

Adapting Boyer-Moore-Like Algorithms for Searching Huffman Encoded Texts

Domenico Cantone Simone Faro Emanuele Giaquinta

Department of Mathematics and Computer Science, University of Catania, Italy

String matching in compressed texts

Compressed matching problem (Amir & Benson, 1992):

- Alphabet Σ
- Pattern P
- Compression system $(\mathcal{E}, \mathcal{D})$
- Encoded Text $\mathcal{E}(T)$

Find all the shifts of P in T using $\mathcal{E}(P)$ and $\mathcal{E}(T)$

String matching in compressed texts

A static compression method is characterized by a system $(\mathcal{E}, \mathcal{D})$ of two complementary functions:

- $\mathcal{E} : \Sigma \rightarrow \{0, 1\}^+$
 - $\mathcal{E}(\varepsilon) = \varepsilon$
 - $\mathcal{E}(T[1..l]) = \mathcal{E}(T[1..l-1]).\mathcal{E}(T[l]), \forall l : 1 \leq l \leq |T|$
- $\mathcal{D}(\mathcal{E}(c)) = c, \forall c \in \Sigma$

String matching in compressed texts

A static compression method is characterized by a system $(\mathcal{E}, \mathcal{D})$ of two complementary functions:

- $\mathcal{E} : \Sigma \rightarrow \{0, 1\}^+$
 - $\mathcal{E}(\varepsilon) = \varepsilon$
 - $\mathcal{E}(T[1..l]) = \mathcal{E}(T[1..l-1]).\mathcal{E}(T[l]), \forall l : 1 \leq l \leq |T|$
- $\mathcal{D}(\mathcal{E}(c)) = c, \forall c \in \Sigma$

- Prefix property - $\nexists c_1, c_2 : \mathcal{E}(c_1) \sqsubseteq \mathcal{E}(c_2)$
- **Canonical Huffman coding**

String matching in compressed texts

t :	00		
e :	01		
w :	100	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
a :	101	twenty	$\bar{0}\bar{0}\bar{1}0\bar{0}\bar{0}1\bar{1}\bar{1}10\bar{0}\bar{0}\bar{1}110$
n :	110		
y :	1110		
b :	1111		

String matching in compressed texts

t :	00		
e :	01		
w :	100	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
a :	101	twenty	$\bar{0}\bar{0}\bar{1}0\bar{0}\bar{0}\bar{1}\bar{1}10\bar{0}\bar{0}\bar{1}110$
n :	110	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
y :	1110	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
b :	1111		

String matching in compressed texts

t :	00		
e :	01		
w :	100	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
a :	101	twenty	$\bar{0}\bar{0}\bar{1}00\bar{0}\bar{1}\bar{1}10\bar{0}\bar{0}\bar{1}110$
n :	110	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
y :	1110	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
b :	1111		

Problem: false positives, occurrences of $\mathcal{E}(P)$ in $\mathcal{E}(T)$ which do not correspond to occurrences of P in T .

An occurrence of $\mathcal{E}(P)$ which does not start on a codeword boundary is a **false positive**.

A binary prefix code can be represented with an ordered binary tree, whose leaves are labeled with characters in Σ and whose edges are labeled by 0 (left) and 1 (right).

The codeword of a given character is the word labeling the branch from the root to the leaf labeled by the same character.

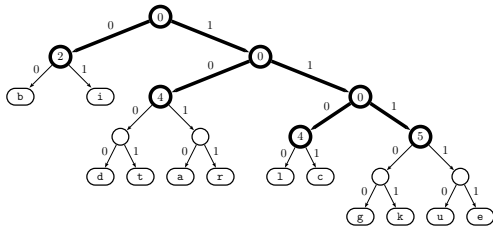
The minimal prefix of a codeword which allows one to unambiguously determine the codeword length corresponds to the minimum depth node in the path from the root induced by the codeword, which is the root of a complete subtree.

The skeleton tree (Klein, 2000) is a pruned canonical Huffman tree whose leaves correspond to minimum depth nodes in the Huffman tree which are roots of complete subtrees.

It is useful to maintain at each leaf of a skeleton tree the common length of the codeword(s) sharing the prefix which labels the path from the root to it.

skeleton tree

b : 00
i : 01
d : 1000
t : 1001
a : 1010
r : 1011
l : 1100
c : 1101
g : 11100
k : 11101
u : 11110
e : 11111



- SK-KMP (Daptardar & Shapira, 2006): modified KMP, prefix alignments always respect codeword boundaries. Decoding uses the skeleton tree. If no prefix matches the border of the current window, the algorithm can skip the remaining bits for the current codeword.

