Total marks: 10                                               Time: 12 minutes

Name: Rahul Vigneswaran K                                    (X)

Roll number: CS23MTECHO2002

-------------------------------------------------------------------------------------------------------

1.  In CNN accelerators, the area occupied by compute units is significantly larger than that of local storage (buffers or SRAMs): [1]
    a.  True
    b.  False

2.  Huffman encoding is a lossless compression scheme: [1]
    a.  True
    b.  False

3.  Briefly explain how/why loop tiling improves the execution time. [2]
    → Loop tiling would reuse that specific tile at hand.
    → This way we process those tiles individually instead of processing as whole.

    *But why is processing tile-wise more efficient?*

    0

4.  Briefly explain the importance of using a dual-port SRAM in accelerators. [2]

    → When there are 2 port, we can use the 1 port to write things onto the SRAM while reading out through another.

    → So we do have to wait for the reading to be done in order to start writing. Both can happen at same time, saving wasted waiting time.

    2

5. Consider the 4x4 matrix as shown below and assume that we have 4 bins 0, 1, 2, 3 with centroids as -1.0, 0.0, 1.0, and 2.0, respectively. Encode the matrix using these bins in a similar manner as discussed in the Deep Compression paper. [2]

| 2.10 | -1.10 | 2.09 | 1.05 |
|------|-------|------|------|
| 0.07 | -1.02 | 1.31 | -0.01 |
| 2.03 | -0.06 | 1.02 | 1.99 |
| -1.01 | 2.01 | 1.12 | -0.98 |

3 — 0 — 3

```
3 0 3 2
1 ③ 2 1
3 1 2 3
0 3 2 0
```

```
-1.0
0.0
1.0
2.0
```

1.5

Assuming that each value in the original matrix uses 32-bits, calculate the storage requirements (in bits) for the original and the encoded representation. [2]
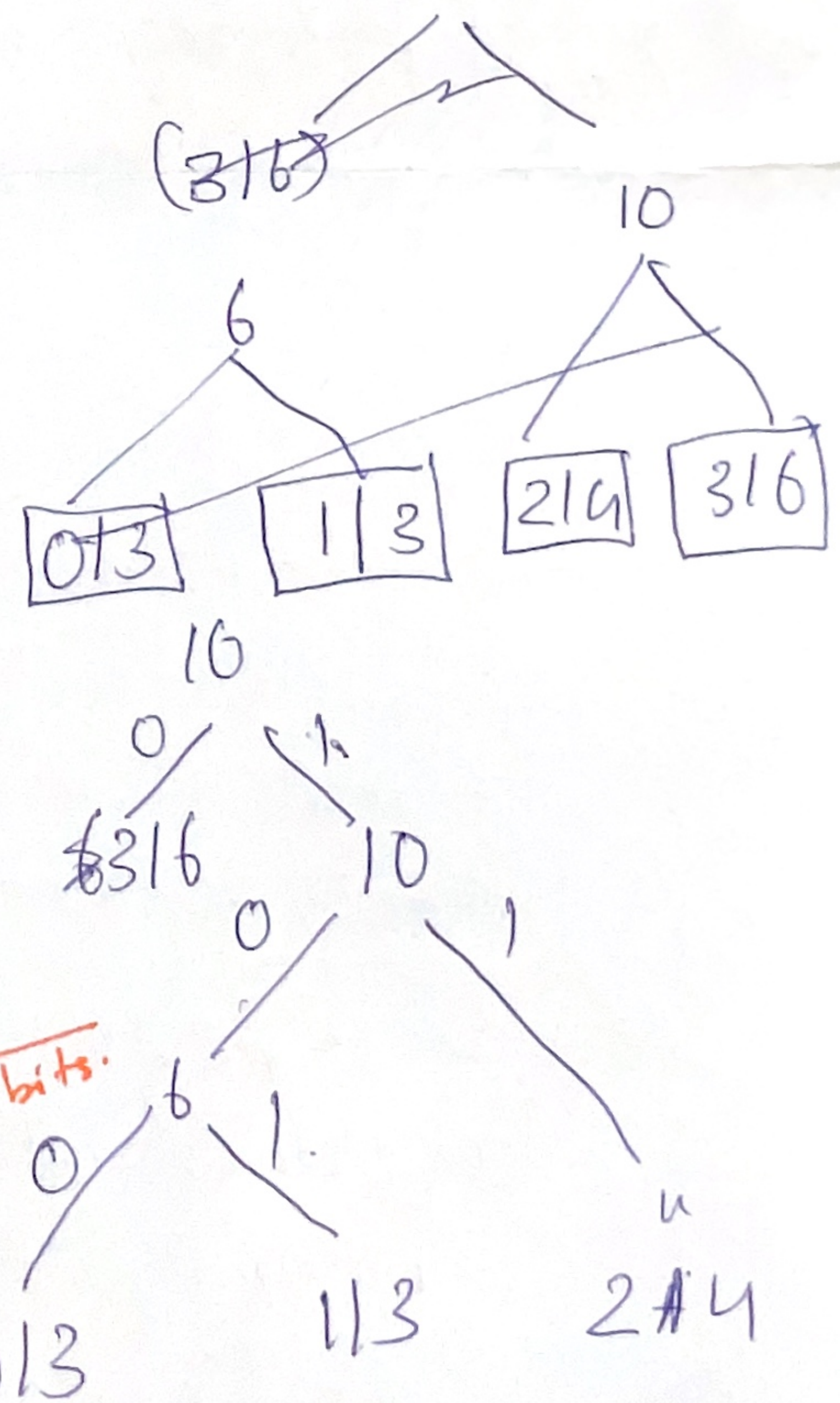
Original : 16 × 32 - bit ✓

Encoded : 0 — 3 times
          1 — 3 times
          2 — 4 times
          3 — 6 times

1.5

0 — 10        0 - 3 bit ×3
1 — 10        1 - 3 bit ×3
2 — 11        - 2 bit ×4
3 — 0         - 1 bit ×6
              _____
              32 bits.

Indent

4 centroid × 32 bit + 8 bits
128 + 3 = 131 bits >> original

(316)
10
6
| 0|3 |   | 1|3 |   | 2|4 |   | 3|6 |

10
0 / 1
$316   10
0 / 1
6 / 1.
0 / 1.
013        1|3        2 A 4