

Privacy Protection in Interactive Content Based Image Retrieval with Copy deterrence

Suhaibathul Aslamiya M P

Department of Computer Science and Engineering, AWH Engineering College kuttikkatoor, Calicut, Kerala

Abstract - Content-based image retrieval (CBIR) is one of the fundamental image retrieval primitives. Privacy protection in Content Based Image Retrieval is a new research topic in cyber security and privacy. The state-of-art CBIR systems usually adopt interactive mechanism, namely relevance feedback, to enhance the retrieval precision. How to protect the user's privacy in such Relevance Feedback based CBIR is a challenge problem. The new CBIR system consists of three stages: 1) private query; 2) private feedback; 3) local retrieval. Private query performs the initial query with a privacy controllable feature vector; private feedback constructs the feedback image set by introducing confusing classes; local retrieval finally re-ranks the images in the user side. In addition, considering the case that the authorized query users may illegally copy and distribute the retrieved images to someone unauthorized. In this paper propose a watermark-based protocol to deter such illegal distributions. In this watermark-based protocol, a unique watermark is directly embedded into the encrypted images before images are sent to the query user. Hence, when image copy is found, the unlawful query user who distributed the image can be traced by the watermark extraction.

Key Words: CBIR, Image privacy, Copy deterrence, Watermark, Relevance feedback etc.

1. INTRODUCTION

With the rapid development of Multimedia and Internet, massive images are generated and distributed, how to store and share such large amount of data efficiently becomes an important issues. Currently, two image retrieval approaches dominate: text based and content-based. In text-based image retrieval, images are searched based on the text descriptions associated with the images. In content-base image retrieval (CBIR), images are retrieved according to their visual similarities measured on low level visual features. CBIR still works when textual annotations are not available, and has been implemented in the state-of-art image search engines, such as Google Image Search1, and Bing Image Search2.

Recently, with the emergence of new applications, an issue with content-based search has arisen – sometimes the query or the database contains privacy-sensitive information. In a networked environment, the roles of the database owner, the database user, and the database service provider can be taken by different parties, who do not necessarily trust each other. A privacy issue arises when an untrusted party wants to access the private information of another party. In that case, measures should be taken to protect the corresponding information. The main challenge is that the search has to be performed without revealing the original query or the database. This motivates the need for privacy-preserving CBIR (PCBIR) systems.

Despite the tremendous benefits, image privacy becomes the main concern with CBIR outsourcing. For example, patients may not want to disclose their medical images to any others except to a specific doctor in medical CBIR applications. To formulate the problem, this paper considers two types of privacy threats. Firstly, the user's search intention could be learned by the service provider.. Secondly, after receiving the retrieved images, the query user may illegally distribute these images to someone unauthorized for benefits.

Several partial-encryption based commutative encryption and watermarking (CEW) methods have been proposed. In these methods, the image data is divided into two parts. One part is encrypted to protect the image content, and the other is used to carry the watermark. The encryption and watermarking operations in these methods do not interfere with each other, which is suitable for our application scenario. However, the watermarked part has not been protected well and will leak information about the images.

The major contributions of this work are summarized as follows.

- A new Private Relevance Feedback based CBIR scheme (PRF-CBIR), which consists of three stages: private query, private feedback and local retrieval.
- A new private query method, which performs the initial query with a privacy controllable feature vector.
- A new private feedback method, which introduces confusing classes into the feedback image set, to protect user intention.
- The existing searchable encryption schemes usually consider that the query users are fully trustworthy. This is not necessarily true in real-world applications. Considering the dishonest query users who may distribute the retrieved images to those who are unauthorized. A watermark-based protocol is designed for the copy deterrence. Specifically, after completing the search operation requested by an image user, a unique watermark associated with the image user is imperceptibly embedded into the retrieved images. Then, the watermarked images are sent to the image user. When an illegal copy of the image is found, the unlawful query user who made the illegal distribution can be traced by the watermark extraction. This will help to deter the illegal distribution.
- A watermark-based protocol in the encryption domain is designed for copy-deterrence. Different from common watermarking techniques, the proposed protocol needs to embed the watermark directly into the encrypted images. After receiving the encrypted and watermarked images, the query user needs to decrypt the images directly. And the decryption should not affect the watermark in the images.

