# CMPSCI 546 (590R)
# Applied Information Retrieval

Indexing

# V-Byte Encoder

```java
public void encode ( int [] input, ByteBuffer output) {
  for (int i : input ) {
    while ( i >= 128 ) {
      output.put ( i & 0x7F ) ;
      i >>>= 7 ; // logical shift, no sign bit extension
    }
    output.put( i | 0x80 );
  }
}
```

# V-Byte Decoder

```
public void decode ( byte [] input, IntBuffer output ) {
    for ( int i = 0; i < input.length; i++ ) {
        int position = 0;
        int result = ((int) input[i] & 0x7F);

        while ( (input[i]] & 0x80) == 0 ) {
            i += 1;
            position += 1;
            int unsignedByte = ((int) input[i] & 0x7F);
            result |= (unsignedByte << (7 * position));
        }
        output.put(result);
    }
}
```

# Auxiliary Structures

- Inverted lists usually stored together in a single file for efficiency
  - *Inverted file*
- *Vocabulary* or *lexicon*
  - Contains a lookup table from index terms to the byte offset of the inverted list in the inverted file
  - Either hash table in memory or B-tree for larger vocabularies
- Term statistics stored at start of inverted lists
- Collection statistics stored in separate file

# Index Construction

- Simple in-memory indexer

**procedure** BUILDINDEX($D$)                    $\triangleright$ $D$ is a set of text documents
 $I \leftarrow$ HashTable()                    $\triangleright$ Inverted list storage
 $n \leftarrow 0$                    $\triangleright$ Document numbering
 **for all** documents $d \in D$ **do**
  $n \leftarrow n + 1$
  $T \leftarrow$ Parse($d$)                    $\triangleright$ Parse document into tokens
  Remove duplicates from $T$
  **for all** tokens $t \in T$ **do**
   **if** $I_t \notin I$ **then**
    $I_t \leftarrow$ Array()
   **end if**
   $I_t$.append($n$)
  **end for**
 **end for**
 **return** $I$
**end procedure**

*Figure 5.8*

# Document-At-A-Time

**procedure** DOCUMENTATATIMERETRIEVAL$(Q, I, f, g, k)$
    $L \leftarrow \text{Array}()$
    $R \leftarrow \text{PriorityQueue}(k)$
    **for all** terms $w_i$ in $Q$ **do**
        $l_i \leftarrow \text{InvertedList}(w_i, I)$
        $L.\text{add}(\ l_i\ )$
    **end for**
    **for all** documents $d \in I$ **do**
        $s_d \leftarrow 0$
        **for all** inverted lists $l_i$ in $L$ **do**
            **if** $l_i.\text{getCurrentDocument}() = d$ **then**
                $s_d \leftarrow s_d + g_i(Q)f_i(l_i)$           ▷ Update the document score
            **end if**
            $l_i.\text{movePastDocument}(\ d\ )$
        **end for**
        $R.\text{add}(\ s_d, d\ )$
    **end for**
    **return** the top $k$ results from $R$
**end procedure**

$$R(Q, D) = \sum_i g_i(Q)f_i(D)$$

*Figure 5.16*

# Term-At-A-Time

**procedure** $\textsc{TermAtATimeRetrieval}(Q,\ I,\ f,\ g\ k)$
$\quad A \leftarrow \text{HashTable}()$
$\quad L \leftarrow \text{Array}()$
$\quad R \leftarrow \text{PriorityQueue}(k)$
$\quad$ **for all** terms $w_i$ in $Q$ **do**
$\quad\quad l_i \leftarrow \text{InvertedList}(w_i,\ I)$
$\quad\quad L.\text{add}(\ l_i\ )$
$\quad$ **end for**
$\quad$ **for all** lists $l_i \in L$ **do**
$\quad\quad$ **while** $l_i$ is not finished **do**
$\quad\quad\quad d \leftarrow l_i.\text{getCurrentDocument}()$
$\quad\quad\quad A_d \leftarrow A_d + g_i(Q)f(l_i)$
$\quad\quad\quad l_i.\text{moveToNextDocument}()$
$\quad\quad$ **end while**
$\quad$ **end for**
$\quad$ **for all** accumulators $A_d$ in $A$ **do**
$\quad\quad s_d \leftarrow A_d \qquad\qquad \triangleright$ Accumulator contains the document score
$\quad\quad R.\text{add}(\ s_d, d\ )$
$\quad$ **end for**
$\quad$ **return** the top $k$ results from $R$
**end procedure**

