

# **MINING VIDEO DATA ON THE CLOUD COMPUTING PLATFORM FOR HIGHER EFFICIENCY**

*Fen Chen, Department of Information Management, Nanjing University of  
Science and Technology, No. 200 Xiao Ling Wei Street, Nanjing, China  
210094, 86-15190491057, pku.lanyan@gmail.com*

*Yue Peng, Department of Information Management, Nanjing University of  
Science and Technology, No. 200 Xiao Ling Wei Street, Nanjing, China  
210094, 86-13851931851, pengyue0321@126.com*

*Duanwu Yan, Department of Information Management, Nanjing University of  
Science and Technology, No. 200 Xiao Ling Wei Street, Nanjing, China  
210094, 86-13951851356, yanwu\_nju@163.com*

## **ABSTRACT**

Rapid growth of video resources makes video mining being one of the most popular research topics in the world. Mass video information has become a big challenge for current information processing technique, especially the efficiency issue. In recent several years, cloud computing has attracted much researchers' attention; its advantage of large scale, high flexibility and low cost provides a potential efficient way for video organization. In this paper, the authors will conduct an experiment of video classification based on the cloud computing platform which focuses on the improvement of automatic cartoon detection, and analyze the running efficiency between the local computer and remote servers.

## **INTRODUCTION**

The rapid developing IT technique leads to the fast growth of video resources. Video provides rich and intuitive information, forming one of the most effective communication ways in our lives. Videos are significantly different from the other resources, such as huge amount of data, complicated structure. The ever-growing video data and its characteristics make the organization and effective retrieval of videos being of great importance. At present, video mining has been one of the hottest topics in the world.

In different video applications and various video categories, cartoon is usually involved in most of them, such as Baidu, Google Video, etc. Thus, cartoon detecting is an important issue in online video organization or the other related field, like digital library. In this paper, the authors will focus on the cartoon detection.

In video mining, the processing efficiency has been a big issue that puzzles the researchers for a long time. In the past several years, the cloud computing technique has been paid a lot of attention. This technique uses large amounts of PCs to substitute the expensive servers used before, thus reducing the operation costs for users. What's more important is that the large scale gives the integrated system

super-computing capability. If researchers can use this technique to organize the video resources, then the advantage of parallel and distributed computing can help improve the video processing efficiency.

In this paper, the authors will try mining features extracted from cartoon videos and analyze the efficiency difference between local computer and remote servers deployed on the clouding platform.

The rest of this paper is organized as follows. Section 2 presents the literature review. The experiment in details is to be introduced in Section 3. In the 4th section, the future work will be explained.

## **LITERATURE REVIEW**

According to the literature review, the authors find that little work has been done in the research field of video mining on the cloud computing platform. In this section, the respective research work in related fields will be introduced.

### **Cartoon Detection**

By far, there is not very much work conducted on classification of cartoons, and cartoon detection is still at its preliminary stage [1].

Rama[2] presented a flexible scheme based on a non-linear classifier called Fuzzy Integral, using a Choquet integral for the data fusion process and extracted the following visual and audio features: pitch, roll-off, frequency-centroid mean, etc. Then, results using this operator for cartoon detection are presented and compared with other well known statistical classification methods like K-NN.

Ianeva [3] introduced new metrics based on the pattern spectrum of parabolic size distributions derived from parabolic granulometries and the complexity of the image signal approximated by its compression ratio. They also computed average color saturation. They also evaluated the effectiveness of proposed features using SVM on a set of keyframes from the TREC-2002 video track and a set of web images.

Gao [1] extracted groups of features including MPEG-7 descriptors to describe the videos, such as HSV color histogram, ratio of color area to the whole image, color structure histogram, edge direction histogram, etc. A content-based video classifier was designed by introducing the active relevance feedback technique into SVM for the cartoon video detection.

Zhao [4] used the following features: (1) color: RGB color histogram, and average of these three color channels, etc.; (2) texture: Gabor wavelet transformation; (3) caption area: detecting if the caption area existed; (4) geometrical size: image width, image height, etc.; (5) inter-frame difference: the average of difference of continuous frames. The following classification stage used SVM as the classifier.

### **Data Mining on the Cloud Computing Platform**

Cloud computing is the development of virtualization, parallel, and grid computing techniques.

