

GIS KRISTAL-II

# KRISTAL - II

1998. 7.

# 1. (Vector Space Model)

- $d$ : (document)
- $f_{d,t}$ :  $t$  (tf)
- $f_t$ :  $t$  (df)
- $M$ : 가 (Measure)
- $N$ : collection
- $n$ : collection
- $t$ : (term)
- **TF\*IDF** : 가 \* (Term Frequency \* Inverse Document Frequency)
- $w_{d,t}$ : - 가 (document - term weight)
- $w_{q,t}$ : - 가 (query - term weight)
- $w_t$ : 가 (term weight)

$d$	Document $D_d$
1	apple balloon balloon elephant apple apple
2	Chocolate balloon balloon chocolate apple chocolate duck
3	Balloon balloon balloon balloon elephant balloon
4	Chocolate balloon elephant
5	Balloon apple chocolate balloon
6	Elephant elephant elephant chocolate elephant

1-1. ( 6, 5)

## Inner product similarity

1-2(a) 1-1 Inner product (binary vector)

1-2(b) inner product 2

$$M(\text{"chocolate, duck"}, D_2) = (0,0,1,1,0) \cdot (1,1,1,1,0) = 2.$$

가

1.  $D_6$  elephant 4 1



the -

t 가  $w_t$   $f_t$  N

$$w_t = \log(N/f_t)$$

가  $w_t$  t  $f_{d,t}$

$$w_{d,t} = f_{d,t} * w_t = f_{d,t} * \log(N/f_t)$$

가 \* (Term Frequency \* Inverse Document Frequency)

**TF\*IDF**

1-1 duck elephant duck 6 1

$$w_{duck} = \log(6/1) = 2.58$$

elephant 6 4

$$w_{elephant} = \log(6/4) = 0.58$$

가 가 . ( 가  
?)

가 가 가  $D_d$   $f_{d,t}$   
가 가  
가 . 가 .

$$M(Q, D_d) = \frac{1}{|D_d|} \sum_{t=1}^n W_{q,t} \bullet W_{d,t}$$

$$|D_d| = \sum_i f_{d,i}$$

$D_d$

## (Vector Space Model)

가 .  
 가 .  
 가 가 . (vector)  
 , .  
 (cosine measure) . (cosine)  
 , 가 0 1, 0  
 가 .

$$\cos(Q, D_d) = \frac{Q \bullet D_d}{|Q||D_d|} = \frac{1}{W_q W_d} \sum_{t=1}^n w_{q,t} \bullet w_{d,t}$$

.  $d$   $W_d$  가

$$W_d = \sqrt{\sum_{t=1}^n w_{d,t}^2}$$

가  $W_q$

$$W_q = \sqrt{\sum_{t=1}^n w_{q,t}^2}$$

. **TF\*IDF** ( $w_{d,t} = f_{d,t} * \log(N/f_t)$ )

$$\cos(Q, D_d) = \frac{1}{W_q W_d} \sum_{t=1}^n w_{q,t} \bullet f_{d,t} \log \frac{N}{f_t}$$

가 .  $f_{q,t}$  1, 0

$$\cos(Q, D_d) = \frac{1}{W_q W_d} \sum_{t \in Q} f_{d,t} \left( \log \frac{N}{f_t} \right)^2 \dots \dots \dots (1)$$

(1)  $W_q$  가

1-1

(a)

$d$	$t$					$W_d$
	a	B	C	d	e	
1	3	2	0	0	1	3.10
2	1	2	3	1	0	3.31
3	0	5	0	0	1	1.42
4	0	1	1	0	1	0.86
5	1	2	1	0	0	1.27
6	0	0	1	0	4	2.39
$F_t$	3	5	4	1	4	
$w_t$	1.00	0.26	0.58	2.58	0.58	

(b)

$d$					
	D $W_q = 2.58$	C $W_q = 0.58$	c, d $W_q = 2.64$	a, b, e $W_q = 1.18$	a, b, c, d, e $W_q = 2.90$
1	0.00	0.00	0.00	0.95	0.39
2	0.78	0.53	0.88	0.29	0.92
3	0.00	0.00	0.00	0.40	0.16
4	0.00	0.67	0.15	0.40	0.30
5	0.00	0.46	0.10	0.76	0.40
6	0.00	0.24	0.05	0.48	0.24
	2	4	2	1	2

1-3. 가 : (a)  $f_{d,t}$  가  $W_{d_i}$  (b) 가 1-1

(1) DB 가 CPU 가?  
 가 . 가

$f_t$  ,  $f_{d,t}$  .  $W_q$   
 가 . (1)  $W_d$  .  $W_d$   
 가 .  
 (accumulator) 가 1-1  
 $\cosine(Q, D_d)$   $r$  가 .

