# ESE532 Project Report

## Deduplication and Compression

**Akhil Gunda, Ruturaj Nanoti, Siddhant Mathur**
**December 2023**

# Contents

# 1 Single ARM Processor Mapped Design

a) The Deduplication and Compression application on a single ARM processor mapped design is described step wise as shown below:

1. The data, ethernet input, is received in real time using ethernet cable from the client machine and stored in local storage, a buffer, using APIs provided in the project.

2. The ethernet input data is then moved into a local input buffer of size 16 KB (Number of Packets * Size of each packet = 16 * 1KB).

3. When the 16 KB buffer is filled, the compression pipeline function is invoked. This compression pipeline includes sub functions (CDC, SHA, DEDUP, LZW and final bit packing) required for compression application.

4. The operations inside compression pipeline take places sequentially which includes:

   - The input data of 16 KB is passed through the content defined chunking function. The CDC function starts calculating hash value with a window size of 16 (WIN_SIZE) on the incoming packets. If the hash value matches with the target hash value of 1024, (MODULUS), a chunk is generated. The end index of the chunk is stored and starts for the next chunk. Finally all the chunk boundaries are stored in a vector and moved to the next function.

   - The SHA256 function starts calculating the fingerprints for the chunks generated. Next the DEDUP function uses an `unordered_map`. It checks if the SHA fingerprint is already present in the unordered map and if not, it assigns each unique chunk a chunk ID and calls LZW on it. If found, just returns the index of the map where it was previously added.

   - If a chunk was not found, the LZW function is called which accepts 16 KB input data along with the chunk boundaries. The LZW sends out the compressed data. Then the compressed data is bit packed and stored in a file.

   - Iterate until all the chunks are traversed.

5. Steps 3 and 4 are repeated for all the packets received via Ethernet. At the end the final compressed binary file is generated.

b) Following are the parameters that were used in our initial implementation.

- **CDC:** The table below shows the default parameters used for CDC.

| WINDOW SIZE | 16 |
|---|---|
| MODULUS | 1024 |
| PRIME | 3 |
| TARGET | 0 |
| MAX CHUNK SIZE | 8192 |

Table 1: CDC Parameters

- **SHA:** The following table shows the default parameters used for SHA/D-EDUP.

| SHA256 BLOCK SIZE | 32 |
|---|---|

Table 2: SHA Parameters

This is set to 32 since the SHA256 digest is 32 byte/256 bit long.

- **DEDUP:** For Insertion and lookup in Deduplication, we used an unordered_map for mapping SHA fingerprints to chunk IDs. The benefit of using unordered_map was that they gave us O(1) lookup time.
Datatype - The SHA fingerprints were generated as a byte array. When this was used as a key to the ordered_map, we were getting a lot of collisions in the Deduplication stage. To counter this, we converted the byte array to a string since an unordered_map needs a key that can be represented as a single data type to hash into the table.

- **LZW:** Initially we used a hash table of size 65536 and 1 bucket. We used a random hash function initially that consisted of some random bit shifts and manipulations. However, we decided to change the hash function in an attempt to decrease the size of the hash table.

  Hash function - We eventually moved to the murmur hash function using which we were able to reduce the size of the hash table to 8192. However, for bigger chunks we were still getting collisions, hence, we increased the buckets to 2 that were used as a 2D array.

- **Overall Application:** Number of Packets, We were initially receiving 16 packets since we were using a block size of 1024.

c) Key Performance Achieved on Single ARM Processor Design: The throughput of our application is calculated based on the compression latency. The throughput calculation follows the following formula:

$$ApplicationThroughput = \frac{\frac{BytesReadFromEthernet \times 8}{1000000}}{TotalCompressionLatency}$$

In determining the compression latency, we utilized the stopwatch class by encapsulating the start and stop calls around our top function. This top function encompasses all the sub stages, such as CDC, SHA, DEDUP, LZW and includes writing data to the compressed file. Consequently, the calculation throughput is 6.815 Mb/s, equivalent to 0.006815 Gb/s.

d) The following table shows the command line parameters:

| Block Size (-b) | 1024 |
|---|---|
| Sleep Parameter (-s) | 5 |

Table 3: Command Line Parameters

e) Compression Achieved: For the test case - LittlePrince.txt file, we achieved a compression ratio of,

$$\frac{CompressedFileSize}{OriginalFileSize} = \frac{11825}{14247} = 0.83$$

e) Breakdown of Time spent on each component:

| Task | Time (Days) |
|---|---|
| CDC Implementation | 2 |
| SHA/DEDUP Implementation | 2 |
| LZW (on ARM Processor) Implementation | 5 |
| Main Application and Integration | 3 |
| Debugging | 4 |
| Testing and Validation | 3 |

Table 4: Breakdown of Time Spent

# 2 Final Ultra96 Mapped Design

a) Performance Achieved on Ultra96 Mapped Design: The throughput of our application is calculated based on the compression latency. The throughput calculation follows the following formula:

$$ApplicationThroughput = \frac{\frac{BytesReadFromEthernet \times 8}{1000000}}{TotalCompressionLatency}$$

In determining the compression latency, we utilized the stopwatch class by encapsulating the start and stop calls around our top function. This top function encompasses all the sub stages, such as CDC, SHA, DEDUP, LZW and includes writing data to the compressed file.

