



# MSMT-CNN for Solar Active Region Detection with Multi-Spectral Analysis

Majedaldein Almahasneh<sup>1</sup> · Adeline Paiement<sup>2</sup> · Xianghua Xie<sup>1</sup> · Jean Aboudarham<sup>3</sup>

Received: 13 May 2021 / Accepted: 11 March 2022  
© The Author(s) 2022

## Abstract

Precisely detecting solar active regions (AR) from multi-spectral images is a challenging task yet important in understanding solar activity and its influence on space weather. A main challenge comes from each modality capturing a different location of these 3D objects, as opposed to more traditional multi-spectral imaging scenarios where all image bands observe the same scene. We present a multi-task deep learning framework that exploits the dependencies between image bands to produce 3D AR detection where different image bands (and physical locations) each have their own set of results. Different feature fusion strategies are investigated in this work, where information from different image modalities is aggregated at different semantic levels throughout the network. This allows the network to benefit from the joint analysis while preserving the band-specific information. We compare our detection method against baseline approaches for solar image analysis (multi-channel coronal hole detection, SPOCA for ARs (Verbeeck et al. *Astron Astrophys* 561:16, 2013)) and a state-of-the-art deep learning method (Faster RCNN) and show enhanced performances in detecting ARs jointly from multiple bands. We also evaluate our proposed approach on synthetic data of similar spatial configurations obtained from annotated multi-modal magnetic resonance images.

**Keywords** Object detection · Solar images · Active regions · Multi-spectral images · Deep neural networks

## Introduction

Active regions (ARs) detection is essential in studying solar behaviors and space weather. The solar atmosphere is monitored on multiple wavelengths, as seen in Fig. 1. However, unlike traditional multi-spectral scenarios such as Earth imaging, e.g. [11, 17, 34], where multiple imaging bands reveal different aspects (e.g. composition) of a same scene, different bands image the solar atmosphere at different temperatures, which correspond to different altitudes

[24]. Therefore, imaging the sun using different wavelengths shows different 2D cuts of the 3D objects that span the solar atmosphere. This makes handling the multi-spectral nature of the data not straightforward. Moreover, the variety of shapes and brightness, and fuzzy boundaries, of ARs also introduce a high complexity in precisely localising them.

Very few solutions were presented to the AR detection problem. Most of these methods exploited single image bands only. Benkhalil et al. [2] proposed a method for single-band images from Paris-Meudon Spectroheliograph (PM/SH) and SOHO/EIT. In [24], ARs were segmented from a single band at a time, which the authors justify by the fact that they each provide information from a different solar altitude, and they showed how the area of ARs differs between the different bands. While we also aim at getting specialised results for each image band, we argue that inter-dependencies exist between bands, which can be exploited for increased robustness. The SPOCA method [32], used in the Heliophysics Feature Catalogue (HFC),<sup>1</sup> segments ARs and coronal holes from SOHO/EIT 171 Å and 195 Å combined images. These two bands image overlapping (but

---

This article is part of the topical collection “Pattern Recognition Applications and Methods” guest edited by Ana Fred, Maria De Marsico and Gabriella Sanniti di Baja.

---

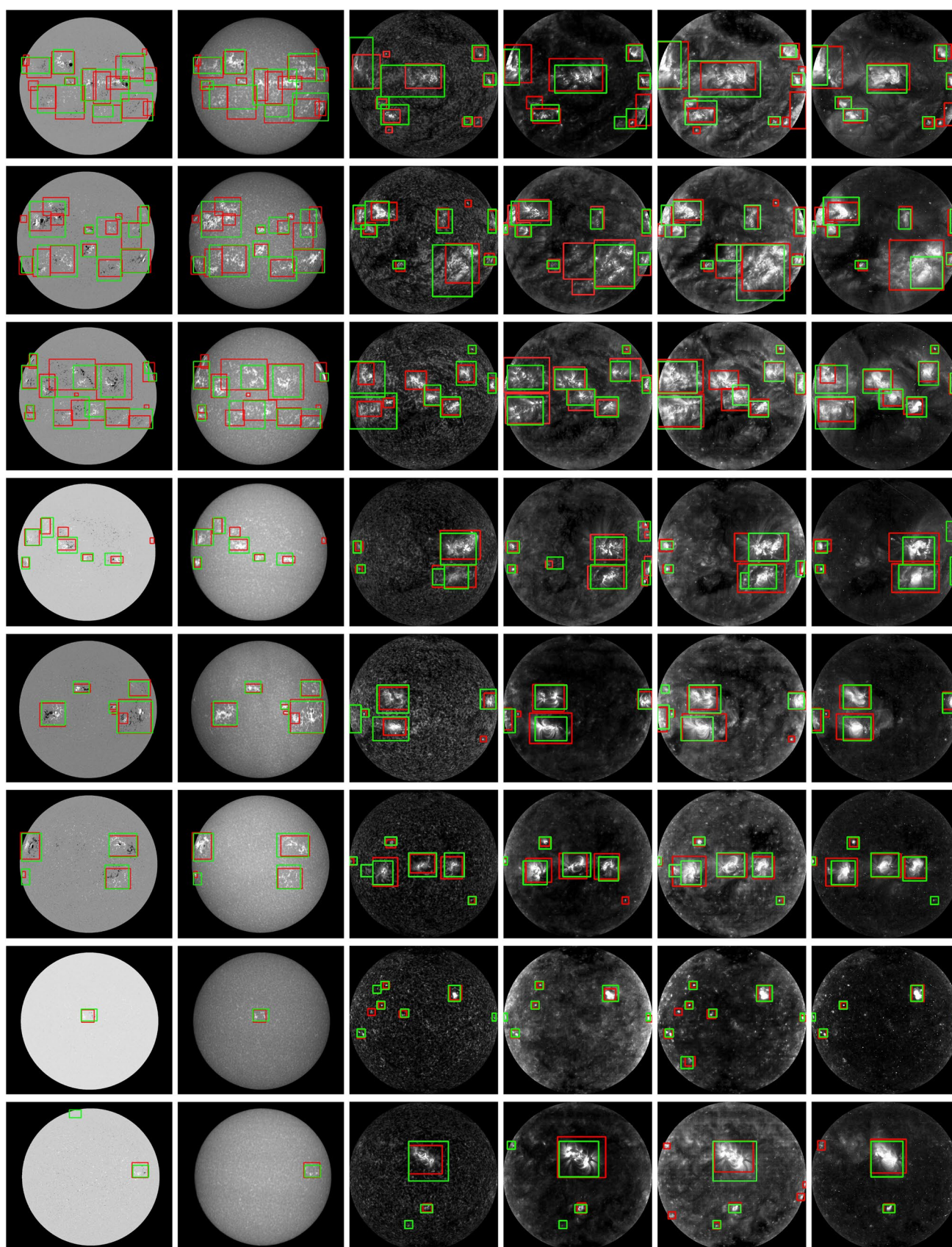
✉ Majedaldein Almahasneh  
809508@swansea.ac.uk

<sup>1</sup> Department of Computer Science, Swansea University, Swansea, UK

<sup>2</sup> Université de Toulon, Aix Marseille Univ, CNRS, LIS, Marseille, France

<sup>3</sup> Observatoire de Paris/PSL, Paris, France

<sup>1</sup> <http://voparis-helio.obspm.fr/hfc-gui/>.



