

# Индексация документов

## Лекция 7

БГУ ФПИИ, 2018

# План

Компоненты и жизненный цикл индекса

Словарь

Списки словопозиций

- Оптимизация posting lists

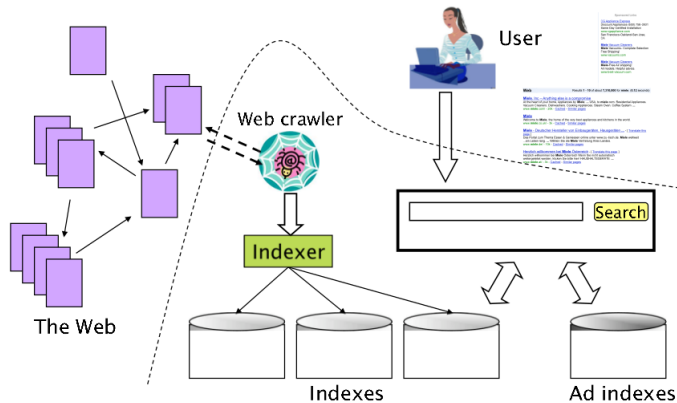
- Static index pruning

Построение статического индекса

- Sort-based indexing

- Single-pass in-memory indexing

- Merge-based indexing



# План

## Компоненты и жизненный цикл индекса

### Словарь

### Списки словопозиций

Оптимизация posting lists

Static index pruning

### Построение статического индекса

Sort-based indexing

Single-pass in-memory indexing

Merge-based indexing

Словарь		Словопозиции (postings)
Терм	$N_t$	
беда	3	3, 10, 11
верста	2	4, 5
друг	7	11, 14, 18, 21, 25, 34, 40
семь	10	1, 3, 4, 11, 15, 23, 37, 45, 51, 56
.....		.....

## Поиск в индексе

- ▶ Термы запроса ищутся в словаре.
- ▶ Списки docId для термов запроса пересекаются аналогично алгоритму merge.
- ▶ Для ускорения работы обработку списков нужно производить в порядке возрастания их длин.

## Операции

Inverted index как структура данных поддерживает

- ▶  $\text{first}(t)$
- ▶  $\text{last}(t)$
- ▶  $\text{next}(t, \text{current})$
- ▶  $\text{prev}(t, \text{current})$

## Типы индексов

### По детализации словопозиции:

- ▶ schema-independent
- ▶ schema-dependent:
  1. docId индекс
  2. частотный (frequency) индекс
  3. позиционный (positional) индекс

### По построению:

- ▶ статический
- ▶ динамический



Document ID	Document Content		
1	Do you quarrel, sir?		
2	Quarrel sir! no, sir!		
3	If you do, sir, I am for you: I serve as good a man as you.		
4	No better.		
5	Well, sir.		

Term	Docid list	Positional list	Schema-independent
a	1; 3	1; (3, 1, ⟨13⟩)	1; 21
am	1; 3	1; (3, 1, ⟨6⟩)	1; 14
as	1; 3	1; (3, 2, ⟨11, 15⟩)	2; 19, 23
better	1; 4	1; (4, 1, ⟨2⟩)	1; 26
do	2; 1, 3	2; (1, 1, ⟨1⟩), (3, 1, ⟨3⟩)	2; 1, 11
for	1; 3	1; (3, 1, ⟨7⟩)	1; 15
good	1; 3	1; (3, 1, ⟨12⟩)	1; 20
i	1; 3	1; (3, 2, ⟨5, 9⟩)	2; 13, 17
if	1; 3	1; (3, 1, ⟨1⟩)	1; 9
man	1; 3	1; (3, 1, ⟨14⟩)	1; 22
no	2; 2, 4	2; (2, 1, ⟨3⟩), (4, 1, ⟨1⟩)	2; 7, 25
quarrel	2; 1, 2	2; (1, 1, ⟨3⟩), (2, 1, ⟨1⟩)	2; 3, 5
serve	1; 3	1; (3, 1, ⟨10⟩)	1; 18
sir	4; 1, 2, 3, 5	4; (1, 1, ⟨4⟩), (2, 2, ⟨2, 4⟩), (3, 1, ⟨4⟩), (5, 1, ⟨2⟩)	5; 4, 6, 8, 12, 28
well	1; 5	1; (5, 1, ⟨1⟩)	1; 27
you	2; 1, 3	2; (1, 1, ⟨2⟩), (3, 3, ⟨2, 8, 16⟩)	4; 2, 10, 16, 24

