAP, ASSA, MOSWOC, NICT, and NJIT.



**Figure 5.** C-, M-, X-class flare conditional probabilities as function of different time intervals as used for the calculations of exceedance. Time interval extends *back* in time from 2015.12.31.

## REFERENCES

- Aschwanden, M. J., Crosby, N. B., Dimitropoulou, M., et al. 2016, *Space Science Reviews*, 198, 47, doi: [10.1007/s11214-014-0054-6](https://doi.org/10.1007/s11214-014-0054-6)
- Barnes, G., & Leka, K. D. 2008, *ApJL*, 688, L107, doi: [10.1086/595550](https://doi.org/10.1086/595550)
- Barnes, G., Leka, K. D., Schrijver, C. J., et al. 2016, *ApJ*, 829, 89, doi: [10.3847/0004-637X/829/2/89](https://doi.org/10.3847/0004-637X/829/2/89)

- Berghmans, D., van der Linden, R. A. M., Vanlommel, P., et al. 2005, *Annales Geophysicae*, 23, 3115, doi: [10.5194/angeo-23-3115-2005](https://doi.org/10.5194/angeo-23-3115-2005)
- Bloomfield, D. S., Higgins, P. A., McAteer, R. T. J., & Gallagher, P. T. 2012, *ApJL*, 747, L41, doi: [10.1088/2041-8205](https://doi.org/10.1088/2041-8205)
- Bobra, M. G., & Couvidat, S. 2015, *ApJ*, 798, 135, doi: [10.1088/0004-637X/798/2/135](https://doi.org/10.1088/0004-637X/798/2/135)
- Centeno, R., Schou, J., Hayashi, K., et al. 2014, *Sol. Phys.*, 289, 3531, doi: [10.1007/s11207-014-0497-7](https://doi.org/10.1007/s11207-014-0497-7)
- Colak, T., & Qahwaji, R. 2008, *Sol. Phys.*, 248, 277, doi: [10.1007/s11207-007-9094-3](https://doi.org/10.1007/s11207-007-9094-3)
- . 2009, *Space Weather*, 7, 6001, doi: [10.1029/2008SW000401](https://doi.org/10.1029/2008SW000401)
- Crown, M. D. 2012, *Space Weather*, 10, 6006, doi: [10.1029/2011SW000760](https://doi.org/10.1029/2011SW000760)
- Devos, A., Verbeeck, C., & Robbrecht, E. 2014, *Journal of Space Weather and Space Climate*, 27, A29, doi: [10.1051/swsc/2014025](https://doi.org/10.1051/swsc/2014025)
- Domingo, V., Fleck, B., & Poland, A. I. 1995, *Sol. Phys.*, 162, 1, doi: [10.1007/BF00733425](https://doi.org/10.1007/BF00733425)
- Falconer, D., Barghouty, A. F., Khazanov, I., & Moore, R. 2011, *Space Weather*, 9, 4003, doi: [10.1029/2009SW000537](https://doi.org/10.1029/2009SW000537)
- Falconer, D. A., Moore, R. L., Barghouty, A. F., & Khazanov, I. 2014, *Space Weather*, 12, 306, doi: [10.1002/2013SW001024](https://doi.org/10.1002/2013SW001024)
- Florios, K., Kontogiannis, I., Park, S.-H., et al. 2018, *solphys*, 293, 28, doi: [10.1007/s11207-018-1250-4](https://doi.org/10.1007/s11207-018-1250-4)
- Gallagher, P., Moon, Y. J., & Wang, H. 2002, *Sol. Phys.*, 209, 171
- Georgoulis, M. K., & Rust, D. M. 2007, *ApJL*, 661, L109, doi: [10.1086/518718](https://doi.org/10.1086/518718)
- Hoeksema, J. T., Liu, Y., Hayashi, K., et al. 2014, *Sol. Phys.*, 289, 3483, doi: [10.1007/s11207-014-0516-8](https://doi.org/10.1007/s11207-014-0516-8)
- Hong, S., Kim, J., Han, J., & Kim, Y. 2014, AGU Fall Meeting Abstracts, SH21A
- Jolliffe, I. T., & Stephenson, D. 2012, *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, 2nd Edition (The Atrium, Southern Gate, Chichester, West Sussex PO19 8SQ, England: Wiley), doi: [10.1002/9781119960003](https://doi.org/10.1002/9781119960003)
- Kubo, Y., Den, M., & Ishii, M. 2017, *Journal of Space Weather and Space Climate*, 7, A20, doi: [10.1051/swsc/2017018](https://doi.org/10.1051/swsc/2017018)
- LaBonte, B. J., Georgoulis, M. K., & Rust, D. M. 2007, *apj*, 671, 955, doi: [10.1086/522682](https://doi.org/10.1086/522682)
- Lee, K., Moon, Y.-J., Lee, J.-Y., Lee, K.-S., & Na, H. 2012, *Sol. Phys.*, 281, 639, doi: [10.1007/s11207-012-0091-9](https://doi.org/10.1007/s11207-012-0091-9)
- Lee, S., Lee, J., & Hong, S. 2013, ASSA GUI User Manual, Version 1.07, [http://www.spaceweather.go.kr/images/assa/ASSA\\_GUI\\_M](http://www.spaceweather.go.kr/images/assa/ASSA_GUI_M)
- Leka, K. D., & Barnes, G. 2003, *ApJ*, 595, 1277
- . 2017, in *Extreme Events in Geospace: Origins, Predictability, Consequences*, 1st edn., ed. Buzulukova, N. (Cambridge, MA, USA: Elsevier), 65–98
- Leka, K. D., Barnes, G., & Wagner, E. L. 2017, *Sol. Phys.*, 292, 36, doi: [10.1007/s11207-017-1057-8](https://doi.org/10.1007/s11207-017-1057-8)
- . 2018, *Journal of Space Weather and Space Climate*, 8, A25, doi: [10.1051/swsc/2018004](https://doi.org/10.1051/swsc/2018004)
- Leka, K. D., Park, S. H., Kusano, K., et al. 2019, *ApJ*, accepted
- McCloskey, A. E., Gallagher, P. T., & Bloomfield, D. S. 2018, *Journal of Space Weather and Space Climate*, 8, A34, doi: [10.1051/swsc/2018022](https://doi.org/10.1051/swsc/2018022)
- McIntosh, P. S. 1990, *Sol. Phys.*, 125, 251
- Murphy, A. H. 1996, *Wea. Forecasting*, 11, 3
- Murray, S. A., Bingham, S., Sharpe, M., & Jackson, D. R. 2017, *Space Weather*, 15, 577, doi: [10.1002/2016SW001579](https://doi.org/10.1002/2016SW001579)
- Murray, S. A., Guerra, J. A., Zucca, P., et al. 2018, *Sol. Phys.*, 293, #60, doi: [10.1007/s11207-018-1287-4](https://doi.org/10.1007/s11207-018-1287-4)
- Nishizuka, N., Sugiura, K., Kubo, Y., et al. 2017, *ApJ*, 835, 156, doi: [10.3847/1538-4357/835/2/156](https://doi.org/10.3847/1538-4357/835/2/156)
- Park, S.-H., Chae, J., & Wang, H. 2010, *ApJ*, 718, 43, doi: [10.1088/0004-637X/718/1/43](https://doi.org/10.1088/0004-637X/718/1/43)
- Park, S. H., Leka, K. D., Kusano, K., et al. 2019, *ApJ*, in preparation
- Pesnell, W. 2008, in *COSPAR, Plenary Meeting*, Vol. 37, 37th COSPAR Scientific Assembly, 2412
- Sawyer, C., Warwick, J. W., & Dennett, J. T. 1986, *Solar Flare Prediction* (Boulder, CO: Colorado Assoc. Univ. Press)
- Scherrer, P. H., Bogart, R. S., Bush, R. I., et al. 1995, *Sol. Phys.*, 162, 129, doi: [10.1007/BF00733429](https://doi.org/10.1007/BF00733429)
- Scherrer, P. H., Schou, J., Bush, R. I., et al. 2012, *Sol. Phys.*, 275, 207, doi: [10.1007/s11207-011-9834-2](https://doi.org/10.1007/s11207-011-9834-2)