1-1.  $r$

Cosine measure  $r$

1.  $A \leftarrow \{ \}$  . ( $A = 0$  )
2.  $t$ 
  - a.  $t$
  - b.  $t$   $f_t$   $I_t$
  - c.  $w_t \leftarrow (\log(N/f_t))^2$
  - d.  $I_t$
  - e.  $f_{d,t}$   
 $A = A_d$   
 $A_d \leftarrow 0$   
 $A \leftarrow A + \{A_d\}$   
 $A_d \leftarrow A_d + f_{d,t} * w_t$
3.  $A_d$   
 $A_d \leftarrow A_d / W_d$  ( $W_q$  ,  $A_d = \cosine(Q, D_d)$ )
4.  $1 \leq i \leq r$ 
  - a.  $A_d = \max\{A\}$   $d$
  - b.  $d$  가 .
  - c.  $A \leftarrow A - \{A_d\}$

가  $W_d$  가  $f_{q,t}$  가 1, 0)가 .

$r(r \ll M)$  가  $1-1$  4 ,  $W_d$  ,  $A$  ,  $A$  .

가  $W_d$

1 가  $W_d$  4Mb(= sizeof(int) \* 1,000,000 = 4byte \* 1M = 4Mb) 가 . , 가 .

가 가 .

가 Cosine 가 가 . (CPU, memory) !

$A$

$f_t$   $f_t$  가 가 ( $w_t$ 가 )  $N$  가 . cosine cosine

가 가 . 가 가 . cosine (continue)

. (continue) quit 가

.) Continue 가 5,000 1,000 , 가 30,000 .





## 2. KRISTAL-II Vector Space Model

가

$$M(Q, D) = (\sum F_{TF} \cdot F_{IDF}^2) / F_{DL}$$

KRISTAL-II

가

(*FTF*),

(*FIDF*),

(*FDL*)

가

method.h

3.

### *FTF*

1.

$$FTF = tf$$

2. Harman 1986

$$FTF = \log(1 + tf)$$

### *FIDF*

1.

$$FIDF = \log(N/df)$$

2. Witten et al. 1994

$$FIDF = \log((N + 1)/df)$$

3. Croft and Harper 1979

$$FIDF = \log((N - df)/df)$$

---

3

가

KRISTAL-II

KRISTAL-II

가

3 4

Cosine

Cosine

4. Sparck Jones 1972

$$FIDF = \log(N/df) + 1$$

**FDL**

1.

$$FDL = Wd$$

$$Wd$$

2.

$$FDL = \log(Wd)$$

3.

$$FDL = Ud$$

$$Ud$$

4. Harman 1986

$$FDL = \log(Ud)$$

**가**

4,142 BLUE  
 2 , 4 , 2 4 가  
 M(method )

M242 M244가 가  
 2-1  
 4 가 M242

M244  
 “Unique Term Conunt” 가 -

4 4,142 BLUE 2  
 (OR) 2 ,  
 4 , 2 1 , 1 , 1  
 { , }, { , }, { , } 28%, 67%, 133%  
 가 가 2  
 KRIST set Recall Precision

$$M(Q, D) = (\sum \log(1 + tf) \cdot (\log(N/df) + 1)^2) / \log(Wd)$$

가

method.h

가

가

가

KRIST Set

(a) Precision

	{ , }	{ , }	{ , }
M111	18/50=0.36	36/100=0.36	12/50=0.24
M242	23/50=0.46	60/100=0.60	28/50=0.56
(M242/M111)	1.28	1.67	2.33

