

Article

Not peer-reviewed version

---

# An Automatic Near-Duplicate Video Data Cleaning Method Based on a Consistent Feature Hash Ring

---

[Yi Qin](#)<sup>\*</sup>, [Ou Ye](#)<sup>\*</sup>, Yan Fu

Posted Date: 1 April 2024

doi: 10.20944/preprints202404.0008.v1

Keywords: Video cleaning; consistent feature hash ring; feature distance means; mountain peak function; multi-head attention mechanism; near-duplicate videos



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

# An Automatic Near-Duplicate Video Data Cleaning Method Based on a Consistent Feature Hash Ring

Yi Qin <sup>1</sup>, Ou Ye <sup>1,\*</sup> and Yan Fu <sup>1</sup>

<sup>1</sup> College of Computer Science & Technology, Xi'an University of Science and Technology, Xi'an 710054, China

\* Correspondence: oye0928@xust.edu.cn

**Abstract:** In recent decades, with the ever-growing scale of video data, near-duplicate videos continue to emerge. Data quality issue caused by near-duplicate videos is becoming more and more prominent, which has affected the application of normal videos. Although current studies on near-duplicate video detection can be helpful to uncover data quality issues for videos, they still lack a process of automatic merging for the video data represented by high-dimensional features, which are difficult to automatically clean the near-duplicate videos to improve data quality for video datasets. At present, there are few studies on near-duplicate video data cleaning. The existing studies have the sensitive problems of video data orderliness and clustering initial center under a condition that prior distribution is unknown, which seriously affect the accuracy of near-duplicate video data cleaning. To address the above issues, an automatic near-duplicate video data cleaning method based on a consistent feature hash ring is proposed in this paper. First, a residual network with convolutional block attention modules, a long short-term memory deep network, and an attention model are integrated to construct an RCLA deep network with the multi-headed attention mechanism to extract spatiotemporal features of video data. Then, a consistent feature hash ring is constructed, which can effectively alleviate the sensitivity of video data orderliness while providing a condition of near-duplicate video merging. To reduce the sensitivity of the initial cluster centers to results of near-duplicate video cleansing, an optimized feature distance-means clustering algorithm is constructed by utilizing a mountain peak function on a consistent feature hash ring, which can implement automatic cleaning of near-duplicate video data. Finally, experiments are conducted based on a commonly used dataset named CC\_WEB\_VIDEO and a coal mining video dataset. Compared with some existing works, simulation results demonstrate that the performance of the proposed method.

**Keywords:** video cleaning; consistent feature hash ring; feature distance means; mountain peak function; multi-head attention mechanism; near-duplicate videos

## 1. Introduction

In recent years, with the innovative progress of video editing, 5G communication, and other related technologies, and the popularity of video-related applications and services has led to the continuous expansion of the scale of video data and shows a continuous exponential growth trend [1]. Take short videos as an example, the data from iiMedia Research shows that the scale of short video users in China has an obvious growth momentum, which exceeded 700 million in 2020. These users deliver vivid and diverse information resources through video creation, video sharing, and video recommendation to enrich daily lives by using short video platforms.

In fact, as the scale of video data increases, many similar videos continue to emerge after video editing, reissue a new version of the modified original, and other operations for the videos, which are also referred to as Near-Duplicate Videos (NDVs) [2]. In [3], near-duplicate videos are defined as identical or approximately identical videos are close to each other and hard to distinguish, but they are different in some detail. In general, near-duplicate videos are derived from the original video, which not only the copyright of the video producer has been illegally infringed illegally [4], but also affects the data quality of video data sets. For example, Wu et al. [5] use 24 keywords to search for

videos on common Web sites of YouTube, Yahoo!Video, and Google Videos, the results show that there are a lot of near-duplicate videos on the afore-mentioned Web sites. In individual cases, the redundancy can reach 93%. These near-duplicate videos will not only cause copyright issues and affect normal applications of video surveillance, video recommendation, etc., but also significantly reduce the data quality of video datasets, making the maintenance and management of video data more and more challenging. At present, some salient problems caused by near-duplicate videos have obtained increasingly attention in academia and industry.

From the perspective of data quality [6], video data quality a great deal of attention is paid to the overall quality of the video data set, and lay stress on the degree of data consistency, data correctness, data completeness, and data minimization are satisfied in the information systems. The emergence of near-duplicate videos will reduce the degree of data consistency and minimization for video data sets. These near-duplicate videos can be taken to be a kind of dirty data, which has a wide coverage, rich and diverse forms. Concretely, regardless of the stage of video collection, video integration, or video processing, it is possible to generate near- duplicate videos. For instance, in the video collection stage, they can be collected from different angles within the same scene; in the video integration stage, there may be near-duplicate videos with different formats and video lengths from different data sources; in the video processing stage, video copy, video editing, and other operations will produce mass near-duplicate videos. Since near-duplicate videos have the characteristics of concealment, dissemination, and harmfulness in video data sets, the tasks of how to identify and reduce these near-duplicate videos are difficult challenges.

Research on near-duplicate video detection can help us discover hidden near duplicate-videos in video datasets. Currently, various kinds of methodologies have been proposed in literature, and the implementation process mainly includes feature extraction, feature signature, and signature index. In either of these methodologies, feature extraction can be regarded as a key component of near-duplicate video detection. From the perspective of video feature representation, near-duplicate video detection methodologies can be categorized into the hand-crafted features-based methodology and high-level features-based methodology.

As mentioned before hand-crafted features-based methodology, local feature descriptors (such as SIFT, SURF, MSER, etc.) and global feature descriptors (such as a bag of the word and HOG) can be utilized to represent the spatial features of video data, and the computational cost of these feature descriptors is small. However, hand-crafted features-based methodology needs to rely on prior knowledge, which is difficult to accurately represent video features for video data with the characteristics of uneven illumination, local occlusion, and significant noise.

In recent years, through in-depth research on deep learning theories and methodologies, a variety of deep neural network modeling techniques have been applied to the detection of near-duplicate videos. The features of videos can be more accurately portrayed and represented to identify video data semantically similar by adopting convolutional neural network models [7], recurrent neural network models [8], and other deep neural network models [9]. Nevertheless, near-duplicate video detection methodologies can only identify the near-duplicate videos in a video dataset, which lack a process of feature sorting and automatic merging for the video data represented by high-dimensional features. Therefore, it is very challenging for them to automatically clean up redundant near-duplicate videos to reduce video copyright infringement and related issues caused by video copying, video editing, and other manual operations.

Currently, major studies focus on near-duplicate video detection [10] and retrieval [11], and less consideration is given to how to automatically clean near-duplicate videos from the perspective of data quality. Data cleaning modeling techniques are important technical ways to effectively reduce near-duplicate data and improve data quality. By using this kind of modeling technique, near-duplicate data existing in the datasets can be automatically cleaned, so the datasets meet data consistency, completeness, correctness, and minimization, and achieve high data quality. At present, data cleaning modeling techniques have been studied more deeply in big data cleaning [12], stream data cleaning [13], contextual data cleaning [14], etc., which can effectively address the data quality issues at the instance layer and schema layer. However, there is still a significant gap in research on

near duplicate video cleaning, and there is a sensitive problem of video data orderliness and clustering initial center sensitive problem under a condition that prior distribution is unknown in the existing studies, which seriously affect the accuracy of near-duplicate video cleaning.

In this paper, an automatic near-duplicate video cleaning method based on a consistent feature hash ring (denoted as RCLA-HAOPFDMC) is proposed to address the above-mentioned issues, which consists of three parts: high-dimensional feature extraction of video data, consistent feature hash ring construction, and cluster cleaning modeling based on a consistent feature hash ring. First, a residual network with convolutional block attention modules, a long short-term memory (LSTM) deep network model, and an attention model are integrated to extract temporal and spatial features from videos by constructing a multi-head attention mechanism. Then, a consistent feature hash ring is constructed, which can effectively alleviate the sensitivity of video data orderliness while providing a condition of near-duplicate video merging. Finally, to reduce the sensitivity of the initial cluster centers to results of near-duplicate video cleansing, an optimized feature distance-means clustering algorithm is constructed by utilizing a mountain peak function on a consistent feature hash ring to implement automatic cleaning of near-duplicate videos. A commonly used dataset named CC\_WEB\_VIDEO [5] and a coal mining video dataset [15] are used to confirm the practical effect of our proposed method. The contributions are summarized as follows: (1) a novel consistent feature hash ring is constructed, which can alleviate the sensitivity issue of videos data orderliness while providing a condition of near-duplicate videos merging; (2) an optimized feature distance-means clustering algorithm is constructed by utilizing a mountain peak function on consistent feature hash ring to merge and clean up near-duplicate videos; (3) the method presented in this paper is successful on the highly difficult CC\_WEB\_VIDEO and coal mining video dataset, where the coal mining video dataset has complex context scenes.

The following is the organizational structure of the remaining parts of this article. In Section 2, a brief review of related works on near-duplicate video detection and data cleaning is presented; an automatic near-duplicate video cleaning method based on a consistent feature hash ring is proposed in Section 3. The experimental results validate the performance of the method presented in Section 4. Finally, the paper is summarized.

## 2. Related Work

In this section, the previous near-duplicate video detection methodologies and image/video cleaning methodologies are briefly reviewed. First, some hand-crafted features-based methodologies and high-dimensional features-based methodologies for near-duplicate video detection are recalled, then some data cleaning methodologies for image cleaning and video cleaning are reviewed. Finally, the shortcomings of the above-mentioned methodologies are analyzed.

