

# COMP6714 Assignment 1

Hongxiao Jin (z5241154)

## Question 1

---

**Algorithm 1:**  $Q1(p_1, p_2, p_s)$ 

---

```
1 answer  $\leftarrow \emptyset$ ;
2 while  $p_1 \neq \text{nil} \wedge p_2 \neq \text{nil}$  do
3   if  $\text{doc}(p_1) = \text{doc}(p_2)$  then
4      $l \leftarrow []$ ;
5      $pp_1 \leftarrow \text{positions}(p_1)$ ;  $pp_2 \leftarrow \text{positions}(p_2)$ ;
6     while  $pp_1 \neq \text{nil}$  do
7        $\text{end} \leftarrow \text{skipTo}(p_s, \text{docID}(pp_1), \text{pos}(pp_1))$ 
8       while  $pp_2 \neq \text{nil}$  do
9         if  $\text{pos}(pp_2) < \text{end}$  then
10          if  $\text{pos}(pp_2) > \text{pos}(pp_1)$  then
11             $\text{add}(l, \text{pos}(pp_2))$ ;
12          else
13            break;
14           $pp_2 \leftarrow \text{next}(pp_2)$ ;
15        while  $l \neq [] \wedge l[1] < \text{pos}(pp_1)$  do
16           $\text{delete}(l[1])$ ;
17        for each  $ps \in l$  do
18           $\text{answer} \leftarrow \text{answer} \cup [\text{docID}(p_1), \text{pos}(pp_1), ps]$ ;
19         $pp_1 \leftarrow \text{next}(pp_1)$ ;
20       $p_1 \leftarrow \text{next}(p_1)$ ;  $p_2 \leftarrow \text{next}(p_2)$ ;
21   else
22     if  $\text{docID}(p_1) < \text{docID}(p_2)$  then
23        $p_1 \leftarrow \text{next}(p_1)$ ;
24     else
25        $p_2 \leftarrow \text{next}(p_2)$ ;
26 return answer
```

---

## Question 2

(1)

At first, the smallest( $\mathbf{Z}_0$ ) is in memory. If  $\mathbf{Z}_0$  gets to the upper bound of the memory, since the memory is full, it will be written into disk as  $\mathbf{I}_0$ . We can treat this  $\mathbf{I}_0$  as **Generation 0**. However, if  $\mathbf{I}_0$  already exists,  $\mathbf{Z}_0$  and  $\mathbf{I}_0$  will be merged into a new sub-index  $\mathbf{Z}_1$ . At this time, if there is no  $\mathbf{I}_1$ ,  $\mathbf{Z}_1$  will be written to disk as  $\mathbf{I}_1$  (**Generation 1**). While if  $\mathbf{I}_1$  is in disk,  $\mathbf{Z}_1$  will merge with  $\mathbf{I}_1$  to form  $\mathbf{Z}_2$  and so on.

From the algorithm of **logarithmic merge**, we can also know that this strategy maintain a series of indexes and each twice as large as the previous one.

Let  $\mathbf{n}$  represent the number of sub-indexes.

- The 1<sup>st</sup> layer can record  $\mathbf{M}$  pages;
- The 2<sup>nd</sup> layer can record  $2\mathbf{M}$  pages;
- .....
- The  $n^{\text{th}}$  layer can record  $2^{n-1} \cdot \mathbf{M}$  pages;

Thus,  $M \cdot t = M \cdot (1 + 2 + 2^2 + 2^3 + \dots + 2^{n-1})$

We can figure out  $t = 2^n - 1$  and  $n = \log_2(t + 1)$ .

**Therefore**, *logarithmic merge* will result in at most  $\lceil \log_2 t \rceil$  sub-indexes.

(2)

First of all, writing all pages into disk will cost  $\mathbf{M} \cdot \mathbf{t}$ .

Then, we should consider **merging cost**.

According to (1), the number of sub-indexes  $\mathbf{n} = \lceil \log_2 t \rceil$ , so sub-indexes can be record as  $\mathbf{I}_0, \mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_{n-1}$ .

Since each index is twice as large as the previous one, we can conclude that

$$2^1 \cdot I_0 = 2^0 \cdot I_1$$

$$2^2 \cdot I_0 = 2^1 \cdot I_1 = 2^0 \cdot I_2$$

$$2^3 \cdot I_0 = 2^2 \cdot I_1 = 2^1 \cdot I_2 = 2^0 \cdot I_3$$

.....

$$2^{n-1} \cdot I_0 = 2^{n-2} \cdot I_1 = 2^{n-3} \cdot I_2 = \dots = 2^1 \cdot I_{n-2} = 2^0 \cdot I_{n-1}$$

Every time when one  $I_{n-1}$  Generating, two  $I_{n-2}$  need to be read. Before reading two  $I_{n-2}$ , they have to be written into the disk.

Therefore, except every  $I_0$  is read once and every  $I_{n-1}$  is written once, the parameters generated in this process are read once and written once respectively.

