

The author thanks the referees for their supportive words and helpful comments that are individually addressed below. The relevant changes to the paper have been marked in bold. I believe it is a better piece of work in light of this feedback.

Referee 1:

This is a very well written and presented paper on the automatic detection of CMEs in SOHO/LASCO data, which is also available as an online catalogue (CORIMP). CORIMP particularly goes beyond previous catalogues in showing the kinematics along many position angles. The previous research in this field is very well discussed, the new results are presented well and I particularly like a section on the application of a machine learning algorithm to distinguish between problem events where CMEs are launched close to one another in time and space. I think that this type of methodology will be used in many future studies on space weather because the amount of available solar and heliospheric data is already so big, and adequate artificial intelligence applications in space weather have almost not been introduced yet. In this respect the paper and the catalogue form a basis for likely many future studies on CMEs and I can definitely recommend the paper for publication with minor revision.

Minor points:

- the abstract is too long (About 300 words) - please make it more concise and to the point (about 200 should be fine). For instance, it is not needed to discuss the advantages of the drawbacks of the previous catalogues in the abstract, the introduction is good enough for that. It would also be good to write all the abbreviations in full the first time they are introduced, throughout the text.

The abstract has been shortened and edited accordingly.

- line 64: a reference for the Hough transform, please.

Reference added.

- line 77: if possible add a note on how to get the realtime alerts.

Footnote added.

- please add some more details on the Savitzky-Golay filter because it is very often mentioned, and I do not think many readers are familiar with it; e.g. whats the difference to a spline?

Details added where the Savitzky-Golay filter is first mentioned in Sect. 2.1.

- there are 6 case studies of events, compared between the different catalogues. There is a table for every one of them. I think it would streamline the presentation to put all these into 1 table, divided by horizontal lines so the events are separated, along with an indicator of the event date and what the event should demonstrate.

That is much better, thank you.

- it might also be good to revisit if every of the 6 events that are discussed in a subsection in section 3 are actually needed, and show something important. When reading the paper, its not immediately clear at the beginning of section 3 why the events have been chosen to demonstrate a particular point. A discussion at the beginning of section 3 on this point would be helpful.

At the beginning of section 3 it is stated that these events were chosen based on their varying styles of eruption and appearance, and I add another line to indicate this means they simply have different morphologies to be characterised and different kinematics to be derived. There is not much more significance to the choice of events than this, simply serving as an example of typical CMEs of interest.

- at the end of each of the subsections 3.1 3.2 etc. the wording is almost completely similar. Please rephrase this so it does not look like its copy-pasted, or leave the concluding sentences for the end of that section or the conclusion section 5; or summarize for each event what it should demonstrate, but in its own words.

I have edited the final lines of each subsection now to avoid this repetition.

- Concerning the section 3.5, it should be emphasized that there is a very large range of detected speeds: SEEDS gives 703 and CDAW 2393 km/s. For outside-users this is rather surprising.

This is true, though I do state in that section that SEEDS cannot be trusted as the CME front is only visible in two C2 frames.

- Section 4: When reading this section, I wondered how k is chosen, and then it is explained at the very end of the section. Please indicate in the beginning somewhere that choosing k=3 or 4 is actually given by the human observation of the number of CMEs.

I have added to an earlier part of this section that the k clusters are manually prescribed by the user.

Are there any ideas on how to figure out k automatically? I am aware that this

is beyond the scope of the current work and I leave it to the author to include a comment on this or not.

I have been finding it challenging to implement any form of reasonable determination of k automatically. It is certainly something to aspire to, as the manual prescription of k is not ideal (which is why I am not currently using it in the catalog). I think this would need to be investigated with supervised machine learning; if one could build a large and dedicated training set for example. I have added a comment on this in the final conclusions.

- It took me a while to figure out exactly what is shown in Figures 10 and 11. Similar to figures 1 and 2, please describe for each panel what is seen (or just use (a), (b), ...) . Alternatively, I suggest to present these 2 plots in this way: show the CORIMP height-time measurements used as input on top, then a vertical arrow pointing to the the current top panel of Figures 10 / 11 as the new middle panel, and then another arrow to the height-time measurements with the clustering results indicated by colors and symbols as the bottom panel. In this way, it would be immediately clear to see how the data were treated and what the output is.

Figures 10 and 11 and their captions have been changed in line with these suggestions and should now be more readily understood.

Referee 2:

The author begins by comparing various CME catalogs for the LASCO data set; CDAW, CACTus, SEEDS, and CORIMP. All of these catalogs use different methods for determining the line-of-sight velocity. All these methods have various advantages and disadvantages. With the expectation of CORIMP, all these catalogs use running difference images that are subject "spatiotemporal crosstalk." The author then goes on to compare the kinematics from CORIMP and the other catalogs for six example CMEs with various kinematic profiles. These comparisons highlight the advantages and disadvantages of using the Savitzky-Golay filter and the CORIMP detection methods. One of the main issues of the CORIMP detections method (and other automatic methods) is the

separation of CMEs that are close in space and time.

This paper presents an interesting comparison of the kinematics produced with CORIMP catalog to other catalogs and a new technique for separating CMEs a common problem in automated methods. This article is essential to validate and document the kinematic data generated by the CORIMP catalog. The current version of this article is acceptable for publication. I have provided some comments on the clarity of the text for the author to consider. I hope these comments are helpful to the author. I do not need to review this article further.

Figure 1 (Top left): This image is a background model subtracted image. The caption implies that these images are used to calculate the velocity in CDAW, which is not correct. CDAW uses running difference movies at the cadence taken by the instrument the same as CACTus.

The caption has been edited for clarification to remove this implication.

Line 135: While CACTus cannot calculate the acceleration of the CME front, it does have the advantage of calculating the velocity of the front over a range of PAs, which can also be an important aspect of the CME kinematics.

A line had been added to this section noting this advantage.

Line 154: How does this assumption affect the derived velocity? I assume that on average the SEEDS velocities are lower since all CME fronts become dimmer as they propagate out.

The referee is correct in understanding that the CMEs become dimmer and so the SEEDS velocities may be lower, however in the text it is explained that it is the resulting increased error on the height-time profile that affects the accuracy of the derived velocity. There is probably some semantics at play here, so I leave the text unchanged.

Section 2: In the introduction you mention that running difference images suffer from “spatiotemporal crosstalk.” However, you never discuss how this can affect the velocity measurements. The time between the images limits the velocities that can be detected in the images. Additional, changes in the image cadence changes the width and intensity increase created by the CME. Do you see a change in the detected CME velocities for these catalogs when the LASCO image cadence changed in 2010? (i.e. similar to the problem in detections found by Wang & Colaninno 2014 ApJL)

This is an interesting question and worth investigating in a larger-scale study of the overall trends in the CORIMP catalog rather than the specific case-studies used for the purposes of this paper.

On the issue of spatiotemporal cross-talk, some details have been added to the text in the introduction.

Fig 4: The black lines in the HT plots are not explained. The lines in the velocity and acceleration plots are impossible to see when printed on paper (old fashion).

I have edited the caption to explain the black lines in the HT plots. And yes, since these plots are quite busy they can be challenging to read, but the black lines are mostly just to imply a trend in the data (which print okay for me albeit without the ability to zoom-in as preferred).

Table 1: The abbreviation AW is ambiguous. Define the abbreviation in the text or add another line to the table headings to avoid abbreviations.

AW changed to Width, and explanations for the Table parameters have been added to the Table caption.

Line 201: Here you report accelerations from 50 to 0 m s⁻² but in the table you give a value of 1 m s⁻² with a range of 14 to -17 m s⁻². Please explain in the text what values from CORIMP are listed in the tables.

The added explanations in the Table caption go towards distinguishing the Table values (of linear or second-order fitted speeds) from those of the

Savitzky-Golay filter values discussed at that point.

Line 222: When was this data gap, in which instrument?

A comment on the data gap (in the LASCO/C3 images between 19:42 and 21:24 UT on 2000 Apr. 18) has been added to the text.

Line 262: Again you need to explain how the values in the tables are generated from the data presented in the figures. Especially for this event where the table values are very different from the “real” values for the CME.

Resolved, for the point above, by the added explanation in the Table caption.

Line 310: I assume that the deceleration seen for this event is the algorithm no longer identifying the leading edge of the CME as it becomes fainter. Are the data points after 15 Rsun unreliable?

By inspection of the kinematic plots for this event, there is a deceleration occurring before the CME reaches 15 Rsun, so the more unreliable HT points thereafter are not creating an artificial deceleration (although they can contribute to such a measure, which is part of the overall message of the paper to not believe single-value quoted measurements without inspecting the kinematic profiles to be sure such effects are not completely to blame).

Investigating the Kinematics of Coronal Mass Ejections with the Automated CORIMP Catalog

Jason P. Byrne

RAL Space, Rutherford Appleton Laboratory, Harwell Oxford, OX11 0QX, UK.
e-mail: jason.byrne@stfc.ac.uk

ABSTRACT

Studying coronal mass ejections (CMEs) in coronagraph data can be challenging due to their diffuse structure and transient nature, compounded by the variations in their dynamics, morphology, and frequency of occurrence. The large amounts of data available from missions like the Solar and Heliospheric Observatory (SOHO) makes manual cataloging of CMEs tedious and prone to human error, and so a robust method of detection and analysis is required and often preferred. A new coronal image processing catalog called CORIMP was developed in an effort to achieve this, through the implementation of a dynamic background separation technique and multiscale edge detection. These algorithms together isolate and characterise CME structure in the field-of-view of the Large Angle Spectrometric Coronagraph (LASCO) onboard SOHO. CORIMP also applies a Savitzky-Golay filter, along with quadratic and linear fits, to the height-time measurements for better revealing the true CME velocity and acceleration profiles across the plane-of-sky. Here we present a sample of new results from the CORIMP CME catalog, and directly compare them with the other automated catalogs of Computer Aided CME Tracking (CACTus) and Solar Eruptive Events Detection System (SEEDS), as well as the manual CME catalog at the Coordinated Data Analysis Workshop (CDAW) Data Center and a previously published study of the sample events. We further investigate a form of unsupervised machine learning by using a k -means clustering algorithm to distinguish detections of multiple CMEs that occur close together in space and time. While challenges still exist, this investigation and comparison of results demonstrates the reliability and robustness of the CORIMP catalog, proving its effectiveness at detecting and tracking CMEs throughout the LASCO dataset.

Key words. Sun – Coronal Mass Ejection (CME) – Space weather – Solar image processing – Machine learning

1. Introduction

Coronal mass ejections (CMEs) represent the largest, most dynamic phenomena that originate from the Sun (Chen, 2011; Webb and Howard, 2012). Propagating at speeds of up to thousands of kilometres per second, with energies on the order of 10^{32} ergs, they can drive adverse space weather throughout the solar system (Howard and Tappin, 2005; Pulkkinen, 2007). Given their potentially

hazardous impact on Earth's geomagnetic environment, the physics governing their eruption and propagation needs to be understood so that their effects may be predicted in the guise of space-weather forecasting. To this end, observations of CMEs must be rigorously inspected in order to determine their dynamics, and this is most generally undertaken with the use of coronagraph instruments (e.g., Koomen et al., 1975; Sheeley et al., 1980; MacQueen et al., 1980; Illing and Hundhausen, 1985; Hundhausen, 1993; Brueckner et al., 1995; Howard et al., 2008).

CMEs tend to be faint, transient phenomena, observed in white-light images that are prone to noise and user-dependent biases in their interpretation. During solar minimum they can occur every few days, but at solar maximum there can be several per day (St. Cyr et al., 2000; Yashiro et al., 2004). They exhibit a wide variety of morphologies, moving in unpredictable directions and speeds in the solar wind (Kilpua et al., 2009; Byrne et al., 2010; Liu et al., 2014). They can drive shocks in the solar atmosphere and interplanetary space (Howard and Tappin, 2005; Carley et al., 2013), and exhibit various levels of geo-effectiveness (Plunkett et al., 2001; Schwenn et al., 2005; Davis et al., 2009; Lugaz and Kintner, 2012). A wealth of image processing techniques have been explored to study CMEs in remote-sensing image data provided by such instruments as the Large Angle Spectrometric Coronagraph (LASCO; Brueckner et al., 1995) onboard the Solar and Heliospheric Observatory (SOHO; Domingo et al., 1995). These techniques generally rely on some form of image differencing to highlight moving features in the observed intensities, but this introduces spatiotemporal crosstalk and scaling issues that affect the accuracy of CME characterisations. **For example, the distance a CME moves between frames of varying cadence, along with its inherent morphological and brightness changes during that time, directly affects the calculations of running-difference images - to the point of changing how a user or algorithm characterises the CME structure and consequently its dynamics.** More advanced image processing methods have thus been explored, such as optical flow techniques (Colaninno and Vourlidas, 2006), supervised segmentation techniques (Goussies et al., 2010), wavelets (Stenborg and Cobelli, 2003) and curvelets (Gallagher et al., 2011). The large volume of data available has made it necessary to automate such techniques for detecting and tracking CMEs across images, with a view to cataloguing their kinematics and morphologies. This allows for more robust CME detections by avoiding the troublesome effects of standard image differencing techniques. It is therefore possible to maintain a non-biased characterisation of the CME structure in every event, since automated techniques are self-consistent.

