

BLG 372E – Analysis of Algorithms, Spring 2010

Assignment 2 – Huffman Coding

Handed out: 09.04.2010

Due: 22.04.2010

Problem:

Find the Huffman Coding for a given text. See the supplementary file “a2-supp.pdf” at Ninova for details of the algorithm.

A text consists of lowercase letters in the Latin alphabet(a-z), space characters() and periods(.).

Submit your source code and a report through Ninova.

In your report:

- Explain how Huffman Coding algorithm works step by step. (10 pts)
- Analyse the time and space complexity of the algorithm. (10 pts)
- Discuss if the time or space complexity of your algorithm could be lowered by using other data structures. (10 pts)
- Discuss the optimality of the algorithm. (10 pts)
- Compress the sample files one by one with huffman coding and gzip respectively. Compare the compression ratios of these algorithms. Explain the differences of gzip and huffman coding algorithms. (10 pts)

Implementation Details:

- All your code must be written in C++ and able to compile and run on linux/unix using g++.
- You are not allowed to use any library except STL.
- You should write a Makefile with a defined “all” and “test” targets. “all” target will compile all the necessary files in your project and “test” target will run your project with the test input.
- Your submission should only contain the necessary files (sources, headers and a makefile)

Example run:

Input format: Input file contains the text to be coded:

lobortis tation ludus causa iaceo camur ibidem te immitto olim iusto os. adsum quis
duis ullamcorper amet laoreet capio regula loquor indoles quae praesent sit. usitas
gilvus humo ludus vero brevitass facilisis vero torqueo sed velit tristique saepius
iriure. illum esse augue virtus jugis vel eum patria. qui plaga suscipit esse
damnum praemitto huic mara nimis blandit minim.

Output format: Output file “freq.txt” should contain the frequencies of the characters in the original text, “hcodes.txt” should contain the Huffman Coding, “encoded.txt” should contain the encoded text, and “stats.txt” file should include the maximum bits required for coding, length of the encoded text in bits, and the compression ratio. Original text is in ASCII format (each character is represented with 8 bits).

freq.txt:	hcodes.txt	encoded.txt	stats.txt
a 27	a 1001	0001010100001001010011100	8 // Max bits required
b 4	b 000010	0011111010110001001100001	1523 // length of the
c 8	c 110011	1010111000110100011111010	encoded text
d 9	d 01000	0011111110101110011100111	0.50 // compression ratio
e 29	e 1101	1111101001101011100111001	
f 1	f 00001110	1110101011011100111001001	
g 5	g 010011	0111100111010110000100110	
h 2	h 0000110	1000110100101011000110110	
i 40	i 011	1011001000100111000100001	
j 1	j 00001111	0110101010001011001010101	
l 18	l 0001	1111111101000010110101011	
m 18	m 0010	1100100101011001010001110	
n 7	n 110001	1111001010111000011110111	
o 20	o 0101	1101010100011110111110101	
p 8	p 00000	1111000100011001001011001	
q 6	q 110000	1010100110000011010011101	
r 19	r 0011	1001001011011000101000110	
s 31	s 1110	0101010011110111011000101	
t 23	t 1000	1100111001000000110101101	
u 31	u 1111	0011110101001111110001100	
v 7	v 110010	1101000101011100001111010	
space 57	space 101	1001110101111000101000010	
. 5	. 010010	1000111011110101110000...	