Grossman [5] described the Sector storage cloud in the Sphere computing platform. Sector and Sphere used clustering method, connecting with the WAN to analyze large data sets. The author used distributed data mining algorithm and compared it with Hadoop.

Ji [6] discussed data mining based on GAE cloud platform. The author designed the data constraint function to solve heterogeneous and non-existing data access problem. Then, layered thinking model is used and a new multilayered plug-in architecture is induced to enhance the scalability of the system.

Wang [7] introduced the Hadoop integrated framework and the Sprint classification algorithm; then, the authors described in detail about the Sprint method and its flow of execution on MapReduce model, and used the analyzed decision tree model for data classification.

Sun [8] analyzed in-depth the key techniques on GAE platform, and compared it with traditional IT system. The authors found that GAE has three basic characteristics: integration of large amounts of cheap servers; coordination of infrastructure and applications to achieve best use of hardware; fault tolerance of nodes by software. All these characteristics are in sharp contrast with traditional IT systems based on high performance Unix server clusters.

As mentioned above, though data mining on the cloud computing platform and video information mining are both hot topics in the world, little work has been conducted in the combined research field. In the following section, the authors will combine these two techniques together, do an experiment of cartoon detection on the cloud computing platform, and verify the efficiency issue.

## **EXPERIMENT AND THE RESULT ANALYSIS**

### **Video Feature Extraction**

Cartoons have their own characteristics; for example, roles in cartoons usually have very simple shapes; cartoons generally have less, simpler and richer color, texture is often simplified to a single color [3]. The features used in this paper are extracted using the similar algorithm the authors have proposed in paper [9]. The following shows the process.

#### **(1) Preprocessing**

Preprocessing contains two steps: shot boundary detection (SBD), and keyframe extraction. In our experiment, SBD is based on color histogram. When difference of two continuous frames' color histograms exceeds the threshold, shot boundary will be considered existing. Keyframe extraction in our test is an automatic procedure, called color histogram average.

#### **(2) Region Segmentation**

Morphologic watershed algorithm is adopted to segment the frames into different logical parts, and every pixel is designated an added value of the region number for further use, thus forming the following matrix:

$$\begin{bmatrix} RN_{00}, RN_{01}, RN_{02}, \dots, RN_{0j}, \dots, RN_{0n} \\ \dots \\ RN_{i0}, RN_{i1}, RN_{i2}, \dots, RN_{ij}, \dots, RN_{in} \\ \dots \\ RN_{m0}, RN_{m1}, RN_{m2}, \dots, RN_{mj}, \dots, RN_{mn} \end{bmatrix}$$

Figure 1: Region Number Matrix

Here, every element in the matrix is a pixel, and pixels of the same value belong to the same region. The largest number indicates the region number after segmentation.

### (3) Visual Features

The following two visual features are extracted after the above region number computing: one is *average color component*, and the other is *average color number* based on different regions.

*Average color component (ACC)* is related to three components computing, namely H, S, and V, where H means hue, S indicates saturation, and V is brightness. Compared to RGB color space, HSV is closer to the way that people feel about color; thus in this paper, the authors apply this color space.

Take *H component* as the example to interpret the extracting procedure. The steps to extract H component include: Firstly, count the sum of hue of every region; secondly, count the average hue value of each region; then, the sum hue data of all regions is calculated; and at last, average hue value of all regions is the result the authors hope to get. In this procedure, the above region number matrix in figure 1 is needed to determine the region id of a pixel. *S* and *V* components are counted similarly. Figure 2 shows the *H component's* value to distinguish the cartoon and non-cartoon videos.

