

LZW Compression

Outline

- 0 Theory part of LZW compression/decompression
- 0 About assignment 2 and some details

Introduction

0 What is LZW?

- 0 Lossless compression method
- 0 Lempel-Ziv-Welch
- 0 Based on LZ78

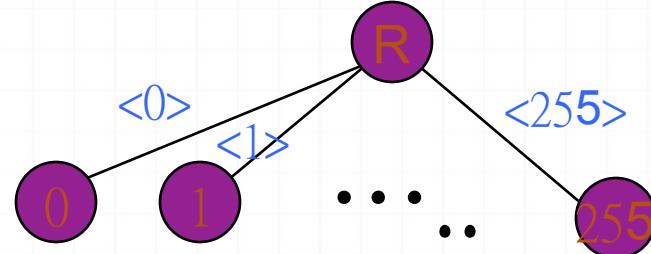
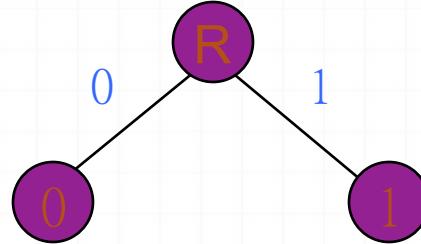
Basic idea

- 0 Assume repetition of phase usually occurs
- 0 Use a code to represent one phase
- 0 Build a dictionary of phases that we met
 - 0 If a phase is found in dictionary, use the code
 - 0 If not found, add it to dictionary and give it a code
- 0 Very high compression ratio if lots of repetition

Skip

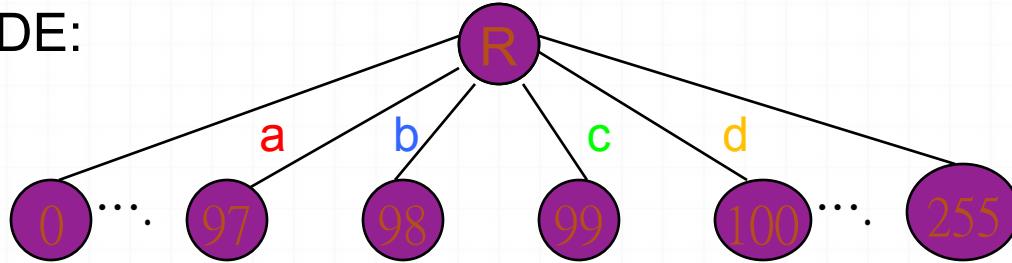
Algorithm

- 0 But how to handle byte steam?
- 0 The algorithm is similar to lecture note!
- 0 However, each node have maximum number of 256 child nodes.



Tree Structure Example

ENCODER SIDE:

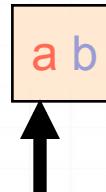
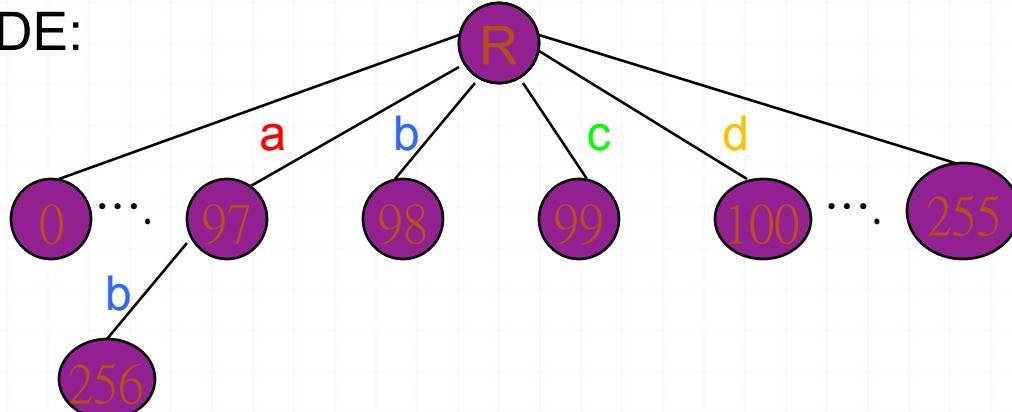


