

# OKSAT at NTCIR-11 Temporalia

## - Plural Sets of Search Terms for a Topic -

Takashi SATO  
 Information Processing Center  
 Osaka Kyoiku University  
 Kashiwara Osaka Japan  
 +81-72-978-3823  
 sato@cc.osaka-kyoiku.ac.jp

Shingo AOKI  
 Graduate School of Education  
 Osaka Kyoiku University  
 Kashiwara Osaka Japan  
 +81-72-978-3823  
 aoki@ss.osaka-kyoiku.ac.jp

### ABSTRACT

Our group submitted runs for Temporal Information Retrieval (TIR) subtask of NTCIR-11 Temporalia. For one of our runs, we prepared plural sets of search terms for a subtopic. Analyzing experimental results, we observe the effectiveness of using plural sets of search terms for a subtopic.

### Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval – *Information filtering, Query formulation, Retrieval models, Search process, Selection Process.*

### General Terms

Experimentation, Performance, Measurement.

### Team Name

OKSAT

### Subtasks

Temporal Information Retrieval (TIR)

### Keywords

Information retrieval, Temporal search, Plural sets of search terms, Gram base index.

## 1. INTRODUCTION

The need of information retrieval including temporal information is increasing nowadays. NTCIR-11 Temporal Information Access (Temporalia) [1] focuses on this problem. Our group participated in Temporal Information Retrieval (TIR) subtask of NTCIR-11 Temporalia. In this paper, we describe our system and techniques used. Following experimental results, we show an example to discuss the effectiveness of our methods.

## 2. OUR APPROACH

### 2.1 Removal of Tags from Corpus

From temporalia corpus [2], we extracted the text surrounded by title tag (`<tag name="title">...</tag>`) in `<meta-info>` tag. Using 'temporalia\_solrify.pl' prepared by task organizer, the text surrounded by text tag (`<text>...</text>`) in which all tags were removed was extracted. In addition, the val of T tags (`<T val="...">`) in the text tag were extracted also.

### 2.2 How to Make Search Terms

We prepared the following three types of search terms.

- (1) From topic file, we extracted words from title and each subtopic of the topic, and then they were filtered by stop word list.
- (2) We searched the internet (Wikipedia and Google) by words from (1), and then we simply added most common words in the search results.
- (3) We classified words from internet search of (2) if possible. Then we added each classified group to (1). As a result, we got plural sets of search terms for a subtopic.

### 2.3 Searching and Scoring

We used only `<text>` tag among tags extracted by 2.1. Because titles in `<title>` tag are short and we observed that search words were not often used in the title even if a document related to the topic. About `<T>` tags in `<text>`, we could not use these values in order to distinguish time factor of subtopics.

Using index made by text in `<text>` tag, we retrieved the search terms (1), (2) and (3) prepared by 2.2. Then we scored and ranked output document id by using probabilistic model [3].

## 3. STRUCTURE OF GRAM BASE INDEX

Gram base indices are known as index structures, which enable arbitrary string search. Grams consist of strings (character sequences), which start every character in a text. The length of strings is less than 3 (such as 1-gram or 2-gram) in common cases. Our grams are longer than 4-grams in average by encoding grams in  $w_g$  byte.  $w_g$  is set 6 in this task. At first, characters are coded in varying length bit in accordance with their frequency of appearance. Then they are stuffed within  $w_g$ . Figure 1 shows the stuffing of coded characters.

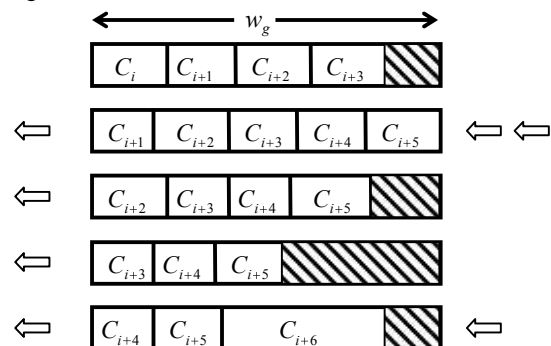


Figure 1. Stuffing of Coded Characters.