- Filtering/verification paradigm
- The filtering phase searches the occurrences of $\mathcal{E}(P)$ in $\mathcal{E}(T)$
- The verification phase uses an algorithm based on the skeleton tree to verify candidate matches

Suppose that we have found a candidate shift s . We need to know if s is codeword aligned.

We maintain an offset ρ pointing to the start of the last window where a verification has taken place.

We update - using the skeleton tree - ρ to a minimal position $\rho^* \geq s$ which is codeword aligned. If $\rho^* = s$, s is a valid shift.

SK-ALIGN(*root*, *t*, ρ , *s*)

```
1  $x \leftarrow \text{root}, \ell \leftarrow 0$ 
2 while TRUE
3     do  $B \leftarrow B_t[\lfloor \rho / k \rfloor] \ll (\rho \bmod k)$ 
4         if  $B < 2^{k-1}$  then  $x \leftarrow \text{LEFT}(x)$  else  $x \leftarrow \text{RIGHT}(x)$ 
5         if  $\text{KEY}(x) \neq 0$ 
6             then  $\rho \leftarrow \rho + \text{KEY}(x) - \ell, \ell \leftarrow 0, x \leftarrow \text{root}$ 
7                 if  $\rho \geq s$  then break
8             else  $\rho \leftarrow \rho + 1, \ell \leftarrow \ell + 1$ 
9 return  $\rho$ 
```

skeleton tree verification

t :	00		
e :	01		
w :	100		
a :	101	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
n :	110	twenty	$\bar{0}\bar{0}\bar{1}0\bar{0}\bar{0}\bar{1}\bar{1}10\bar{0}\bar{0}\bar{1}110$
y :	1110		
b :	1111		

skeleton tree verification

```
t : 00
e : 01
w : 100
a : 101    ten    0̄0̄0̄1̄1̄10
n : 110    twenty 0̄0̄1̄00̄0̄1̄1̄100̄0̄1̄110
y : 1110
b : 1111
```

$\bar{0}\bar{0}\bar{1}00\bar{0}\bar{1}\bar{1}10\bar{0}\bar{0}\bar{1}110$ $j = 3$

$\text{SK-ALIGN}(\text{root}, t, 0, 3) = 5 \neq 3$

skeleton tree verification

t :	00		
e :	01		
w :	100		
a :	101	ten	$\bar{0}\bar{0}\bar{0}1\bar{1}10$
n :	110	twenty	$\bar{0}\bar{0}\bar{1}0\bar{0}\bar{0}1\bar{1}10\bar{0}\bar{0}1110$
y :	1110		
b :	1111		

$\bar{0}\bar{0}\bar{1}0\bar{0}\bar{0}1\bar{1}10\bar{0}\bar{0}1110$ $j = 3$

$\text{SK-ALIGN}(\text{root}, t, 0, 3) = 5 \neq 3$

$\bar{0}\bar{0}\bar{1}0\bar{0}\bar{0}1\bar{1}10\bar{0}\bar{0}1110$ $j = 9$

$\text{SK-ALIGN}(\text{root}, t, 5, 9) = 10 \neq 9$

Pros

- Lazy decoding: in the best case (the pattern does not occur) no decoding is performed
- For every codeword the algorithm reads the minimum number only of bits necessary to infer its length

Cons

- The number of bits processed depends on the position of candidate matches

```
HUFFMAN-MATCHER( $P, m, T, n$ )
1  PRECOMPUTE_GLOBSALS( $P$ )
2   $n \leftarrow |T|$ 
3   $m \leftarrow |P|$ 
4   $s \leftarrow 0$ 
5   $\rho \leftarrow 0$ 
6   $root \leftarrow$  BUILD-SK-TREE( $\phi$ )
7  for  $i \leftarrow 0$  to  $n - 1$ 
8      do  $j \leftarrow$  CHECK_SHIFT( $s, P, T$ )
9          if  $j = s + m - 1$ 
10             then  $\rho \leftarrow$  SK-ALIGN( $root, T, \rho, j$ )
11                 if  $\rho = j$  then PRINT( $j$ )
12              $s \leftarrow s +$  SHIFT_INCREMENT( $s, P, T, j$ )
```

