

COMP9319 Web Data Compression and Search

LZW,
Adaptive Huffman

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

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

<https://eduassistpro.github.io/>

3

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

4

Hempel-Ziv-Welch (LZW) Algorithm

- Most popular modification to LZ78
 - Unix compress, TIFF, etc.)
- regarding its patents

s (12bit 4096
reached

Patent issues again

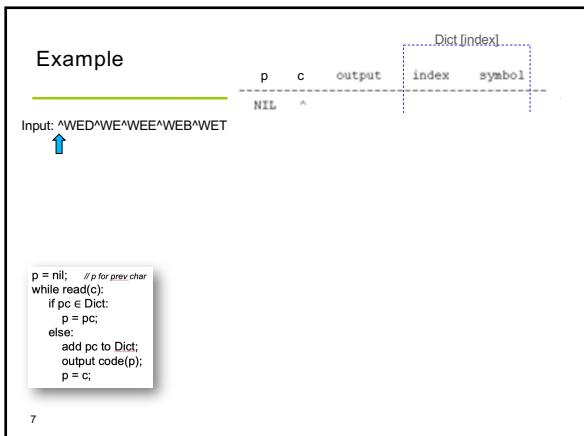
From Wikipedia: "In 1993–94, and again in 1999, Unisys Corporation received widespread condemnation when it attempted to enforce licensing fees for LZW in GIF images. The 1993–1994 Unisys-Compuserve (Compuserve being the creator of the GIF format) controversy engendered a Usenet comp.graphics discussion *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