$$R(Q, D) = \sum_i g_i(Q)f_i(D)$$

*Figure 5.18*

# Reminder about using lists

- "Align" inverted lists / traverse in parallel
- Allows
  - Simple Boolean AND, OR, NOT
  - Proximity operators (if has position information)
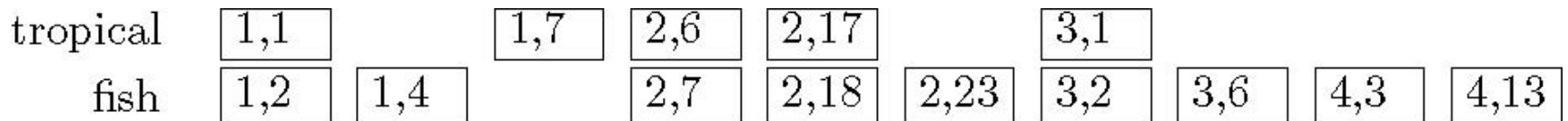  - Combining scores across terms

| tropical | 1,1 | | | 1,7 | 2,6 | 2,17 | | | 3,1 | | | |
|----------|-----|-----|-----|-----|-----|------|-----|-----|-----|-----|-----|------|
| fish | 1,2 | 1,4 | | | 2,7 | 2,18 | 2,23 | 3,2 | | 3,6 | 4,3 | 4,13 |

*Figure 5.6*

```
 1: procedure TERMATATIMERETRIEVAL(Q, I, f, g, k)
 2:     A ← Map()
 3:     L ← Array()
 4:     R ← PriorityQueue(k)
 5:     for all terms w_i in Q do
 6:         l_i ← InvertedList(w_i, I)
 7:         L.add( l_i )
 8:     end for
 9:     for all lists l_i ∈ L do
10:         d_0 ← −1
11:         while l_i is not finished do
12:             if i = 0 then
13:                 d ← l_i.getCurrentDocument()
14:                 A_d ← A_d + g_i(Q)f(l_i)
15:                 l_i.moveToNextDocument()
16:             else
17:                 d ← l_i.getCurrentDocument()
18:                 d' ← A.getNextAccumulator(d)
19:                 A.removeAccumulatorsBetween(d_0, d')
20:                 if d = d' then
21:                     A_d ← A_d + g_i(Q)f(l_i)
22:                     l_i.moveToNextDocument()
23:                 else
24:                     l_i.skipForwardToDocument(d')
25:                 end if
26:                 d_0 ← d'
27:             end if
28:         end while
29:     end for
30:     for all accumulators A_d in A do
31:         s_d ← A_d                    ▷ Accumulator contains the document score
32:         R.add( s_d, d )
33:     end for
34:     return the top k results from R
35: end procedure
```

# Conjunctive Term-at-a-Time

*Figure 5.20*

```
 1: procedure DOCUMENTATATIMERETRIEVAL(Q, I, f, g, k)
 2:     L ← Array()
 3:     R ← PriorityQueue(k)
 4:     for all terms w_i in Q do
 5:         l_i ← InvertedList(w_i, I)
 6:         L.add( l_i )
 7:     end for
 8:     d ← −1
 9:     while all lists in L are not finished do
10:         s_d ← 0
11:         for all inverted lists l_i in L do
12:             if l_i.getCurrentDocument() > d then
13:                 d ← l_i.getCurrentDocument()
14:             end if
15:         end for
16:         for all inverted lists l_i in L do
17:             l_i.skipForwardToDocument(d)
18:             if l_i.getCurrentDocument() = d then
19:                 s_d ← s_d + g_i(Q)f_i(l_i)            ▷ Update the document score
20:                 l_i.movePastDocument( d )
21:             else
22:                 d ← −1
23:                 break
24:             end if
25:         end for
26:         if d > −1 then R.add( s_d, d )
27:         end if
28:     end while
29:     return the top k results from R
30: end procedure
```

# Conjunctive Document-at-a-Time

*Figure 5.21*