**Application Throughput: Linux.tar (Size: 200273920 B)(Binary File)**
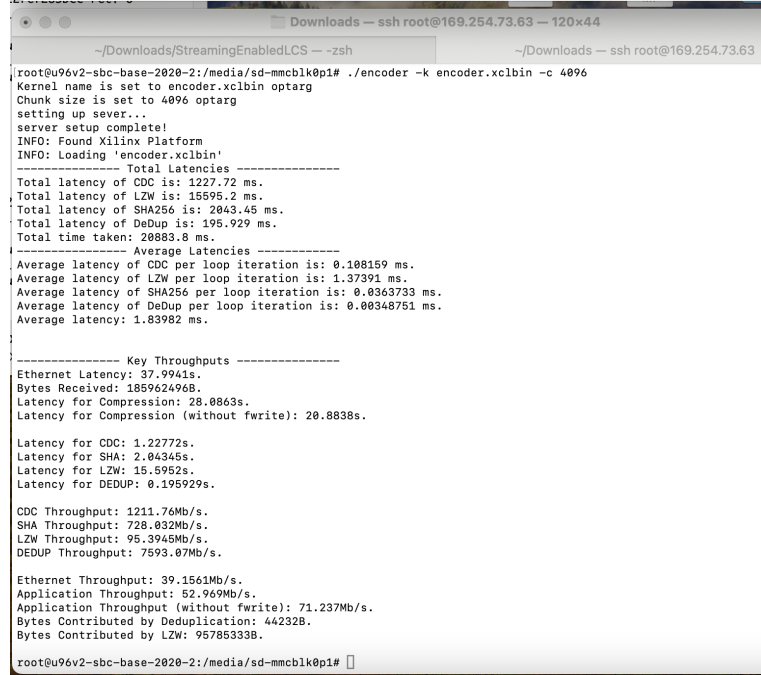The throughput achieved is **71.237 Mb/s**. The intermediate throughputs were 1211.76 Mb/s for CDC, 728.032 Mb/s for SHA, 7593.07 Mb/s for Dedup, and 95.3945 Mb/s for LZW. Our output file is compressed to 95829565 Bytes. These values were determined from the "Linux.tar" test case.

**Application Throughput: Franklin.txt (Size: 399054 B)(Text File)**
Throughput achieved is **67.443 Mb/s**. The intermediate throughputs were 1178.01 Mb/s for CDC, 713.09 Mb/s for SHA, 8704.41 Mb/s for Dedup, and 92.6546 Mb/s for LZW. Our output file is compressed to 291731 Bytes. These values were determined from the "Franklin.txt" test case.

Terminal Output Linux Tar and Franklin respectively:

```
[akhilgunda@Akhils-MacBook-Air Primary Files % ../client_mac -i 169.254.73.63 -b 8192 -f linux.tar -s 2000
ip is set to 169.254.73.63
blocksize is 8192
filename is linux.tar
sleep is set to 2000 optarg
bytes_read 200273920
akhilgunda@Akhils-MacBook-Air Primary Files %
```

```
Downloads — ssh root@169.254.73.63 — 120×44
   ~/Downloads/StreamingEnabledLCS — -zsh          ~/Downloads — ssh root@169.254.73.63
[root@u96v2-sbc-base-2020-2:/media/sd-mmcblk0p1# ./encoder -k encoder.xclbin -c 4096
Kernel name is set to encoder.xclbin optarg
Chunk size is set to 4096 optarg
setting up sever...
server setup complete!
INFO: Found Xilinx Platform
INFO: Loading 'encoder.xclbin'
--------------- Total Latencies ---------------
Total latency of CDC is: 1227.72 ms.
Total latency of LZW is: 15595.2 ms.
Total latency of SHA256 is: 2043.45 ms.
Total latency of DeDup is: 195.929 ms.
Total time taken: 20883.8 ms.
---------------- Average Latencies ------------
Average latency of CDC per loop iteration is: 0.108159 ms.
Average latency of LZW per loop iteration is: 1.37391 ms.
Average latency of SHA256 per loop iteration is: 0.0363733 ms.
Average latency of DeDup per loop iteration is: 0.00348751 ms.
Average latency: 1.83982 ms.


--------------- Key Throughputs ---------------
Ethernet Latency: 37.9941s.
Bytes Received: 185962496B.
Latency for Compression: 28.0863s.
Latency for Compression (without fwrite): 20.8838s.

Latency for CDC: 1.22772s.
Latency for SHA: 2.04345s.
Latency for LZW: 15.5952s.
Latency for DEDUP: 0.195929s.

CDC Throughput: 1211.76Mb/s.
SHA Throughput: 728.032Mb/s.
LZW Throughput: 95.3945Mb/s.
DEDUP Throughput: 7593.07Mb/s.

Ethernet Throughput: 39.1561Mb/s.
Application Throughput: 52.969Mb/s.
Application Throughput (without fwrite): 71.237Mb/s.
Bytes Contributed by Deduplication: 44232B.
Bytes Contributed by LZW: 95785333B.

root@u96v2-sbc-base-2020-2:/media/sd-mmcblk0p1#
```
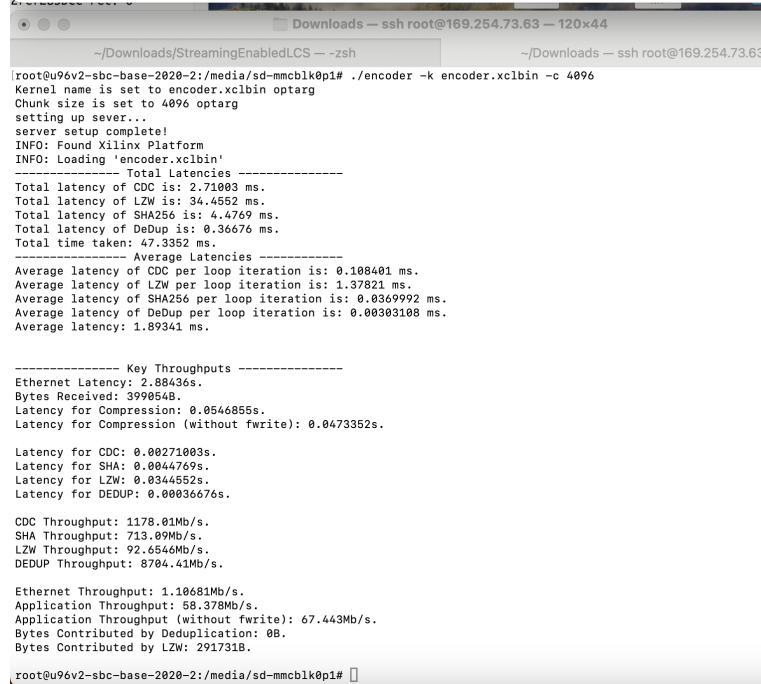
Figure 1: Linux Tar

```
[akhilgunda@Akhils-MacBook-Air Primary Files % ../client_mac -i 169.254.73.63 -b 8192 -f Franklin.txt
ip is set to 169.254.73.63
blocksize is 8192
filename is Franklin.txt
bytes_read 399054
akhilgunda@Akhils-MacBook-Air Primary Files %
```

```
Downloads — ssh root@169.254.73.63 — 120×44
   ~/Downloads/StreamingEnabledLCS — -zsh          ~/Downloads — ssh root@169.254.73.63
[root@u96v2-sbc-base-2020-2:/media/sd-mmcblk0p1# ./encoder -k encoder.xclbin -c 4096
Kernel name is set to encoder.xclbin optarg
Chunk size is set to 4096 optarg
setting up sever...
server setup complete!
INFO: Found Xilinx Platform
INFO: Loading 'encoder.xclbin'
--------------- Total Latencies ---------------
Total latency of CDC is: 2.71003 ms.
Total latency of LZW is: 34.4552 ms.
Total latency of SHA256 is: 4.4769 ms.
Total latency of DeDup is: 0.36676 ms.
Total time taken: 47.3352 ms.
---------------- Average Latencies ------------
Average latency of CDC per loop iteration is: 0.108401 ms.
Average latency of LZW per loop iteration is: 1.37821 ms.
Average latency of SHA256 per loop iteration is: 0.0369992 ms.
Average latency of DeDup per loop iteration is: 0.00303108 ms.
Average latency: 1.89341 ms.


--------------- Key Throughputs ---------------
Ethernet Latency: 2.88436s.
Bytes Received: 399054s.
Latency for Compression: 0.0546855s.
Latency for Compression (without fwrite): 0.0473352s.

Latency for CDC: 0.00271003s.
Latency for SHA: 0.0044769s.
Latency for LZW: 0.0344552s.
Latency for DEDUP: 0.00036676s.

CDC Throughput: 1178.01Mb/s.
SHA Throughput: 713.09Mb/s.
LZW Throughput: 92.6546Mb/s.
DEDUP Throughput: 8704.41Mb/s.

Ethernet Throughput: 1.10681Mb/s.
Application Throughput: 58.378Mb/s.
Application Throughput (without fwrite): 67.443Mb/s.
Bytes Contributed by Deduplication: 0B.
Bytes Contributed by LZW: 291731B.

root@u96v2-sbc-base-2020-2:/media/sd-mmcblk0p1#
```