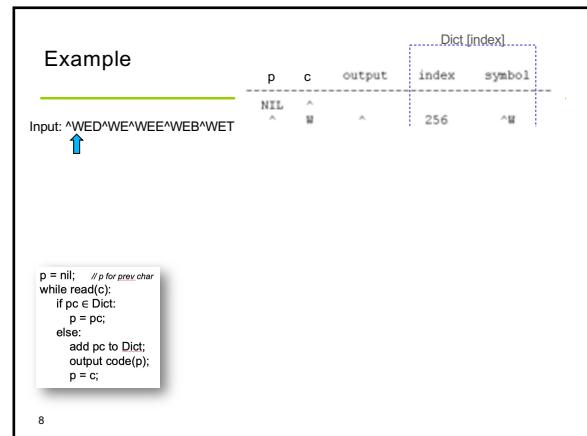
LZW Compression

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

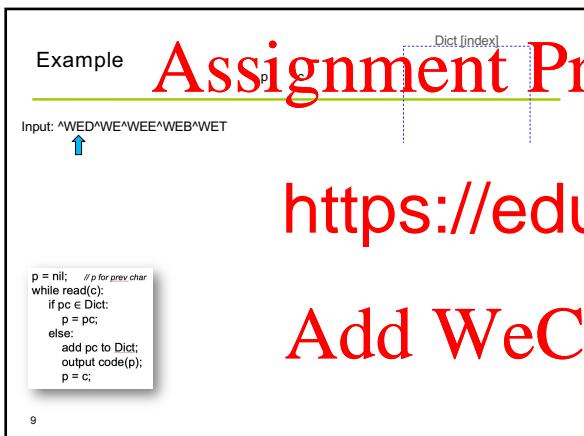
6



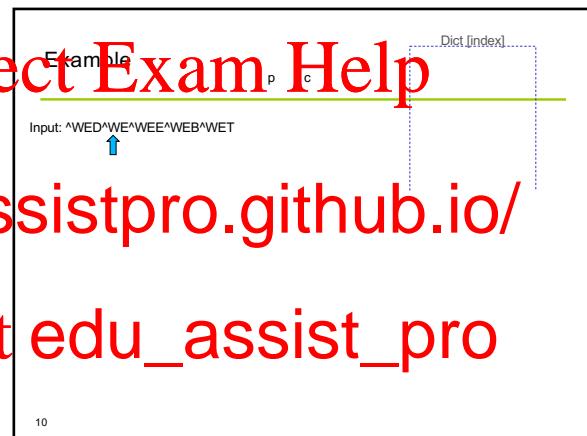
7



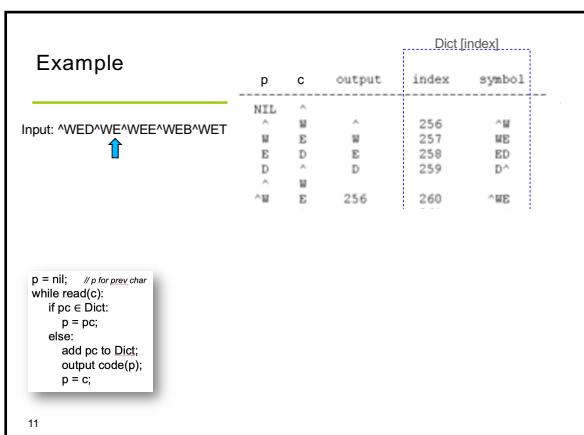
8



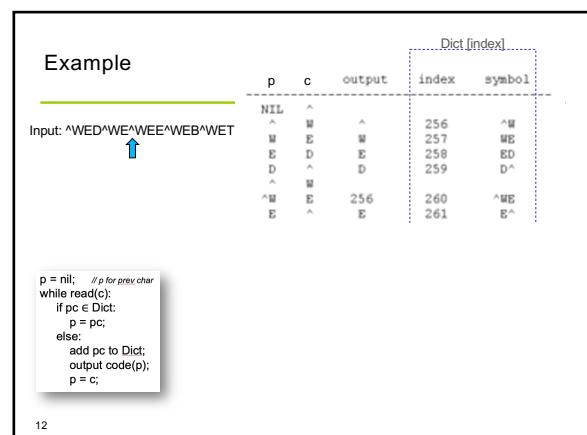
9



10



11



12

Example

p	c	output	index	symbol
NIL	^		256	^W
W	E	W	257	WE
E	D	E	258	ED
D	^	D	259	D^
^	W		260	^WE
W	^	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;
```

13

Example

p	c	output	index	symbol
NIL	^		256	^W
W	E	W	257	WE
E	D	E	258	ED
D	^	D	259	D^
^	W		260	^WE
W	^	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;
```

14

13 14

Example

p	c	output	index	symbol
NIL	^		256	^W

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

```
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

Example

p	c	output	index	symbol
NIL	^		256	^W

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

```
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

15 16

Example

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

```
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

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

17 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

Example

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

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

20

Example

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

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

21

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

22

Example

Input: ^WED<260>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
			<260>	

23

Example

Input: ^WED<260>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		<260>	259	D^

24

Example

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

```

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

```

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

Example

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

```

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

```

p	c	output	index	symbol	Dict [index]
^	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^	

25

26

Example

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

