



Journal of Applied Sciences

ISSN 1812-5654

science
alert

ANSI*net*
an open access publisher
<http://ansinet.com>

A Two-Layer Dictionary Organization for Quickly Retrieval

¹Jau-Ji Shen and ²Chun-Ta Li

¹Department of Information Management,

National Formosa University, 64 Wen-Hwa Road, Huwei, Yunlin, Taiwan 632, R.O.C.

²Department of Computer Science, National Chung Hsing University,

250 Kuo Kuang Road, Taichung, Taiwan 402, R.O.C.

Abstract: Present goal in this study was to develop a dictionary searching technique designed specifically for World Wide Web, which is called “Two-layer dictionary retrieval”. Two-layer dictionary retrieval, through “signature extraction process” which extracts the “signature” from every vocabulary in the dictionary, generates the dictionary signature file and distributes to every user. The user want to looks up a vocabulary in the dictionary, through the same process that extracts the signature from which the vocabulary is desired, such that the user can search this vocabulary on the downloaded dictionary signature file. Our technique is not only able to minimize the searching area in the dictionary, but is also able to minimize the space required to establish an index so as to greatly improve the efficiency of dictionary retrieval over the Internet. According to the experimental results, our method is five times better than the other methods available in terms of efficiency.

Key words: Dictionary organization, false drop, indexing, internet, spelling check

INTRODUCTION

Nowadays more and more data files have been quickly created each day since 1995, the year in which the World Wide Web has begun to advance so rapidly. To quickly retrieve the desired information for users from such enormous data files is an important issue of dictionary retrieval.

There are currently a number of dictionary searching techniques that have been proposed (Angell *et al.*, 1983; Burkhard and Keller, 1973; Doster, 1977; Hall and Dowling, 1980; Kohonen, 1987; Morgan, 1970; Owolabi, 1996; Owolabi and McGregor, 1988; Salton and Wong, 1978; Szanzer, 1969; Szanzer, 1973), which primarily engage themselves in dictionary organization. That is, they place their strategy upon organization for the index establishment. For example, the length of characters contained in a vocabulary is taken as a basis for classification (Morgan, 1970; Szanzer, 1969), or the first letter in a vocabulary is taken as a classification (Hall and Dowling, 1980; Szanzer, 1973). The above two classification methods, after undergoing the experiment (Owolabi, 1996) that O. Owolabi pointed out, result in an uneven distribution that impedes the searching efficiency. Thus O. Owolabi combined the two methods in which he

first classifies the primary index by the length of characters in a vocabulary, then moves on to the secondary index and classifies by the first letter from the length of characters in the vocabulary. This way he can amend the uneven distribution, but this way brings up another issue leading to larger space required for the indexing. For example, assume in the primary index containing the length of characters from 2 to 50 and there are a total of 49 classifications, if he wants to classify each vocabulary by its first letter in the subordinate index and each length of characters contains 26 classifications (from A to Z), there will be a total of 1274 classifications needed to establish the secondary index and thus the resulting large file will make it difficult to be transferred over the Internet. Therefore our goal in this study is to make possible the even distribution in the dictionary and minimize the additional space required to establish the indexes. Furthermore, our technique is not only able to greatly improve the efficiency of dictionary retrieval over the Internet, but also to be applied on the spelling check that it can quickly retrieve stored lexicon to spell check the vocabulary does really correctness in the dictionary.

The concept of two-layer dictionary retrieval: Present study proposes a two-layer dictionary organization

Table 1: Statistical table of the selected medical dictionary containing the number of occurrences for variant characters

