

A clustering-based method for identifying and tracking squall lines

Zhao Shi^{1,2}, Yuxiang Wen^{1,2}, Jianxin He^{1,2}

5 ¹College of Atmospheric Sounding, Chengdu University of Information Technology, Chengdu 610225, China

²Key Laboratory of Atmospheric Sounding, China Meteorological Administration, Chengdu 610225, China

Correspondence to: Yuxiang Wen (3220305009@stu.cuit.edu.cn)

Abstract. The squall line is a type of convective system that is characterized by storm cells arranged in a line or band pattern and is usually associated with disastrous weather. The identification and tracking of squall lines thus play important roles in early warning systems for meteorological disasters. Here, a clustering-based identification and tracking algorithm for squall lines is presented based on weather radar data. Clustering analysis is designed to distinguish the strong echo area and estimate the feature value, including the reflectivity value, length, width, area, endpoints, central axes, and centroid. The linearly arranged clusters are merged to improve the identification of squall line development. The three-dimensional structure and movement tracking of the squall line are obtained using the centroid and velocity of the squall lines identified in each layer. The results demonstrate that the method can effectively identify and track one or more squall lines across the radar surveillance area. The results also show that the recognition accuracy rate for the single scan elevation of this method is 95.06%, and the false-positive rate is 3.17%. This method improves the accuracy of squall line identification in the development stage of squall lines and still works efficiently even when high interference contamination occurs.

1 Introduction

20 The squall line is a prevalent convective weather commonly occurring in mid-latitude regions during spring and summer. Squall lines are arranged by storm cells in a linear structure, which can extend over one hundred or even hundreds of kilometres. The typical life cycle of a squall line is approximately 6-12 hours (Rotunno et al., 1988; Wanghong et al., 2009). The squall lines occurring in coastal areas can bring large amounts of precipitation inland from coastal areas (Oliveira and Oyama, 2020). Squall lines are also associated with severe disastrous weather events, including rainstorms, lightning, hail, downbursts, and even tornadoes (Trapp et al., 2005; Xiaohong et al., 2021). Therefore, identifying and tracking squall lines are crucial for early warning of meteorological disasters.

Weather radar is an effective meteorological remote sensing instrument with high spatiotemporal resolution and has been widely used in monitoring and nowcasting mesoscale convection. To enhance the understanding of radar meteorology on squall lines, since the 1980s, the National Oceanic and Atmospheric Administration (NOAA) has conducted many studies on squall lines using Doppler weather radar data (Smull and Houze, 1985; Srivastava et al., 1986; Smull and Houze, 1987;

Bluestein et al., 1987). The evolution mechanism of the squall line is analysed with radar observations. By combining Vertical Integrated Liquid Water (VIL), Echo Top (ET), Composite Reflectivity (CR) and Velocity–Azimuth Processing (VAP) calculated wind field data and autoweather station observation data, the relationships between squall lines and rainstorms, strong winds, hail, and other disastrous weather processes are revealed. The occurrence of squall lines may lead to clockwise vertical wind shear at low altitudes and counterclockwise vertical wind shear at high altitudes. This shear favoured the generation and strengthening of unstable weather and provided a favourable environment for the development of convection. Simultaneously, the dramatically changing VIL and high ET often heralded hail and strong winds (Wanghong et al., 2009). One of the characteristics of squall lines in weather radar data is the formation of a strong echo band on the radar reflectivity image (Ma, 2022), making these lines visually identifiable. However, their suddenness and wide-ranging impact make it difficult to improve real-time forecasting capabilities via manual identification. Moreover, the automatic identification of squall lines is a complex task (Chengling et al., 2017). Over the past few decades, numerous researchers have conducted studies on the automatic identification of squall lines and have employed various algorithms, such as the two-dimensional Fourier transform (Kelly, 2003), wavelet transform, Hu moment theory (Chengling et al., 2017), Hough transform (Wang et al., 2021; Chengling et al., 2017), and machine learning (Ziqi et al., 2021). Rinehart and Garvey developed pattern recognition schemes of correlation coefficient techniques, Fourier analysis, and Gaussian curve fitting to detect the movement, merging and splitting of storms (Rinehart and Garvey, 1978). Johnson et al. proposed the Storm Cell Identification and Tracking Algorithm (SCIT) (Johnson et al., 1998), which develops centroid techniques to identify and track individual storms (including isolated storms, clustered storms, supercells, and squall lines). The SCIT algorithm has been proven to have an accuracy of more than 90% in identifying storms and has been applied to the WSR88d radar operational system. Dixon and Wiener developed the 'TITAN' algorithm (Dixon and Wiener, 1993), which defines a "storm" as a continuous area that exceeds the reflectivity and size thresholds (adjacent areas with a reflectivity exceeding 35 dbz and a volume exceeding 50 km³). Storms are tracked using the results of the comparison of the previous scan data with subsequent scans with the storm's movement characteristics and maximum horizontal movement speed. Morphological methods were used to identify the merging and splitting of storms. Kelly et al. proposed a method using a two-dimensional Fourier transform and morphological algorithm to identify the characteristics of disaster weather radar images (Kelly, 2003). The above methods identify squall lines as storm cells. However, due to the length and large area of the squall line, as well as the special arrangement (especially when the squall lines are arranged in a shape such as 'L'), identifying the squall line using methods similar to storm cells will result in inadequate identification accuracy (Gangqiang et al., 2021). Promoting the technological development of automatic identification, tracking and prediction of severe weather is a long-pursued research topic. In recent years, weather radar networks have been widely deployed in densely populated areas around the world for severe weather monitoring and early warning, and the use of radar meteorological data has increased explosively. With the development of digital image processing, big data mining, artificial intelligence, and other technologies, the squall

line recognition algorithm has improved. The maximal margin detection method based on wavelet transform patterns and the Hu moment principle can extract the echo characteristics of squall lines(Chengling et al., 2017); however, using a limited number of thresholds may result in false-positives and false-negatives. A convolutional neural network (CNN) was also used for squall line identification (Ziqi et al., 2021). The proposed model effectively identifies the presence of squall lines during the early development stage and the mature convective stage, even when the reflectivity is lower than that in the exuberant stage. However, the size of the dataset of atypical squall lines used for training the model is limited, which may lead to false-positives and false-negatives.

When machine learning techniques are used, physical and morphological features and characteristics cannot be neglected in identification methods.

The clustering algorithm is a type of unsupervised learning algorithm that is commonly used in machine learning and data mining (Gower, 1967). When using supervised machine learning to identify squall lines, it is necessary to label the existing data in advance. However, the appearance of squall lines is random in time and space, making data labelling a complicated project. The clustering algorithm can classify the data points by using some characteristics of the data points without presetting the labels, which is very effective for the squall line process with strong randomness. The clustering algorithm can thus be used to classify the points in the radar scan results, and the data points can be classified based on certain characteristics (such as distance or density). Each individual scan result of the radar sample can correspond to the Euclidean space. Clustering classifies these points into multiple irrelevant sets according to certain information and uses other features of the set to identify the weather system. When used in conjunction with squall line features, the clustering algorithm can identify clusters that meet certain criteria in the radar reflectivity factor data and extract data points that are associated with squall lines.

2 Materials and Methods

The squall line carries a large number of precipitation particles, so its reflectivity factor is significantly greater than that of the surrounding area. The spatial characteristics of the squall line are that the convective system is linearly distributed and covers a large area, which makes it appear in the radar image as a high reflectivity echo band with a large area and long length. This algorithm uses the spatial and temporal evolution characteristics of squall line echoes to extract points with reflectivities that are significantly greater than those in the surrounding area, and through the distribution of points, the noise is filtered out, and the points are divided into clusters. The areas that met the spatial structural characteristics of squall lines were filtered out, and the results of squall line identification were ultimately obtained.

This method mainly comprises the following steps: data preprocessing, threshold calculation, clustering analysis, and target identification and tracking.

The data preprocessing step aims to correspond the radar data to a real geographic location so that the radar data reflect the real weather conditions, and the following steps are based on the preprocessed data. Four thresholds are set to ensure that the

95 algorithm can accurately identify the squall lines during operation. The clustering analysis step separates the points with high reflectivity extracted by threshold 1 into multiple regions and removes noise and isolated points, and the regions' features are calculated after being separated. The regions' features (including area, length, maximum value of the reflectivity, centroid, and moving speed) are compared with the threshold of 2-4 to identify the squall lines. The squall lines are tracked by the centroid of the squall lines, based on their moving speed.

100 **2.1 Data sources**

China and the United States are two of many countries that are significantly impacted by natural meteorological disasters. It is crucial for both nations to undertake meteorological observation and weather prediction and to facilitate scientific assistance in disaster prevention, mitigation, and response to climate change. Both China and the US boast expansive land areas and comprehensive weather radar networks; therefore, both countries have rich weather radar data, including squall line
105 observation results. The method is based on volume scan data from NEXRAD and CINRAD, and these radars have already been used for operation. These radars were well calibrated, and the data obtained by these radars were preliminarily quality controlled, including ground clutter suppression, velocity dealiasing, attenuation correction, and so on. These data have been in use for years and have proven to be reliable most of the time.