(b) Recall

	{ , }	{ , }	{ , }
M111	18/23=0.78	36/71=0.51	12/33=0.36
M242	23/23=1.00	60/71=0.85	28/33=0.85
(M242/M111)	1.28	1.67	2.33

2-1. 가

Precision(P<sub>50</sub>)    50    Precision(P<sub>100</sub>)    100  
 Recall(R<sub>50</sub>)    23    Recall(R<sub>100</sub>)    71

```

int          SetNum, MemCnt; /*          Set          DocID */
T_RESULT     DocIDList;
USERSECLIST  SecList;
USERTERMLIST TermList;

/* ..... DB          ..... */
while(          ) {
    VECTOR_SecList_Init( &SecList );
    for (          ) {
        VECTOR_SecList_Insert( &SecList,          ,          가          );
    }

    VECTOR_TermList_Init( &TermList );
    for (          ) {
        VECTOR_TermList_Insert( &SecList,          ,          가          );
    }

    FIRE_VectorRank( SecList, TermList, &SetNum, &MemCnt );

    DocIDList = (T_RESULT *) malloc( MemCnt * sizeof(T_RESULT) );
    FIRE_GetDocIDList( SetNum, DocIDList );
    /* DocIDList          . */
}
/* ..... DB          ..... */

```

## KRISTAL - II

1. ( , )
2. ( )
3. DocID
4. DocID
5. (tf)가 가
6. MAX\_SORT\_SIZE DocID

- 7. 가
- 8. MAX\_SET\_SIZE
- 9. MAX\_SET\_SIZE 가
- 10. DocIDList

KRISTAL - II

가 . ,  
 5 .

(USERSECLIST)

USERSECLIST  
 MAX\_SEC\_NUM

```
#define MAX_SEC_NUM 50 /* 가 */
#define MAX_TERM_NUM 100 /* 가 */
#define MAX_TERM_LEN 100 /* */
```

```
typedef struct
{
    char Name[MAX_TERM_LEN]; /* */
    double Wgt; /* 가 , default = 1.0 */
} USERSECINFO;
```

```
typedef struct {
    int Cnt;
    USERSECINFO Sec[MAX_SEC_NUM];
} USERSECLIST;
```

## API

1. `VECTOR_SecList_Init( USERSECLIST SecList );`

`USERSECLIST SecList` (SecList.Cnt = 0) .

2. `VECTOR_SecList_Insert( &SecList, char * , double 가 );`

`SecList` 가 가 . 1 가 .  
가  
가 1.0 .

3. `VECTOR_SecList_Display( &SecList );`

, 가 .

## (USERTERMLIST)

`USERTERMLIST` .  
`MAX_TERM_NUM` .

`typedef struct`

```
{  
    char term[MAX_TERM_LEN]; /* */  
    double Wgt; /* 가 , default = 1.0 */  
} USERTERMINFO;
```

`typedef struct {`

```
    int Cnt; /* */  
    USERTERMINFO Term[MAX_TERM_NUM]; /* */  
} USERTERMLIST;
```

## API

1. `VECTOR_TermList_Init( USERTERMLIST TermList );`

`USERTERMLIST TermList` (TermList.Cnt = 0) .

2. `VECTOR_TermList_Insert( &TermList, char * , double 가 );`

`TermList` 가 가 . 1 가 .

3. `VECTOR_TermList_Display( &TermList );`

, 가 .

```

가
가 SECTIONINFO
vs_FillSecInfo() . 가 SECTIONINFO
2 , DocID

```

```

typedef struct
{
    char SecName[MAX_TERM_LEN]; /* */
    int IndexType; /* */
    double Wgt; /* 가 */

    int DIDCnt; /* DocID */
    R_POSTINFO *RPostList; /* DocID List */
} SECTIONINFO;

```

```

가 vqTERMINFO vs_FillTermInfo()
. vqTERMINFO 2 ,
( ) .