**Fig. 1** Ground-truth (red) and MSMT-CNN's (green) detection of ARs in randomly selected images from (left to right) SOHO/MDI Magnetogram and PM/SH 3934 Å, SOHO/EIT 304, 171, 195, and 284 Å. Contrast has been increased for convenience of visualisation

different) regions of the solar atmosphere. SPOCA considers that they should yield identical detections. This approximation may result in a bad analysis of at least one of these bands. We provide separate but related results for these bands. We also exploit more bands for richer information on the solar atmosphere, with separate results for each band.

SPOCA's segmentation is based on clustering, with Fuzzy and Possibilistic C-means [13], followed by morphological operations. The method of [2] uses local thresholding and morphological operations followed by region growing. The method was evaluated against manual detections (synoptic maps) produced at PM and National Oceanic and Atmospheric Administration (NOAA), and detected similar numbers of ARs as PM, and  $\sim 50\%$  more than NOAA. In [24], ARs were segmented by computing the pixel-wise fractal dimension (a measure of non-linear growth that reflects the degree of irregularity over multiple scales) in a convolutional fashion, and feeding the resulting feature map to a Fuzzy C-means [3] algorithm. Overall, these methods, mainly based on clustering and morphological operations, are very pre- and post-processing dependant. This makes them difficult to adapt to new image domains. We address this limitation using deep learning (DL).

Object detection has evolved dramatically in the past two decades, from hand-crafted features based detection (e.g. Haar [33], and HOG [6]) to deep neural networks (DNN) such as YOLO [22], SSD [15], R-FCN [5], Cornernet [14], and Faster RCNN [23].

Deep learning methods generally aim at analysing 2D images or dense 3D volumes, while the sparse 3D nature of the solar imaging data requires designing a specialised DL framework. Furthermore, multi-spectral images are commonly treated in a similar manner to RGB images, by stacking different bands into composite multi-channel images, extracting a common feature map, and producing a single detection result for the composite image, e.g. [8, 9, 11, 18, 19, 21]. This multi-channel strategy is ill-suited to the solar imaging scenario, since different images show different scenes and should have their own detection results.

Another common approach is to aggregate information from different bands at different levels, e.g. feature level and image level, [7, 9, 10, 28, 29, 31, 34]. This feature fusion strategy demonstrates the potential for DNNs to improve localisation by exploiting the multi-spectral aspect of the data.

In [34], a feature-fusion approach was proposed where HOG features extracted separately from the RGB+thermal images were concatenated before performing the final analysis by a fully connected layer. This strategy obtained better results than the previously mentioned image-level fusion. Authors discussed that the network may optimise the learned features for each band. Moreover, they reckon that small misalignments may be overcome as spatial

information gets less relevant in late network stages. This may be an advantage in our case of images showing different parts of a scene.

However, when comparing image-level and feature-level fusion, [9] found on the contrary that image fusion worked best when segmenting soft tissue sarcomas in multi-modal medical images. These different results suggest that there is no universal best fusion strategy, and it needs to be adapted to each case. In our detection scenario, we investigate the best stages where to apply fusion.

Another feature-fusion strategy was used in [12] to segment coronal holes from 7 SDO bands and a magnetogram. The method relies on training a CNN to segment coronal holes from a single band, followed by fine-tuning the learned CNN over the other bands consecutively. The feature maps of each specialised CNN are used in combination as input to a final segmentation CNN, resulting in a unique final prediction. The production of a unique localisation result for all multi-spectral images is a common limitation to all cited works, which we address in this study with a multi-task network. In this work, continuing from our previous work [1] on AR detection, we introduce MultiSpectral-MultiTask-CNN (MSMT-CNN), a multi-tasking DNN framework, as a robust solution for solar AR detection that takes into consideration the multi-spectral aspect of the data and the three-dimensional spatial dependencies between image bands. This multi-spectral and multi-tasking concept may be applied to any CNN backbone.

The 3D nature of our multi-spectral imaging scenario, which differs from previous multi-spectral applications, requires a new benchmark. We introduce two annotated datasets comprised of solar images from both ground and space, and which cover evenly all phases of solar activity, which follows an 11-year cycle. To the best of our knowledge, no detection ground-truth was previously available for such data. A labelling tool was hence designed to cope with its temporal and multi-spectral nature and will be also released.

## Methodology

While some existing works were developed for analysing multi-spectral images, to our best knowledge, the problem of detecting objects over sparse 3D multi-spectral imagery, in which different bands show different scenes, was not yet addressed. Our framework exploits jointly several time-matched image bands in parallel, to predict separate, although related, detection results for each image. This framework is general and may be used with any DNN backbone, we demonstrate it using Faster RCNN [23].

The intuition behind our framework manifests in three key principles:



1. Extracting features from different image bands individually using parallel feature extraction branches. This allows the network to learn independent features from each band, according to their specific modality.
2. Aggregating the learned features from the different branches using some appropriate fusion operator. This assists the network to jointly analyse the extracted features from different bands and thus learn their interdependencies. In this work, we test fusion by addition and concatenation, at different feature levels (i.e. early and late fusion).
3. Generating a set of results per image band, based on a multi-task loss, allowing the detection of different sections or layers of 3D objects.

Points 1 and 3 are motivated by the nature of the multi-spectral data, where different bands image different locations in a 3D scene, each providing a unique information. Our multi-tasking framework aims at getting specialised results for each image band, in contrast to most existing works where focus is on producing a unique prediction to all image bands. This is crucial since the localisation information may differ from one band to another in solar (sparse) multi-spectral images. Yet, all bands are correlated, which motivates point 2. Our framework exploits the inter-dependencies between the different bands by its joint analysis strategy, increasing the robustness of its performance in individual bands.

Furthermore, our framework emulates how experts manually detect ARs (see also Sect. “Data”), where a suspected region’s correlation with other bands is evaluated prior to its final classification. This demonstrates the usefulness and importance of accounting for (spatially and temporally) neighbouring slices in robustly detecting ARs.

The MSMT-CNN framework is very modular and flexible. It may accommodate any number of available multi-spectral images. Additionally, since different scenarios may require different fusion strategies (as suggested by existing works), the modularity of our framework allows it to be easily adapted to different types and levels of feature fusion (e.g. addition and concatenation, early and late). This modular design also allows our framework to adopt different backbone architectures (e.g. Faster RCNN in our experiments). Indeed, its three key principles are applicable to any backbone, as they are not architecture dependent.