2. LITERATURE SURVAY

Service provider owns an enormously huge number of images and is willing to provide the CBIR service on the dataset. Service provider can provide CBIR service via its own servers or outsourcing to cloud provider. Facing an untrustworthy service provider, the problem is how to prevent it from inferring the user's search intention. A few works have been reported on image privacy issues in CBIR.

In [1], Lu et al. randomized the search index using secure inverted index and secure min-Hash with approximate distance preserving property. Yuan et al. [2] proposed a lightweight secure image search scheme, namely SEISA. In SEISA, matrix based encryption is used to encrypt image features, which can achieve high search efficiency and accuracy. Xia et al. [3] presented a privacy-preserving CBIR approach which supports local-feature based CBIR with the earth movers distance (EMD) as similarity metric. In their approach, the EMD problem is transformed to a linear programming problem. The cloud provider then solves the linear programming problem without learning the sensitive information. Most recently, Weng et al. [4] studied this problem in the scenario of the near duplicate detection based on robust hashing and piece-wise inverted indexing. The basic idea is to randomly omit certain bits in the query to prevent the service provider from predicting the user's search intention accurately. However, this randomly query feature omitting approach is appropriate for near duplicate detection, but not suitable for general CBIR task. Majhi et al. [5] used the texture based feature along with the color histogram feature to exploit both color and texture properties of images which are effectively participating in the image retrieval. Then compare color histogram of two images for visual similarity. Chun et al. [6] proposed two texture based operator known as BDIP and BVLC technique to measure smoothness and extract sketch features respectively. This texture and quantized HSV color space histogram features are exploited to formulate the feature vector which is encrypted by performing XORed operation with its sliced biplanes by a random binary bit pattern to preserve the hamming distance. Yanyan Xu et al. [7]'s work is based on orthogonal decomposition. Orthogonal decomposition is a kind of vector representation method, any vector can be expressed as a sum of a set of component coefficients through orthogonal decomposition. The image is divided into two different components, for which encryption and feature extraction are executed separately. As a result, service provider can extract features from an encrypted image directly and compare them with the features of the queried images, so that users can thus obtain the image. There is no special requirements to encryption algorithms, which makes it more universal and can be applied in different scenarios.

None of these schemes consider the dishonest query users who may illegally distribute the retrieved images. Actually, it is difficult to design a method to completely prevent illegal distributions. However, it is possible to design certain techniques to deter such illegal behaviors. Watermarking techniques have been widely studied for the copy deterrence.

3. PRIVACY ISSUES

Privacy-preserving CBIR (PCBIR) systems face the common problem that the server is not trusted by the database owner or the user. According to where privacy is emphasized, a PCBIR technique is typically used in the following scenarios:

- The database contains private information, e.g., photo sharing websites;
- The query contains private information, e.g., remote Diagnosis.
- The CBIR technique contains private information, e.g. proprietary technology.

The major privacy concern is that the user’s search intention could be learned by the service provider. The service provider can learn the search intention through the following attacks:

- *Query Attack:* The service provider acquires the user’s search intention directly with the query image Q . For example, if the user queries with a tumor image, the service provider can easily infer the user’s search intention is tumor.
- *Result Attack:* By analysing the results returned to user, R_0 and R_1 , the service provider may be very likely to infer the search intention. For instance, if the majority of the result are tumor images, it is likely that the user’s search intention is tumor.
- *Feedback Attack:* The user’s feedback image set F , also can be utilized by the service provider to learn the search intention. Suppose the user labels some tumor images as relevant, the service provider can acquire that the search intention is tumor.