Adaptation of existing algorithms for binary strings

- BINARY-HASH-MATCHING (Lecroq & Faro, 2009)
- FED (J. Kim & E. Kim & Park, 2007)

string matching in binary strings

- Scanning T with bit granularity is slow
- $T[0, \dots, n - 1] \rightarrow B_T[0, \dots, \lceil n/k \rceil - 1]$, $B_T[i]$ is a block of k bits
- An occurrence of P can start at bit $1, 2, \dots, k$ of a block
- $P_i \leftarrow P \gg i$
- $P_i[0, \dots, m + i - 1] \rightarrow Patt_i[0, \dots, \lceil (m + i)/k \rceil - 1]$

string matching in binary strings

- Scanning T with bit granularity is slow
- $T[0, \dots, n - 1] \rightarrow B_T[0, \dots, \lceil n/k \rceil - 1]$, $B_T[i]$ is a block of k bits
- An occurrence of P can start at bit $1, 2, \dots, k$ of a block
- $P_i \leftarrow P \gg i$
- $P_i[0, \dots, m + i - 1] \rightarrow Patt_i[0, \dots, \lceil (m + i)/k \rceil - 1]$

```
ent      10001110
twenty  00100011-10001110
```

string matching in binary strings

- Scanning T with bit granularity is slow
- $T[0, \dots, n-1] \rightarrow B_T[0, \dots, \lceil n/k \rceil - 1]$, $B_T[i]$ is a block of k bits
- An occurrence of P can start at bit $1, 2, \dots, k$ of a block
- $P_i \leftarrow P \gg i$
- $P_i[0, \dots, m+i-1] \rightarrow Patt_i[0, \dots, \lceil (m+i)/k \rceil - 1]$

ent	10001110
twenty	00100011-10001110
P_2	00100011-10000000

string matching in binary strings

$$\mathcal{E}(P) = 110010110010110010110$$

(A) <i>Patt</i>	0	1	2	3
0	<u>11001011</u>	<u>00101100</u>	<u>10110000</u>	
1	<u>01100101</u>	<u>10010110</u>	<u>01011000</u>	
2	<u>00110010</u>	<u>11001011</u>	<u>00101100</u>	
3	<u>00011001</u>	<u>01100101</u>	<u>10010110</u>	
4	<u>00001100</u>	<u>10110010</u>	<u>11001011</u>	<u>00000000</u>
5	<u>00000110</u>	<u>01011001</u>	<u>01100101</u>	<u>10000000</u>
6	<u>00000011</u>	<u>00101100</u>	<u>10110010</u>	<u>11000000</u>
7	<u>00000001</u>	<u>10010110</u>	<u>01011001</u>	<u>01100000</u>

(B) <i>Mask</i>	0	1	2	3
0	<u>11111111</u>	<u>11111111</u>	<u>11111000</u>	
1	<u>01111111</u>	<u>11111111</u>	<u>11111100</u>	
2	<u>00111111</u>	<u>11111111</u>	<u>11111110</u>	
3	<u>00011111</u>	<u>11111111</u>	<u>11111111</u>	
4	<u>00001111</u>	<u>11111111</u>	<u>11111111</u>	<u>10000000</u>
5	<u>00000111</u>	<u>11111111</u>	<u>11111111</u>	<u>11000000</u>
6	<u>00000011</u>	<u>11111111</u>	<u>11111111</u>	<u>11100000</u>
7	<u>00000001</u>	<u>11111111</u>	<u>11111111</u>	<u>11110000</u>

The pattern $\mathcal{E}(P)$ is aligned with the s -th bit of the text if

$$Patt_i[h] = B_T[j + h] \& Mask_i[h], \text{ for } h = 0, 1, \dots, m_i - 1$$

where $j = \lfloor s/k \rfloor$, $i = s \bmod k$, $m_i = \lceil (m + i)/k \rceil$

string matching in binary strings

- Bad character rule does not work well with binary alphabets,
 $bc(c) \rightarrow 1, \forall c \in \Sigma$
- Super alphabet: $2 \rightarrow 2^k$
- P is logically divided into $m - k$ overlapping grams

Shift with bit granularity

$$Hs(B) = \min \left(\{0 \leq u < m \mid p[m - u - k .. m - u - 1] \supseteq B\} \cup \{m\} \right), 0 \leq B < 2^k$$