To date, a point-and-click catalog of CMEs in LASCO data has been undertaken at the Coordinated Data Analysis Workshop (CDAW) Data Center (Gopalswamy et al., 2009), which operates by tracking CMEs in running difference images to produce information on the dynamics of each event. It is a wholly manual procedure and is therefore subject to user bias in interpreting the data. Automated catalogs have since been developed to overcome this bias and tedium. The Computer Aided CME Tracking routine (CACTus; Robbrecht and Berghmans, 2004) is the first such automated catalog, that works by using a Hough transform (Hough, 1962) to detect intensity ridges corresponding to CME tracks in time-height stacks (J-maps) of polar-unwrapped running-difference LASCO images. The Solar Eruptive Events Detection System (SEEDS; Olmedo et al., 2008) is another automated catalog that similarly uses polar-unwrapped running-difference LASCO images but with a form of threshold segmentation to approximate the shape of the CME leading edge in every image. A new automated catalog has recently been developed from a unique set of coronal image processing techniques, called CORIMP, that overcomes many of the limitations of

76 current catalogs in operation (Morgan et al., 2012; Byrne et al., 2012). An online database has been
 77 produced for the SOHO/LASCO data and event detections therein; providing information on CME
 78 onset time, position angle, angular width, speed, acceleration, and mass, along with kinematic plots
 79 and observation movies. By investigating the catalog output it is intended that this work will lead
 80 to an improved understanding of the dynamics of CMEs. Furthermore, a realtime version of the
 81 algorithm has been implemented to provide CME detection alerts to the interested space weather
 82 community¹.

83 In Sect. 2 the CME catalogs are discussed in greater detail to highlight their methodologies and
 84 drawbacks. A sample of CMEs is then investigated in Sect. 3 in order to compare the outputs of each
 85 catalog, paying particular interest to the robustness and reliability of the new CORIMP catalog. In
 86 Sect. 4 a first effort is made to use a form of unsupervised machine learning to isolate spatially
 87 and temporally overlapping CME detections. The conclusions of this investigation are presented in
 88 Sect. 5.

89 2. Cataloging CMEs

90 In coronagraph images CMEs are observed as outwardly moving regions of stronger brightness in-
 91 tensities than the background corona (see the example in Fig. 1). Different methods for thresholding
 92 the CME intensity in such images have been employed by different catalogs. This is in order to de-
 93 tect their appearance and track their motion through the field-of-view, leading to a determination of
 94 the CME kinematics and morphology. However, each method suffers from drawbacks and, as such,
 95 the resulting CME catalogs can vary significantly in their characterisations and measurements of
 96 each event.

97 2.1. CORIMP Automated Catalog

98 The CORIMP² catalog was developed with a method of dynamic signal separation and multiscale
 99 edge detection to overcome certain drawbacks of previous detection and tracking methods that rely
 100 on running-difference images. The CME signal is separated from the more quiescent streamers and
 101 coronal structures in LASCO/C2 and C3 images, such that a multiscale filtering technique may be
 102 used to suppress noise in order to characterise the CME structure and track its motion. A spread
 103 of heights is then measured across the angular span of the CME. A type of cleaning algorithm is
 104 applied to the height-time measurements before the kinematics are determined. This is required to
 105 overcome cases where pixels along the CME front edges at a specific position angle are not cor-
 106 rectly identified, but rather a pixel corresponding to core material behind the CME front is measured,
 107 which causes unnecessary scatter in the height-time datapoints. The cleaning algorithm works by
 108 stepping along the height-time measurements at each position angle and requiring that only se-
 109 quentially increasing heights are plotted. This helps remove the detections of trailing CME material
 110 from the height-time plots such that only the kinematics of the main CME front are determined.
 111 The kinematics are then derived in three ways: using a Savitzky-Golay filter (Savitzky and Golay,
 112 1964), a quadratic fit (second-order polynomial), and a linear fit (straight-line). The importance of

¹ Automated email alerts may be requested from the author, and are also published on social media at twitter.com/CMEcatalog and facebook.com/CMEcatalog

² <http://alshamess.ifa.hawaii.edu/CORIMP>

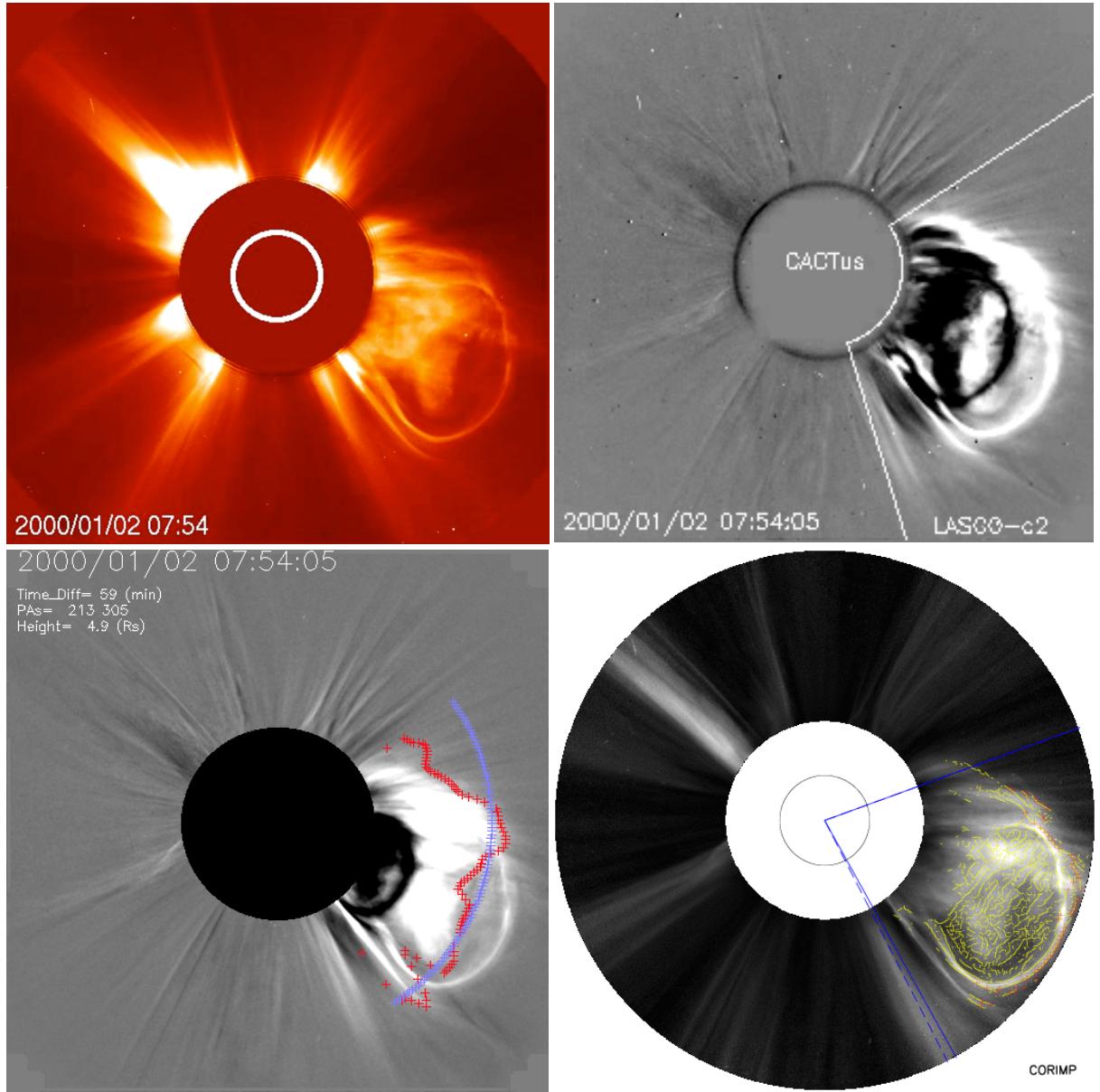


Fig. 1. LASCO/C2 observations of a CME on 2000 Jan. 02 at 07:54 UT. *Top left:* A level 2 background-model subtracted image of the event. *Top right:* Running difference image reproduced from the CACTus catalog with the angular span of the CME detection indicated by the white lines. (**The manual CDAW catalog uses similar such running-difference images.**) *Bottom left:* Running difference image reproduced from the SEEDS catalog with the CME front detection highlighted in red (and the extended ‘half-max lead’ in purple). *Bottom right:* Normalised radial-gradient filtered (NRGF) image taken from the CORIMP catalog with the angular span of the CME detection indicated in blue, the pixel-chained CME structure in yellow, and the CME front in red.

¹¹³ a robust method for determining the kinematics of a transient event is discussed in [Byrne et al.](#)

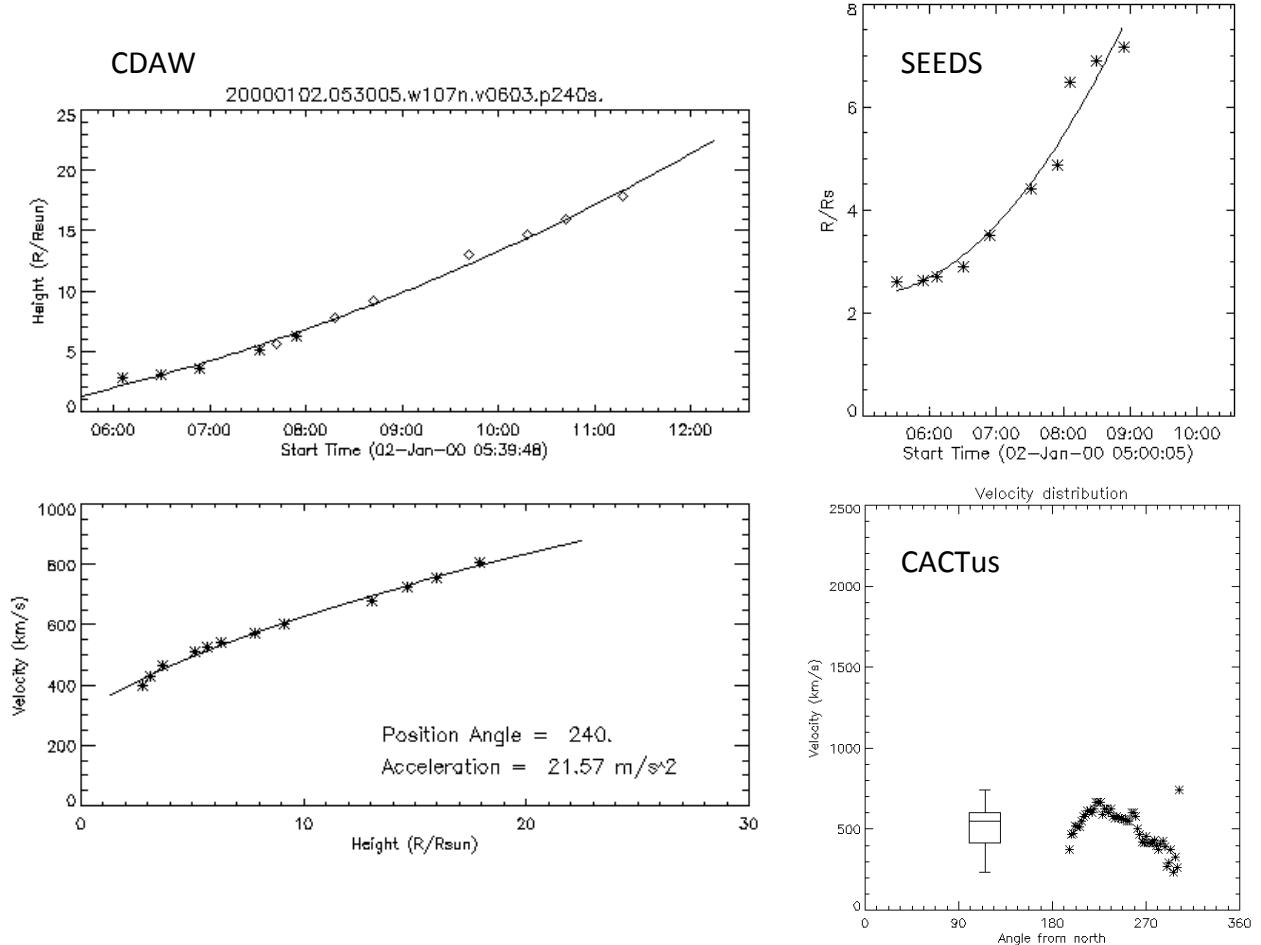


Fig. 2. The kinematic outputs for the 2000 Jan. 02 CME reproduced from the CDAW, SEEDS and CACTus catalogs. *Top left:* The CDAW catalog height-time measurements of the CME chosen manually along the running-difference bright front (at position angle 240) with a second-order fit. *Bottom left:* The corresponding CDAW velocity profile plotted against height, showing an acceleration of 21.57 m s^{-2} . *Top right:* The automated SEEDS height-time measurements and second-order fit resulting in an acceleration of 18.6 m s^{-2} in the LASCO/C2 field-of-view. *Bottom right:* The automated CACTus velocities determined along the angular span of the CME, with a corresponding box-and-whisker plot to highlight the median (548 km s^{-1}) and interquartile range.