Figure 2: Franklin Text File

Decode Diff Outputs for Linux.tar and Franklin.txt:



Figure 3: Linux Tar Decoder Diff



Figure 4: Franklin Decoder Diff

**Vivado Analysis:**

- Power Consumption for this Design: The estimated total on-chip power is **2.803 W**.



Figure 5: Power Consumption

- Resource Utilization:



Figure 6: Power Consumption

b) Task Decomposition: For our final implementation, we receive data in buffers of 16KB from the client. Since we have moved to a block size of 8192, we are receiving 16KB in 2 packets.

Once we receive all of the 16KB data from the client, this is sent to the CDC function which makes multiple chunks of 4KB. We are using a different version of CDC called the Fast CDC which ensures that none of the chunks are greater than 4KB. This implementation of Fast CDC uses the GEAR hash. Chunks could be smaller than 4KB but we perform a hard chunk at 4KB. This version of Fast CDC performs a modulus in a different way as opposed to the regular CDC functionality. This works by skipping sub-minimum chunk cut-points. In this implementation, if the current chunking position is less that the average chunk size the modulus operation with the hash is performed with a value higher than the average chunk size and if the current chunking position is greater than the average chunk size the modulus operation with the hash is performed with a value smaller than the average chunk size. This ensures an even distribution of chunk sizes. This is how, a chunk is declared. If not, it declares a hard chunk of 4KB by looking at the previous chunk boundary.

Once all these chunks are made, SHA generates a 256 bit digest for each of these chunks. We have mapped the SHA computation onto the NEON registers which provided a substantial speedup over the software SHA implementation. This was implemented using the MPSOC library that was given in the project handout. This

implementation uses the SHA cryptographic intrinsics as well as the NEON intrinsics.

These chunk fingerprints are used as keys in the hash table in the Deduplication stage. Fingerprints associated with duplicate chunks are only added to the hash table once. As of our current implementation, we are performing LZW on each of these chunks irrespective of it being a duplicate or not. This is because branch statements in our kernel were proving to slow down the overall application. Hence, we chose to perform LZW compression on every chunk regardless of it being a duplicate or not. Then on the basis of the output of the Deduplication stage, we decide whether to write just the header or the compressed codes in the form of packets to the file.



