



International Journal for Innovative Engineering and Management Research

A Peer Reviewed Open Access International Journal

www.ijiemr.org

COPY RIGHT



ELSEVIER
SSRN

2019IJIEMR. Personal use of this material is permitted. Permission from IJIEMR must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. No Reprint should be done to this paper, all copy right is authenticated to Paper Authors

IJIEMR Transactions, online available on 6th Mar 2018. Link

[:http://www.ijiemr.org/downloads.php?vol=Volume-08&issue=ISSUE-03](http://www.ijiemr.org/downloads.php?vol=Volume-08&issue=ISSUE-03)

Title: **IDENTICAL SEARCH BY USING STRING SIMILARITY BY USING HASHING**

Volume 08, Issue 03, Pages: 41–45.

Paper Authors

MR.A.JANARDHAN RAO, N.NAGARJUNA

Vignan's Lara Institute of Technology & Science



USE THIS BARCODE TO ACCESS YOUR ONLINE PAPER

To Secure Your Paper As Per **UGC Guidelines** We Are Providing A Electronic Bar Code

IDENTICAL SEARCH BY USING STRING SIMILARITY BY USING HASHING

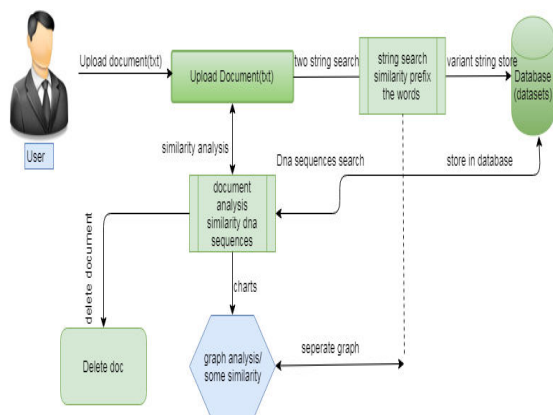
MR.A.JANARDHAN RAO¹, N.NAGARJUNA²

Assistant Professor¹, Department of M.C.A , Vignan's Lara Institute of Technology & Science

M.C.A Student², Department of M.C.A, Vignan's Lara Institute of Technology & Science

Abstract: String similarity search is a fundamental query that has been widely used for DNA sequencing, error-tolerant query auto-completion, and data cleaning needed in database, data warehouse and data mining. In this paper, we study string similarity search based on edit distance that is supported by many database management systems such as Oracle and PostgreSQL. Given the edit distance, $ed(s, t)$, between two strings, s and t , the string similarity search is to find every string t in a string database D which is similar to a query string s such that $ed(s, t) \leq \theta$ for a given threshold. In the literature, most existing work take a filter-and-verify approach, where the filter step is introduced to reduce the high verification cost of two strings by utilizing an index built offline for D . The two up-to-date approaches are prefix filtering and local filtering. In this paper, we study string similarity search where strings can be either short or long. Our approach can support long strings, which are not well supported by the existing approaches due to the size of the index built and the time to build such index. We propose two new hash-based labeling techniques, named OX label and XX label, for string similarity search. We assign a hash-label, H_s , to a string s , and prune the dissimilar strings by comparing two hash-labels, H_s and H_t , for two strings s and t in the filter step. The key idea behind is to take the dissimilar bit-patterns between two hash-labels. We discuss our hash-based approaches, address their pruning power, and give the algorithms. Our hash-based approaches achieve high efficiency, and keep its index size and index construction time one order of magnitude smaller than the existing approaches in our experiment at the same time.

ARCHITECTURE:



EXISTING SYSTEM:

Most of the existing string similarity search algorithms take a filter-and-verify approach. The filter step is introduced to reduce the verification cost of two strings, s and t , which is costly when two strings are long. In order to find similar strings in a string dataset D for a given query string s with a threshold, they first prune strings, t , that cannot be possibly similar with s such that $ed(s, t) > \theta$ using an index built offline for D in the filter step, and then verify those

strings that are possibly similar one by one in the verification step. The performance of an approach is measured by the query cost and the index cost. The query cost is the sum of the filter cost (the total running time in the filter step) and the verification cost (the total running time in the verification step). The index cost is the index construction time and the index space needed. To efficiently process string similarity search, the existing work attempts to prune strings in D as many as possible based on the index built offline. Almost all the existing work needs to know the edit distance threshold beforehand, in order to construct the index for a string dataset D , except for BitTree. Behm et al. propose a hierarchical structure containing different filters, e.g., the length and charsum filter, in Flamingo package. Gravano et al. propose to partition a string into a set of q -grams and prune a string pair (s, t) that have less than a certain number of common q -grams. The chunk-based approaches share the similar idea but partition the string using disjoint q -grams, called chun. Instead of using fixed-length q -grams, Li et al. selectively choose high-quality grams of variable length in index construction.