4. PROPOSED SYSTEM

This work proposed a new privacy preserving CBIR system. Compared with existing solutions, proposed model is more attractive. This system is a retrieval system with a given image as a query image; the system returns relevant images from the database and preserve the users search intention at the same time. Also consider the dishonest image users who correctly follow the protocol specification, but may distribute the retrieved images to the unauthorized others for benefits. The watermarking technology is adopted to deter the illegal distribution. Fig.1 shows the system model of new CBIR.

It has three stages — private query, private feedback and local retrieval.

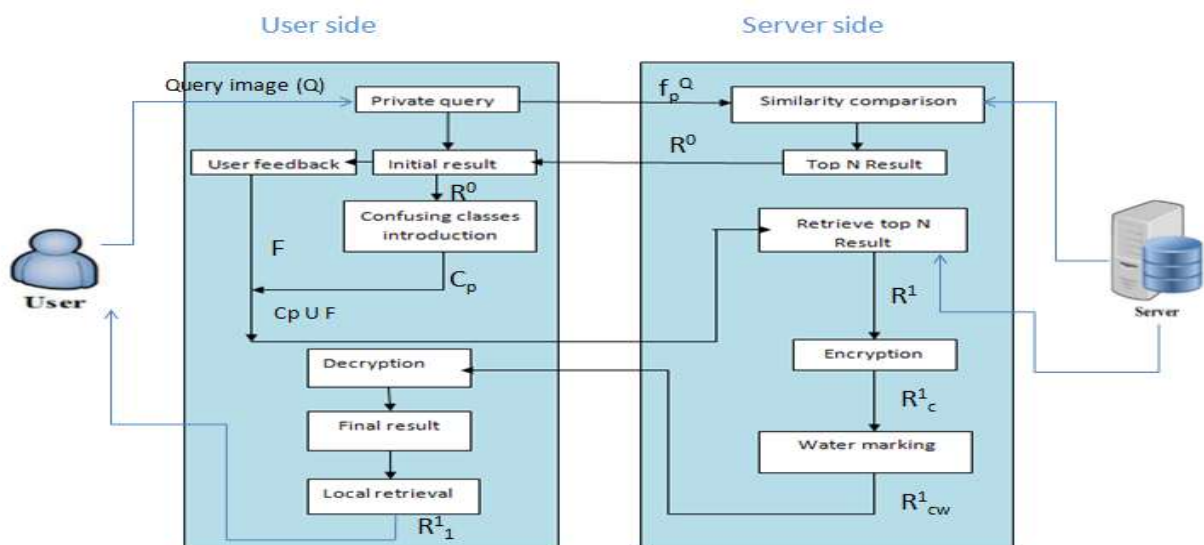


Fig-1: System Architecture

4.1 Private Query

Private query is proposed to address the query attack and the result attack on R^0 . The basic idea is, instead of query image Q , we use a part of the Q 's feature f_p^Q as query. In this way, the query attack can be avoided, and the result attack on R^0 can be alleviated by adjusting the privacy information contained in f_p^Q . Private query is developed by utilizing the Percentage of Variance (PoV) [8] defined in the Principal Component Analysis (PCA) [8] reduced feature space. In this work, PCA is performed prior to retrieval for two goals. Firstly, PCA is adopted for dimension reduction [9] to avoid the curse of dimensionality. PCA can perform dimensionality reduction while preserving as much of the variance in the high dimensional space as possible. Secondly, we make use of PoV, defined in the PCA reduced feature space, to measure privacy and generate query vector f_p^Q . This is based on the observation that PoV can characterize the importance of different feature components. Fig.2 shows the flow chart of the proposed private query method, which consists of offline stage and online stage.

