

COMP9319 Web Data Compression and Search

LZW,
Adaptive Huffman

1

1

Dictionary coding

- Patterns: correlations between part of the data
- Idea: replace recurring patterns with references to dictionary
- LZ algorithms are adaptive:
 - Universal coding (the prob. distr. of a symbol is unknown)
 - Single pass (dictionary created on the fly)
 - No need to transmit/store dictionary

2

2

LZ77 & LZ78

- LZ77: referring to previously processed data as dictionary
- LZ78: use an explicit dictionary

3

3

Lempel-Ziv-Welch (LZW) Algorithm

- Most popular modification to LZ78
- Very common, e.g., Unix compress, TIFF, GIF, PDF (until recently)
- Read <http://en.wikipedia.org/wiki/LZW> regarding its patents
- Fixed-length references (12bit 4096 entries)
- Static after max entries reached

4

4

Patent issues again

From Wikipedia: "In 1993–94, and again in 1999, Unisys Corporation received widespread condemnation when it attempted to enforce licensing fees for LZW in GIF images. The 1993–1994 Unisys-Compuserve (Compuserve being the creator of the GIF format) controversy engendered a Usenet comp.graphics discussion *Thoughts on a GIF-replacement file format*, which in turn fostered an email exchange that eventually culminated in the creation of the patent-unencumbered Portable Network Graphics (PNG) file format in 1995. Unisys's US patent on the LZW algorithm expired on June 20, 2003 ..."

5

5

LZW Compression

```
p = nil;    // p for prev char
while read(c):
    if pc ∈ Dict:
        p = pc;
    else:
        add pc to Dict;
        output code(p);
        p = c;
```

6

6

Example

	p	c	output	Dict [index]	
				index	symbol
Input: ^WED^WE^WEE^WEB^WET	NIL	^			

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

7

7

Example

	p	c	output	Dict [index]	
				index	symbol
Input: ^WED^WE^WEE^WEB^WET	NIL	^			
	^	W		256	^W

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

8

8

Example

	p	c	output	Dict [index]	
				index	symbol
Input: ^WED^WE^WEE^WEB^WET	NIL	^			
	^	W		256	^W
	W	E	W	257	WE

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

9

9

Example

	p	c	output	Dict [index]	
				index	symbol
Input: ^WED^WE^WEE^WEB^WET	NIL	^			
	^	W		256	^W
	W	E	W	257	WE
	E	D	E	258	ED
	D	^	D	259	D^
	^	W			

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

10

10

Example

	p	c	output	Dict [index]	
				index	symbol
Input: ^WED^WE^WEE^WEB^WET	NIL	^			
	^	W		256	^W
	W	E	W	257	WE
	E	D	E	258	ED
	D	^	D	259	D^
	^	W			
	^W	E	256	260	^WE

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

11

11

Example

	p	c	output	Dict [index]	
				index	symbol
Input: ^WED^WE^WEE^WEB^WET	NIL	^			
	^	W		256	^W
	W	E	W	257	WE
	E	D	E	258	ED
	D	^	D	259	D^
	^	W			
	^W	E	256	260	^WE
	E	^	E	261	E^

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

12

12

Example

Input: ^WED^WE^WEE^WEB^WET

p	c	output	Dict [index]	
			index	symbol
NIL	^		256	^W
^	W		257	WE
W	E		258	ED
E	D		259	D^
D	^			
^W	E	256	260	^WE
E	^		261	E^
^	W			

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

13

Example

Input: ^WED^WE^WEE^WEB^WET

p	c	output	Dict [index]	
			index	symbol
NIL	^		256	^W
^	W		257	WE
W	E		258	ED
E	D		259	D^
D	^			
^W	E	256	260	^WE
E	^		261	E^
^	W			
^W	E			

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

14

13

14

Example

Input: ^WED^WE^WEE^WEB^WET

p	c	output	Dict [index]	
			index	symbol
NIL	^		256	^W
^	W		257	WE
W	E		258	ED
E	D		259	D^
D	^			
^W	E	256	260	^WE
E	^		261	E^
^	W			
^W	E			
^WE	E	260	262	^WEE

