

# Sensing Identicals Files and Eliminating SHA Base Algorithm (SIFE-SBA) for WhatApps

SonaSahu, Prof. AdityaSinha

Department of CSE, BITS, Bhopal, Madhya Pradesh, India

## ABSTRACT

An extreme transaction of the WhatsApp's which is more duplicate or replicate files. where whole documents may be served in different formats like JPG, JPEG, PNG., GIF, PDF, Doc, DOCX and more. That Files may get copied to avoid delays/data consumption to provide fault acceptance. large databases of WhatsApp's files store in the declassification effort with receiving folders. copies of the files are abundant in images, video, Audio and short text in the databases. In our examination on On-line Social Networks (WhatsApp's) of images, more than 60% of total data are exact duplicate. due to forwarding by friends with linked each other. In this algorithms for detecting replicated files before store in local database by checking with updated hash index which are more critical in WhatsApp's applications where data is received from many friends links sources or groups. The deduction of duplication files are very necessary, to reduce data consumption in runtime, to improve data storing capacity with higher accuracy.

**Keywords :** Unique Documents, Detecting Duplicate, Copying, Database Storing, Data Consumptions, Whatsapp's.

## I. INTRODUCTION

In an Online Social Networking applications WhatsApp's are increasingly becoming the medium for people to keep in touch with every one, share their information about their photos, daily activities, political upraising, and travels. Social networking has attracting hundreds of millions of users, spending billions of minutes on such services.

Duplicates are abundant in photo's, Video's in WhatsApp's databases. For example, popular mobile phone images base messages, Videos messages, and many more may be forwarded by lots of people, and billions of people may express their feelings on the same hot topic by WhatsApp's and images base messages. In our examination on mobile phone more than 60% Images base messages have duplicate.

In this paper, we propose a algorithm that automatically seeks and identifies copying data in WhatsApp's. The important concept behind its logic is that it services

user-specific or user-identifying data, collected from the user's original WhatsApp's profile to locate similar profiles across WhatsApp's

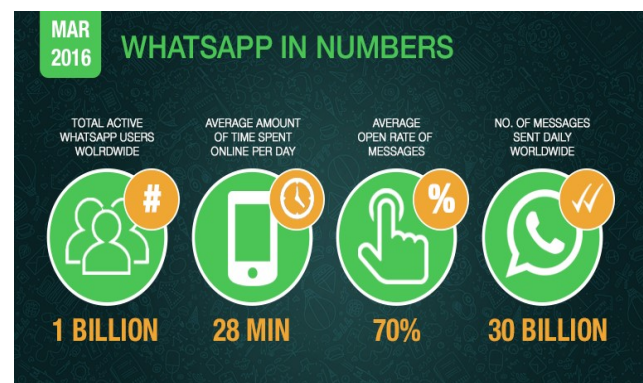
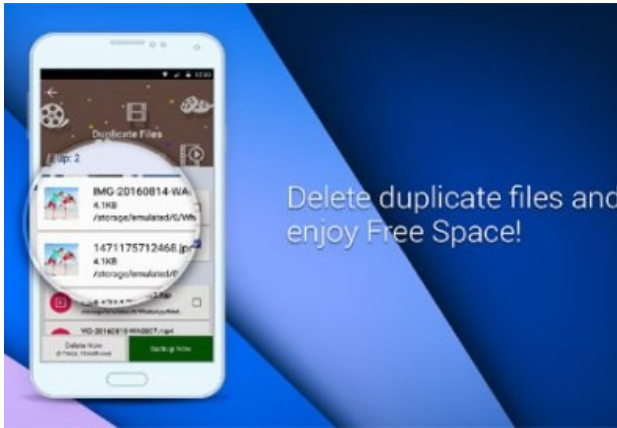


Figure 1: Actual User in WhatsApp's in March 2016

acknowledged and send folders. Any repaid results, depending on how rare the joint profile data is considered to be, are deemed distrustful and further inspection is performed. Finally, the user is presented with a list of likely profile duplication data.



**Figure 2:** Actual User in WhatsApp's Messages

In data folder of WhatsApp's contains a process stands for extract, transform and load. During the removal phase, scores of the data with hash values come to the data folder send or received. Some of these portals are of over-all attention some are extremely area exact. Independent of focus, the vast majority of portals obtain to the data, loosely called the documents, from multiple sources [1]. Obtaining data from multiple input sources typically results in the duplication show in figure 2,3, and 4. The discovery of duplicate documents within a collection has recently become an area of excessive interest [2] and is the focus of our described effort of the application.