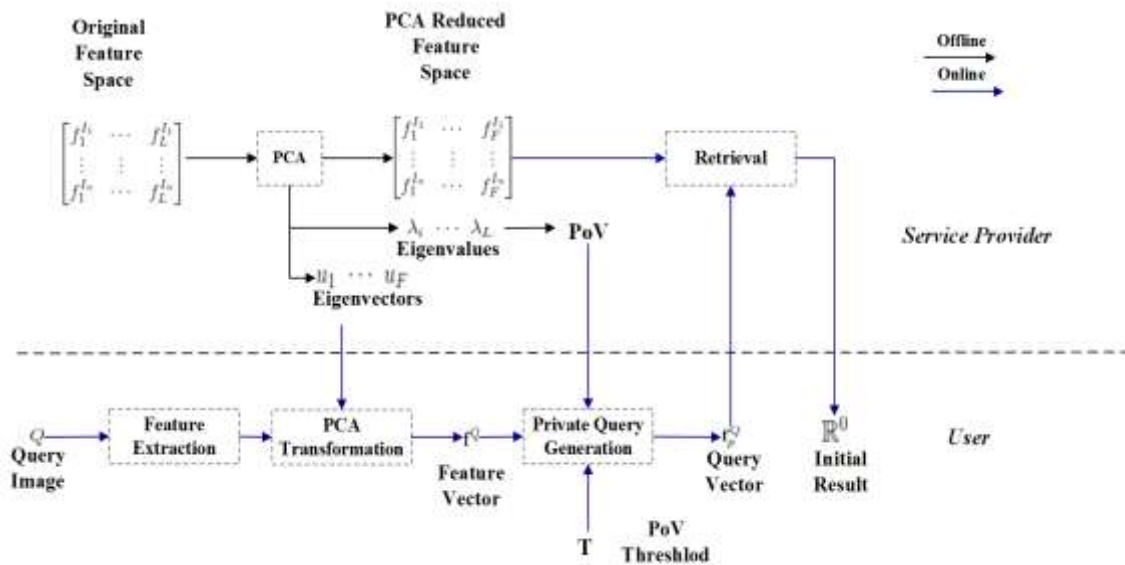


Fig-2: The flow chart of private query

4.1.1 Offline Stage

In this stage, the service provider performs PCA on the image dataset, and computes the PoV values of PCA feature components. Suppose the image dataset is $\Omega = \{I_1, \dots, I_i, \dots, I_n\}$, the feature vector of image $I_i \in \Omega$ can be denoted as $f^{i_i} = (f^{i_i}_1, \dots, f^{i_i}_j, \dots, f^{i_i}_F)^T$. And the feature vectors of all images in Ω can be represented as a matrix $F = (f^1, \dots, f^i, \dots, f^n)^T$. PCA is then employed to learn a linear transformation as,

$$F_L = F(u_1, \dots, u_i, \dots, u_F)$$

where F_L keeps only the first L ($L < F$) important feature components. Those feature components are sorted in the descending order of importance. The learned weight vectors $u_1, \dots, u_i, \dots, u_F$ have two properties. Firstly, $u_1, \dots, u_i, \dots, u_F$ form an orthogonal basis for the L feature components. Secondly, they can preserve as much variability as possible in the original data. $u_1, \dots, u_i, \dots, u_F$ are also the top L eigenvectors of F 's covariance matrix, and their eigenvalues $\lambda_1, \dots, \lambda_i, \dots, \lambda_F$ characterize the variance degree explained by the corresponding eigenvectors respectively.

For each PCA feature component f_i , we convert the eigenvalues to PoV as,

$$PoV(f_i) = \frac{\lambda_i}{\sum_{j=1}^L \lambda_j}$$

4.1.2 Online Stage

In this stage, the user generates the query vector f^Q_p , and retrieves images from the service provider. Suppose that the user's query image is Q , after feature extraction and PCA transformation, the PCA features of Q can be represented as f^Q . The PoV threshold T and the query vector length l are set according to the user's privacy policy. According to above equation, the PoV of f^Q_p can be calculated as,

$$PoV(f^Q_p) = \sum_{f_i \in f^Q_p} PoV(f_i).$$

Given PoV threshold T and l , the goal of private query generation is to choose a continuous segment from f^Q satisfying T , which can be formulates as,

$$\text{Min } ||T - PoV(f^Q_p)||,$$

$$\text{s.t. } f^Q_p \in \{f^Q[i : (i + l - 1)] \mid 1 \leq i \leq (L + 1 - l)\},$$

where $f^Q[i : (i + l - 1)]$ denotes a continuous segment from f^Q , with the index from i to $(i + l - 1)$. Consider that the length of f^Q is small, we traverse f^Q sequentially to find appropriate f^Q_p satisfying above equation. Finally, the service provider performs retrieval with respect to f^Q_p , and returns the N most similar images as the result R^0 .