Sort	Character	Number of occurrences	Sort	Character	Number of occurrences	Sort	Character	Number of occurrences
1	e	202760	17	y	45348	33	3	386
2	a	191554	18	g	39335	34	4	329
3	l	182916	19	b	28847	35	5	289
4	o	159617	20	f	25856	36	(252
5	r	149568	21	v	19116	37)	252
6	s	138910	22	x	9334	38	6	224
7	t	136899	23	-	9232	39	7	141
8	l	117922	24	,	6594	40	9	95
9	Space	103379	25	w	6268	41	8	90
10	c	102408	26	z	4267	42	/	61
11	u	71942	27	'	4228	43	:	51
12	p	70857	28	j	2498	44	.	31
13	m	70257	29	q	2070	45	+	24
14	h	56494	30	l	749	46	&	8
15	d	56454	31	2	489	47	-	3
16			32			48		

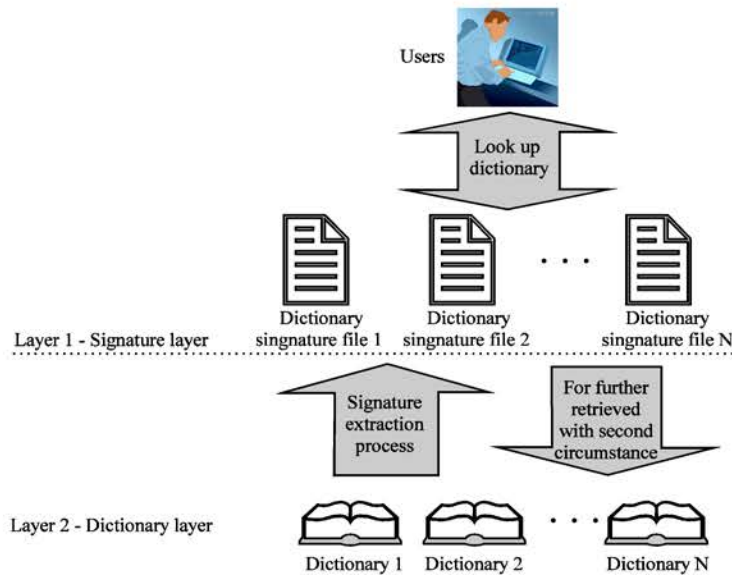


Fig. 1: The architecture of two-layer dictionary retrieval

method also called two-layer dictionary retrieval, in which its architecture is shown in Fig. 1.

The architecture in Fig. 1. can be applied to any current types of dictionaries such as Medical Dictionary, Webster Dictionary, etc. During the pre-process for each dictionary, after the signature extraction process that extracts a dictionary signature file in the first step, the dictionary signature file can now be downloaded depending on user's specific needs. When the user wants to look up the vocabulary, he will directly make his initial retrieval on the specific dictionary signature file at layer 1. Since the dictionary signature file after undergoing the signature extraction process is considered small-sized, this will greatly reduce the cost for sending the information back to the user over the Internet.

There are only two circumstances of retrievals. First circumstance is that the signature for which the vocabulary is to be retrieved after the same signature extraction process does not exist in the dictionary

signature file. At this time, the user will receive a definite reply that the vocabulary to be retrieved does not exist in the original dictionary. Second circumstance is that the signature for which the vocabulary is to be retrieved after the signature extraction process does exist in the dictionary signature file, implying that the vocabulary to be retrieved may exist in the original dictionary. It will then remotely log onto the dictionary database at the back end for further retrieval at layer 2. Because the dictionary is pre-classified by the signature, that the searching can be achieved on a few specific areas in the dictionary without going through the entire dictionary. It is this way that makes our dictionary retrieval better prepared and applied to current networks.

The signature extraction process

The encoding table of dictionary: From the architecture of two-layer dictionary retrieval mentioned in the last Section, each vocabulary in the dictionary must be