Figure 7: Flowchart

8

c) Compression Achieved:

For the test case - Franklin.txt file, we achieved a compression ratio of,

$$\frac{CompressedFileSize}{OriginalFileSize} = \frac{291731}{399054} = 0.73$$

For the test case - Linux.tar file, we achieved a compression ratio of,

$$\frac{CompressedFileSize}{OriginalFileSize} = \frac{95829565}{200273920} = 0.48$$

d) Parallelism:

1. For our CDC implementation, we tried to precalculate and store the power values in a local buffer. Since the hash calculation was using powers of 3 up to 17 repetitively for each iteration, we decided to hand calculate these and store them in an array from which they can be read for hash calculations. This made our CDC significantly quicker. To further optimize our CDC, we changed to FastCDC which used the GEAR hash and skips sub-minimum chunk cutpoints. This made our final CDC implementation much faster as opposed to regular CDC.
CDC throughput before changing to FastCDC: 377.822 Mb/s
CDC throughput after adapting the FastCDC implementation : 1185.7 Mb/s
Speedup: 3.138x

2. For the SHA implementation, we changed our software implementation of SHA 256 with the NEON intrinsics version of the implementation. This was implemented using the MPSOC library that was given in the project handout. This implementation uses the SHA cryptographic intrinsics as well as the NEON intrinsics which gives us significant speedup for the SHA since it uses 8 vector lanes.
SHA latency with the software implementation: 0.00058189 s
SHA latency with the NEON intrinsics implementation: 0.00028 s
Speedup: 2.078x

3. Pipelining and parallelism can also be exploited within the kernel by enabling streams and data flow between the load, compute and store functions of the LZW compression. These streams act as a FIFO buffer from which data can be written to or read from each iteration. Functions within the dataflow region can execute concurrently.
We tried to use the dataflow pragma to enable data flow between the load, compute and store function to speed up the kernel. Although, we weren't able to successfully complete it due to issues with synchronization in the output stream generated by the compute function. We noticed that with the dataflow pragma the C simulation worked perfectly fine, and the synthesis also passed, but the post C checking in co-simulation failed due to the output array being filled to zeros.

4. To exploit parallelism using threads in the host code, we tried to map the data collection from the ethernet on one thread and the compression_pipeline

function to another thread. This way we could use double buffering to fill data up in a buffer while another buffer is being processed on. Due to this our data processing pipeline is never waiting for data. Although this too didn't work out, and we were facing issues with the two buffers that were setup, which resulted in a segmentation fault at the end of the application.

e) Zynq Resource Mapping:

| Application Stage | Hardware | Memory Utilization | FPGA Resource Utilization |
|---|---|---|---|
| Encoder | ARM Processor | 145KB | 0 |
| CDC | ARM Processor | 17KB | 0 |
| SHA | NEON Vector Processor + ARM Processor | 10KB | 0 |
| DEDUP | ARM Processor | 500 KB (Worst Case) | 0 |
| LZW | FPGA | Any memory used in the host is included in the Encoder | LUTs: 48550 FFs: 11269 DSPs: 3 BRAM_18Ks: 187 |

Table 5: Zynq Resource Mapping



Figure 8: Vitis Resource Utilization

The following screenshot shows the breakdown of the resource utilization on the FPGA for our LZW kernel.

Figure 9: Vitis Resource Profile

f) Performance Model:
Our final implementation consists of only 1 kernel/ compute unit. Hence, that will be our LZW throughput. The throughput is going to be defined by the bottleneck or the slowest operation. Since, LZW is our bottleneck, it is going to define our overall performance/ throughput to a large extent. Since everything is currently running on 1 thread, our overall application throughput is going to be min (CDC, SHA, Deduplication , LZW Compression).
The overall application throughput does not match LZW's throughput because our current implementation does not overlap the computation of LZW with that of the other functions. As a result of which, the overall application throughput is a little lesser than the throughput of LZW.

| Stage | Latency (ms) |
|---|---|
| CDC | 2.71 |
| SHA | 4.47 |
| DEDUP | 0.36 |
| LZW | 34.45 |

Table 6: Latencies

Total Latency: 41.99 ms

Total Application Latency(including migrating memory back and forth from the host): 47.33 ms

$$ModeledThroughput = \frac{SizeOfFranklin}{TotalLatency} = \frac{\frac{399054\times8}{1000000}}{47.33} = 67.45Mb/s$$

The reason this does not match our actual throughput is because of memory transfer overheads.

g) Current Bottleneck Preventing Higher Performance:
The bottleneck in our final implementation continues to be the LZW compression. To improve its performance, we moved the LZW onto the FPGA and tried a variety of methods to improve it. The screenshot below shows the Application timeline of our kernel in Vitis Analyzer.

11

Figure 10: Timeline Trace

# 3   Validation and Real-Time Input

a) Design Validation:

At the beginning we focused on developing functionality for each of the different stages in the compression pipeline which included CDC, SHA, DEDUP and LZW. To enable seamless integration, during this stage we identified the way in which interfaces needed to be designed for each of the pieces to come together. This meant coming up with a high-level design of the complete control flow of our application, i.e., what data is passed to what stage and how it will process it before proceeding to the next.

Following is a detailed description of the design and validation that we performed for each of the aforementioned stages.

- **CDC:** Initially we started off by using our validated implementation from HW2 that used the rolling hash. This helped us make sure that we were working with a functional CDC and that it would not break the application when integrating it with the other components. As a part of our validation process we made sure that our CDC produced the same output as the python code in HW2 for varying modulus values as a part of further validation process. This was done by printing out chunk indices and redirecting them into a file, and the files created from our implementation and the python code were compared using the diff tool, to validate the functionality.

  We did this at every stage whenever the CDC was optimized as a sanity check and also to preemptively catch any bugs that would arise in the future. Apart from using text files we also used binary data to validate the functionality.

- **SHA:** From the beginning we focused on getting a fully functional SHA implementation since it was a critical part of the whole application. The output from SHA would decide if our chunk was a duplicate and whether it needed to be compressed, and if our SHA produced the same digest for different chunks it would severely harm the overall functionality.