```

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

```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat `edu_assist_pro`

Example

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

27

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.

LZW implementation

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

29

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

So far

We have covered:

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

32

Assignment Project Exam Help

More online readings:
<http://www.ics.uci.edu/~dan/pubs/DC-Sec1.html>
<http://marknelson.us/1991/02/01/arithmetic-coding/>

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)

33

34

Huffman Coding (revisit)
and then
Adaptive Huffman

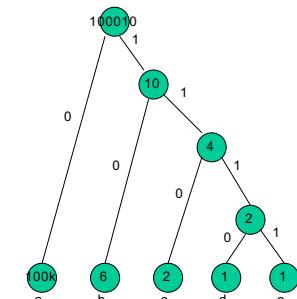
35

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

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

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

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

Problems of Huffman coding

Need statistics & static: e.g., single pass over the data just to unchanged during To decode, the stat table transmitted. Table for small msg.

=> Adaptive compression: e.g., adaptive huffman

39

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Adaptive compression

Encoder

Initialize the model
Repeat for each input char
 Encode char
 Update the model

Update the model

Decoder

Initialize the model
Repeat for each input char
 Decode char
 Update the model
 me Initialize & update model

40

39

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

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

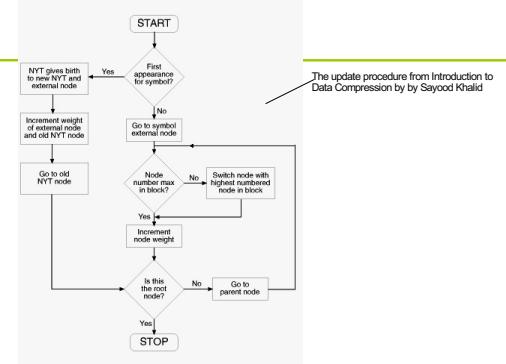
41

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



44

43

44

Adaptive Huffman

abbbbba: 01100001011000100110001001100010011000100110001
abbbbba: 01100001001100010011110

<https://eduassistpro.github.io/>
Add WeChat [edu_assist_pro](#)

a: 01100001
b: 01100010

From an old Wikipedia page

Adaptive Huffman

abbbbba: 0110000101100010011000100110001001100010011000100110001

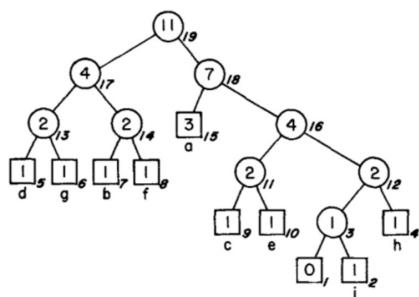
101

<https://eduassistpro.github.io/>
Add WeChat [edu_assist_pro](#)

a: 01100001
b: 01100010

