

Deep learning reconstruction of sunspot vector magnetic fields for forecasting solar storms

DATTARAJ B. DHURI,^{1,2} SHAMIK BHATTACHARJEE,¹ SHRAVAN M. HANASOGE,^{1,2} AND SASHI KIRAN MAHAPATRA¹

¹*Department of Astronomy and Astrophysics, Tata Institute of Fundamental Research, Mumbai, India 400005*

²*Center for Space Science, New York University Abu Dhabi, Abu Dhabi, UAE*

ABSTRACT

Solar magnetic activity produces extreme solar flares and coronal mass ejections, which pose grave threats to electronic infrastructure and can significantly disrupt economic activity. It is **therefore important to appreciate the triggers of explosive solar activity and develop reliable space-weather forecasting**. Photospheric vector-magnetic-field data capture sunspot magnetic-field complexity and can therefore **improve the quality of space-weather prediction**. However, state-of-the-art vector-field observations are consistently only available from Solar Dynamics Observatory/Helioseismic and Magnetic Imager (SDO/HMI) since 2010, with most other current and past missions and observational facilities such as Global Oscillations Network Group (GONG) only recording line-of-sight (LOS) fields. Here, **using an inception-based convolutional neural network, we reconstruct HMI sunspot vector-field features from LOS magnetograms of HMI as well as GONG with high fidelity (~ 90% correlation) and sustained flare-forecasting accuracy**. We **rebuild vector-field features during the 2003 Halloween storms, for which only LOS-field observations are available, and the CNN-estimated electric-current-helicity accurately captures the observed rotation of the associated sunspot prior to the extreme flares, showing a striking increase**. Our study thus **paves the way for reconstructing three solar cycles worth of vector-field data from past LOS measurements, which are of great utility in improving space-weather forecasting models and gaining new insights about solar activity**.

Keywords: Sun: magnetic fields — methods: data analysis — methods: miscellaneous — methods: statistical

1. INTRODUCTION

Sunspot magnetic fields are generated within the solar interior, become buoyant through the solar convection zone and emerge at the photosphere and the corona as large-scale structures of sunspots and active regions (ARs) in the form of giant loops (Cheung & Isobe 2014). Coronal loops are dynamic, driven by emerging magnetic flux, electric current, and turbulent flows. Free magnetic energy stored in these loops is occasionally released via magnetic reconnection in the form of explosions such as flares and coronal mass ejections (CMEs) (Shibata & Magara 2011; Su et al. 2013). Radiation and charged particles emitted in these explosions can lead to severe space weather, disrupting our life on Earth significantly (Pulkkinen et al. 2005; Eastwood et al. 2017; Boteler 2019). In the past, the geomagnetic storm of 1989, resulting from a X15-class flare and subsequent CME, tripped circuit breakers in Hydro-Quebec power-grid causing

a widespread blackout in Quebec (Boteler 2019). The Halloween storm of 2003 produced extreme flares causing transformer malfunction and blackouts in Sweden, and damaging multiple science-mission satellites (Pulkkinen et al. 2005). In today's society, a high-magnitude solar storm can potentially lead to trillions of US dollars worth economic losses, with up to a decade of recovery time (Eastwood et al. 2017). Improving our understanding of AR magnetic-fields is therefore important for identifying triggers of these explosions and achieving reliable space-weather forecasting.

Coronal and photospheric AR magnetic fields are non-potential, comprising twisted flux-tubes as revealed by high-resolution, high-cadence observations of the SDO (Pesnelli et al. 2012) since 2010. Large ARs and their complex dynamics, e.g. twisting and rotation, are known to be associated with solar explosive activity (Toriumi & Wang 2019). The SDO/HMI photospheric vector-magnetic-field observations facilitate the calculation of AR features (Leka & Barnes 2007), such as total unsigned magnetic flux, free energy density, electric current helicity and Lorentz forces, characterizing the AR magnetic-field dynamics. These features are publicly available as the HMI data-product Space-weather HMI

Active Region Patches (SHARPs) (Bobra et al. 2014). The SHARPs features are extensively used for statistical studies of pre-flare magnetic-field evolution and energy build up (Dhuri et al. 2019) and improving space-weather forecasting using Machine Learning (ML) (Bobra & Couvidat 2015; Bobra & Ilonidis 2016; Chen et al. 2019). HMI observations are limited to only one full solar cycle (cycle 24) and therefore, statistical space-weather forecasting models based on SHARPs are restricted.

Various difficulties are associated with the measurement of the transverse component of the photospheric magnetic-field (Stenflo 2013) and therefore, ground- and space-based instruments monitoring the Sun since the 1970s, provide observations of only the longitudinal, i.e., line-of-sight (LOS) component. Continuous full-disk LOS field observations are available through the ground-based NASA/National Solar Observatory (NSO) Kitt Peak Telescope (1974 - present) (Livingston et al. 1976), space-based Michelson and Doppler Imager (MDI, 1996 - 2011) (Scherrer et al. 1995) and ground-based Global Oscillations Network Group (GONG, 1995-present). These LOS-field measurements, although not sufficient for quantifying sunspot complexity to non-potential energy and helicity, have been useful for providing a qualitative assessment of AR morphology via sunspot classification schemes, such as the McIntosh classification (McIntosh 1990) and Mount Wilson classification which form the basis of operational space-weather forecasts (Crown 2012).

Improving on these qualitative AR classifications and formally devising a method to quantify vector-field properties from LOS fields is of great utility — (i) because it allows for “improving” past datasets of LOS observations and understanding how vector-field features have evolved over multiple solar cycles, (ii) a reliable estimation of vector-field features over the past few decades can be used to build more robust statistical models for space-weather forecasting, and (iii) for future missions acquiring only LOS data, vector-field features and even full vector-field construction can be part of an on-ground data-processing pipeline. ML methods such as convolutional neural networks (CNN) developed through the past decade have proven to be hugely successful in identifying patterns and correlations in large, high-dimensional datasets and particularly images (LeCun et al. 2015; Goodfellow et al. 2016). Here, we explore dependencies between LOS magnetograms and the corresponding full vector-field of ARs through a CNN model developed to estimate vector-field features SHARPs using the LOS magnetograms measurements from space-based SDO/HMI as well as ground-based GONG.

2. DATA

We use photospheric LOS-magnetogram data provided by HMI and GONG. GONG provides only LOS magnetograms.

	Data	
	Train & Val	Test
	May'10 - Sep'15	Oct'15 - Aug'18
# HMI ARs	848	194
# HMI Samples	124633	26820
# GONG ARs	848	145
# GONG Samples	114443	13454

Table 1. The Helioseismic Magnetic Imager (HMI) and Global Oscillation Network Group (GONG) data used for training a CNN to obtain vector-field features from Space-Weather HMI Active Region Patches (SHARPs) (Bobra et al. 2014).

HMI-derived SHARPs (the hmi.sharp_cea_720s data series (Bobra et al. 2014)) include vector and LOS magnetograms of AR patches that are automatically detected and tracked as they rotate across the visible solar disk (Bobra et al. 2014). HMI magnetograms are available at a plate scale of 0.5 arcsecs, i.e., ~ 380 km at the disk center. GONG magnetograms are available at a plate scale of 2.5 arcsecs. The magnetograms available in the hmi.sharp_cea_720s series are on a cylindrical equal-area (CEA) grid, thus eliminating the projection effects. We similarly remap the GONG AR magnetograms to a CEA grid. We train a CNN to obtain SHARPs features directly from LOS magnetograms of HMI as well as GONG. We only consider top SHARPs features that produce maximum flare forecasting accuracy for a machine learning (ML) model (Bobra & Couvidat 2015). These are listed in Table 2.

HMI measurements are sensitive to the observation conditions as well as the relative velocity between SDO and the Sun (Hoeksema et al. 2014). Observation conditions are indicated by the QUALITY flag and we consider measurements for which the Stokes vectors are reliable (QUALITY ≤ 10000 in hexadecimal) and when the relative velocity between SDO and the Sun is < 3500 m/s (Bobra & Ilonidis 2016). Data closer to the limb are noisier because the higher relative velocities as well as projection effects. Therefore, we limit observations to within $\pm 45^\circ$ of the central meridian. Further, we only include ARs from the SHARPs data series that grow to a maximum area of > 25 Mm². This eliminates a significant number of small ARs that do not produce major (M- or X-class) flares. The SHARPs feature calculation using HMI vector-field observations considers those pixels in the AR magnetograms for which the 180° ambiguity resolution is reliable (Bobra et al. 2014).

Observations between May 2010 and Aug 2018 are used to train the CNN — approximately 80% of the data are used to train and validate the CNN, while the remaining is the unseen or test data. We chronologically split the ARs into training and validation and test partitions: ARs in the period May 2010 - Sep 2015 for training and validation and Oct 2015 -

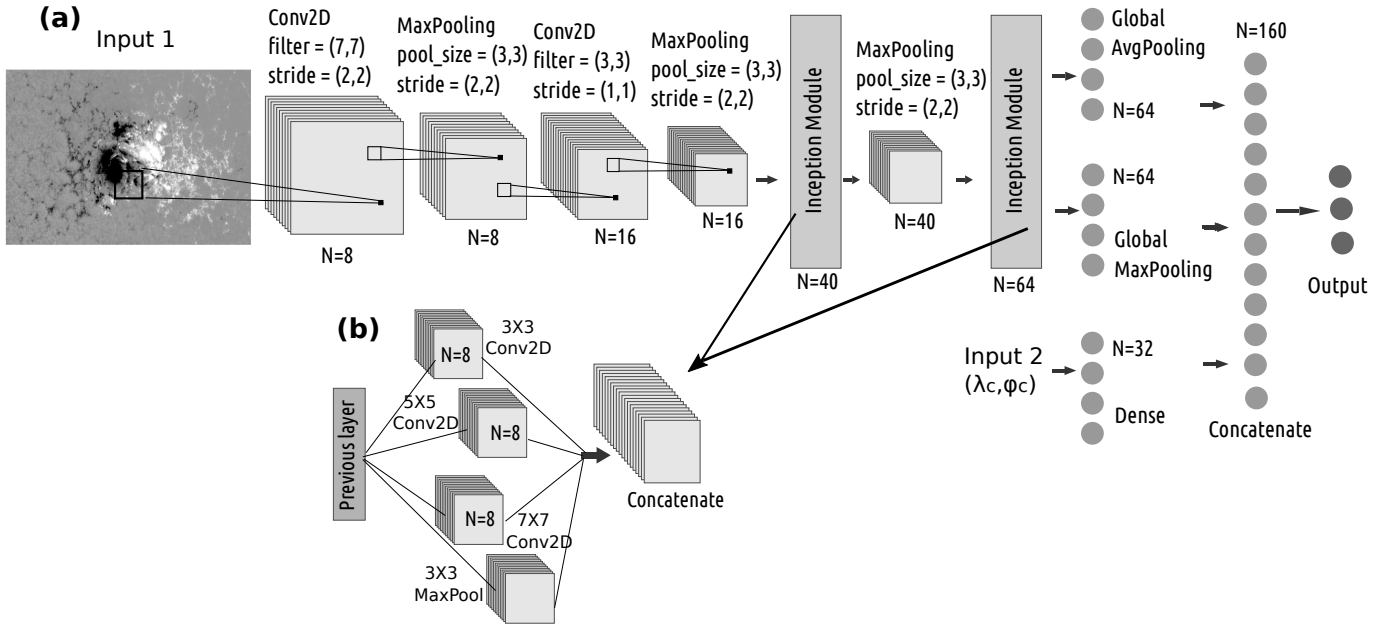


Figure 1. The CNN architecture. (a) Convolutional neural network (CNN) architecture used for obtaining vector-field features from LOS magnetograms. The architecture incorporates inception modules similar to GoogleNet [Szegedy et al. \(2015\)](#). The CNN takes in two inputs — i) LOS magnetograms ii) AR center latitude (λ_c) and longitude (ϕ_c). There are no fully connected layers that directly process the magnetogram input and therefore the CNN can process magnetogram patches of variable sizes. (b) Inception module used in the CNN.