  Our implementation for SHA-256 was based on the mpsoc-crypto library, which contained both the software and NEON versions of SHA-256. Initially aiming for functionality we went with implementing the software version. This library

12

had its own test test files that we first used to make sure that the SHA implementation was valid and legitimate. Once that was done we integrated that library into our implementation. Once that was done we tested it with our interface to validate that the functionality was preserved. The testing was done in two parts. Firstly we used the linux command line utility called sha256sum as our gold standard, and compared its output against it. Following that to further validate our design we used the NIST examples to test out the SHA implementation, provided here. We tested against the inputs and compared the generated digest with the digest provided in the reference.

We made sure that our design was validated both on text and binary data to ensure full functionality.

- **DEDUP:** Our implementation for DEDUP used unordered_maps to lookup SHA fingerprints and check if an entry for it already existed, if it did we would return the value, i.e., the chunk index mapped to that fingerprint or return -1 after inserting the digest into the table with a new chunk index.

  To validate this we made sure that previously seen SHA fingerp rints were producing the same chunk index, upon lookup into the map. This was done by testing the output from the same SHA fingerprint and validating that both the calls to dedup returned the same value. Apart from this we also tested our design by providing unique SHA fingerprints and checking the return value from the function (which should be -1, after the new mapping from SHA to chunk id is inserted into the map). Finally we tested if our design would produce the right chunk index after providing a new SHA fingerprint and then calling dedup again to check if the new fingerprint was inserted into the map with the right chunk index.

- **LZW:** For this part we started off with the LZW implementation from geeks for geeks and modified it so that it did not use any unsupported C++ constructs like strings and vectors. We made sure that all our parts could be easily ported to the FPGA if they were to turn out to be the bottleneck.

  Our interface for LZW required modifying the function to accept start and end indices into a big buffer, which represented the "chunk". Hence, we needed to make sure that the implementation worked for the chunk being somewhere in between the buffer, at the beginning of the buffer and at the end of the buffer. This was done by randomly generating start and end indices and passing them as arguments to the LZW function. This also helped us validate our loop bounds and ensured we were processing all the data that was passed into the function.

  Initially our validation process for this stage consisted of comparing the output from our implementation to the one produced by the geeks for geeks version. We started out by testing small strings like "I am Sam, Sam I am, am I Sam" and made sure that our output produced the same results as compared with the original. We then moved onto LittlePrince and finally moved onto small binary files. Once we were sure that our implementation was functioning correctly,

we tested with larger text files like Franklin and the Lord of the Rings, and after this passed. After this we moved to testing large files like a zip archive for a font named FiraCode (20MB) and tar file for a graphical editor named neovide (24 MB).

Now for the application as a whole once we had integrated everything we started out with a simple packet emulation program, that would read blocksize amount of bytes from a file at a time and copy them into a 16KB buffer in a loop. Once this buffer was filled up we would call our compression_pipeline function to process that buffer. This helped us quickly weed out any problems that occurred during the integration phase since we were able to quickly test out our design on the PC itself instead of migrating all the binary files to the board and testing it there. It also aided in simplifying the original implementation so that we could focus on squashing any bugs in the standalone implementations of CDC, SHA, DEDUP and LZW.

After the packet emulation was fully functional we moved onto using the data packets transmitted via ethernet and running the whole application on the board (everything was mapped to a single ARM processor at this point). This helped us validate our design with all the intricacies involved in packet transfer. Apart from that this phase aided in identifying the bottleneck since we used stopwatch class to time and profile the different stages of our application. This data backed our decision of mapping LZW onto the FPGA.

Finally, when LZW was to be mapped to the FPGA we needed to focus on verifying its hardware implementation. This is when we wrote our testbench in Vitis that would compare the output from our kernel with the golden implementation from geeks for geeks. Initially, the testbench would read chunk_size data into a buffer and pass that to the LZW kernel for processing.

This worked for our initial implementations but when we optimized our host code to just call LZW once and process a batch of chunks we realized that to get a good idea of the functionality we had to expand the scope of testing, and emulate the conditions that the kernel is actually going to be running in, i.e., receiving start and end indices to chunks in a large buffer. So, we created an array that would hold randomly generated chunk indices and tested the kernel with them. Moving forward we decided that to emulate and generate realistic chunk indices we needed to use our CDC implementation and hence we included a function call to CDC in our testbench.

Finally, based on some bugs that we faced (which are outlined in the section below) we decided to enhance our testing routine. We would now read in the whole file 16KB at a time in a loop, pass that buffer through CDC, generate chunk indices, and then pass that data to our kernel for testing. Following which the output from the kernel would be compared against the golden implementation at the granularity of a chunk. The packet length and the actual data was compared against the golden implementation to make sure that our kernel was not producing any incorrect values. This helped us understand the level of current functionality and the points at which our implementation would break, and in this way we could fix our design before deploying it onto the FPGA.

b) Real-Time Guarantee:
Based on our extensive testing and validation we can say with confidence that our design is completely functional and stable with a throughput of approximately 70 Mb/s as per the "linux.tar" test for files up to 200 MB with an -s parameter of 2000 at a block size of 8192B. The -s parameter is only required in case of large files like the "linux.tar" and the "FiraCode" archive (which required an -s parameter of 1000). Whereas for smaller files like "Franklin.txt" and "LittlePrince.txt" we noticed that the -s parameter was not required. We observed this trend of a smaller or no -s parameter as the size of the file reduced. As per our reasoning this may be due to the fact that our CDC creates chunks based on the input data, and the number of chunks may vary depending on the content. As the number of chunks increase the processing time for each 16KB data buffer that we maintain increases, and this may change dynamically based on the type of the input, and since the probability of having this worst case delay goes up with a higher number of packets (for large files), the -s parameter also needs to go up.

