

Inverted File Index

以单词或短语作为关键词，将他们出现在文档中的位置记录在索引中。

例如：{

No Terms Times; (Doc ID; Places)

↑

Term dictionary

↑

Posting list

Index generator

```
while (read a document D) {
```

```
    while (read a Term T in D) {
```

```
        if (Find (Dictionary, T) == false)
```

```
            Insert (Dictionary, T);
```

```
        Get T's Posting list;
```

```
        Insert a node to T's posting list;
```

```
}
```

```
}
```

Write the inverted index to disk;

read a term

word stemming 词干分析

process, processing, processed, processes \Rightarrow process.

stop words 停用词

"a" "the" "it". 但会有问题..?

access a term

search trees (B-tree, B+tree, Tries...)

Hashing

⇒ 如果没有足够的内存:

if (out of memory) {

 write blockIndex[i] to disk;

 cnt++

 Free Memory block;

}

然后 merge (Inverted index, BlockIndex[i])
↑
o - cnt

distributed indexing 分布式索引

Term-partitioned index : a-c, d-f, ...

Document-partitioned index : 文档/机器分布

Dynamic indexing

当有文档 come in, 有 term 需添加或插入, 那么这一块 posting list 且会被
更新 (dynamic). 另外一部分做上标记, 如果不常用 in db → 相对静态,
分块来减少操作复杂度. Docs deleted.

Main Index 主存

Auxiliary index 辅助存储

地圖： 西北亞文學

英語文學.