**Figure 3:** Actual User in WhatsApp's Messages

Naturally inverted indexes are used to support efficient query processing in the information search and retrieval engines. Storing fake files affects both accuracy and the efficiency of Phones. Recovering virtual files in response to the user's query evidently drops the number of the valid responses provided to user, hence lowering

the accurateness of user's response set. Furthermore, processing duplicates necessitates additional calculation.



**Figure 4:** Actual User in WhatsApp's Messages

When WhatsApp's documents, 1 of the capacity think that, at minimum, matching URL's would recognize to exact matches. However, many WhatsApp's user are used dynamic presentation where in content changes depending on region or other variables. In the calculation, data providers often make several names for one of the site in an attempt to attract the users with unlike interests or viewpoints.

Fake where they come from single of the main problems with current geospatial databases is that they are known to contain many duplicate points (e.g., [3, 4, 5]). The main reason why geospatial databases contain fake is the databases are rarely formed totally. from scratch., and instead are built by combining measurements from numerous sources. Since some of the measurements are characterized in data from the several sources.

Duplicate morals can corrupt results of the statistical data processing and investigation. when instead of a actual measurement result, here see several dimension results confirming each other, and may get an erroneous impression so this measurement result is more consistent than it actually is Sensing and eliminating fake is therefore an important part of the assuring and refining quality of geospatial data, as recommended by the US Federal Standard [6].

## II. LITERATURE SURVEY

In this research article author [7], write documents that are 100% similar are termed to be duplicate documents and near duplicate documents (NDD) are not bitwise

identical but strikingly similar. If the NDD papers are clustered then they almost share the same cluster. The existence of near duplicate web pages are due to exact replica of the original site, mirrored sites, versioned sites, and multiple representations of the same physical object and plagiarized documents. Author proposed algorithm comprises of three phases Transliterate phase, filtering and Location Sensitive Bitwise Similarity method (LSBSM). It is to identify the query page is how similar to all the records in the repository. author have analyzed the system using the parameters like precision, recall, f-measure and efficiency, the results showed improvement in the values when compared with systems using existing weighting schemes which clarifies the efficiency of the proposed system. Mainly the elapsed time for the identification of near duplicate web pages has reduced and accuracy has increased.

In this research paper [8], online social networking has become one of the most important forms of today's communication. While an online social network can be attractive for many socially interesting features, its competitive edge will diminish if it is not able to keep pace with increasing user activities. Deploying more servers is an intuitive way to make the system scale, but for the best performance one needs to determine where best to put the data, whether replication is needed, and, if so, how. This paper is focused on replication; specifically, they propose S-CLONE, a socially-aware data replication scheme which can significantly improve a social network's efficiency by taking into account social relationships of its data. S-CLONE's performance is substantiated in our evaluation study.

In this research paper author [9], say that the Text mining, also known as Intelligent Text Analysis is an important research area. It is very difficult to focus on the most appropriate information due to the high dimensionality of data. Feature Extraction is one of the important techniques in data reduction to discover the most important features. Processing massive amount of data stored in a unstructured form is a challenging task. Several pre-processing methods and algorithms are needed to extract useful features from huge amount of data. The survey covers different text summarization, classification, clustering methods to discover useful features and also discovering query facets which are multiple groups of words or phrases that explain and summarize the content covered by a query thereby reducing time taken by the user.

In this paper [10] Social networking is one of the most popular Internet activities, with millions of users from around the world. The time spent on sites like Facebook or LinkedIn is constantly increasing at an impressive rate. At the same time, users populate their online profile with a plethora of information that aims at providing a complete and accurate representation of themselves. Attackers may duplicate a user's online presence in the same or across different social networks and, therefore, fool other users into forming trusting social relations with the fake profile. By abusing that implicit trust transferred from the concept of relations in the physical world, they can launch phishing attacks, harvest sensitive user information, or cause unfavourable repercussions to the legitimate profile's owner. In this paper they propose a methodology for detecting social network profile cloning. they present the architectural design and implementation details of a prototype system that can be employed by users to investigate whether they have fallen victims to such an attack. Our experimental results from the use of this prototype system prove its efficiency and also demonstrate its simplicity in terms of deployment by everyday users. Finally, they present the findings from a short study in terms of profile information exposed by social network users.