2.2 Data Preprocessing

110 Radar volume scan data are stored in polar coordinates (azimuth and radial distance). To facilitate spatial correspondence and further processing, the data need to be converted to Cartesian coordinates and interpolated into a regular spatial grid. The distribution of radar echo information in real space is calculated using the elevation ((El)), azimuth ((Az)) and radial distance ((r)).

$$x = \sin(Az) \sin(El)r ; y = \cos(Az) \sin(El)r \quad (1)$$

Finally, nearest-neighbour interpolation is applied to interpolate the spatial points into the corresponding grid, enabling the
115 gridded data to reflect the actual weather conditions from weather radar scanning.

2.3 Threshold calculation

In previous studies, one approach to extract these points used a threshold value, for which various threshold settings have been proposed, as listed in Table 1.

Num	Threshold settings
1	The length of echo bands with reflectivity greater than 12 dBZ should be no less than 150 km, and the length-to-width ratio of bands with reflectivity greater than 36 dBZ should be no less than 3:1 (Chen and Chou, 1993).

2	The length of echo bands with reflectivity greater than 20 dBZ should be no less than 100 km, and the length-to-width ratio of bands with reflectivity greater than 40 dBZ should be no less than 5:1 (Geerts, 1998)
3	The length of echo bands with reflectivity greater than 40 dBZ should be no less than 100 km (Parker and Johnson, 2000)

Table 1. Threshold settings of the previous studies

120 However, different threshold settings will lead to different identification results. Without proper threshold settings, the method may output incorrect identification results. Therefore, a reasonable selection of thresholds is necessary. Previous studies have focused on the following thresholds: the reflectivity value threshold used to extract storms from weather radar data and the length and width threshold used to identify storms as squall lines.

125 Combined with the study on South China, which used raindrop spectra combined with polarimetric radar (Wang et al., 2019), the following thresholds were set for this method, as listed in Table 2.

Num	Threshold settings
Threshold 1	The minimum values of the reflectivity in the region are not less than 40 dBZ.
Threshold 2	The maximum values of the reflectivity in the region are not less than 52 dBZ.
Threshold 3	The length of the region is not less than 100 kilometres.
Threshold 4	Length to width ratio of the region not less than 3:1

Table 2. Threshold settings of this research

2.4 Clustering analysis

130 The data points that satisfy threshold 1 are extracted through the minimum reflectivity threshold, but it should be noted that the regions extracted using this threshold contain points that do not belong to squall lines, such as storm cells, noise, and clutter points. Therefore, further analysis of these points is required to differentiate points of the convective system from other points and identify the squall lines. The main steps of the cluster analysis process include region clustering, region characterization, and region combination.

2.4.1 Region clustering

135 The area of the convective system that is considered in this method is a high-density cluster point with high reflectivity, where the density is defined as the number of points in a certain area. In this step, a clustering method based on point coordinates, density, and searching distance features is proposed to classify the points extracted by threshold 1. Three parameters are needed in region clustering (shown in Figure 1, the distances mentioned in this method are referred to as the Euclidean distance): the

radius of the field (Eps), the minimum number of points required to judge the core points (MinPts), and the search condition to form clusters (the searching distance). In this method, points in radar data are categorized into core points and noise points according to whether the core conditions are met. The core condition means that there are at least MinPts points within the Eps distance of the point itself. Points that satisfy the core condition are core points; otherwise, they are noise points. If there is a series of core points and the distance between each core point is within the searching distance, these core points are assigned to the same cluster. The method iterates over all the points in the extracted region above. A set of clusters and a composition of unclassified noise points are finally obtained. This method improves the classification ability of the method, especially when the squall line identification process and the line-arranged convective cells are not fully merged or when there is occlusion or interference in the radar data itself.

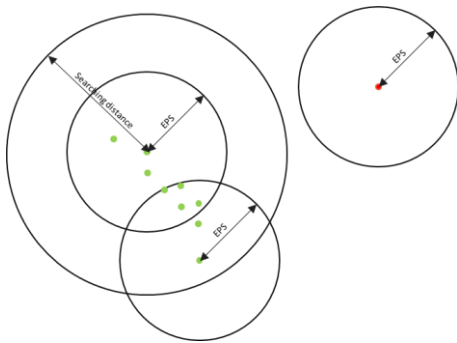


Figure 1 Region clustering parameters

The detailed steps to realize region clustering are as follows:

Parameter Selection: Determine three parameters: Eps, MinPts, and the searching distance. According to Orlanski's classification of convective scales (Orlanski, 1975), squall lines are β -size convective systems; however, in the development stage, squall lines consist of several smaller scale convective systems (γ -size convective systems with a range of 2-20 km) in a linear and tightly packed formation. When searching the points using the Eps range, the presence of γ -scale convection cannot be neglected, so Eps should be in the range of 1-10 km in this algorithm (Eps is insensitive to the threshold in this range).

Missing radar data usually occur in 1-2 radial data points, and the searching distance and EPS take the following values:

$$\text{Searching Distance} = \text{EPS} + 2 \times \text{RadialGate} \quad (2)$$

where RadialGate is the length of the weather radar radial gate. The MinPts is set to the maximum area of the isolated point or noise to be detected. In a previous study (Wang et al., 2021), the number of valid points present in a rectangular box of $n \times n$ (the number of valid points threshold is $0.25 \times$ number of the point) was used to determine the presence of isolated points, so the MinPts was approximately $0.25 \times$, which is the area of the circle with the searching distance as the radius

(round towards negative infinity and retaining one valid digit).

Initialization of labels: Initialize the clustering labels by assigning labels to all points as 'not labelled'.

Core point identification: For a point labelled as 'not labelled', determine whether it meets the core condition; if it meets the core point condition, mark it as a core point, and classify it into the cluster of the current label; otherwise, mark it as a noise point.

165 **Cluster formation:** For other points within the searching distance of the current core point, determine whether they meet the core condition, and if so, mark them as the core point and classify them into the current clustering; otherwise, they are marked as noise points. This step is repeated until all points within the searching distance are labelled.

Clustering iteration: When there is no core point in the searching distance assigned as 'not labelled', update the label to traverse the 'not labelled' points without the searching distance, and repeat the above steps until all points are labelled.

170 Output results: output point coordinates and corresponding labels.

2.4.2 Region characterization

Region clustering divides the extracted points into clusters. To determine whether squall lines exist in these clusters, further analysis of the features of the clusters is needed. The features of the clusters include the central axis, endpoints, area, intensity (maximum reflectivity value), velocity, and position of the centroid.

175 The velocity, intensity, and area of the clusters can be easily obtained via spatial transformation and correspondence. However, determining the central axis and endpoints of clusters is difficult because convective systems are unstable, which results in clusters with irregular shapes. Therefore, it is necessary to determine an efficient and accurate method to estimate the central axis of the clusters. The Hough transform is an image processing algorithm published by Hough et al. in 1959 (Hough, 1959) that has been widely used to recognize lines or circles in complex images (Duda and Hart, 1972), and it has also been used for
180 squall line identification in previous research (Wang et al., 2021). The Hough transform is a method that utilizes a voting-based approach to transform a collection of lines into a collection of points. This method transforms the point space (X, Y) into parameter space (ρ, θ) to form a series of voting accumulations. The resulting parameter space consists of two parameters: ρ and θ . The point coordinates X and Y are converted to (ρ, θ) by the following equation:

$$\rho = X * \cos(\theta) + Y * \sin(\theta) \tag{3}$$

The Hough peak is the partial maximum in the point set voting results. These partial peaks correspond to the most voted (ρ, θ)
185 and represent potential lines in the original data. The ρ, θ values of the Hough peak are associated with a straight line in the gridded radar data. The ρ, θ obtained at this point corresponds to straight lines in the gridded radar data (as shown in Figure), and these straight lines are considered central axes for clustering.