## Pre-processing

The input of our system are time-matched observations, possibly acquired by different instruments or at different orientations of the same instrument. As such they need to be spatially aligned. We harmonise the radius and center location of the solar disk, either using SOHO image preparation routines, or through Otsu [20] thresholding of the solar

disc of PM/SH images followed by minimum enclosing circle fitting and re-projection into a unified center and radius. Orientation is normalised by SOHO and PM routines to a vertical north-south solar axis. Although this process does not correct a possible small time difference and resulting east-west rotation of the Sun between two acquisitions, it ensures a sufficient alignment for our purpose of AR detections from spatially (and temporally) correlated solar disks.

The SOHO/EIT images are prepared by EIT routines. We eliminate any prominences or solar eruptions by masking out all areas outside the solar disk. The contrast of SOHO/MDI Magnetograms is enhanced by intensity rescaling. Contrast enhancement was not used on SOHO/EIT and PM/SH images, as it was found to have minimal effect on our detection results.

Both datasets are augmented using north-south flipping, east-west flipping, and a combination of the two. Augmentation with arbitrary rotations of the images is a popular way of augmenting astronomy datasets. However, such rotations are ruled out from our study because ARs tend to appear predominantly alongside the solar equator.

Finally, a single-channel solar image was repeated along the depth axis resulting in a 3-channel image to match the pre-trained CNN’s input depth.

## MSMT-CNN

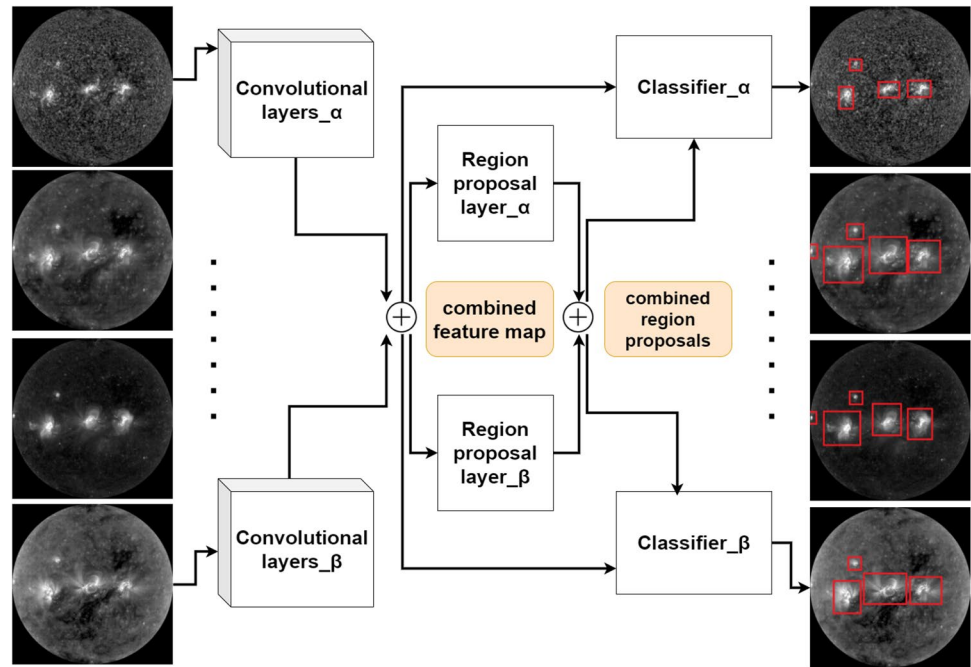
Our detection DNN is presented in Fig. 2. A CNN (ResNet50 or VGG16 in our experiments) is first used as a feature extraction network. Parallel branches (subnetworks) produce a feature map per image band, following the late feature fusion strategy. This allows the subnetworks’ filters to be optimised for their input bands individually. The feature maps are then concatenated across the bands.

The combined feature map is jointly analysed by one parallel module per image band that performs Faster RCNN’s region proposal network (RPN). The RPN stage uses three aspect ratios ([1:1], [1:2], [2:1]) and four sizes of anchor (32, 64, 128, and 256 pixel width), found empirically to match well the typical size and shape of ARs. One specialised RPN per image band is trained.

At training, for each band, the correspondent region proposals along with the combined feature map are used by a Faster RCNN’s detector module to perform the final prediction for the band. However, at testing time, the band-specialised detector modules use the region proposals from all bands. This combination of region proposals helps finding potential AR locations in bands where they are more difficult to identify.

It is good to note that during training, the RPN’s proposals for a band are filtered (i.e. labeled as positive or negative) with respect to their overlap with the band’s ground-truth. Hence, combining them in the training time would mean

**Fig. 2** MSMT for detection using the Faster-RCNN backbone. ‘Plus’ sign denotes concatenation of the feature maps, or of the lists of region proposals (at testing time). Each image band is analysed independently using a band-specific convolutional branch to extract band-specific features. These are then fused and are jointly analysed by band-specific RPN modules such that each RPN produces region proposal for its correspondent band. Region proposals from different RPNs are then aggregated and passed onto band-specific detection heads, where each band gets a separate (but related) set of predictions



implicitly inheriting the ground-truth of a band to another, in contradiction with the band-specific ground-truth used for training the detector module. This may hinder the learning of both the RPN and detector modules. Therefore, region proposals are only combined at testing time to ensure a better learning of the final detection modules.

Using the combined feature map for both RPN prediction and classification helps the network learn the relationship between the image bands and hence provide more consistent region proposals and final predictions. We demonstrate in Sect. “Experiments” that this is particularly helpful in cases where an AR is difficult to detect in a single band.

The network is trained in the same way as the original Faster RCNN, using all input bands and branches according to a combined loss function:

$$L = \sum_b \left( \frac{1}{N_{\text{cls}}} \sum_i L_{\text{cls}}(p_{b_i}, p_{b_i}^*) + \lambda \frac{1}{N_{\text{reg}}} \sum_i p_{b_i}^* L_{\text{reg}}(t_{b_i}, t_{b_i}^*) \right) \quad (1)$$

where  $b$  and  $i$  refer to the image band and the index of the bounding box being processed, respectively. The terms  $L_{\text{cls}}$  and  $L_{\text{reg}}$  are the bounding-box classification loss and the bounding-box regression loss defined in [23].  $N_{\text{cls}}$  and  $N_{\text{reg}}$  represent the size of the mini batch being processed and the number of anchors, respectively.  $\lambda$  balances the classification and the regression losses (we set  $\lambda$  to 10 as suggested in [23]).  $p$  and  $p^*$  are the predicted anchor’s class probability and its actual label, respectively. Lastly,  $t$  and  $t^*$  represent the predicted bounding box coordinates and the ground-truth

coordinates, respectively. It is worth noting that our proposed framework is not limited to using Faster RCNN’s loss and may be trained with using other task-suitable loss functions.