PROPOSED SYSTEM:

In this paper, we study string similarity search, when the query string s and the average string t in D can be long. The up-to-date approaches cannot efficiently process long string similarity search for the following main reasons. For the prefix filtering approaches, the main idea is to use a small number of q -grams for filtering. When strings become long, the pruning

power of such a small number of q -grams will reduce significantly. In addition, the prefix filtering approaches need to know before the index construction. However, when the average strings become long, users want to use different for string similar search: a small for short strings and a large for long strings. It cannot be easily handled by the prefix filtering approaches. For the local filtering approach, the BitTree index will be extremely large to be stored and it is time consuming to construct such an index. Different from the existing work in the literature, we propose new hash-based labeling for string similar search. Let H_s and H_t be two hash-labels for strings, s and t . We show that s and t are definitely dissimilar for a given using H_s and H_t . We propose two hash-based approaches, namely OX label and XX label. Both are in the scheme of $(\sim, \aleph, \}, \#)$. Here, \sim and \aleph are two functions to create a hash-label H_s for a string s , and $\}$ and $\#$ are two functions to compare two hash-labels, H_s and H_t for two strings, s and t . The key idea behind is to take the dissimilar bit-patterns between two hash-labels. We discuss our hash-based approaches, address their pruning power, and give the algorithms. New optimizations to the verification algorithm are proposed for efficiently verifying whether a candidate string is an answer. We have conducted extensive performance studies and confirm the efficiency of our hash-based approaches in both datasets of long strings and datasets of short strings with much smaller index size.

ALGORITHM:

String Similarity Search: Given a string dataset D of n strings, a query string s and an edit distance threshold ϵ , the string similarity search problem is to find all strings $t \in D$ such that $ed(s, t) \leq \epsilon$. A well-known algorithm to compute the edit distance between two strings s and t is to fill an edit distance matrix of size $(|s| + 1) \times (|t| + 1)$ using dynamic programming. However, it requires $O(|s| \cdot |t|)$ time complexity which is costly for long strings. The *filter-and-verify* framework adopted by the existing work builds an index to prune the dissimilar strings of the query string in the dataset D in the filter step, and verifies the remaining candidates to get the real result in the verification step. The filter step is important to reduce the cost of computing the edit distance between two strings, by pruning the strings that cannot be possibly in the final results as many as possible using the index built offline. There are some simple heuristics that can be applied in the filter step. The length-filter is such an example, which prunes the string t if $||s| - |t|| > \epsilon$. The index built will further prune strings that cannot be simply pruned by such simple heuristics.

$$|Q_s \cap Q_t| \geq \max\{|s|, |t|\} + q - 1 - q\epsilon$$

MODULES:

1. UPLOAD PRODUCT

The registered users are authorized to upload the product. The product owners have ability to change or even delete the product from the application at

any point of time. The products can be viewed to other users and product owners can only access the details.

2. STRING SIMILARITY SEARCH

The uploaded products are listed in the users' view. There are lot of products are listed and in order to avoid congestion, the search can be available to make utilize the products in effective way. The searches have more number of details. In order to avoid the congestions searches can be utilized and give suggestion.

3. UPLOAD DOCUMENT ANALYSIS

According to user search it shows the suggestion of product can be shown to the user. The products are shown to user according to most searches and have different types of search to get the details and better retrieval of product in order to implement and make use of the search.

4. GRAPH ANALYSIS

Graph analysis of details can be taken from the data which are utilized in flow of project. The graph can be utilized to showcase the products maximum retrieval by users search and how effective to user while they are searching in the system..

Conclusion:

In this paper, we study two new hash-based approaches, OX label and XX label, for string similarity search based on edit distance, where $OX = (\sim, \vee, \oplus, \#)$ and $XX = (\sim, \oplus, \oplus, \#)$. Both OX and XX label use the same last two functions, \oplus and $\#$, to compare two hash-labels for pruning. But they take a different way to create the hash-labels. Here, OX label uses two functions, \sim and \vee , to create a hash-label for a string, whereas XX label uses two functions, \sim and \oplus , to create a hash-label for a string. We prove that both OX and XX label can be used to prune dissimilar strings, s and t , when $ed(s, t) > .$ The index size for OX label and XX label is determined by L , and the hash-label for string of any length has the same L (the number of bits). We analyze the pruning power by OX label and XX label. We show that OX label is effective when L is sufficiently large comparing to the sum of the lengths of two strings, s and t . We also show that the pruning power of XX label only depends on the number of different q -grams between the q -gram set Q_s and the q -gram set Q_t for s and t , and can be effectively used for both short and long string similarity pruning. We conducted extensive performance studies using 6 real string datasets.