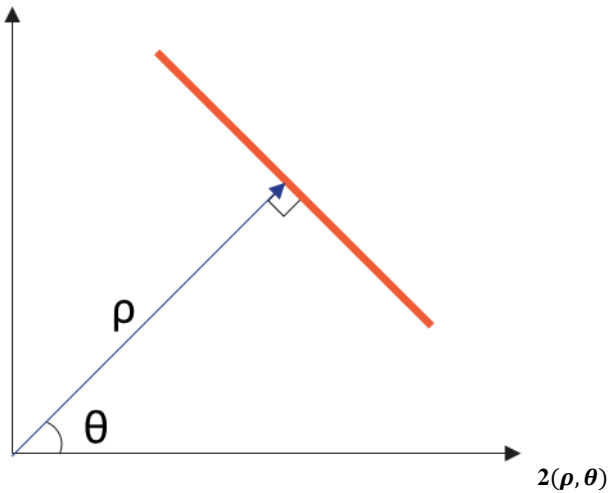


Figure Correspondence between straight lines and parameters

190 **Moreover**, the straight line intersects with the edge of the clusters, and these intersection points are considered to be the endpoints of the clusters. Thus, the length of the central axis can be estimated using the endpoints. The clusters associated with storms are approximately considered ellipses, and the widths of the clusters are calculated using the area and length of the central

axes obtained from the above steps:

$$width = \frac{4 * aera}{\pi * length} \quad (4)$$

195 where, *aera* is the area of the cluster, and *length* is the estimated length of the cluster.

2.4.3 Region combination

To enable the algorithm to identify the squall lines when they are fully formed (convective systems that are not fully merged but are linearly arranged in the development stage of the squall line), in this step, the linearly arranged clusters are merged. The traditional methods use the centroid distance to determine whether to merge storms. However, for the squall line, due to
 200 the special characteristics of its linear arrangement and large length-to-width ratio, the determination method of simply determining the distance between centroids in accordance with the distance circle has several drawbacks. When the distance is large, the nonlinearly arranged clusters in the 'width' direction (which does not line up with other clusters but within the distance) may be merged, and when the value is small, the clusters in the 'Length' direction cannot be searched.

This algorithm determines whether to merge two clusters by the distances between the obtained endpoints above. If the two
 205 clusters' nearest endpoints are within 10 km, then the two clusters are combined into the same cluster. After the two clusters are merged, the features change, so the recalculation of the features is needed. The length and area of the clusters are added

together as the area and length of the newly merged clusters, and the maximum reflectivity factor and width are taken as the larger values of the two clusters. The coordinates of the centroid in the horizontal direction are calculated as follows:

$$X = \frac{\sum j * Zh(i, j)}{\sum Zh(i, j)}, Y = \frac{\sum i * Zh(i, j)}{\sum Zh(i, j)} \quad (5)$$

2.5 Target identification

210 The results obtained from the cluster analysis step are compared with the threshold Conditions 2-4 in Table 2. If there is a cluster that satisfies the threshold, at least one squall line exists in this layer of the current volume scanning data, and this cluster is considered to be an identified squall line.

2.6 Target Tracking

215 The above steps enable the squall lines to be identified and the locations of squall lines to be obtained in a single layer of volume-scan data. However, in practice, the vertical structure plays a more important role than the horizontal structure during strong convective weather, which is prone to cause major meteorological disasters (Ma, 2022). Therefore, it is necessary to obtain the three-dimensional structure of the radar-scanned information of squall lines. Moreover, convective storms are characterized by rapid structural evolution and movement. Over the life cycle of the squall line, it may undergo multiple splits, regenerations, and reorganizations (Ye-Qing et al., 2008). Therefore, it is also necessary to track the changes in the shape and
220 location of squall lines in practical applications.

The wind field and velocity information are calculated by the VAP method from radar radial velocity data. The traditional VAP method assumes that the wind field is uniform in the region, and the calculation ability is good in the case of a uniform wind field, but the error is larger in the case of a nonuniform wind field, so the extended VAP (EVAP) (Zhouzhenbo et al., 2006) is used to calculate the wind field. The EVAP inversion method is as follows:

$$\cos(\Delta\theta + \Delta\beta) = \frac{V_{r1} + V_{r2}}{V_r} \quad (6)$$

$$\tan \beta = \left(\frac{V_{r1} - V_{r2}}{V_{r1} + V_{r2}} \right) \cot(\Delta\theta + \Delta\beta) \quad (7)$$

$$V = \left| \frac{V_r}{\cos \beta} \right| \quad (8)$$

225 where V_{r1}, V_{r2} are the radial velocities in the azimuthal angle adjacent to V_r on the equidistant circle. The position of the centroid of the squall line is combined with the velocity data obtained beforehand, and the maximum wind field velocity is considered to be the maximum moving speed of the squall line. The squall lines in different layers of the scanning data whose centroid converges within a distance (R) are considered to be the same squall line. (The process of calculating R takes into

account the fact that the shape of the squall line might change in the process of moving and evolving, which results in centroid shifting, so the calculation method introduces the width of the squall line to improve the searching ability):

$$R = V_{max} * \Delta time + width/2 \quad (9)$$

where V_{max} is the maximum wind field velocity obtained by inversion in the squall line region, $\Delta time$ is the scanning data time interval, and $width$ is the squall line width estimated above.

By applying the above method to data with different elevation angles from the same volume-scanning process, the three-dimensional structure of the squall line can be obtained. Applying this method to data from different volume-scanning processes enables the squall line to be tracked.

The resulting flowchart of the algorithm is as follows:

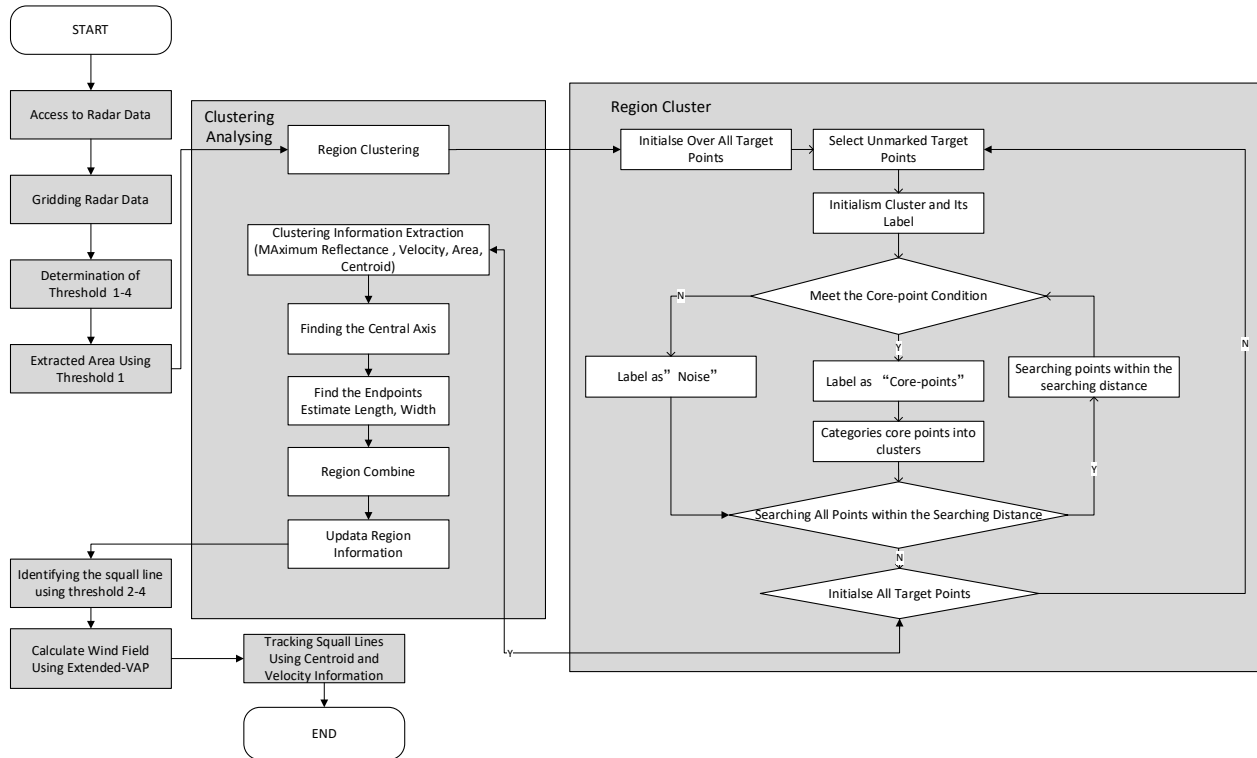


Figure 3 Overall flowchart of the proposed algorithm

3 Results

240 3.1 Experimental design

The proposed method is based on NEXRAD and CINRAD; thus, to test the efficacy of this method, radar data from both weather radar networks are employed for the experiment. A typical squall line was observed by the CINRAD Z9762 weather radar in He Yuan, Guangdong Province, on June 4, 2016. Partial thunderstorms and gusty weather occurred that day. The radar

recorded the developmental stage and exuberant stage of the squall line. Through visual observation by meteorologists, a strong
245 echo band with a length of approximately 200 km was found in the radar echo image. The volume-scan data of the radar on
that day are selected to demonstrate the algorithm identification process and to show the results of each step in detail. Moreover,
three-dimensional structure merging and dynamic tracking are demonstrated. Based on this example, the anti-interference
ability is verified by adding artificial noise interference.