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

15

15

Example

Input: ^WED^WE^WEE^WEB^WET

p	c	output	Dict [index]	
			index	symbol
NIL	^		256	^W
^	W		257	WE
W	E		258	ED
E	D		259	D^
D	^			
^W	E	256	260	^WE
E	^		261	E^
^	W			
^W	E			
^WE	E	260	262	^WEE
E	^			

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

16

16

Example

Input: ^WED^WE^WEE^WEB^WET

p	c	output	Dict [index]	
			index	symbol
NIL	^		256	^W
^	W		257	WE
W	E		258	ED
E	D		259	D^
D	^			
^W	E	256	260	^WE
E	^		261	E^
^	W			
^W	E			
^WE	E	260	262	^WEE
E	^			
E^	W	261	263	E^W
W	E		264	WEB
WE	B	257	265	B^
B	^			
^W	E			
^WE	T	260	266	^WET
T	EOF			

p = nil; // p for prev char
 while read(c);
 if pc ∈ Dict:
 p = pc;
 else:
 add pc to Dict;
 output code(p);
 p = c;

17

17

LZW Compression

- Original LZW used dictionary with 4K entries, first 256 (0-255) are ASCII codes.
- In the above example, a 19 symbols reduced to 7 symbols & 5 code. Each code/symbol will need 8+ bits, say 9 bits.
- Reference: Terry A. Welch, "A Technique for High Performance Data Compression", IEEE Computer, Vol. 17, No. 6, 1984, pp. 8-19.

18

18

LZW Decompression

```
read(c); // c is likely to be > 8 bits
output c;
p = c;
while read(c):
    output Dict[c];
    add p + Dict[c][0] to Dict;
    p = Dict[c];
```

19

19

Example

Input: ^WED<256>E<260><261><257>B<260>T

```
test(c); // c is likely to be > 8 bits
output c;
p = c;
while read(c):
    output Dict[c];
    add p + Dict[c][0] to Dict;
    p = Dict[c];
```

p	c	output	Dict [index]	
			index	symbol

20

20

Example

Input: ^WED<256>E<260><261><257>B<260>T

```
test(c); // c is likely to be > 8 bits
output c;
p = c;
while read(c):
    output Dict[c];
    add p + Dict[c][0] to Dict;
    p = Dict[c];
```

p	c	output	Dict [index]	
			index	symbol
			256	^W

21

21

Example

Input: ^WED<256>E<260><261><257>B<260>T

```
test(c); // c is likely to be > 8 bits
output c;
p = c;
while read(c):
    output Dict[c];
    add p + Dict[c][0] to Dict;
    p = Dict[c];
```

p	c	output	Dict [index]	
			index	symbol
			256	^W
			257	WE

22

22

Example

Input: ^WED<256>E<260><261><257>B<260>T

```
test(c); // c is likely to be > 8 bits
output c;
p = c;
while read(c):
    output Dict[c];
    add p + Dict[c][0] to Dict;
    p = Dict[c];
```

p	c	output	Dict [index]	
			index	symbol
			256	^W
			257	WE
			258	ED

23

23

Example

Input: ^WED<256>E<260><261><257>B<260>T

```
test(c); // c is likely to be > 8 bits
output c;
p = c;
while read(c):
    output Dict[c];
    add p + Dict[c][0] to Dict;
    p = Dict[c];
```

p	c	output	Dict [index]	
			index	symbol
			256	^W
			257	WE
			258	ED
			259	D^

24

24

Example

Input: ^WED<256>E<260><261><257>B<260>T

Dict [index]				
p	c	output	index	symbol
^	W	W	256	^W
W	E	E	257	WE
E	D	D	258	ED
D	<256>	^W	259	D^
<256>	E	E	260	^WE