```

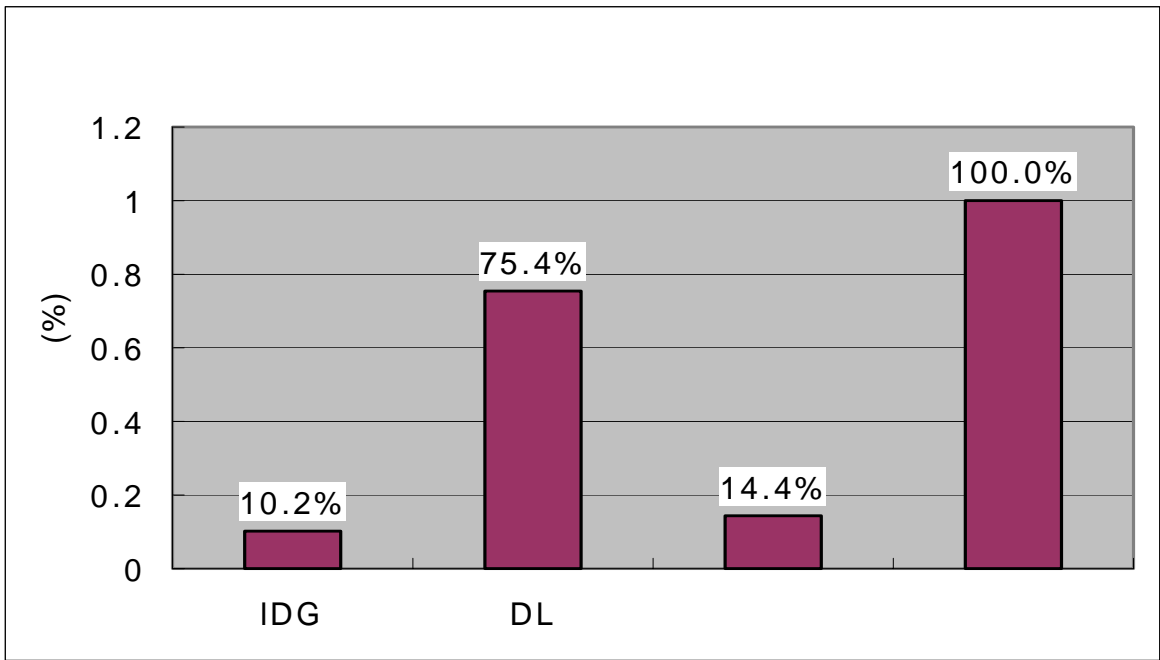
```

typedef struct
{
    char term[MAX_TERM_LEN]; /* index term itself */
    int IndexType; /* Index type for this term */
    int df; /* document frequency */
    T_SEARCHOP op;
    double idf; /* Inverse Document Ferequency (IDF) */
    double idf_2; /* (IDF)^2 */
    double Wgt; /* Term Weight for Relevance feedback */

    int SecDIDCntMax; /* max size of RPostList retrieved each section */
    int DIDCnt; /* Total Size of DocIDList for this term */
    R_POSTINFO *DocIDList; /* DocID list containing this term */
} vqTERMINFO; /* vector query term information */

```





2-1.

). IDG:

DL:

( DocID

가

가

**DocID (vs\_GetDocIDList())**

SecInfo

RPostList

가

DocID

가

while ( i)

SecInfo[i].RPostList <-

i

DocID

가

가

DocID

가

가

가

DocID

10%

" ", " ", " ", " "

가

DocID

33

40%

60%

가

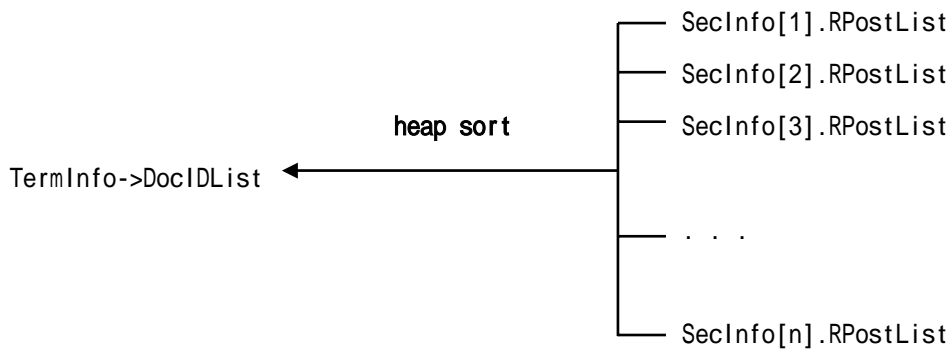
DocID

RPostList

TermInfo

(vs\_MergeSecDIDList\_Heap())