To ensure the performance of the algorithm in atypical situations, the ability to simultaneously identify more than one squall
250 line needs to be considered. A special squall line process was observed in the United States on July 1, 2014. Hurricanes cause
large areas of severe wind damage. Two squall lines were sounded in the same volume-scan data. Observations from the
Chicago-based KLOT radar were chosen to validate the algorithm's ability to identify, 3D merge, and track multiple storm
lines in the presence of multiple squall lines at the same moment in the volumetric scan data.

To better simulate the actual situation and to verify the performance of the method objectively, the radar volume scanning data
255 related to tornadoes are selected to verify the algorithm's performance in the identification process by comparing the algorithm
identification results with the manual identification results and TITAN identification results.

3.2 Example from Z9762 in HeYuan, Guangdong, China

3.2.1. Static identification

Static identification refers to the process of identifying single layers of radar volume scan data. The third layer (data with an
260 elevation angle of 1.36) of the volume scan data of the Z9762 radar in Heyuan city, Guangdong Province, China, in 2016 at
7:00:00 UTC is selected as a typical example to analyse the identifying capability of the algorithm in a single layer of volume
scan data.

First, the radar data are gridded, and the nearest neighbour interpolation method is used to obtain the information of the echo
data in real space (shown in Figure 4 (a)).

265 Using a threshold of 1, the points with a reflectivity factor that is significantly greater than that of the surrounding area are
extracted, as shown in Figure 4 (b). A squall line can be visually observed in the image, but it should be noted that there are
also some points that do not correspond to the squall lines. Further analysis is needed to extract the squall line from these
points.

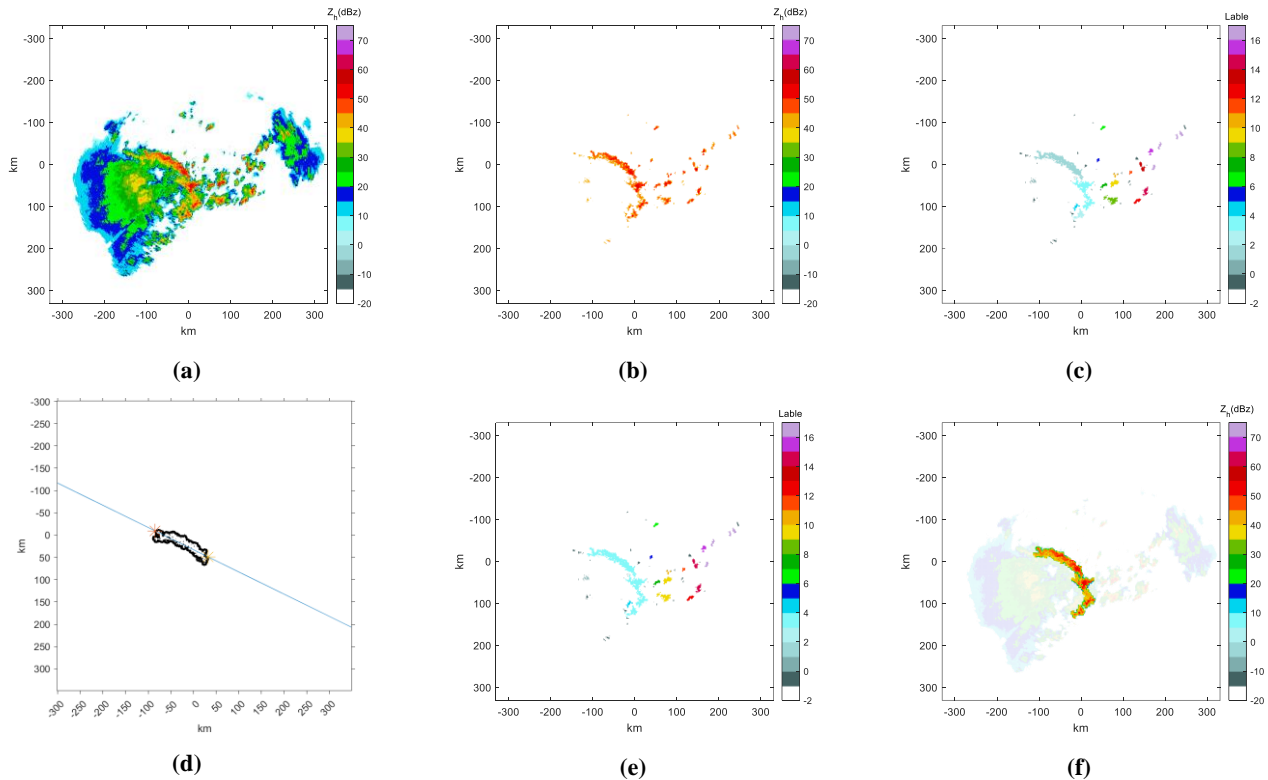
The above data points can be classified based on the point density and searching distance via the clustering method. In the
270 clustering method, the density clustering parameters are set as follows: Eps = 5, MinPts = 30, and the searching distance = 7.
Moreover, the parameters are not sensitive to the threshold and can be adjusted to suit different conditions of use. The
classification results obtained in the region clustering step of the clustering analysis are shown in Figure 4 (c).

Cluster characterization is based on the clusters in the above results. The cluster labelled 1 is used to show the details of the
region characterization step. The cluster area is calculated to be 1,239 km². The central axis of the cluster obtained through the

275 Hough transform and the two endpoints are shown in Figure 4 (d). The cluster is estimated to have a length of 120.7 km and a width of 15.4 km. The central intensity of the cluster is 61 dBZ.

All the clusters are iterated, and the satisfactory clusters are merged according to their endpoint distances. The results shown in Figure 4 (e) are obtained by recalculating the clusters' features and verifying whether they meet thresholds 2-4 (mentioned in Table 2). The clusters that satisfy the thresholds are considered to be squall lines (Figure 4(f)).

280

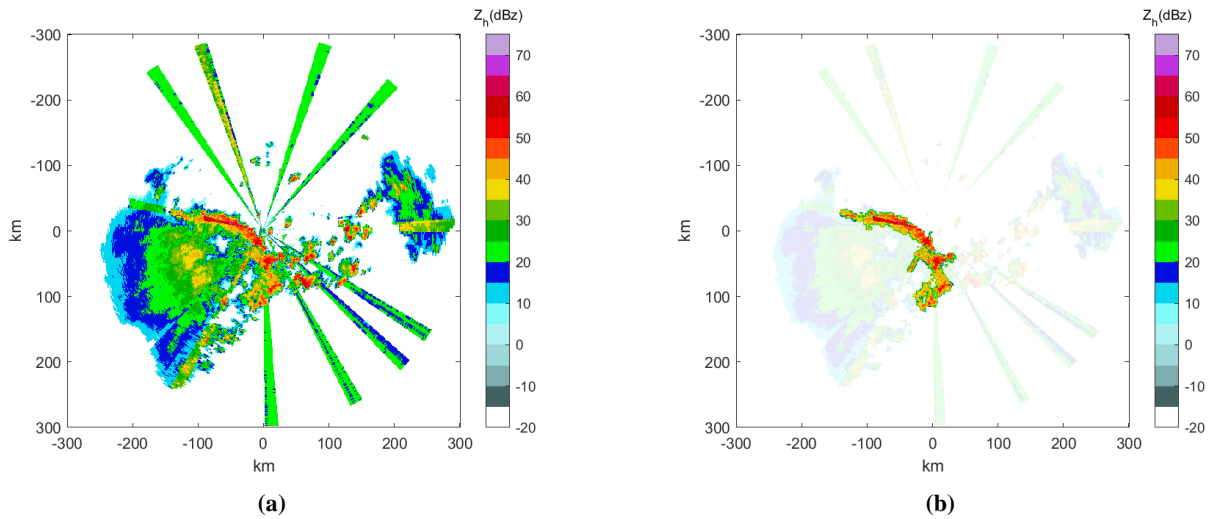


285 **Figure 4 (a) The weather system detected at Z9762 at 7:00:00 UTC on June 04, 2016, with a squall line that can be seen in the figure. (b) For the points extracted using the threshold, the reflectivities of these points are significantly greater than those of the surrounding points. (c) The results of region clustering; the colour bar indicates different clusters. (d) An example of obtaining the central axis. The edge of the cluster and the identified central axis and the endpoints of the obtained cluster central axis are marked as ‘*’, and the length of the cluster can be calculated using these two points. (e) The results of region combination. (f) The result of squall line identification (areas with a hundredth white opacities are identified as squall lines).**

290 The above results show that the algorithm is effective in identifying typical squall lines in single-layer radar data, can effectively identify the existence of squall lines and mark the location of squall lines, and can differentiate squall lines between convective cells and clutter points.