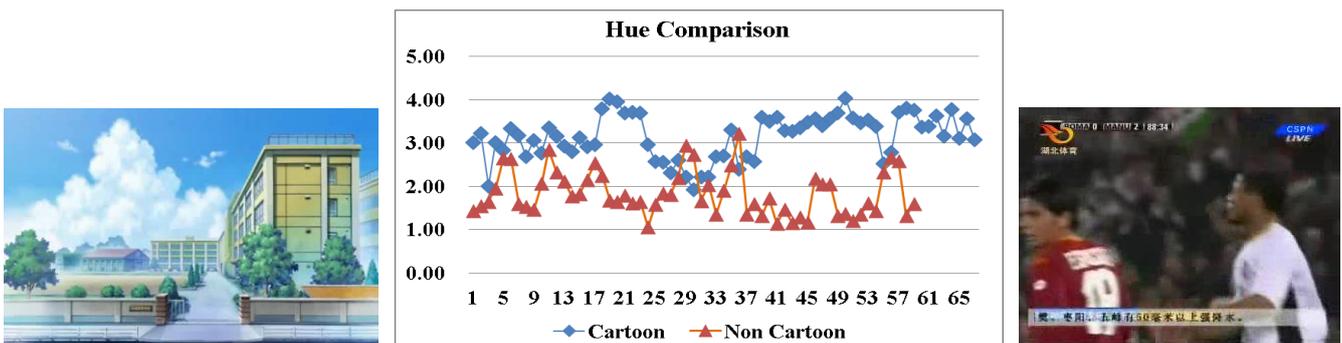


Figure 2: Hue Comparison of Cartoon and Non-cartoon

*Average color number (ACN)* also comes from the results of region segmentation. Figure 3 shows the extraction process: Firstly count the average region color kinds (color types divided by pixel number of a region), then count the average color number of all regions in a video frame. The idea behind the algorithm is that cartoons usually have less color number than the ordinary videos.

Besides the above visual features, some traditional features are also been extracted to compare the

classification accuracy, including color histogram (72 Dimension), and texture based on gray-level co-occurrence matrix (GLCCM, 8 Dimension).