Commands to run our design:

- **On a Mac (Host):** Run the following command on your board (by default the file name for the compressed file is *compressed_file.bin*,

```
$ ./encoder -k encoder.xclbin -f <Filename>
```

Run the following command in your host computer,

```
$ ./client_mac -i <board_ip_address> -b 8192 -f <Filename> -s 1000
```

Now to decode the file run,

```
$ ./decoder <compressedFilename> <outputFilename>
```

- **On a Windows (Host):** Run the following command on your board (by default the file name for the compressed file is *compressed_file.bin*,

```
$ ./encoder -k encoder.xclbin -f <Filename>
```

Run the following command in your host computer,

```
$ ./client -i <board_ip_address> -b 1024 -f <Filename> -s 1000
```

Now to decode the file run,

```
$ ./decoder <compressedFilename> <outputFilename>
```

c) Challenges and Debugging:
During the development of this project, we faced a lot of challenges and bugs which were uncovered during the thorough testing that we carried out for each of the different stages. Following is a detailed description of all the bugs that we fixed and

the challenges that we faced broken down based on the different stages of the data pipeline.

- **CDC:** When we were integrating our CDC into the whole application, we noticed that CDC was creating chunk boundaries at every character. After carefully analyzing the code and stepping through the code in GDB we noticed that the offset that we used for filling up the 16KB buffer was invalid and this resulted in CDC operating on garbage data. Once this was identified we fixed our offset calculation to use the length from the header of the ethernet packet and this rectified our issue.

- **SHA/DEDUP:** In this stage we faced more of a challenge than a bug, which was related to storing the SHA fingerprints as keys into an unordered_map that was used to map SHA digests to chunk indices. Now, the unordered_map in C++ can't efficiently map keys which are represented as containers like vectors, arrays, etc. This meant that we couldn't store our SHA digest as a byte array but had to convert it into a form that could be represented by a single data type, which was a string. This conversion of the byte array to a string proved to be really challenging since we had to take special care while converting "0"s presents in the SHA fingerprint since NULL (0x00) characters are considered to be string terminators in C. To mitigate this problem we used the stringstream class in C++. Our solution included creating a hexadecimal stringstream and iterating over the byte array. At each iteration we would cast the current byte into an int using the static_cast directive, and append the int value as a two character wide hexadecimal string to the output string. Finally after every byte was processed we would convert the stringstream object into a string and return that from SHA so that DEDUP could then utilize an unordered_map.

  Apart from that in the SHA library we noticed that the implementation used strlen to find the length of the string and compute the SHA based on that, but we knew that this wouldn't work on binary data since as described earlier NULL (0x00) characters are considered to be string terminators in C, and this would result in the string getting cut out prematurely which in-turn would give out an incorrect SHA. To pre-emptively catch this we adapted the library to use start and end indices into a buffer representing the chunk to ensure that every chunk is processing completely.

- **LZW:** The first bug that we faced in LZW was related to the last character not showing up properly in the decoded output for a chunk. This happened in all the files that we tested, and it pointed us to our incorrect handling of the loop bounds for the last character. When we analyzed the code using gdb we found out that the loop terminated prematurely because our loop bound was set by subtracting 1 from the end index whereas the input passed to LZW already took that into account and this led to an off by one error. Once it was identified we fixed it by removing the subtraction by 1.

  After we had moved to the FPGA mapping of the LZW we noticed that the implementation was failing when working with binary data. This was due to

16

the fact that the first 256 entries in the hash table were initialized with the key being the ASCII values of all the characters followed by a zero to the respective ASCII values of the character itself, i.e., 'a0' was mapped to 97 since the ASCII value of 'a' is 97. This meant that a sequence of characters represented by '00' was already mapped to a value of 0, although actually we haven't seen that mapping before and hence the mapping should really be of an invalid character followed by a valid character to ensure that we aren't mapping any sequences accidentally. To ensure this we can just remove the initialization of the hash table and keep the next code set to 256. This way we ensure that there are no invalid mapping present in the hash table which will prevent from getting any false positive lookups.

Another problem that we faced while processing binary data with LZW was that our implementation would fail at random places within the binary file and work perfectly fine most of the time. It took us a lot of time trying to understand where the problem was occurring, and since we were using Vitis there wasn't a good debugging tool available. So, firstly we included some debugging print statements in the code to find the point of failure (this included the iteration index, offset, etc.). After that we decided to compile the testbench and our kernel using g++ outside vitis and used two gdb instances, to stop at the exact iteration and offset in both our kernel and the golden implementation parallelly. There we first examined the prefix and the next character to ensure that the problem didn't occur in a previous iteration. We found that both of them were the same, then we advanced the code in the golden implementation to check if the sequence was present in the map and was found or was it inserted into the map this iteration, and surprisingly our implementation was also following the same pattern. This suggested that there was a problem with the value returned by the lookup. We checked the current code value and found that the value was exceeding our set maximum of 4096, and due to this our hash lookup returned a value that wasn't associated with that particular key, since our hash table was configured to have a key consisting of a 12-bit value and an 8-bit character. Due to this our complete implementation broke down and we were getting random points of failure.

Once we found this out, we decided that we will need to move to 13-bit hash table configuration and 13-bit packed output since our lzw codes were exceeding the value represented by 12-bits. Hence, we modified the code accordingly. After that we encountered a similar problem. We were still getting random failures. After setting up the parallel gdb instances again we found that the problem was with overflow, the insert function in the associative memory didn't set the proper bit in memory when the fill value exceeded 32, since $1 << 32$ exceeded the int maximum value. So, after identifying that issue, we modified the code to use 1UL instead of 1 to increase the length of the data type and prevent overflow, and the bug was fixed.

# 4  Key Lessons

a) Design and Optimization:

- Gaining insight into the functionalities of various sub-stages in the system design and algorithmic aspects of CDC for processing Ethernet input data. This includes the incorporation of SHA-256 for deduplication, the utilization of LZW compression, and implementing bit-packing techniques.

- Systematically improving function performance by addressing hashing function requirements, selecting data types compatible with both text and binary inputs, and eliminating unnecessary memory copy functions.