Aug 2018 is for the test. Six-hourly samples are drawn from the time series of each AR. All samples from a given AR are exclusively part of either the training set, validation set or the test set, to avoid biases arising from temporal coherence of observations from an AR ([Ahmadzadeh et al. 2021](#)). The number of ARs and magnetogram observations used for training, validation and test are listed in Table 1. Since solar activity depends on the phase of the cycle, the chronological splitting may introduce a bias for training the CNN. Indeed, the ratio of flaring-to-nonflaring ARs in the test data is approximately half its value in the training and validation dataset ([Bhattacharjee et al. 2020](#)). However, chronological splitting is appropriate for operational space-weather forecasting tools.

3. METHODS

CNNs are neural networks with convolution filters (kernels) to scan over the input data, typically 2D data of images, and detect spatial patterns for tasks such as classification and identification ([LeCun et al. 2015](#); [Goodfellow et al. 2016](#)). The convolution filters are $K \times K$ neurons that slide over the images and detect different patterns. Convolution filters have free parameters — each neuron has weight w and each convolution filter has bias b . Neurons process pixels of the inputs (or the outputs from previous layers) x_i by performing the operation $f(\sum_i w_i x_i + b)$, where f is the activation function ([Hastie et al. 2001](#)). CNNs also have pooling layers which are used to reduce the input size as it progresses to deeper levels of the CNN. A max-/average-pooling filter picks out

the maximum or average value from the given $N \times N$ feature map. Pooling layers typically follow a convolutional layer in a CNN to reduce the dimensionality.

We use a CNN architecture with inception modules similar to inception V1 modules from GoogleNet ([Szegedy et al. 2015](#)). Typically, in a convolutional layer, we use filters of fixed size that work best for the particular problem. However, inception modules are designed to detect patterns over a variety of length scales that may be present in the input. They involve convolution filters of different sizes in a single layer. The outputs from all the convolutional layers in an inception module are concatenated and supplied as an input to the following layer. The inception module used here comprises three convolution filters of sizes 3×3 , 5×5 and 7×7 and one 3×3 max-pooling filter.

The CNN architecture is shown in Figure 1. The CNN takes in two inputs — i) LOS magnetograms of AR patches and ii) latitude (λ_c) and longitude (ϕ_c) of the center of AR patches. The CNN consists of two regular convolutional layers and followed by the two inception modules that process the LOS magnetograms. The latitude and longitude are processed by a fully connected layer of neurons. The output of the two regular convolutional layers and the first inception module are reduced by a max-pooling layer. The output of the final inception module is reduced by a global max-pooling layer and also a global-average pooling layer. These are concatenated with the output from the fully connected layer that processes the longitude and latitude. The concatenated layer is connected to the output layer of neurons. The

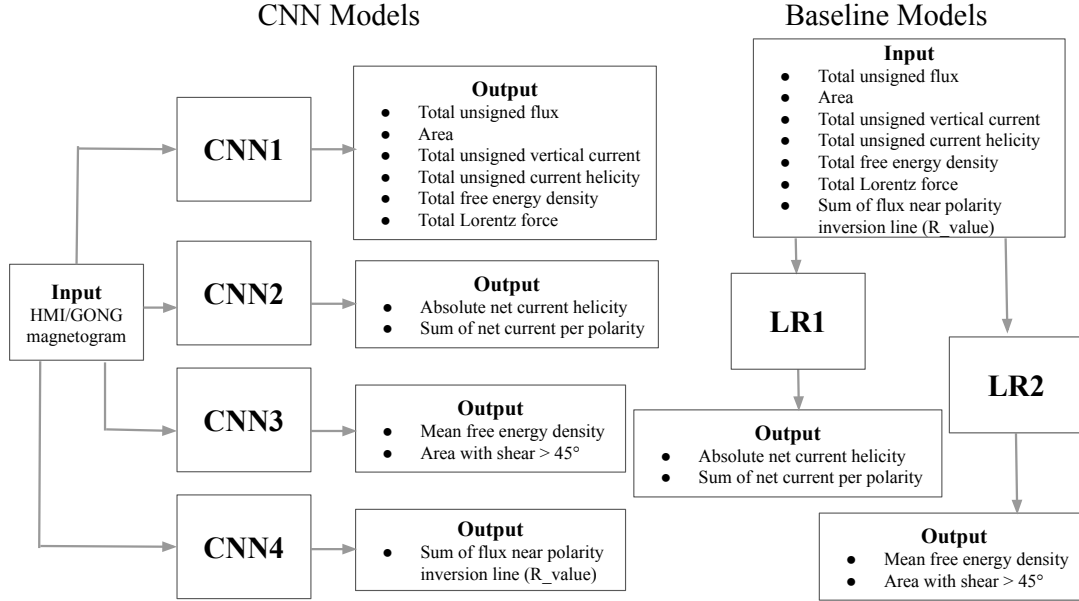


Figure 2. Schematic of ML models. (left:) Convolutional Neural Network (CNN) models process LOS magnetograms as input and produce vector-field features SHARPs as output. SHARPs features group together in four groups based on their mutual correlations as shown. We develop four different CNN models to estimate SHARPs from four different groups. All CNN models have identical architecture described in Figure 1 except the final output layer, where, the number of output neurons is equal to the number of SHARPs features to be estimated from the respective group. (right:) Baseline models using Linear Regression (LR) for estimation of two groups of SHARPs features which depend on electric current and free energy respectively using extensive SHARPs features and Schrijver’s R.value (Schrijver 2007) as an input.

number of neurons in the output layers is equal to the number of SHARPs features being estimated (see Figure 2). We use a *linear* activation function for the convolutional layers, which explicitly treats the positive and negative pixel values from LOS magnetograms symmetrically. Also, the fully connected layer of neurons has a *tanh* activation to explicitly treat positive and negative values of latitude and longitude, that are normalised between ± 1 , symmetrically. The final output layer of neurons have *sigmoid* activation (Han & Moraga 1995; Hastie et al. 2001) to yield the normalised value of the estimated SHARPs features between 0 and 1.

The absence of fully connected layers in the network that processes the LOS-magnetogram input implies that the CNN architecture can analyze LOS magnetograms of arbitrary sizes. Since AR patches are of varied dimensions, magnetograms in the training, validation and test data are also correspondingly differently sized. As such, our CNN does not require pre-processing to convert magnetograms to a fixed size and thus it is free from biases that may arise as a result of resizing (Bhattacharjee et al. 2020).

We use 10-times repeated-holdout validation for training the CNN (Hastie et al. 2001). We randomly split ARs in the training and validation sets into three parts and use data from two parts for training and the remaining part for validation. This process is repeated nine times while ensuring that the data from an AR is part of either the training or the validation and not both. The output from the CNN is compared

to the original SHARPs feature values. The *sigmoid* output layer of the CNN lies in a continuous range between 0 to 1. The original SHARPs features are normalised by dividing by their respective maximum values. We partition the normalised features (over range 0 to 1) in the training set into ten bins of equal width (0.1) and oversample the data in each bin to match the number of samples in the maximally populated bin. The input magnetograms are standardised, i.e., a mean is subtracted and the resultant magnetogram is divided by a standard deviation of the magnetic field values. The mean and standard deviation used for standardisation are calculated over all pixels of all magnetograms in the training and validation data of the respective instrument. The CNN output is compared to the original SHARPs values and the loss function — defined to be the mean squared error — is computed. We train the CNN to minimize the mean squared error over different epochs using stochastic gradient descent (Bottou 1991; Hastie et al. 2001) with a learning rate of 0.00007. The CNN is developed using the Python library *keras*.

4. RESULTS

4.1. Estimation of AR vector-magnetic-field features using CNN

The SHARPs features considered (listed in Figure 2 and Table 2) are correlated among each other and are divided into four groups based on mutual Pearson correlations (Dhuri et al. 2019): (i) features that depend on the area of ARs, i.e.,

SHARPs Features	10-times Repeated-Holdout Validation		Test	
	HMI	GONG	HMI	GONG
Total unsigned flux	95.14 ± 00.62	90.87 ± 01.96	89.73 ± 02.70	87.42 ± 01.39
Area	95.87 ± 00.49	95.06 ± 00.84	92.00 ± 01.70	92.88 ± 00.88
Total unsigned vertical current	94.78 ± 00.71	91.80 ± 01.69	88.86 ± 02.57	89.00 ± 01.76
Total unsigned current helicity	95.74 ± 00.50	91.65 ± 01.76	88.33 ± 02.65	83.31 ± 02.28
Total free energy density	96.19 ± 00.80	92.60 ± 01.60	90.17 ± 02.37	91.22 ± 01.25
Total Lorentz force	96.64 ± 00.47	94.94 ± 00.98	90.63 ± 02.46	92.71 ± 00.87
Absolute net current helicity	90.37 ± 03.28	63.76 ± 03.65	57.83 ± 08.84	57.60 ± 06.97
Sum of net current per polarity	89.51 ± 02.53	64.58 ± 03.08	61.93 ± 07.63	59.09 ± 06.96
Mean free energy density	95.10 ± 01.00	89.92 ± 00.79	92.13 ± 01.80	91.73 ± 00.54
Area with shear > 45°	95.02 ± 00.81	90.00 ± 01.19	90.59 ± 01.57	90.48 ± 00.46
Flux near polarity inversion line	90.54 ± 00.56	76.28 ± 01.83	77.11 ± 00.79	70.43 ± 00.79

Table 2. Pearson correlations between the CNN-estimated vector-field features SHARPs and their true values. SHARPs (Space-weather HMI Active Region Patches) features are calculated from HMI vector-field observations (Bobra et al. 2014). The p-values for all correlations are ~ 0.0 . SHARPs features are mutually correlated (Dhuri et al. 2019) and accordingly arranged in the four groups as features depending on (i) AR area (ii) electric current (iii) mean free energy density (iv) R-value i.e. the sum of flux near polarity inversion line (Schrijver 2007).

SHARPs Features	10-times Repeated-Holdout Validation		Test	
	HMI	GONG	HMI	GONG
Total unsigned flux	86.92 ± 01.13	86.27 ± 01.75	81.19 ± 01.74	77.11 ± 02.20
Area	89.57 ± 01.10	92.11 ± 00.75	87.10 ± 01.33	86.91 ± 01.21
Total unsigned vertical current	86.82 ± 01.41	87.34 ± 01.38	81.22 ± 02.10	79.07 ± 02.52
Total unsigned current helicity	87.37 ± 01.41	87.49 ± 01.48	82.61 ± 02.24	79.40 ± 02.69
Total free energy density	84.23 ± 01.46	85.53 ± 02.17	81.32 ± 03.27	80.11 ± 02.42
Total Lorentz force	90.63 ± 01.04	92.78 ± 00.98	86.96 ± 01.51	86.49 ± 01.66
Absolute net current helicity	59.61 ± 05.26	59.35 ± 03.36	57.70 ± 06.29	2.27 ± 02.73
Sum of net current per polarity	60.02 ± 03.95	65.75 ± 02.95	57.51 ± 02.58	7.85 ± 02.79
Mean free energy density	92.02 ± 00.97	91.59 ± 00.92	93.02 ± 01.42	92.58 ± 00.35
Area with shear > 45°	93.19 ± 00.98	90.08 ± 00.54	91.63 ± 00.65	89.96 ± 00.37
Flux near polarity inversion line	91.69 ± 00.47	76.63 ± 02.60	83.00 ± 00.85	75.32 ± 00.92

Table 3. Spearman correlations between the CNN-estimated vector-field features SHARPs and their true values. The p-values for all correlations are ~ 0.0 .

extensive features that include AR area, total unsigned flux, total unsigned vertical current, total unsigned current helicity, total free energy density and total Lorentz force, (ii) features that depend on the electric current in ARs, i.e., absolute net current helicity and sum of net current per polarity, (iii) features that depend only on the non-potential energy in ARs, i.e., mean free energy density and area with shear > 45°, and finally, (iv) Schrijver R-value (Schrijver 2007) viz. the sum of flux on the polarity inversion line. Overall, we develop four different CNNs (Figure 2) to estimate SHARPs features from these respective four groups. For each CNN, the output layer comprises K neurons to estimate K SHARPs features corresponding to each of the four groups.