4.2 Private Feedback

Relevance feedback can significantly improve the retrieval performance. However, since user labels relevant images during feedback, the user's search intention is completely exposed to the service provider. In order to leverage the performance gain and preserve the users' privacy at the same time, private feedback is developed here. Private feedback can deal with the feedback attack and result attack on R_1 . The new idea of private feedback is introducing confusing classes into feedback image set F_p .

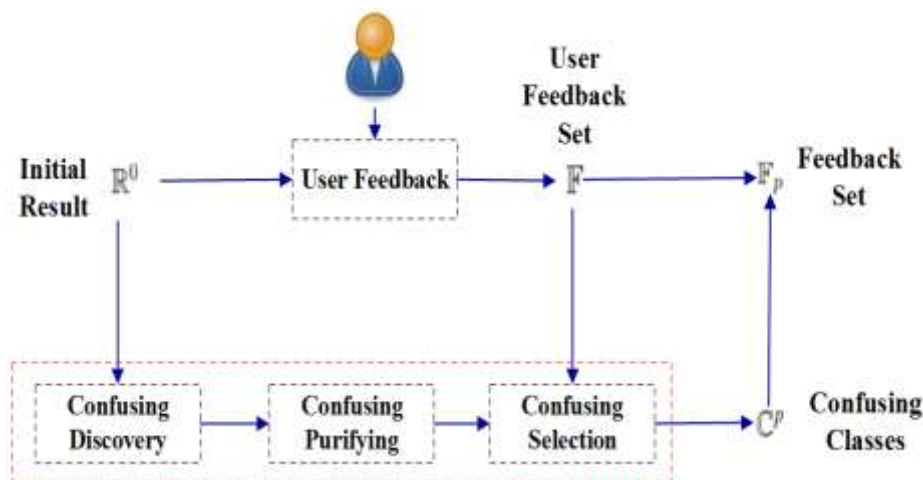


Fig-3: The flow chart of confusing classes introduction.

During the relevance feedback, the user randomly labels some relevant images in R_0 as F . During private feedback, the user introduces confusing classes into feedback image set F_p , and submits F_p to the service provider. A three-step approach is employed to introduce confusing classes, including confusing discovery, confusing purifying and confusing selection. The overall flow of confusing classes introduction is presented in Fig.3.

4.2.1 Confusing discovery

The objective of this step is to discover confusing classes candidates in R^0 . In this paper, we use k-means [10] to cluster the images in R^0 , and regard the result clusters of k-means as confusing classes candidates. The k-means relies on distances among images. Considering that R_0 are retrieved with private query f^0_p , the distances among images on feature components in f_p are little. Therefore, we compute the distances using remained feature components in f . The output of confusing discovery is k clusters, $C = \{C_1, \dots, C_k\}$, which will be used as confusing classes candidates in next step

4.2.2 Confusing purifying

The output of confusing discovery is k clusters, $C = \{C_1, \dots, C_k\}$. Ideally, each cluster should be in the same class. However, in practice, those clusters may be unclear. Therefore, we employ a consensus filters approach to further purify those clusters. The basic idea is, based on random subspace strategy [11], to construct an ensemble of SVMs [12] as filters. Then the outputs of SVMs are aggregated according to the consensus rules [13, 14].

