

Project 2

CDA 4630/5636: Embedded Systems

Total: 10 points

Due: March 31, 2021, 11:30 PM

In this project, you need to **implement both code compression and decompression using C, C++, Java or Python**. Assume that the dictionary can have sixteen entries (index 4 bits) and the sixteen entries are selected based on frequency (the most frequent instruction should have index 0000). If two entries have the same frequency, priority is given to the one that appears first in the original program order. The original code consists of 32-bit binaries. You are allowed to use only eight possible formats for compression (as outlined below). Note that if one entry (32-bit binary) can be compressed in more than one way, choose the most beneficial one i.e., the one that provides the shortest compressed pattern. If two formats produce exactly the same compression, choose the one that appears earlier in the following listing (e.g., *run-length encoding* appears earlier than *direct matching*). If a 32-bit binary can be compressed using multiple dictionary entries by any specific compression format (e.g., *bitmask-based compression*), please use the dictionary entry with the smallest index value. Please count the starting location of a mismatch from the leftmost (MSB) bit of the pattern – the position of the leftmost bit is 00000.

Format of the *Original Binaries*

000	Original Binary (32 bits)
-----	---------------------------

Format of the *Run Length Encoding (RLE)*

001	Run Length Encoding (3 bits)
-----	------------------------------

Format of *bitmask-based compression* – starting location is counted from left/MSB

010	Starting Location (5 bits)	Bitmask (4 bits)	Dictionary Index (4 bits)
-----	----------------------------	------------------	---------------------------

Please note that a bitmask location should be the first mismatch point from the left. In other words, the leftmost bit of the 4-bit bitmask pattern should be always '1'.

Format of the *1-bit Mismatch* – mismatch location is counted from left/MSB

011	Mismatch Location (5 bits)	Dictionary Index (4 bits)
-----	----------------------------	---------------------------

Format of the *2-bit consecutive mismatches* – starting location is counted from left/MSB

100	Starting Location (5 bits)	Dictionary Index (4 bits)
-----	----------------------------	---------------------------

Format of the *4-bit consecutive mismatches* – starting location is counted from left/MSB

101	Starting Location (5 bits)	Dictionary Index (4 bits)
-----	----------------------------	---------------------------

Format of the *2-bit mismatches anywhere* – Mismatch locations (ML) are counted from left/MSB

110	1 st ML from left (5 bits)	2 nd ML from left (5 bits)	Dictionary Index (4 bits)
-----	---------------------------------------	---------------------------------------	---------------------------

Format of the *Direct Matching*

111	Dictionary Index (4 bits)
-----	---------------------------

Run-Length Encoding (RLE) can be used when there is consecutive repetition of the same instruction. The first instruction of the repeated sequence will be compressed (or kept uncompressed if it is not part of the dictionary) as usual. The remaining ones will be compressed using RLE format shown above. The three bits in the RLE indicates the number of occurrences (000, 001, 010, 011, 100, 101, 110 and 111 imply 1, 2, 3, 4, 5, 6, 7 and 8 occurrences, respectively), excluding the first one. A single application of RLE can encode up to 8 instructions. In other words, up to 9 repetitions can be covered by RLE (first one non-RLE compression followed by up to 8 RLE compression). If you have more than 9 repetitions, you can apply RLE multiple times as long as it is profitable compared to other available options. While multiple combinations are possible, please cover 9 repetitions followed by the remaining ones. For example, if you have 21 repetitions, you should apply RLE three times as <1+8> + <1+8> + <1+2> (instead of <1+6> + <1+6> + <1+6>, or other possible compositions).

You need to show a working prototype that will take any 32-bit binary (0/1 text) file and compress it to produce a output file that shows compressed patterns arranged in a sequential manner (32-bit in each line, last line padded with 0's, if needed), a separation marker "xxxx", followed by sixteen dictionary entries. Your program should also be able to accept a compressed file (in the above format) and decompress to generate the decompressed (original) patterns. Please see the sample files posted in the webpage.

Command Line and Input/Output Formats: The simulator should be executed with the following command line. Use parameters "1" and "2" to indicate compression and decompression, respectively.

./SIM 1 (or **java SIM 1** or **python3 SIM.py 1**) for compression

./SIM 2 (or **java SIM 2** or **python3 SIM.py 2**) for decompression

Please **hardcode** the input and output files as follows:

1. Input file for your compression function: **original.txt**
2. Output produced by your compression function: **cout.txt**
3. Input file for your decompression function: **compressed.txt**
4. Output produced by your decompression function: **dout.txt**

Submission Policy:

Please follow the submission policy outlined below. There will be up to 10% **score penalty** based on the nature of submission policy violations.

1. Please prepare your project submit only in one source file. **Please add ".txt" at the end of your filename.** Your file name must be SIM (e.g., SIM.c.txt or SIM.cpp.txt or SIM.java.txt or SIM.py.txt). On top of the source file, please include the following sentence:

/* On my honor, I have neither given nor received unauthorized aid on this assignment */

2. Please test your submission. These are the exact steps we will follow too.
 - Download your submission from eLearning (ensures your upload was successful).
 - Remove ".txt" extension (e.g., SIM.c.txt should be renamed to SIM.c). We know that eLearning adds a number at the end of the file if you submit multiple times. My script will take care of it (so do not worry about that number inserted by eLearning).

- Login to storm.cise.ufl.edu. If you are not a CISE student, please find a linux machine in your department and use the following commands. “It is working in my laptop” is not an acceptable argument. If you run your code in a linux machine, I should be able to test it in storm.cise.ufl.edu.
 - Please compile to produce an executable named **SIM**.
 - gcc SIM.c -o SIM **or** javac SIM.java **or** g++ SIM.cpp -o SIM **or**
g++ -std=c++0x SIM.cpp -o SIM
 - Please do not print anything on screen.
 - Assume hardcoded input/output files as outlined in the project description.
 - Compress the input file (original.txt) and check with the expected output (compressed.txt)
 - ./SIM 1 (or java SIM 1 or python3 SIM.py 1)
 - diff -w -B cout.txt compressed.txt
 - Decompress the input file (compressed.txt) and check with the expected output (original.txt)
 - ./SIM 2 (or java SIM 2 or python3 SIM.py 2)
 - diff -w -B dout.txt original.txt
3. *In previous years, there were many cases where output format was different, filename was different, command line arguments were different, or e-Learning submission was missing. All of these led to un-necessary frustration and waste of time for the instructor and students. **Please use the exact same commands as outlined above to avoid 10% score penalty.***
4. **You are not allowed to take or give any help in completing this project.** *In previous years, some students violated academic honesty (giving help or taking help in completing this project). We were able to establish cheating in several cases - those students received “0” in the project and their names were reported to Dean of Students Office (DSO). If your name is already in DSO for violation in another course, the penalty for the second offence is determined by DSO. In the past, two students from my class were suspended for a semester due to repeated academic honesty violation (implies deportation for international students).*

Grading Policy

The project assignment has the sample input and output files. Correct handling of the sample input will be used to determine 60% of credit awarded. The remaining 40% will be determined from other input test cases that you will not have access prior to grading. The other test cases can have different types and number of 32-bit binaries (0/1 text). It is recommended that you construct your own sample input files with which to further test your compression and decompression functions. You can assume that we will use less than 1024 32-bit binary (0/1 text) patterns in the new test file. **Please note that the new test case will NOT test any exceptional scenarios that are not described in this document.**