## Статический индекс

- ▶ Индексируем коллекцию один раз, затем ищем
- ▶ Индексация:
  - ▶ in-memory
  - ▶ sort-based
  - ▶ merge-based.

# План

Компоненты и жизненный цикл индекса

Словарь

Списки словопозиций

Оптимизация posting lists

Static index pruning

Построение статического индекса

Sort-based indexing

Single-pass in-memory indexing

Merge-based indexing

## Реализация

- ▶ Хэш-таблица
- ▶ Дерево поиска (В-дерево)
- ▶ Сортированный массив

Выбор решения:

- ▶ сколько ключей требуется?
- ▶ изменяется ли словарь?
- ▶ какова частота обращения к различным ключам?

# План

Компоненты и жизненный цикл индекса

Словарь

Списки словопозиций

Оптимизация posting lists

Static index pruning

Построение статического индекса

Sort-based indexing

Single-pass in-memory indexing

Merge-based indexing

# Posting

$$(d, f_{t,d}, \langle p_0, \dots, p_{f_{t,d}} \rangle)$$

- ▶  $d$  — docId
- ▶  $f_{t,d}$  — tf
- ▶  $\langle p_0, \dots, p_{f_{t,d}} \rangle$  — offsets, sorted

## Skip lists

Во время поиска задействованы списки словопозиций разных длин:

- Это неэффективно
- Хорошо бы перепрыгивать через некоторое количество документов
- Компрессия делает это проблематичным

## Skip lists

Skip list — структура данных для поддержки прыжков по словопозициям, представляет собой список пар (posting, position).

Оптимальное количество skip pointers  $\sim \sqrt{P}$ .



## Пример

- ▶ Постинг лист:

5,11,17,21,26,34,36,37,45,48,51,52,57,80,89,91

- ▶ Дельта кодирование:

5,6,6,4,5,8,2,1,8,3,3,1,5,23,9,2

- ▶ Skip list:

(17,3),(34,6),(45,9),(52,12),(89,15)

INTERSECTWITHSKIPS( $p_1, p_2$ )

```
1  answer  $\leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4      then ADD(answer,  $\text{docID}(p_1)$ )
5           $p_1 \leftarrow \text{next}(p_1)$ 
6           $p_2 \leftarrow \text{next}(p_2)$ 
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8      then if  $\text{hasSkip}(p_1)$  and  $(\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2))$ 
9          then while  $\text{hasSkip}(p_1)$  and  $(\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2))$ 
10             do  $p_1 \leftarrow \text{skip}(p_1)$ 
11             else  $p_1 \leftarrow \text{next}(p_1)$ 
12  else if  $\text{hasSkip}(p_2)$  and  $(\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1))$ 
13      then while  $\text{hasSkip}(p_2)$  and  $(\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1))$ 
14          do  $p_2 \leftarrow \text{skip}(p_2)$ 
15          else  $p_2 \leftarrow \text{next}(p_2)$ 
16  return answer
```

## Static index pruning

top  $k$  prune( $I, k, \epsilon$ )  
  for each term  $t$  in  $I$   
    retrieve the posting list  $P_t$  from  $I$   
    if  $|P_t| > k$   
      for each entry  $d \in P_t$   
        compute  $A(t, d)$  (e.g. by Equation 1)  
        let  $z_t$  be the  $k$ th best entry in row  $t$  of  $A$   
         $\tau_t \leftarrow \epsilon \cdot z_t$   
        for each entry  $d \in P_t$   
          if  $A(t, d) \leq \tau_t$   
            remove entry  $d$  from  $P_t$   
  Save  $(t, P_t)$  in the pruned inverted file

$$A(t, d) = \frac{\frac{\log(1+f_{t,d})}{\log(1+avgf_d)} \log \frac{N}{N_t}}{|d|}$$