Table 1 shows code table of character we used. The column 'C' is the character encoded; 'L' is the bit length of encoding. The

column 'CODE' is the code in hexadecimal. In this table, '\_' stands for a space character. The lowercase is converted to uppercase. Other than numeric and alphabet letters are converted into a space.

**Table 1. Code Table for Gram Coding**

C	L	CODE	C	L	CODE
_	11	0	I	4	3
0	8	1	J	9	1
1	8	41	K	7	21
2	8	81	L	5	d
3	8	82	M	5	11
4	8	83	N	4	2
5	8	84	O	4	7
6	8	85	P	6	3
7	8	86	Q	11	1
8	8	87	R	4	d
9	8	40	S	4	5
A	4	1	T	3	5
B	6	30	U	5	19
C	5	9	V	7	1
D	5	c	W	6	31
E	3	7	X	8	80
F	6	2	Y	6	11
G	6	1	Z	10	1
H	4	9			

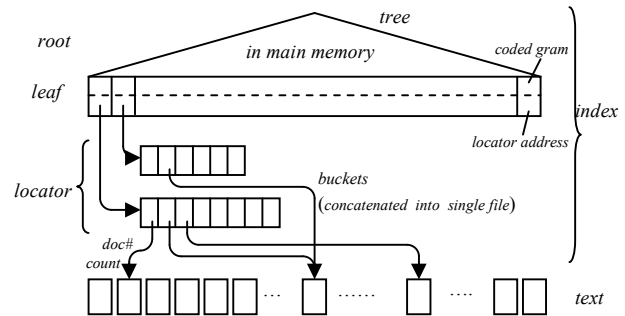
Figure 2 shows an example of gram coding (each row is gram length in bit and characters in a gram) when text is 'KFWB\_NEWS\_980\_' a part of the title of the first document of the corpus.

```

43: KFWB_NE
42: FWB_NEW
40: WB_NEWS
34: B_NEWS
39: _NEWS_
47: _NEWS_TALK
43: EWS_TALK
40: WS_TALK
45: S_TALK_
41: _TALK_
46: TALK_98
43: ALK_98
47: LK_980
42: K_980
46: _980_
35: _980_
27: 80_
19: 0_
11: -
    
```

**Figure 2. Example of Grams.**

Figure 3 shows a data structure of gram index [4][5][6]. An index has three parts called *root*, *leaf* and *locator*. Grams are sorted and stored in secondary storage as leaf. There is a pointer from a gram in leaf to a bucket, which stores the document numbers where the string corresponding to the gram is found and the count of the gram appeared in the document. Locator, which is stored in secondary storage, is collection of buckets. Root, which is put in main memory when searching, is a wide range map of grams.



**Figure 3. Structure of Gram Index.**

The search algorithm is explained in terms of three cases according to the relation between the length of search key word ( $l_k$ ) and gram length ( $l_g$ ). Figure 4 (a), (b) and (c) show how to follow the pointers in leaves and locators when  $l_k = l_g$ ,  $l_k < l_g$  and  $l_k > l_g$  respectively. Since the buckets of the locator are stored sequentially, they are drawn in one box and separated by double lines.

Our index has tree structure, which has sorted gram and wide range map of them. So, not only search terms whose length is equal to gram length, but also shorter or longer terms can be searched efficiently. When we search a longer term, every gram in the term is searched. Then retrieved sets of document numbers are intersected.

