

Efficient Autocompletion with Error tolerance

Sheng HU[†] Chuan XIAO[‡] and Yoshiharu ISHIKAWA[†]

[†] Graduate School of Information Science, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, 466-8601 Japan

[‡] Institute for Advanced Research, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, 466-8601 Japan

E-mail: [†] hu@db.ss.is.nagoya-u.ac.jp, [‡] {chuanx, y-ishikawa}@nagoya-u.ac.jp

Abstract Most popular search engines have provided query autocompletion features in their systems. When a user issues a query, without inputting complete keywords, the search engine will prompt several options which complete the remaining part of the query automatically. This research will focus on investigating novel techniques for error-tolerant query autocompletion by addressing several key issues of this feature, and eventually improve the efficiency and the accuracy of the query processing. In our work, we proposed an improve expansion algorithm based on current state-of-the-art work. Moreover, we also improve the index reduction method to reduce the index size as much as possible. Experiment results show that our new method can improve the algorithm efficiency and reduce the index size very well.

Keyword Query Autocompletion, Query Processing, Textual Database

1. Introduction

Recently, the user experience of search engine becomes more and more important. Most popular search engines have provided query autocompletion features in their systems. When a user issues a query, without inputting complete keywords, the search engine will prompt several options which complete the remaining part of the query automatically.

Moreover, search engines like Google, have equipped their system with an advanced autocompletion technique, which is called autocompletion with error tolerance. Even a user issues an uncompleted query which may contain typos, the search engine can still identify the intended query and give the correct answers. As more and more users are using smartphones as search engine interface, this feature is considered of great importance and prospective.

This research will focus on investigating novel techniques for error-tolerant query autocompletion by addressing several key issues of this feature, e.g., improving the efficiency and the accuracy of query processing. Figure.1 shows an example of the result of autocompletion with inputting “coff”.

```
coff
coffee
coffeescript
coffee bar cowboy
coffee kajita
coffee break
coffin
```

Figure 1. An example of autocompletion of “coff”

2. Related Work

There are various ways to implement the autocompletion technique[1][6]. Currently, trie index[2] becomes very popular and is adopted as the most efficient way among the current research works. The ICPAN algorithm put forward by Guoliang Li[3] is considered the most reliable trie traversal algorithm while IncNGTrie which is put forward by Chuan Xiao[4] is considered the most efficient trie traversal algorithm. The ICPAN[3] gives a definition about the Pivotal Active Node which is seemed as the basic composition unit of the trie index and can support 2-length characters error tolerance which means the search system can give the correct answer despite of two characters typos. The IncNGTrie[4] uses the Neighbor Generation method[5] so as to include the possibility of typos in advance. In the experiment, it can support 3-length characters error and achieves up to two orders of magnitude speedup over Li’s ICPAN method. However, IncNGTrie achieves the efficiency at the price of sacrificing large memory space, and under the special circumstances the overhead of the space may be unaffordable.

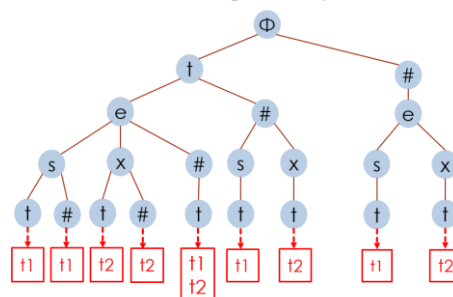


Figure 2. An example of IncNGTrie with keyword {test, text} with $\tau=1$

Example 2 Figure 2. shows a simple example of IncNTrie Index Structure which consists of two strings t1 and t2. T1 represents “test” and t2 represents “text”. Single path from the root to one leaf represents one string. For example, from left side, we have Φ -t-e-s-t path, which represents string “test”. Meanwhile we use ‘#’ as the placeholder to represent the deletion conducted on the data string side. Below we give the definition of ‘#’.