# План

Компоненты и жизненный цикл индекса

Словарь

Списки словопозиций

Оптимизация posting lists

Static index pruning

Построение статического индекса

Sort-based indexing

Single-pass in-memory indexing

Merge-based indexing

## Пример

1. ['Семь раз отмерь, один раз отрежь.', 'Один за всех, все за одного', 'Семь бед — один ответ.', 'Семь вёрст до небес и все лесом.']
2. [['семь', 'раз', 'отмерить', 'один', 'раз', 'отрезать'], ['один', 'за', 'все', 'все', 'за', 'один'], ['семь', 'беда', 'один', 'ответ'], ['семь', 'верста', 'до', 'небеса', 'и', 'все', 'лес']]
3. [('беда', 3), ('верста', 4), ('все', 2), ('все', 2), ('все', 4), ('до', 4), ('за', 2), ('за', 2), ('и', 4), ('лес', 4), ('небеса', 4), ('один', 1), ('один', 2), ('один', 2), ('один', 3), ('ответ', 3), ('отмерить', 1), ('отрезать', 1), ('раз', 1), ('раз', 1), ('семь', 1), ('семь', 3), ('семь', 4)]

Терм	Словопозиции
'беда'	[3]
'верста'	[4]
'все'	[2, 4]
'до'	[4]
'за'	[2]
'и'	[4]
'лес'	[4]
'небеса'	[4]
'один'	[1, 2, 3]
'ответ'	[3]
'отмерить'	[1]
'отрезать'	[1]
'раз'	[1]
'семь'	[1, 3, 4]

**buildIndex\_sortBased** (*inputTokenizer*)  $\equiv$

```
1   position  $\leftarrow$  0
2   while inputTokenizer.hasNext() do
3       T  $\leftarrow$  inputTokenizer.getNext()
4       obtain dictionary entry for T; create new entry if necessary
5       termID  $\leftarrow$  unique term ID of T
6       write record  $R_{position} \equiv (termID, position)$  to disk
7       position  $\leftarrow position + 1$ 
8   tokenCount  $\leftarrow position$ 
9   sort  $R_0 \dots R_{tokenCount-1}$  by first component; break ties by second component
10  perform a sequential scan of  $R_0 \dots R_{tokenCount-1}$ , creating the final index
11  return
```

## Алгоритм блочного индексирования

1. Сегментируем коллекцию на равные части.
2. Для каждой части строим обратный индекс в памяти одним из способов.
3. Сохраняем inverted index каждой части на диск.
4. Объединяем все промежуточные результаты в окончательный индекс.



```
mergeIndexPartitions ( $\langle I_0, \dots, I_{n-1} \rangle$ )  $\equiv$   
1   create empty inverted file  $I_{\text{final}}$   
2   for  $k \leftarrow 0$  to  $n - 1$  do  
3       open index partition  $I_k$  for sequential processing  
4    $\text{currentIndex} \leftarrow 0$   
5   while  $\text{currentIndex} \neq \text{nil}$  do  
6        $\text{currentIndex} \leftarrow \text{nil}$   
7       for  $k \leftarrow 0$  to  $n - 1$  do  
8           if  $I_k$  still has terms left then  
9               if ( $\text{currentIndex} = \text{nil}$ )  $\vee$  ( $I_k.\text{currentTerm} < \text{currentTerm}$ ) then  
10                   $\text{currentIndex} \leftarrow I_k$   
11                   $\text{currentTerm} \leftarrow I_k.\text{currentTerm}$   
12           if  $\text{currentIndex} \neq \text{nil}$  then  
13                $I_{\text{final}}.\text{addPostings}(\text{currentTerm}, \text{currentIndex}.\text{getPostings}(\text{currentTerm}))$   
14                $\text{currentIndex}.\text{advanceToNextTerm}()$   
15   delete  $I_0 \dots I_{n-1}$   
16   return
```