During training, the weights of each stage (i.e. feature extraction, region proposal, and detection) are stored independently whenever the related Faster RCNN loss decreases. At testing time, the best performing set of weights is retrieved per stage. We refer to this practice as ‘Multi-Objective Optimisation’ (MOO). The improved performance that we observe in Sect. “Experiments” may be explained by each stage having a different objective to optimise, which may be reached at different times.

In this paper, we experiment with a 2, 3, and 4-band pipeline. However, the approach may generalise straightforwardly to  $n$  bands and new imaging modalities. Similarly, our framework may exploit any DNN architecture, and may be updated with new state-of-the-art DL architectures easily.

## Experiments

Our framework was implemented using Tensorflow and run on an NVIDIA GeForce GTX 1080 Ti. We evaluate our detection stage using precision, recall, and F1-score. Solar ARs are dynamic structures that are constantly changing (e.g. merging and splitting, from and to multiple regions) during their lifetime [26, 27, 30, 32]. Accordingly, we design a dynamic criterion to match the characteristics of AR. Particularly, we address two prevalent scenarios, (1) an AR area is detected by multiple neighbouring bounding boxes, or (2)

a cluster of neighboring ARs is detected in a single bounding box. See Fig. 3. Therefore, a detection is considered a true positive if its intersection with a ground-truth box is greater or equal to 50% of either the predicted or ground-truth area, and is of an area that lies within the AR area distribution of the annotated dataset. We empirically found that this provides a good trade-off of precision over recall. Additionally, we evaluate our detection using IoU (intersection over union) thresholds of 25% and 50%, and compare it to baseline methods. Note that both our annotation and evaluation processes were validated by a solar physics expert. Non maximum suppression (NMS) is used to discard any redundant detections.

All tested CNN were initialised with pre-trained ImageNet [25] weights. Indeed, we demonstrated in [4] that CNNs pre-trained on RGB images may fine-tune and adapt well to other modalities such as depth images, provided that the image's gain and contrast are suitably enhanced to match those of the pre-training RGB images.

## Data

We consider data from SOHO (space-based) and PM observatory (ground-based). We use the bands 171 and 195 Å (transition region), 284 Å (corona), and 304 Å (chromosphere and base of transition region) from SOHO/EIT, 3934 Å (chromosphere) from PM/SH, and line-of-sight magnetograms (photosphere) from SOHO/MDI, as illustrated in Fig. 1. To account for the regular solar cycle, we select images evenly from each of three periods of varying solar activity level, with years 2002–03, 2004–05, and 2008–10 for high, medium, and low activity respectively.

We publish two new datasets with detection annotations (i.e bounding boxes): the Lower Atmosphere Dataset (LAD) and Upper Atmosphere Dataset (UAD). All annotations were validated by a solar physics expert.

Localising ARs in multi-spectral images can be challenging due to inconsistent numbers of polarity centers, shapes, sizes, and activity levels. To mitigate this issue, when manually annotating, we exploit neighboring image bands, magnetograms, as well as temporal information where we examine the evolution of a suspected AR to validate its detection and localisation. We designed a new multi-spectral labelling tool<sup>2</sup> which displays, side by side, images from an auxiliary modality and from a sequence of three previous and three subsequent time steps.

Auxiliary imaging modalities may have different observation frequencies and times, therefore we work with *time-matched* images, i.e. the time-closest image, if any, in a 12-h

window during which ARs may not undergo any significant change.

ARs have a high spatial coherence in 3934 Å and magnetogram images due to the physical proximity of the two imaged regions. Hence, when annotating 3934 Å images, time-matched magnetograms were used as an extra support. Furthermore, the 3934 Å bounding boxes could be considered to be good approximations of magnetograms' annotations. Table 1 presents an overview of annotated images for both datasets. We split the datasets into training and testing sets in the following proportions. For LAD, we use 213 images (1380 bounding box) for training, and 53 images (406 bounding box) for testing. For UAD, we use 283 images for training, and 40 images for testing. This amounts to 2205, 1919, 2341, and 2016 training bounding boxes in the 304, 171, 195, and 284 Å bands respectively, and 287, 262, 330, and 263 testing bounding boxes.

To compare against SPOCA, we consider a subset of the UAD testing set for which SPOCA detection results are available in HFC: the SPOCA subset. It consists of 26 testing images (181, 168, 213, and 166 bounding boxes in the 304, 171, 195 and 284 Å images respectively).

## Independent Detection on Single Image Bands

We first compare detection results of Faster RCNN over individual image bands analysed independently (Table 2). This aims to evaluate different feature extraction DNNs, and will further serve as baseline to assess our proposed framework.

The ResNet50 architecture consistently produces better results than the VGG one with higher F1 scores in all experiments. Based on these results, we adopt the ResNet50 architecture as backbone for our framework in the next experiments.

When comparing the detection results per image band, we notice that 304 Å images are repeatedly amongst the most difficult to analyse in UAD, having the lowest F1-scores in all tests. On the other hand, 171 Å has the best results of UAD bands, followed by 284 and 195 Å. This may be explained by ARs having a denser or less ambiguous appearance in 171, 195, and 284 Å image bands than in 304 Å since they are higher in the corona. A similar observation can be made in the LAD dataset when comparing the magnetogram results to 3934 Å, where magnetograms observe a lower altitude than 3934 Å.

We also notice a strong contrast between a same detector's precision and recall on the different UAD bands. This further demonstrates that these bands are not equal in how easily they may be analysed, even though they were acquired at the same time with same size and resolution.

Detections are visually verified to be poorer for small ARs and for spread and faint ones with more ambiguous boundaries. Visual inspection also confirms the different

<sup>2</sup> Our labelling tool will be released on the project's website.

**Table 1** Technical summary of the annotated datasets (top: UAD, bottom: LAD)

Modality	Image res.	Obs. frequency	Activity level	# Img.	# BBoxes
SOHO/EIT 284 Å	1024 × 1024	12 min	High	84	801
			Medium	93	610
			Low	146	868
			All	323	2279
SOHO/EIT 171 Å	1024 × 1024	12 min	High	84	873
			Medium	93	635
			Low	146	673
			All	323	2181
SOHO/EIT 195 Å	1024 × 1024	12 min	High	84	875
			Medium	93	678
			Low	146	1118
			All	323	2671
SOHO/EIT 304 Å	1024 × 1024	12 min	High	84	807
			Medium	93	614
			Low	146	1071
			All	323	2492
PM/SH 3934 Å	1500 × 1340	~1 day	High	47	618
			Medium	91	697
			Low	128	471
			All	266	1786
SOHO/MDI Magn.	1024 × 1024	96 min	High	47	618
			Medium	91	697
			Low	128	471
			All	266	1786

**Table 2** Detection performance of the single image band detectors. For each band, the highest scores are highlighted in bold