- Applying vectorization techniques with the NEON accelerator to enhance bandwidth, informed by lessons learned from Homework 4. Identifying application bottlenecks through initial calculations and a simplified software version.

- Leveraging Vitis Analyzer for debugging, port assignment, and kernel execution tracing. Utilizing VITIS pragmas to optimize FPGA functions in terms of initiation interval (II) and resource utilization.

- Throughout this submission, we found Git to be an invaluable tool, enabling us to maintain tagged code versions. This proved essential for identifying major commits that significantly contributed to functionality or signaled upcoming substantial code revisions.

b) Debugging:

- We found GDB to be quite useful while debugging the code. It enabled us to use breakpoints, to jump to specific areas in the code and examine everything carefully to pinpoint the issue. It also helped us achieve full functionality of the kernel when working with Vitis, since we used it outside Vitis by compiling our kernel and the testbench with g++.

- The packet emulation using file read was also really helpful when testing out the functionality of the complete application.

- Including print statements helped examine the state of the output from the kernel since we couldn't access the host buffer mapped to the kernel using OpenCL in GDB.

- Vitis Analyzer proved to be helpful when we were facing a synchronization problem in the host code, and looking at the timeline trace helped us understand the issue and rectify it.

- The data flow diagram and the scheduler view in Vitis helped us visualize the control flow and the operations scheduled at various cycles enabling us to pinpoint the source code that caused a low II and fix it.

c) Teamwork and Collaboration:

- We gained valuable insights into effective collaboration by using version control systems like Git and GitHub, where we honed our skills in branching, merging, and resolving conflicts seamlessly within the team.

- Task distribution among team members for design, optimization, and debugging was executed with efficiency, ensuring a well-coordinated effort.

- Our regular team meetings became a platform for constructive feedback, focusing on aspects like code quality, efficiency, and addressing any issues promptly.

- Emphasizing the significance of investing time in thorough testing and debugging to prevent integration issues during the final stages of the project.

- Communication played a pivotal role; engaging in discussions about individual tasks, timelines, and challenges faced by team members led to more efficient problem-solving.

- Navigating changes in project requirements was approached with flexibility, adapting to unforeseen challenges and proactively adjusting the project plan accordingly.

- Recognition and celebration of team achievements and milestones became integral, contributing to the fostering of a positive team culture and maintaining high motivation levels.

# 5 Design Space Exploration, Graphs and Models

## 5.1 CDC

**Before moving to FastCDC:**
The main design space that we explored was the chunk size that gets set by the modulus value. Initially this value was set to 1024 but we had to eventually increase this to 4096. A smaller chunk size would imply more deduplication but the compression and overall throughput increases with larger chunks. As mentioned in the project handout, we wanted to chunk with at least a size of 4096 bytes. Another thing that drastically helped improve our CDC throughput was pre calculating powers of 3 (which was our chosen PRIME) and storing them in a local buffer. Since the hash calculation was using powers of 3 up to 17 (this was based on our WIN_SIZE used for the hash) repetitively for each iteration, we decided to pre-calculate these and store them in an array from which they can be read for hash calculations. This made our CDC significantly quicker.

**After moving to FastCDC:**

To further optimize our CDC, we changed to FastCDC which used the GEAR hash and skipped sub-minimum chunk cut-points. Fast CDC theorizes that some of the hash calculations can be skipped when defining chunk boundaries in favor of a low probability that a chunk will occur at a position less than the desired average chunk size. To explain this let's consider that we want chunks of size 4096, now the FastCDC says that the probability of drawing a chunk boundary at say 1024 (defined as the minimum chunk size) is very low, and hence we can skip hash calculation until the first 1024 bytes in a chunk. This drastically reduces the latency of CDC while efficiently drawing good chunk boundaries.

Now the GEAR hash consists of essentially a lookup table of 256 values and the hash value for a certain character is the value stored at the index represented by the ASCII value of that character. This means that the hash calculation only requires an array lookup. This along with the skipping of hash calculations, made our final CDC implementation much faster as opposed to regular CDC.

**Data: Test Case - Franklin Text File (File Size: 399054 B)**

| Chunk Size | Total Throughput (Mb/s) | CDC Throughput (Mb/s) | SHA Throughput (Mb/s) | LZW Throughput (Mb/s) | DEDUP Throughput (Mb/s) |
|---|---|---|---|---|---|
| 2048 | 58.4843 | 1781.58 | 551.707 | 79.0904 | 5561.34 |
| 4096 | 67.443 | 1178.01 | 713.09 | 92.6546 | 8704.41 |
| 8192 | 72.8656 | 1000.7 | 851.229 | 100.217 | 13057 |

Table 7: Chunk Size Versus Throughput - Franklin

| Chunk Size | Contributed DEDUP (B) | Contributed LZW (B) | Compression Ratio |
|---|---|---|---|
| 2048 | 0 | 327310 | 0.82 |
| 4096 | 0 | 291731 | 0.73 |
| 8192 | 0 | 259191 | 0.65 |

Table 8: Chunk Size Versus Compression Ratio - Franklin

The graph below was generated by using the data points that we have and extrapolating them using the power rule, i.e., $y = Ax^B$. This enabled us to estimate the values of throughput and compression ratios for values of chunk size ranging from 256 to 65536, and then plot them.