- Schou, J., Scherrer, P. H., Bush, R. I., et al. 2012, Sol. Phys., 275, 229, doi: [10.1007/s11207-011-9842-2](https://doi.org/10.1007/s11207-011-9842-2)
- Schrijver, C. J. 2007, ApJL, 655, L117, doi: [10.1086/511857](https://doi.org/10.1086/511857)
- Sharpe, M. A., & Murray, S. A. 2017, Space Weather, 15, 1383, doi: [10.1002/2017SW001683](https://doi.org/10.1002/2017SW001683)
- Steward, G., Lobzin, V., Cairns, I. H., Li, B., & Neudegg, D. 2017, Space Weather, 15, 1151, doi: [10.1002/2017SW001595](https://doi.org/10.1002/2017SW001595)
- Steward, G. A., Lobzin, V. V., Wilkinson, P. J., Cairns, I. H., & Robinson, P. A. 2011, Space Weather, 9, S11004, doi: [10.1029/2011SW000703](https://doi.org/10.1029/2011SW000703)
- Strugarek, A., Charbonneau, P., Joseph, R., & Pirot, D. 2014, Sol. Phys., 289, 2993, doi: [10.1007/s11207-014-0509-7](https://doi.org/10.1007/s11207-014-0509-7)
- Wheatland, M. S. 2000, ApJL, 536, L109, doi: [10.1086/312739](https://doi.org/10.1086/312739)
- . 2005, Space Weather, 3, 7003, doi: [10.1029/2004SW000131](https://doi.org/10.1029/2004SW000131)
- Woodcock, F. 1976, Monthly Weather Review, 104, 1209, doi: [10.1175/1520-0493\(1976\)104](https://doi.org/10.1175/1520-0493(1976)104)