3. Index reduction

In IncNGTrie research, they used two techniques to reduce the index size as much as possible, so as to make it fit the computer memory. One is Common Data String Merge, the other is Common Subtree Merge. After my experiments, I found some disadvantages of these two techniques. For the former one, this work simply does the reduction exhaustively without any consideration of the reduction extent parameterization. Thus I come up with the idea that I can introduce the parameter p to control the reduction extent. In the original work, this parameter is set to 1, that is, only when the parent node share completely the identical data strings, it'll be merge into 1 string, which drastically limits the benefit of the reduction. And after I introduce this parameter ***BRC***, we can generally control the benefit and gain from the common data string merge technique. Here we give the definition of Branch Reduction Control.

We use BRC as the parameter to control how much we will reduce the index size. First we will traverse the whole trie in a Depth-First order, and check every nodes to see if its descendants amount is lower than BRC, if so, we will

Example 3 Figure 3, Figure 4, Figure 5 give an example of the trie variation under the condition when we increase BRC from 1 to 3 consecutively. If we set the BRC parameter to 1, the leaf nodes' parent node will be merged when the descendant identical data string number less equals than 1. Numbers in the bracket represent the amount of the descendant identical data string under each node. As showed in Figure 3, note that in the left subtree under the node "a", as all the descendants under "a" only represent the same string "t3", all the paths start from "a" and contain "#" will be removed from the index. In Figure 3, we change the nodes' color into grey if they are removed. In Figure 4 and Figure 5, we increase BRC to 2 and 3, as the nodes satisfying our reduction condition increase quite a lot, we can find the index size will be reduced incrementally. Especially when BRC=3, we successfully prune almost a half part in our example. In our experiment, this change will reduce the size of our tree drastically when τ is larger than 1 and decrease a little bit on the search time side because this is an inevitable tradeoff between the storage and efficiency. Moreover, I change the underlying data storage to make all the identical data point to just one copy instead of store a copy for every redundant record. The parameter we introduce could be adjusted according to the requirement of the real environment.

Figure 3. Index of {test, text, taxi} when BRC=1

Figure 4. Index of {test, text, taxi} when BRC=2

and the larger tau is, the more the total size is reduced. When BRC is increased to 9999999, the whole IncNGTrie is degenerated to a direct trie, which has a very small size.