Purifying is performed on clusters with the size a little larger than F . That is, $|C_i| \geq (1 + \alpha)|F|$, where α is an adjustment factor. Regarding each cluster C_i as positive training dataset, and remained examples $R_0 - C_i$ as negative training dataset, an ensemble of SVMs are constructed by using the random subspace method. Random subspace method is an ensemble learning method. It constructs base classifier on random sampled features instead of the entire feature set. Given the original feature space f , a base SVM classifier C_t is constructed on the subspace f_t randomly sampled from f . Suppose T SVM classifiers are constructed, those SVMs are then used to classify the images in C_i . According to the consensus rule, only the images labelled as positive by all SVMs are retained in C_i . Algorithm 1 presents the procedure of consensus filters based confusing purifying.

```

Input : The clusters set :  $C = \{C_1, \dots, C_k\}$ ,
          the initial retrieval result :  $R^0$ .
Output: The purified clusters set :  $C^P$ .

1  $C^P \leftarrow \emptyset$ ;
   // Choose the clusters with certain size.
2  $C^{temp} \leftarrow \{C_i | (C_i \in C) \wedge (|C_i| \geq (1 + \alpha)|F|)\}$ ;
3 foreach Cluster  $C_i \in C^{temp}$  do
   // Regard  $C_i$  as positive training dataset,
   // and  $(R^0 - C_i)$  as negative training
   // dataset.
4    $S^+ \leftarrow C_i$ ;
5    $S^- \leftarrow (R^0 - C_i)$ ;
   // Train  $T$  SVMs using the random subspace
   // strategy.
6   for  $t \leftarrow 1$  to  $T$  do
   | // Bootstrap  $f_t$  of length  $L$  from
   | // original feature space  $f$ .
7   |  $f_t \leftarrow \text{Bootstrap}(f, L)$ ;
8   |  $C_t \leftarrow \text{TrainSVM}(f_t, S^+, S^-)$ ;
9   end
   // Classify images in  $C_i$  using the
   // constructed SVMs.
10  for  $t \leftarrow 1$  to  $T$  do
11  |  $L_t \leftarrow \text{SVMPredict}(C_t, C_i)$ ;
12  end
   // Aggregate the outputs of SVMs with
   // consensus rule.
13   $L \leftarrow \text{ConsensusAggregation}(\{L_t\}_{t=1}^T)$ ;
   // Purify the cluster based on aggregation
   // result.
14   $C_i \leftarrow \{I_m | (I_m \in C_i) \wedge (L(I_m) == \text{positive})\}$ ;
15   $C^P \leftarrow C^P \cup C_i$ ;
16 end
17 return  $C^P$ .
    
```

Algorithm 1: Consensus filters based confusing purifying

4.2.3 Confusing Selection

The goal of this step is to select $K - 1$ clusters to form image set. Suppose the cluster set after purifying is C_p , the valid clusters firstly should be of certain size, that is, $|C_i| \geq (1 + \alpha)|F|$, where α is the adjustment factor. Secondly, the clusters should not contain any images in F , as $C_i \cap F = \emptyset$. We randomly select $K - 1$ clusters from those valid clusters. For each of those selected clusters, we randomly choose $(1 + \alpha)|F|$ images. Finally, those $K - 1$ clusters with the size of $(1 + \alpha)|F|$ are added to the user feedback image set F , forming the private feedback set F_p . Finally, regarding F_p as positive examples, and randomly selected images from $R^0 - F_p$ with the same size of F_p as negative examples, the service provider constructs SVM to rank the images in Ω , and returns the top N ranked images, as result R^1 . The overall process of private feedback is summarized in Algorithm 2.

```

Input : The initial retrieval result :  $R^0$ .
Output: The feedback retrieval result :  $R^1$ .

// The user side.
// The user randomly labels some relevant
// images in  $R^0$  as  $F$ .
1  $F \leftarrow \text{UserFeedback}(R^0)$ ;
// Introduce confusing classes through three
// steps: discovery, purifying and
// selection.
2  $C \leftarrow \text{Discovery}(R^0)$ ;
3  $C^p \leftarrow \text{Purifying}(C)$ ;
4  $C^p \leftarrow \text{Selection}(C^p)$ ;
// Construct the feedback image set
// following the  $K$ -anonymity feedback
// principle.
5  $F_p = \cup_{C_i \in C^p} C_i \cup F$ ;

// The service provider side.
// Construct SVM using  $F_p$  and randomly
// selected  $N$  to rank the images in  $\Omega$ .
6  $C \leftarrow \text{TrainSVM}(f, F_p, N)$ ;
7  $R^1 \leftarrow \text{SVMPredict}(C, \Omega)$ ;
8 return  $R^1$ .
    
```