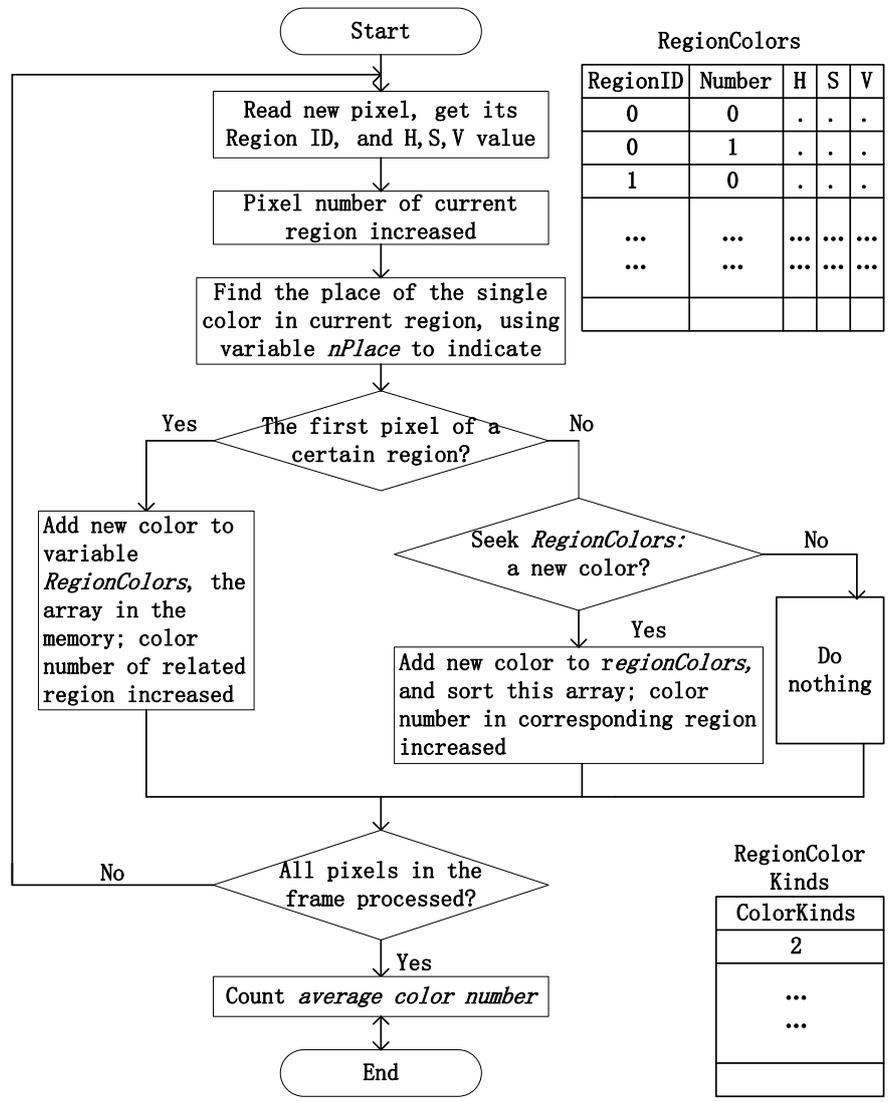


Figure 3: Average Color Number Counting

## Cloud Platform Selection

At present, several clouding platforms have been developed in the world, such as Google App Engine (GAE), Amazon Elastic Compute Cloud, Sun Hydrizine, IBM Blue Cloud, etc. In this paper, the authors will select GAE as the experiment platform, based on the following reasons:

- GAE is a new generation of cloud-based web development platform, and programs built on this platform are easy to build, maintain, and extend. In addition, Google cloud computing platform provides rich API functions
- From the other papers, GAE is one of the most often used platforms in research. For example, Hussain [10] proposed and demonstrated the usage of GAE cloud computing platforms as a possible solution for mining

and analyzing large amounts of data from Twitter; Ji [6] explored the methods of data mining based on Google App Engine as talked above; Huang [11] analyzed the three core technologies of Google Cloud, GFS, Map-Reduce and BigTable, and deployed a video surveillance system on the GAE platform.

- The authors once conducted preliminary comparative analysis research on GAE platform about simple application [12]. The current research result indicates that whether from stability, or running efficiency, GAE can get better effect than local system. This is another important reason for the authors to select GAE.

## Experiment Configuration

### (1) Data Collection

The data collection in this experiment is shown in table 1. Among the videos, 51 videos are selected randomly as the training set, and the extra serves as the test set.

Table 1: Data Collection in the Experiment

Video Type	Number	Comment
Cartoon	94	American cartoons; Chinese cartoons; Japanese Hayao Miyazaki animation; etc.
Advertising	6	Various kinds of advertising
Sport	30	All kinds of sports: Soccer, basketball, bicycle, Wrestling, etc.
Entertainment	43	Teleplay, movie, entertainment
Documentary	26	Various documentary videos
News	41	News videos from multiple channels

### (2) Evaluation Measure

The following formula is used to evaluate the classification accuracy, as follows:

$$Accuracy = \frac{\sum_{i=1}^k C_i}{n} \quad (1)$$

This formula is used to evaluate the percent of videos in the test set that have been classified into the right class. Here, k is the class number, Ci is the number of videos that overlap between our result and the Gold Standard, and n is the video number in test set. The higher the value is, the better the effect.

### (3) Classifier

Researchers have developed many kinds of classifiers, such as Decision Tree, Naïve Bayes, and SVM, etc. All these algorithms have been used in video classification for automatic video processing. In this paper, Decision Tree and SVM will both be used to check the mining effect on different platforms.

## Results and Analysis

In this section, the authors will demonstrate the mining result using different classifiers on different platforms. Table 2 shows the experiment's overall cases.

Table 2: Overall Cases in the Experiment

	Decision Tree		SVM	
	Color Histogram + GLCCM (80-D)	ACC+ ACN+ GLCCM (12-D)	Color Histogram+ GLCCM (80-D)	ACC+ ACN+ GLCCM (12-D)
GAE	10 times	10 times	10 times	10 times
Local	10 times	10 times	10 times	10 times

Table 3 shows the classification of Decision Tree and SVM; the results are the same whether on the local system or on the GAE platform, indicating the stability of these algorithms.

Table 3: Mining Accuracy of Different Algorithms

	Decision Tree		SVM	
	Color Histogram + GLCCM (80-D)	ACC+ ACN+ GLCCM (12-D)	Color Histogram+ GLCCM (80-D)	ACC+ ACN+ GLCCM (12-D)
GAE	68.3%	72%	71.4%	80.4%
Local				

Table IV and V show the computing efficiency using Decision Tree and SVM on local and GAE platforms.