The calculation process is as follows:

$$\begin{aligned}
merge\_cost &= (2^{n-1} \cdot 2^0 + 2 * (2^{n-2} \cdot 2^1 + 2^{n-3} \cdot 2^2 M + \dots + 2^1 \cdot 2^{n-2}) + 2^0 \cdot 2^{n-1}) \cdot M \\
&= M \cdot (2^n + 2^n \cdot (n - 2)) \\
&= M \cdot 2^n \cdot (n - 1) \\
&= M \cdot t \cdot (\log_2 t - 1) \\
total\_cost &= write\_cost + merge\_cost \\
&= M \cdot t + M \cdot t \cdot (\log_2 t - 1) \\
&= O(M \cdot t \cdot \log_2 t)
\end{aligned}$$

Therefore, the total I/O cost of the logarithmic merge is  $O(t \cdot M \cdot \log_2 t)$ .

### Question 3

Decode: 01000101 11110001 01110000 00110000 11110110 11011

0 : **1**

Let 10 0 0 represents number **k<sub>1</sub>**

$$\because k_{1dd} = \lfloor \log_2(k_{1d} + 1) \rfloor = 1$$

$$k_{1dr} = (k_{1d} + 1) - 2^{k_{1dd}} = 0$$

$$\therefore k_{1d} = 1$$

$$\because k_{1r} = k_1 - 2^{k_{1d}} = 0$$

$$\therefore k_1 = 2$$

Thus, 10 0 0 represents **2**.

Let 10 1 11 represents number **k<sub>2</sub>**

$$\because k_{2dd} = \lfloor \log_2(k_{2d} + 1) \rfloor = 1$$

$$k_{2dr} = (k_{2d} + 1) - 2^{k_{2dd}} = 1$$

$$\therefore k_{2d} = 2$$

$$\because k_{2r} = k_2 - 2^{k_{2d}} = 3$$

$$\therefore k_2 = 7$$

Thus, 10 1 11 represents **7**.

Similarly, we can calculate:

110 00 101: **13**

110 00 000: **8**

110 00 011: **11**

110 11 011011: **91**

Therefore, the result of decoding is 1, 2, 7, 13, 8, 11, 91.

Since these numbers The document IDs are 1, 3, 10, 23, 31, 42, 133.

## Question 4

(1)

**Line 21** causes the bug in the algorithm in Figure 2.

The main problem is the **pickTerm()** method. According to the original paper, we can know that **pickTerm()** selects the term with the maximal idf. So if the idf of a term is maximum among all terms and is the **pTerm** in one cycle, this algorithm will end up in an infinite loop.

(2)

The table below is from lecture 5b. Infinite loop would appear in the process of this algorithm.

	A	B	C
UB	4	5	8
List	$\langle 1, 3 \rangle$	$\langle 1, 4 \rangle$	$\langle 1, 4 \rangle$
	$\langle 2, 4 \rangle$	$\langle 2, 1 \rangle$	$\langle 2, 2 \rangle$
	$\langle 7, 1 \rangle$	$\langle 7, 2 \rangle$	$\langle 5, 1 \rangle$
		$\langle 8, 5 \rangle$	$\langle 7, 7 \rangle$
		$\langle 9, 2 \rangle$	$\langle 10, 1 \rangle$
		$\langle 11, 5 \rangle$	$\langle 11, 8 \rangle$

**In the 1<sup>st</sup> cycle**

$\theta = 0$ , curDoc = 0

Sorted Term	A	B	C
Doc	1	1	1
Cumulative Upper Bound	4	9	17

- pTerm = A
  - pivot = posting[0].DID = 1
  - pivot > curDoc (**Line14**)
    - posting[0].DID = pivot (**Line15**)
- curDoc = 1,  $\theta = 11$
- return (curDoc, posting)

**In the 2<sup>nd</sup> cycle**

$\theta = 11$ , curDoc = 1

Sorted Term	A	B	C
Doc	1	1	1
Cumulative Upper Bound	4	9	17

- pTerm = C
- pivot = posting[2].DID = 1
- pivot = curDoc (**Line10**)
  - aterm = pickTerm(terms[0..pTerm]) (**the idf of A is biggest, choose A**)

posting[A].DID = 2

**In the 3<sup>rd</sup> cycle**

$\theta = 11$ , curDoc = 1

Sorted Term	B	C	A
Doc	1	1	2
Cumulative Upper Bound	5	13	17

- pTerm = C
- pivot = posting[1].DID = 1
- pivot = curDoc (**Line10**)
  - aterm = pickTerm(terms[0..pTerm]) (**assume choose C**)
  - posting[C].DID = 2

**In the 4<sup>th</sup> cycle**

$\theta = 11$ , curDoc = 1

Sorted Term	B	A	C
Doc	1	2	2
Cumulative Upper Bound	5	9	17

- pTerm = C
- pivot = posting[2].DID = 2
- pivot > curDoc (**Line14**)
  - curDoc  $\neq$  posting[0].DID (**Line19**)
    - aterm = pickTerm(terms[0..pTerm]) (**the idf of A is biggest, choose A**)
    - posting[A].DID = 7

**In the 5<sup>th</sup> cycle**

$\theta = 11$ , curDoc = 1

Sorted Term	B	C	A
Doc	1	2	7
Cumulative Upper Bound	5	13	17

- pTerm = C
- pivot = posting[C].DID = 2
- pivot > curDoc (**Line14**)
  - posting[0].DID  $\neq$  pivot (**Line19**)
    - aterm = pickTerm(terms[0..pTerm]) (**the idf of A is biggest, choose A**)
    - posting[A].DID = 7

In this case, **pickTerm()** will always choose A because of the maximum idf.  $\theta$ , curDoc and other parameters would not change as well.

Therefore, this algorithm will end up in an infinite loop.