Extensive features are strongly correlated with AR total unsigned flux, that depends only on the radial component of the magnetic field. The radial component is tradition-

ally estimated from AR LOS magnetic field using a potential field approximation (Leka et al. 2017). Using a CNN, we directly estimate these extensive features without first requiring to estimate the radial magnetic field. The R-value depends only on LOS magnetic field and can be directly calculated using GONG LOS magnetograms. However, to match HMI SHARPs R-value, GONG LOS magnetic fields require a cross-calibration. CNN models are expected to implicitly learn the cross-instrument calibration during training (Munoz-Jaramillo et al. 2022) and estimated SHARPs values are also expected to be automatically cross-calibrated.

Unlike the extensive features, an accurate estimation of SHARPs depending on electric current and mean free energy requires explicit knowledge of the full vector-magnetic fields. Such features are important for understanding triggers of solar storms and are typically estimated assuming magnetic-

SHARPs Features	10-times Repeated-Holdout Validation		Test	
	Pearson	Spearman	Pearson	Spearman
Absolute net current helicity	70.99 ± 06.08	68.21 ± 02.25	67.10 ± 03.94	70.32 ± 00.56
Sum of net current per polarity	71.83 ± 03.79	70.51 ± 02.23	57.97 ± 02.90	67.07 ± 00.55
Mean free energy density	72.11 ± 04.87	79.63 ± 04.49	74.41 ± 01.90	81.91 ± 00.66
Area with shear > 45°	67.89 ± 02.93	69.16 ± 03.41	70.83 ± 00.39	77.86 ± 00.20

Table 4. Correlation between SHARPs features estimated using the baseline models and their true values. Two baseline Linear Regression (LR) models are developed (Figure 2) that take extensive SHARPs features and Schrijver’s R_value as input and produce SHARPs features that depend on electric current and mean non-potential energy as output. The p-values for all correlations are ~ 0.0 .

		Total unsigned flux		Absolute net current helicity		Mean free energy density	
		Spline Fit	Time Derivative	Spline Fit	Time Derivative	Spline Fit	Time Derivative
Validation	HMI	97.41 ± 00.37	84.46 ± 02.75	94.99 ± 02.02	77.70 ± 12.16	98.54 ± 00.38	89.68 ± 02.50
	GONG	91.62 ± 01.94	25.51 ± 08.68	66.12 ± 03.73	33.23 ± 05.40	92.12 ± 07.77	75.52 ± 02.29
Test	HMI	90.50 ± 02.60	75.18 ± 06.81	60.21 ± 09.21	32.44 ± 18.33	94.14 ± 01.52	82.06 ± 03.06
	GONG	88.63 ± 01.48	18.07 ± 06.22	59.72 ± 07.14	27.03 ± 07.81	93.69 ± 00.49	70.52 ± 01.91

Table 5. Trend comparison of CNN-estimated and true SHARPs features. Pearson correlations between the time derivative of the true and the CNN-estimated values of total unsigned flux, absolute net current helicity and mean free energy density. The time derivative is obtained for the true and the CNN-estimated features of each AR after fitting the respective time series to a cubic spline. For reference, the Pearson correlations between the spline fit values of the true and the CNN-estimated features are also shown. Note that these are consistent with the values in Table 2.

field models, e.g., linear and non-linear force-free models (Régner & Priest 2007). Here, we provide a purely data-driven estimation of these features using a CNN. In order to assess the performance of the CNN, we use Linear Regression (LR) models as a baseline. We develop two separate LR models, one each for features that depend on electric current and free energy, respectively. As input, the LR models have extensive features and R_Value. The first LR model produces absolute net current helicity and sum of net current per polarity as the output while the second produces mean free energy density and area with shear > 45° as the output. Figure 2 shows a schematic of the CNN models as well as the baseline models.

We use Pearson and Spearman correlations for measuring the performance of the CNN and baseline models. Pearson correlation measures a linear correlation between the true and estimated values of the vector-magnetic-field features. Spearman correlation is a rank correlation that captures the monotonic relationship between the true and estimated values in addition to the linear relationship measured by the Pearson correlation. Pearson and Spearman correlations for the CNN-estimated vector-magnetic-field features are listed in Tables 2 and 3 respectively. For the baseline models, these correlations are listed in Table 4.

From Table 2, the Pearson correlations of CNN-estimated vector-field features is higher for HMI than GONG and thus appear to be dependent on the spatial resolution of LOS magnetograms. For HMI, the CNN-estimated extensive features yield a Pearson correlation of $\sim 95\%$ for the validation

and $\sim 90\%$ for the test data. For GONG data, the corresponding correlation is $\sim 90\%$. The Pearson correlations of CNN-estimated values of extensive features are not a perfect $\sim 100\%$, since the SHARPs calculation does not consider all pixels, rather, only taking into account those for which disambiguation of the azimuthal component of the magnetic field is reliable (Bobra et al. 2014). From Table 3, the Spearman correlations for the extensive features are only slightly lower than the corresponding Pearson correlations implying that the ranking of the estimated features is generally consistent with the true ranking.

The Pearson and Spearman correlation values for features that depend on the non-potential energy are significantly high > 90% across the validation and test datasets. These correlation values are also $\sim 10\% - 20\%$ higher compared to the linear regression baseline model (Table 4). These features, namely mean free energy density and area with shear > 45°, explicitly depend on the full vector-magnetic-field.

For features that depend on electric current — i.e., absolute net current helicity and sum of net current per polarity — the CNN does not perform better than the baseline model. While the Pearson correlation for the validation are 20% higher compared to the baseline of 70%, Spearman correlations are approximately equal (up to the error bars) at 60%. Also, the CNN fails to generalise to the test data with a low Pearson and Spearman correlation scores of 60% each.

Figure 3 shows scatter plot visualisations of the correlation between true and CNN-estimated SHARPs features for HMI and GONG from the ten validation sets. For HMI, the

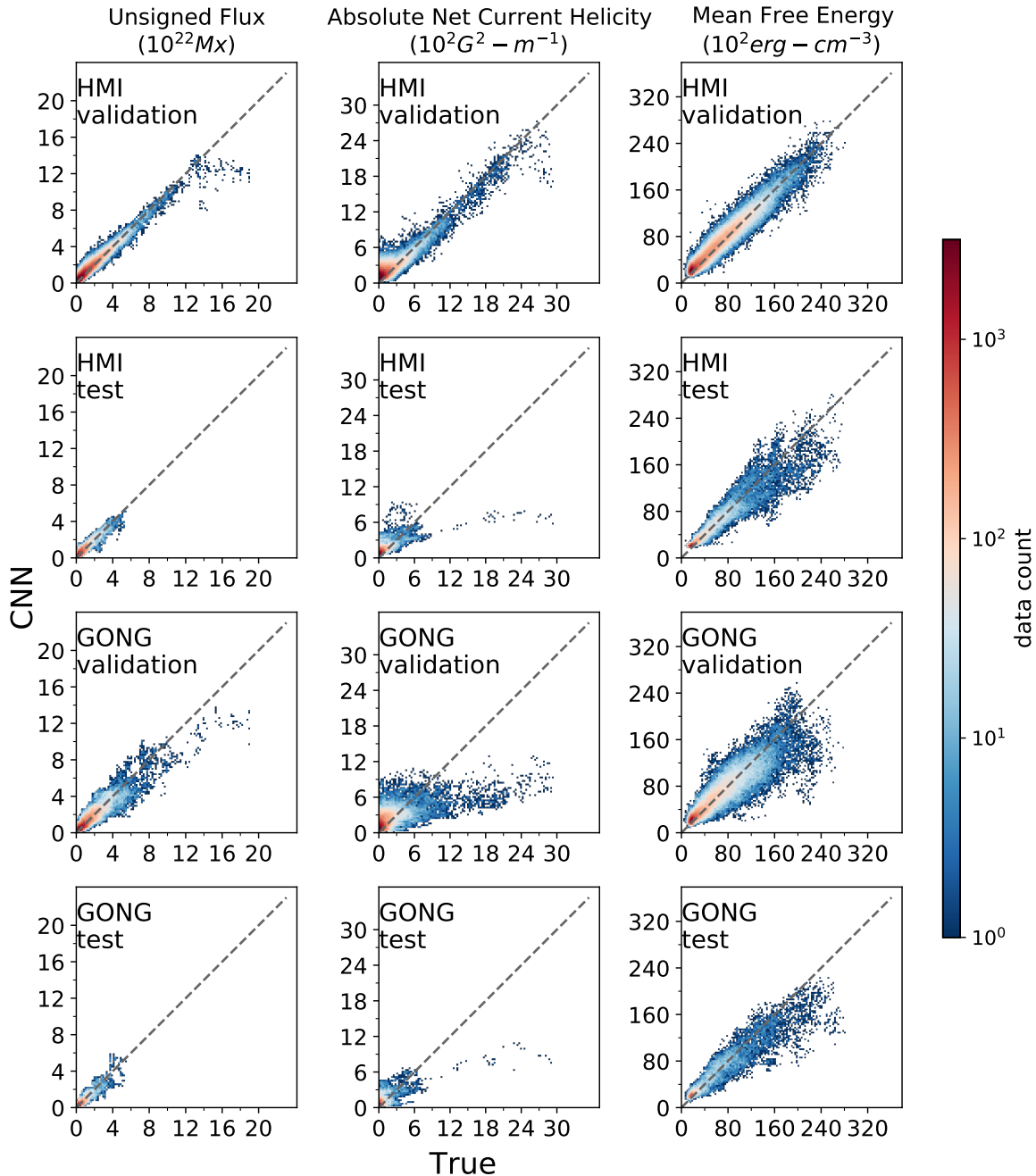


Figure 3. Comparison of true and CNN-estimated SHARPs vector-field features. Scatter plots for the distribution of CNN-estimated and true HMI values of total unsigned flux, absolute current helicity and mean free energy density obtained from HMI and GONG LOS magnetograms. The 45° line is marked for reference.

true and CNN-derived values mostly match relatively closely, except at only very small values ($< 200 G^2 - m^{-1}$) of absolute net current helicity where the CNN estimates are significantly larger. For the HMI test data as well as GONG data, the CNN-estimation of absolute net current helicity for large values ($> 1000 G^2 - m^{-1}$) is consistently on the lower side ($\sim 500 G^2 - m^{-1}$). Figure 4 explicitly shows mean absolute errors in the CNN estimation as a function of the true values for HMI and GONG. Mean absolute er-

rors in CNN-estimated values from GONG magnetograms show higher dependence on true values compared to HMI and increase significantly with increasing true values of the respective features, particularly for the validation data. For total unsigned flux, mean absolute errors of CNN-estimated features of both HMI and GONG are significantly higher for the extreme values $15 - 20 \times 10^{22} Mx$. For the HMI test data and GONG data, mean absolute errors are more than 12 times higher at large magnitudes ($> 1000 G^2 - m^{-1}$) of