Shift with bit granularity

$$Hs(B) = \min \left(\{0 \leq u < m \mid p[m - u - k .. m - u - 1] \supseteq B\} \cup \{m\} \right), 0 \leq B < 2^k$$

- if $Hs[B] = 0$, the alignment $(s - m + 1) \bmod k$ is checked, where s is the current position in T in bits

Shift with byte granularity

$$\delta_i(c) = \min(\{m_i - 2 + 1\} \cup \{m_i - 2 + 1 - k \mid \text{Patt}_i[k] = c \text{ and } 1 \leq k \leq m_i - 2\})$$

$$\delta(c) = \min\{\delta_i(c), 0 \leq i < k\}$$

Shift with byte granularity

$$\delta_i(c) = \min(\{m_i - 2 + 1\} \cup \{m_i - 2 + 1 - k \mid \text{Patt}_i[k] = c \text{ and } 1 \leq k \leq m_i - 2\})$$

$$\delta(c) = \min\{\delta_i(c), 0 \leq i < k\}$$

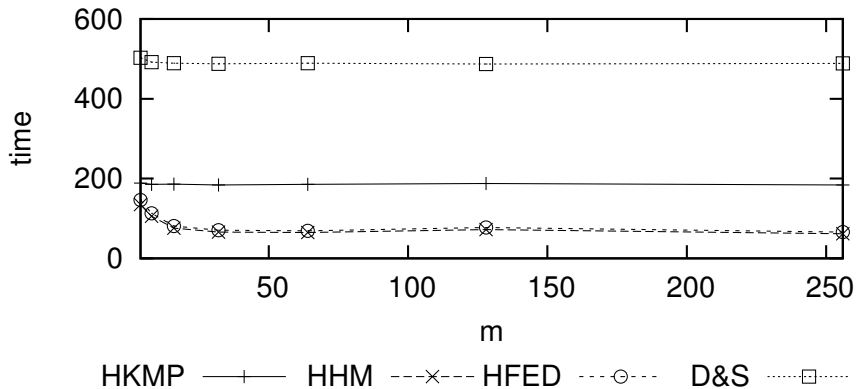
- hash table λ to find candidate alignments, each entry in the table pointing to a linked list of patterns
- pattern Patt_i is inserted into the slot with hash value $\text{Patt}_i[m_i - 2]$

- $\mathcal{O}(\lceil m/k \rceil n)$ worst case time complexity
- $\mathcal{O}(m + 2^k)$ space complexity

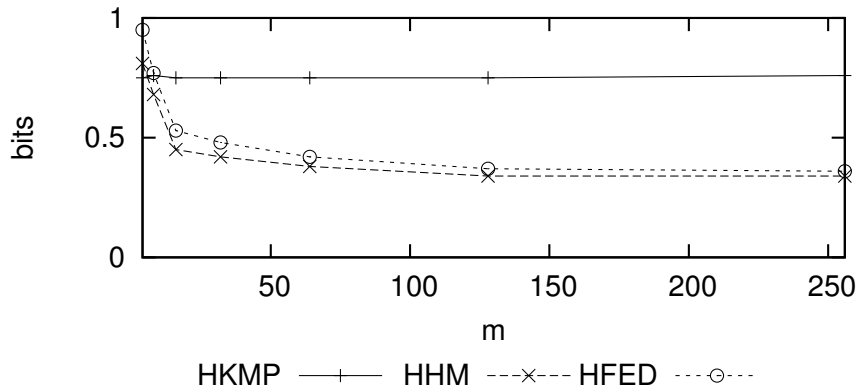
$k \leq 8$ works well for a single pattern, space overhead negligible

- Implementation in C, compiled with gcc 4, run on PowerPC G4 1.5 GHz
- Natural language texts
- Set of 100 patterns of fixed length $m \in \{4, 8, 16, 32, 64, 128, 256\}$, randomly extracted from the text
- Comparison between the following algorithms:
 - HUFFMAN-KMP
 - HUFFMAN-HASH-MATCHING
 - HUFFMAN-FED
 - Decompress and Search with *3-Hash* algorithm

Experimental results - running times



Experimental results - processed bits



- Our generic algorithm can skip many bits when decoding
- Sublinear on average when used with Boyer-Moore like algorithms
- Fast, especially when the pattern frequency is low