25

25

Example

Input: ^WED<256>E<260><261><257>B<260>T

Dict [index]				
p	c	output	index	symbol
^	W	W	256	^W
W	E	E	257	WE
E	D	D	258	ED
D	<256>	^W	259	D^
<256>	E	E	260	^WE
E	<260>	^WE	261	E^

26

26

Example

Input: ^WED<256>E<260><261><257>B<260>T

Dict [index]				
p	c	output	index	symbol
^	W	W	256	^W
W	E	E	257	WE
E	D	D	258	ED
D	<256>	^W	259	D^
<256>	E	E	260	^WE
E	<260>	^WE	261	E^
<260>	<261>	E^	262	^WEE

27

27

Example

Input: ^WED<256>E<260><261><257>B<260>T

Dict [index]				
p	c	output	index	symbol
^	W	W	256	^W
W	E	E	257	WE
E	D	D	258	ED
D	<256>	^W	259	D^
<256>	E	E	260	^WE
E	<260>	^WE	261	E^
<260>	<261>	E^	262	^WEE
<261>	<257>	WE	263	E^W
<257>	B	B	264	WEB
B	<260>	^WE	265	B^
<260>	T	T	266	^WET

28

28

Note: LZW decoding

- There is one special case that the LZW decoding pseudocode presented is unable to handle.
- This is your exercise to find out in what situation that happens, and how to deal with it.
- I'll go through this at the live lecture.

29

29

LZW implementation

- Parsing fixed number of bits from input is easy
- Fast and efficient

30

30

Types (revision)

- Block-block
 - source message and codeword: fixed length
 - e.g., ASCII
- Block-variable
 - source message: fixed; codeword: variable
 - e.g., Huffman coding
- Variable-block
 - source message: variable; codeword: fixed
 - e.g., LZW
- Variable-variable
 - source message and codeword: variable
 - e.g., Arithmetic coding

31

31

So far

We have covered:

- Course overview
- Background
- RLE
- Entropy
- Huffman code
- Arithmetic code
- LZW

32

32

More online readings

<http://www.ics.uci.edu/~dan/pubs/DC-Sec1.html>
<http://marknelson.us/1991/02/01/arithmetic-coding-statistical-modeling-data-compression/>

33

33

Lossless compression revisited

- Run-length coding
 - Huffman coding
 - Arithmetic coding
- Dictionary methods
 - Lempel Ziv algorithms

Static (Huffman, AC) vs Adaptive (LZW)

34

34

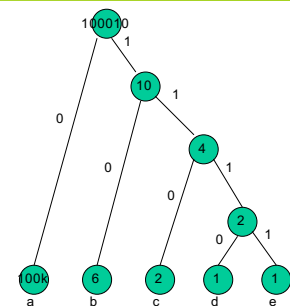
Huffman Coding (revisit) and then Adaptive Huffman

35

35

Huffman coding

S	Freq	Huffman
a	100000	0
b	6	10
c	2	110
d	1	1110
e	1	1111



36

36

Huffman not optimal

$$H = 0.9999 \log 1.0001 + 0.00006 \log 16668.333 \\ + \dots + 1/100010 \log 100010 \\ \approx 0.00$$

$$L = (100000 \cdot 1 + \dots) / 100010 \\ \approx 1$$

37

37

Problems of Huffman coding

Huffman codes have an integral # of bits.

E.g., $\log(3) = 1.585$ while Huffman may need 2 bits

Noticeable non-optimality when prob of a symbol is high.

=> Arithmetic coding

38

38

Problems of Huffman coding

Need statistics & static: e.g., single pass over the data just to collect stat & stat unchanged during encoding

To decode, the stat table need to be transmitted. Table size can be significant for small msg.

=> Adaptive compression e.g., adaptive huffman

39

39

Adaptive compression

Encoder

Initialize the model

Repeat for each input char

```
(
  Encode char
  Update the model
)
```

Decoder

Initialize the model

Repeat for each input char

```
(
  Decode char
  Update the model
)
```

Make sure both sides have the same Initialize & update model algorithms.