heap sort



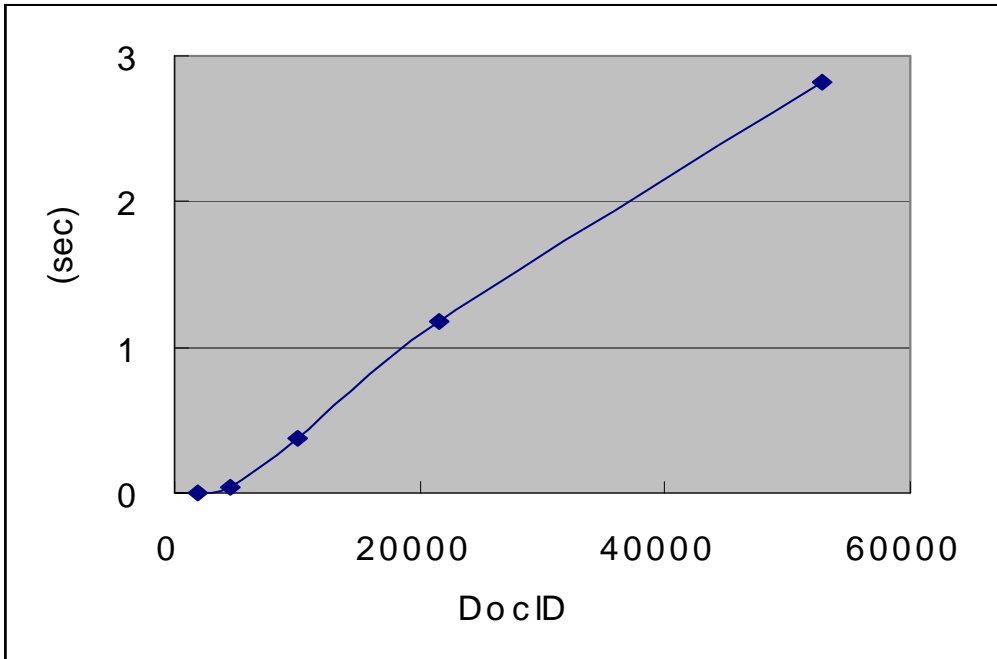
(tf) 6  
 가 DocID (vqTERMINFO) 가 가  
 , IDF , . 가 가  
 IDF .  
 가 DocID  
 DocID  
 df 가 가 DocID DocID  
 . 가 DocID 가  
 DocID DocID  
 가 .  
 가 .

**DocIDList (Accumulation)**

DocIDList 2가 .  
 DocID MAX\_SORT\_SIZE  
 가 ( 2-2  
 MAX\_SORT\_SIZE ) .  
 MAX\_SET\_SIZE 가

---

6 2 heap sort ( DocIDList  
 (accumulator) ),  
 C 가 qsort ( 가 가 DocIDList  
 ). 2-2 quick sort .



2-2. 가 qsort  
 Sun SparcCenter200E( 1.2GB) DocID 5000  
 0.4 , 60000 3 가  
 Sort (qsort() ) Quick

. MAX\_SORT\_SIZE MAX\_SET\_SIZE  
 DocID

(DocLen)  
 75.4% ( 2-1 ) . TF IDF 가  
 KRISTAL-II UNIX

가

가

DocID

" ", " ", " "

가

DocID 가 1 , MAX\_SORT\_SIZE DocID  
MAX\_SET\_SIZE

가

가

가

가

가

MAX\_SORT\_SIZE  
Pruning continue  
가

MAX\_SORT\_SIZE DocID가  
가

MAX\_SET\_SIZE DocID 가

34

(

450MB)

가

7

1 3 ( 2.58 )

/

"

"