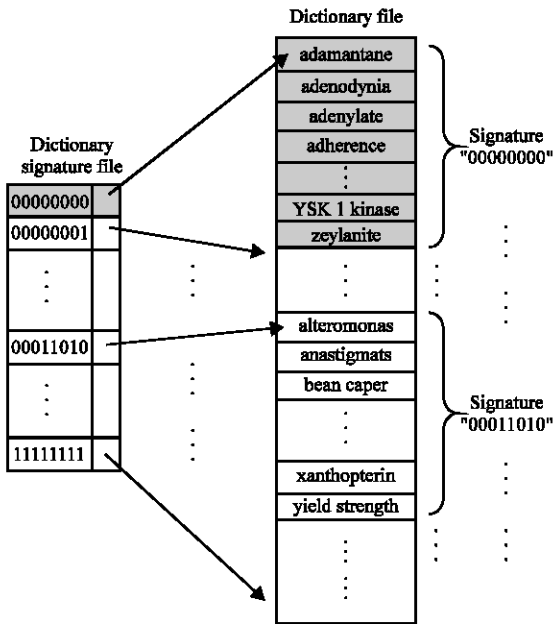


Fig. 4: Establishment of indexing for C = 8

lower distinguishability; extravagant number of partition take up too much space and thus lose the meaning of signature extraction process. Since the length of group codes of every vocabulary is different, the length of each partitioned sub-string from every vocabulary is not the same. A signature extraction algorithm is designed as below to extract the signature from each partitioned sub-string. The signature extraction version of dictionary is then composed by all signatures which are extracted from all group coded vocabularies. Our paper designs the signature extraction process as the following algorithm shown in Fig. 3. Example 2. is a simple illustration that explains the progress of the signature extraction algorithm.

Example 2: Let TX = "001001100110110101", set the number of partitions C = 4.

- Step 1: |TX| = 18.
- Step 2: $d = (4 - (|TX| \bmod 4)) \bmod 4 = 2$,
TX = "001001100110110101" + "00"
= "00100110011011010100".
- Step 3: TX is evenly partitioned into four sub-strings, which are respectively: TX₁ = "00100", TX₂ = "11001", TX₃ = "10110" and TX₄ = "10100".
- Step 4: Determine the number of occurrences of 0 and 1 separately for the four sub-strings in TX and take the one with more occurrences as signature code for each sub-string. If the number of occurrences of 0 and 1 are the same, 1 is preferred. The result is K = "0 1 1 0".
- Step 5: Signature of TX is "0 1 1 0".

Figure 4 shows how a dictionary signature file is generated from the incorporation of signature of every vocabulary extracted from the signature extraction process and organizes these signatures containing the vocabulary to establish the index. When the signature for which the vocabulary the user wants to look up exists in the dictionary signature file, the user can retrieve from the established index according to the retrieved signatures; in this case, the overall dictionary searching efficiency is improved.

RESULTS AND DISCUSSION

The experimental data is collected from a medical dictionary that contains 142,709 vocabularies, 2,156,957 characters and the average length of characters of every vocabulary is 16. Determination of number of partitions, C, can be made by the size of dictionary or the average length of vocabulary considering the size of which the dictionary signature file is to be sent to the user over the Internet is appropriate. Here C ranges from 8 to 16 and the observation of various partitions which result in different outcomes is shown as Table 3 in the following. We randomly pick 3000 vocabularies for our experiment where 2000 vocabularies do exist and the other 1000 vocabularies do not exist in the dictionary, then we take the average value for these two circumstances of retrievals.

From Table 3, when C gets bigger, not only will the practical range to be searched be much less with bigger number of signatures indices but will the time be much less also and the search time is not even 1 millisecond. Yet this brings up another disadvantage-the memory space occupied is increased. How to determine the size of C value for the best coordination between the number of signature indices and memory space occupied is the critical part to which we should pay special attention in our method. From these experimental results, we obtain the best value for C is 14 (shadow part) considering that this value results in little memory space occupied under the minimum number of words searched and the number of words searched only take up 0.09% of overall vocabulary in the dictionary as well as that the search time is below 1 millisecond. Another reason for picking C = 14 is due to that 60% of the average length of vocabulary in the dictionary is between the length of characters being 7 and 16. When the number of partitions is set as 14, the partition point for the vocabulary exactly locates in the space between the characters to extract exactly the signature in each character so as to enhance the representation to which the index it belongs and the overall search efficiency is increased hence.