40

40

Adaptive Huffman Coding (dummy)

Encoder

Reset the stat

Repeat for each input char

```
(
  Encode char
  Update the stat
  Rebuild huffman tree
)
```

Decoder

Reset the stat

Repeat for each input char

```
(
  Decode char
  Update the stat
  Rebuild huffman tree
)
```

41

41

Adaptive Huffman Coding (dummy)

Encoder

Reset the stat

Repeat for each input char

```
(
  Encode char
  Update the stat
  Rebuild huffman tree
)
```

Decoder

Reset the stat

Repeat for each input char

```
(
  Decode char
  Update the stat
  Rebuild huffman tree
)
```

This works but too slow!

42

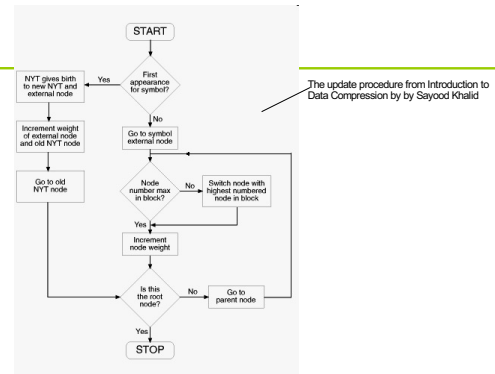
42

Adaptive Huffman (Algorithm outline)

1. If current symbol is NYT, add two child nodes to NYT node. One will be a new NYT node the other is a leaf node for our symbol. Increase weight for the new leaf node and the old NYT and go to step 4. If not, go to symbol's leaf node.
2. If this node does not have the highest number in a block, swap it with the node having the highest number
3. Increase weight for current node
4. If this is not the root node go to parent node then go to step 2. If this is the root, end.

43

43

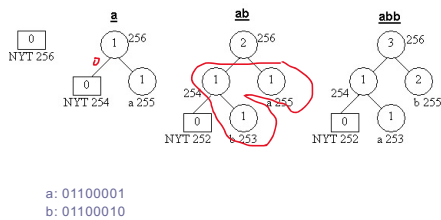


44

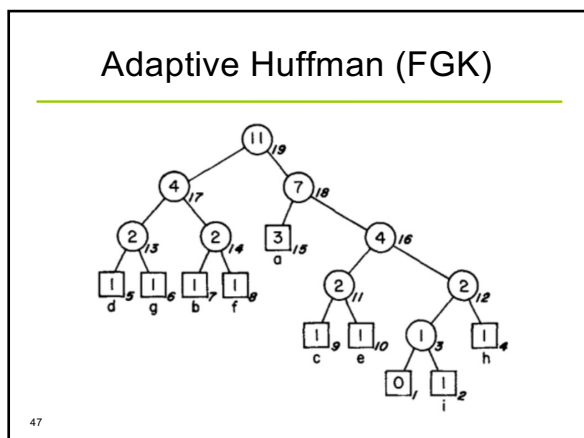
44

Adaptive Huffman

abbbba: 011000010110001001100010011000100110001001100001
 abbbba: 011000010011000100111101



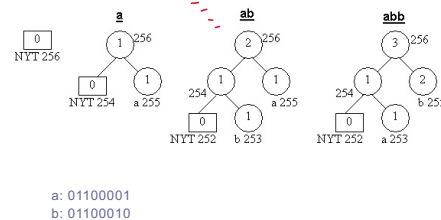
45



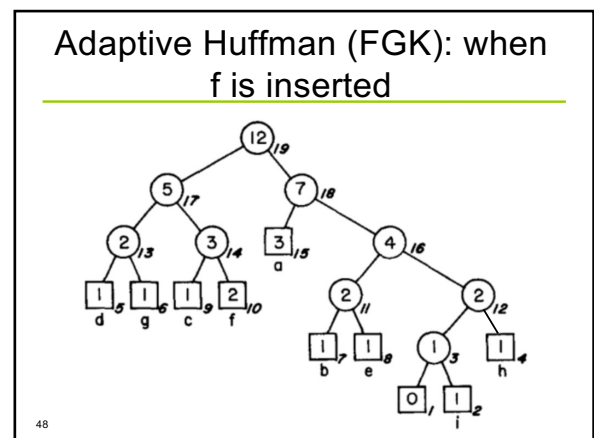
47

Adaptive Huffman

abbbba: 011000010110001001100010011000100110001001100001
 abbbba: 011000010011000100111101



46



48

Adaptive Huffman (FGK vs Vitter)

1.