3.2.2. Anti-interference capability

The proposed algorithm uses quality-controlled weather radar data in principle. However, as mentioned in a prior study (Wang et al., 2021), in the case of interference that is not eliminated during quality control. There will be interference (with a small probability of occurrence) that may not be eliminated in the quality-control process, traditional algorithms are not able to overcome this interference (The prior method would recognize the interference as squall lines). Compared to traditional methods, the proposed method has a greater degree of interference resistance. Cochannel interference data are added to the reflectivity data by random replacement or addition to test the anti-interference performance of the algorithm, and the simulated interference and the identification of the results are shown.



300

Figure 5 The radar data with random cochannel noise added (a) and the result of the squall line identifying the squall line (b)

The results show that this method is more robust than the traditional method in the process of static identification when encountering interference. The method can identify the squall line in the data with cochannel interference if the interference does not cover the weather information, while the method in the prior study is impossible to use for effective identification.

305 3.2.3 Three-dimensional structures

The 3D structure of the squall line obtained by merging the different layers of volume-scan data using centroid and EVAP is shown in Figure 6:

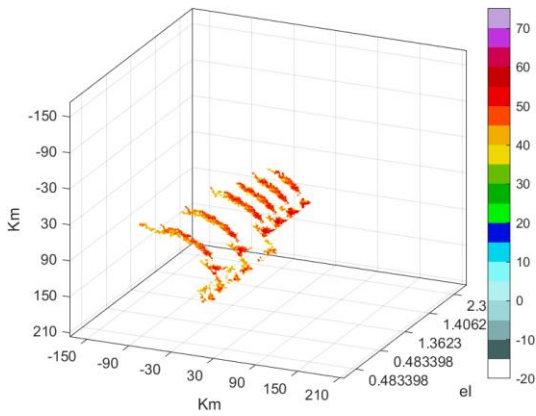
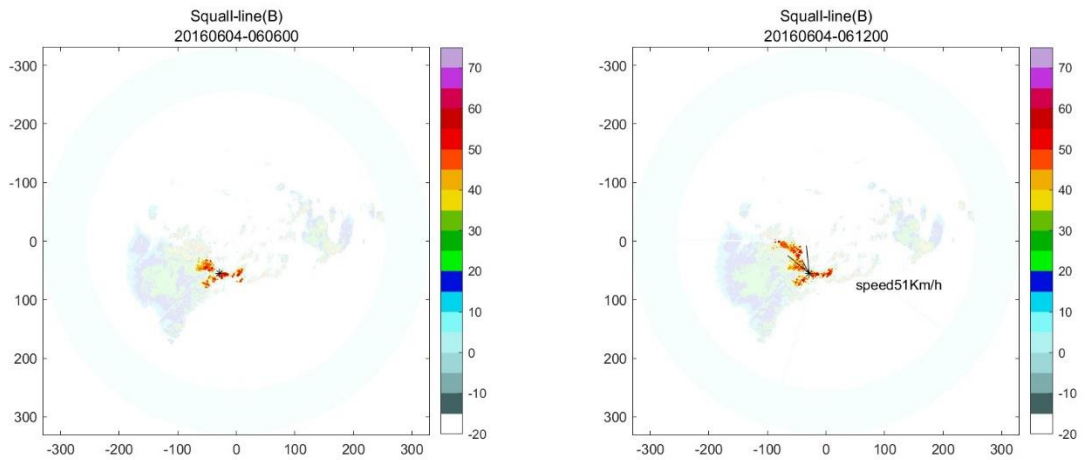


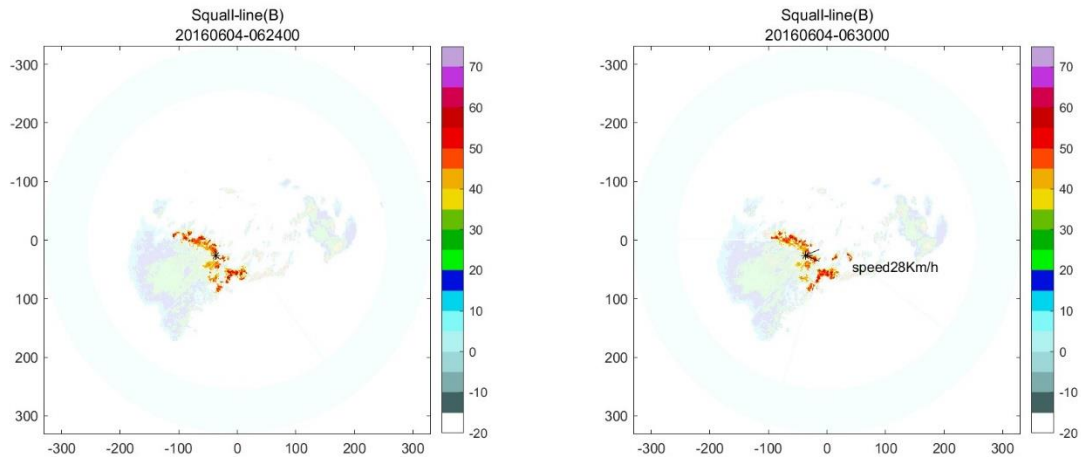
Figure 6 The 3D strut of the squall line of the reflectivity data.

310 **3.2.4 Target racking**

The results of squall line tracking using the above method are as follows (the third layer of volume-scan data for the period 05:54:00-07:00:00 is chosen to demonstrate the tracking effect).



(a)



(b)

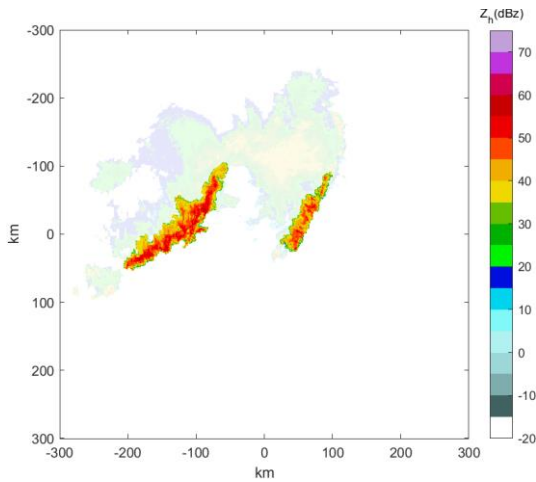
315 **Figure 7 Results of squall line tracking at 6:06:00 UTC on June 04, 2016 (a), and at 6:24:00 UTC on June 04, 2016 (b). The figure shows the squall lines' movement of the latter moment relative to the former moment. '*' shows the location of the centroid of the squall lines, and the arrows show the direction of movement of the squall lines. The moving speed of the squall lines is given on the right side of the arrow.**

The results show that the method can be used to effectively track squall lines in the development and exuberant stages of squall development.

3.3 Example from KLOT in the US

320 3.3.1. Static identification

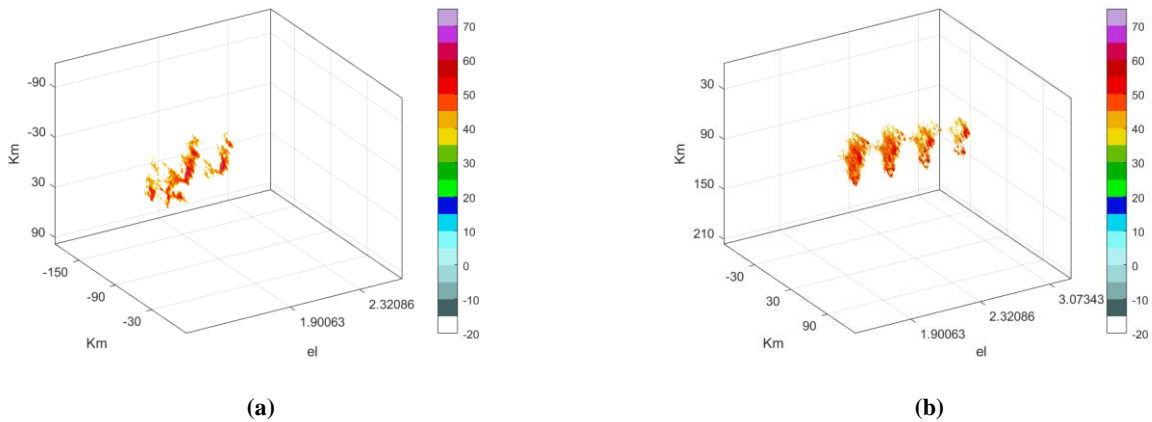
The static identification result of the KLOT's lowest elevation level III data observed at 01:13:00 UTC Jul 01, 2017, is shown in Figure 8 . Two squall lines are identified by the algorithm.



325 **Figure 8** The result of squall line identification on KLOT at 01:13:00. UTC July 01, 2014. (Areas with a hundredth white opacity are identified as squall lines.)

3.3.2. Three-dimensional structures

The volume scan data from KLOT on 01:16:55 UTC Jul 01, 2017, level II data were selected for the validation of the radar data 3D merging capability of the double squall line process, and the 3D structure is shown in Figure 9.