114 (2013), wherein the often-used method of 3-point Lagrangian interpolation and associated error
 115 propagation were shown to behave counter-intuitively and provide misleading kinematic results. It
 116 was shown that the Savitzky-Golay smoothing filter performs well on CME height-time data,
 117 through its use of a kernel-based estimation of a local polynomial regression, that also directly
 118 computes the derivatives of the parameters to provide the kinematics. Given the large variety
 119 of CMEs that can occur with differing kinematic profiles, the inspection of these three automated
 120 fitting techniques in CORIMP can better reveal the true trends of CME motion.

121 **2.2. CACTus Automated Catalog**

122 The CACTus³ catalog was the first automated CME detection algorithm, in operation since 2004.
 123 It is based upon the detection of CMEs as bright ridges in time-height slices at each angle around a
 124 coronagraph image. A running-difference technique is applied and each image is transformed into
 125 Sun-centred polar coordinates, re-binned, and the C2 and C3 fields-of-view are combined. These are
 126 then stacked in time, and for each angle the corresponding time-height slice undergoes a modified
 127 Hough transform for detecting intensity ridges across it. Thresholding the most significant ridges
 128 filters out the progression of CMEs, with the variables for each ridge characterised by onset time,
 129 velocity, and position angle. A median velocity across the angular span of each event is quoted as
 130 the CME speed.

131 The running-difference cadence, the ridge intensity threshold, and the imposed limit on how
 132 many frames a CME may exist, all affect how successful the automated detection can be. However,
 133 [Robbrecht and Berghmans \(2004\)](#) show the algorithm to be robust in reproducing the detections of
 134 a human user by directly comparing with the CDAW catalog. The main drawback of the CACTus
 135 catalog for studying CMEs is the imposed zero acceleration of the detection algorithm, since the
 136 Hough transform thresholds the ridges as straight lines whose slopes provide a constant velocity.
 137 The velocity itself may also be an underestimate since it is a median across the span of the CME,
 138 **although having a spread of velocity measurements across position angles is an advantage of**
 139 **CACTus** (see the example velocity plot at the bottom right of Fig. 2). However, it is sometimes
 140 possible that the angular spans are over-estimated since side outflows in the images are enhanced
 141 by the running-difference and may include streamer deflections. It is also difficult to distinguish
 142 when one CME has fully progressed from the field-of-view and another CME has entered it, so in
 143 some cases trailing portions of a CME are detected as separate events.

144 **2.3. SEEDS Automated Catalog**

145 The SEEDS⁴ catalog employs an automated CME detection algorithm for tracking an intensity
 146 thresholded CME front in running-difference images from LASCO/C2. The images are unwrapped
 147 into Sun-centred polar coordinates, and a normalised running-difference technique is applied (such
 148 that the mean intensity of the new image is effectively zero). The pixel intensities (positive values
 149 only) are then summed along angles and thresholded at a certain number of standard deviations
 150 above the mean intensity. This determines the “core angles” of the CME, and a region growing
 151 technique based on a secondary threshold of intensities in the rest of the image is applied to open
 152 the angular span to include the full CME. An issue arises when streamer deflections occur that
 153 offset the region growing technique and overestimate the CME angular width. An intensity average
 154 across the angles within the span of the CME is then determined, and where the forward portion
 155 of this intensity profile equals half its maximum value is taken as the CME height. The velocity
 156 and acceleration are determined from the heights through consecutive images and these results are
 157 output with the CME position angle and angular width in the SEEDS catalog (see the example
 158 height-time plot at the top right of Fig. 2).

³ <http://sidc.oma.be/cactus/>

⁴ <http://spaceweather.gmu.edu/seeds/>

Along with the issues of streamer deflections and the tracking being limited to only the C2 field-of-view, the choice of the “Half-Max-Lead” as the CME height is dependant on the overall CME brightness, and thus any brightness change during its propagation will affect this measurement. This adds to the error on the height-time profile, which affects the accuracy of the derived velocity and acceleration.

2.4. CDAW Manual Catalog

The CME catalog hosted at the CDAW Data Center⁵ grew out of a necessity to record a simple but effective description and analysis of each event observed with LASCO (Gopalswamy et al., 2009). The catalog is wholly manual in its operation, with a user tracking the CME through C2 and C3 running-difference images and producing a “point-&-click” height-time plot of each event. A linear fit to the height-time profiles provides a 1st-order estimate for the plane-of-sky velocity, and a quadratic fit provides a 2nd-order velocity fit and an acceleration for the event. The central position angle and angular width of the CME are also deduced from the images, and the event is flagged as a halo if it spans 360° , partial halo if it spans $\geq 120^\circ$, and wide if it spans $\geq 60^\circ$. The catalog itself lists each CME’s first appearance in C2, central position angle, angular width, linear speed, 2nd-order speed at final height, 2nd-order speed at $20 R_\odot$, acceleration, mass, kinetic energy, and measurement position angle (the angle along which the heights of the CME are determined; see the example height-time and velocity-height plots in the left of Fig. 2). While the human eye is supremely effective at distinguishing CMEs in coronagraph images, errors may be introduced to the manual cataloging procedure through the biases of different operators; for example, in deciding how the images are scaled, where along the CME the heights are measured, or whether a CME is worth including in, or discarding from, the catalog.

3. CME Event Sample and Catalog Results

A selection of CMEs from the SOHO/LASCO data was chosen in the analysis of Byrne et al. (2009), wherein multiscale methods of edge detection and a resulting ellipse characterisation of the CME front were used to track its apex. These events were chosen based on their varying styles of eruption and appearance, in order to compare with the measurements of the manual CDAW catalog (see images of each CME in Fig. 5 of Byrne et al. 2009). **They exhibit various forms that typical CMEs in coronagraph observations can take, to serve as examples of how the detection and characterisation algorithms fare on each, and how well their varying kinematic trends are revealed.** The images were not differenced to avoid spatiotemporal cross-talk and associated scaling issues. The uncertainties on the height measurements were quantified by the multiscale filter size and subsequent ellipse-fitting, and propagated into the kinematics via numerical differentiation using 3-point Lagrangian interpolation. However, it has since been demonstrated by Byrne et al. (2013) that this method for deriving kinematics is not wholly reliable, and other approaches must be considered as discussed in Sect. 2.1 above. Following the development of the automated CORIMP algorithms, these events are now revisited in this new catalog and directly compared with the other automated CACTus and SEEDS catalogs, the manual CDAW catalog, and the results of Byrne et al.

⁵ http://cdaw.gsfc.nasa.gov/CME_list

CME Date & Start Time	Catalog	CPA [deg.]	Width [deg.]	Lin. Speed [$km\ s^{-1}$]	Accel. [$m\ s^{-2}$]
2000 Jan. 02 ~06:06 UT (Arcade eruption)	CORIMP	250	81 ⁸³	454 ⁷⁴³	1 ¹⁴ ₋₁₇
	CACTus	250	106	548 ⁷⁴⁴ ₂₃₁	
	SEEDS	257	96	292	18.6
	CDAW	253	107	603	21.6
2000 Apr. 18 ~14:54 UT (Gradual CME)	CORIMP	210	98	431 ⁵³⁷	4 ¹⁵ ₋₁₁
	CACTus	198	102	463 ⁷⁴⁴ ₂₂₇	
	SEEDS	195	108	338	17.7
	CDAW	195	105	668	23.1
2000 Apr. 23 ~12:54 UT (Impulsive CME)	CORIMP	287	119 ¹²⁵	836 ¹⁷⁰⁶	-11 ⁵⁰ ₋₁₅₄
	CACTus	144	360	1114 ¹⁸⁴⁹ ₂₄₅	
	SEEDS	275	130	594	-8.5
	CDAW	281	360	1187	-48.5
2001 Apr. 23 ~12:39 UT (Faint CME)	CORIMP	232	72 ⁷⁴	187 ²⁸³	3 ¹⁵ ₋₁₃
	CACTus	231	88	459 ⁶⁰² ₃₁₅	
	SEEDS	224	77	408	-46.6
	CDAW	228	91	530	-0.7
2002 Apr. 21 ~01:26 UT (Fast CME)	CORIMP	235	154 ¹⁷⁷	1129 ²³⁰⁰	61 ³⁴⁵ ₋₆₁₉
	CACTus	322	352	1103 ¹⁹¹³ ₂₉₈	
	SEEDS	250	186	703	31.8
	CDAW	282	360	2393	-1.4
2004 Apr. 1 ~23:04 UT (Slow CME)	CORIMP	58	42 ⁴⁴	401 ⁵⁰²	2 ¹⁸ ₋₂₂
	CACTus	60	70	485 ⁸²⁹ ₂₄₄	
	SEEDS	60	59	261	19.7
	CDAW	59	79	460	7.1

Table 1. Catalog measurements of a sample of CMEs observed by SOHO/LASCO. **CME Date & Start Time** refers to the first observation of the CME in LASCO. **CPA** refers to the central position angle of the CME. **Width** refers to the angular span, or opening angle, of the CME. **Lin. Speed** is the derived speed of the CME using a linear fit to the height-time measurements. **Accel.** is the derived acceleration of the CME using a second-order fit to the height-time measurements. Note, some values have a corresponding maximum and/or minimum (x_{min}^{max}) as specified in the respective catalogs.

¹⁹⁷ (2009). In each case below, the tabled information and kinematic plots are reproduced directly from
¹⁹⁸ the online catalogs, and not rescaled or otherwise manipulated, for a fair comparison.

¹⁹⁹ 3.1. Arcade eruption: 2000 January 2

²⁰⁰ The CME that erupted off the southeast limb of the Sun on 2000 Jan. 02 from ~06:06 UT in LASCO
²⁰¹ exhibited an arcade-type structure consisting of multiple bright loops. CORIMP identified the bulk
²⁰² of the CME through the LASCO field-of-view to $\sim 24 R_\odot$. However, this CME may be deemed the
²⁰³ third in a series of four CMEs that occurred in succession off the southeast limb, that CORIMP failed
²⁰⁴ to separate due to their spatial and temporal overlap (essentially a smaller CME in between two large