## 4. EXPERIMENTAL RESULTS

### 4.1 Indexing

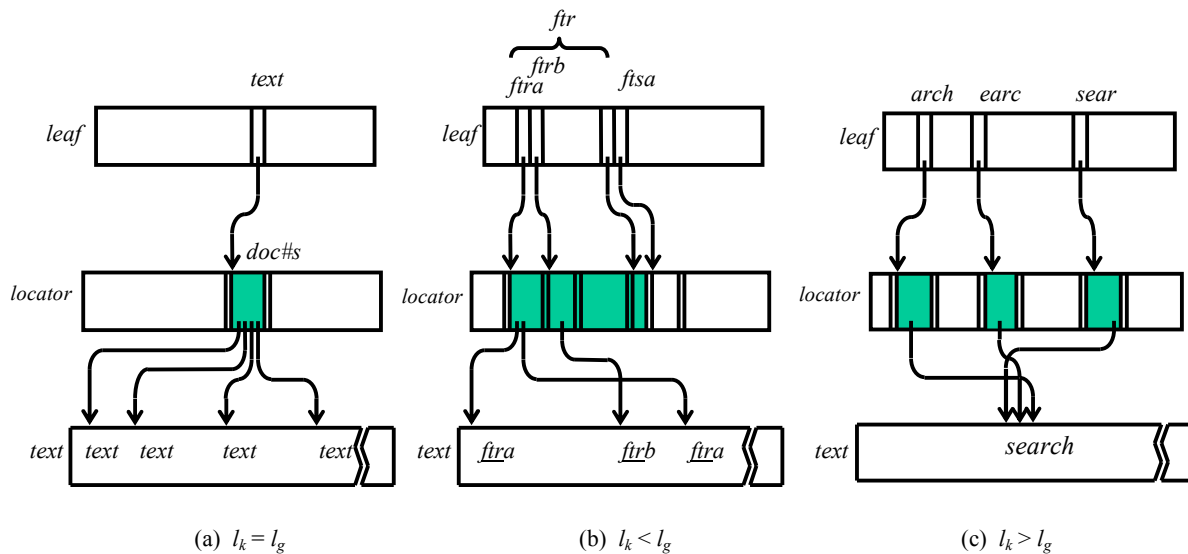
We made gram based indices from text surrounded by <text> and <title> tag of corpus. Table 2 shows specifications of computer we used. And Table 3 shows statistics of indices.

**Table 2. Specifications of Computer**

<b>CPU</b>	Intel Core i5-4430@3.0GHz, 4C/4T
<b>MEM</b>	8GB, DDR3-1600
<b>O S</b>	FreeBSD 8.4, 64bit
<b>HDD</b>	1TB, SATA 6GB/s, 64MB Cache

**Table 3. Statistics of Indices**

	title	text
data size (GB)	0.170	8.90
index size (GB)	0.489	18.9
time (min.)	2.08	153


**Figure 4. Index Search.**

For val of  $\langle T \rangle$  tag in  $\langle \text{text} \rangle$ , we made a word base index using Berkeley NDBM library.

## 4.2 Runs

Using three types of search terms of 2.2 we made the following runs.

- OKSAT-TF01 : type (1) of 2.2
- OKSAT-TF02 : type (2) of 2.2
- OKSAT-TF03 : type (3) of 2.2

Table 4 shows time (searching and scoring in minutes) and AP (mean average precision) of our submitted runs. Because task organizers have a shallow pooling at document 20,  $P@20$  and  $nDCG@20$  are shown in this table also.

**Table 4. Time and MAP of Submitted Runs**

RUN ID	Time	AP	$P@20$	$nDCG@20$
OKSAT-TF01	54	0.2138	0.5828	0.4566
OKSAT-TF02	83	0.2063	0.5843	0.4551
OKSAT-TF03	25	0.1971	0.5723	0.4391

The CPU of our computer has 4 cores and can process 4 threads simultaneously. OKSAT-TF01 and OKSAT-TF02 were executed simultaneously by 2 threads each. On the other hand OKSAT-TF03 was executed alone by 4 threads. So, OKSAT-TF03 was executed twice as fast as OKSAT-TF01 and OKSAT-TF02.

As a result, AP and  $nDCG@20$  dropped slightly by the addition of search terms.

## 4.3 Temporal Class Analysis

Table 5 shows  $P@20$  of each temporal query classes. The atemporal query class was better than other query classes. The past query class was difficult for our group.

**Table 5.  $P@20$  of Each Temporal Query Classes**

RUN ID	atemporal	future	past	recency
OKSAT-TF-01	0.6460	0.5740	0.5050	0.6060
OKSAT-TF-02	0.6260	0.6030	0.5070	0.6010
OKSAT-TF-03	0.5830	0.6190	0.4870	0.6000

The future class of OKSAT-TF-02 and OKSAT-TF-03 were better than that of OKSAT-TF-01. We added ‘future’ to search terms of a future class query of OKSAT-TF-02 and OKSAT-TF-03 when those of OKSAT-TF-01 did not have ‘future’.