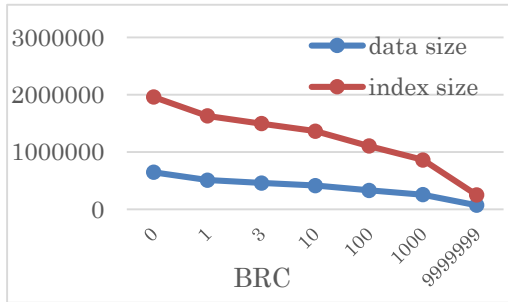


Figure 6. Data size and index size reduction effect when $\tau=1$

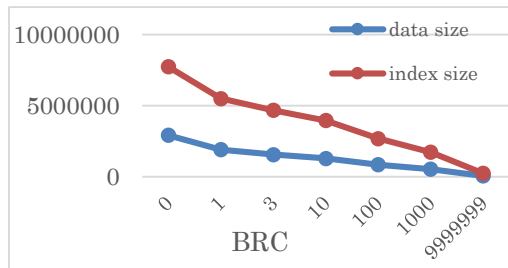


Figure 7. Data size and index size reduction effect when $\tau=2$

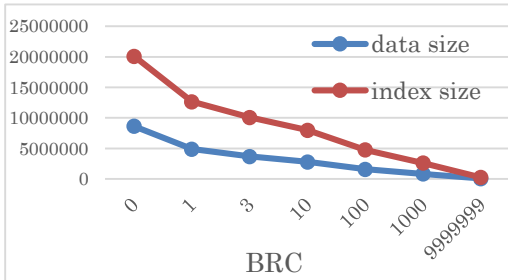


Figure 8. Data size and index size reduction effect when $\tau=3$

5.2 Effect of expansion improvement

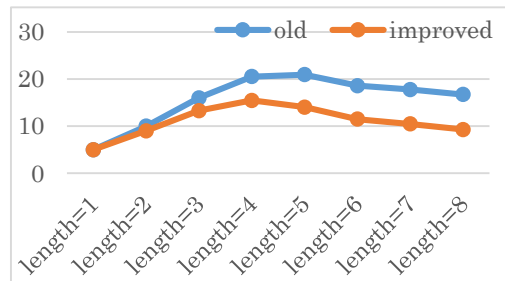


Figure 9. Returned active nodes amount when $\tau=1$

The three figures Figure 9, Figure 10, Figure 11 show the expansion improvement effect under the condition $\tau=1$, 2 and 3 . X-axis represents the query length and Y-axis represents the active node states amount. The common point of the three is that, as the length of query word in-

creases, the improved expansion algorithm has much less active states, which can remarkably save the memory and improve the time efficiency compared with the old algorithm. When query length equals 8, the improved algorithm's active nodes states number is almost a half of the old one, which means our method reduce over 50% active node states when query is long enough.

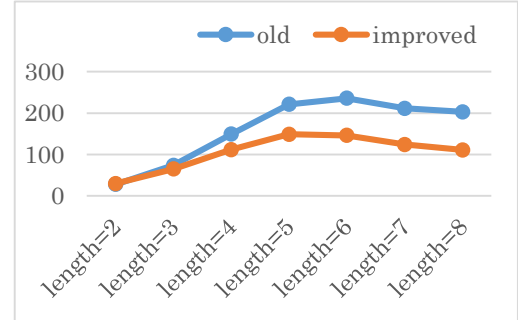


Figure 10. Returned active nodes amount when $\tau=2$

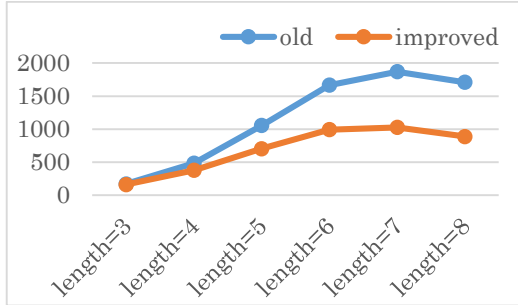


Figure 11. Returned active nodes amount when $\tau=3$

References

- [1] L. Boytsov. Indexing methods for approximate dictionary searching: Comparative analysis. ACM Journal of Experimental Algorithmics, 16(1), 2011.
- [2] D. Deng, G. Li, and J. Feng. An efficient trie-based method for approximate entity extraction with edit-distance constraints. In ICDE, pages 762-773, 2012.
- [3] G. Li, S. Ji, C. Li, J. Feng. Efficient Fuzzy Full-text Type-ahead Search. VLDB Journal, 2011.
- [4] C. Xiao, J. Qin, W. Wang, Y. Ishikawa, K. Tsuda, K. Sadakane. Efficient Error-tolerant Query Autocompletion. Proceedings of the VLDB Endowment, Vol. 6, No. 6, 2013.
- [5] T. Bocek, E. Hunt, and B. Stiller. Fast Similarity Search in Large Dictionaries. Technical Report ifi-2007.02, Department of Informatics, University of Zurich, April 2007.
- [6] Inci Cetindil, Jamshid Esmaelnezhad, Taewoo Kim, and Chen Li. Efficient instantfuzzy search with proximity ranking. In Data Engineering (ICDE), 2014 IEEE 30th International Conference on, pages 328-339. IEEE, 2014.
- [7] Thomas Bocek, Ela Hunt, Burkhard Stiller, and Fabio Hecht. Fast similarity search in large dictionaries. University of Zurich, 2007.
- [8] Moshe Mor and Aviezri S. Fraenkel. A hash code method for detecting and correcting spelling errors. Communications of the ACM, 25(12):935-938, 1982.