Detector	Dataset	Band (Å)	Precision	Recall	F1
Faster RCNN (ResNet50)	LAD	3934	<b>0.93</b>	<b>0.82</b>	<b>0.87</b>
	LAD	Magn.	<b>0.89</b>	<b>0.78</b>	<b>0.83</b>
	UAD	304	<b>0.73</b>	<b>0.83</b>	<b>0.78</b>
	UAD	171	<b>0.84</b>	<b>0.89</b>	<b>0.86</b>
	UAD	195	<b>0.81</b>	<b>0.75</b>	<b>0.78</b>
	UAD	284	<b>0.86</b>	<b>0.82</b>	<b>0.84</b>
	SPOCA	304	<b>0.72</b>	<b>0.82</b>	<b>0.77</b>
	SPOCA	171	<b>0.87</b>	<b>0.87</b>	<b>0.87</b>
	SPOCA	195	<b>0.82</b>	<b>0.73</b>	<b>0.77</b>
	SPOCA	284	<b>0.86</b>	<b>0.82</b>	<b>0.84</b>
Faster RCNN (VGG16)	UAD	304	0.67	0.78	0.72
	UAD	171	<b>0.84</b>	0.81	0.82
	UAD	195	0.79	0.73	0.76
	UAD	284	0.83	0.81	0.82
	SPOCA	304	0.68	0.80	0.74
	SPOCA	171	0.85	0.80	0.82
	SPOCA	195	0.78	0.72	0.75
	SPOCA	284	0.84	<b>0.82</b>	0.83

performances in various bands being caused by differing visual complexities of ARs. These observations suggest that detecting ARs using information provided by a single band may be an under-constrained problem.

### Joint Detection on Multiple Image Bands

We now present the results of our framework when detecting ARs over the LAD/UAD bands jointly. Joint detection results are summarised in Table 3.

In our first experiment, we compare early fusion (pixel level concatenation) against late fusion (feature level concatenation or addition), on the LAD dataset. Overall, the three approaches show an enhanced performance in contrast to single band based detection. However, we find that late fusion with concatenation shows higher performance than early fusion, having 0.90 F1-score versus 0.88 for magnetograms, while both scored 0.89 over 3934 Å. We further test late fusion using element wise addition and observe a decrease of 1% and 3% in the F1-score over 3934 Å and Magnetogram, respectively. Hence, we choose the late fusion with concatenation approach for all the following experiments.

We also evaluate the benefit of our MOO strategy using our 2-band architecture over the UAD. This approach

**Table 3** AR detection performance of the MSMT-CNN detectors. For each band, the highest scores are highlighted in bold

Detector	Fusion	Dataset	Bands (Å)	Prec.	Recall	F1			
MSMT-CNN (ResNet50—MOO)	Early—concat.	LAD	3934	0.96	<b>0.82</b>	<b>0.89</b>			
			Magn.	0.95	0.82	0.88			
	Late—concat.		3934	<b>0.97</b>	<b>0.82</b>	<b>0.89</b>			
			Magn.	<b>0.96</b>	<b>0.85</b>	<b>0.90</b>			
	Late—addition		3934	0.95	<b>0.82</b>	0.88			
			Magn.	0.94	0.80	0.87			
MSMT-CNN (ResNet50)	Late—concat.	UAD	171	0.92	0.77	0.84			
			284	0.90	0.81	0.85			
			171	0.82	0.85	0.83			
			195	0.86	0.72	0.78			
			195	0.88	0.67	0.77			
			284	0.84	0.78	0.81			
			304	0.82	0.79	<b>0.80</b>			
			195	0.87	0.75	0.80			
			MSMT-CNN (ResNet50—MOO)	Late—concat.	UAD	171	0.90	0.83	<b>0.87</b>
						284	<b>0.93</b>	0.80	<b>0.86</b>
SPOCA	171	<b>0.89</b>			0.83	<b>0.86</b>			
	284	<b>0.92</b>			<b>0.80</b>	<b>0.86</b>			
UAD	171	0.86			0.77	0.82			
	195	0.89			0.75	0.81			
SPOCA	171	0.83			0.77	0.80			
	195	0.86			0.73	0.79			
MSMT-CNN (ResNet50—MOO)	Late—concat.	UAD			195	0.88	0.68	0.77	
					284	0.84	0.78	0.81	
		SPOCA	195	<b>0.87</b>	0.67	0.75			
			284	0.81	0.78	0.80			
		UAD	304	0.82	0.78	<b>0.80</b>			
			195	0.88	<b>0.78</b>	<b>0.83</b>			
		SPOCA	304	<b>0.79</b>	<b>0.78</b>	<b>0.79</b>			
			195	0.85	0.77	<b>0.81</b>			
		UAD	304	0.78	0.74	0.76			
			171	0.76	0.76	0.76			
MSMT-CNN (ResNet50—MOO)	Neighbour bands—late—concat.	UAD	284	0.79	0.78	0.78			
			304	<b>0.93</b>	0.69	0.79			
			171	<b>0.94</b>	0.66	0.78			
			195	<b>0.91</b>	0.72	0.80			
			284	<b>0.93</b>	0.66	0.77			
			304	0.72	0.76	0.74			
			171	0.74	0.79	0.76			
			195	0.81	0.73	0.77			
			284	0.68	<b>0.84</b>	0.75			
			SPOCA	Early—concat.	SPOCA	171	0.54	<b>0.93</b>	0.68
195	0.58	<b>0.82</b>				0.68			
[12] using Faster-RCNN (ResNet50)	Sequential fine-tuning	UAD	304	0.73	<b>0.83</b>	0.78			
			171	0.80	<b>0.90</b>	0.84			
			195	0.83	0.72	0.77			
			284	0.86	0.80	0.83			



generally improves the F1-scores in most bands comparing to the non-MOO architectures. This behaviour may indicate that the two feature extraction stages were indeed more effectively optimised for their different tasks at different epochs. We retain this MOO approach for all other experiments.

On the UAD dataset, with various combinations of two bands, we notice a general improvement over single band detections. In addition, the performance varies in correspondence to the bands being used. Combining bands that are difficult to analyse (304 or 195 Å that have lowest F1-scores in the single band analyses) with easier bands (171 and 284 Å) unsurprisingly enhances their respective performance. More interestingly, combining the difficult 304 and 195 Å bands together also improve on their individual performance. Similarly, when combining bands that are easier to analyse (171 and 284 Å), in contrast to using combinations of difficult and easy bands in the analysis, performances are also improved over their individual analyses. Following these settings, our two-band based approach was able to record higher or similar F1-scores in contrast to the best performing single-band detector. This supports our hypothesis that joint detection may provide an increased robustness through learning the inter-dependencies between the image bands.

Moreover, the most dramatic improvement in F1-scores across both LAD and UAD datasets is for the 3934 Å images when magnetograms are added to the analysis. This is in line with the current understanding of AR having strong magnetic signatures.

Generally, in the UAD dataset, we find that using a combination of 2 bands produces the best results in comparison to using 3 or 4 bands. This may be caused by the fact that optimising the network for multiple tasks (2, 3, or 4 detection tasks) simultaneously increases the complexity of the problem. While the network successfully learned to produce better detections in the case of 2 bands, it was difficult to find a generalised yet optimal model for 3 or 4 bands at the same time.