FGK: (Explicit) node numbering

Vitter: Implicit numbering

2

Vitter's Invariant:

- (*) For each weight w , all leaves of weight w precede (in the implicit numbering) all internal nodes of weight w .

49

49

Adaptive Huffman (Vitter'87)

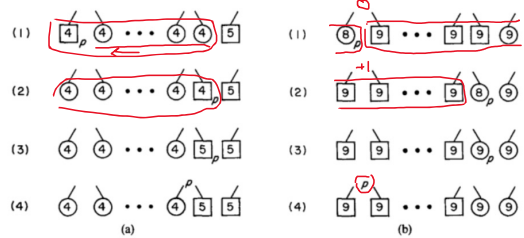
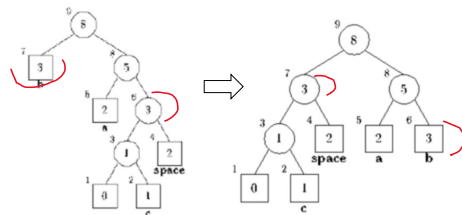


FIG. 6. Algorithm A's *SlideAndIncrement* operation. All the nodes in a given block shift to the left one spot to make room for node p , which slides over the block to the right. (a) Node p is a leaf of weight 4. The internal nodes of weight 4 shift to the left. (b) Node p is an internal node of weight 8. The leaves of weight 9 shift to the left.

50

50

Adaptive Huffman (Vitter's Invariant)



51

51

Issues with Wikipedia

[illegible]