a b a b c d



Tree Structure Example

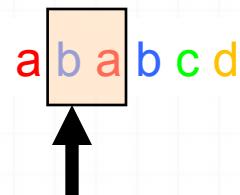
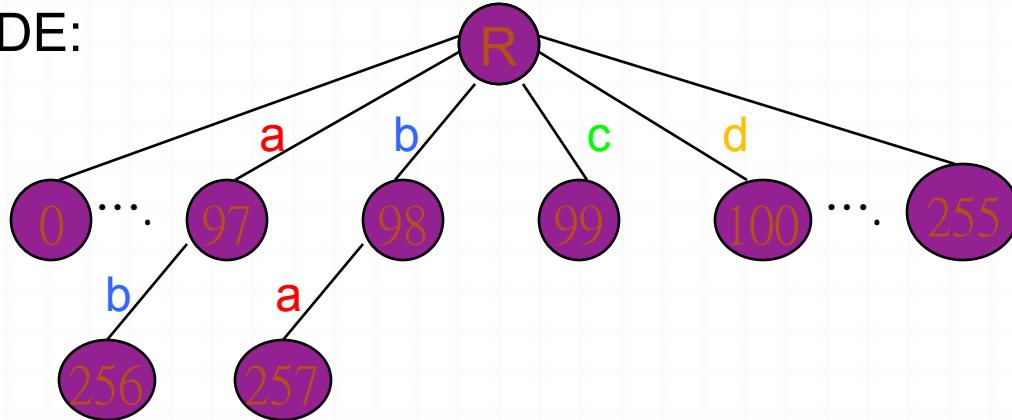
ENCODER SIDE:



output: 97,

Tree Structure Example

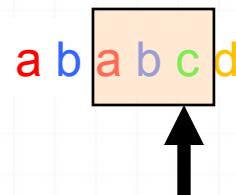
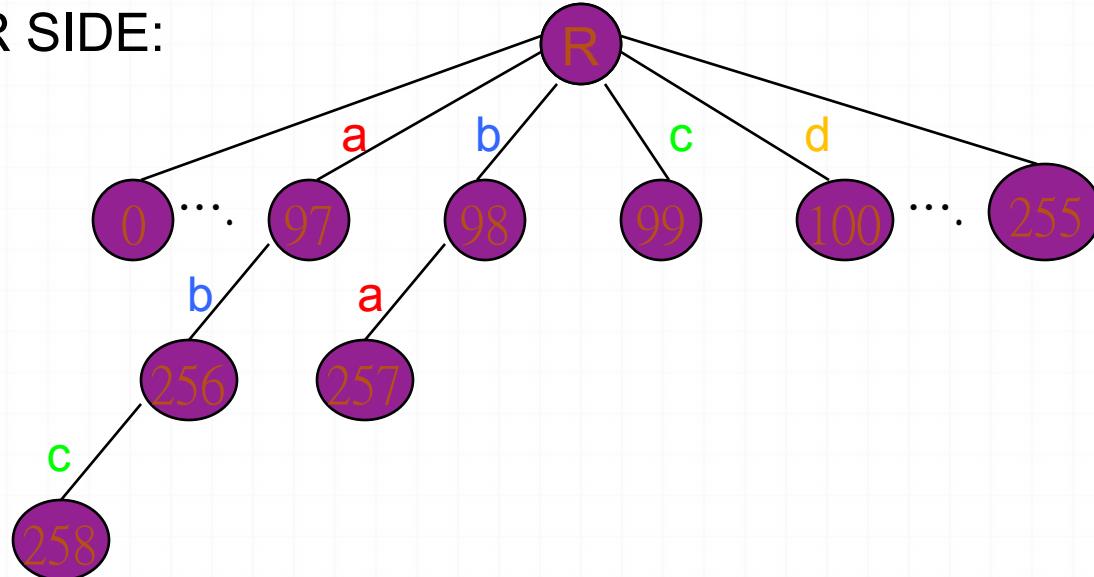
ENCODER SIDE:



output: 97,98

Tree Structure Example

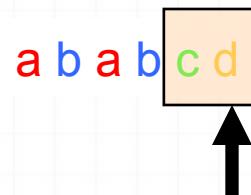
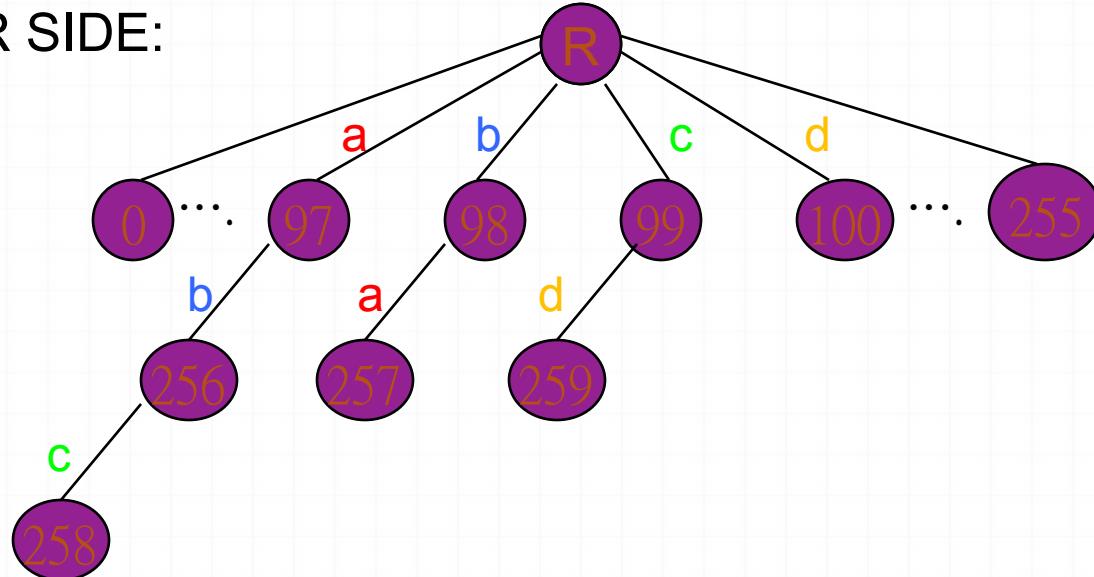
ENCODER SIDE:



output: 97,98,256

Tree Structure Example

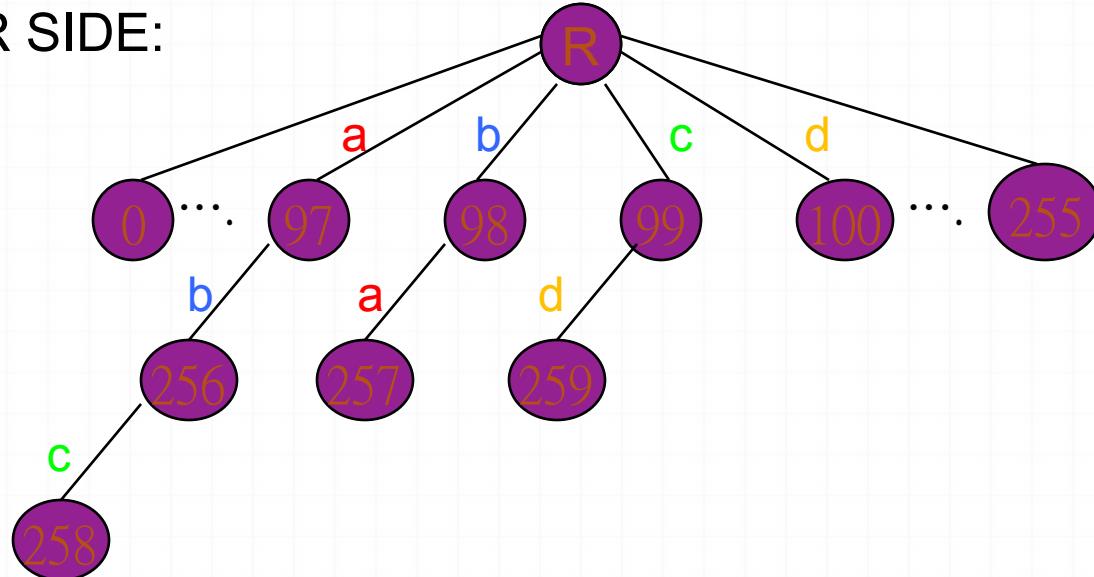
ENCODER SIDE:



output: 97,98,256,99

Tree Structure Example

ENCODER SIDE:



a b a b c d

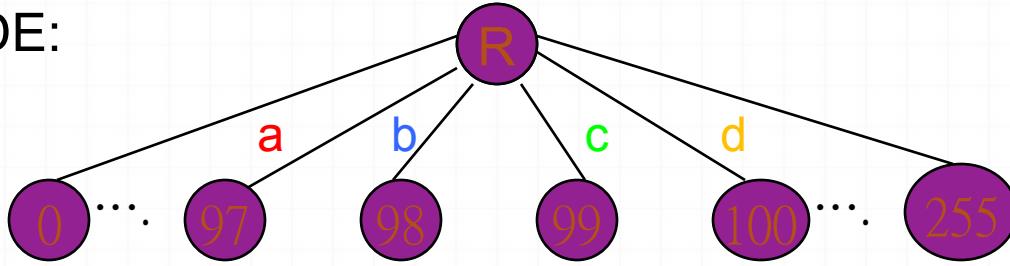


output: 97,98,256,99,100

Encode complete!

Tree Structure Example

DECODE SIDE:



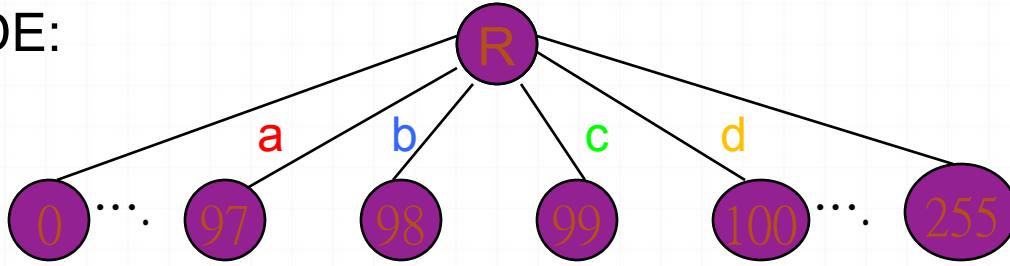
input: 97,98,256,99,100



output:

Tree Structure Example

DECODE SIDE:



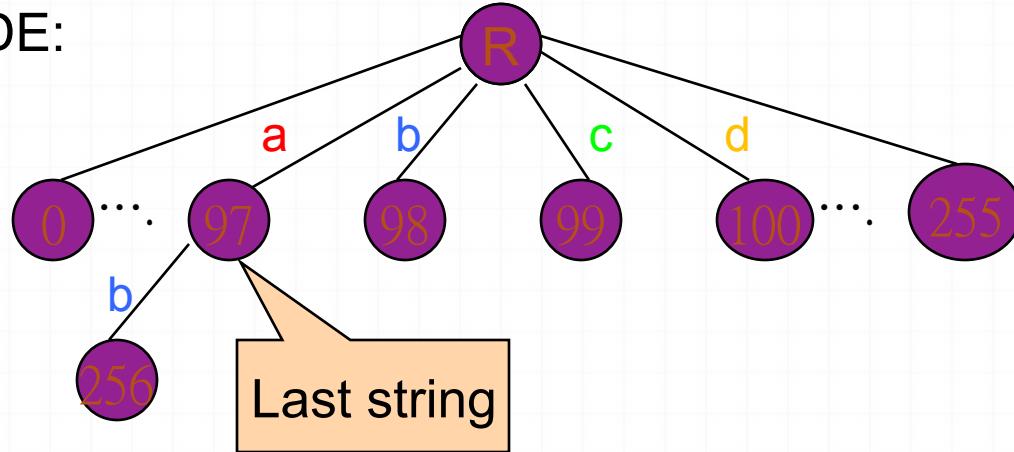
input: 97,98,256,99,100



output: a

Tree Structure Example

DECODE SIDE:



input: 97,98,256,99,100

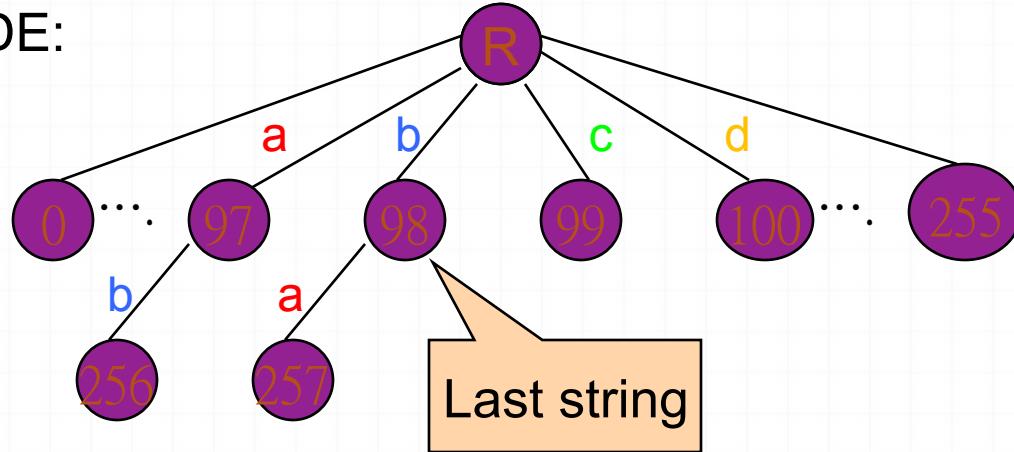


output: a b

Last string = a

Tree Structure Example

DECODE SIDE:



input: 97,98,256,99,100

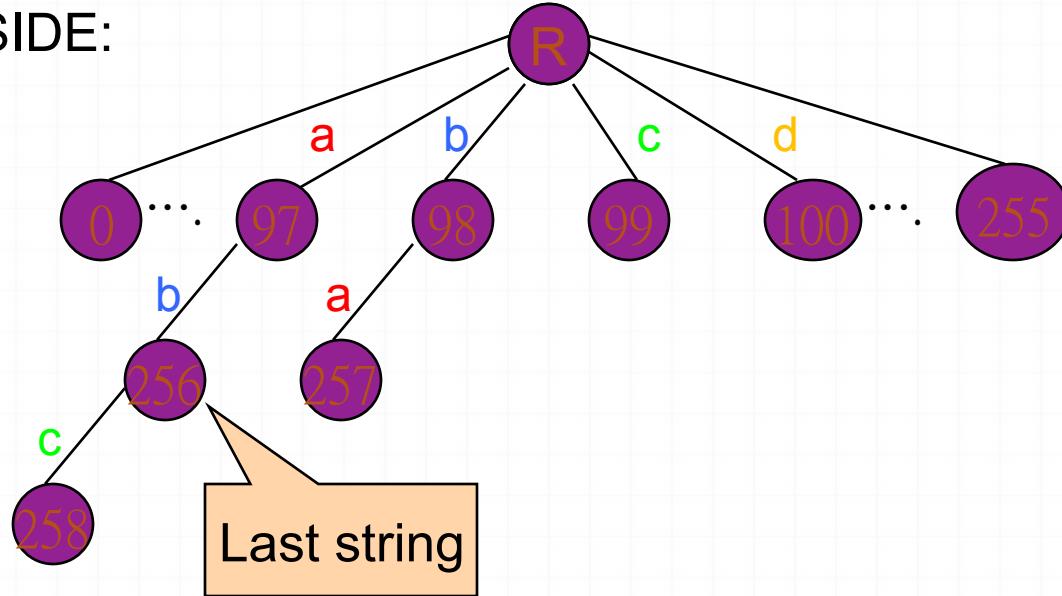


output: a b a b

Last string = b

Tree Structure Example

DECODE SIDE:



input: 97,98,256,99,100

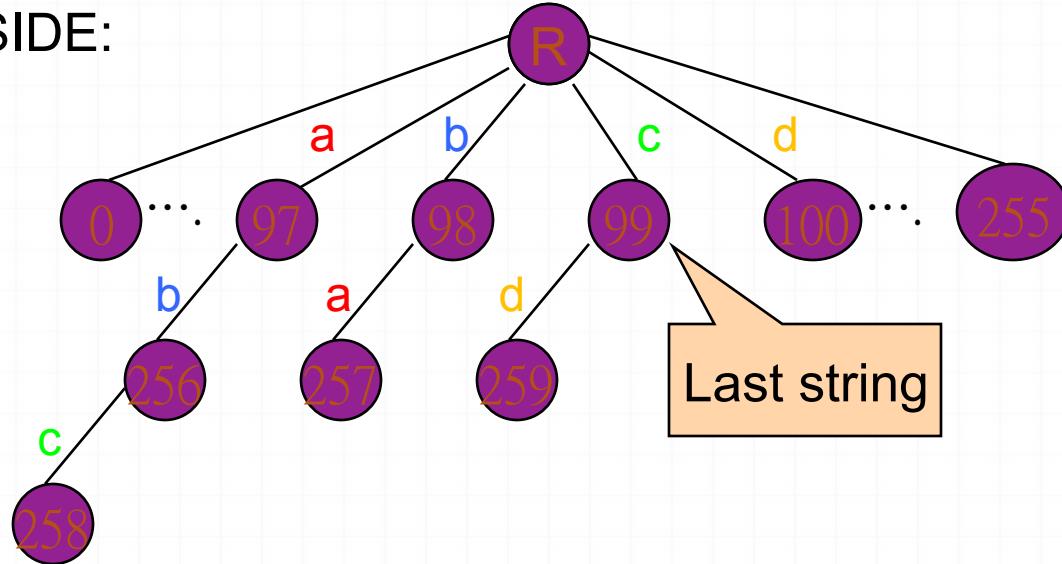


output: a b a b c

Last string = a b

Tree Structure Example

DECODE SIDE:



input: 97,98,256,99,100



output: a b a b c d

Last string = c

Algorithm

0 Now let's see a table structure example

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL			
L			

OUT:

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL	'a'	"a"	97
"a"	'b'	"ab"	

OUT: 97

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
	,"

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL	'b'	"b"	98
'b'	'c'	"bc"	

OUT: 97, 98

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
	,

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL	'c'	"c"	99
'c'	'c'	"cc"	

OUT: 97, 98, 99

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	" cc "
	" "
	" "
	" "

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL	'c'	"c"	99
'c'	'c'	"cc"	258
'cc'	'd'	"ccd"	