10

DocID 12 x10 = 1.2MB

2 2.4MB가

MAX\_SORT\_SIZE DocID

MAX\_SORT\_SIZE가 10,000

120kb가

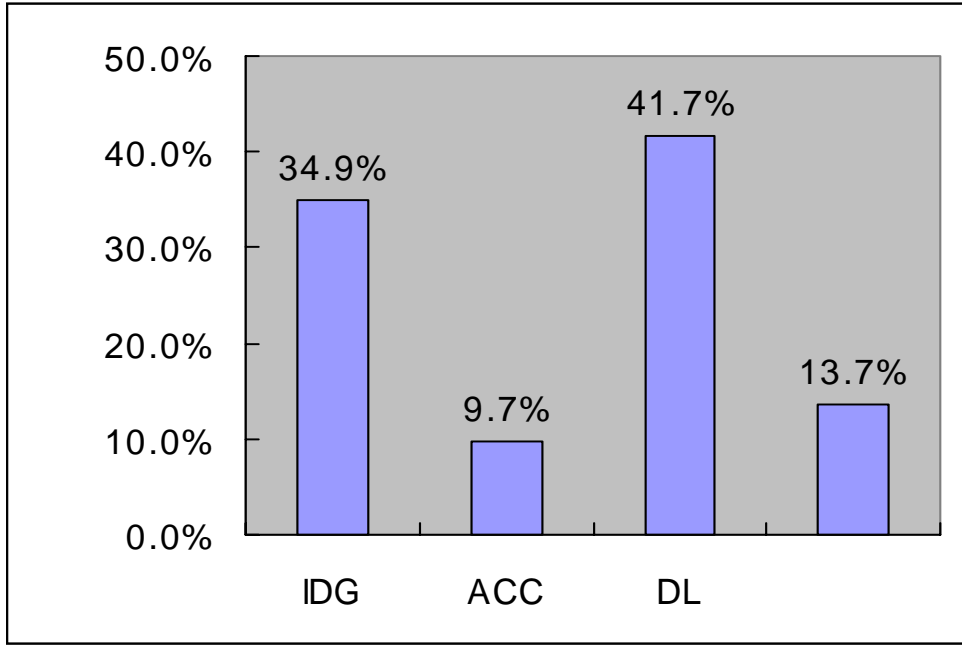
PC

가

DB 가 344,869

가

78,649



2-3. KRISTAL-II

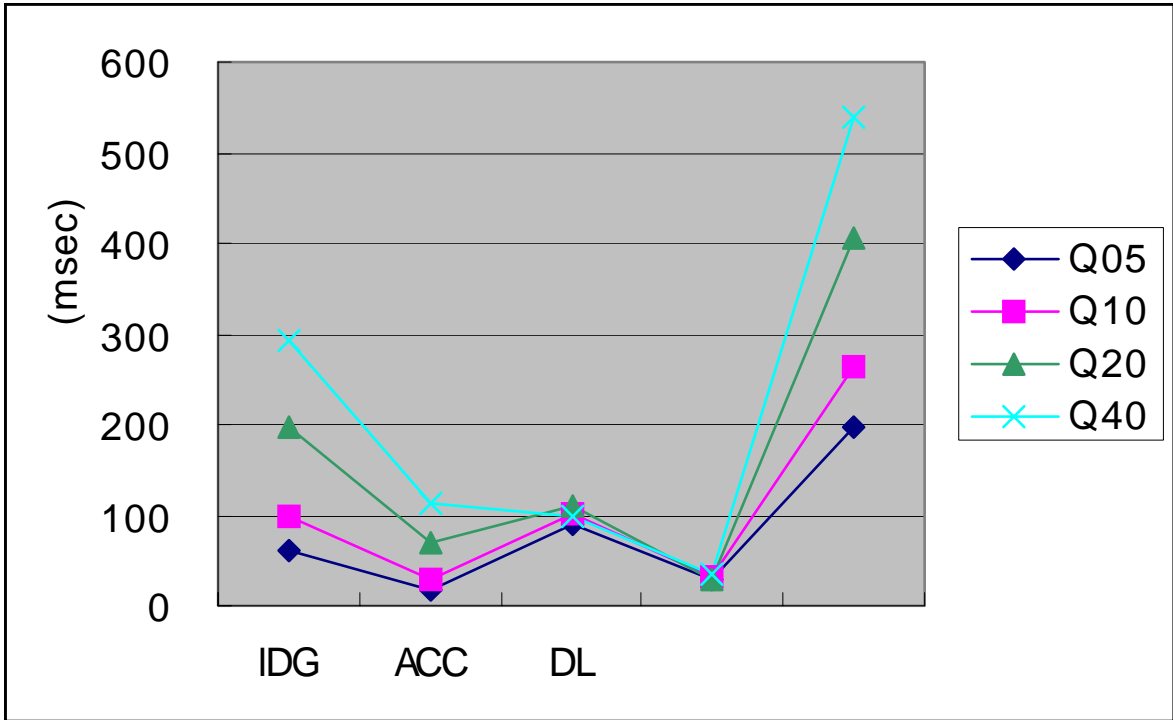
IDG: DocID, ACC: , DL: , :IDG, ACC, DL  
 . 20 , IDG가 0.90 , ACC가 0.25 , DL  
 2.58 , 1.08 , 0.35 가 ( ).