52

52

COMP9319 students correcting lots of Wiki pages

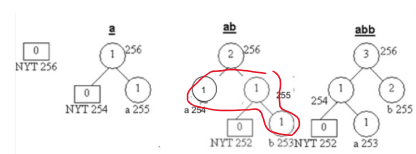
- `(cor >= 0.95 && March 2016)` **Redundant task (removed)** (`h` = 0.8620909) (`h1` = 0.8620909) (`h2` = 0.8620909) (`h3` = 0.8620909) (`h4` = 0.8620909) (`h5` = 0.8620909) (`h6` = 0.8620909) (`h7` = 0.8620909) (`h8` = 0.8620909) (`h9` = 0.8620909) (`h10` = 0.8620909) (`h11` = 0.8620909) (`h12` = 0.8620909) (`h13` = 0.8620909) (`h14` = 0.8620909) (`h15` = 0.8620909) (`h16` = 0.8620909) (`h17` = 0.8620909) (`h18` = 0.8620909) (`h19` = 0.8620909) (`h20` = 0.8620909) (`h21` = 0.8620909) (`h22` = 0.8620909) (`h23` = 0.8620909) (`h24` = 0.8620909) (`h25` = 0.8620909) (`h26` = 0.8620909) (`h27` = 0.8620909) (`h28` = 0.8620909) (`h29` = 0.8620909) (`h30` = 0.8620909) (`h31` = 0.8620909) (`h32` = 0.8620909) (`h33` = 0.8620909) (`h34` = 0.8620909) (`h35` = 0.8620909) (`h36` = 0.8620909) (`h37` = 0.8620909) (`h38` = 0.8620909) (`h39` = 0.8620909) (`h40` = 0.8620909) (`h41` = 0.8620909) (`h42` = 0.8620909) (`h43` = 0.8620909) (`h44` = 0.8620909) (`h45` = 0.8620909) (`h46` = 0.8620909) (`h47` = 0.8620909) (`h48` = 0.8620909) (`h49` = 0.8620909) (`h50` = 0.8620909) (`h51` = 0.8620909) (`h52` = 0.8620909) (`h53` = 0.8620909) (`h54` = 0.8620909) (`h55` = 0.8620909) (`h56` = 0.8620909) (`h57` = 0.8620909) (`h58` = 0.8620909) (`h59` = 0.8620909) (`h60` = 0.8620909) (`h61` = 0.8620909) (`h62` = 0.8620909) (`h63` = 0.8620909) (`h64` = 0.8620909) (`h65` = 0.8620909) (`h66` = 0.8620909) (`h67` = 0.8620909) (`h68` = 0.8620909) (`h69` = 0.8620909) (`h70` = 0.8620909) (`h71` = 0.8620909) (`h72` = 0.8620909) (`h73` = 0.8620909) (`h74` = 0.8620909) (`h75` = 0.8620909) (`h76` = 0.8620909) (`h77` = 0.8620909) (`h78` = 0.8620909) (`h79` = 0.8620909) (`h80` = 0.8620909) (`h81` = 0.8620909) (`h82` = 0.8620909) (`h83` = 0.8620909) (`h84` = 0.8620909) (`h85` = 0.8620909) (`h86` = 0.8620909) (`h87` = 0.8620909) (`h88` = 0.8620909) (`h89` = 0.8620909) (`h90` = 0.8620909) (`h91` = 0.8620909) (`h92` = 0.8620909) (`h93` = 0.8620909) (`h94` = 0.8620909) (`h95` = 0.8620909) (`h96` = 0.8620909) (`h97` = 0.8620909) (`h98` = 0.8620909) (`h99` = 0.8620909) (`h100` = 0.8620909) (`h101` = 0.8620909) (`h102` = 0.8620909) (`h103` = 0.8620909) (`h104` = 0.8620909) (`h105` = 0.8620909) (`h106` = 0.8620909) (`h107` = 0.8620909) (`h108` = 0.8620909) (`h109` = 0.8620909) (`h110` = 0.8620909) (`h111` = 0.8620909) (`h112` = 0.8620909) (`h113` = 0.8620909) (`h114` = 0.8620909) (`h115` = 0.8620909) (`h116` = 0.8620909) (`h117` = 0.8620909) (`h118` = 0.8620909) (`h119` = 0.8620909) (`h120` = 0.8620909) (`h121` = 0.8620909) (`h122` = 0.8620909) (`h123` = 0.8620909) (`h124` = 0.8620909) (`h125` = 0.8620909) (`h126` = 0.8620909) (`h127` = 0.8620909) (`h128` = 0.8620909) (`h129` = 0.8620909) (`h130` = 0.8620909) (`h131` = 0.8620909) (`h132` = 0.8620909) (`h133` = 0.8620909) (`h134` = 0.8620909) (`h135` = 0.8620909) (`h136` = 0.8620909) (`h137` = 0.8620909) (`h138` = 0.8620909) (`h139` = 0.8620909) (`h140` = 0.8620909) (`h141` = 0.8620909) (`h142` = 0.8620909) (`h143` = 0.8620909) (`h144` = 0.8620909) (`h145` = 0.8620909) (`h146` = 0.8620909) (`h147` = 0.8620909) (`h148` = 0.8620909) (`h149` = 0.8620909) (`h150` = 0.8620909) (`h151` = 0.8620909) (`h152` = 0.8620909) (`h153` = 0.8620909) (`h154` = 0.8620909) (`h155` = 0.8620909) (`h156` = 0.8620909) (`h157` = 0.8620909) (`h158` = 0.8620909) (`h159` = 0.8620909) (`h160` = 0.8620909) (`h161` = 0.8620909) (`h162` = 0.8620909) (`h163` = 0.8620909) (`h164` = 0.8620909) (`h165` = 0.8620909) (`h166` = 0.8620909) (`h167` = 0.8620909) (`h168` = 0.8620909) (`h169` = 0.8620909) (`h170` = 0.8620909) (`h171` = 0.8620909) (`h172` = 0.8620909) (`h173` = 0.8620909) (`h174` = 0.8620909) (`h175` = 0.8620909) (`h176` = 0.8620909) (`h177` = 0.8620909) (`h178` = 0.8620909) (`h179` = 0.8620909) (`h180` = 0.8620909) (`h181` = 0.8620909) (`h182` = 0.8620909) (`h183` = 0.8620909) (`h184` = 0.8620909) (`h185` = 0.8620909) (`h186` = 0.8620909) (`h187` = 0.8620909) (`h188` = 0.8620909) (`h`