References

1. Sharing in MULTICS. In Proceedings of the Fourth Symposium on Operating System Principles, SOSP 1973, Thomas J. Watson, Research Center, Yorktown Heights, New York, USA, October 15-17, 1973.
2. Robert Morris and Ken Thompson. Password Security: A Case History, 1979. <http://cs-www.cs.yale.edu/homes/arvind/cs422/doc/unix-sec.pdf>.
3. Philippe Oechslin. Making a faster cryptanalytic time-memory trade-off. In Dan Boneh, editor, Advances in Cryptology – CRYPTO 2003, 23rd Annual International Cryptology Conference, Santa Barbara, California, USA, August 17-21, 2003, Proceedings, volume 2729 of Lecture Notes in Computer Science, pages 617–630. Springer, 2003.
4. Password Hashing Competition (PHC), 2014. <https://passwordhashing.net/index.html>.
5. Donghoon Chang, Arpan Jati, Sweta Mishra, and Somitra Kumar Sanadhya. Rig: A simple, secure and flexible design for password hashing. In Dongdai Lin, Moti Yung, and Jianying Zhou, editors, Information Security and Cryptology - 10th International Conference, Inscrypt 2014, Beijing, China, December 13-15, 2014, Revised Selected Papers, volume 8957 of Lecture Notes in Computer Science, pages 361–381. Springer, 2014.
6. Ari Juels and Ronald L. Rivest. Honeywords: making passwordcracking detectable. In 2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13, Berlin, Germany, November 4- 8, 2013, 2013.
7. Fred Cohen. The Use of Deception Techniques: Honeypots and Decoys. <http://all.net/journal/deception/DeceptionTechniques.pdf>.



8. Lance Spitzner. Honeytokens: The Other Honey-pot, 2003. <http://www.symantec.com/connect/articles/honeytokens-other-honey-pot>.
9. Hristo Bojinov, Elie Bursztein, Xavier Boyen, and Dan Boneh. Kamouflage: Loss-resistant password management. In Computer Security - ESORICS 2010, 15th European Symposium on Research in Computer Security, Athens, Greece, September 20-22, 2010. Proceedings, pages 286–302, 2010.
10. Wikipedia contributors. 2012 LinkedIn hack. Wikipedia, The Free Encyclopedia, Date retrieved: 29 May 2016. Available at: https://en.wikipedia.org/w/index.php?title=2012_LinkedIn_hack&oldid=722095159.
11. Bruce Schneier. Cryptographic Blunders Revealed by Adobe's Password Leak. Schneier on Security, 2013. Available at: https://www.schneier.com/blog/archives/2013/11/cryptographic_b.html.
12. Swati Khandelwal. Hacking any eBay Account in just 1 minute, 2014. Available at: <http://thehackernews.com/2014/09/hacking-ebay-accounts.html>.
13. Wikipedia contributors. Ashley Madison data breach. Wikipedia, The Free

Encyclopedia, Date retrieved: 29 May 2016. Available at: https://en.wikipedia.org/w/index.php?title=Ashley_Madison_data_breach&oldid=721001290.

14. Troy Hunt. Observations and thoughts on the LinkedIn data breach, 2015. Available at: <https://www.troyhunt.com/observations-and-thoughts-on-the-linkedin-data-breach/>.

15. Michael Gilleland. Levenshtein Distance, in Three Flavors. Available at: <http://people.cs.pitt.edu/~kirk/cs1501/assignments/editdistance/Levenshtein%20Distance.htm>.

15. Bibliography

- (1) Java Complete Reference by Herbert Shield
- (2) Database Programming with JDBC and Java by George Reese
- (3) Java and XML By Brett McLaughlin
- (4) Wikipedia, URL: <http://www.wikipedia.org>.
- (5) Answers.com, Online Dictionary, Encyclopedia and much more, URL: <http://www.answers.com>
- (6) Google, URL: <http://www.google.co.in>
- (7) Project Management URL: <http://www.startwright.com/project.html>