### 2.1. Near-Duplicate Video Detection Methodologies

In the past decade, hand-crafted features are widely used to near-duplicate video detection, such as SIFT, HOG, and MSER. For example, the study in [16] adopts SIFT local feature to encode temporal information, generates temporal set SIFT features by tracking SIFT, and combines local sensitive hash algorithms to detect near-duplicate videos. Despite SIFT features of a video frame can maintain invariance to rotation, scaling, and illumination changes, as well as maintain a certain degree of affine transformation and noise, some of the invariance of SIFT will be damaged to a certain extent during strong Camcording. Henderson et al. [17] adopt key point features from a Harris corner, SURF, BRISK, FAST, and MSER descriptors to detect video frame duplication. This method incorporates different local features to represent video frames, but ignores the global feature and spatiotemporal features that video frames have. Zhang et al. [18] integrate Harris 3D spatiotemporal feature and HOG/HOF global feature descriptors to detect near-duplicate news web videos, and the Jaccard coefficient is applied to similarity metric, but there exists the issue of inefficient detection in this method. In [19], a new near-repeat video detection system, CompoundEyes, is proposed, which combines seven hand-made features (such as color consistency, color distribution, edge direction, motion direction, etc.) to improve the efficiency of near-repeat video detection. However, this method

is susceptible to feature changes. The work in [3] adopts an unsupervised cluster algorithm based on temporal and spatial key points to automatically identify and classify near-duplicate videos, but the results of near-duplicate video detection are sensitive to initialize the cluster center.

In general, the combination of spatial and temporal features can be more comprehensively and accurately to represent the spatiotemporal information contained in video data than the representation of a single low-level feature, hence the methodologies based on spatiotemporal features can identify near-duplicate videos more accurately. However, the methodologies based on low-level features need to know prior knowledge, and the results of near-duplicate video detection are easily affected by disturbances from illuminations, occlusions, distortion, etc.

Recently, various deep network models have been utilized to detect near-duplicate videos, which have more excellent representational capacity than the methodologies based on hand-crafted features. For instance, the study in [7] presents a survey on the utilization of deep learning techniques and frameworks. Nie et al. [20] use a pre-trained convolutional neural network model to extract high-dimensional features of videos, and a simple but efficient multi-bit hash function is proposed to detect near-duplicate videos. This is a supervised joint view hashing method, which can improve the performance of accuracy and efficiency. However, the distribution of near-duplicate video data in the video dataset is usually unknown in practical applications, so the supervised joint view hashing method is limited. In [21], the near-duplicate video is detected by combining the two-stream network model of RGB and optical flow, multi-head attention and Siamese network model. The limitation of this method is that it adopts a cosine distance function to measure the similarity between every two videos, which results in relatively low efficiency. Moreover, a neighborhood attention mechanism is integrated into an RNN-based reconstruction scheme to capture the spatial-temporal features of videos, which is used to detect near-duplicate videos [22]. The work in [23] uses a temporal segment network model to detect near-duplicate video data. These above-mentioned models based on temporal networks can detect near-duplicate videos by capturing the temporal features of videos. In [24,25], a Parallel 3D ConvNet model and a spatiotemporal relationship neural network model are adopted to extract spatiotemporal features to detect near-duplicate videos.

In summary, high-dimensional features-based methodologies can achieve better performance than hand-crafted features-based methodologies, they can reduce the impacts of disturbances from illuminations, occlusions, and distortion on the model results. However, near-duplicate videos can be identified directly by either hash mapping or similarity metric, but the abovementioned methodologies lack a process to automatically merge data with high-dimensional features, which are more difficult to clean up near-duplicate videos automatically.

## *2.2. Data Cleaning Methodologies*

Data duplication may have the following reasons: data maintenance, manual input, device errors, and so on [26], data cleaning modeling techniques are effective ways to automatically clean and reduce near-duplicate data. Recently, the amount of literature on the topic of data cleaning [27] has shown a rapid growth trend, most of the existing works are concentrated stream data cleaning and spatiotemporal data cleaning.

In a line works of stream data cleaning, for instance, the literature [28] proposes a stream data cleaning method named SCREEN, which can clean up the steam data by finding the repair sequence with the smallest difference from input, construct an online cleaning model and calculate the local optimal of the data point. However, this method does not guarantee that near-duplicate data are exactly adjacent to each other in the same sliding window. The work in [29] proposes a streaming big data system, which is based on an efficient, compact, and distributed data structure to maintain the necessary state for repairing data. Additionally, it improves cleaning accuracy by supporting rule dynamics and utilizing sliding window operations. The limitation of this method is that the fixed size of the sliding window has a significant impact on the cleaning results. In [30], a sliding window and K-means cluster algorithm are adopted to clean stream data, but the result of this method is sensitive to the initialization of cluster centers.

To sum up, although methodologies of stream data cleaning can clean up the stream data cleaning effectively, there are limitations in that the results of models are sensitive to the fix-size sliding window and pre-defined the initialized clustering center.

In line works of time series data cleaning, for example, Ranjan et al. [31] utilize a k-nearest neighbor algorithm and a sliding window prediction approach to clean time series data on a set of nonvolatile and volatile time series. This method can optimize the width of the sliding window to enhance the performance, but to optimize parameters that affect performance, a general scheme needs to be developed. In [32], a top-N keyword query processing method is proposed, which is based on real-time entity parsing to clear data sets with duplicate tuples. The limitation of this method is that the selection of keywords has a salient impact on the results. The study in [33] proposes an approach of real-time data cleaning and standardization, which clarified the workflows of data cleaning and data reliability, and it can be adapted to clean up near-duplicate time series data. However, this approach is not enough to describe the details of real-time data cleaning.

Through the studies of time series data cleaning, it is found that there is a prevalence of erroneous data in the industrial field, hence the studies mainly focus on cleaning the erroneous data, and less on cleaning near-duplicate data. In the existing studies, k-nearest neighbor and top-K algorithms are widely adopted to clean near-duplicate time series data since they do not rely on prior knowledge of the distribution of time series data. Nevertheless, the pre-setting of the k parameters has a significant impact on the cleaning results.

In a line works of spatiotemporal data cleaning, take the literature [34] as an example, a probabilistic system named CurrentClean is presented, which uses a spatiotemporal probabilistic model and a set of inferences to identify and clean stale data in the database. This probabilistic system applies to data cleaning with spatiotemporal features in relational databases, but it is difficult to apply to unstructured data cleaning with high dimensional spatiotemporal features. To address this issue, the study in [15] proposes a method to clean up near-duplicate videos by using locality-sensitive hashing and a sorted neighborhood algorithm. However, this method is challenging to use SIFT and SURF hand-crafted features accurately to portray video features, and the use of a sorted neighborhood algorithm causes a more prominent orderliness-sensitive problem of video data. The work in [35] achieved an improvement in the quality of the image dataset by automatically clearing minority images using a convolutional neural network model. However, the completeness of image datasets may be destroyed when cleaning images of the minority classes. Fu et al. [36] propose a near-duplicate video cleaning method based on VGG-16 deep network model and feature distance means clustering fusion to improve the data quality of video datasets, which takes less account of the temporal feature representation of videos and suffers from clustering initial center sensitive problem under a condition that prior distribution is unknown. Moreover, a novel content based on the video segmentation identification scheme is proposed to reduce the mass of near-duplicate video clips [37]. H. Chen et al. [38] utilize the similarity measurement to clean the duplicate annotations of video data in the MSR-VTT dataset.

In summary, there are few studies on data cleaning for unstructured data with spatiotemporal features, such as video and audio data, due to less consideration from the perspective of data quality. Recently, the studies on near-duplicate video cleaning mainly have the following issues: (1) it is more difficult to be able to arrange all near-duplicate videos nearby sorting algorithms, so the accuracy of cleaning is more sensitive to the data orderliness; (2) by utilizing the idea of clustering, a video data with the most significant features can be retained in all near-duplicate data and the rest deleted, but the setting of clustering initial center is more sensitive to cleaning results under a condition that prior distribution is unknown.

To address these two issues, we consider constructing a novel consistent feature hash ring based on optimizing video feature extraction to map video data to low-dimensional space, which is used to reduce data orderliness sensitivity issues caused by data sorting, and provide a condition of near-duplicate video merging. On this basis, an optimized feature distance-means clustering algorithm is constructed, which merges a mountain peak function on a consistent feature hash ring to overcome

the clustering initial center-sensitive problem under a condition that the prior distribution is unknown.

### 3. The Proposed Method

In this section, how to utilize the proposed novel automatic near-duplicate video cleaning method based on a consistent feature hash ring is described, which can automatically clean up near-duplicate videos to improve the data qualities of video datasets. It is important to note here that the concepts between the data quality of video datasets and video quality are different. Normally, video quality is concerned with the clarity of videos and involves performance metrics such as resolution and frames per second. The data quality of the video dataset is concerned with the degree to which the video dataset satisfies data consistency, data completeness, data correctness, and data minimization. The goal is to remove redundant near-duplicate videos from video datasets by the method proposed in this paper, so that the video datasets consistently maintain high data quality.

To achieve this goal, three main stages need to be accomplished: feature representations of video data, identification of near-duplicate videos, and deletion of near-duplicate videos. It should be noted that if all the identified near-duplicate videos are removed, the data completeness and data correctness of video datasets will be affected. If only some the near-duplicate videos are removed, the data consistency and data minimization of video datasets will be affected. Therefore, it is a difficult challenge to retain video data with the most salient features to ensure the data completeness and data correctness of video datasets, while removing the rest near-duplicate video data to ensure data consistency and data minimization of video datasets.

At present, considering the time cost of near-duplicate video cleaning, the key insight of existing works is to overcome the above challenges by exploiting data ordering and clustering to retain video data with the most salient features and remove the rest of near-duplicate videos. However, the sensitivities of cleaning effect to data orderliness and initial clustering center setting are major issues. Inspired by the distributed big data storage processing, a consistent feature hash ring is constructed in this paper. The advantage of utilizing a consistent feature hash ring is to reduce the impact of data sorting on the cleaning results by mapping video data with high-dimensional features to a feature hash ring, while providing a condition for removing near-duplicate videos. Figure 1 outlines our approach, which consists of three parts: high-dimensional feature extraction of videos, construction of a consistent feature hash ring, clustering, and cleaning near-duplicate videos based on a consistent feature hash ring. Each of these sections is explained next.