Furthermore, since bands imaging consecutive layers of solar atmosphere are expected to be highly correlated, we test our framework by combining directly neighbouring bands together, such that a prediction for a band is performed using the band's own feature map combined with the feature map(s) of its available (1 or 2) direct neighbour(s). This approach gets the highest recall score on the UAD band 284 Å of all tests, where it is combined with the 195 Å band. However, it does not improve the performance on the other bands comparing to the single-band and 2-band based experiments.

We compare against state-of-the-art SPOCA [32] on the SPOCA subset, and against the first stage of [12] (sequentially fine-tuned networks) by adapting their approach to

Faster RCNN and testing it on UAD. SPOCA detections were obtained from 171 Å and 195 Å images only, combined as two channels of an RGB image, and SPOCA produces a single detection for both bands. We compare this detection against the ground truth of each of the bands individually. To prove the robustness and versatility of our detector, we also experiment with a combination of chromosphere, transition region, and corona bands on the SPOCA subset in addition to the whole UAD.

On the SPOCA subset, over the bands 171 Å and 195 Å for which it is designed, SPOCA gets the poorest performance of all multi-band and single-band experiments. It is worth noting that this method relies on manually tuned parameters according to the developers' own definition and interpretation of AR boundaries, which may differ from the ones we used when annotating the dataset. While supervised DL-based methods could integrate this definition during training, SPOCA could not perform such adaptation. This may have had a negative impact on its scores. Furthermore, visual inspection shows a poor performance for SPOCA on low solar activity images, see Fig. 4. This may be due to the use of clustering in SPOCA, since in low activity periods the number of AR pixels (if any) is significantly smaller than solar background pixels, which makes it hard to identify clusters.

Moreover, the fine-tuned networks of [12] suffer from a high rate of false positives, and show a close performance to single band detection using Faster RCNN with an identical precision, recall and F1-score over the band 304 Å and a slight decrease over the other 3 bands. This may be due to the fact that its transfer learning does not incorporate the inter-dependencies directly when analysing the different bands.

We further evaluate our method using IoU and compare our results to both, single band based detection by Faster-RCNN and joint detection from SPOCA and the fine tuned network of [12] (see Table 4). When using an IoU threshold of 50%, we find that our method produces the highest F1 score over the 195 Å band amongst all methods on both UAD and SPOCA datasets, with a comparable performance over 171 Å and 284 Å. On the other hand, our method shows a drop in performance on 304 Å comparing to Faster RCNN and the sequentially fine-tuned network. Generally, all methods show a significant decrease in the F1 scores when using IoU based criterion, with SPOCA being the lowest amongst all methods. A similar pattern is observed when using a less strict IoU threshold of 25%. During our visual inspection, we notice that in some cases, detected ARs are regarded false positives when evaluated against the IoU criterion. Particularly, when a cluster of neighbouring active region areas (i.e. neighbouring ARs) is detected as a single AR structure, or when an AR is detected by multiple neighbouring boxes. See Fig.

**Table 4** AR detection performance of MSMT-CNN and baseline detectors using two IoU (intersection over union) threshold values, 0.5 and 0.25. The rows with italic background indicate results produced using single-band based detectors. For each band, the highest scores are highlighted in bold

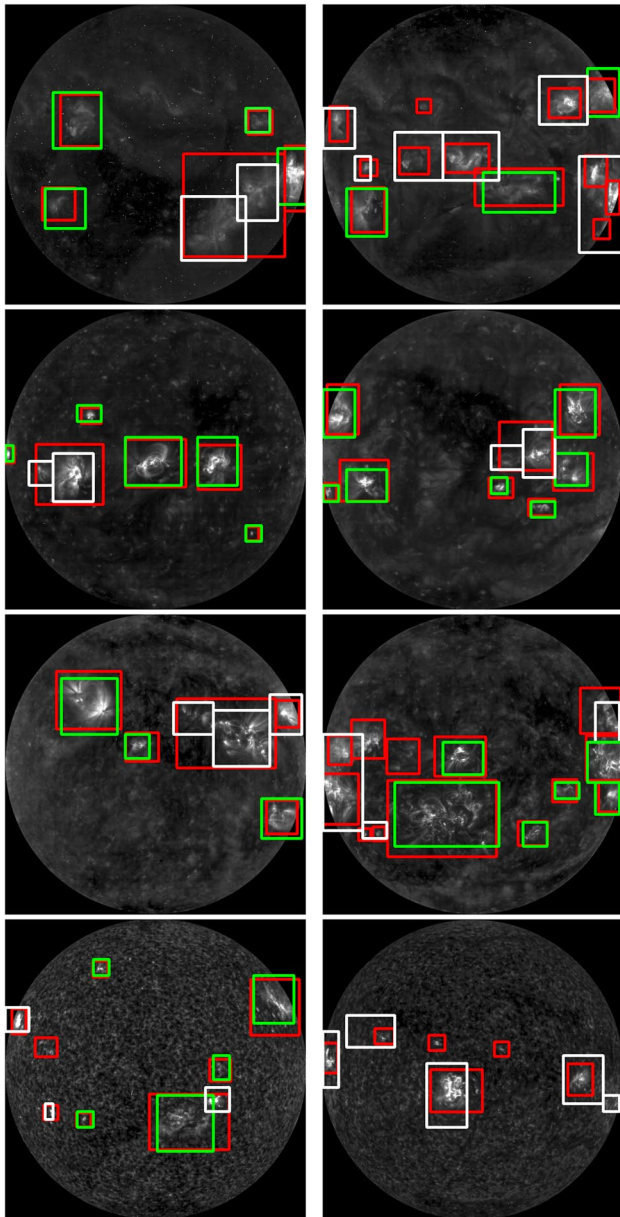
Detector	Dataset	Bands (Å)	Prec. IoU 0.5 (0.25)	Recall	F1
Faster-RCNN (ResNet50)	UAD	304	<b>0.56</b> (0.70)	<b>0.65</b> (0.82)	<b>0.60</b> (0.75)
		171	<b>0.65</b> (0.78)	0.68 (0.82)	<b>0.66</b> (0.80)
		195	0.61 (0.76)	0.52 (0.57)	0.56 (0.65)
		284	0.66 (0.82)	<b>0.63</b> (0.79)	0.64 (0.81)
	SPCOA	304	<b>0.55</b> (0.70)	<b>0.65</b> (0.82)	<b>0.59</b> (0.75)
		171	0.66 (0.81)	0.67 ( <b>0.83</b> )	0.67 ( <b>0.82</b> )
		195	0.60 (0.76)	0.54 (0.69)	0.57 (0.73)
		284	0.68 (0.83)	<b>0.65</b> (0.80)	<b>0.66</b> (0.82)
MSMT-CNN (Late concat. – ResNet50)	UAD	304	0.51 ( <b>0.77</b> )	0.48 (0.73)	0.49 ( <b>0.75</b> )
		171	0.69 ( <b>0.86</b> )	0.61 (0.76)	0.65 ( <b>0.81</b> )
		195	<b>0.67</b> ( <b>0.83</b> )	<b>0.59</b> ( <b>0.73</b> )	<b>0.62</b> ( <b>0.77</b> )
		284	<b>0.69</b> ( <b>0.88</b> )	0.59 (0.76)	0.64 ( <b>0.81</b> )
	SPOCA	304	0.52 ( <b>0.75</b> )	0.50 (0.73)	0.51 (0.74)
		171	<b>0.70</b> ( <b>0.84</b> )	0.64 (0.76)	0.66 (0.80)
		195	<b>0.64</b> ( <b>0.80</b> )	<b>0.58</b> ( <b>0.73</b> )	<b>0.61</b> ( <b>0.76</b> )
		284	<b>0.71</b> ( <b>0.88</b> )	0.62 (0.78)	<b>0.66</b> ( <b>0.83</b> )
SPOCA	SPOCA	171	0.16 (0.34)	0.26 (0.56)	0.19 (0.42)
		195	0.16 (0.33)	0.20 (0.43)	0.18 (0.37)
[12] using Faster-RCNN (ResNet50)	UAD	304	<b>0.56</b> (0.70)	<b>0.65</b> ( <b>0.82</b> )	<b>0.60</b> ( <b>0.75</b> )
		171	0.62 (0.74)	<b>0.70</b> ( <b>0.84</b> )	<b>0.66</b> (0.79)
		195	0.62 (0.79)	0.55 (0.70)	0.58 (0.74)
		284	0.68 (0.83)	0.62 (0.75)	<b>0.65</b> (0.79)
	SPOCA	304	<b>0.55</b> (0.70)	<b>0.65</b> ( <b>0.82</b> )	<b>0.59</b> ( <b>0.75</b> )
		171	0.66 (0.77)	<b>0.71</b> (0.82)	<b>0.68</b> (0.79)
		195	0.62 (0.76)	0.54 (0.67)	0.58 (0.71)
		284	0.69 (0.83)	0.63 (0.77)	<b>0.66</b> (0.80)