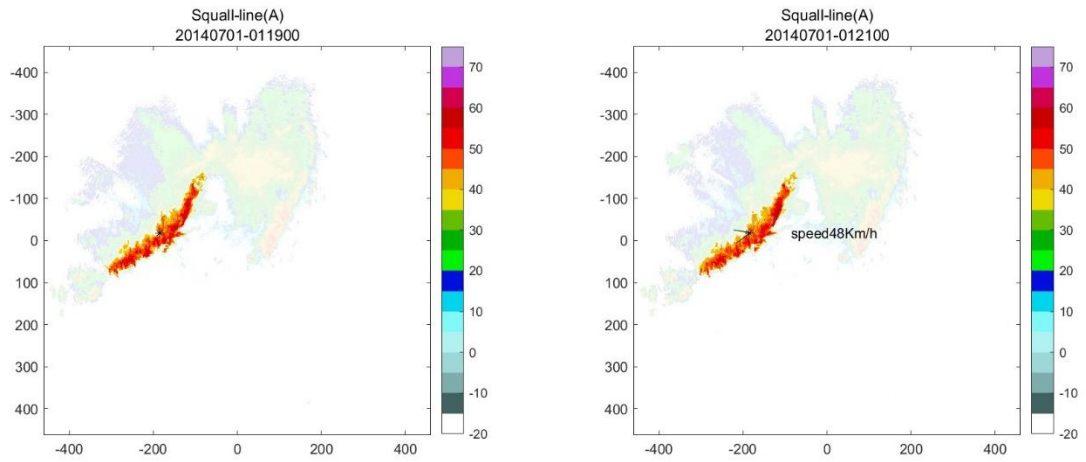


330 **Figure 9** The 3D strut of the reflectivity data of the squall line on the left in the radar image (a) and the squall line on the right (b).

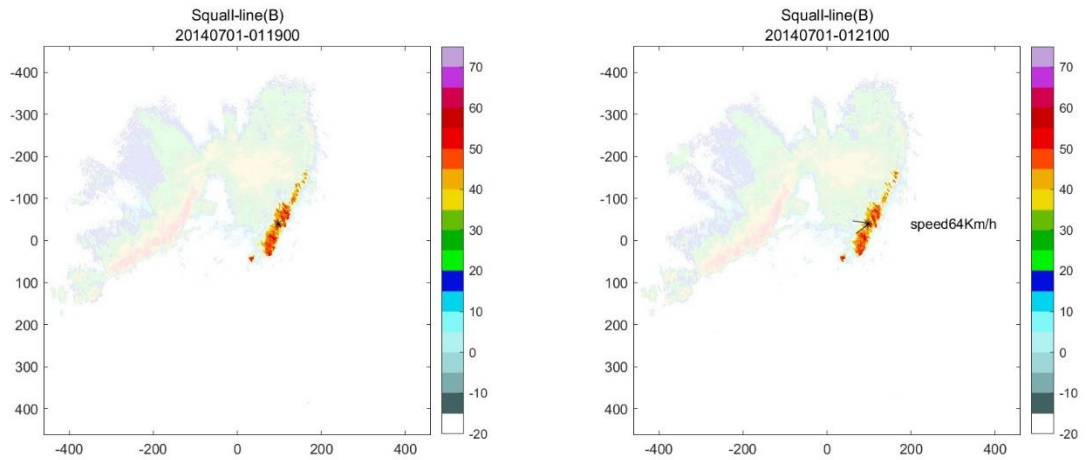
3.3.3 Target racking

The two squall lines appearing in the above examples are tracked separately, and the tracking results are shown in Figure 10

335 The above results show that when two squall lines appear in the volume scanning process of the same radar, the algorithm can track the two squall lines separately. At the same time, the method can provide the direction and speed of movement.



(a)



(b)

340 **Figure 10** The tracking results for the two squall lines at 01:13:00 UTC Jul 01, 2017, are shown in (a) and (b). The figure shows the movement of the squall lines in the latter moment relative to those in the former moment. ‘*’ shows the location of the centroid of the squall lines, and the arrows show the direction of movement of the squall lines. The moving speed of the squall lines is given on the right side of the arrow.

3.4 Quantitative analysis

To better represent physical scenarios and to obtain a more objective verification of the performance of the algorithm, a significant number of verification experiments are needed. Due to the limited observational range of a single radar, the appearance of squall lines is characterized by randomness. The amount of available radar data is limited for finding data containing squall lines. There is a certain connection between the existence of squall lines and the occurrence of tornadoes. Weather radar data associated with tornadoes observed in Jiangsu Province (the selected data are shown in Table 3) are selected as the dataset for quantitative analysis. Additionally, precipitation data from 2022/11 to 2023/05 obtained from the RXM25 radar in Chengdu were also selected. The performance of the algorithm was verified using two approaches: manually identifying data for comparison with the algorithm results and using the TITAN algorithm identifying results.

355 By comparing the results of manual identification or TITAN with those of algorithmic identification, the confusion matrix is obtained as follows. The corresponding results are shown in Table 4 and Table 5. In manual identification, computers are primarily utilized to search the data to meet the following request: the highest reflectivity is greater than 50 dBZ, and the number of points greater than 40 dBZ is not less than 1,000. Then, the reflectivity data are manually selected to meet the following request: there is a region in the radar-girded data with a reflectivity of not less than 40 dBZ, the area of this region is not less than 2000 km², the length is not less than 100 km, and the maximum reflectivity of the region is more than 50 dBZ. The data thus selected are considered to be manually identified squall lines.

360 The following events are defined by taking the manual (or TITAN) identifying result as the true result of the sample (if a squall line is manually (or using TITAN) observed in the radar echo image, the squall line is considered to actually exist within the radar observation range) in the quantitative analysis process: true positive (TP): the true result of the sample is positive, and the algorithm predicts that the result is positive. True Negative (TN): the true result of the sample is negative, and the algorithm predicts that the result is negative. False-positive (FP): the true result of the sample is negative, and the algorithm predicts that the result is positive. False-negative (FN): the true result of the sample is positive, and the algorithm predicts that the result is negative. Based on the above samples, the following parameters are defined to reflect the algorithm performance:

Radar Name	Observation time	Radar Name	Observation time
Z9250	2007/07/03	Z9517	2016/06/23
	2011/07/12		2017/08/01
Z9513	2009/08/27	Z9518	2008/07/04
	2011/07/13	Z9516	2006/07/03
	2016/07/06		2008/07/29
	2006/07/03		2008/07/30
2008/07/29	2008/08/17		
Z9515	2008/07/30	Z9519	2011/07/11
	2008/08/17		2012/08/10
	2011/07/11		2011/08/02
	2012/08/10		2016/07/06
	2013/07/07		Z9527
Z9523	2015/07/24	Z9527	2018/08/18
	2017/07/02		

365 Table 3 Data source information

Algorithm	Manual	Y	N
	Y		1040
N		54	7409

'Y' indicates that there are squall lines identified, 'N' indicates that there are no squall lines identified

Table 4 The manual identification and algorithmic results

Algorithm	TITAN	Y	N
	Y		953
N		46	7417

'Y' indicates that there are squall lines identified, 'N' indicates that there are no squall lines identified

Table 5 The TITAN identification and algorithmic results

370

Accuracy of algorithmic recognition:

$$ACC = \frac{TP + TN}{TP + FP + TN + FN} \quad (10)$$

Successful squall line identification rate of the algorithm:

$$PRE = \frac{TP}{TP + FP} \quad (11)$$

False identification rate of the algorithm:

$$FAR = FN / (TN + FN) \quad (12)$$

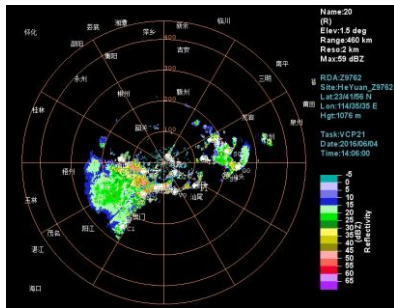
Missed identifying rate of the algorithm:

$$NAR = FP / (TP + FP) \quad (13)$$

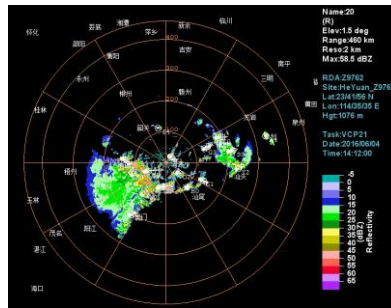
375 In the proposed method, the manually identified result is taken as the true result, and the calculated result is $ACC = 98.25\%$, $PRE = 95.06\%$, $FAR = 3.17\%$, $NAR = 4.93\%$. Meanwhile, taking the TITAN identification result as the true result, the calculated result is $ACC = 97.31\%$, $PRE = 95.40\%$, $FAR = 16.33\%$, $NAR = 4.84\%$.