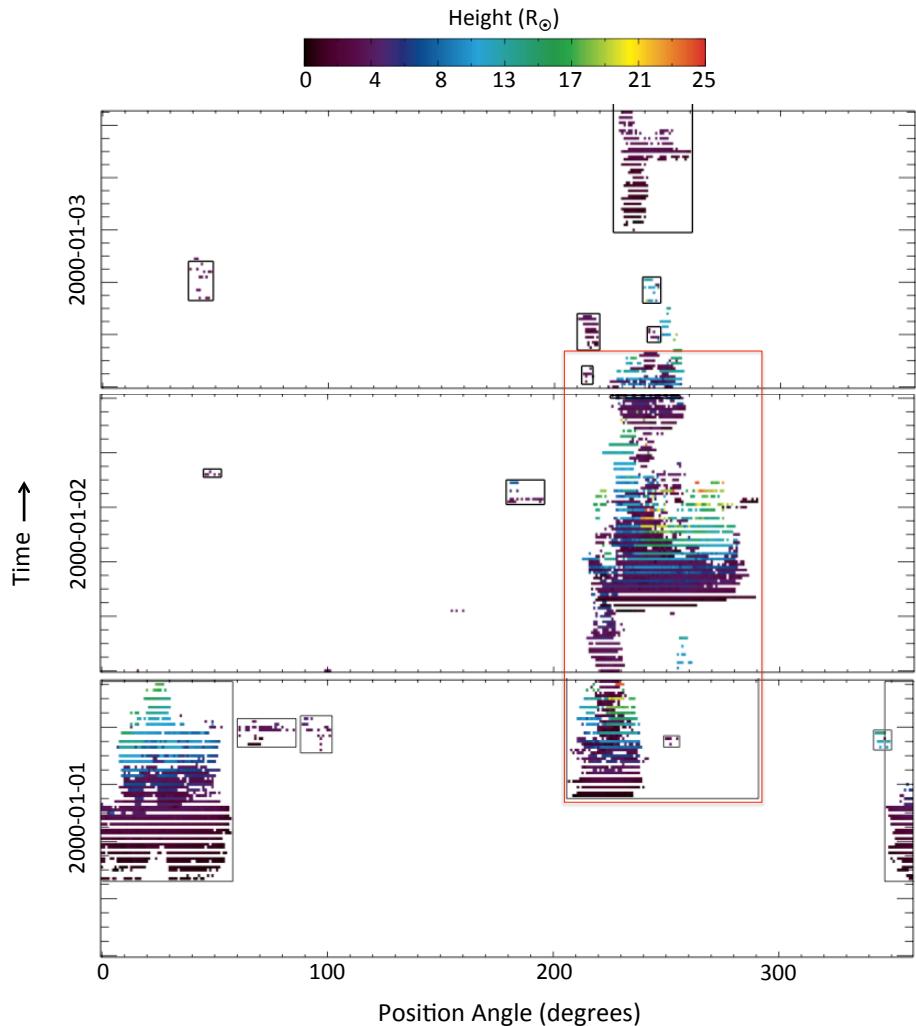


Fig. 3. Daily detection stacks reproduced from the CORIMP CME catalog for the SOHO/LASCO observations in the date range 2000 Jan. 01–03. These stacks are generated from the automatically measured CME front heights at every position angle in an image, stacked in time. CMEs appear as groupings of colour-graded pixels, as indicated by the boxed regions. The overlapping CMEs in this time interval are boxed in red, corresponding to the height-time profiles in Fig. 4 (that have been put through the cleaning algorithm).

ones connects their detections along with a fourth smaller one afterwards, as seen in the CORIMP CME detection stacks reproduced in Fig. 3 - see [Byrne et al. 2012](#) for details). This therefore serves as an example of the need to inspect the catalog output before trusting the quoted values listed in Table 1. The CORIMP height-time measurements (in the time range ~06:00 – 16:00 UT in Fig. 4) reveal a non-linear trend indicative of an early acceleration that the Savitzky-Golay filter determines decreases from a maximum of $\sim 50 \text{ m s}^{-2}$ to 0 m s^{-2} , as the maximum velocity levels off in the range of ~ 500 – 600 km s^{-1} . This is consistent with the measurements of [Byrne et al. \(2009\)](#) shown in their Fig. 6. The quadratic (and linear) fits in CORIMP agree with a maximum velocity in this range of ~ 500 – 600 km s^{-1} and an acceleration in the range of approximately $\pm 20 \text{ m s}^{-2}$. CACTus

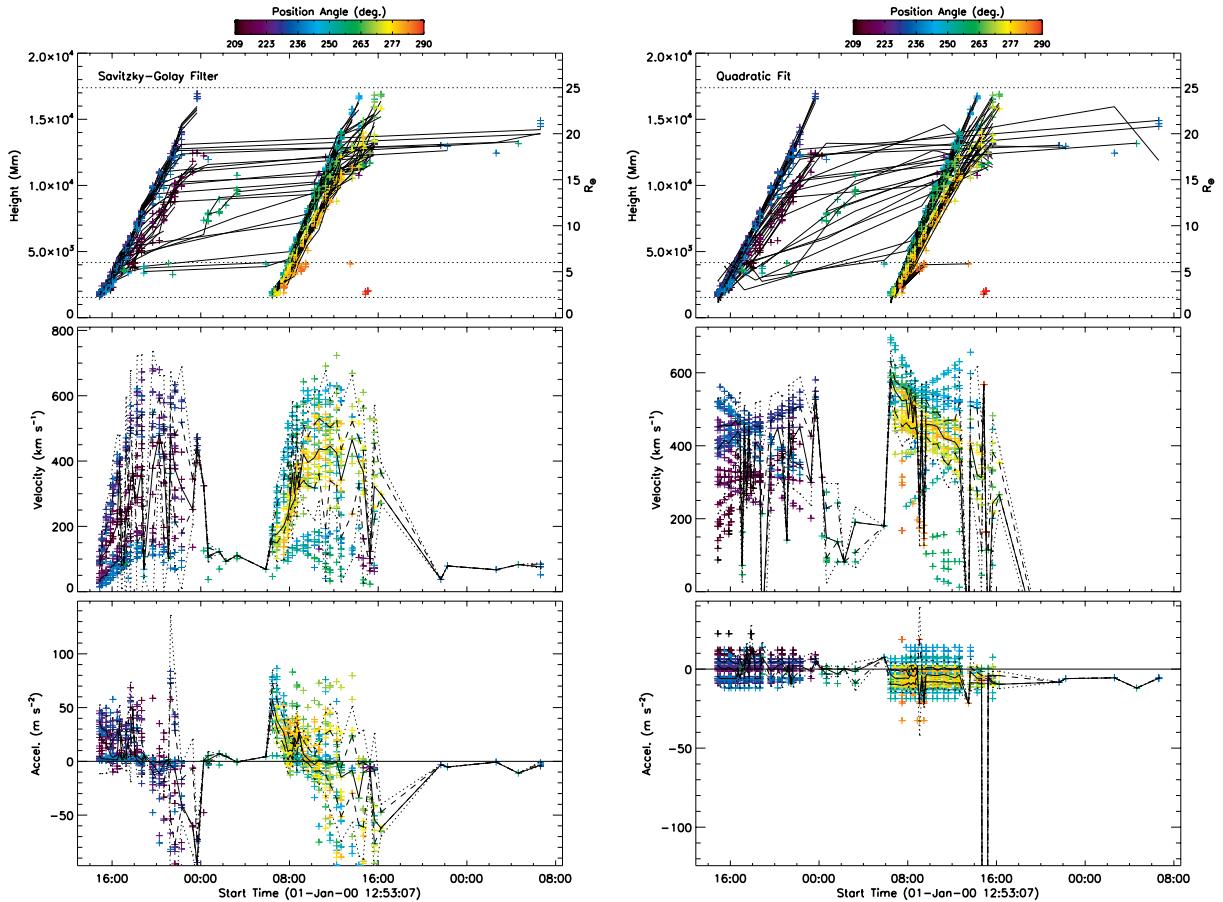


Fig. 4. Kinematic plots of the 2000 Jan. 02 CME from the automatic detection and tracking in the CORIMP catalog. The top plots show the height-time measurements with a colorbar to indicate the angular span of the data points, **and solid black lines to indicate the fitting**. The middle and bottom plots show the velocity and acceleration profiles of the CME with the median (solid line), interquartile range (inner dashed lines) and upper and lower fences (outer dashed lines) over-plotted. The left plots are determined by a Savitzky-Golay filter applied to the height-time measurements with a 7-point moving window, while the right plots are determined with a second-order quadratic fit.

214 determined a linear velocity of 548 km s^{-1} (in the range $231 - 744 \text{ km s}^{-1}$). SEEDS determined a
 215 linear velocity of 292 km s^{-1} and an acceleration of 18.6 m s^{-2} (in the C2 field-of-view). And CDAW
 216 determined a linear velocity of 603 km s^{-1} and an overall acceleration of 21.6 m s^{-2} . These catalog
 217 measurements are listed in Table 1. (Note that the slightly lower angular width in CORIMP is due
 218 to the exclusion of part of the questionable streamer deflection/interaction along the southern flank
 219 of the CME.) **Therefore, by inspection, the results of the CORIMP CME catalog are in relative**
 220 **agreement with the other catalogs and manual analysis of this event, and CORIMP is deemed**
 221 **robust albeit unreliable at separating overlapping events.**

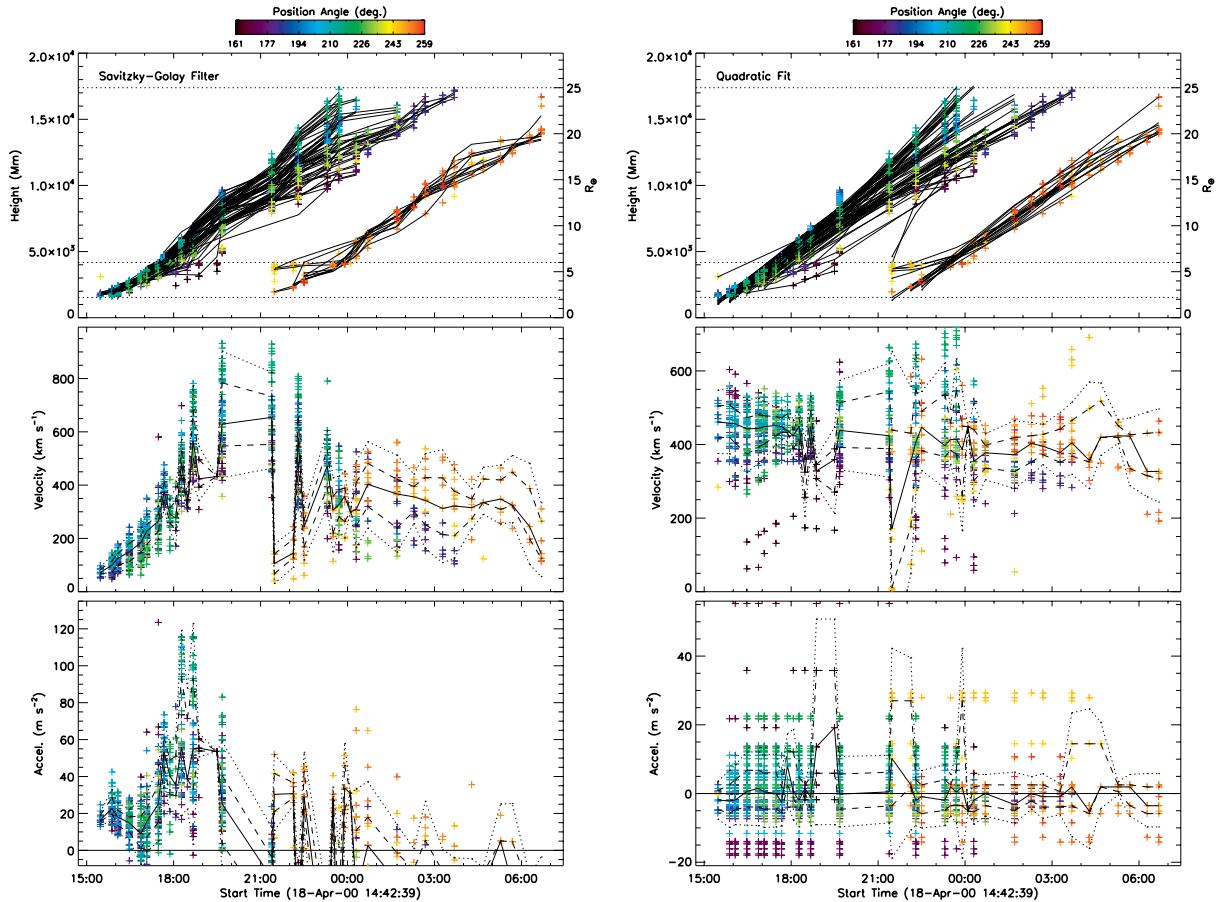


Fig. 5. Kinematic plots of the 2000 Apr. 18 CME from the automatic detection and tracking in the CORIMP catalog, as in Fig. 4.