Algorithm 2: Private feedback

4.3 Local Retrieval

Since some confusing classes are introduced for private feedback, the feedback retrieval result R^1 would not be satisfactory. Therefore, we re-rank the images in R^1 in the user side locally, namely local retrieval, to enhance the retrieval performance. In local retrieval, the user feedback image set F is regarded as positive training dataset, and the introduced confusing classes are regarded as negative training dataset. The SVM is then trained to re-rank the images in R^1 , producing the final result R^1_L to the user.

4.4 Copy Right Problem

The watermarking technology is employed for the copy deterrence in the proposed scheme. At the beginning, a unique watermark associated with the query user is embedded into the encrypted images. Then, the encrypted and watermarked images are sent to the query user. After receiving the images, the query user can directly decrypt these images. The watermark

is still preserved after the decryption. When an illegal copy of the image is found, the unlawful user who made the illegal distribution can be traced by examining the watermark in the image. This will deter the illegal distributions.

The images retrieved by the service provider are encrypted with a standard stream cipher. The keys are generated by a one-way pseudorandom number generator which takes as input a secret key and the unique image identity. The keys for different images are different with each other. In this process each pixel in an grayscale image is composed of 8 binary bits. The pixel bits of image are encrypted into random bits through the exclusive-or operation with a standard stream cipher. Then, the encrypted image is segmented into non overlapping blocks and a part of them are randomly chosen to carry watermark bits. Next, the pixels in each of chosen blocks are randomly divided into two sets S0 and S1 according to a secret key. If the watermark bit is 0, flip the 3 least significant bits (LSBs) of the pixels in S0. Otherwise, flip the 3 LSBs of the pixels in S1. In this way, an encrypted and watermarked image is generated. The encrypted and watermarked image can be decrypted by the same stream cipher. The decrypted image still contains the watermark in it.

If the image owner finds his image is exposed to someone unauthorized, the owner can submit the illegal copy and the corresponding original image. Then service provider takes the responsibility to extract the watermark from the doubtful image. Then, the extracted watermark is used to identify the illegal user. In the extraction process, extract the watermark bits by comparing the watermarked image and its corresponding original version.

5 Experiments and Results

In this paper, Corel database of 1000 images is used for the retrieval evaluation of CBIR system. The database consists of 10 categories which are classified as African people, Beaches, Building, Buses, Dinosaurs, Elephant, Flowers, Horses, Mountains and Food. This database has been widely used as ground-truth for evaluating color image retrieval each category consists of 100 images. These images are in JPEG format with resolution of 384×256 or 256×384 pixels. PYTHON 3.7 with Intel(R) Core(TM) i3 processor and 4GB RAM is used to simulate the algorithms. Precision and Recall are used as evaluation search metrics. The similarity measures between query feature vector set with each image feature vector set in the database is computed using Euclidean distance. Smaller value of Euclidean distance indicates better similarity between images.

5.1 Evaluation metrics

This paper use success rate to measure the performance of privacy, and precision to measure the performance of retrieval. Success rate is used to measure the service provider's attack success, which is the ration of the number of successful attacks over the number of all the attacks.

$$\text{Success Rate} = \frac{\text{successful attack}}{\text{all the attacks}}$$

Precision is used to measure the accuracy of image retrieval. For a query q, precision is defined as the fraction of retrieved images that are relevant.

$$\text{Precision (q)} = \frac{\# \text{ relevant images}}{\# \text{ retrieved images}}$$

Traditional RF-CBIR scheme suffers from the query attack, result attack and feedback attack. The success rate of feedback attack drops dramatically from 100% in normal Relevance Feedback -CBIR to 41% in proposed-CBIR, while the average precision of top 20 decreases slightly from 46% to 42%.