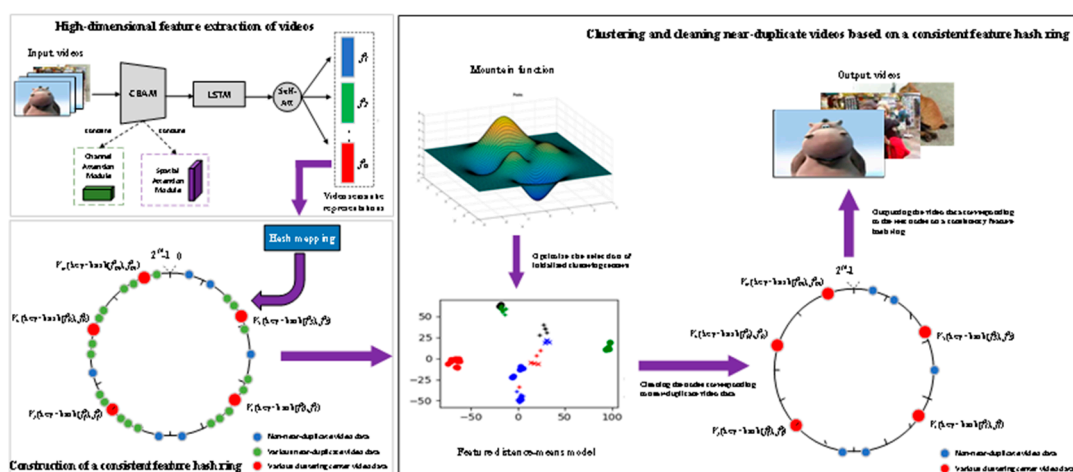


Figure 1. The overall framework of our proposed method.

#### 3.1. High-Dimensional Feature Extraction of Videos

The feature representation of video data is an important stage in the data cleaning process of videos. Currently, several convolutional neural network models have been adopted in the study of

image and video data cleaning for feature representation of image or video data, such as AlexNet [35], GoogleNet [35], VGG-16 [36], and ResNet50 [39]. Due to the spatiotemporal features of video data, it is not enough to rely on the above-mentioned convolutional neural network models for spatial feature representations of videos, and the extracted video features are less likely to highlight the spatiotemporal features of salient regions in video data, which will affect the accurate representation of video semantics. To overcome such a limitation, a residual network with convolutional block attention modules is adopted firstly to extract spatial features of video data, the channel attention and spatial attention modules in this convolutional block attention modules can effectively improve the representation capability of spatial features in the saliency region of video data. Then, the above network and a long short-time memory model are integrated to extract the spatiotemporal features of video data. Finally, to highlight the role of key information in video data on the video semantic representations, an attention model based on the above network models is introduced, thereby along three independent dimensions of channel, space, and time series to construct a video spatiotemporal feature extraction model with multi-head attention mechanism, which is named as RCLA model in this paper. The concrete architecture of the RCLA model is shown in Figure 2.

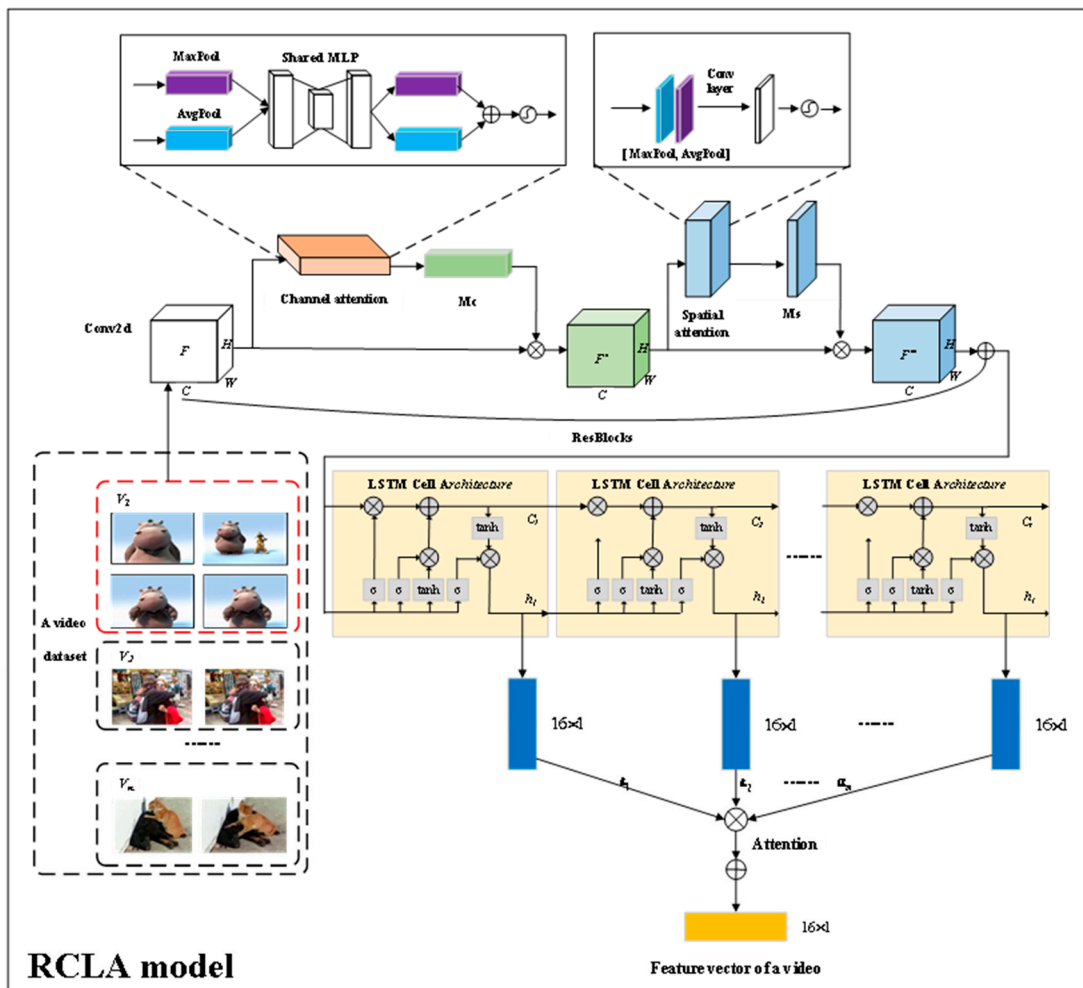


Figure 2. The architecture of the RCLA model.

Considering the scale of data used for training and testing in this paper, several video data as the training samples are firstly selected to input a residual network with convolutional block attention modules. Let the size of the above video data be  $w \times h \times c \times l$ , where  $w \times h$  represents the size of a video frame,  $c$  represents the number of channels per frame, and  $l$  represents the number of frames of the video data [40]. Before training the 34-layer residual network, set the values of  $w$  and  $h$  to 224, and the value of  $c$  to 3. In addition,  $7 \times 7$  the convolutional kernel with stride 2 in the convolution layer are firstly fixed, and then the pooling window with stride 2 in the pooling layer is

fixed  $3 \times 3$  to implement the convolutional operation and max-pooling process. During the above process, a BatchNorm2d method is used to normalize the input data, and a Rectified Linear Units (ReLU) activation function is used to alleviate the problem of gradient dispersion. Thus, a feature map  $F$  of size  $56 \times 56 \times 3$  for a video frame can be obtained. Then, the  $F$  into the middle part of the residual network with 34 layers is input. This part is composed of 4 blocks, which respectively include 3, 4, 6, and 3 residual blocks, and each residual block contains a convolutional block attention module (CBAM) [41]. In this module, the details of spatial features in a video frame can be profiled by constructing a channel attention map  $M_c$  and a spatial attention map  $M_s$ . Specifically, the spatial information of feature map  $F$  is aggregated by first performing maximum pooling and average pooling operations on  $F$ , respectively, and this spatial information is input to a multilayer perceptron (MLP) network model with 2 layers, whose two layers share the weights  $W_0 \in \mathbb{R}^{c/r \times c}$  and  $W_1 \in \mathbb{R}^{c \times c/r}$  ( $r$  is the reduction ratio,  $r$  is set to 16 in this paper), respectively, so that the channel attention map  $M_c$  can be obtained by using Eq. (1) and an intermediate feature map  $F'$  is computed by using Eq. (2):

$$M_c(F) = \sigma(W_1(W_0(F_{avg}^c)) + W_1(W_0(F_{max}^c))) \quad (1)$$

$$F' = M_c(F) \otimes F \quad (2)$$

where  $\sigma(\cdot)$  denotes the sigmoid function; The meanings of  $F_{avg}^c$  and  $F_{max}^c$  are average-pooled features and max-pooled features respectively; " $\otimes$ " is used for element-by-element multiplication.

Then, the average pooling and maximum pooling operations are performed on  $F'$  to generate features  $F_{avg}^s \in \mathbb{R}^{1 \times h \times w}$  and  $F_{max}^s \in \mathbb{R}^{1 \times h \times w}$ , respectively. On this basis, the spatial attention map  $M_s$  is calculated by Eq. (3), and the final refined feature map  $F''$  can be obtained by Eq. (4) as follows:

$$M_s(F') = \sigma(conv([F_{avg}^s \otimes F_{max}^s])) \quad (3)$$

$$F'' = M_s(F') \otimes F' \quad (4)$$

where  $conv(\cdot)$  denotes a convolution operation.