3.2. Gradual CME: 2000 April 18

The CME that erupted off the south limb of the Sun on 2000 Apr. 18 from ~14:54 UT in LASCO exhibited a typical 3-part structure of leading CME front, cavity and bright core. CORIMP identified the bulk of the CME through the LASCO field-of-view to ~ $25 R_\odot$, though it did not detect a southern portion of the faint CME front in the latter C3 observations. A western portion of material also erupted as a delayed part of the northern flank of the CME, that appears as a somewhat secondary height-time profile in the CORIMP kinematic plots in Fig. 5 (at position angles ~ 250° in the redder end of the colorbar). The CORIMP height-time measurements reveal a non-linear trend indicative of an early acceleration that the Savitzky-Golay filter determines to be approximately $20 m s^{-2}$ as the velocity increases to over $400 km s^{-1}$ before the data gap in the LASCO/C3 images between 19:42 and 21:24 UT causes a large scatter in the derived kinematics (e.g., an artificial acceleration peak of $>100 m s^{-2}$). The initial increasing velocity profile up to a maximum in the range $\sim 600 - 800 km s^{-1}$ by ~20:00 UT agrees with that of Byrne et al. (2009) as shown in their Fig. 7. The quadratic (and linear) fits in CORIMP are not as prone to the scattering effects of the data gap, and thus derive a slightly lower maximum velocity range of $\sim 500 - 550 km s^{-1}$ and an acceleration in the range of

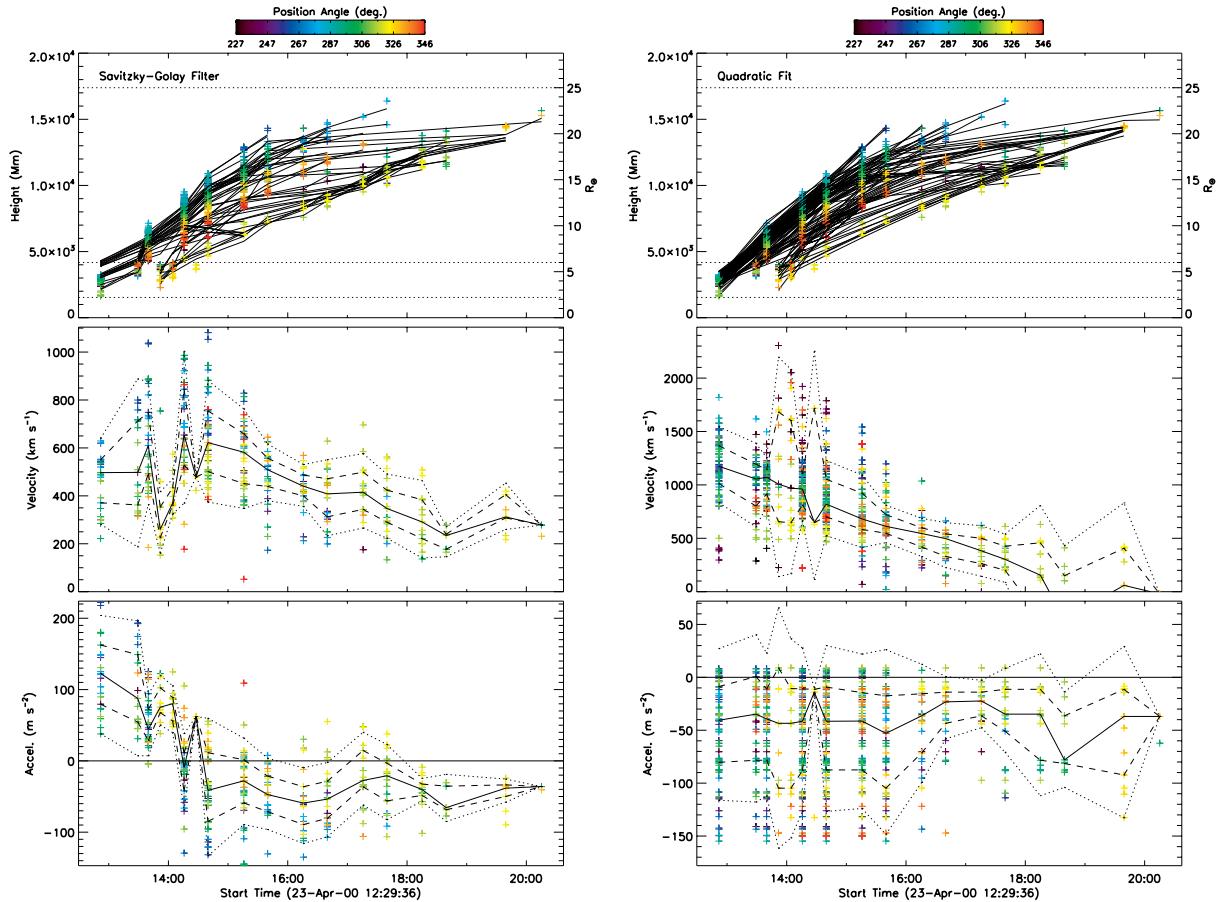


Fig. 6. Kinematic plots of the 2000 Apr. 23 CME from the automatic detection and tracking in the CORIMP catalog, as in Fig. 4.

approximately $\pm 15 \text{ m s}^{-2}$. CACTus determined a linear velocity of 463 km s^{-1} (in the range $227 - 744 \text{ km s}^{-1}$). SEEDS determined a linear velocity of 338 km s^{-1} and an acceleration of 17.7 m s^{-2} (in the C2 field-of-view). And CDAW determined a linear velocity of 668 km s^{-1} and an overall acceleration of 23.1 m s^{-2} . Therefore, by inspection, all sets of results are in relative agreement for this event.

3.3. Impulsive CME: 2000 April 23

The large and fast CME that erupted off the west limb of the Sun on 2000 Apr. 23 from $\sim 12:54$ UT in LASCO underwent a hugely impulsive acceleration as it exploded into the corona. CORIMP identified the bulk of the CME through the LASCO field-of-view to $\sim 20 R_\odot$ after which the CME front became too faint. Strong streamer deflections occurred to the north and south flanks of the CME, with very faint material visible as a full halo or shock around the east limb separate to the bulk flux-rope structure in the west. The CORIMP height-time measurements (Fig. 6) reveal an initial acceleration that the Savitzky-Golay filter determines to be $\gtrsim 150 \text{ m s}^{-2}$ dropping quickly to a range of approximately -100 to 0 m s^{-2} , as the velocity decreases from ~ 1000 to 500 km s^{-1} ; though this

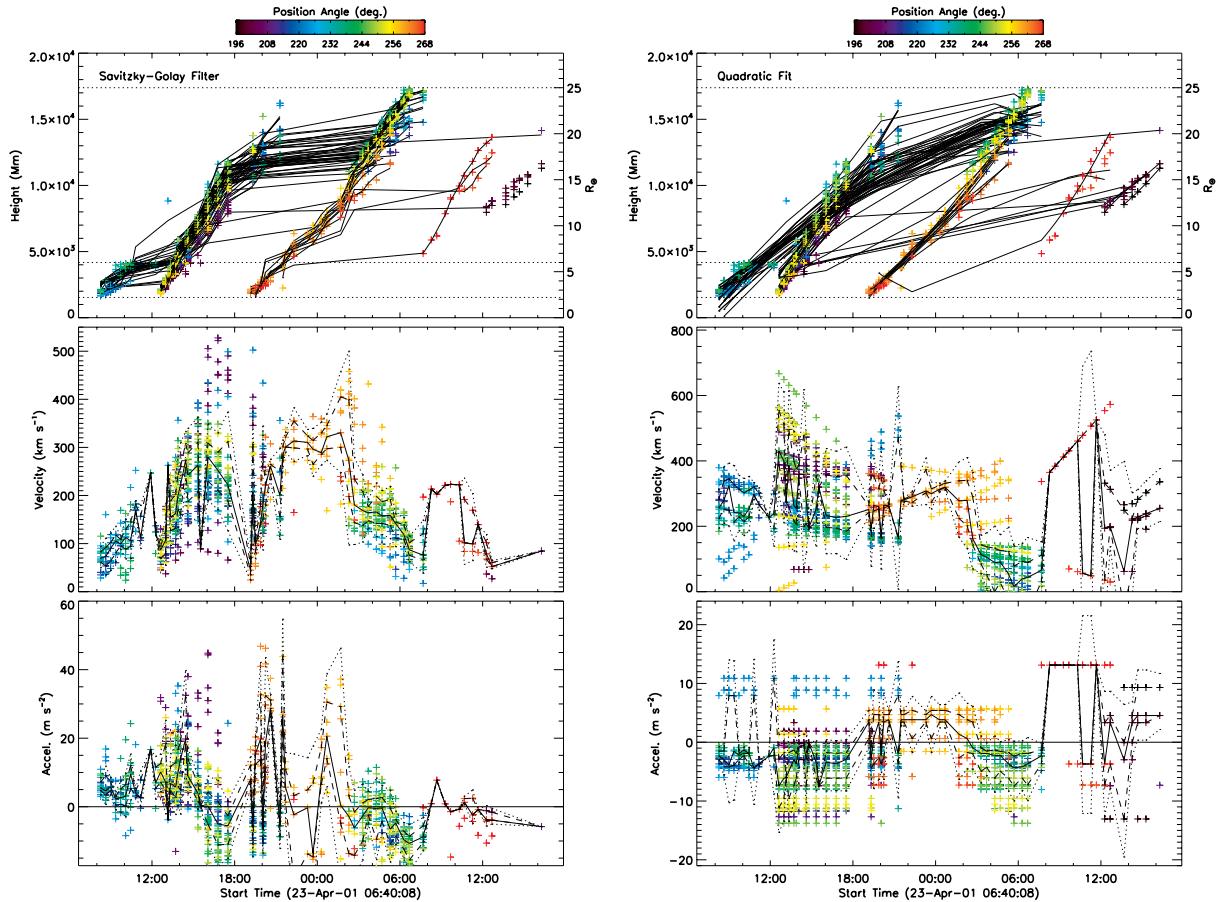


Fig. 7. Kinematic plots of the 2001 Apr. 23 CME from the automatic detection and tracking in the CORIMP catalog, as in Fig. 4.

is an underestimate since the filter overly smoothes the relatively under-sampled height-time measurements. The quadratic fits in CORIMP better handle this data and derive an initial velocity range of $\sim 1200 - 1500 \text{ km s}^{-1}$, while the linear fits derive an initial velocity range of $\sim 1000 - 1200 \text{ km s}^{-1}$, which are consistent with the measurements of Byrne et al. (2009) shown in their Fig. 8. The resulting deceleration is determined to have a median of approximately -50 m s^{-2} , reaching as low as -150 m s^{-2} . CACTus determined a linear velocity of 1114 km s^{-1} (in the range $245 - 1849 \text{ km s}^{-1}$). SEEDS determined a linear velocity of 594 km s^{-1} and a deceleration of -8.5 m s^{-2} (in the C2 field-of-view). And CDAW determined a linear velocity of 1187 km s^{-1} and an overall deceleration of -48.5 m s^{-2} . **Therefore, by inspection and careful consideration of the low sampling of the event, the results of the CORIMP CME catalog are in relative agreement with the corresponding results of the other catalogs and manual analysis.**

3.4. Faint CME: 2001 April 23

The CME that erupted off the southwest limb of the Sun on 2001 Apr. 23 from $\sim 12:54$ UT in LASCO appeared relatively faint behind multiple streamers in the line-of-sight, some of which de-

flected especially along the southern flank of the CME. CORIMP identified the bulk of the CME through the LASCO field-of-view to $\sim 20 R_{\odot}$ after which the CME front became too faint. However, this CME was the first of two that occurred in close succession off the southwest limb, that CORIMP failed to separate due to their spatial and temporal overlap (plus some ejecta ahead of this CME was detected from $\sim 08:16$ UT). Therefore the kinematic profiles must be inspected before trusting the quoted catalog values listed in Table 1. Investigating the relevant portion of the plots in Fig. 7, in the time interval $\sim 12:00 - 18:00$ UT, the CORIMP height-time measurements reveal an initial acceleration that the Savitzky-Golay filter determines to be $\sim 10 m s^{-2}$ dropping to scatter about zero as the maximum velocity levels off at $\sim 350 km s^{-1}$; though this is an underestimate since the measurements are dominated by the material ahead of the CME front that the algorithm detected as part of the main event (from $\sim 08:00 - 12:00$ UT). The quadratic fits to the measurements are more dominated by the overall deceleration of the CME (approx. $-10 m s^{-2}$ from 12:00 UT onwards) as the velocity drops from an initial range of $\sim 550 - 650 km s^{-1}$ (consistent with Byrne et al. 2009 shown in their Fig. 9) to $\sim 400 km s^{-1}$ by 18:00 UT; though this appears biased to lower values by the overlapping measurements of the second CME. The linear fits are less trustworthy as they tend to fit across the two CMEs and preceding ejected material, resulting in the underestimated CORIMP linear speed in Table 1. CACTus determined a linear velocity of $459 km s^{-1}$ (in the range $315 - 602 km s^{-1}$). SEEDS determined a linear velocity of $408 km s^{-1}$ and a deceleration of $-46.6 m s^{-2}$ (in the C2 field-of-view), however it failed to detect the CME front in the final frames which accounts for this erroneously large deceleration. And CDAW determined a linear velocity of $530 km s^{-1}$ and an overall deceleration of $-0.7 m s^{-2}$. **Therefore, while all sets of results are found to be in relative agreement, there is again the issue of separating overlapping event kinematics.**