The above tests show that the accuracy and recognition rate of this method are greater than 95%. Using the manually identified results as a benchmark, the FAR and NAR were greater than 5%. However, the FAR is greater than 15% when using the
380 TITAN result as a benchmark. By comparing the TITAN results with the manual identification results, we observe that TITAN identified linear storm cells as multiple independent cells during the squall line development stage.

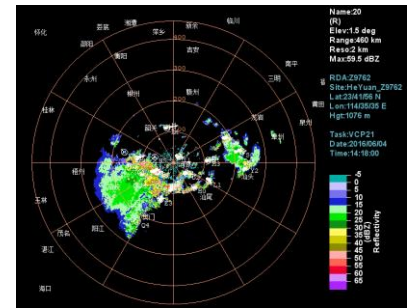
Figure 11 shows the products generated by the TITAN algorithm in the operations of the CINRADSA radar network. The product shows the squall line identifying the result of Z9762 from 6:06:00 to 6:24 UTC on June 04, 2016 (Figure 11(a)-(f)). The product shows the storm cells' location by the centroid and uses an arrow to show the movement of the cells. The squall
385 lines can also be displayed separately as an independent product. The identification results during this time period of the method designed in this experiment are shown in Figure 11 (g) –(i).



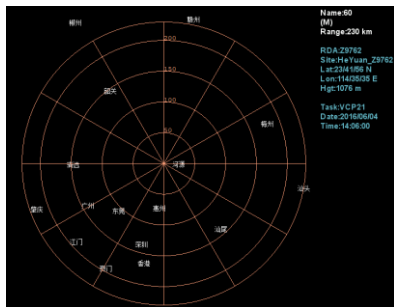
(a)



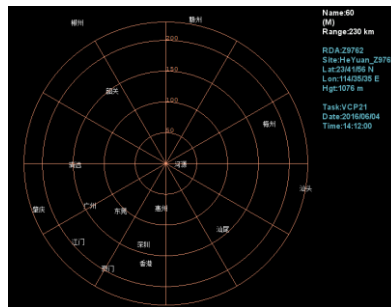
(b)



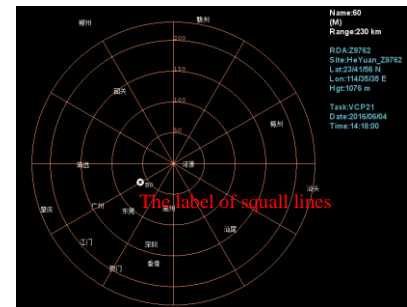
(c)



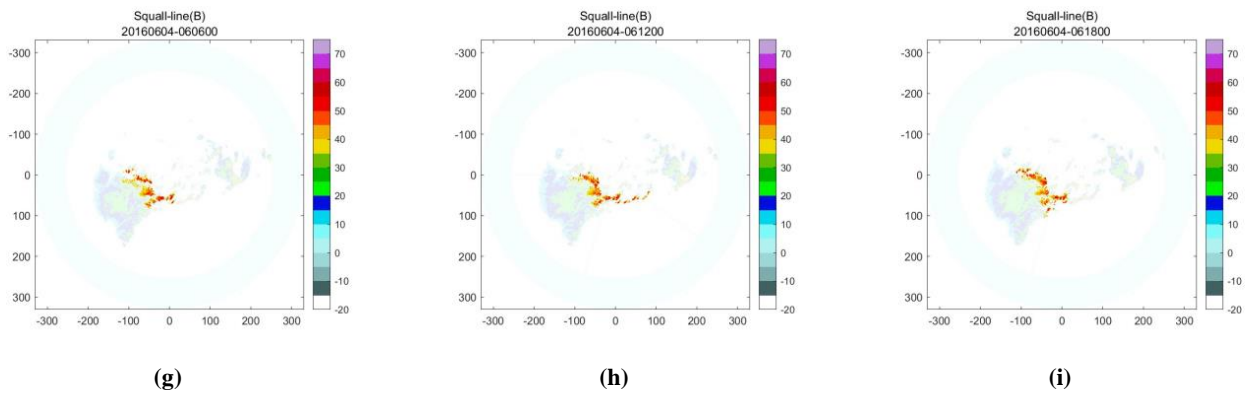
(d)



(e)



(f)



390 **Figure 11 (a)-(f) The products generated by the TITAN algorithm. (g)-(i) The identification results of the method in this study. Figures (a)-(c) show the identification results of TITAN for all storm cells at 6:06:00 UTC on June 04, 2016 UTC (a), 6:12:00 UTC on June 04, 2016 UTC (b), 6:18:00 UTC on June 04, and 2016 UTC (c). The mesoscale convection identification results of TITAN or the corresponding moments of (a), (b), and (c) are specifically labelled in (d)-(f), and the label is near the text in (f). Figures (g)-(i) show the identification results for the corresponding moments of the method designed in this study. The time shown in (a)-(f) refers to Beijing time (UTC +8).**

As shown in Figure 11, the proposed method first identifies the squall line at 6:06:00 UTC, while the TITAN identifies the squall line at 6:18:00 UTC. According to the above figures, when the cells are not fully merged, TITAN identifies them as independent cells. The method of this study involves identifying squall lines before they are fully merged. It is demonstrated that the squall warning can be improved compared to the TITAN algorithm. By statistically analysing the results with manual identification, the algorithm can advance the squall line warning time by approximately 15 minutes. This result indicates that the combination of storm cells in a linear arrangement allows early identification of squall lines.

400 Overall, the proposed method can be used to effectively identify squall lines in selected weather radar data. This method can also be used to effectively identify squall lines in both the development stage and the exuberant stage. It can also provide the three-dimensional structure of the squall line and track the squall line before it is fully merged. The proposed algorithm can improve the timeliness of weather forecasting through the above advantages.

4 Conclusion and discussion

In this paper, an automatic squall line identification and tracking method for weather radar echo data is presented. Doppler weather radar data are used as the data source. The points that are significantly higher than the surrounding area are extracted by threshold from the preprocessed radar data. The points are distinguished into clusters by the clustering method. The cluster features, including the reflectivity value, length, width, area, endpoints, central axis, and centroid, are obtained by clustering characterization. The linearly arranged clusters are merged to improve the identification ability in the squall line development stage. The movement tracking and three-dimensional structure of the squall line are obtained using the centroid and velocity of the squall lines identified in a single layer.

Analysing the weather processes from two different radars proved that the proposed method can identify one or more squall lines in the radar data effectively. In the process of quantitative analysis of the algorithms, the manual identification or TITAN identification results were identified as true samples. Both analyses show that the method has high accuracy. The analysis of the TITAN results shows that the method in this study can advance the early warning of squall lines. However, the manual
415 identification process is still somewhat subjective; therefore, further optimization experiments using other polarization parameters, weather station data, etc., are needed.

Compared with traditional methods, this method does not rely on manual observation, so the identification and tracking process can be automated through computers to improve the accuracy and timeliness of weather warning operations. This method can identify the squall lines earlier than the traditional method. The squall lines are often associated with tornadoes, downburst
420 storms, and other catastrophic weather events, and the earlier identification of squall lines in the early warning system can be used to send out warning information of the related meteorological events. Meanwhile, in the usage of the collaborative radar system, combined with the identification of squall lines, refined structural detection results will be carried out and more accurate analysis results will be derived, leading to more precise warning results.

The identification of the squall line in this method is mainly based on the radar reflectivity data, which may not be very accurate
425 for the edge marking of the squall lines. Therefore, other radar parameters will be used in combination with machine learning algorithms to obtain accurate edges of the squall lines in subsequent studies. Another limitation is that the identification and tracking process of this method only works in the scanning results of a single radar, which requires the radar to scan the complete squall lines, and in the subsequent research, the identification of squall lines in a larger coverage will be realized by the girded data of multiple radars.

430 Using this algorithm together with other convective identification and tracking algorithms, the information of squall lines, storm cells, supercells, and other targets can be used simultaneously, and the ability to predict catastrophic weather, including tornadoes, can be greatly improved. In combination with traditional weather warning algorithms, this approach can further improve the reliability of catastrophic weather warning work. A finer vertical structure of squall may be obtained with deep learning technology, and a 3D structure can be obtained via this method.

435

Author Contributions: Conceptualization, Y.W., Z.S., and J.H.; methodology, Z.S. and Y.W.; software, J.H. and Z.S.; formal analysis, Z.S., Y.W. and J.H.; writing—original draft preparation, Y.W.; writing—review and editing, J.H., Y.W.; visualization, Y.W. All authors have read and agreed to the published version of the manuscript.