### III. PROPOSED ALGORITHM

The recognition of duplicates in WhatsApp's is performed during the storage of the each documents, ensuring that the individually distinct file is the stored in a single folder within case. When a client requests to the storage of a documents, the system performs a sequence of tasks:

1. Creates a sign  $G_i$  for the each document.
2. Marks in sign-location to sign and their obtains the location  $L_i$  of the consistent folder .
3. Examinations for a folder in location  $L_i$  within the  $n-1$  capacities that compose the case multicasting requests to capacities.
4. If a folder is found on 1 of capacities then the document is considered to be a duplicate and its reference counter is incremented by 1. Otherwise the document is stored in the new folder with location  $L_i$  in the capacity is identified by  $G_i \bmod n-1$ .
5. Finally, a *content key or links* referencing their folder is returned to the client.

The mod-based policy used to determine the capacity where to store a document divides the load equally among the sizes. However, clients can define the capacity where to store each document implementing alternative load-balancing policies. This way, in the presence of heterogeneous nodes, clients can impose higher workloads on the ones with higher throughput.

#### IV. ALGORITHM

Check WhatsApps Database  $WA_i$   
 For each record from  $WA_i$ ,  
 System Create a identicals record with arrange it in lexicographics order.  
 System Regularly Sance new data in receiving folder  
 Compute the New receiving data files directories  $WA(SHA)q_i = WA_i$   
 Compute the directories  $WA(SHA)p_i = WA_i$   
 If  $(SHA)p_i = R(SHA)q_i$  Then  
 Delete the identicals  $WA(SHA)q_i$   
 Link that identicals detected file link with  $WA_i$  index postion.  
 Else  
 Add in index  $WA_i$   
 As a result, we get an index-lexicographically ordered list of records:  
 $WA(1)..... :: .....WA(n-1)$ ,  
 where  $1,2,3,4,.....n-1$

In every WhatsApp's account client retrieves a document/File which is either image or document sending by their friends. In the WhatsApp's firstly check Web store decomposes that content key, identifies the signature of the document and the size that hosts the correspondent folder. The location of the folder in the volume is obtained by applying sign-location to the signature.

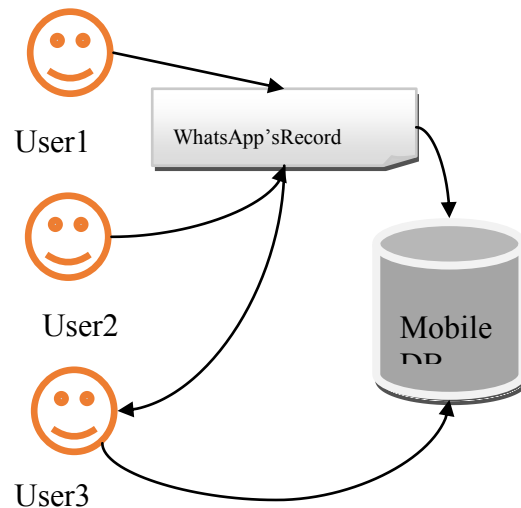


Figure 5. In WhatsApp's duplication of data file

Finally, all the document is kept in the folder is decompressed using this algorithm to detailed in folder's header, and the document is returned to WhatsApp's client. The delete operation is also invoked as soon as received duplicate file or documents as the argument. The location of folder is executed by the following same process as for retrieve operation. If the index counter is contained in file of in the folder with value 1, the folder is deleted the file and link that location with current account holders. Otherwise, the place of the counter is decremented. Since the location of document is determined by content key or links, the capacity where document is stored is directly accessed, both for the save and the delete operations. Therefore, performance of these operations is independent from the number of the capacities that compose an occasion.

In this proposed algorithm author use SHA for detecting duplication, This standard specifies four secure hash base algorithms that are, SHA-1 [11], SHA-256, SHA-384, and SHA-512. All four of the algorithms are iterative, one-way hash functions that can procedure a message to produce a reduced representation called a MD. These algorithms enable the determination of a message's integrity without any change to the message with a identical probability with outcome in a different MD. This property is useful in the generation and confirmation of digital signatures and the message authentication codes and in generation of the random numbers.

Our algorithm defined into the two stages preprocessing and hash calculation. Preprocessing involves the padding

with message parsing padded message into  $m$ -bit folders and setting initialization the values to be used in hash computation only. The hash computation is generated from a message schedule from padded message and uses that schedule along with the functions that coefficients the word operations to iteratively generate series of the hash values. The final hash value is generated by hash computation which is used to determine MD.