## 4.4 Topic-based Analysis

We show poorly and better performing topic examples about expansion of search terms (i.e. from OKSAT-FT-01 to OKSAT-FT-02).

026a: We expanded search term from ‘passive smoke’ to ‘smoke’ and ‘smoking’. This made  $P@20$  fall down from 0.9500 to 0.5500.

045f: We expanded search term from ‘Papacy’ to ‘Pope’. This raised  $P@20$  from 0.0500 to 0.4000.

However, the effect of search term expansion was not similar about other temporal classes. So the effect of word expansion was sensitive to temporal classes.

## 4.5 Examples of Plural Sets of Search Terms

In temporal information retrieval, there are cases when plural events should be searched. Searching by only one set of search terms is not enough, in those cases. So, we made our system to handle plural sets of search terms for a subtopic. Concretely, we explain an effect of plural sets of search terms about subtopic id 001p (subtopic of type past of topic id 001). We added the words

extracted from title of topic id 001 and its subtopic 001p to the words after '+' below.

001p-1 (+ 'Tokyo subway,' Tokyo', 'subway', 'sarin')

001p-2 (+ 'Hiroshima nuclear')

001p-3 (+ 'Tohoku earthquake')

Then three sets of search terms were made. Table 6 shows relations of the term sets above and the AP (average precision), P@20 and nDCG@20 of their runs.

We merged these three runs into one by two ways. One is merging by score order of document and the other is merging by rotation (first we gather top of three runs, next we gather second of three runs and so on.)

Table 7 shows merge type and AP, P@20 and nDCG@20. In this table, 'base words' stands for no words from the internet added i.e. words extracted from title of topic id 001 and its subtopic 001p only. From Table 7, we observe that the AP, P@20 and nDCG@20 of merging by rotation is higher than those of base words only and merging by rotation.

**Table 6. Term Set and AP**

Term Set	AP	P@20	nDCG@20
001p-1	0.0807	0.2355	0.2328
001p-2	0.0911	0.3009	0.3716
001p-3	0.1849	0.5319	0.4958

**Table 7. Merge Type and AP**

Merge Type	AP	P@20	nDCG@20
base words	0.1849	0.5139	0.4958
by score	0.1368	0.4205	0.4465
by rotation	0.2669	0.6362	0.5589

## 5. CONCLUSIONS

Our group submitted three runs for Temporal Information Retrieval (TIR) subtask of NTCIR-11 Temporalia). In third run, we prepared plural sets of search terms for a subtopic using words added by the internet search. Analyzing experimental results, we observe the effectiveness of using plural sets of search terms for a subtopic.

## 6. REFERENCES

- [1] H. Joho, A. Jatowt, R. Blanco, H. Naka, and S. Yamamoto, Overview of NTCIR-11 Temporal Information Access (Temporalia) Task, in *Proceedings of the NTCIR-11 Conference*, Tokyo, Japan, 2014.
- [2] M. Matthews, P. Tolchinsky, R. Blanco, J. Atserias, P. Mika, and H. Zaragoza, Searching Through Time in the New York Times, in *Proceedings of the 4th Workshop on Human-Computer Interaction and Information Retrieval*, NJ, USA, pp.41-44, 2010.
- [3] S.E. Robertson and S. Walker, Some simple effective approximations to the 2-Poisson model for probabilistic weighted retrieval, in *Proceedings of the 17th International Conference Research and Development in Information Retrieval*, pp. 232-241, 1994.
- [4] T. Sato, Fast full text retrieval using gram based tree structure, in *Proceedings of the ICCPOL '97*, Vol.2, pp.572-577, 1997.
- [5] T. Sato and K. Han, NTCIR-3 CLIR Experiments at Osaka Kyoiku University - Compression of Gram-based Indices -, in *Proceedings of the NTCIR-3*, Tokyo, December 2002.
- [6] T. Sato, T. Satomoto, and K. Han, NTCIR-3 PAT Experiments at Osaka Kyoiku University -Long Gram-based Index and Essential Words -, in *Proceedings of the NTCIR-3*, Tokyo, December 2002.