OUT: 97, 98, 99, **258**

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	"cc"
259	"ccd"

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL	'd'	"d"	100
'd'	'c'	"dc"	

OUT: 97, 98, 99, 258, **100**

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	"cc"
259	"ccd"
260	"dc"

Compression

IN: a b c c c d c c d

Prefix	Char.	Search	Code Saved
NUL	'c'	"c"	99
'c'	'c'	"cc"	258
'cc'	'd'	"ccd"	259

OUT: 97, 98, 99, 258, 100, **259**

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	"cc"
259	"ccd"
260	"dc"

Reminder

1. Initialize the dictionary, and set P to an empty string
2. While there are still characters to read,
 - a. Get the next character, C, from the input text
 - b. Search for the string $\langle P, C \rangle$ in the dictionary
 - c. if found,
 - i. Append C to the end of the string P
 - d. if not found,
 - i. Output the code that corresponds to P
 - ii. Add to the dictionary the new entry $\langle P, C \rangle$
 - iii. $P := C$
3. Output the code that corresponds to P

Reminder

- 0 N-bit code → $2^N - 1$ phases
- 0 First 0-255 entries are all the single characters
- 0 Last phase $2^N - 1$ is reserved
 - 0 Denote “End-Of-File”
 - 0 Must be written as the last code of one file

Decompression

- 0 Reverse of the compression process
- 0 Look up the phase according to the code that we read
- 0 Update the dictionary on-the-fly

Decompression

IN: 97, 98, 99, 258, 100, 259

cW	pW	C	dict(pW)+C
97			

OUT: a

Decompression

IN: 97, 98, 99, 258, 100, 259

cW	pW	C	dict(pW)+C
98	97	'b'	"ab"

OUT: a b

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"

Decompression

IN: 97, 98, 99, 258, 100, 259

cW	pW	C	dict(pW)+C
99	98	'c	"bc"

OUT: a b **c**

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"

Decompression

IN: 97, 98, 99, 258, 100, 259

cW	pW	C	dict(pW)+C
258	99	'c	"cc"



OUT: a b c **cc**

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	"cc"

Decompression

IN: 97, 98, 99, 258, 100, 259

cW	pW	C	dict(pW)+C
100	258	'd'	"ccd"

OUT: a b c c c d

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	"cc"
259	"ccd"

Decompression

IN: 97, 98, 99, 258, 100, 259

cW	pW	C	dict(pW)+C
259	100	'c	"dc"

OUT: a b c c c d **c c d**

CODE	Entry
0	NUL
:	
97	a
98	b
99	c
100	d
:	:
255	ASCII-255
256	"ab"
257	"bc"
258	"cc"
259	"ccd"
260	"dc"

Reminder

1. Initialize the dictionary,
2. Get the first code, cW
3. Find cW in dictionary and output the string dictionary(cW)
4. While there are still codes to read,
 - a. pW := cW
 - b. Get the next code, cW, and find it in the dictionary
 - c. if found,
 - i. Output the string dictionary(cW)
 - ii. Let C be the first character of dictionary(cW)
 - iii. P := dictionary(pW)
 - iv. Add the string <P, C> to the dictionary
 - d. if not found, **(exception)**
 - i. P := dictionary(pW)
 - ii. Let C be the first character of P
 - iii. Output the string <P, C>
 - iv. Add the string <P, C> to the dictionary

Reminder

- 0 Exception that a code from the lzw file refers to an entry not present yet
- 0 Decompression “lags” behind compression
 - 0 in compression, after producing the first code, the dictionary now contains the 256-th entry
 - 0 in decompression, the 256-th entry is added after reading the second code

What have to do in Ass2 ?

0 Compress and Decompress a list of files to a single archive

0 In command line:

0 For compression :

0 >lzw –c <lzw filename> < list of file names>

0 For decompression:

0 >lzw –d <lzw filename>

What have to do in Ass2 ?