The four algorithms vary most significantly in number of the bits of safety that provided data being hashed this is directly related to MD length. When a secure hash algorithm is used in the combination with another algorithm, there may be requirements specified elsewhere that require the use of a secure hash algorithm with a certain number of bits of security. For example, if a message is being signed with a digital signature algorithm that provides 128 bits of security, then that signature algorithm may require the use of a secure hash algorithm that also provides 128 bits of security (e.g., SHA-256).

The performance numbers above were for a single negotiated implementation on an Intel Core 2 with 2.3 GHz processor with Windows 7 in 32-bit mode, and serve only as a rough point of the general comparison. [12].

Function is rapidly compares with large numbers of files to matching content by the computing the SHA hash of the each file and detecting fake with as soon as possible. The probability of the two non-identical files is hash will be same, even in a theoretical directory they containing millions of the files is remarkably remote or locally. Thus hashes rather than the file contents are compared with the process of the detecting fake is greatly accelerated.

## V. CONCLUSION

In this research paper author estimated its performance using various data collections of WhatsApp's users. In terms of human usability, The ultimate purpose of how the similar files which must be considered a fake or duplicated, files on human judgment bases. Therefore here our solution necessity to be easy to usage here author planned a new fake discovery algorithm which assessed its performance using multiple data collections of the WhatsApp's users. In this research paper here collections of data for experiment are used by a group in which too many user use WhatsApp's and share their

data for communication and take backup daily basis. After the collection of data here our algorithm is applied and analysis the result. Whole experiment was done in Matlab tools. In the terms of social usability is more similar to the document so our detection approach become helpful and it is perfect. This paper try to save data of user because data is more costly and its limited availability become data more important. And also this paper proof how our algorithm save storage which is current era most popular problem of the world in data storage field. This research intends aid to upcoming researchers in this field of fake data detection in the WhatsApp's to know available methods and their help to complete research in further road.

## VI. REFERENCE

- [1]. BRODER, A., GLASSMAN, S., MANASSE, S., AND ZWEIG, G. "Syntactic clustering of the web", In Proceedings of the Sixth International World Wide Web Conference (WWW6'97) (Santa Clara, CA., April). 391-404.
- [2]. SHIVAKUMAR, N. AND GARICA-MOLINA, H. 1998. Finding near-replicas of documents on the web. In Proceedings of Workshop on Web Databases (WebDB'98) (Valencia, Spain, March). 204-212.
- [3]. McCain, M., and William C., 1998. Integrating Quality Assurance into the GIS Project Life Cycle, Proceedings of the 1998 ESRI Users Conference. <http://www.dogcreek.com/html/documents.html>
- [4]. Good child, M., and Gopal, S. (Eds.), "Accuracy of Spatial Databases", Taylor & Francis, London. 1989.
- [5]. Scott, L. "Identification of GIS Attribute Error Using Exploratory Data Analysis", Professional Geographer 46(3), 378.386, 1994.
- [6]. FGDC Federal Geographic Data Committee, Content standard for digital geospatial metadata (revised June 1998), Federal Geographic Data Committee, Washington, D.C., <http://www.fgdc.gov/metadata/contstan.html>, 1998. FGDC-STD- 001-1998.
- [7]. LavanyaPamulaparty Dr. C.V. Guru Rao Dr. M. SreenivasaRao, "LSBSM: A Novel Method for Identification of Near Duplicates in Web Documents", International Journal of Computer Science and Information Security (IJCSIS), Vol. 15, No. 2, February 2017
- [8]. Mauro Conti, RadhaPoovendran, Marco Secchiero, "FakeBook: Detecting Fake Profiles in On-line Social Networks", ACM International Conference on Advances in Social Networks Analysis and Mining, IEEE 2012.

- [9]. Ramya R S, Venugopal K R, Iyengar S S&Patnaik L, “Feature Extraction and Duplicate Detection for Text Mining: A Survey”, Global Journal of Computer Science and Technology: C Software & Data Engineering, Volume 16 Issue 5 Version 1.0 Year 2016
- [10]. GeorgiosKontaxis, IasonasPolakis, Sotiris Ioannidis and Evangelos P. Markatos, “Detecting Social Network Profile Cloning”, 3rd International Workshop on Security and Social Networking, IEEE, 2011.
- [11]. The SHA-1 algorithm specified in this document is identical to the SHA-1 algorithm specified in FIPS 180-1.
- [12]. <http://en.wikipedia.org/wiki/SHA-2>.