“Experiments”. These observations suggest that an IoU based criterion does not perfectly capture the dynamic characteristics of solar ARs, even when incorporating an IoU threshold as low as 25%. Unlike generic object detection tasks, where object morphology and boundaries are well defined, ARs are dynamic structures that are continuously evolving (e.g. merging and splitting, emerging and dying out) [26, 27, 30, 32]. Therefore, when designing our evaluation criterion (Sect. “Experiments”), we take into account the aforementioned scenarios associated with solar ARs.

To further demonstrate the benefits of our joint analysis based approach, we create a synthetic dataset from the BraTS multi-modal dataset [16] of similar spatial configurations to the solar imaging bands. BraTS consists of full 3D MR image volumes of brain in four modalities (T1GD, T1, T2, and Flair) and three classes: enhancing tumour (ET), necrotic and non-enhancing tumour core (NCR/NET), and peritumoural edema (ED). We create the synthetic dataset by selecting one 2D slice of each

image modality separated by a spatial gap of size 1 voxel. This emulates the solar images scenario where each band shows ARs in a different solar altitude.

Although this gap size may seem much lower than for solar images, they are justified by the speed of change of the imaged brain from one slice to another neighbouring one being much larger than for the generally smoother ARs. For each modality, we use a total of 11,533 and 190 training and testing images, respectively. We evaluate our detection approach with different fusions, over the 4 bands of BraTS-prime dataset, and compare it against single band based detection using Faster RCNN. See Tables 5 and 6. All fusion strategies significantly outperform single band detectors, with an average F1-score increase of 18%, 28%, and 32% for early addition, early concatenation, and late concatenation fusion, respectively. This confirms our hypothesis that exploiting inter-dependencies between the image bands by the joint analysis may provide a superior performance in contrast to single band based detection.



**Fig. 3** Detections of solar active regions visualized when an IoU based criterion is applied during evaluation. Ground-truth (red), and MSMT-CNN detections (true positive in green and false positive in white) of ARs in images from SOHO/EIT (top to bottom) 284, 171, 195, and 304 Å. We observe that using IoU as an evaluation criterion causes some detected ARs to be regarded as false positives. Particularly, when an AR is detected by multiple boxes, or when multiple neighboring ARs are detected as a single structure

## Conclusion

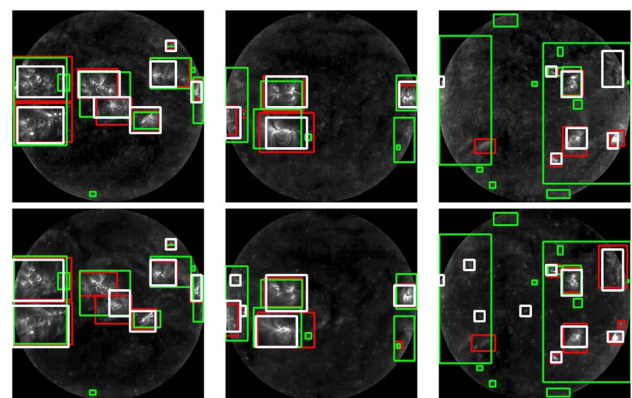
We presented MSMT-CNN, a multi-branch and multi-tasking framework to tackle the 3D solar AR detection problem from multi-spectral images that observe different cuts of the 3D solar atmosphere. MSMT-CNN analyses multiple image bands jointly to produce consistent detection across them.

**Table 5** F1-scores of single band based detectors on the BraTS-prime dataset (with 1 slice gap). All detectors are based on ResNet50

Detector	Image band	Precision	Recall	F1
Faster RCNN (ResNet50)	T1Gd	0.57	0.55	0.56
	T1	0.56	0.52	0.54
	T2	0.50	0.47	0.48
	Flair	0.59	0.98	0.73

**Table 6** F1-scores of MSMT-CNN with different fusion strategies on the BraTS-prime dataset (with 1 slice gap). All detectors are based on ResNet50. For each band, the highest scores are highlighted in bold

Detector	Image band	Precision	Recall	F1
MSMT-CNN (Early—addition)	T1Gd	0.62	0.97	0.76
	T1	0.65	<b>0.98</b>	0.78
	T2	0.61	<b>0.98</b>	0.75
	Flair	0.60	<b>0.99</b>	0.74
MSMT-CNN (Early—concat.)	T1Gd	0.78	<b>0.98</b>	0.86
	T1	0.82	<b>0.98</b>	0.89
	T2	0.77	0.97	0.86
	Flair	0.73	0.98	0.84
MSMT-CNN (Late—concat.)	T1Gd	<b>0.83</b>	0.96	<b>0.89</b>
	T1	<b>0.85</b>	<b>0.98</b>	<b>0.91</b>
	T2	<b>0.85</b>	<b>0.98</b>	<b>0.91</b>
	Flair	<b>0.81</b>	0.98	<b>0.89</b>



**Fig. 4** Ground-truth (red), proposed (white), and SPOCA's (green) detection of ARs over the three solar activity periods (left to right: high, medium, low) in randomly selected images from SOHO/EIT 171 Å (top) and 195 Å (bottom)

It is a flexible framework that may use any CNN backbone, and may be straightforwardly generalised to any number and modalities of images. Our findings suggest that jointly analysing information from different image bands, at different feature levels, can elevate the detection performance across the image bands. MSMT-CNN showed competitive

results against baseline and state-of-the-art detection methods. Future research could investigate incorporating multi-level feature fusion in multi-view imaging systems such as cross sectional imaging in the medical domain which demonstrates an interesting direction for future research.