Through the above-mentioned different residual blocks with convolutional block attention modules, the feature vectors  $f_{rc}$  of size [512, 1] can be obtained to represent the spatial features exhibited by video frames.

Then, the spatial features  $f_{rc}$  are input into the long short-term memory network (LSTM) to further extract the temporal features of video data. Considering the size of datasets used in this paper, a one-layer LSTM with  $N$  (we set  $N=16$  in this paper) hidden layer nodes is employed, which is consist of an input gate, forget gate, and output gate. The hidden state  $h_t^s$  at  $t$  moment is calculated as shown in Eq. (5):

$$h_t^s = \text{LSTM}(f_{rc}, W_{ls}, h_{t-1}^s; \theta_{ls}) \quad (5)$$

where  $\text{LSTM}(\cdot)$  denotes the formalized function of an LSTM network model;  $W_{ls}$  denotes a parameter matrix learned during training of the LSTM network model;  $h_{t-1}^s$  denotes the hidden state at  $t-1$  moment; and  $\theta_{ls}$  denotes the hyperparameters of an LSTM network model. Through the output layer of an LSTM network model, the spatial-temporal features  $f_{st}$  of size  $5 \times 16 \times 1$  can be obtained to represent video data.

To focus on the visual features of different video frames to highlight the semantic contents of video data, an attention module based on the above models is introduced, as shown in Eq. (6):

$$f = \text{Att}(h_t^s, f_{st}; W_{\text{Att}}) \quad (6)$$

where  $\text{Att}(\cdot)$  denotes a standard additive attention function;  $\mathbf{W}_{\text{Att}}$  denotes the weight vectors in an attention module;  $\mathbf{f}$  denotes the semantic features of video data obtained by using the multi-head attention mechanism, and their size are  $16 \times 1$ .

When training the RCLA model, the goal is to minimize the learning loss of each deep neural network model. Considering the cascade relationship between each of the above neural network models, the output of the above attention module as the input of a loss function is used to perform optimization of an RCLA model, the above-mentioned loss function is constructed as shown in Eq. (7) and Eq. (8):

$$\mathbf{y}_s = \arg \min(\|\mathbf{y}_i - \mathbf{y}_{\text{video-seed}}^j\|_2^2) \quad \text{s.t.} \quad i \in [1, N_v], j \in [1, N_{v-c}] \quad (7)$$

$$L(y'_i, \mathbf{y}_i) = -\sum_{i=1}^{16} y'_i \times \log(\text{sim}(\mathbf{y}_i, \mathbf{y}_s)) \quad (8)$$

where  $\mathbf{y}_i$  denotes a feature vector of the  $i^{\text{th}}$  video data in a video dataset with  $N_v$  video data,  $\mathbf{y}_{\text{video-seed}}^j$  denotes a feature vector of the  $j^{\text{th}}$  seed video in the above video dataset,  $N_{v-c}$  denotes the total number of preset seed videos,  $y'_i$  denotes the label corresponding to  $\mathbf{y}_i$ ,  $\text{sim}(\cdot)$  denotes a function of the similarity measurement.

### 3.2. The Construction of a Consistent Feature Hash Ring

To efficiently identify near-duplicate video data with high-dimensional features and reduce the impact of high-dimensional feature sorting on near-duplicate video data cleaning, a consistent feature hash ring is constructed, which is inspired by the way of using distributed hash tables to address the problems of load balancing and high scalability in big data storage, as shown in Figure 3.

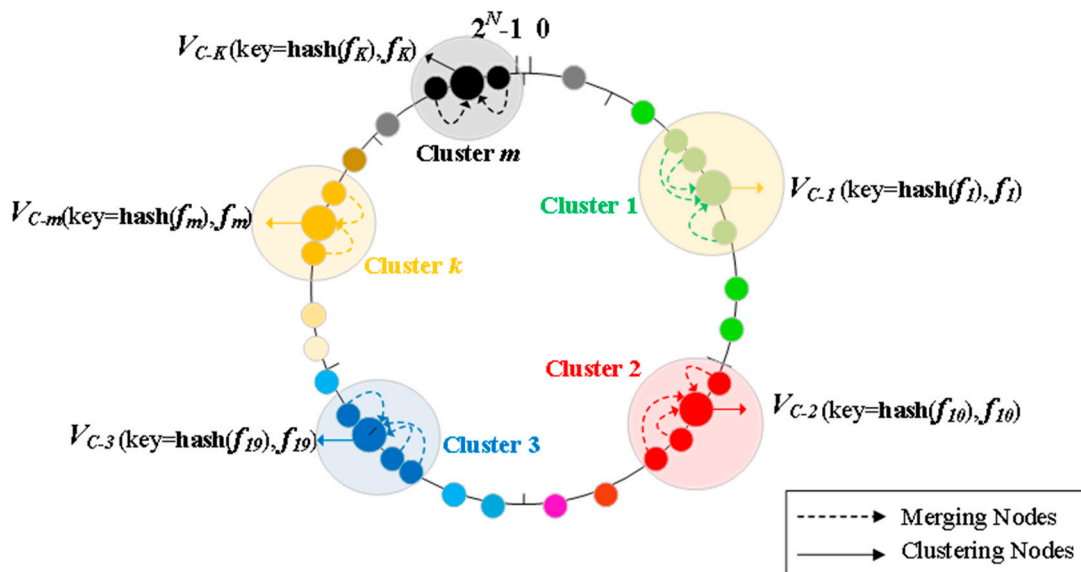


Figure 3. The structure of a consistent feature hash ring.

Assuming that a hash function can be used to map the high-dimensional features of all video data into a set of the hash values, and the largest hash value is less than  $2^N$ , then a consistent feature hash ring is a virtual ring constructed through this set of hash value, and the range of hash value space of this ring is  $[0, 2^N - 1]$  (hash value is a 32-bit unsigned integer, and  $N=11$  is set in this paper). An entire spatial distribution of the consistent feature hash ring is organized in a clockwise direction. The hash value at the top of this ring as the starting point is 0, the hash value at the first point on the left of 0 is  $2^N - 1$ , and the hash value at the first point on the right of 0 is 1. By analogy, all the hash values are distributed from small to large clockwise until they return to the starting point.

When constructing a consistent feature hash ring, a hash function is designed as shown in Eq. (9), which can be used to calculate the hash values corresponding to each video and map all video data to the consistent feature hash ring.

$$\text{hash}(\mathbf{f}_i) = [\text{binary}(\mathbf{f}_i) + 1] \bmod (2^N - 1) \quad (9)$$

where  $\text{binary}(\cdot)$  denotes a binary encoding function, and  $\mathbf{f}_i$  denotes a high-dimensional feature vector of the  $i$ th video data.

By utilizing the binary function as shown in Eq. (10), the high-dimensional features of all video data can be mapped into compact binary codes, which can not only decrease the dimensions of video features, but also perform metric operations in the low-dimensional space.

$$\text{binary}(\mathbf{f}_i) = \text{MD5}(\|\tanh(\mathbf{f}_i)\|) \quad (10)$$

where  $\text{MD5}(\cdot)$  denotes an encryption function;  $\tanh(\cdot)$  is an activation function.

Specifically, a  $\tanh(\cdot)$  activation function is utilized to perform a nonlinear variation of the video data with a high-dimensional feature. On this basis, to improve the sparsity of the hash distribution on the consistent feature hash ring, a  $\ell_1$  norm is used to calculate a value to represent video data. Subsequently, this value is encrypted with the MD5 encryption algorithm and converted into a fixed-length binary code as a hash value of the video data (the length of a binary code is set to 16 in this paper). Finally, a linear detection method is used when storing multiple hash values to avoid the same hash addresses being preempted by different hash values, i.e., an address is assigned by adding 1 to the back, and the modulus of the maximum value on a consistent feature hash ring is taken as the upper bound of the address range until there is a free address. Here, the modulo operation is to ensure that the location found is in the effective space of  $2^N - 1$ . Thus, the  $i$ th video data can be mapped as video hash feature points  $x_i$  on a consistent feature hash ring in the form of two-dimensional coordinates  $(\text{hash}(\mathbf{f}_i), \mathbf{f}_i)$ .

### 3.3. FD-Means Clustering Cleaning Optimization Algorithm with Fused Mountain Peak Function

The efficiency and accuracy of partitioning clustering algorithm are closely related to the selection of initial clustering center. The FD-Means clustering algorithm [36] randomly selects several initial cluster centers for multiple clustering, and finally selects the optimal clustering centers as the initial clustering center. However, it has a large amount of calculation, and poor effect leads to the volatility of clustering results. To address this issue, a mountain peak function is fused with the FD-Means clustering algorithm to optimize the selection of the initial clustering centers to automatically clean the near duplicate video data accurately.

Specifically, assuming that a set of video features  $S_f = \{\|\mathbf{f}_1\|_2, \|\mathbf{f}_2\|_2, \dots, \|\mathbf{f}_i\|_2, \dots, \|\mathbf{f}_N\|_2\}$ , where  $\mathbf{f}_i$  denotes a vertical ordinate of  $x_i$ , and  $N_v$  denotes the total number of video data. To select several initial clustering centers of video data, all the data samples on the consistent feature hash ring are first divided into a finite grid, and all the cross points of a  $K \times K$  grid can be used as the candidate centroids of the clustering centers, as shown in Eq. (11) and Eq. (12):

$$T_{\text{interval}} = \frac{\max(S_f) - \min(S_f)}{N_v} \quad (11)$$

$$T_{\text{interval}} = \frac{\max(S_f) - \min(S_f)}{N_v} \quad (12)$$

where  $T_{\text{interval}}$  denotes the partition interval of a grid;  $i_c$  denotes an index of the  $i$ th cluster center  $V_C^{i_c}$ , and  $i_c \in [1, K]$ ;  $P_C^{i_c}$  denotes the value of the  $i$ th cluster center  $V_C^{i_c}$  in a grid. Moreover,  $K$  can be determined by inequality  $(K-1)^2 < T_{\text{interval}} \leq K^2$ .