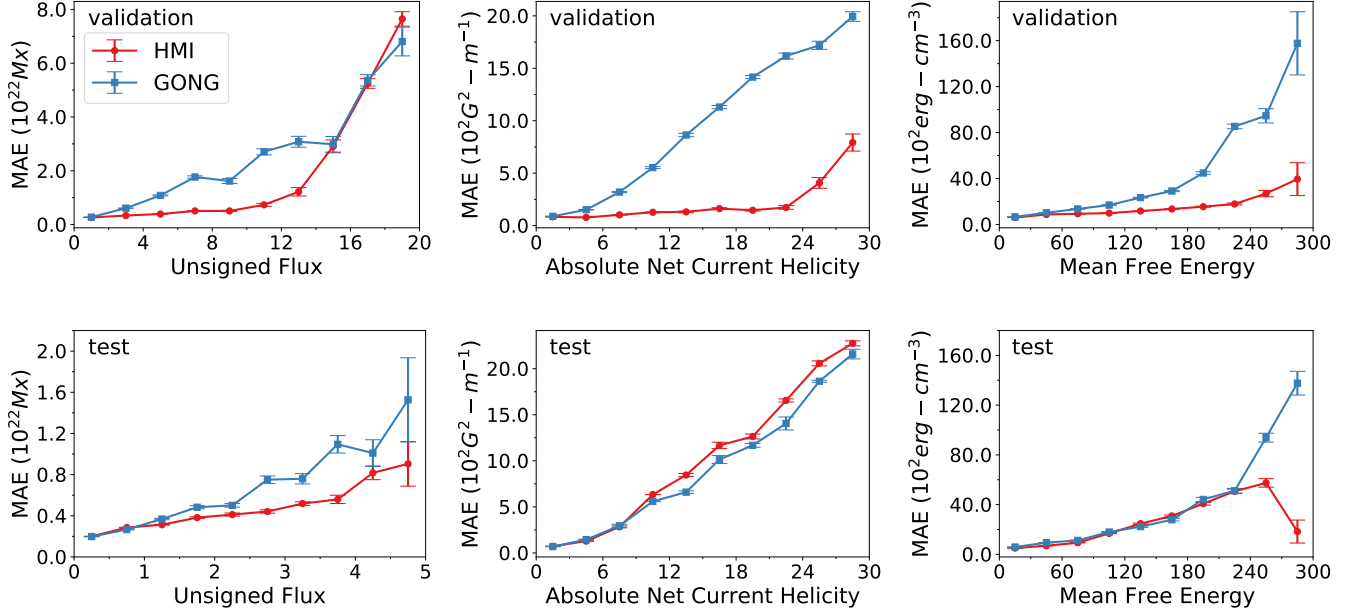


Figure 4. Mean absolute error of the CNN-reconstructed vector-field features. The CNN-estimated values of total unsigned flux, absolute current helicity and mean free-energy density for HMI and GONG are binned into 10 uniform bins as per the respective true values. Mean absolute error as a function of the mean true values of each population bin are shown. $1\text{-}\sigma$ error bars for each bin are also shown. The legend in the top left panel applies to all panels.

absolute net current helicity compared to the HMI validation data. The average relative errors for GONG and HMI are comparable at $\approx 80 \pm 10\%$, $900 \pm 100\%$ and $25 \pm 2\%$ for total unsigned flux, absolute net current helicity and mean free energy density respectively. The high average relative errors imply that the CNN estimates are far off from true values, particularly for SHARPs features with low true values. SHARPs features from ARs which produce at least one major flare (M5 or greater) show a significant drop in average relative errors, at approximately $30 \pm 15\%$, $300 \pm 80\%$ and $16 \pm 2\%$ respectively.

4.2. Time evolution of the CNN-derived features on flaring active regions

For understanding AR magnetic-field dynamics and improving forecasting of solar storms, it is important that temporal variations of the CNN-estimated SHARPs is faithful to the true SHARPs. We measure trends in the time evolution of SHARPs features of an AR by fitting the observed and the CNN-estimated values with smooth spline curves and calculate numerical time derivatives. Table 5 lists Pearson correlations between time derivatives of splines, fitted to the true and the CNN-estimated values of total unsigned flux, absolute net current helicity and mean free energy density. We find that the Pearson correlations are high, $\sim 80\%$, for HMI, with the exception of absolute net current helicity values from the test data. For GONG, only the Pearson correlations for mean free energy density are high enough, $\sim 70\%$, to suggest that the

corresponding trends are captured reasonably accurately in the CNN-estimated features. These discrepancies between trends of the CNN-estimated GONG and true values appear to be a consequence of the lower resolution of GONG magnetograms.

A comparison of the time evolution of true and CNN-estimated features obtained from HMI and GONG for individual ARs that produce at least one major flare (M5 or greater) is shown in Figure 5. The true and CNN-estimated values of total unsigned flux, absolute net current helicity and mean free energy density are in agreement, particularly for HMI, capturing evolution of these features before and after flares. Disagreements between the true and CNN-estimated features occur only at the extreme values of these features. E.g., for X9.3 flare in NOAA 12673 in September 2017, which was the largest flare in cycle 24, the CNN accurately estimates the rise of the total unsigned flux and also mean free energy density prior to the flare. The absolute net current helicity rises to unusually high values prior to the X9.3 flare (30% more than the maximum values encountered in the training data) and therefore the corresponding CNN estimates are inaccurate. More examples of comparisons of the time evolution between true and CNN-estimated features from ARs that produce at least one M5 or greater flare are included in the Appendix Figure 9.

The CNN estimation of SHARPs features on flaring ARs is thus useful for understanding AR magnetic-field evolution leading to particularly violent solar storms in the past. The

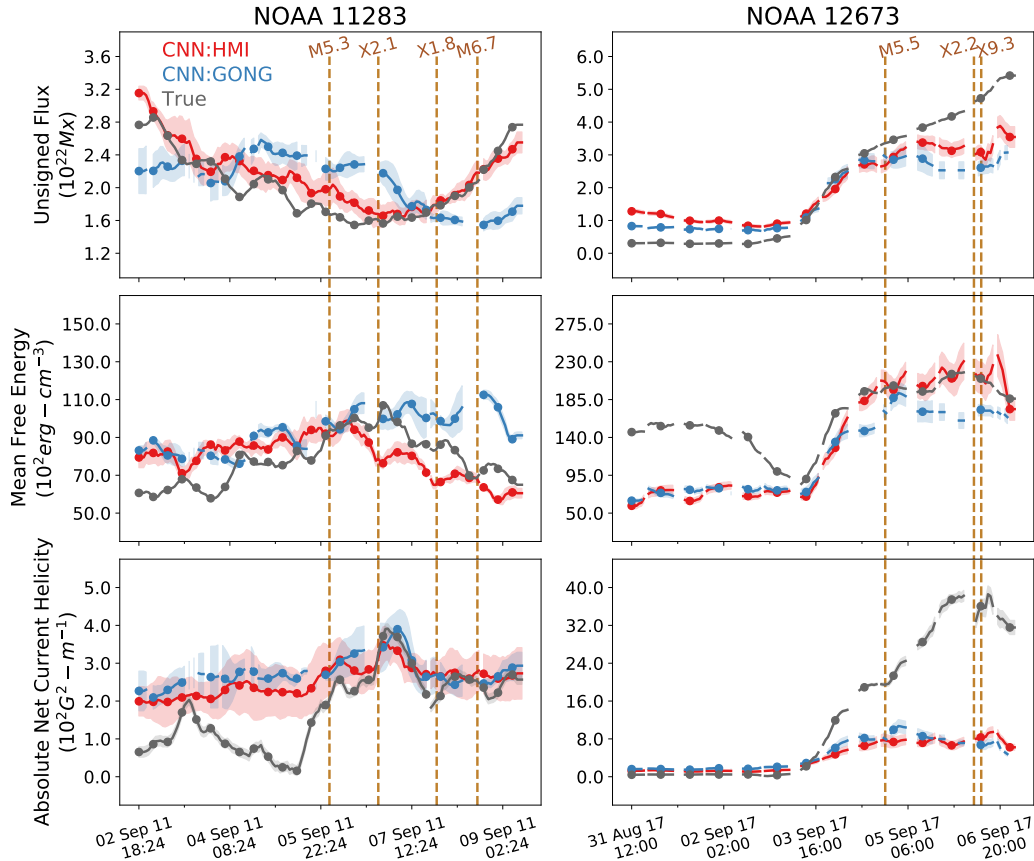


Figure 5. Comparison of true and CNN-estimated SHARPs vector-field features. Comparisons of time evolution of the CNN-estimated total unsigned flux, absolute net current helicity and mean free energy density with true values for ARs that produce M5 or greater flares. Only AR observations within $\pm 45^\circ$ of the central meridian are considered. The left plot shows a typical result (see Appendix Figure 9 for all ARs with major flares). The right plot shows an extreme event with the largest flare observed in cycle 24. The gaps correspond to the missing observations and $1-\sigma$ error bars are shown. The legend in top left applies to all plots. The plots are smoothed with a six hourly running average.

Halloween storms of October 2003 produced extreme flares from AR NOAA 10486 of magnitudes X17.0, X10.0 and the largest recorded flare X28.0 (Pulkkinen et al. 2005). The magnetic-field evolution leading to these extreme flares was characterised by rotation of a major positive polarity of the delta sunspot as shown in the top panel of Figure 6 (Zhang et al. 2008). Without the knowledge of vector-magnetic-fields, free energy and current helicity during these storms are previously modelled based on the magnetic virial theorem (Metcalf et al. 2005; Régnier & Priest 2007), linear/non-linear force-free field extrapolation (Régnier & Priest 2007), and a Minimum Current Corona model (Kazachenko et al. 2010). We obtain a model-free and purely data-driven CNN-estimates of total unsigned flux, absolute net current helicity and mean free energy density during these storms using LOS magnetograms. However, the HMI observations are not available for this period. We therefore use the CNN trained with HMI magnetograms to process LOS observations from MDI during the Halloween storms to estimate time evolution of total unsigned flux, absolute net current helicity and mean free energy density. Flare X28.0 is excluded as it occurred

outside 45° of the central meridian. In particular, the CNN-estimated absolute net current helicity of NOAA 10486 rises continuously by 25% between X1.2 flare and X17.0 flare corresponding to the observed sunspot rotation. A similar gradual rise of a modelled helicity flux by 50% between X1.2 and X17.0 flare has been reported (Kazachenko et al. 2010), caused primarily by helicity injection from the rotation of the sunspot. The CNN estimates show that the absolute net current helicity stays high leading to the X10.0 flare and falls thereafter. The CNN-estimated mean free energy density also rises leading to the X17.0 and X10.0 flares. Note that these CNN-estimated values from MDI magnetograms are not expected to be corrected for the instrument cross-calibration between the MDI and HMI since the CNN is trained with only HMI magnetograms. Table 8 in the Appendix lists the Pearson and Spearman correlations between the true values and the CNN-estimated values using MDI line-of-sight magnetograms, during the overlap period of MDI and HMI. These correlation values are significantly lower compared to those estimated from HMI magnetograms (Tables 2 and 3). Therefore, a rigorous estimation first requires standardisation of