2-3

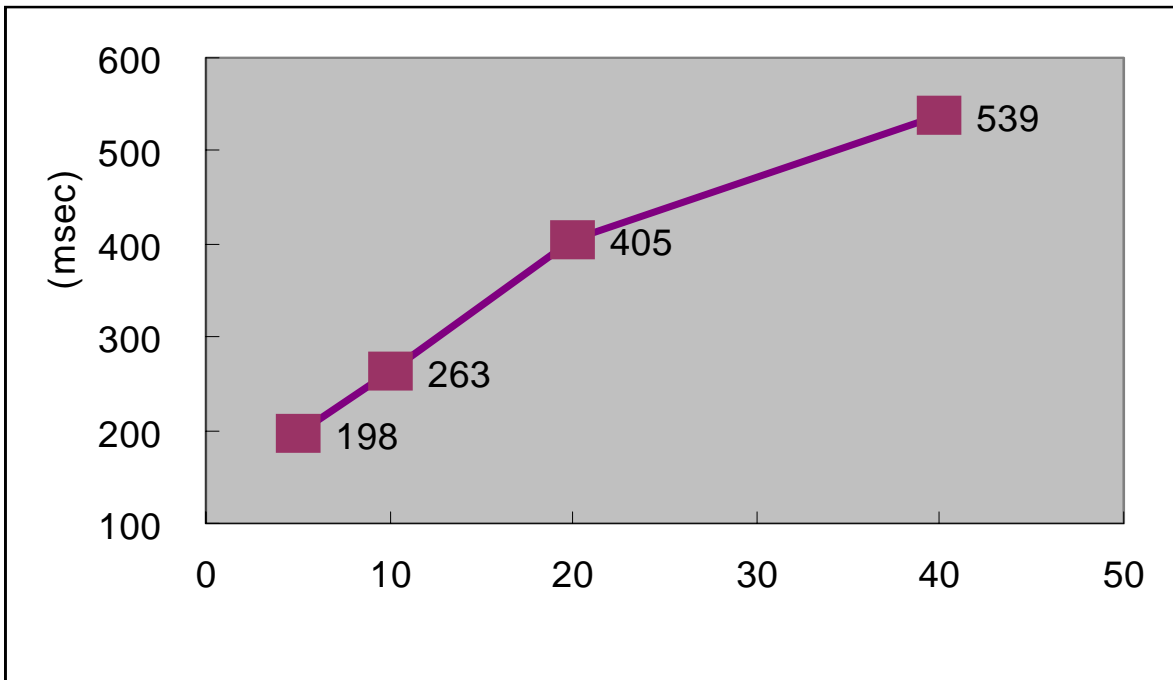
344,869 6  
 2-3 DocID (DL), (ACC),  
 (DL) 5 20  
 20  
 (MAX\_SORT\_SIZE) 5,000 ,  
 (MAX\_SET\_SIZE) 500 , 20

2.58 가

2-1 2-3 (DL)  
 70% 40%  
 ( 가 500 ,  
 (http://altavista.digital.com) 200 .)  
 가 DocID (IDG)  
 IDG  
 10% 35% 가  
 60% 가



(a)



(b)

2-4. Q05 - 5, Q10 - 10, Q20 - 20, Q40 - 40. (a) IDG: DocID, ACC: (accumulator), DL: ( ) (b)

2-4      2-2  
 5 (Q05), 10 (Q10), 20 (Q20), 40 (Q40)      10  
 2-4(a)      2-2      가  
 DocID      (IDG)      가  
 가      DocID      가  
 가      (ACC)      가      가  
 DocID      (DL)      MAX\_SET\_SIZE  
 가

( : msec)

set \	IDG	ACC	DL		
Q05	60	17	90	30	198
Q10	100	28	102	33	263
Q20	198	69	110	28	405
Q40	294	112	97	35	539
	163	57	100	32	351

2-2.

set      10      40  
 3.51 가

KRISTAL-II

MAX\_SET\_SIZE

KRISTAL-II

가

- test set 가
- DocID
- 
- ( )
- .....

similar to this) 2 가 (Query Refinement) (Find documents)

1

, 가 가 가

가  $w_t = TF * IDF^2$

TF IDF

가

( 2-5), 가

가

가

### 3.

Croft, W. B., and D. J. Harper. 1979. "Using Probabilistic Models of Document Retrieval Without Relevance Information." *Documentation*, 35(4), 285-95

Harman, D. 1986. "An Experimental Study of Factors Important in Document Ranking." Paper presented at ACM Conference on Research and Development in Information Retrieval, Cambridge, England.



Sparck Jones, K. 1972. "A Statistical Interpretation of Term Specificity and Its Application in Retrieval." *J. Documentation*, 28(1), 11-20.

Witten, I. H., Moffat, A., and Bell, T. C. 1994. "Managing Gigabytes", 141-148, Van Nostrand Reinhold, New York