Subsequently, the points with higher density in grid space are found by constructing a mountain peak function for each cross point to calculate its peak value  $H_{V_c}$ . For example, a peak value of the  $i$ th cluster center  $V_c^i$  is calculated as shown in Eq. (13):

$$H_{V_c}^i = \sum_{i=1}^{N_v} \exp\left(-\frac{(P_c^i - \|\mathbf{f}_i\|_2)^2}{2\sigma^2}\right) \quad (13)$$

where  $\sigma$  is a constant value.

On this basis, the video hash feature points corresponding to the cross points of  $K$  maximum peak values are selected sequentially as the initialized clustering centers on a consistent feature hash ring, denoted as  $V_c = \{V_{c-1}, V_{c-2}, \dots, V_{c-K}\}$ .

Since the FD-Means clustering algorithm seldom considers the curse of dimensionality when using Euclidean distance to measure the similarities between video features, a new similarity measurement function is constructed, as shown in Eq. (14):

$$\text{Dist}_{\text{VFP}}(\mathbf{f}_i, \mathbf{f}_j) = \alpha \times |FD(\mathbf{f}_i, \mathbf{f}_j)| + |\text{hash}(\mathbf{f}_i) - \text{hash}(\mathbf{f}_j)| \quad (14)$$

where  $\alpha$  denotes a weight factor;  $FD(\mathbf{f}_i, \mathbf{f}_j)$  denotes the Euclidean distance between any  $i$ th and  $j$ th video hash feature points, as shown in Eq. (15):

$$FD(\mathbf{f}_i, \mathbf{f}_j) = \begin{pmatrix} \|\mathbf{f}_1 - \mathbf{f}_1\|_2 & \|\mathbf{f}_1 - \mathbf{f}_2\|_2 & \cdots & \|\mathbf{f}_1 - \mathbf{f}_m\|_2 \\ \|\mathbf{f}_2 - \mathbf{f}_1\|_2 & \|\mathbf{f}_2 - \mathbf{f}_2\|_2 & \cdots & \|\mathbf{f}_2 - \mathbf{f}_m\|_2 \\ \vdots & \vdots & \ddots & \vdots \\ \|\mathbf{f}_q - \mathbf{f}_1\|_2 & \|\mathbf{f}_q - \mathbf{f}_2\|_2 & \cdots & \|\mathbf{f}_q - \mathbf{f}_m\|_2 \end{pmatrix} \quad (15)$$

where  $q$  and  $m$  denote the dimensions of two feature vectors.

When updating an FD-Means cluster center, the sum of distances between a video hash feature point and other video hash feature points in a cluster is first calculated, and the video hash feature point with the smallest sum of distances is selected as a cluster center through iteration, as shown in Eq. (16):

$$\text{Dist}(x_i, d) = \sum_{j \in d - x_i} \text{Dist}_{\text{VFP}}(\mathbf{f}_i, \mathbf{f}_j), \text{ s.t. } x_i \in d \quad (16)$$

$$V_{c-i}^* = \text{argminDist}(x_i, d)$$

where  $d$  denotes a cluster;  $V_c^{i*}$  denotes a new cluster center after the initial cluster center of  $d$  is updated.

When updating a cluster, the distance between the video hash feature points of the non-cluster centers and all the cluster centers are first calculated. In the K-Means clustering algorithm, all points of non-cluster are divided into the nearest clusters according to the nearest neighbor principle. Unlike the K-Means clustering algorithm, the optimized FD-Means clustering algorithm compares the distances between the video hash feature points of non-cluster centers and all cluster centers. If the minimum distance is not less than the given threshold  $\delta$ , the video hash feature points are used as the cluster centers of new clusters; otherwise, according to the nearest neighbor principle, they are grouped into the nearest clusters, and the automatic clustering of video hash feature points on a consistent hash ring is finally achieved, as shown in Eq. (17) and Eq. (18):

$$\min \text{Dist}_c(x_i, V_c; \delta) = \text{Dist}(x_i, V_c) - \delta \quad (17)$$

$$d^* = \begin{cases} d \cup x_i, & \text{s.t. } \min \text{Dist}_c(x_i, V_c; \delta) < 0 \\ x_i, & \text{otherwise} \end{cases} \quad (18)$$

where  $d^*$  denotes an updated cluster.

In the automatic cleaning of near-duplicate video data, all the near-duplicate videos obtained by clustering cannot be eliminated to ensure the data consistency, completeness, correctness, and minimization of a video set. Therefore, a representative video can be retained from the near-duplicate videos, and other near-duplicate video data can be automatically deleted to improve the video data quality.

At present, the traditional image or video data de-duplication methods tend to retain the first detected data and remove other near-duplicate data in an ordered sequence, but such methods are randomized at different settings of sequence length, which will lead to fluctuations in the results of data cleaning. Hence, a video in a cluster cannot be arbitrarily selected as a seed video or representative video to be reserved. In this paper, the video data corresponding to the cluster centers are as the representative video data to be retained, and others in the clusters are deleted, as shown in Eq. (19):

$$D^* = F(V_c) \quad (19)$$

where  $F(\cdot)$  denotes a mapping function between a set  $V_c$  of cluster centers and a set  $V$  of the corresponding video data, and the mapping relationship is expressed as  $F: V_c \rightarrow V$ ,  $D^*$  denotes the original video data set is updated after the mapping function  $F(\cdot)$ .

Through the above processes, the automatic cleaning of near-duplicate video data is finally achieved.

#### 4. Experimental Evaluation

The extensive experiments on a commonly used dataset named CC\_WEB\_VIDEO and a coal mining video dataset are conducted to evaluate the performance of RCLA-HAOPFDMC (the proposed method in this paper) and compare them with other representative advanced methods, such as the CBAM-Resnet [42] and BS-VGG16 [36] methods. All experiments were conducted on the same machine, which had an 8-Inter Xeon processors with 2.10GHz and a graphics card NVIDIA Corporation GP102 with 16G memory, the programs are implemented based on Python version 3.6.5 and PyTorch version 0.4.0. Next, we will provide a detailed explanation of the experiment and results.

##### 4.1. Dataset and Evaluation Criteria

In this paper, the CC\_WEB\_VIDEO and coal mining video datasets are used to carry out the comparative experiment of the proposed method. The CC\_WEB\_VIDEO dataset contains 24 scenes and a total of 13129 video data. This paper randomly selects 63 videos from scenes ("The Lion Sleeps Tonight", "Evolution of dance", "Folding Shirt", "Cat Massage", and "ok go-here it goes again" scenes) to verify the effectiveness of the proposed method. The coal mining video dataset includes 125 video data with 10 scenes, which are all used to test the performance of the method presented in this paper.

We use common metrics such as accuracy, precision, recall, and F1-score to evaluate video data cleaning, in order to evaluate the performance of the method proposed in this paper. The expression is as follows:

$$precision = \frac{TP}{TP + FP} \quad (17)$$

$$Recall = \frac{TP}{TP + FN} \quad (18)$$

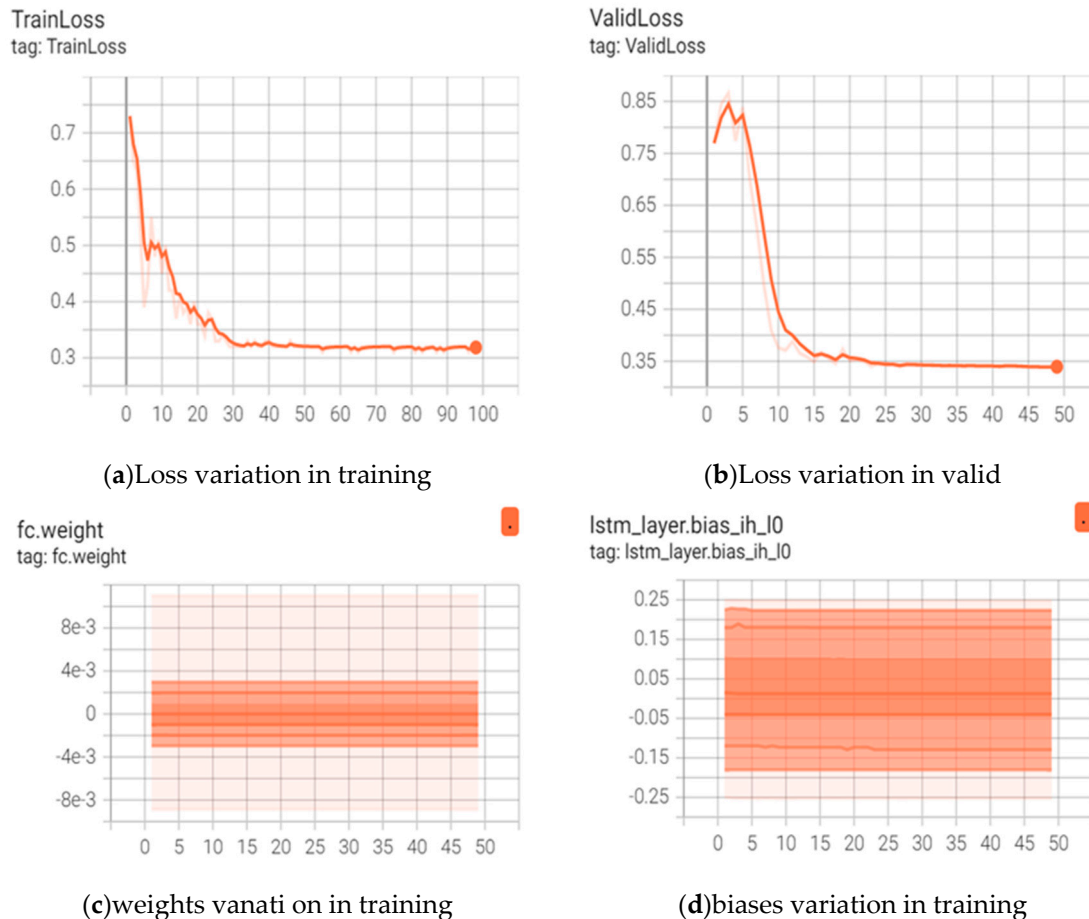
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (17)$$

$$F1-score = \frac{2 \times precision \times recall}{precision + recall} \quad (18)$$