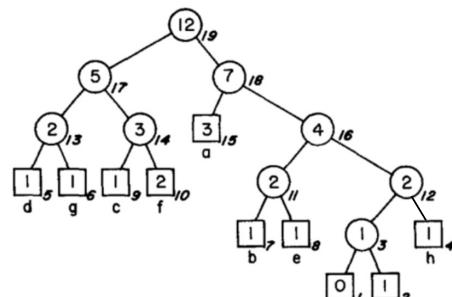
From an old Wikipedia page

Adaptive Huffman (FGK)



47

Adaptive Huffman (FGK): when f is inserted



48

47

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

Adaptive Huffman (Vitter'87)

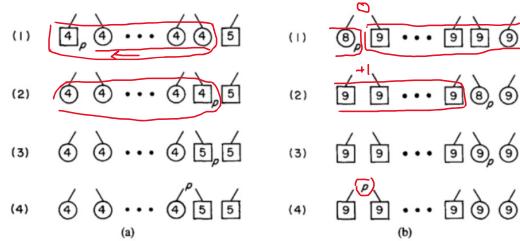


FIG. 6. Algorithm A's *SlideAndIncrement* operation. All the nodes in a given block shift to the left one spot to make room for node p , which slides over the block to the right. (a) Node p is a leaf of weight 4. The internal nodes of weight 4 shift to the left. (b) Node p is an internal node of weight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

50

Adaptive Huffman (Vitter's Invariant)

<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

51

Issues with Wikipedia

52

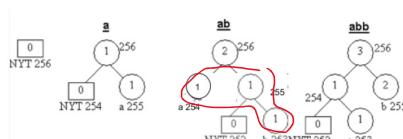
COMP9319 students correcting lots of Wiki pages

A screenshot of a GitHub commit history showing many contributions from COMP9319 students to correct issues on Wikipedia. The commits range from March 2016 to March 2018, addressing various bugs and improving descriptions of Vitter's algorithm.

53

Adaptive Huffman (Vitter's)

abbbbba: 0110001011000100110001001100010011000100110001
abbbbbba: 0110001001100010**1111101**



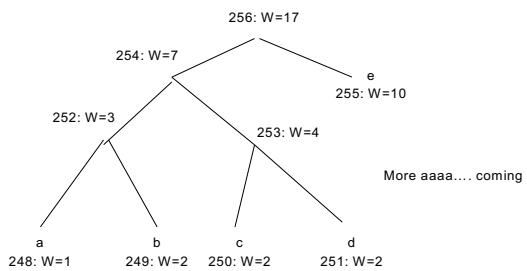
a: 0110001

b: 01100010

Corrected fr Wikipedia

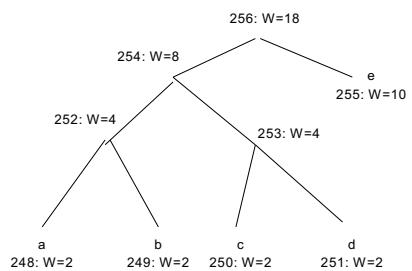
54

More example on Vitter's



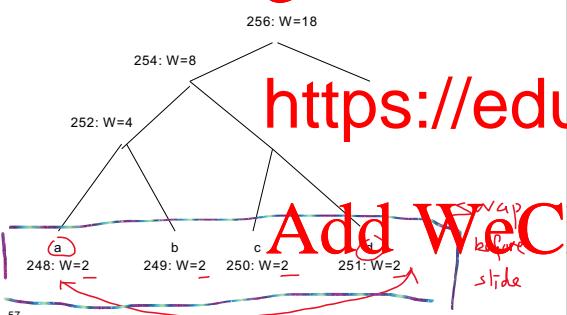
55

More example



56

More example

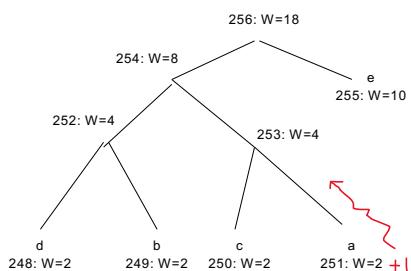


57

The Vitter's algo: swaps then slides

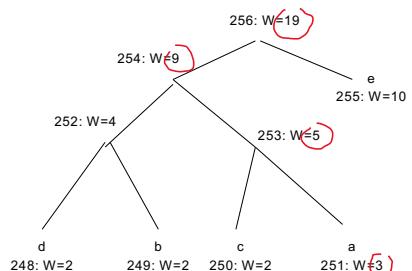
58

More example



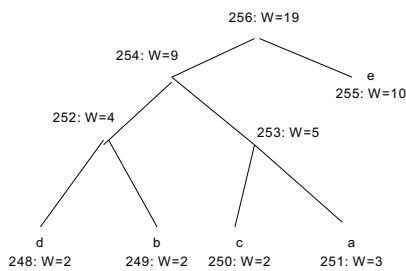
59

More example



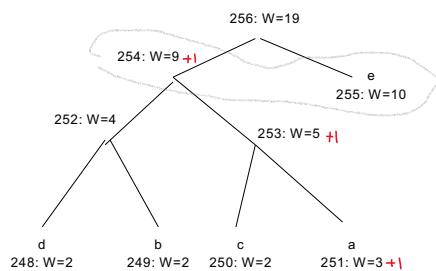
60

More example



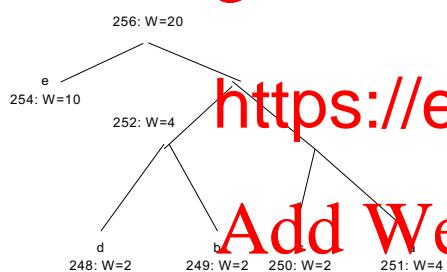
61

More example



62

More example



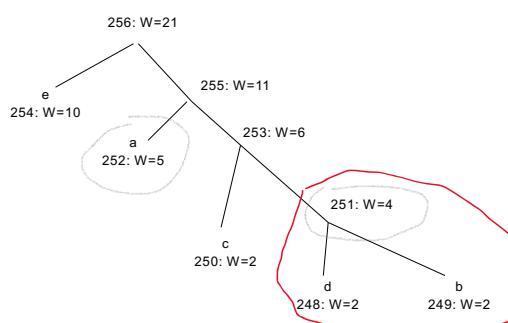
63

More example



64

More example



65

Adaptive Huffman

Question: Adaptive Huffman vs Static Huffman

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

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.

68

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

More experiments

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

69

69