Query images from different categories of image database with its corresponding first five relevant retrieval are shown in Fig.4.



Fig-4: Retrieval of Relevant Images

5.2 Watermark Extraction Accuracy

After receiving the watermarked image, the query user may change it by regular image processing operations, e.g., JPEG compression. In addition, with the knowledge of watermarking algorithm, the user may try to remove the watermark bits by flipping the bits of image pixel. In this case, the watermark bits could not be extracted with the 100% accuracy. In this paper, the watermark extraction accuracy is defined as $p = N_e/N_w$, where N_e is the number of correctly extracted watermark bits and N_w is the total number of the watermark bits. A larger block size s helps to enhance the robustness of our watermarking algorithm.

3. CONCLUSIONS

This work propose a privacy-enhancing method for large-scale content-based image retrieval. A new Private Relevance Feedback-CBIR (PRF-CBIR) scheme is proposed to protect the user’s search intention and leverage the performance gain of relevance feedback. PRF-CBIR consists of three stage: private query, private feedback and local retrieval. PRF-CBIR can deal with query attack, result attack and feedback attack existing in Relevance Feedback-CBIR. This scheme can effectively control privacy leakage and significantly reduce the attack success probability. An elaborate watermark-based protocol in the encryption domain is designed for copy-deterrence . Different from common watermarking techniques, the proposed protocol needs to embed the watermark directly into the encrypted images. After receiving the encrypted and watermarked images, the query user needs to decrypt the images directly. And the decryption should not affect the watermark in the images.

REFERENCES

- [1] W. Lu, A. L. Varna, A. Swaminathan, and M. Wu, “Secure image retrieval through feature protection,” in The 2009 IEEE International Conference on Acoustics, Speech and Signal Processing, Taipei, 2009, pp. 1533–1536.
- [2] J. Yuan, S. Yu, and L. Guo, “SEISA: Secure and efficient encrypted image search with access control,” in 2015 IEEE Conference on Computer Communications (INFOCOM), Kowloon, 2015, pp. 2083–2091
- [3] Z. Xia, Y. Zhu, X. Sun, Z. Qin, and K. Ren, “Towards privacy-preserving content-based image retrieval in cloud computing,” IEEE Transactions on Cloud Computing, no. 99, 2015.

- [4] L. Weng, L. Amsaleg, A. Morton, and S. Marchand-Maillet, "A privacy-preserving framework for large-scale contentbased information retrieval," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 1, pp. 152–167, 2015.
- [5] Mukul Majhi and Sushila Maheshkar, "Privacy Preserving in CBIR Using Color and Texture Features," in 2016 international conference on parallel, distributed and grid computing (PDGC).
- [6] Young Deok Chun, Sang Yong Seo, and Nam Chul Kim, "Image Retrieval Using BDIP and BVLC Moments," *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 13, No. 9, 2003.
- [7] Y. Xu, J. Gong, L. Xiong, Z. Xu, J. Wang, Y-q. Shi, "A Privacy-Preserving Content-based Image Retrieval Method in Cloud Environment," *Journal of Visual. Communication and Image Retrieval.*, volume 43, 2017, pp 164-172.
- [8] I. Jolliffe, *Principal component analysis*. Wiley Online Library, 2002.
- [9] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, 2000.
- [10] A. K. Jain, "Data clustering: 50 years beyond k-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651–666.
- [11] T. K. Ho, "The random subspace method for constructing decision forests," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 8, pp. 832–844, 1998.
- [12] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [13] C. E. Brodley and M. A. Friedl, "Identifying mislabeled training data," *Journal of Artificial Intelligence Research*, vol. 11, pp. 131–167, 1999.