0 Compression

```
LZW.exe -c lzwfile Ephesians.txt Matthew.txt
```



Ephesians.txt

3/9/2013 13:50



LZW.exe

3/8/2013 23:25



lzwfile

3/9/2013 13:53



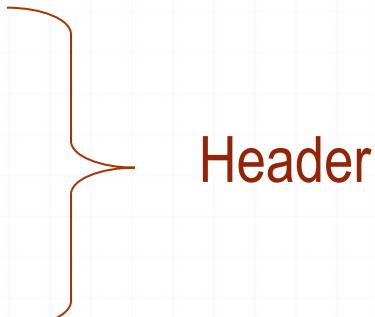
Matthew.txt

3/9/2013 13:49

File format

0LZW file format

```
<filename1>\n
<filename2>\n
<filename3>\n
...
\n
<Compressed file1><end-of-file>
<Compressed file2><end-of-file>
<Compressed file3><end-of-file>
...
```



The diagram shows the 0LZW file format structure. It consists of a header section containing multiple filenames followed by compressed files. The header section is highlighted with a yellow background and includes '<filename1>\n', '<filename2>\n', '<filename3>\n', '...', and '\n'. A red curly brace groups these lines, with the word 'Header' written in red to its right. Below the header, the compressed files are listed with their respective end-of-file markers: '<Compressed file1><end-of-file>', '<Compressed file2><end-of-file>', '<Compressed file3><end-of-file>', and '....'.

What have to do in Ass2 ?

0 Decompression

```
LZW.exe -d lzwfile
```

 Ephesians.txt
 LZW.exe
 lzwfile
 Matthew.txt

3/9/2013 14:03
3/8/2013 23:25
3/9/2013 13:53
3/9/2013 14:03

Details of Implementation

- 0 Full dictionary
 - 0 Code size = 12 in this assignment
 - 0 At most 4096 entries
 - 0 Create a new dictionary in compression
 - 0 Similar steps in decompression

Discussion

0 File end

- 0 Denote $2^N - 1$ as EOF code of one file
- 0 Output this EOF defined by us when finishing the compression of one file
- 0 Stop or open another file when detects the EOF in decompression

Discussion

0 File end

- 0 Please do not make yourselves confused with “ 2^N-1 as EOF code” and EOF returned from reading file.
- 0 “ 2^N-1 as EOF code” is for separating the files from one to one. EOF, means that there is no any character to be read.

Discussion

0 You are provided with sketch file. And six functions should be paid attention to.

0 write_code(lzw_file, code, code_size)

This function is for writing the code to file. It is not necessary to modify this function.

Discussion

- 0 Be careful on using the write_code() routine, since the write_code() routine stores things in a buffer before actual writing
- 0 Remember to flush the buffer at last!
 - 0 `write_code(lzw_file, 0, 12);`

Discussion

0 read_code(lzw_file, code_size)

This function is for reading the code from file. It is not necessary to modify this function as well.

Discussion

0 readfileheader(lzw_file, outputfile, numberoffile)

This function is for reading the information of header from file. It is not necessary to modify this function.

Discussion

0 writefileheader(lzw_file, outputfile, numberoffile)

This function is for writing the information of header to file. It is not necessary to modify this function as well.

Discussion

0 compress(inputfile, lzw_file)

This function is for compression. You are required to implement this function.

Discussion

0 decompress(lzw_file, ouputfile)

This function is for decompression. You are required to implement this function.

Reminder

- 0 You are provided with a sample program
 - 0 Your output should be the same as the ones from sample program.
 - 0 “fc /w /c file1 file2” for file comparison in Windows Command Line.

Reminder

- 0 Try to speed up look-up process
 - 0 Tree structure
 - 0 Hashing function
 - 0 There are bonus for advanced data structure accelerating the speed of compression and decompression.
 - 0 Reference page: <http://en.wikipedia.org/wiki/LZW>