## Declarations

**Conflict of Interest** The authors declare that they have no conflict of interest.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Almahasneh M, Paiement A, Xie X, Aboudarham J. Active region detection in multi-spectral solar images. In: International conference on pattern recognition applications and methods; 2021.
- Benkhalil A, Zharkova V, Zharkov S, Ipson S. Active region detection and verification with the solar feature catalogue. *Sol Phys*. 2006;235:87–106.
- Bezdek J. Objective function clustering. In: Pattern recognition with fuzzy objective function algorithms. Advanced applications in pattern recognition; 1981.
- Crabbe B, Paiement A, Hannuna S, Mirmehdi M. Skeleton-free body pose estimation from depth images for movement analysis. In: IEEE international conference on computer vision workshops; 2015.
- Dai J, Li Y, He K, Sun J. R-fcn: Object detection via region-based fully convolutional networks. In: Advances in neural information processing systems. Curran Associates, Inc.; 2016.
- Dalal N, Triggs B. Histograms of oriented gradients for human detection. In: IEEE computer society conference on computer vision and pattern recognition; 2005.
- Eitel A, Springenberg JT, Spinello L, Riedmiller M, Burgard W. Multimodal deep learning for robust rgb-d object recognition. In: IEEE/RSJ international conference on intelligent robots and systems; 2015.
- Gani MO, Kuiry S, Das A, Nasipuri M, Das N. Multispectral object detection with deep learning; 2021.
- Guo Z, Li X, Huang H, Guo N, Li Q. Deep learning-based image segmentation on multimodal medical imaging. *IEEE Trans Radiat Plasma Med Sci*. 2019;3(2):162–9.
- Hwang S, Park J, Kim N, Choi Y, Kweon I. Multispectral pedestrian detection: Benchmark dataset and baselines. In: Proceedings of the IEEE conference on computer vision and pattern recognition; 2015.
- Ishii T, Simo-Serra E, Iizuka S, Mochizuki Y, Sugimoto A, Ishikawa H, Nakamura R. Detection by classification of buildings in multispectral satellite imagery. In: International conference on pattern recognition; 2016.
- Jarolim R, Veronig A, Hofmeister S, Temmer M, Heinemann S, Podladchikova T, Dissauer K. Multi-channel coronal hole detection with a CNN. In: Machine learning in heliophysics; 2019.
- Krishnapuram R, Keller J. The possibilistic C-means algorithm: insights and recommendations. *IEEE Trans Fuzzy Syst*. 1996;4(3):385–93.
- Law H, Deng J. Cornernet: Detecting objects as paired keypoints. In: European conference on computer vision; 2018.
- Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu C-Y, Berg A. SSD: single shot multibox detector. In: European conference on computer vision; 2016.
- Menze BH, Jakab A, Bauer S, Kalpathy-Cramer J, Farahani K, Kirby J, Burren Y, Porz N, Slotboom J, Wiest R et al. The multimodal brain tumor image segmentation benchmark (BRATS). *IEEE Trans Med Imaging* 2015.
- Mohajerani S, Krammer T, Saeedi P. A cloud detection algorithm for remote sensing images using fully convolutional neural networks. In: IEEE 20th international workshop on multimedia signal processing; 2018.
- Mohajerani S, Saeedi P. Cloud-Net: an end-to-end cloud detection algorithm for Landsat 8 imagery. In: IEEE international geoscience and remote sensing symposium; 2019.
- Osorio K, Puerto A, Pedraza C, Jamaica D, Rodríguez L. A deep learning approach for weed detection in lettuce crops using multispectral images. *AgriEngineering*. 2020;2(3):471–88.
- Otsu N. A threshold selection method from gray-level histograms. *IEEE Trans Syst Man Cybernet*. 1979;9(1):62–6.
- Penatti OAB, Nogueira K, dos Santos JA. Do deep features generalize from everyday objects to remote sensing and aerial scenes domains? In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops; 2015.
- Redmon J, Divvala S, Girshick R, Farhadi A. You only look once: Unified, real-time object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition; 2016.
- Ren S, He K, Girshick R, Sun J. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Trans Pattern Anal Mach Intell*. 2015;39:1137–49.
- Revathy K, Lekshmi S, Prabhakaran Nayar S. Fractal-based fuzzy technique for detection of active regions from solar images. *Solar Phys*. 2005;228:43–53.
- Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg A, Fei-Fei L. ImageNet large scale visual recognition challenge. *International journal of computer vision*. 2015.
- Schunker H, Birch AC, Cameron RH, Braun DC, Gizon L, Burston RB. Average motion of emerging solar active region polarities-i two phases of emergence. *Astron Astrophys*. 2019;625:A53.
- Schunker H, Braun D, Birch A, Burston R, Gizon L. Sdo/hmi survey of emerging active regions for helioseismology. *Astronomy & Astrophysics*. 2016;595:A107.
- Simonyan K, Zisserman A. Two-stream convolutional networks for action recognition in videos. *Neural Inf Process Syst*. 2014.
- Song X, Gao S, Chen C. A multispectral feature fusion network for robust pedestrian detection. *Alexandria Eng J*. 2021.
- Sun X, Bobra MG, Hoeksema JT, Liu Y, Li Y, Shen C, Couvidat S, Norton AA, Fisher GH. Why is the great solar active region 12192 flare-rich but cme-poor? *Astrophys J Lett*. 2015;804(2):L28.
- Takumi K, Watanabe K, Ha Q, Tejero-De-Pablos A, Ushiku Y, Harada T. Multispectral object detection for autonomous vehicles. In: Proceedings of the on thematic workshops of ACM multimedia; 2017.
- Verbeeck C, Delouille V, Mampaey B, Visscher RD. The SPoCA-suite: software for extraction, characterization, and tracking of



- active regions and coronal holes on EUV images. *Astron Astrophys.* 2013;561:16.
33. Viola P, Jones MJ. Robust real-time face detection. *Int J Comput Vision.* 2004;57(2):137–54.
  34. Wagner J, Fischer V, Herman M, Behnke S. Multispectral pedestrian detection using deep fusion convolutional neural networks. In: *European symposium on artificial neural networks*; 2016.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.