53

53

Adaptive Huffman (Vitter's)

```
abbbba: 011000010110001001100010011000100110001001100001
abbbba: 011000010011000101111101
```

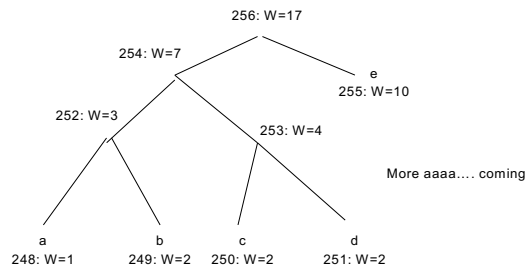


a: 01100001
b: 01100010

Corrected fr Wikipedia

54

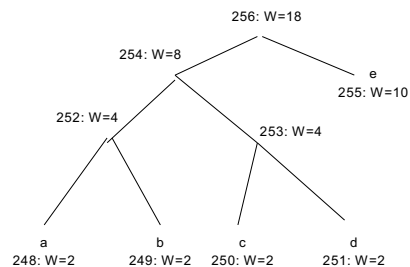
More example on Vitter's



55

55

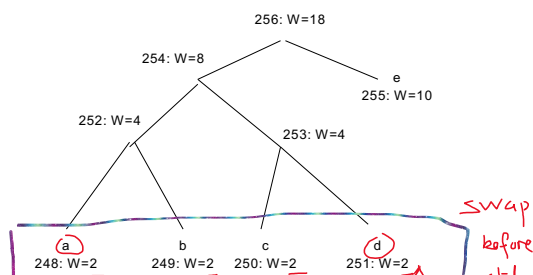
More example



56

56

More example



57

57

The Vitter's algo: swaps then slides

Definition 4.1. Blocks are equivalence classes of nodes defined by $v = x$ iff nodes v and x have the same weight and are either both internal nodes or both leaves. The leader of a block is the highest numbered (in the implicit numbering) node in a block.

The blocks are linked together by increasing order of weight; a leaf block always precedes an internal block of the same weight. The main operation of the algorithm needed to maintain invariant (*) is the SlideAndIncrement operation, illustrated in Figure 6. The version of Update we use for Algorithm A is outlined below:

```

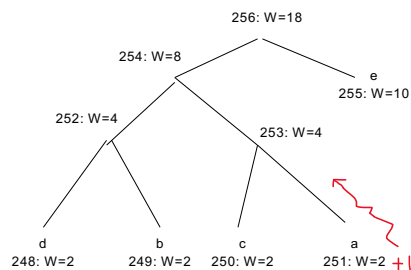
procedure Update;
begin
  leafToIncrement := 0;
  q := leaf node corresponding to  $a_{n-1}$ ;
  if (q is the 0-node) and ( $k < n-1$ ) then
    begin
      [Special Case #1]
      Replace q by an internal 0-node with two leaf 0-node children, such that the right child
        corresponds to  $a_{n-1}$ ;
      q := internal 0-node just created;
      leafToIncrement := the right child of q
    end
  else begin
    Interchange q in the tree with the leader of its block;
    if q is the sibling of the 0-node then
      begin
        [Special Case #2]
        leafToIncrement := q;
        q := parent of q
      end
    end;
    while q is not the root of the Huffman tree do
      [Main loop: q must be the leader of its block]
      SlideAndIncrement(q);
      if leafToIncrement  $\neq 0$  then
        SlideAndIncrement(leafToIncrement)
      end;
  end;
end;