3.5. Fast CME: 2002 April 21

The CME that erupted off the west limb of the Sun on 2002 Apr. 21 from $\sim 01:27$ UT in LASCO propagated very fast through the field-of-view. CORIMP identified the bulk of the CME through the LASCO field-of-view to $\sim 17 R_{\odot}$ after which the CME front became too faint, and only the southern flank material continued to be detected. Figure 8 shows the CORIMP height-time measurements, which the Savitzky-Golay filter struggles to fit appropriately due to the small window-size available at each position angle (as the filter requires a minimum of 7 data points). The quadratic fits to the data reveal a high initial acceleration of $\gtrsim 1000 m s^{-2}$ followed by a deceleration in the range of approximately -500 to $0 m s^{-2}$. The velocity shows an initial range of $\sim 2000 - 2500 km s^{-1}$ possibly reaching $\sim 3000 km s^{-1}$ before dropping to $\sim 1000 km s^{-1}$, which is consistent with the measurements of Byrne et al. (2009) shown in their Fig. 10. The linear fits also reveal an initial velocity range of $\sim 2000 - 2500 km s^{-1}$ dropping to $\sim 1000 km s^{-1}$. CACTus determined a linear velocity of $1103 km s^{-1}$ (in the range $298 - 1913 km s^{-1}$). SEEDS determined a linear velocity of $703 km s^{-1}$ and an acceleration of $31.8 m s^{-2}$ (in the C2 field-of-view), however these cannot be trusted as the CME front is only visible in two C2 frames. And CDAW determined a linear velocity of $2393 km s^{-1}$ and an overall deceleration of $-1.4 m s^{-2}$. **Therefore, CORIMP is deemed reliable albeit problematic at handling low-sampled events (with the Savitzky-Golay filter anyway, as the quadratic and linear fits remain robust).**

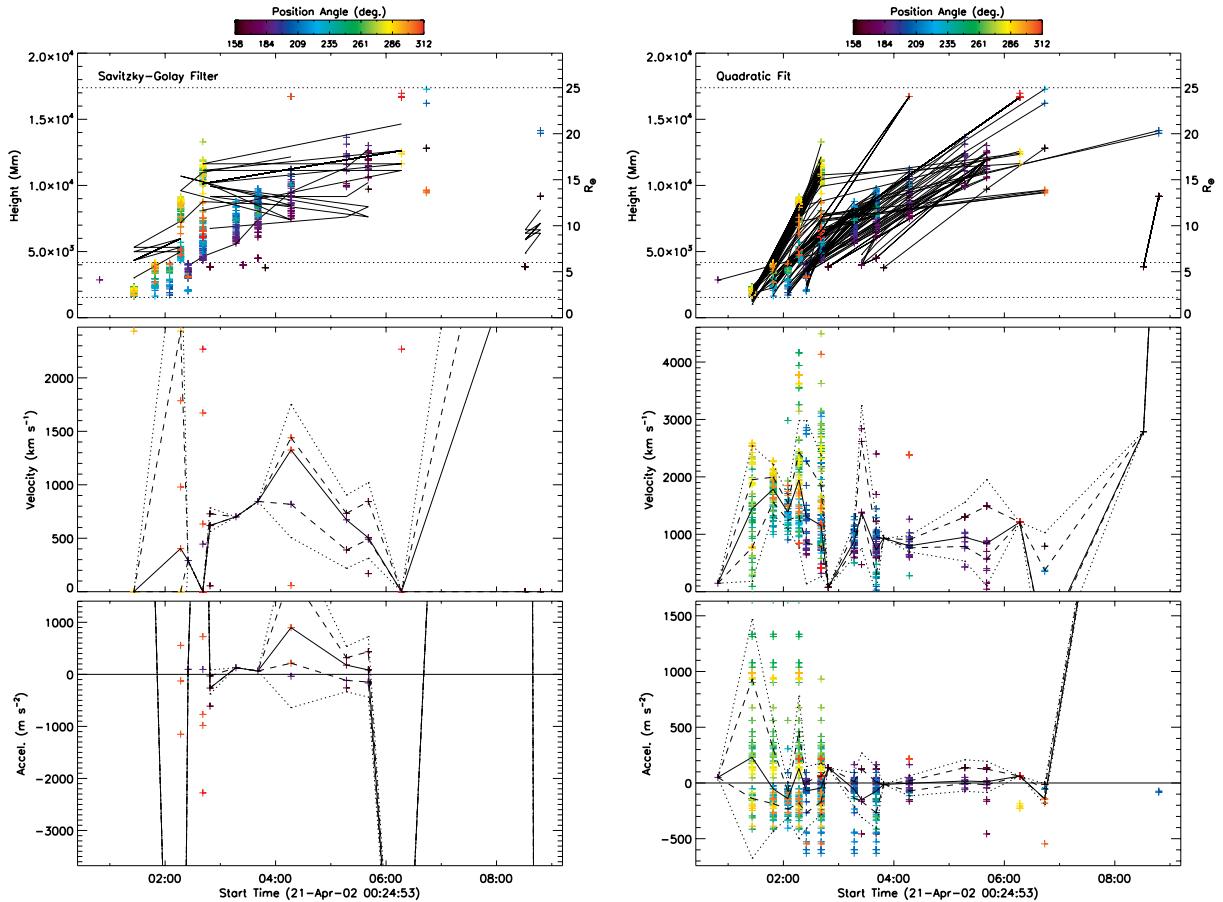


Fig. 8. Kinematic plots of the 2002 Apr. 21 CME from the automatic detection and tracking in the CORIMP catalog, as in Fig. 4.

305 3.6. Slow CME: 2004 April 1

306 The CME that erupted off the northeast limb of the Sun on 2004 Apr. 01 from $\sim 23:05$ UT in
 307 LASCO, exhibited a clear flux-rope structure and propagated relatively slowly. CORIMP identi-
 308 fied the bulk of the CME through the LASCO field-of-view to $\sim 20 R_\odot$ after which the CME front
 309 became too faint. Figure 9 shows the CORIMP height-time measurements, which are plentiful given
 310 the slow motion and clean detection of the event. These measurements reveal an initial acceleration
 311 that the Savitzky-Golay filter determines to be $\gtrsim 25 m s^{-2}$ dropping to $0 m s^{-2}$ by the time the CME
 312 reaches $\sim 15 R_\odot$ and the maximum velocity levels off in the range $\sim 500 - 600 km s^{-1}$. The quadratic
 313 fits to the data reveal a bulk velocity in the range $\sim 400 - 600 km s^{-1}$, with an overall decelera-
 314 tion of the CME of approximately $-5 m s^{-2}$. The linear fits also produce a velocity in this range.
 315 These results are consistent with the measurements of Byrne et al. (2009) shown in their Fig. 11,
 316 though without reproducing the “staggered” velocity profile. CACTus determined a linear velocity
 317 of $485 km s^{-1}$ (in the range $244 - 829 km s^{-1}$). SEEDS determined a linear velocity of $261 km s^{-1}$
 318 and overall acceleration of $19.7 m s^{-2}$ (in the C2 field-of-view). And CDAW determined a linear ve-

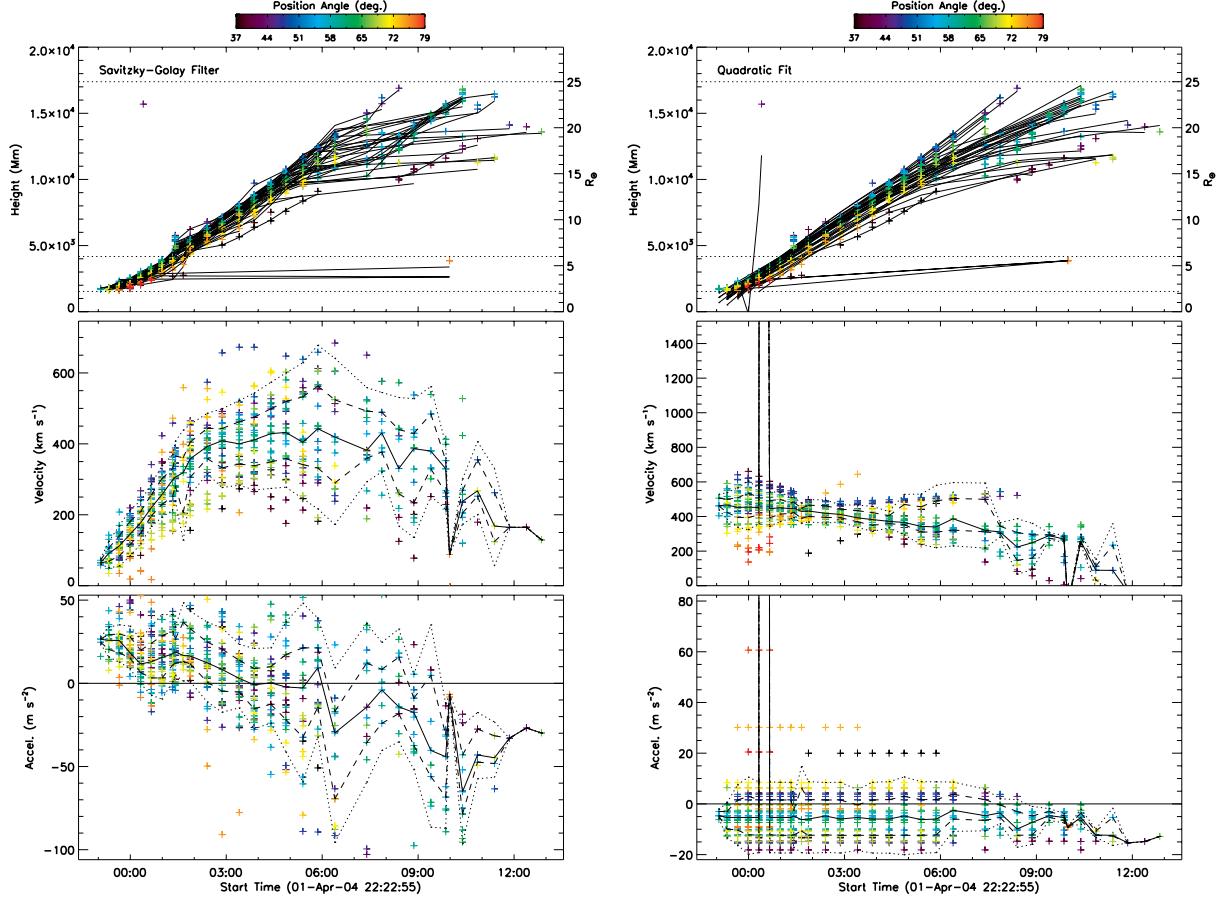


Fig. 9. Kinematic plots of the 2004 Apr. 01 CME from the automatic detection and tracking in the CORIMP catalog, as in Fig. 4.

319 locity of 460 km s^{-1} and an overall acceleration of 7.1 m s^{-2} . Therefore, by inspection, CORIMP
 320 and the other CME characterisations are in relative agreement for this event.

321 4. Separating Multiple CME Detections via K-means Clustering

322 The case studies in the previous section highlight a key issue in the automatic detection, tracking
 323 and cataloging of CMEs: namely the difficulties in distinguishing between multiple events that oc-
 324 cur close together in space and time. The events of Sect. 3.1 and 3.4 demonstrate how the CORIMP
 325 catalog can fail to separately characterise CMEs that a human user would label as distinct events
 326 (although even this can be a non-trivial task since projection effects can make it hard to determine
 327 if two CMEs are truly merging in space or simply overlapping on the plane-of-sky). The reason for
 328 this, is that the thresholds in place to identify the beginning of a new CME detection cannot read-
 329 ily determine the end of a previous CME detection whose trailing material overlaps the subsequent
 330 CME material on the plane-of-sky. Other observational factors can help a human user distinguish the
 331 two, such as the differing speeds, densities, brightness and cohesiveness of the structure. However,

these differences are too subtle to employ in an automated algorithm that must be able to characterise all manner of CMEs with a large variety of such properties; and so the overlapping CMEs are classified as a single event. While it is still possible to investigate their separate kinematic trends, as in Sect. 3.1 and 3.4, this is not ideal for accurately counting CMEs nor for reliably producing independent CME detection alerts. It therefore remains a challenge to employ some form of machine intelligence in cataloging CMEs, which can use the CME detection parameters to reveal instances when multiple CMEs occur together.

An initial effort to do this has been made using a clustering algorithm in the field of unsupervised machine learning. Specifically, the method of k -means clustering was investigated, which works by partitioning n observations into k clusters that are distinguished by minimising the within-cluster sum of squares, i.e., using Euclidean distance as a metric on the parameter space. It is well suited to generating globular, non-hierarchical, non-overlapping clusters, and may be computationally fast if k remains small. This approach could work with the parameters available in the CME detection analysis, such as the time, location/direction, size and speed of a CME. For example, the bulk of the overlapping CMEs may have different average propagation times as one may be proceeding slightly later in time than another. Or if the CMEs are propagating at different speeds, then a linear fit to each of their height-time profiles would have a different slope. Choosing these parameters of “average time” and “slope of a linear fit”, determined at every position angle in the span of the event detection, it is possible to cluster the height-time measurements into separate CMEs.