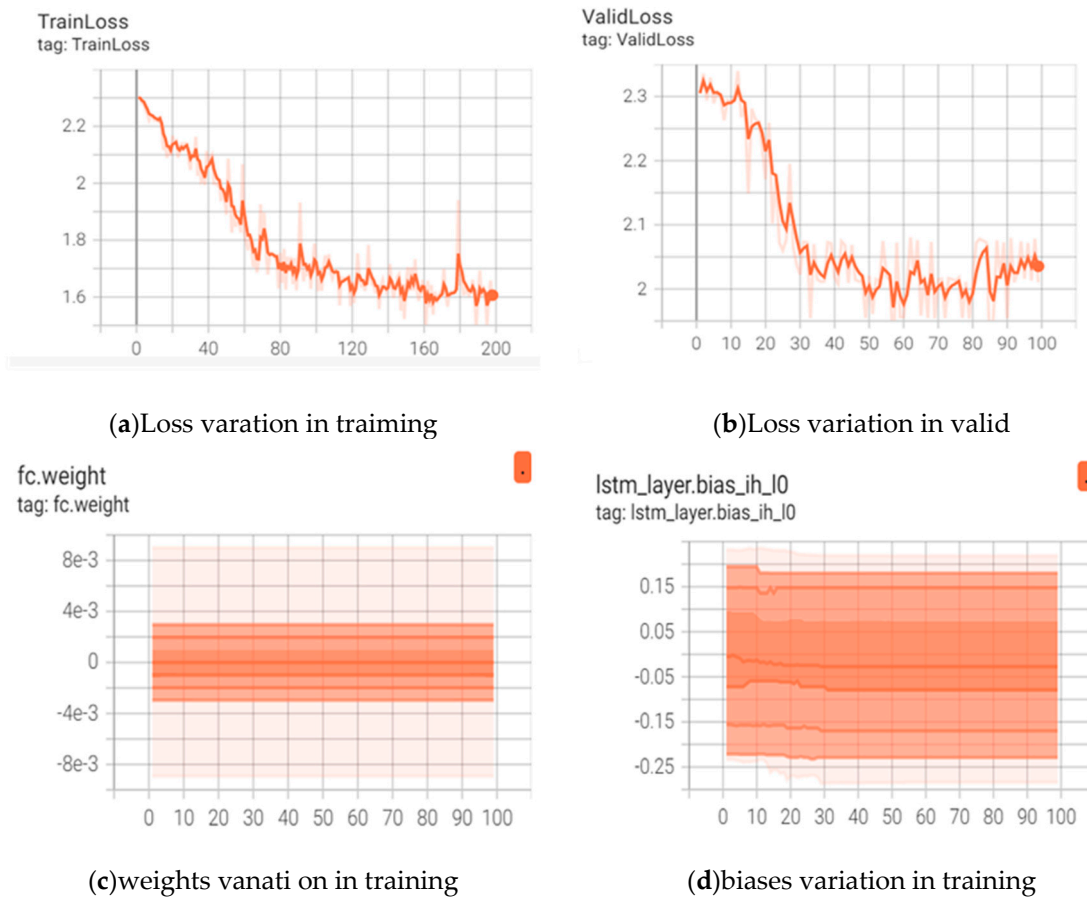
where TP is the number of true positive samples; FN is that of the false negative samples; FP is that of the false positive samples, and TN is that of the true negative samples.

#### 4.2. Experimental Results and Analysis

For the CC\_WEB\_VIDEO and coal mining video datasets, when training the RCLA model, the initialization of weights and bias variables is randomly generated. Moreover, The overfitting problem is solved by dropout function and parameter-sharing methods, and the loss function is optimized by Adam algorithm. Figs. 4 and Figs. 5 show the weight variable changes and training losses of CC\_WEB\_VIDEO and coal mine video datasets. It can be seen from Figs. 4 and 5 that the designed loss function in this paper is converged, and the different weights in the full connection layer of the RCLA model change in the range of -0.003 to 0.003 in Figs. 4 and 5. Moreover, Figure 4 shows that the range of different deviation values is -0.2 to 0.2 and Figure 5 shows -0.24 to 0.2. The above experimental results show that there is no overfitting issues during the training and validation process of the experiment.



**Figure 4.** The visualization of loss and weight variation during the CC\_WEB\_VIDEO dataset.



**Figure 5.** The visualization of loss and weight variation during the coal mining video dataset.

Since the number of hidden layers in the LSTM network model and attention size in the attention module have a great impact on the performance of the proposed method, the CC\_WEB\_VIDEO public dataset is used to evaluate the performance indicators of different number of hidden layers and attention size. The experimental results are shown in Tables 1 and 2.

**Table 1.** The experimental results of different parameter settings for the number of hidden layers in the LSTM network model.

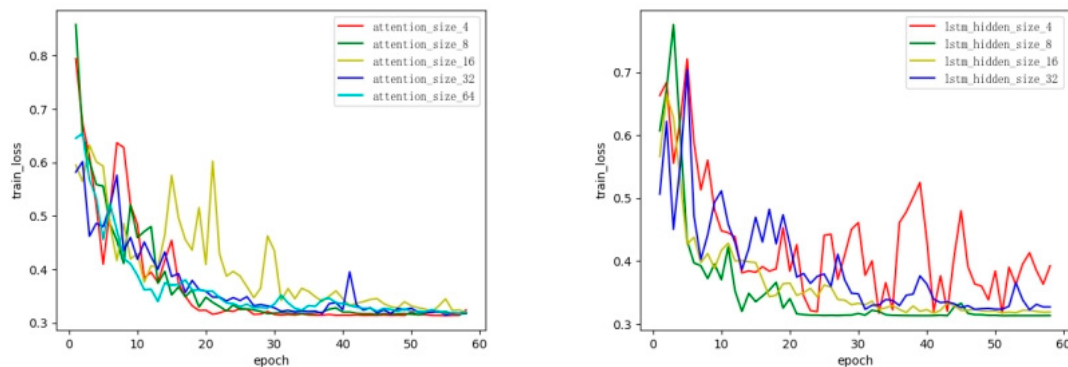
The number of hidden layers	precision	recall	F1-score	Accuracy
4	0.5577	0.611	0.583	0.7
8	0.9375	0.9375	0.9375	0.95
<b>16</b>	<b>0.9375</b>	<b>0.944</b>	<b>0.941</b>	<b>0.975</b>
32	0.8375	0.8	0.818	0.925

**Table 2.** The experimental results of different parameter settings for attention size.

Attention size	precision	recall	F1-score	Accuracy
4	0.7944	0.86	0.826	0.9
8	0.8375	0.9	0.868	0.925
16	0.8819	0.944	0.912	0.95
<b>32</b>	<b>0.9375</b>	<b>0.978</b>	<b>0.941</b>	<b>0.975</b>
64	0.85	0.9	0.874	0.925

It can be seen from Table 2 that it is not more hidden layers, the higher the performance indicators. In the above experiments, the performance indicators are the highest when the number of

hidden layers is set to 16. It is considered that since the limited scale of video data used in the experiments, when the number of hidden layers is small, it is challenging to ensure the accuracy of video feature representation. Therefore, all indicators are significantly lower when the number of hidden layers is 4. If there are many hidden layers, the number of nodes in the RCLA model is large, which is easy to fall into the local optimization. For example, when the number of hidden layers is 32, all indicators are relatively low. Besides, when the parameter of attention size is set to be small, the method proposed in this paper focuses on portraying the local features of a small region in a video keyframe. Since the near-duplicate video data have similar but different visual features in the same local area, amplifying the difference in feature representation will affect the performance of this proposed method. When the attention size is large, there is a confusion problem with near-duplicate video recognition causing incorrect video data cleaning. In addition, as shown in Figure 6, although the number of hidden layers and attention size are set to different, the loss function converges during training, indicating no over-fitting. Finally, according to the experimental results in Tables 2 and 3, the numbers of hidden layers and the attention size in LSTM are set to 16 and 32 in the proposed method.



**Figure 6.** The visualization of loss variation under different parameter settings during training.

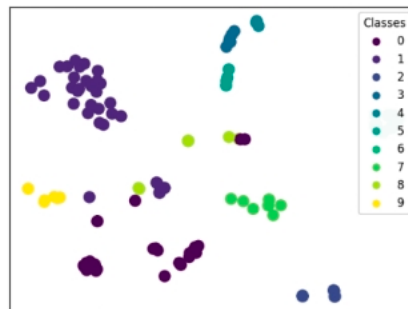
**Table 3.** The experimental results of different parameter settings for attention size.

Models	CC_WEB_VIDEO dataset			
	precision	recall	F1-score	Accuracy
Spatiotemporal Keypoint [3]	0.61	0.96	0.75	0.64
BS-VGG16 [36]	0.79	0.92	0.85	0.85
LBoW [43]	0.63	0.85	0.72	0.66
MLE-MRD [44]	0.82	0.91	0.86	0.87
CBAM-Resnet [42]	0.77	0.92	0.84	0.88
3D-CNN [24]	0.88	0.76	0.84	0.93
RCLA	0.93	0.94	0.94	0.95

In addition, to evaluate the performance of the RCLA model for feature representation of video data, the softmax function is used to achieve the comparison of different feature representation models through the detection of near-duplicate video data. The experimental results are shown in Table 3.