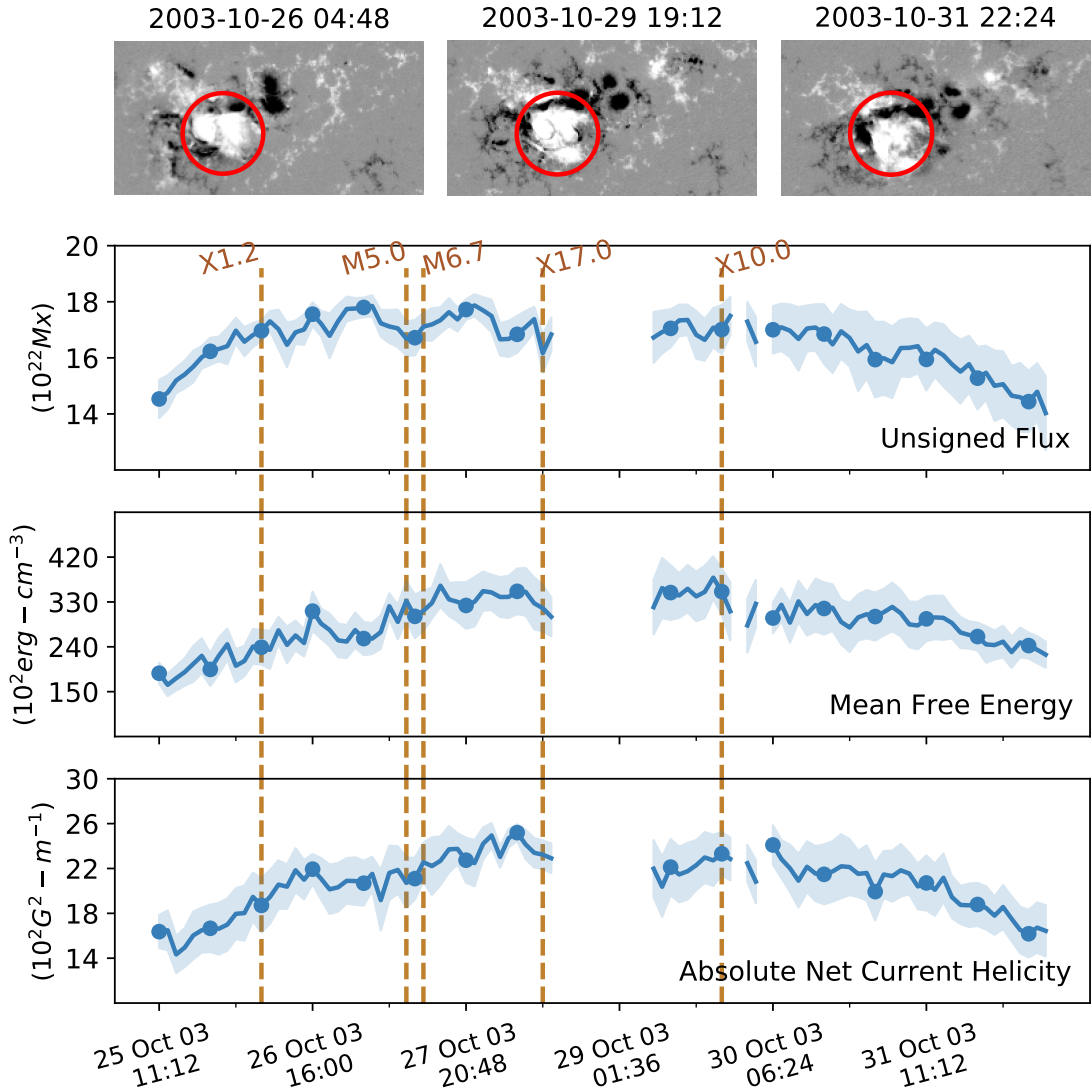


Figure 6. The 2003 Halloween Storms. (top) AR NOAA 10486 MDI magnetograms show the AR dynamics leading to extreme flare events X17.0, X10.0 and X28.0 (outside $\pm 45^\circ$ of the central meridian, not shown) characterised by rotation of the positive polarity spot (Zhang et al. 2008; Kazachenko et al. 2010), highlighted within the red circle. Positive (white) and negative (black) polarities are saturated at 1000 G. (bottom) The CNN-estimated total unsigned flux, mean-free-energy density and absolute net current helicity during these storms showing a systematic rise leading to the flares. The CNN accurately captures helicity injection due to the sunspot rotation showing a significant increase in the mean-free-energy density and electric current helicity. The gaps correspond to the missing observations and $1-\sigma$ error bars are shown. The plots are smoothed with a six hourly running average.

MDI and HMI magnetograms (e.g., with other approaches such as super-resolution (Munoz-Jaramillo et al. 2022)). We also used the GONG magnetograms to estimate the vector-field features during the storms using the CNN trained with GONG (see Appendix Figure 10). The values of the vector-field features estimated using GONG magnetograms are in the extreme range, as expected during the storms. However, the sensitivity of these estimated values to pre- and post-flare magnetic-field variations is lower compared to the features estimated from MDI.

4.3. Flare forecasting using CNN-derived features

The SHARPs features have been extensively used for building flare-forecasting models using ML (Bobra & Couvidat 2015; Bobra & Itonidis 2016; Nishizuka et al. 2017; Dhuri et al. 2019; Chen et al. 2019; Ahmadzadeh et al. 2021). In order to assess the utility of CNN-estimated SHARPs for flare-forecasting tasks, we compare their flare-forecasting performance to the true SHARPs. We set up the problem of forecasting M-/X-class flares with 24h warning similar to Bobra & Couvidat (2015). We use two approaches for the comparison. First, we build Linear Discriminant Analysis (LDA) classification models using one SHARPs feature at a time. This allows for direct comparing of the true and

Feature	accuracy	recall(+)	recall(-)	TSS
True SHARPs				
Total unsigned flux	87.04 ± 02.36	51.16 ± 18.79	88.99 ± 02.95	40.16 ± 17.15
Area	85.27 ± 02.16	58.78 ± 17.32	86.67 ± 02.71	45.45 ± 15.55
Total unsigned vertical current	88.19 ± 02.75	64.39 ± 07.67	89.50 ± 03.00	53.89 ± 07.14
Total unsigned current helicity	89.81 ± 02.82	66.17 ± 05.76	91.15 ± 02.99	57.32 ± 06.20
Total free energy density	89.77 ± 02.04	52.17 ± 11.96	91.84 ± 02.33	44.01 ± 11.07
Total Lorentz force	87.68 ± 02.47	52.39 ± 16.75	89.60 ± 02.90	41.99 ± 15.55
Absolute net current helicity	91.01 ± 02.17	58.29 ± 08.63	92.92 ± 02.40	51.21 ± 08.32
Sum of net current per polarity	90.96 ± 01.99	60.64 ± 08.29	92.72 ± 02.30	53.37 ± 07.56
Mean free energy density	76.56 ± 02.94	70.16 ± 04.95	76.92 ± 03.09	47.08 ± 05.98
Area with shear > 45°	67.15 ± 03.32	74.55 ± 04.39	66.72 ± 03.54	41.27 ± 05.54
Log of flux near polarity inversion line	67.57 ± 03.83	97.98 ± 01.63	65.83 ± 04.05	63.81 ± 05.06
CNN:HMI				
Total unsigned flux	85.01 ± 02.69	69.52 ± 08.66	85.87 ± 02.99	55.39 ± 07.65
Area	83.54 ± 01.84	70.88 ± 06.98	84.23 ± 02.17	55.11 ± 05.36
Total unsigned vertical current	85.51 ± 02.59	71.04 ± 07.54	86.30 ± 02.94	57.34 ± 06.02
Total unsigned current helicity	86.39 ± 02.23	69.81 ± 06.28	87.32 ± 02.46	57.13 ± 05.36
Total free energy density	88.02 ± 02.00	61.87 ± 08.08	89.48 ± 02.30	51.35 ± 06.68
Total Lorentz force	86.04 ± 02.74	68.47 ± 09.26	87.03 ± 02.94	55.50 ± 08.97
Absolute net current helicity	88.70 ± 02.03	72.11 ± 06.77	89.62 ± 02.14	61.73 ± 06.52
Sum of net current per polarity	88.82 ± 02.22	73.94 ± 08.13	89.64 ± 02.42	63.58 ± 07.50
Mean free energy density	76.33 ± 02.68	69.17 ± 05.34	76.72 ± 02.87	45.89 ± 05.40
Area with shear > 45°	71.55 ± 02.72	70.73 ± 08.44	71.57 ± 03.00	42.30 ± 07.84
Log of flux near polarity inversion line	76.62 ± 02.80	91.88 ± 03.58	75.72 ± 03.16	67.59 ± 02.28
CNN:GONG				
Total unsigned flux	87.01 ± 02.44	53.17 ± 11.56	88.87 ± 02.84	42.05 ± 10.34
Area	84.32 ± 02.57	60.14 ± 12.87	85.62 ± 02.98	45.76 ± 11.51
Total unsigned vertical current	87.06 ± 02.43	54.33 ± 11.78	88.87 ± 02.83	43.20 ± 10.66
Total unsigned current helicity	88.15 ± 02.42	49.19 ± 13.78	90.30 ± 02.83	39.49 ± 12.85
Total free energy density	89.44 ± 01.69	41.35 ± 17.17	92.11 ± 02.34	33.45 ± 15.85
Total Lorentz force	86.96 ± 02.39	48.25 ± 19.33	89.06 ± 02.82	37.30 ± 18.34
Absolute net current helicity	91.26 ± 02.18	54.50 ± 11.18	93.39 ± 02.34	47.89 ± 11.16
Sum of net current per polarity	90.79 ± 01.77	52.30 ± 10.90	93.02 ± 02.06	45.32 ± 10.61
Mean free energy density	77.19 ± 03.73	67.10 ± 07.69	77.76 ± 04.08	44.86 ± 07.50
Area with shear > 45°	70.50 ± 04.47	73.21 ± 07.29	70.32 ± 04.84	43.53 ± 07.77
Log of flux near polarity inversion line	78.22 ± 04.06	84.98 ± 03.89	77.82 ± 04.38	62.80 ± 04.83

Table 6. A comparison of the CNN-estimated and true SHARPs features for flare forecasting using linear discriminant analysis (LDA) of each feature. 1- σ standard deviation is shown.

the CNN-estimated values of each SHARPs feature for flare forecasting. Second, we use all SHARPs features together to train a support vector machine (SVM) for flare forecasting. We measure the flare-forecasting performance using accuracy, recall and the True Skill Statistics (TSS) score (Peirce 1884). However, only the latter two are robust to the class-imbalance prevalent in the flare-forecasting problem (Bobra & Couvidat 2015; Ahmadzadeh et al. 2021), and therefore reliable for comparison. Our definitions of positive and negative classes are identical to the operational approach de-

scribed in Bobra & Couvidat (2015). In addition, we use the 10-times repeated-holdout validation described in Section 3. Unlike Bobra & Couvidat (2015), we explicitly ensure that the samples from a given AR are not mixed in training and validation sets (Ahmadzadeh et al. 2021). Also, as mentioned in Section 2, we only consider ARs with maximum area > 25 Mm². Both the LDA and SVM are implemented using the scikit-learn library in Python.

Table 6 lists performance metrics for the classification of M-/X-class flares using the LDA of one SHARPs feature at

Flare forecasting using CNN obtained SHARPs features				
Number of observations				
# Positives	338			
# Negatives	6011			
	accuracy	recall(+)	recall(-)	TSS
True SHARPs	0.842 ± 0.030	0.856 ± 0.044	0.841 ± 0.033	0.697 ± 0.045
CNN:HMI	0.812 ± 0.028	0.869 ± 0.056	0.809 ± 0.031	0.677 ± 0.046
CNN:GONG	0.818 ± 0.031	0.801 ± 0.064	0.819 ± 0.035	0.621 ± 0.056
True SHARPs (Bobra & Couvidat 2015)	0.924 ± 0.007	0.832 ± 0.042	0.929 ± 0.008	0.761 ± 0.039

Table 7. Flare-forecasting performance of the CNN-reconstructed vector-field features. Flare-forecasting performance of a Support Vector Machine (Cortes & Vapnik 1995; Hastie et al. 2001) trained using CNN-estimated SHARPs features (Table 2). The SVM is trained to forecast M- and X-class flares 24h in advance, similar to Bobra and Couvidat (2015) (Bobra & Couvidat 2015). $1-\sigma$ standard deviation is shown.