This is demonstrated in Figs. 10 and 11 for the two events discussed in Sects. 3.1 and 3.4, respectively. These figures show top plots of the k -means clustering algorithm (**where k is manually prescribed by the user**) applied to the parameters of “average time” and “slope of a linear fit”, where the means are plotted as black asterisks and the associated groups of points in separate clusters are plotted with different symbols and colours. The bottom plots of these figures show the corresponding height-time profiles that have been separated according to their clusters. In Fig. 10 the results are shown for both $k = 3$ (left plots) and $k = 4$ (right plots), to illustrate the effect of changing k . By inspection, the clustering algorithm works well at distinguishing the separate events. The bulk measurements of the case-study CME beginning at ~06:06 UT on 2000 Jan. 02 are quite well clustered, though some of its later C3 measurements are wrongly determined as part of a separate CME. For this event the $k = 3$ case fares better at grouping the CMEs, while the $k = 4$ case splits apart the profile of the first CME (shown as the green and blue datapoints in the bottom right plot of Fig. 10). Similarly in Fig. 11 the clustering algorithms go some way towards distinguishing multiple CMEs in the event detection on 2001 Apr. 23, but there are datapoints that are wrongly classified, even for the two instances of $k = 4$ shown for this event. These results highlight the difficulty in applying an automatic extraction of separate CME height-time profiles when detected so close together in space and time. Furthermore, there is an inherent limitation to k -means clustering by having to specify the number of clusters required from the data, which is not known a priori - especially not for an automated methodology such as in the CORIMP catalog. Further investigation into the parameters to be clustered, alternative clustering, or different machine learning algorithms, may produce better results.

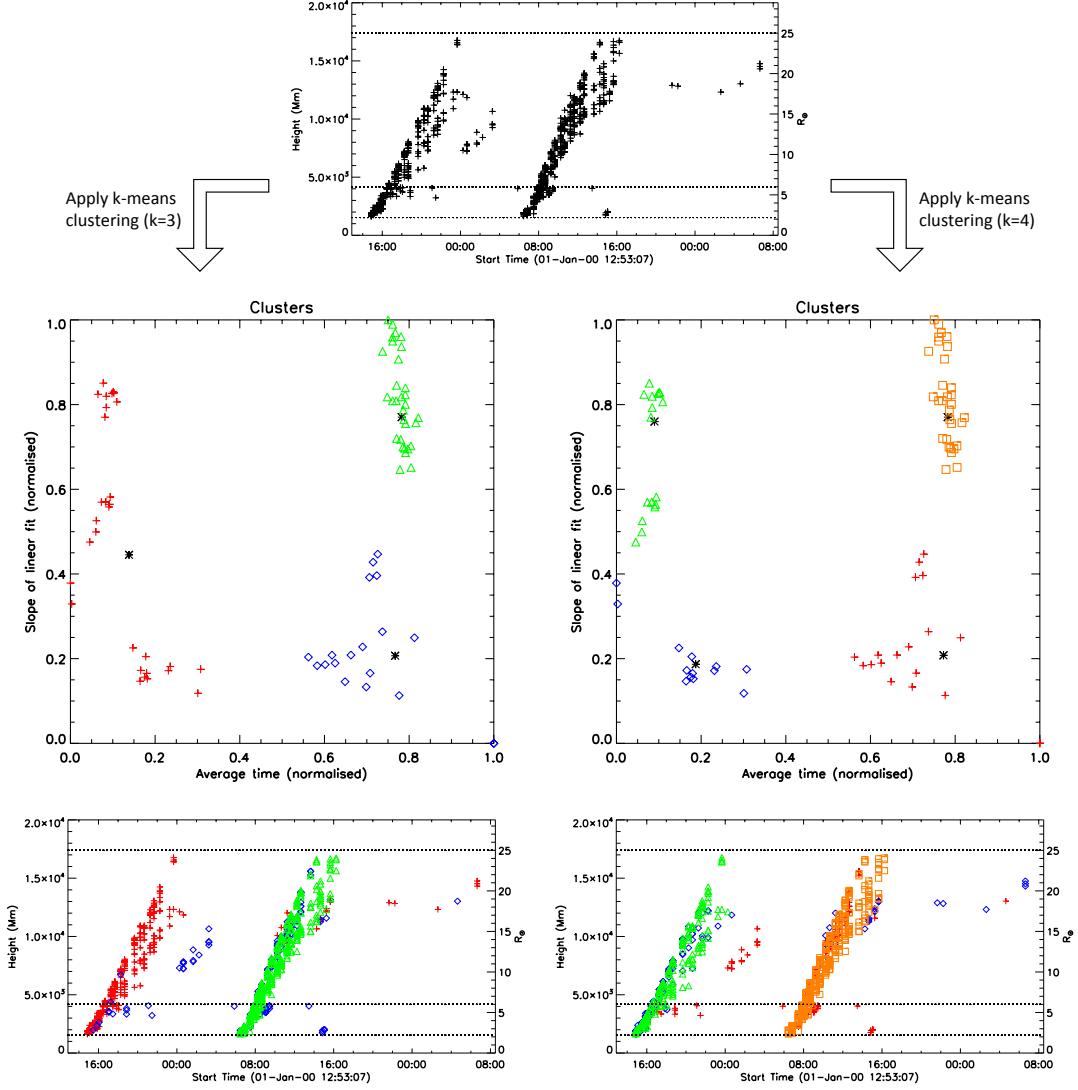


Fig. 10. A k -means clustering algorithm applied to the height-time data on 2000 Jan. 01–03 (from Fig. 4), in an effort to distinguish the multiple CME profiles that were detected as a single event due to their close proximity in space and time. *Top plot:* The height-time measurements of the CME detections from the automated CORIMP catalog. (Note, this dataset has been put through a cleaning algorithm, discussed in Sect. 2.1, that removes a lot of inner-core and trailing-material datapoints, thus making it easier to distinguish their separate profiles.) *Middle plots:* The resulting clusters for the cases of $k = 3$ (left) and $k = 4$ (right), applied to the normalised parameters of the slope of a linear fit to, and the mean time of, the height-time profile at each position angle. The clusters are distinguished by different plot symbols and colours, with black asterisks to indicate the mean of each cluster. *Bottom plots:* The resulting effort at separately distinguishing the CME height-time profiles.

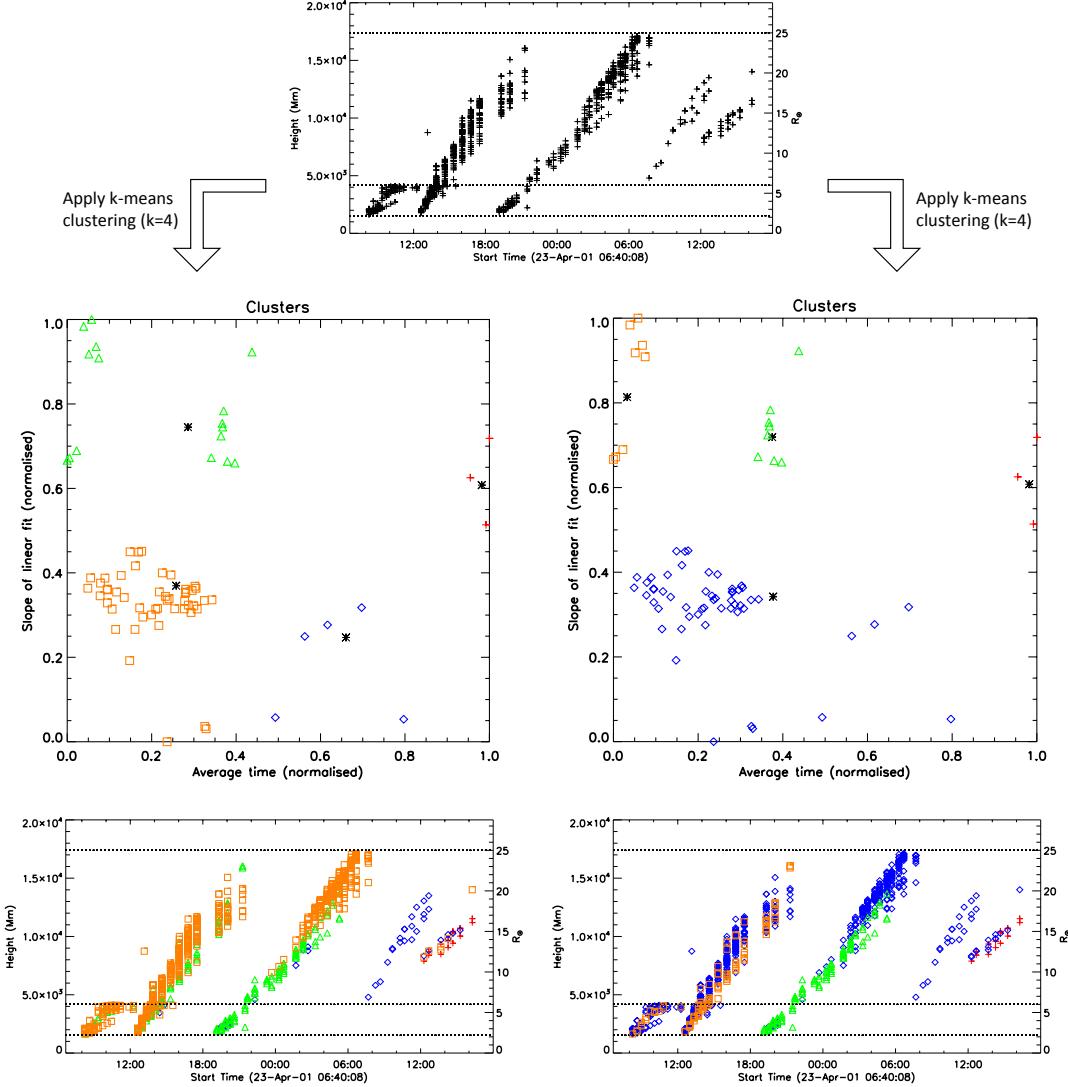


Fig. 11. A k -means clustering algorithm applied to the height-time data on 2001 Apr. 23–24 (from Fig. 7), in an effort to distinguish the multiple CME profiles, as in Fig. 10 but for two cases of $k = 4$.

372 5. Conclusions

373 As the wealth of coronagraph data and CME observations has increased dramatically since the
 374 launch of SOHO in 1995, it has become important to develop robust and reliable methods of de-
 375 tecting and tracking CMEs in white-light images. Since CMEs are faint and transient phenomena
 376 that prove difficult to consistently isolate from the background corona, manual inspection of the
 377 images is open to interpretation and prone to user-specific biases. Similarly, it is challenging to fix
 378 the criteria and thresholds necessary in a computerised methodology for automating this task, al-
 379 though advances have been made to achieve this and provide the benefit of having a self-consistent
 380 catalog of results. Efforts to both manually and automatically catalog CMEs have been discussed in
 381 Sect. 2 with the aim of comparing how each fares in light of the newly developed CORIMP catalog,

which was built to overcome some of the drawbacks of current catalogs. To this end, a selection of CMEs was chosen from a previous study by [Byrne et al. \(2009\)](#), and the new results in the CORIMP catalog were investigated alongside the results of the automated CACTus and SEEDS catalogs and the manual CDAW catalog.

In the previous study of [Byrne et al. \(2009\)](#), the CMEs were characterised with the use of a multiscale edge-detection filter, whereby an ellipse was fitted to the isolated CME front and its apex tracked to produce height-time measurements. Since this approach avoided differencing the images, it was possible to quantify single-image uncertainties for the resulting height-time measurements, to be used for gauging a confidence interval on the derived CME kinematics. However, [Byrne et al. \(2013\)](#) demonstrated that the often-used method of numerical differentiation using 3-point Lagrangian interpolation, and its associated error propagation, is not wholly reliable at deriving the true CME kinematics. This motivated the use of the Savitzky-Golay filter along with quadratic and linear fits to the height-time measurements in CORIMP, across the angular span of the CME such that the statistical spread in the kinematics of each event may better indicate the true underlying trends. It is therefore warranted to compare these new automatically-generated results with the outputs of the other catalogs.

The spread of measurements along the angular span of the CME proves more useful than choosing a single fixed apex of the CME, because it propagates as an impermanent, evolving structure that can undergo various rates of expansion across the plane-of-sky. The variety of events chosen here as a subset of the thousands in the LASCO data is enough to demonstrate this. Having the angular spread of kinematics also provides insight to the bulk motion of the CME as well as its flanks and front: with the angular extent indicating the flanks and the upper values on the velocities indicating the CME front (usually the fastest part of its structure). Therefore a greater amount of information is available on the overall CME motion.