440 *Funding:* This work was supported by the National Key R&D Program of China (2021YFC3090203), the Key Laboratory of Atmospheric Sounding Program of China Meteorological Administration (U2021Z01, U2021Z09), and the CMA Meteorological Observation Centre (CMAJBGS202203).

Data availability statement: Not applicable.

445 *Acknowledgements:* The authors would like to express their sincere thanks to the Guangdong Meteorological Network and Equipment Support Centre for supplying the data used in this manuscript and their viewers for their constructive comments and editorial suggestions, which considerably helped improve the quality of the manuscript.

Conflicts of Interest: The authors declare no conflicts of interest.

References

450

Bluestein, H. B., Marx, G. T., and Jain, M. H.: Formation of Mesoscale Lines of Precipitation: Nonsevere Squall Lines in Oklahoma during the Spring, *Monthly Weather Review*, 115, 2719-2727, [https://doi.org/10.1175/1520-0493\(1987\)115<2719:FOMLOP>2.0.CO;2](https://doi.org/10.1175/1520-0493(1987)115<2719:FOMLOP>2.0.CO;2), 1987.

455 Chen, G. T. J. and Chou, H. C.: General Characteristics of Squall Lines Observed in TAMEX, *Monthly Weather Review*, 121, 726-733, [https://doi.org/10.1175/1520-0493\(1993\)121<0726:GCOSLO>2.0.CO;2](https://doi.org/10.1175/1520-0493(1993)121<0726:GCOSLO>2.0.CO;2), 1993.

CHENGLing, z., HEJian, x., and ZENG, X.-j.: Radar Echo Recognition of Squall Line based on Wavelet and Hu Moment, *JOURNAL OF CHENGDU UNIVERSITY OF INFORMATION TECHNOLOGY*, 32, 369-374, 10.16836/j.cnki.jcuit.2017.04.005, 2017.

460 Dixon, M. and Wiener, G.: TITAN: Thunderstorm Identification, Tracking, Analysis, and Nowcasting—A Radar-based Methodology, *Journal of Atmospheric and Oceanic Technology*, 10, 785-797, [https://doi.org/10.1175/1520-0426\(1993\)010<0785:TTITAA>2.0.CO;2](https://doi.org/10.1175/1520-0426(1993)010<0785:TTITAA>2.0.CO;2), 1993.

Duda, R. O. and Hart, P. E.: Use of the Hough transformation to detect lines and curves in pictures, *Commun. ACM*, 15, 11–15, 10.1145/361237.361242, 1972.

465 Gangqiang, N. A. N., Mingxuan, C., Rui, Q. I. N., Lei, H. A. N., and Weihua, C. A. O.: Identification, tracking and classification method of mesoscale convective system based on radar composite reflectivity mosaic and deep learning, *Acta Meteorologica Sinica*, 79, 1002-1021, 10.11676/qxxb2021.062, 2021.

Geerts, B.: Mesoscale Convective Systems in the Southeast United States during 1994–95: A Survey, *Weather and Forecasting*, 13, 860-869, [https://doi.org/10.1175/1520-0434\(1998\)013<0860:MCSITS>2.0.CO;2](https://doi.org/10.1175/1520-0434(1998)013<0860:MCSITS>2.0.CO;2), 1998.

Gower, J. C.: A Comparison of Some Methods of Cluster Analysis, *Biometrics*, 23, 623-637, 10.2307/2528417, 1967.

470 Hough, P. V. C.: Machine Analysis of Bubble Chamber Pictures,

Johnson, J. T., MacKeen, P. L., Witt, A., Mitchell, E. D. W., Stumpf, G. J., Eilts, M. D., and Thomas, K. W.: The Storm Cell Identification and Tracking Algorithm: An Enhanced WSR-88D Algorithm, *Weather and Forecasting*, 13, 263-276, [https://doi.org/10.1175/1520-0434\(1998\)013<0263:TSCIAT>2.0.CO;2](https://doi.org/10.1175/1520-0434(1998)013<0263:TSCIAT>2.0.CO;2), 1998.

Kelly, W. E.: IMAGE PROCESSING FOR HAZARD RECOGNITION IN ON-BOARD WEATHER RADAR, 2003.

475 MA, J. Z. A. M. W. A. S. G. A. H. H. A. L.: The scattering mechanism of squall lines with C-Band dual polarization radar. Part I: echo characteristics and particles phase recognition, *Front. Earth Sci.*, 16, 2022.

Oliveira, F. P. and Oyama, M. D.: Squall-line initiation over the northern coast of Brazil in March: Observational features, *Meteorological Applications*, 27, 10.1002/met.1799, 2020.

480 Orlandi, I.: A rational subdivision of scales for atmospheric processes, *Bulletin of the American Meteorological Society*, 56, 527-530, 1975.

Parker, M. D. and Johnson, R. H.: Organizational Modes of Midlatitude Mesoscale Convective Systems, *Monthly Weather Review*, 128, 3413-3436, [https://doi.org/10.1175/1520-0493\(2001\)129<3413:OMOMMC>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<3413:OMOMMC>2.0.CO;2), 2000.

Rinehart, R. E. and Garvey, E. T.: Three-dimensional storm motion detection by conventional weather radar, *Nature*, 273, 287-289, 10.1038/273287a0, 1978.

485 Rotunno, R., Klemp, J. B., and Weisman, M. L.: A Theory for Strong, Long-Lived Squall Lines, *Journal of the Atmospheric Sciences*, 45, 463-485, 10.1175/1520-0469(1988)045<0463:atfssl>2.0.co;2, 1988.

Smull, B. F. and Houze, R. A.: A Midlatitude Squall Line with a Trailing Region of Stratiform Rain: Radar and Satellite Observations, *Monthly Weather Review*, 113, 117-133, [https://doi.org/10.1175/1520-0493\(1985\)113<0117:AMSLWA>2.0.CO;2](https://doi.org/10.1175/1520-0493(1985)113<0117:AMSLWA>2.0.CO;2), 1985.

- 490 Smull, B. F. and Houze, R. A.: Dual-Doppler Radar Analysis of a Midlatitude Squall Line with a Trailing Region of Stratiform Rain, *Journal of Atmospheric Sciences*, 44, 2128-2149, [https://doi.org/10.1175/1520-0469\(1987\)044<2128:DDRAOA>2.0.CO;2](https://doi.org/10.1175/1520-0469(1987)044<2128:DDRAOA>2.0.CO;2), 1987.
- Srivastava, R. C., Matejka, T. J., and Lorello, T. J.: Doppler Radar Study of the Trailing Anvil Region Associated with a Squall Line, *Journal of Atmospheric Sciences*, 43, 356-377, [https://doi.org/10.1175/1520-0469\(1986\)043<0356:DRSOTT>2.0.CO;2](https://doi.org/10.1175/1520-0469(1986)043<0356:DRSOTT>2.0.CO;2), 1986.
- 495 Trapp, R. J., Tessendorf, S. A., Godfrey, E. S., and Brooks, H. E.: Tornadoes from Squall Lines and Bow Echoes. Part I: Climatological Distribution, *Weather and Forecasting*, 20, 23-34, 10.1175/waf-835.1, 2005.
- Wang, H., Kong, F., Wu, N., Lan, H., and Yin, J.: An investigation into microphysical structure of a squall line in South China observed with a polarimetric radar and a disdrometer, *Atmospheric Research*, 2019.
- 500 Wang, X., Bian, H.-X., Qian, D.-L., Miao, C.-S., and Zhan, S.-W.: An automatic identifying method of the squall line based on Hough transform, *Multimedia Tools and Applications*, 80, 18993-19009, 10.1007/s11042-021-10689-3, 2021.
- WangHong, MaFeng-lian, and WangWan-jun: Doppler Radar Data Analysis of a SquallLine Process, *Desert and Oasis Meteorology*, 3, 39-43, 2009.
- Xiaohong, L., Wenjuan, Z., and Nengzhu, , A. F.: Lightning Activity in the Pre_TC Squall Line of Typhoon Lekima (2019) Observed by FY4A LMI and Its Relationship with Convective Evolution, *Remote Sensing Technology and Application*, 36, 2021.
- 505 Ye-qing, Y., Xiao-ding, Y., Yijun, Z., Hua, C., Ming, W., and Jin, L.: Analysis on a Typical Squall Line Case with Doppler Weather Radar Data, *Plateau Meteorology*, 27, 373-381, 2008.
- ZHOUZhenbo, MINJinzhong, and PENGXiayun: Extended—VAP Method for Retrieving Wind Field from Single-DopplerRadar.(I):Methods and Contrast Exprement, *PLATEAU METEOROLOGY*, 516-524, 2006.
- 510 Ziqi, J., Xinmin, W., Yansong, B., Han, L., Ming, W., and Mingyue', L.: Squall Line Identification Method Based on Convolution Neural Network, *JOURNAL OF APPLIED METEOROLOGICAL SCIENCE*, 32, 580-591, 2021.