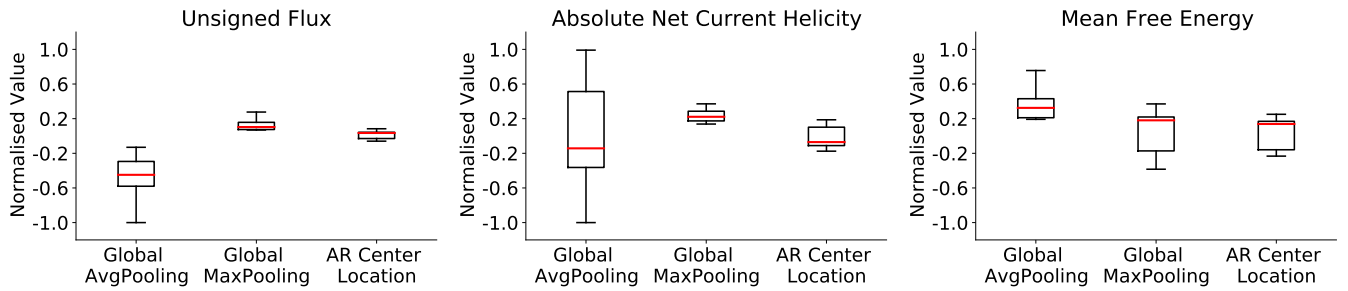


Figure 7. Interpretation of the trained CNN model. Boxplots showing weights in the penultimate layer of the CNN trained to reconstruct total unsigned flux, absolute net current helicity and mean free energy density. Top five weights from the components of the penultimate layer, i.e., global average pooling, global max-pooling and fully connected layer processing AR location (see Figure 1), are shown for the ten validation models. The red line indicates median of the weight populations.

a time. The accuracy, recall and TSS values obtained using each of the CNN-estimated features from HMI and GONG magnetograms are consistent with those of the true SHARPs features up to the validation error bars. We note that Schrijver’s R-value (Schrijver 2007) gives the highest TSS values for flare forecasting using individual features.

Table 7 lists the performance metrics for the SVM classification of M-/X-class flares using all SHARPs features together. TSS ($\sim 68\%$) and recall ($\sim 86\%$) values obtained using an SVM trained with the CNN-estimated features from HMI are consistent with those obtained using the true SHARPs. TSS ($\sim 62\%$) and recall ($\sim 80\%$) values from an SVM trained with the CNN-estimated features from GONG are slightly lower. For a comparison, we list TSS ($\sim 76\%$) and recall ($\sim 83\%$) from Bobra & Couvidat (2015) that are higher. The systematically lower TSS of the SVM in forecasting flares when using true SHARPs values here as compared with Bobra & Couvidat (2015) is due to exclusion of observations from ARs with maximum area $< 25 \text{ Mm}^2$ (all nonflaring) and the explicit restriction that samples from an AR are part of either training or validation sets. Largely consistent performance metrics for flare forecasting with the CNN-estimated SHARPs imply that high relative errors notwithstanding, the CNN-estimated features

can be useful for building space-weather forecasting tools. This is a consequence of (true) SHARPs feature values varying over several orders of magnitudes and thus being significantly different for flaring and nonflaring ARs for forecasting of flares (Dhuri et al. 2019). Accuracy of the CNN-estimated SHARPs features may be improved by significantly increasing the resolution of LOS magnetograms from, e.g., GONG, using techniques such as super-resolution (Munoz-Jaramillo et al. 2022). Our method is thus suitable for reconstructing vector-field features from historical LOS magnetograms, ultimately useful for reliable space-weather forecasting.

4.4. Interpreting the CNN

CNNs and, in general, deep learning are extremely efficient at identifying correlations in the data. In this case, the CNN builds a useful model of AR vector magnetic fields from the observed LOS magnetograms. In particular, the CNN estimated SHARPs features may be reliably used to study energy build up and time evolution of magnetic fields in flaring ARs. Yet it is very challenging to open up the trained network and understand the CNN to uncover the information absorbed. Nevertheless, weights learned by the CNN can shed some light on its working. There are also attribution methods to quantify the contribution of different parts of

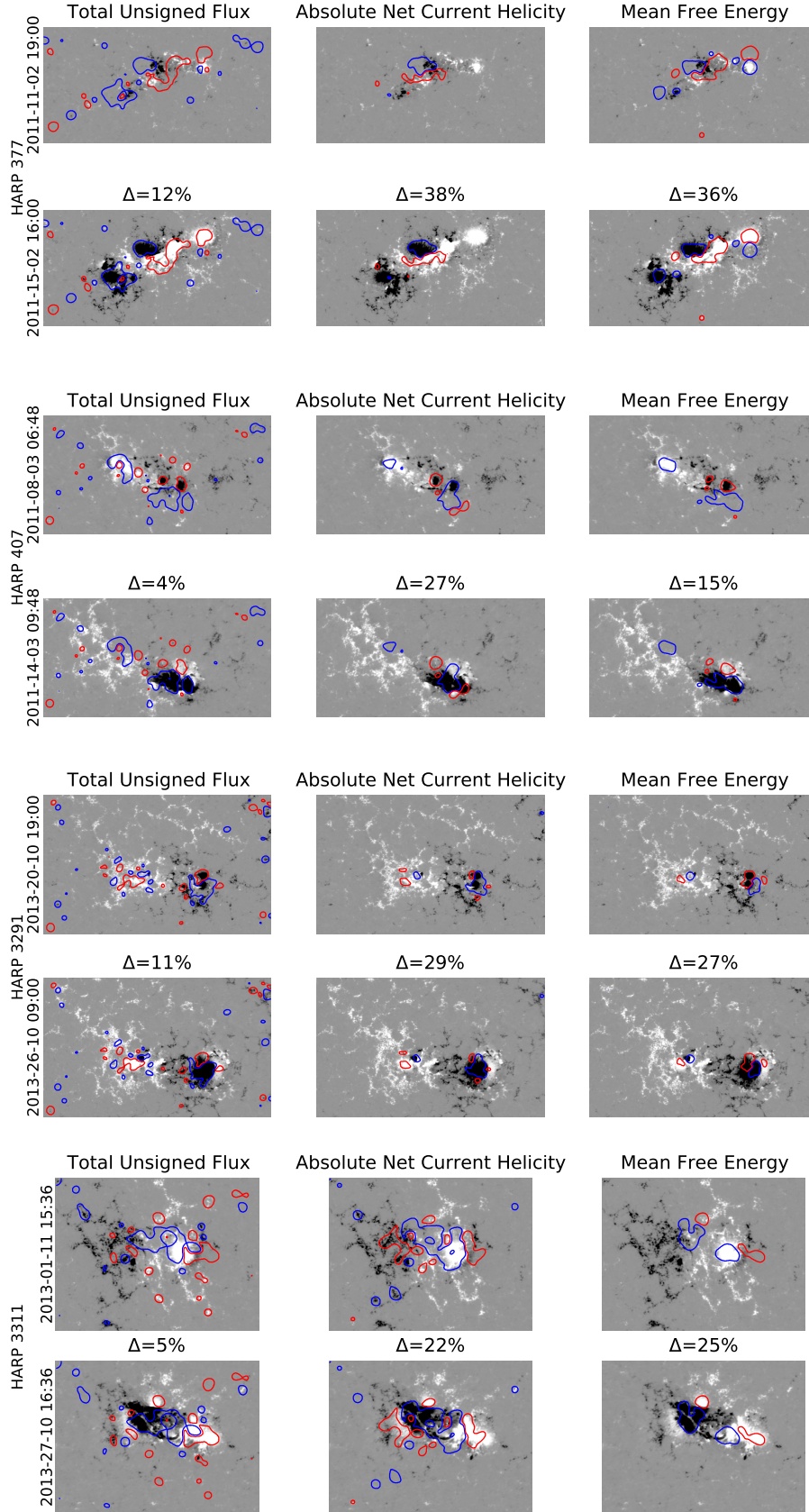


Figure 8. Integrated gradient (IG) attribution maps. The contour plots for the IG attribution maps highlight the regions from magnetograms with the highest attributions for the CNN estimation of SHARPs features — total unsigned flux, absolute net current helicity, and mean free energy density. The red/blue contours correspond to the regions with the net positive/negative attributions for the magnetograms in the bottom rows relative to the reference magnetograms in the top rows. The Δ values are the percentage change in the normalised values of the respective features between the present and the reference observations. The colorbar for the magnetograms is saturated with ± 500 G. The IG attribution maps are smoothed with a Gaussian filter of standard deviation 10 pixels before obtaining the contours.

the input image to the CNN’s output. Here, we analyse the weights of the CNN as well as obtain attribution maps for input magnetograms to interpret the trained CNN.

The CNN architecture (Figure 1) comprises a fully convolutional network for processing LOS magnetograms and a fully connected layer of neurons for processing information about the location of ARs on the solar disk. The penultimate concatenation layer comprises a global-average-pooling layer, a global-max-pooling layer that process the LOS magnetograms and a fully-connected layer that processes the location of the ARs. The global-average-pooling neurons are sensitive to the entire spatial extent of LOS magnetograms, while the global-max-pooling neurons are sensitive to spatially local patterns. The fully-connected neurons are sensitive to AR coordinates on the disk. Figure 7 illustrates the distribution of the top weights of each of the three components in the penultimate layers as their contribution to the output of the CNN that estimates total unsigned flux, absolute net current helicity and mean free energy density. Neurons associated with global-average pooling contribute dominantly to the total unsigned flux and mean free energy, implying that their estimation depends on the consideration of the entire LOS magnetograms. For absolute net current helicity, key contributors are neurons from the global max-pooling layer and its estimation is sensitive to spatially local patterns from LOS magnetograms. Without the global-max-pooling layer, absolute net current helicity and related CNN-estimated SHARPs features show $\sim 30\%$ less Pearson correlation with the true values. Weights from neurons related to AR location on the solar disk are ~ 0 , and thus, the CNN estimation does not strongly depend on the AR location. Indeed, the CNN may be trained equally well without the additional input of the AR location. This may be a consequence of considering AR patches only within $\pm 45^\circ$ where the projection effects are not significant.

While there are many attribution methods, gradient-based methods such as saliency maps (Simonyan et al. 2013), grad-CAMs (Selvaraju et al. 2017), integrated gradients (IG) (Sundararajan et al. 2017) etc., are favoured over perturbation-based methods such as occlusion masks (Zeiler & Fergus 2014) because of computational efficiency and higher resolution attribution maps. IG attribution maps are of the same resolution as the input magnetograms and are thus superior to grad-CAMs obtained from the CNN feature maps. Also, unlike saliency maps, IG attribution maps are calculated using a reference input image that facilitates assigning a cause for the attribution e.g. by comparing the magnetic-field evolution (Sun et al. 2022). Thus, here we use IG attribution maps to identify pixels, and hence the magnetic-field features in the input, that are important for the CNN output. The IG attribution map for a given input image is calculated by integrating

gradients in the CNN output along the path from a reference image. Formally,

$$L^f(x, x_0) = (x - x_0) \times \int_{\alpha=0}^1 \frac{\partial Y^f(x_0 + \alpha \times (x - x_0))}{\partial x} d\alpha, \quad (1)$$

where x_0 is the reference image, Y^f is the CNN output for SHARPs vector-field feature f .