In addition, Table 4 is a comparison of our method with previously mentioned methods. These methods

Table 3: The impact of size of C value chosen

C	No. of signature indices	Memory space occupied (bits)	Average No. of words searched	Average search time (10 ⁻³ sec)
8	256	2048	1212	4.63
9	512	4608	728	3.72
10	1024	10240	610	2.31
11	1928	21208	466	1.48
12	3881	46572	304	1.42
13	5769	74997	259	1.25
14	11345	158830	124	0.78
15	15364	230460	116	0.76
16	25103	401648	107	0.67

Table 4: Summary comparison of four methods

Method	Average dictionary size searched (# of vocabularies)	Dictionary search (%)	Average search time (10 ⁻³ sec)	False drop probability (%)
Partition by string length	2159	1.51%	14.679	33.33
Partition by first letter	6089	4.27%	39.691	33.33
Tow-level indexing	647	0.45%	3.656	32.49
Two-level dictionary retrieval (C = 14)	124	0.09%	0.782	15.42

include Partition by String Length, Partition by First Letter, Two-Level Indexing and Two-Layer Dictionary Retrieval (C = 14). Partition by String Length and Partition by First Letter are commonly seen dictionary classification methods. Two-Level Indexing by Owolabi (1996) shows the best efficiency. That's why we choose these three methods as subjects of comparison with our two-layer dictionary retrieval. In the comparison, we compare the average number of words searched, average search time and even the false drop probability for the four methods under the circumstance that the vocabulary to be retrieved does not exist in the dictionary to see whether it's effective in filtering the dictionary and reducing the search range. The formula for false drop probability is shown as follows:

$$\text{False drop probability (\%)} = \frac{\text{Number of false drops}}{\text{Number of total executions}} * 100\%$$

From Table 4 it may concluded that present method is far better than other three dictionary search methods considering the average range to be searched in the dictionary is less than 0.1% and average search time is not even 1 m sec. As for false drop probability, our method shows a lower value than others. Even when a false drop circumstance occurs, the range to be searched is greatly decreased because of better filtering capability so as to reduce the search time over the Internet.

CONCLUSIONS

In this study, we proposed a dictionary retrieval system appropriate for current distributive networks-Two-Layer Dictionary Retrieval. Through the signature extraction concept, the dictionary retrieval over the Internet will not need as much space for the indexing so

the file to be transferred over the Internet is small, reduce the false drop probability and even reduce the range to be searched in the dictionary by means of dictionary filtering capability for much better overall efficiency of dictionary retrieval in terms of Internet application.

REFERENCES

Angell, R.C., G.E. Freund and P. Willett, 1983. Automatic spelling correction using a trigram similarity measure. *Inform. Process. Manage.*, 19: 255-261.

Burkhard, W.A. and R.M. Keller, 1973. Some approaches to best-match file searching. *Comm. ACM*, 16: 230-236.

Doster, W., 1977. Contextual postprocessing system for cooperation with a multiple-choice character-recognition. *IEEE Trans. Comp.*, 26: 1090-1101.

Hall, P.A.V. and G.R. Dowling, 1980. Approximate string matching. *ACM Comp. Sur.*, 12: 381-402.

Kohonen, T., 1987. *Content Addressable Memories*. Springer-Verlag, Berlin.

Morgan, H.L., 1970. Spelling correction in system programs. *Comm. ACM*, 13: 90-94.

Owolabi, O., 1996. Dictionary organizations for efficient similarity retrieval. *J. Sys. Software*, 34: 127-132.

Owolabi, O. and D.R. McGregor, 1988. Fast approximate string matching. *Software Prac. Exp.*, 18: 387-393.

Salton, G. and A. Wong, 1978. Generation and search of clustered files. *ACM Trans. Databases Sys.*, 3: 321-346.

Szanzer, A.J., 1969. Error-correcting methods in natural language processing. *Inform. Process.* 68, IFIP, pp: 1412-1416.

Szanzer, A.J., 1973. Bracketing techniques in elastic matching. *Comp. J.*, 16: 132-134.