Table 4: Decision Tree – Processing Time - Feature used: ACC + ACN + GLCCM (12-Dimension)

DT	1	2	3	4	5	6	7	8	9	10
GAE	1383	1387	1316	1329	1357	1808	1347	1380	1383	1372
Local	14102	14492	14196	13846	14393	14924	14025	14988	14130	15522

Table 5: SVM – Processing Time - Feature used: ACC + ACN + GLCCM (12-Dimension)

SVM	1	2	3	4	5	6	7	8	9	10
GAE	6620	6682	6892	6506	6890	6682	6776	6790	6880	6890
Local	38444	47520	41098	40036	37685	41098	50703	47400	47345	40764

From table 3 to table 5, we can draw the following conclusions:

- The features of *average color component* and *average color number* can get higher accuracy, and they are suited to be introduced into cartoon video processing.
- When using Decision Tree to detect cartoons, the processing efficiency in GAE platform is as 6 to 10 times high as in the local system. The same thing happens to SVM algorithm, and the processing time on GAE platform is about one-sixth of that in the local computer.
- The processing efficiency of Decision Tree is higher than that of SVM in the same platform. But taking the mining accuracy into consideration, SVM is better to use in video mining. What's more, the efficiency of SVM on GAE platform is about one half comparing to that of Decision Tree in the local system. Thus SVM can be a relatively good choice in video mining in the cloud computing environment.

## CONCLUSION AND FUTURE WORK

In this paper, the authors focus on video feature extraction and the mining result comparison on different platforms. The experiment shows that cloud computing can be a good choice to organize the vast amount of video data to achieve better processing efficiency. In the future, the authors will try to optimize the parameters of mining algorithms and improve the accuracy and efficiency further.

## ACKNOWLEDGEMENT

This paper is supported by Ministry of Education, Humanities and Social Sciences project - Video Information Organization and Mining Research based on the Clouding Platform (Project Number: 10YJC870001), Chinese Postdoc Fund – Video Information Classification based on Multimodal Framework (Project Number: 2012M521061), and the Independent Research Program (Project Number: 2011YBXM94) of Nanjing University of Science and Technology.

## REFERENCES

- [1]Gao Xinbo, Tian Chunna, Zhang Na. A Cartoon Video Detection Method Based on Active SVM Learning. *Journal of Electronics & Information Technology*, 2007, vol. 29, issue.6, pp. 1338-1342
- [2]Rama A., Tarres F., Sanchez L. Cartoon Detection Using Integral. *Eight International Workshop on Image Analysis for Multimedia Interactive Services (WIAMIS'07)*. 2007, pp. 45-45
- [3]Ianeva T. I., Vries A. P. D., Rohrig H. Detecting Cartoons: A Case Study in Automatic Video-Genre Classification. *ICME*. 2003: 449-452
- [4]Zhao Ming. *Content-based Classification of Short Animation*. Shanghai Jiao Tong University, Thesis, 2007, pp.27-32
- [5]Robert L Grossman, Yunhong Gu. *Data Mining Using High Performance Data Clouds: Experimental Studies Using Sector and Sphere*. [J/OL]
- [6]Ji Jun. *The Framework and Implementation of Data Mining based on Cloud Computing*. Shandong: Thesis of Qingdao University, 2009
- [7]Wang E, Li Ming. *Research on Mass Data Mining under Cloud Computing*. *Modern Computer*, 2009(11): 22-25, 50
- [8]Sun Jian, Jia Xiaojing. *Framework of Google Cloud Computing Platform and its Impact on Cost*. *Telecommunication Science*, 2010(1): 38-44
- [9]Chen Fen, Lai Maosheng, Ye Zhining. *Detecting Cartoons: Automatic Video Genre Classification*. *IEEE International Conference on Management and Service Science*, 2010
- [10]Rashid Hussain, Adnan; Hameed, Mohd Abdul; Hegde, Nagaratna P. *Mining Twitter using cloud computing*. *World Congress on Information and Communication Technologies (WICT)*, 2011: 187-190
- [11]Zhenyu Huang. *Research And Implementation Of Cloud Computing Based Network Video Surveillance System*. Shanghai Jiaotong University Thesis, 2011
- [12] Fen Chen. *Performance Analysis of the Cloud Platform from a Micro Perspective*. *IEEE International Conference on Management and Service Science*, 2012