The Savitzky-Golay filter provides an indication of the kinematic trends that a first or second-order fit cannot necessarily produce. Since this filter is applied in a moving-window on the data points, it can be problematic in cases of events with low-sampling (as in Sect. 3.5), but otherwise performs very well at automatically quantifying the different phases of acceleration of a CME. Therefore the dynamics of the eruption may be better quantified and understood.

While the robustness of the CORIMP catalog is clear (in so far as it can demonstrably produce results that are accurate and consistent across the data), there is a reliability issue that arises in cases of multiple CMEs that overlap in space and time. The problem with such cases is that another CME can erupt in the same direction as a previous one, close enough in time that the two detections are merged, as though the second CME were part of the trailing material of the first (as shown in Fig. 3). The opposite problem to this is that harsher thresholds would split apart single CMEs into multiple events, especially large CMEs with substantial trailing material. Indeed such problems can affect all automated catalogs, such that CORIMP appears to suffer from the former issue, while CACTus and SEEDS suffer from the latter. A form of unsupervised machine learning was explored, by applying a k -means clustering algorithm to certain parameters in the CME detections, in an effort to distinguish overlapping events. While these first results shown in Sect. 4 are promising, they highlight the difficulty of the task, and warrant further investigation. **For example, perhaps a form of supervised machine learning would fare better, if a substantial training set of correctly labelled event data was produced from the current database of results and used to train an intelligent algorithm.** For now, this issue is only overcome by a manual inspection of the data,

as highlighted in the events of Sect. 3.1 and 3.4. In conclusion, any catalog should not be quoted blindly, as the thresholds cannot always distinguish the exact eruption that a user would isolate by eye. However, knowing this, CORIMP still offers the most rigorous details on the kinematics and morphologies of CMEs in a catalog to date, from which a user can infer a wealth of information.

Acknowledgements. The SOHO/LASCO data used here are produced by a consortium of the Naval Research Laboratory (USA), Max-Planck-Institut fuer Aeronomie (Germany), Laboratoire d’Astronomie (France), and the University of Birmingham (UK). SOHO is a project of international cooperation between ESA and NASA. The CACTus CME catalog is generated and maintained by the SIDC at the Royal Observatory of Belgium. The SEEDS CME catalog has been supported by NASA Living With a Star Program and NASA Applied Information Systems Research Program. The CDAW Data Center CME catalog is generated and maintained by NASA and The Catholic University of America in cooperation with the Naval Research Laboratory. The author would like to thank the anonymous referees for their helpful comments.

References

- Brueckner, G. E., R. A. Howard, M. J. Koomen, C. M. Korendyke, D. J. Michels, et al. The Large Angle Spectroscopic Coronagraph (LASCO). *Sol. Phys.*, **162**, 357–402, 1995. 10.1007/BF00733434. [1](#)
- Byrne, J. P., P. T. Gallagher, R. T. J. McAteer, and C. A. Young. The kinematics of coronal mass ejections using multiscale methods. *A&A*, **495**, 325–334, 2009. 10.1051/0004-6361:200809811, [0901.3392](#). [3](#), [3.1](#), [3.2](#), [3.3](#), [3.4](#), [3.5](#), [3.6](#), [5](#)
- Byrne, J. P., D. M. Long, P. T. Gallagher, D. S. Bloomfield, S. A. Maloney, R. T. J. McAteer, H. Morgan, and S. R. Habbal. Improved methods for determining the kinematics of coronal mass ejections and coronal waves. *A&A*, **557**, A96, 2013. 10.1051/0004-6361/201321223, [1307.8155](#). [2.1](#), [3](#), [5](#)
- Byrne, J. P., S. A. Maloney, R. T. J. McAteer, J. M. Refojo, and P. T. Gallagher. Propagation of an Earth-directed coronal mass ejection in three dimensions. *Nature Communications*, **1**, 2010. 10.1038/ncomms1077, [1010.0643](#). [1](#)
- Byrne, J. P., H. Morgan, S. R. Habbal, and P. T. Gallagher. Automatic Detection and Tracking of Coronal Mass Ejections. II. Multiscale Filtering of Coronagraph Images. *ApJ*, **752**, 145, 2012. 10.1088/0004-637X/752/2/145. [1](#), [3.1](#)
- Carley, E. P., D. M. Long, J. P. Byrne, P. Zucca, D. S. Bloomfield, J. McCauley, and P. T. Gallagher. Quasiperiodic acceleration of electrons by a plasmoid-driven shock in the solar atmosphere. *Nature Physics*, **9**, 811–816, 2013. 10.1038/nphys2767. [1](#)
- Chen, P. F. Coronal Mass Ejections: Models and Their Observational Basis. *Living Reviews in Solar Physics*, **8**, 1, 2011. [1](#)
- Colaninno, R. C., and A. Vourlidas. Analysis of the Velocity Field of CMEs Using Optical Flow Methods. *ApJ*, **652**, 1747–1754, 2006. 10.1086/507943. [1](#)
- Davis, C. J., J. A. Davies, M. Lockwood, A. P. Rouillard, C. J. Eyles, and R. A. Harrison. Stereoscopic imaging of an Earth-impacting solar coronal mass ejection: A major milestone for the STEREO mission. *Geophys. Res. Lett.*, **36**, 8102, 2009. 10.1029/2009GL038021. [1](#)

- 463 Domingo, V., B. Fleck, and A. I. Poland. The SOHO Mission: an Overview. *Sol. Phys.*, **162**, 1–2, 1995.
464 10.1007/BF00733425. [1](#)
- 465 Gallagher, P. T., C. A. Young, J. P. Byrne, and R. T. J. McAteer. Coronal mass ejection detection using
466 wavelets, curvelets and ridgelets: Applications for space weather monitoring. *Advances in Space Research*,
467 **47**, 2118–2126, 2011. 10.1016/j.asr.2010.03.028, [1012.1901](#). [1](#)
- 468 Gopalswamy, N., S. Yashiro, G. Michalek, G. Stenborg, A. Vourlidas, S. Freeland, and R. Howard. The
469 SOHO/LASCO CME Catalog. *Earth Moon and Planets*, **104**, 295–313, 2009. 10.1007/s11038-008-9282-
470 7. [1](#), [2.4](#)
- 471 Goussies, N. A., M. E. Mejail, J. Jacobo, and G. Stenborg. Detection and Tracking of Coronal Mass
472 Ejections Based on Supervised Segmentation and Level Set. *Pattern Recogn. Lett.*, **31**(6), 496–501, 2010.
473 10.1016/j.patrec.2009.07.011, URL <http://dx.doi.org/10.1016/j.patrec.2009.07.011>. [1](#)
- 474 Hough, P. V. C. A Method and Means for Recognizing Complex Patterns. US Patent: 3,069,654, 1962. [1](#)
- 475 Howard, R. A., J. D. Moses, A. Vourlidas, J. S. Newmark, D. G. Socker, et al. Sun Earth Connection Coronal
476 and Heliospheric Investigation (SECCHI). *Space Science Reviews*, **136**, 67–115, 2008. 10.1007/s11214-
477 008-9341-4. [1](#)
- 478 Howard, T. A., and S. J. Tappin. Statistical survey of earthbound interplanetary shocks, associated coro-
479 nal mass ejections and their space weather consequences. *A&A*, **440**, 373–383, 2005. 10.1051/0004-
480 6361:20053109. [1](#)
- 481 Hundhausen, A. J. Sizes and locations of coronal mass ejections - SMM observations from 1980 and 1984-
482 1989. *J. Geophys. Res.*, **98**, 13,177, 1993. 10.1029/93JA00157. [1](#)
- 483 Illing, R. M. E., and A. J. Hundhausen. Observation of a coronal transient from 1.2 to 6 solar radii.
484 *J. Geophys. Res.*, **90**, 275–282, 1985. 10.1029/JA090iA01p00275. [1](#)
- 485 Kilpua, E. K. J., J. Pomoell, A. Vourlidas, R. Vainio, J. Luhmann, Y. Li, P. Schroeder, A. B. Galvin, and
486 K. Simunac. STEREO observations of interplanetary coronal mass ejections and prominence deflection
487 during solar minimum period. *Annales Geophysicae*, **27**(12), 4491–4503, 2009. URL <http://www.ann-geophys.net/27/4491/2009/>. [1](#)
- 489 Koomen, M. J., C. R. Detwiler, G. E. Brueckner, H. W. Cooper, and R. Tousey. White Light Coronagraph
490 in OSO-7. *Applied Optics*, **14**, 743–751, 1975. URL <http://www.opticsinfobase.org/abstract.cfm?URI=ao-14-3-743>. [1](#)
- 492 Liu, Y. D., J. G. Luhmann, P. Kajdič, E. K. J. Kilpua, N. Lugaz, et al. Observations of an extreme storm in
493 interplanetary space caused by successive coronal mass ejections. *Nature Communications*, **5**, 3481, 2014.
494 10.1038/ncomms4481. [1](#)
- 495 Lugaz, N., and P. Kintner. Effect of Solar Wind Drag on the Determination of the Properties of Coronal Mass
496 Ejections from Heliospheric Images. *Sol. Phys.*, **47**, 2012. 10.1007/s11207-012-9948-1, [1204.3813](#). [1](#)
- 497 MacQueen, R. M., A. Csoeke-Poeckh, E. Hildner, L. House, R. Reynolds, A. Stanger, H. Tepoel, and
498 W. Wagner. The High Altitude Observatory Coronagraph/Polarimeter on the Solar Maximum Mission.
499 *Sol. Phys.*, **65**, 91–107, 1980. 10.1007/BF00151386. [1](#)

- 500 Morgan, H., J. P. Byrne, and S. R. Habbal. Automatically Detecting and Tracking Coronal Mass Ejections.
501 I. Separation of Dynamic and Quiescent Components in Coronagraph Images. *ApJ*, **752**, 144, 2012.
502 10.1088/0004-637X/752/2/144. [1](#)
- 503 Olmedo, O., J. Zhang, H. Wechsler, A. Poland, and K. Borne. Automatic Detection and Tracking of Coronal
504 Mass Ejections in Coronagraph Time Series. *Sol. Phys.*, **248**, 485–499, 2008. 10.1007/s11207-007-9104-
505 5. [1](#)
- 506 Plunkett, S. P., B. J. Thompson, O. C. St. Cyr, and R. A. Howard. Solar source regions of coronal mass
507 ejections and their geomagnetic effects. *Journal of Atmospheric and Solar-Terrestrial Physics*, **63**, 389–
508 402, 2001. 10.1016/S1364-6826(00)00166-8. [1](#)
- 509 Pulkkinen, T. Space Weather: Terrestrial Perspective. *Living Reviews in Solar Physics*, **4**, 1, 2007. [1](#)
- 510 Robbrecht, E., and D. Berghmans. Automated recognition of coronal mass ejections (CMEs) in near-real-
511 time data. *A&A*, **425**, 1097–1106, 2004. 10.1051/0004-6361:20041302. [1](#), [2](#), [2](#)
- 512 Savitzky, A., and M. Golay. Smoothing and differentiation of data by simplified least squares procedures.
513 *Analytical Chemistry*, **36**, 1627–1639, 1964. [2](#), [1](#)
- 514 Schwenn, R., A. dal Lago, E. Huttunen, and W. D. Gonzalez. The association of coronal mass ejections with
515 their effects near the Earth. *Annales Geophysicae*, **23**, 1033–1059, 2005. [1](#)
- 516 Sheeley, N. R., Jr., D. J. Michels, R. A. Howard, and M. J. Koomen. Initial observations with the SOLWIND
517 coronagraph. *ApJ*, **237**, L99–L101, 1980. 10.1086/183243. [1](#)
- 518 St. Cyr, O. C., S. P. Plunkett, D. J. Michels, S. E. Paswaters, M. J. Koomen, et al. Properties of coronal
519 mass ejections: SOHO LASCO observations from January 1996 to June 1998. *J. Geophys. Res.*, **105**,
520 18,169–18,186, 2000. 10.1029/1999JA000381. [1](#)
- 521 Stenborg, G., and P. J. Cobelli. A wavelet packets equalization technique to reveal the multiple spatial-scale
522 nature of coronal structures. *A&A*, **398**, 1185–1193, 2003. 10.1051/0004-6361:20021687. [1](#)
- 523 Webb, D. F., and T. A. Howard. Coronal Mass Ejections: Observations. *Living Reviews in Solar Physics*, **9**,
524 3, 2012. [1](#)
- 525 Yashiro, S., N. Gopalswamy, G. Michalek, O. C. St. Cyr, S. P. Plunkett, N. B. Rich, and R. A. Howard. A
526 catalog of white light coronal mass ejections observed by the SOHO spacecraft. *Journal of Geophysical
527 Research (Space Physics)*, **109**, 7105, 2004. 10.1029/2003JA010282. [1](#)