Figure 8 shows contour plots of typical IG attribution maps for a few example magnetograms from flaring ARs (*bottom rows*). The red/blue contours include regions of net positive/negative contribution towards the CNN output. The IG attribution maps for the three SHARPs features — total unsigned flux, absolute net current helicity, and mean free energy density — are shown separately along with the reference magnetograms (*top rows*) used. In general, increasing/decreasing positive polarity flux corresponds to net positive/negative attribution. For total unsigned flux, almost all magnetic-field regions, even relatively smaller regions with weaker magnetic fields, constitute a positive/negative attribution. In contrast, for absolute net current helicity and mean free energy density, only relatively larger and stronger magnetic-field regions constitute a positive/negative attribution. For mean free energy, positive/negative attribution regions typically correspond to the uniformly increasing/decreasing positive flux. In the case of absolute net current helicity, attributions correspond to regions with “mixed” magnetic fields of the positive-negative polarities closely located. The appearance of a spurious magnetic-field polarity inversion line (PIL) is a known artifact in the line-of-sight magnetograms whenever the magnetic-field inclination relative to the line-of-sight exceeds 90° (Leka et al. 2017). We find that in many cases (e.g. HARP 407, 3291, and 3311) when the PIL artifact exists for magnetic fields within penumbrae, it wrongly constitutes an important attribution. The misattribution results from the failure of the CNN to learn the PIL artifact (Sun et al. 2022) and as a consequence, limits the accuracy of the reconstructed the vector-field-features.

5. DISCUSSION

We have thus developed a CNN model for quantifying vector-field properties — extensive features such as total unsigned flux as well as properties depending explicitly on transverse magnetic-field component such as free-energy density and current helicity — using LOS magnetograms taken from space-based HMI and ground-based GONG instruments. The CNN-estimated features strongly correlate ($> 90\%$) with their true measurements from HMI SHARPs, particularly for high-resolution LOS magnetograms from HMI. Time-evolution of the CNN-estimated features reliably mimic true AR magnetic-field evolution, particularly for ARs

producing major flares (M5 or greater). Prior to HMI, vector-magnetic-field observations available from instruments such as Imaging Vector Magnetograph and Hinode/Spectro Polarimeter (Kosugi et al. 2007) have limited spatial and temporal coverage. In contrast, near-continuous observations of LOS magnetograms are available since the 1970s from missions such as the Kitt Peak telescope (KP), MDI and GONG. LOS magnetograms from these instruments vary in their spatial resolution that are lower than HMI resolution. Nonetheless, these instruments' observation periods overlap with HMI (KP:2010-present, MDI:2010-2011, GONG:2010-present) and the attendant observations may be used to train or fine tune the CNN model to estimate SHARPs vector-field features. We explicitly show that the flare-forecasting performance of the CNN-estimated features is comparable to the true SHARPs. Therefore, vector-fields estimated from past LOS observations of nearly five decades using CNN can provide approximately four times more solar storms' data than currently available, useful for building robust statistical models for space-weather forecasting using ML. A larger sample size of solar storms also facilitates building ML algorithms based on time series of AR observations which may significantly improve forecasting performance (Dhuri et al. 2019). The CNN estimated vector-fields also provide a new perspective to understand and quantify magnetic-field dynamics during the past extreme events such as 2003 Halloween storms as demonstrated here.

Our CNN estimates are reliable for studies of solar storms, yet there is also a significant scope of improvement. Our estimates of vector-field features using HMI magnetograms are consistently more accurate compared to those estimated using lower resolution GONG magnetograms. Using LOS magnetograms from GONG and other instruments that are explicitly cross-calibrated with HMI LOS magnetograms may significantly improve accuracy of the cor-

responding vector-field instruments. Also, Deep-learning-based techniques for improving the resolution of magnetograms, namely super-resolution, are being successfully developed (Rahman et al. 2020; Munoz-Jaramillo et al. 2022). Using super-resolved LOS magnetograms as input to the CNN promises to yield more accurate CNN estimates of the vector-field features. Our estimates are also based only on the training data from the rising phase of cycle 24. Using new data available from HMI and also from newer instruments, a robust CNN regression is achievable. Extending our method, reasonable data-driven estimates of even the full photospheric vector-magnetic-field from only LOS magnetograms may be feasible, which opens up a new approach in studying and modelling AR magnetic-fields using ML.

ACKNOWLEDGMENTS

S.M.H acknowledges funding from Department of Atomic Energy grant RTI4002 and the Max-Planck Partner Group programme. D.B.D and S.M.H. acknowledge discussions with Mark C. M. Cheung and Marc DeRosa. The authors would also like to thank the anonymous reviewer and the scientific editor Manolis K. Georgoulis for their comments and suggestions that helped improve clarity of the manuscript. The authors declare that they have no competing interests. D.B.D. and S.M.H. designed the research. D.B.D., S.B. and S.K.M. analysed data. D.B.D. and S.M.H. interpreted the results. D.B.D. wrote the manuscript with contributions from S.M.H. HMI LOS and vector magnetograms, MDI LOS magnetograms and SHARPs data are publicly accessible on the JSOC data server at <http://jsoc.stanford.edu/>, courtesy the HMI and MDI science teams. The GONG LOS magnetograms are publicly available at <https://gong.nso.edu/> and were acquired by GONG instruments operated by NISP/NSO/AURA/NSF with contribution from NOAA.

REFERENCES

- Ahmadzadeh, A., Aydin, B., Georgoulis, M. K., et al. 2021, *The Astrophysical Journal Supplement Series*, 254, 23, doi: [10.3847/1538-4365/abec88](https://doi.org/10.3847/1538-4365/abec88)
- Bhattacharjee, S., Alshehhi, R., Dhuri, D. B., & Hanasoge, S. M. 2020, *The Astrophysical Journal*, 898, 98, doi: [10.3847/1538-4357/ab9c29](https://doi.org/10.3847/1538-4357/ab9c29)
- Bobra, M. G., & Couvidat, S. 2015, *The Astrophysical Journal*, 798, 135. <http://stacks.iop.org/0004-637X/798/i=2/a=135>
- Bobra, M. G., & Ilonidis, S. 2016, *The Astrophysical Journal*, 821, 127, doi: [10.3847/0004-637X/821/2/127](https://doi.org/10.3847/0004-637X/821/2/127)
- Bobra, M. G., Sun, X., Hoeksema, J. T., et al. 2014, *Solar Physics*, 289, 3549, doi: [10.1007/s11207-014-0529-3](https://doi.org/10.1007/s11207-014-0529-3)
- Bobra, M. G., Wright, P. J., Sun, X., & Turmon, M. J. 2021, *The Astrophysical Journal Supplement Series*, 256, 26, doi: [10.3847/1538-4365/ac1f1d](https://doi.org/10.3847/1538-4365/ac1f1d)
- Boteler, D. H. 2019, *Space Weather*, 17, 1427, doi: [10.1029/2019SW002278](https://doi.org/10.1029/2019SW002278)
- Bottou, L. 1991, in *Proceedings of Neuro-Nimes*, 687–706
- Chen, Y., Manchester, W. B., Hero, A. O., et al. 2019, *Space Weather*, 17, 1404, doi: [10.1029/2019SW002214](https://doi.org/10.1029/2019SW002214)
- Cheung, M. C. M., & Isobe, H. 2014, *Living Reviews in Solar Physics*, 11, 3, doi: [10.12942/lrsp-2014-3](https://doi.org/10.12942/lrsp-2014-3)
- Cortes, C., & Vapnik, V. 1995, *Machine learning*, 20, 273
- Crown, M. D. 2012, *Space Weather*, 10, 1, doi: [10.1029/2011SW000760](https://doi.org/10.1029/2011SW000760)

- Dhuri, D. B., Hanasoge, S. M., & Cheung, M. C. M. 2019, Proceedings of the National Academy of Sciences, 116, 11141, doi: [10.1073/pnas.1820244116](https://doi.org/10.1073/pnas.1820244116)
- Eastwood, J. P., Biffis, E., Hapgood, M. A., et al. 2017, Risk Analysis, 37, 206, doi: [10.1111/risa.12765](https://doi.org/10.1111/risa.12765)
- Goodfellow, I., Bengio, Y., & Courville, A. 2016, Deep Learning (The MIT Press)
- Han, J., & Moraga, C. 1995, in From Natural to Artificial Neural Computation, ed. J. Mira & F. Sandoval (Berlin, Heidelberg: Springer Berlin Heidelberg), 195–201
- Hastie, T., Tibshirani, R., & Friedman, J. 2001, The Elements of Statistical Learning, Springer Series in Statistics (New York, NY, USA: Springer New York Inc.)
- Hoeksema, J. T., Liu, Y., Hayashi, K., et al. 2014, Solar Physics, 289, 3483, doi: [10.1007/s11207-014-0516-8](https://doi.org/10.1007/s11207-014-0516-8)
- Kazachenko, M. D., Canfield, R. C., Longcope, D. W., & Qiu, J. 2010, The Astrophysical Journal, 722, 1539, doi: [10.1088/0004-637x/722/2/1539](https://doi.org/10.1088/0004-637x/722/2/1539)
- Kosugi, T., Matsuzaki, K., Sakao, T., et al. 2007, Solar Physics, 243, 3, doi: [10.1007/s11207-007-9014-6](https://doi.org/10.1007/s11207-007-9014-6)
- LeCun, Y., Bengio, Y., & Hinton, G. 2015, Nature, 521, 436 EP . <https://doi.org/10.1038/nature14539>
- Leka, K. D., & Barnes, G. 2007, The Astrophysical Journal, 656, 1173. <http://stacks.iop.org/0004-637X/656/i=2/a=1173>
- Leka, K. D., Barnes, G., & Wagner, E. L. 2017, Solar Physics, 292, 36, doi: [10.1007/s11207-017-1057-8](https://doi.org/10.1007/s11207-017-1057-8)
- Livingston, W. C., Harvey, J., Pierce, A. K., et al. 1976, Appl. Opt., 15, 33, doi: [10.1364/AO.15.000033](https://doi.org/10.1364/AO.15.000033)
- McIntosh, P. S. 1990, Solar Physics, 125, 251, doi: [10.1007/BF00158405](https://doi.org/10.1007/BF00158405)
- Metcalf, T. R., Leka, K. D., & Mickey, D. L. 2005, The Astrophysical Journal, 623, L53, doi: [10.1086/429961](https://doi.org/10.1086/429961)
- Munoz-Jaramillo, A., Jungbluth, A., Gitiaux, X., et al. 2022, Nature Portfolio, doi: [10.21203/rs.3.rs-713430/v1](https://doi.org/10.21203/rs.3.rs-713430/v1)
- Nishizuka, N., Sugiura, K., Kubo, Y., et al. 2017, The Astrophysical Journal, 835, 156. <http://stacks.iop.org/0004-637X/835/i=2/a=156>
- Peirce, C. S. 1884, Science, ns-4, 453, doi: [10.1126/science.ns-4.93.453-a](https://doi.org/10.1126/science.ns-4.93.453-a)
- Pesnell, W. D., Thompson, B. J., & Chamberlin, P. C. 2012, Solar Physics, 275, 3, doi: [10.1007/s11207-011-9841-3](https://doi.org/10.1007/s11207-011-9841-3)
- Pulkkinen, A., Lindahl, S., Viljanen, A., & Pirjola, R. 2005, Space Weather, 3, 1, doi: [10.1029/2004SW000123](https://doi.org/10.1029/2004SW000123)
- Rahman, S., Moon, Y.-J., Park, E., et al. 2020, The Astrophysical Journal, 897, L32, doi: [10.3847/2041-8213/ab9d79](https://doi.org/10.3847/2041-8213/ab9d79)
- Régnier, S., & Priest, E. R. 2007, The Astrophysical Journal, 669, L53, doi: [10.1086/523269](https://doi.org/10.1086/523269)
- Scherrer, P. H., Bogart, R. S., Bush, R. I., et al. 1995, Solar Physics, 162, 129, doi: [10.1007/BF00733429](https://doi.org/10.1007/BF00733429)
- Schrijver, C. J. 2007, The Astrophysical Journal Letters, 655, L117. <https://iopscience.iop.org/article/10.1086/511857>
- Selvaraju, R. R., Cogswell, M., Das, A., et al. 2017, 2017 IEEE International Conference on Computer Vision (ICCV), 1, 618
- Shibata, K., & Magara, T. 2011, Living Reviews in Solar Physics, 8, 6, doi: [10.12942/lrsp-2011-6](https://doi.org/10.12942/lrsp-2011-6)
- Simonyan, K., Vedaldi, A., & Zisserman, A. 2013, CoRR, abs/1312.6034, 1
- Stenflo, J. O. 2013, The Astronomy and Astrophysics Review, 21, 66, doi: [10.1007/s00159-013-0066-3](https://doi.org/10.1007/s00159-013-0066-3)
- Su, Y., Veronig, A. M., Holman, G. D., et al. 2013, Nature Physics, 9, 489 EP . <http://dx.doi.org/10.1038/nphys2675>
- Sun, Z., Bobra, M. G., Wang, X., et al. 2022, The Astrophysical Journal, 931, 163, doi: [10.3847/1538-4357/ac64a6](https://doi.org/10.3847/1538-4357/ac64a6)
- Sundararajan, M., Taly, A., & Yan, Q. 2017, in Proceedings of Machine Learning Research, Vol. 70, Proceedings of the 34th International Conference on Machine Learning, ed. D. Precup & Y. W. Teh (PMLR), 3319–3328. <https://proceedings.mlr.press/v70/sundararajan17a.html>
- Szegedy, C., Wei Liu, Yangqing Jia, et al. 2015, in 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1–9, doi: [10.1109/CVPR.2015.7298594](https://doi.org/10.1109/CVPR.2015.7298594)
- Toriumi, S., & Wang, H. 2019, Living Reviews in Solar Physics, 16, 3, doi: [10.1007/s41116-019-0019-7](https://doi.org/10.1007/s41116-019-0019-7)
- Zeiler, M., & Fergus, R. 2014, in Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), Vol. 8689 LNCS, Computer Vision, ECCV 2014 - 13th European Conference, Proceedings, part 1 edn. (Springer Verlag), 818–833, doi: [10.1007/978-3-319-10590-1_53](https://doi.org/10.1007/978-3-319-10590-1_53)
- Zhang, Y., Liu, J., & Zhang, H. 2008, Solar Physics, 247, 39, doi: [10.1007/s11207-007-9089-0](https://doi.org/10.1007/s11207-007-9089-0)