```

58

58

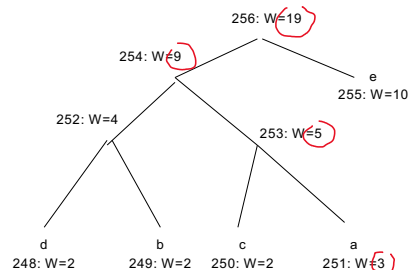
More example



59

59

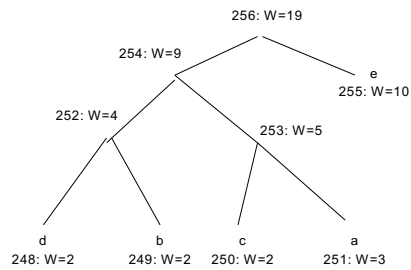
More example



60

60

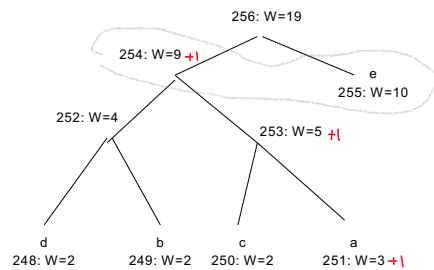
More example



61

61

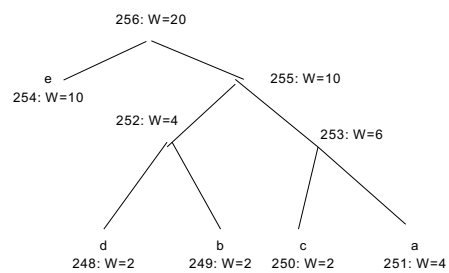
More example



62

62

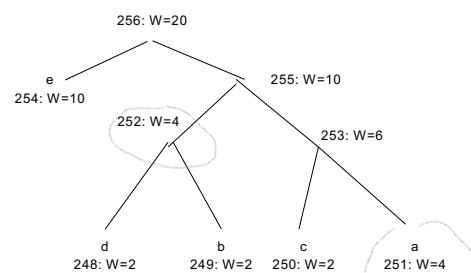
More example



63

63

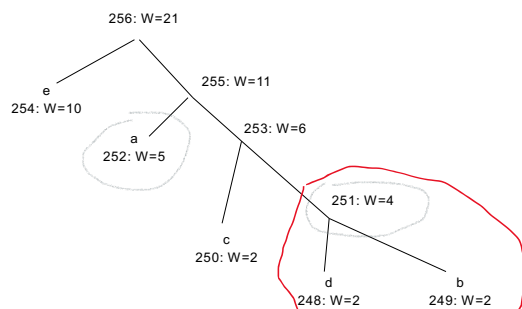
More example



64

64

More example



65

65

Adaptive Huffman

Question: Adaptive Huffman vs Static Huffman

66

66

Compared with Static Huffman

Dynamic and can offer better compression
(cf. Vitter's experiments next)

Works when prior stat is unavailable

Saves symbol table overhead (cf. Vitter's
expt next)

67

67

Vitter's experiments

Include overheads such as symbol tables / leaf node code etc.

t	k	S_t	b/l	D_t^A	b/l
100	96	664	13.1	569	10.2
500	96	3320	7.9	3225	7.4
960	96	6400	7.1	6305	6.8

Exclude overheads such as symbol tables / leaf node code etc.

95 ASCII chars + <end-of-line>

68

From Vitter's paper. You know where it is. ©

68

More experiments

t	k	S_t	b/l	D_t^A	b/l
100	34	434	7.1	420	6.3
500	52	2429	5.7	2445	5.5
1000	58	4864	5.3	4900	5.2
10000	74	47710	4.8	47852	4.8
12280	76	58457	4.8	58614	4.8

69

69