It can be seen from Table 3 that the hand-crafted feature extraction models of spatiotemporal key points and LBoW have limited to represent the video features in the near-duplicate video detection task. Hence, all indicators are low. After introducing the CBAM module in each residual block, the ability of video spatial feature extraction is improved using the channel and spatial attention mechanisms. Therefore, all indicators are improved compared with the hand-crafted feature extraction models. Since the video data has the spatiotemporal feature, not only the LSTM

deep neural network in the RCLA model is used to extract the temporal feature of the video data, but also the standard attention mechanism is used to enhance the feature representation of the local regions of the near-duplicate video data, the near-duplicate video can be accurately identified, and the indicators are generally high, but the recall indicator is lower than that of the spatiotemporal keypoints model. It is considered that the spatiotemporal key points model can extract the spatiotemporal feature of video data, and the number of detection results is enormous, where the number of correct near-duplicate video data is also massive. Hence, the recall indicator of the spatiotemporal key points model is high, but the precision and accuracy indicators are low. On this basis, take the coal mine video dataset as an example, this paper visually shows the clustering results of 10 types of near-duplicate video data, as shown in Figure 7.



**Figure 7.** The clustering results of near-duplicate video data in the coal mine video dataset.

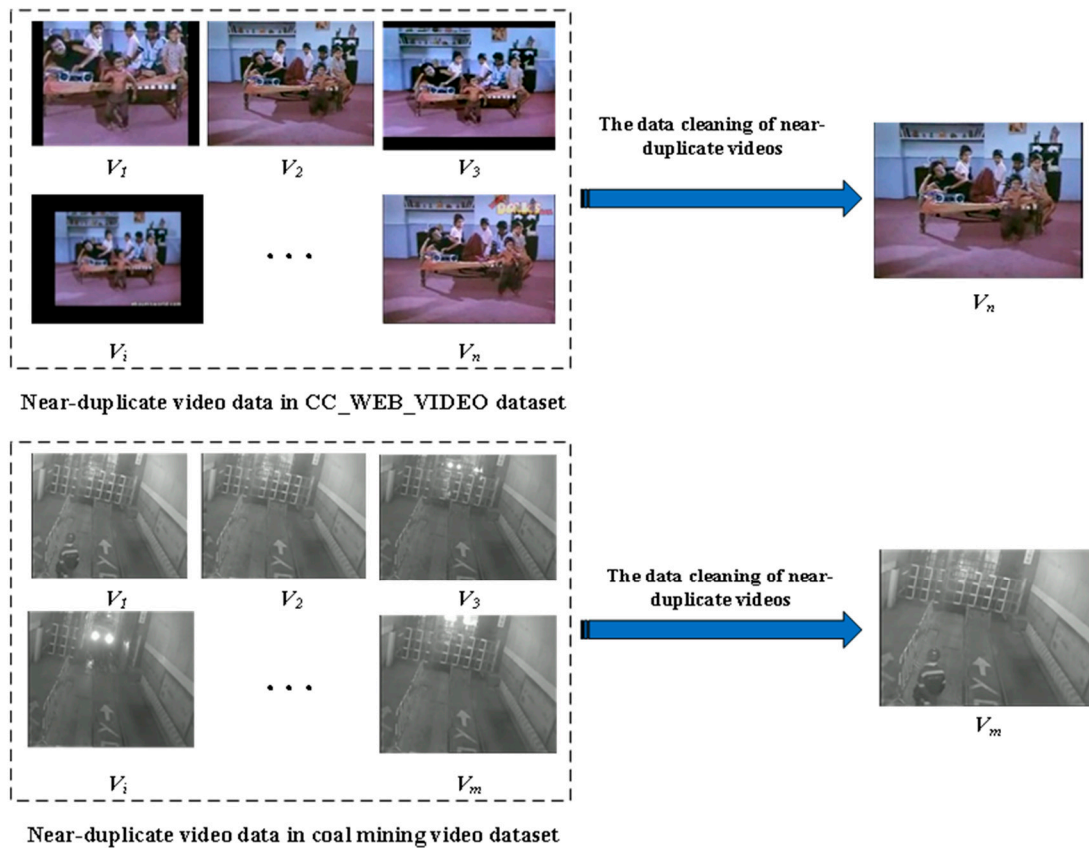
The performance of the proposed method (RCLA-HAOPFDMC) is verified in comparison with the existing research on near-repeat video data cleaning. The experimental results are shown in Table 4.

**Table 4.** The experimental comparison of the different methods.

Methods	cluster cleaning	CC_WEB_VIDEO dataset			The coal mine video dataset		
		acc	rec	F1-score	acc	rec	F1-score
Spatiotemporal Keypoint [3]	K-Means	0.4527	0.451	0.451	0.466	0.5333	0.497
	FD-Means	0.4776	0.538	0.506	0.666	0.5333	0.592
	FD-Means fused with the peak function	0.522	0.835	0.612	0.857	0.6	0.706
LBoW [43]	K-Means	0.453	0.472	0.462	0.5	0.5333	0.516
	FD-Means	0.587	0.615	0.601	0.733	0.733	0.733
	FD-Means fused with the peak function	0.572	0.813	0.632	0.833	0.5	0.625
BS-VGG16 [36]	K-Means	0.275	0.615	0.436	0.465	0.362	0.382
	FD-Means	0.650	0.929	0.76	0.667	0.833	0.74
	FD-Means fused with the peak function	0.49	0.967	0.633	0.5	0.4	0.444
MLE-MRD [44]	K-Means	0.53	0.62	0.57	0.57	0.65	0.61
	FD-Means	0.72	0.79	0.75	0.69	0.76	0.72
	FD-Means fused with the peak function	0.76	0.82	0.79	0.75	0.86	0.80
CBAM-Resnet [42]	K-Means	0.423	0.56	0.481	0.5333	0.666	0.592
	FD-Means	0.587	0.615	0.601	0.733	0.733	0.733

	FD-Means fused with the peak function	0.825	0.681	0.779	0.777	0.7	0.736
3D-CNN [24]	K-Means	0.75	0.75	0.75	0.71	0.91	0.80
	FD-Means	0.80	0.69	0.74	0.87	0.70	0.77
	FD-Means fused with the peak function	0.875	0.7	0.778	0.936	0.723	0.816
	K-Means	0.672	0.67	0.671	0.733	0.733	0.733
RCLA- HAOPFDMC	FD-Means	0.901	0.802	0.848	0.864	0.9333	0.897
	FD-Means fused with the peak function	<b>0.914</b>	<b>0.801</b>	<b>0.854</b>	<b>0.872</b>	<b>0.9333</b>	<b>0.902</b>

This paper compares the proposed method with the existing studies and the different clustering cleaning models, as shown in Table 4. First, it can be seen from Table 4 that the performance indicators of the near-duplicate video cleaning methods based on hand-crafted feature extraction are relatively low, such as spatiotemporal key points and LBoW models, which is due to the limited ability of hand-crafted features to represent video features. Second, the BS-VGG16 model only extracts the spatial features of the video data, and the CBAM-Resnet model introduces the channel and spatial attention mechanisms in the spatial feature extraction. On this basis, the ACNNBN-LSTM model can extract the spatiotemporal features of the video data, and the RCLA-HAOPFDMC method based on the spatiotemporal feature extraction to introduce the standard attention mechanism, which can more accurately depict the features of near-duplicate video data to help to clean the near-duplicate video data accurately and automatically. In addition, by comparing the experimental results of the K-Means, FD-Means, and FD-Means fused with the peak function clustering algorithms, the performance indicators after near-duplicate video cleaning using the K-Means algorithm are low, it is caused by the randomness of the initial clustering center setting. When the FD-Means algorithm is used for near-duplicate video cleaning, the influence of the K value in the K-Means algorithm on the experimental results can be reduced. Thus, the performance indicators are relatively high. This paper constructs a consistent feature hash ring to decrease the impact of data ordering on near-duplicate data cleaning. On this basis, the fusion of the FD-Means algorithm and peak function can further reduce the influence of the random initial cluster center setting on the near-duplicate video cleaning. Therefore, the performance indicators of the proposed method (RCLA-HAOPFDMC) in this paper are higher than the existing methods. Finally, the results of near-duplicate video data cleaning are shown in Figure 8.



**Figure 8.** The results of near-duplicate video data cleaning on the CC\_WEB\_VIDEO and coal mining video datasets.

## 5. Conclusions

In this paper, an automatic near-duplicate video data cleaning method based on a consistent feature hash ring is proposed, which can be utilized to improve data quality for video datasets. In this method, a novel consistent feature hash ring is constructed to alleviate the sensitivity of video data orderliness. On this basis, an optimized feature distance-means clustering algorithm fusing the mountain peak function on a consistent feature hash ring is used to automatically clean the near-duplicate video data. The experiment results on the CC\_WEB\_VIDEO and coal mining video datasets demonstrate the advantages of the proposed method, which can achieve automatic cleaning for near-duplicate video data. However, the method proposed in this paper is not an end-to-end deep neural network model, which needs to be trained separately in the feature extraction and clustering stages. In addition, the computation of cluster cleaning on the consistent feature hash ring is large. In the future, how to construct an end-to-end near-duplicate video data cleaning method will be explored. Moreover, it is of great interest to introduce the swarm intelligence optimization algorithms in the further to improve the accuracy of near-duplicate video data cleaning by optimizing the parameter selection.

**Acknowledgments:** This work was supported in part by the National Natural Science Foundation of China under Grant 61873277, and in part by the Chinese Postdoctoral Science Foundation under Grant 2020M673446.