APPENDIX

A. TIME EVOLUTION OF THE CNN-ESTIMATED FEATURES ON ALL ARS PRODUCING FLARES M5 OR GREATER.

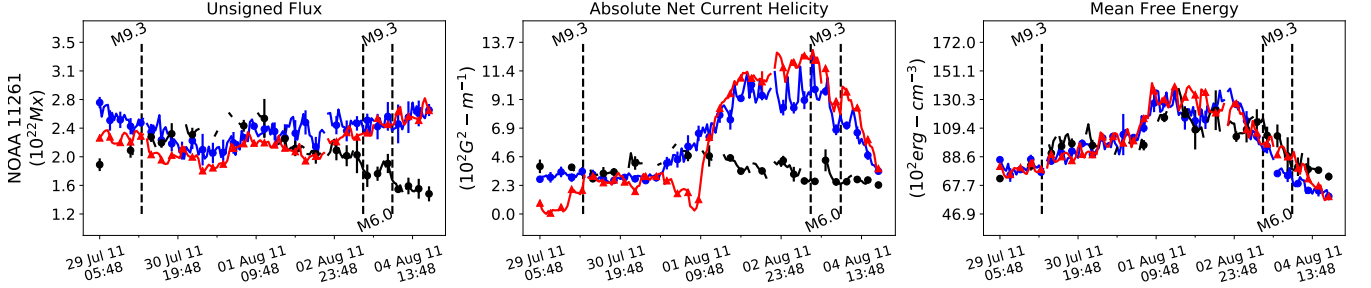


Figure 9. Time-evolution of the CNN-estimated vector-field features on flare-productive active regions. Comparison of the CNN-estimated values, using HMI (blue) and GONG (black) LOS magnetograms, of total unsigned flux, absolute net current helicity and mean free energy density with true values (red) calculated from HMI vector magnetograms are shown for HARP 750. Only observations within $\pm 45^\circ$ are considered. $1-\sigma$ error bars are shown. The gaps indicate missing observations. The complete figure set (28 images) for active regions that produced at least one M5 or greater flare is available in the online journal.

B. COMPARISON OF THE CNN-ESTIMATED FEATURES DURING THE 2003 HALLOWEEN STORMS

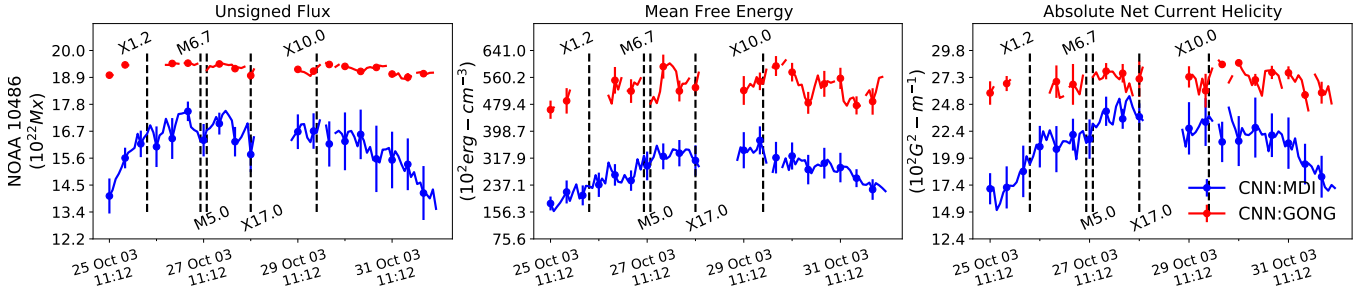


Figure 10. A comparison of the time-evolution of the CNN-estimated vector-field features during the 2003 Halloween storms using MDI and GONG LOS magnetograms. The CNN trained with HMI magnetograms is used for the estimation from MDI magnetograms whereas the CNN trained with GONG magnetograms is used for the estimation from GONG magnetograms. Note that the HMI observations are not available before 2010. The CNN:GONG feature values are generally high compared to the CNN:MDI, showing little variation throughout the storms. The CNN:MDI features appear to capture the variation of these features during the storms expected from the theoretical modelling e.g. (Kazachenko et al. 2010). The $1-\sigma$ errors are shown. The gaps indicate missing observations. The legend in the rightmost panel applies to all panels.

C. MDI CORRELATIONS

SHARPs Features	Pearson correlation	Spearman correlation
Total unsigned flux	82.49 ± 21.01	68.09 ± 37.73
Area	85.55 ± 23.82	74.59 ± 17.72
Total unsigned vertical current	74.83 ± 45.39	64.60 ± 46.31
Total unsigned current helicity	75.22 ± 45.28	65.80 ± 47.14
Total free energy density	84.18 ± 12.15	76.40 ± 19.94
Total Lorentz force	89.13 ± 11.33	75.37 ± 29.10
Absolute net current helicity	51.62 ± 27.28	48.97 ± 23.21
Sum of net current per polarity	42.84 ± 35.65	38.69 ± 32.27
Mean free energy density	92.60 ± 04.02	89.26 ± 07.80
Area with shear $> 45^\circ$	91.78 ± 03.32	89.63 ± 03.32
Flux near polarity inversion line	62.64 ± 14.63	59.09 ± 21.23

Table 8. Pearson and Spearman correlations between the CNN-estimated vector-field features SHARPs using MDI line-of-sight magnetograms and their true values. The SHARPs features are estimated using the CNN trained with the HMI line-of-sight magnetograms. The AR patches of MDI line-of-sight magnetograms are taken from the publicly available data product Space-Weather MDI Active Region Patches (SMARPs) (Bobra et al. 2021). SMARPs and SHARPs data overlap between 1 May 2010 and 28 October 2010 (Bobra et al. 2021).

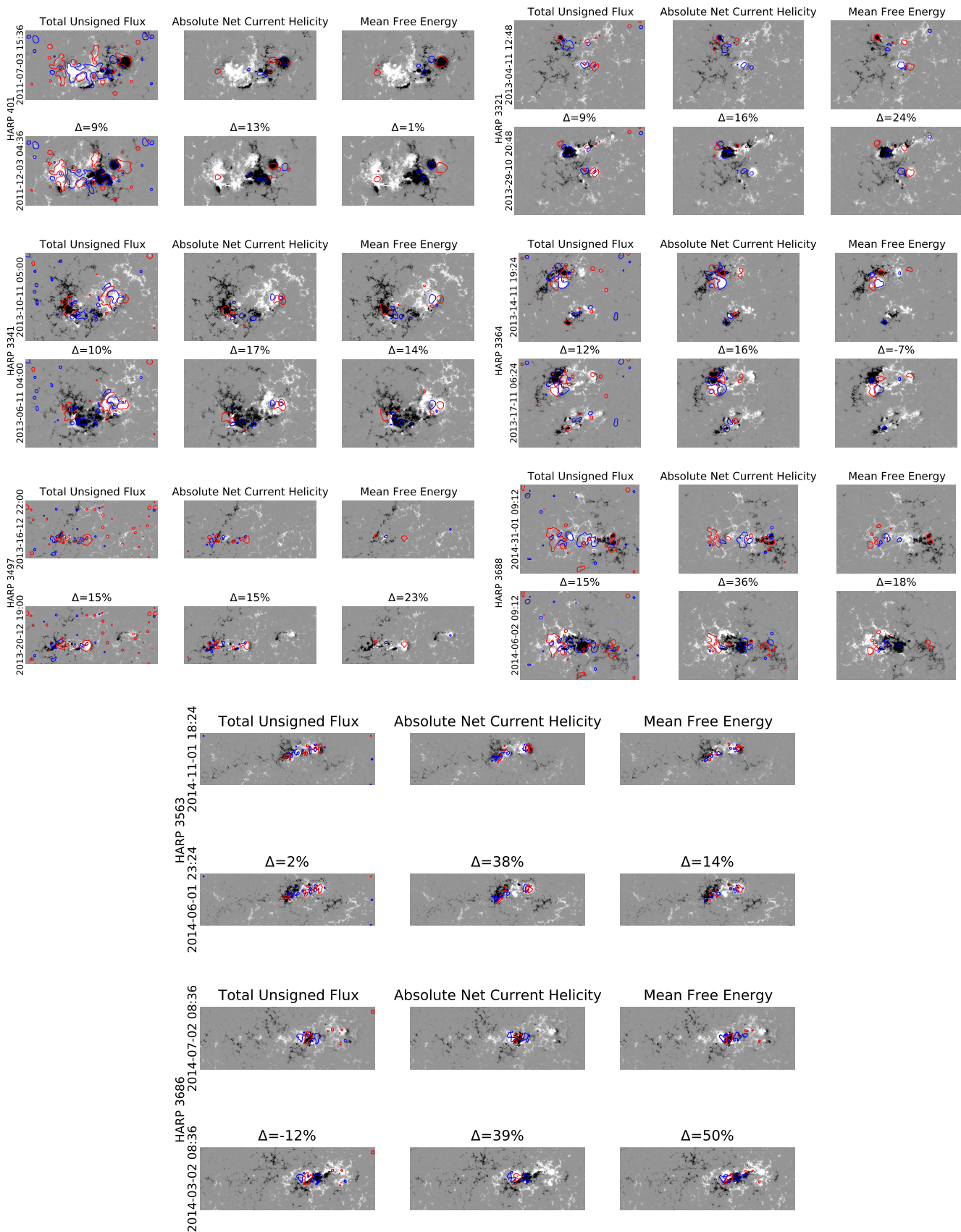


Figure 11. Additional examples of Integrated gradient (IG) attribution maps (Figure 8).