## References

1. H.-K. Tan et al., "Scalable detection of partial near-duplicate videos by visual-temporal consistency," in Proc. 2009 ACM Int. Conf. Multimedia., Beijing, China, Oct. 2009, pp. 145–154.
2. A. Basharat, Y. Zhai, and M. Shah, "Content-based video matching using spatiotemporal volumes," Comput. Vis. Image Understand., vol. 110, pp. 360–377, Jun. 2008.
3. D.-K. Zhang, Z.-H. Sun, and K.-B. Jia, "Near-duplicate video detection based on temporal and spatial key points," Smart Innov. Syst. Technol., vol. 180, pp.129–137, Apr. 2020.

4. X. Nie et al., "LLE-based video hashing for video identification," in Proc. IEEE 10th Int. Conf. Signal Process., Beijing, China, Oct. 2010, pp. 1837–1840.
5. J.-J. Liu et al., "Near-duplicate video retrieval: Current research and future trends," *ACM Comput. Surv.*, vol. 45, no. 4, pp. 1–23, Aug. 2013.
6. C. Carvalho, R S. Moreira, and J M. Jose, "Data quality visual analysis (DQVA) A tool to process and pin spot raw data irregularities," in Proc. IEEE Annu. Comput. Commun. Workshop Conf., Las Vegas, NV, United States, Jan. 2021, pp. 1036–1045.
7. D. -A. Phalke, and S. Jahirabadkar, "A survey on near-duplicate video retrieval using deep learning techniques and framework," in Proc. IEEE Pune Sect. Int. Conf., Pune, India, Dec. 2020, pp. 124–128.
8. Y.-C. Hu, and X.-B. Lu, "Learning spatial-temporal features for video copy detection by the combination of CNN and RNN," *J. Visual Commun. Image. Represent.*, vol. 55, pp. 21–29, Aug. 2018.
9. S. A. Abdu, A. H. Yousef, and S. Ashraf, "Multimodal video sentiment analysis using deep learning approaches, a survey," *Inf. Fusion*, vol. 76, pp. 204–226, Dec. 2021.
10. S. Mohiuddin, S. Malakar, and R. Sarkar, "Duplicate frame detection in forged videos using sequence matching," in Proc. Commun. Comput. Info. Sci., Santiniketan, India, Jan. 2021, pp. 29–41.
11. L. Shen, R.-C. Hong, Y.-B. Hao, "Advance on large scale near-duplicate video retrieval," *Front. Comput. Sci.*, vol. 14, no. 5, pp. 1–24, Oct. 2020.
12. X.-L. Wang et al., "Research on electricity characteristic recognition method of clean heating based on the big data model," in Proc. Lect. Notes Inst. Comput. Sci. Soc. Informatics Telecommun. Eng., Binzhou, China, Sep. 2020, pp. 25–35.
13. H.-C. Zhou, M.-H. Li, and Z.-Q. Gu, "Knowledge fusion and spatiotemporal data cleaning: A review," in Proc. IEEE Int. Conf. Data Sci. Cyberspace, Hong Kong, China, Jul. 2020, pp. 295–301.
14. Z. Zheng, "Contextual data cleaning with ontology FDs," in Proc. ACM SIGMOD Int. Conf. Manage. Data, Xi'an, China, Jun. 2021, pp. 2911–2913.
15. O. Ye, Z.-L. Li, and Y. Zhang, "Near-duplicate video cleansing method based on locality sensitive hashing and the sorted neighborhood method," in Proc. EAI/Springer Inno. Comm. Comp., Kitakyushu, Japan, Jul. 2019, pp. 129–139.
16. Y.-Y. Zhu et al., "Large-scale video copy retrieval with temporal-concentration SIFT," *Neuro computing*, vol. 187, pp. 83–91, Apr. 2016.
17. C. Henderson, and E. Lzquierdo, "Robust feature matching in long-running poor-quality videos," *IEEE Trans Circuits Syst Video Technol*, vol. 26, no. 6, pp. 1161–1174, Jun. 2016.
18. C.-D. Zhang et al., "Near-duplicate segments-based news web video event mining," *Signal Process*, vol. 120, pp. 26–35, Mar. 2016.
19. Y.-X. Chen et al., "Effective and efficient content redundancy detection of web videos," *IEEE Trans. Big Data*, vol. 7, no. 1, pp. 187–198, Mar. 2021.
20. X.-S. Nie et al., "Joint multi-view hashing for large-scale near-duplicate video retrieval," *IEEE Trans Knowl Data Eng*, vol. 32, no. 10, pp. 1951–1965, Oct. 2020.
21. K.-H. Wang et al., "Attention-based deep metric learning for near-duplicate video retrieval," in Proc. Int. Conf. Pattern Recognit., Milan, Italy, Jan. 2021, pp. 5360–5367.
22. S.-Y. Li et al., "Neighborhood preserving hashing for scalable video retrieval," in Proc. IEEE Int. Conf. Comput. Vision, Seoul, Korea, Oct. 2019, pp. 8211–8220.
23. X.-B. Ai et al., "Inter-frame relationship graph based near-duplicate video clip detection method," in Proc. Commun. Comput. Info. Sci., Beijing, China, Apr. 2019, pp. 70–79.
24. H. Chen et al., "A supervised video hashing method based on a deep 3d convolutional neural network for large-scale video retrieval," *Sensors*, vol. 21, no. 9, pp. 3094–, 2021.
25. Y.-Z. Hu, Z.-K. Mu, and X.-B. Ai, "STRNN: End-to-end deep learning framework for video partial copy detection," in Proc. J. Phys. Conf. Ser., Xi'an, China, Jul. 2019, pp. 1–10.
26. Z. H. Mohamed, and J. S. Vinila, "A comparative study on data cleaning approaches in sentiment analysis," in Proc. Lect. Notes Electr. Eng., Thiruvananthapuram, India, Dec. 2019, pp. 421–431.
27. H. -C. Zhou, M. Li, and Z.- Q. Gu, "Knowledge fusion and spatiotemporal data cleaning: A review," in Proc. IEEE Int. Conf. Data Sci. Cyberspace, Hong Kong, China, Jul. 2020, pp. 295–301.
28. S. -X. Song et al., "SCREEN: Stream data cleaning under speed constraints," in Proc. ACM SIGMOD Int. Conf. Manage. Data, Melbourne, VIC, Australia, Mar. 2015, pp. 827–841.
29. Y. C. Tian, P. Michiardi, and M. Vukolic, "Bleach: A distributed stream data cleaning system," in Proc. IEEE Int. Congr. Big Data, Honolulu, HI, United States, Jun. 2017, pp. 113–120.
30. Q. -M. Liu et al., "Cleaning RFID data streams based on l-means clustering method," *J. China Univ. Post Telecom.*, vol. 27, no. 2, pp. 72–81, Apr. 2020.
31. K G. Ranjan, B. R. Prusty, and D. Jena, "Comparison of two data cleaning methods as applied to volatile time-series," in Proc. Int. Conf. Power Electron. Appl. Technol., Surathkal, India, Oct. 2019, pp. 1–6.
32. L. Zhu et al., "Keyword search with real-time entity resolution in relational databases," in Proc. ACM. Int. Conf. Proc. Ser., Macau, China, Feb. 2018, pp. 134–139.

33. A. Aissani, E T J. Yi, and T. Thamilyanan, "End to end real-time data cleansing, standardization and reliability towards a smooth digital oil field deployment," in Proc. Offshore Technol. Conf. Asia, Kuala Lumpur, Malaysia, Nov. 2020, pp. 1–10.
34. Z. Zheng, M. Milani, and F. Chiang, "CurrentClean: Spatio-temporal cleaning of stale data," in Proc. Int. Conf. Data Eng., Macau, China, Apr. 2019, pp. 172–183.
35. Y. Zhang et al., "ImageDC: Image data cleaning framework based on deep learning," in Proc. IEEE Int. Conf. Artif. Intell. Inf. Syst., Dalian, China, Mar. 2020, pp. 748–752.
36. Y. Fu, Z. Han, O. Ye, "FD-means clustering cleaning algorithm for near-duplicate videos," *Comput. Eng. and Appl.*, vol. 58, no. 1, pp. 197–203, Jan. 2022.
37. W. Jia et al., "Scalable hash from triplet loss feature aggregation for video de-duplication," *J. Visual Commun. Image Represent.*, vol. 72, pp.1–9, Oct. 2020.
38. H. Chen et al., "The MSR-video to text dataset with clean annotations," *Comput. Vis. Image Understand.*, vol. 225, pp. 1–7, Dec. 2022.
39. H.-Y. Cheng, C.-C. Yu, "Automatic data cleaning system for large-scale location image database using a multilevel extractor and multiresolution dissimilarity calculation," *IEEE Intell. Syst.*, vol. 36, no. 5, pp. 49–56, May. 2021.
40. O. Ye, J. Deng, Z H. Yu, T. Liu and L H. Dong, "Abnormal event detection via feature expectation subgraph calibrating classification in video surveillance scenes," [J]. *IEEE ACCESS*, 2020, 8: 97564–97575.
41. Woo S. et al., "CBAM: Convolution block attention module," in Proc. 15th Euro. Conf. Comput. Vision, Munich, Germany, Sep. 2018, pp. 3–19.
42. Wang Z, "Research on garbage image classification based on neural network," in Proc. the 2022 International Conference on Computer Network, Electronic and Automation (ICCNEA), 2022, Xi'an, China, pp. 214–217.
43. Giorgos K et al., "Near-Duplicate Video Retrieval by Aggregating Intermediate CNN Layers," in Proc. the 23rd International Conference on MultiMedia Modeling, 2017, Reykjavik, Iceland, pp. 251–263.
44. H. -Y. Cheng, C. -C. Yu. Automatic Data Cleaning System for Large-Scale Location Image Databases Using a Multilevel Extractor and Multiresolution Dissimilarity Calculation[J]. *IEEE Intelligent Systems*, 2021, 36(5): 49-56.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.