Figure 11: Chunk Size VS Throughput VS Compression Ratio Graph

**Data: Test Case - vmlinuz_small.tar (File Size: 399054 B)**

| Chunk Size | Total Throughput (Mb/s) | CDC Throughput (Mb/s) | SHA Throughput (Mb/s) | LZW Throughput (Mb/s) | DEDUP Throughput (Mb/s) |
|---|---|---|---|---|---|
| 2048 | 73.3759 | 1639.94 | 593.419 | 89.3484 | 10673.8 |
| 4096 | 87.4704 | 1099.73 | 776.749 | 109.426 | 16190.4 |
| 8192 | 94.3199 | 965.239 | 911.699 | 119.165 | 20780 |

Table 9: Chunk Size Versus Throughput - vmlinuz_small

| Chunk Size | Contributed DEDUP (B) | Contributed LZW (B) | Compression Ratio |
|---|---|---|---|
| 2048 | 768 | 643 | 0.004 |
| 4096 | 380 | 734 | 0.003 |
| 8192 | 184 | 938 | 0.003 |

Table 10: Chunk Size Versus Compression Ratio - vmlinuz_small

## 5.2 SHA 256

Initially we were using a software implementation from the mpsoc-crypto library provided in the project handout. To optimize this further, we moved the SHA implementation onto the ARM NEON cores since this was an effective method of exploiting parallelism.

For our final implementation, we adapted the NEON intrinsics version of the implementation. The mpsoc -crypto library provided in the project handout was used for this. This uses the SHA cryptographic intrinsics as well as the NEON intrinsics which gives us significant speedup for the SHA since it uses 8 vector lanes.

## 5.3  LZW Kernel

We tried to exploit various design axes while trying to optimize the LZW kernel. We started by removing the for loop from the associative memory lookup (assoc_lookup) function and replaced these with bit shifting and manipulation operations. Since, 'for' loops are computation intensive, this change helped us decrease the latency significantly. Additionally, the main LZW loop initially consisted of a while loop. We changed this to a for loop since they are much faster and to help Vitis HLS know that the number of iterations are already known. We pipelined these loops to help decrease the II and make it more efficient. We also implemented load-store-compute which enabled us to stream data across these functions in the kernel.

We were also getting many hash collisions with the original hash function. To counter this, we also adapted the murmur hash in our LZW kernel implementation to reduce the hash collisions. This also allowed us to decrease the size of our hash table which resulted in reduced lookup times. As mentioned in the previous milestones, we also made use of 2 buckets instead of the one since we were getting multiple hash collisions for bigger binary files. This meant that each index was associated with 2 key-value pairs instead of just one. As a result of which, the associative memory would only come into use when 2 keys have already been previously hashed to the same index and entered into the main hash table. This decreased the entries in our associative memory significantly.

To reduce the number of calls to the kernel and minimize the data transfer overhead we batched all the chunks in a 16KB buffer and sent that as a whole to the kernel in the kernel it would iterate over all the chunks and compute LZW. This gave us a huge speedup in comparison to calling the kernel in a loop based on the output from DEDUP. Once the data was received back from the kernel we would use the output from DEDUP to decide what to write into the compressed file.

Apart from that to minimize the kernel overhead we reduced the amount of input arguments required. Firstly, the chunk_indices array for the batched chunks that needed to be sent to LZW consisted of its first element being set to the length of the array; this way we avoided sending another parameter to the kernel and decreasing the setup time. Similarly for the out_packet_lengths argument of the kernel which populated the length of the encoded LZW chunks had its first element signify if insertion into the associative memory had failed. This helped us further reduce the data transfer time. As described previously in section III under "Challenges and Debugging", we noticed that with binary files the code value in LZW exceeded 4096, which was the configuration for our hash table lookup and bit packing. Due to this we decided to move to a 13-bit packing routine and modified the hash table to support 13-bit values, and updated the size of the key accordingly. This helped us mitigate the issue of the code exceeding 4096 and ensured that our design would be stable and functional for our configured chunk size.

Another optimization that we tried was moving the bit packing inside the kernel and producing the bit-packed output for the batched chunks. We thought that this would give us a significant speedup, instead it slowed down our kernel and the overall performance of the application took a hit. We tried to optimize the bit-packing in Vitis by pipelining and unrolling wherever possible, but it didn't help. Due to this we decided to move the bit-packing routine outside the kernel and back onto the ARM processor.

Finally, another optimization that we tried was to reduce the number of parameters sent to our compression_pipeline function, since the OpenCL buffers were set up in main and then they were passed to the compression_pipeline function along with their host mapped buffers. Apart from this the function took in a lot of other parameters like the OpenCL command queue, the kernel, the file pointer to write into the file, etc. Hence to minimize the amount of arguments, we tried to wrap the OpenCL buffers and their host counterparts into a struct and tried to pass a pointer to that struct, but when we did that we weren't able to properly retrieve data back into our buffer from the kernel and decided to get rid of it, in interest of complete functionality.

# 6 Future Scope

In an attempt to optimize this further, we thought of implementing thread pooling where different stages could be run on different threads. By using multiple threads, we could explore pipelining LZW with other parts of the pipeline such as data collection, CDC, SHA and Deduplication. This would enable us to overlap the computation of LZW with that of other stages.

Additionally, we could try to run multiple compute units concurrently for different chunks since each chunk is independent of one another. We tried implementing this for our final design, however, since we were using almost 65% of the available LUTs, we could not replicate this design to run on multiple compute units.

As of now when the chunks are batched in the kernel, we are sequentially iterating over all the chunks and computing LZW on them individually at a time. Given a better (lower) resource utilization we could unroll the for loop that processes these chunks, and potentially use the dataflow pragma to simulate multiple compute units, i.e., running multiple instances of LZW concurrently.

Apart from that another optimization that we could potentially make is to parallelize the lookup between the buckets in the hash table and the associative memory which would significantly speed up the LZW process and improve the performance overall. Furthermore, we could also use better hash functions and techniques like double hashing to reduce the size of our hash table and improve our resource utilization. This in-turn ties up to the problem of not being able to completely unroll the loop, and effectively leverage the resources present on the FPGA.

# 7 Task Distribution

We allocated tasks among ourselves to design and optimize the project. Akhil Gunda focused on enhancing the CDC implementation, optimizing with Fast CDC implementation and the host code, Siddhant Mathur concentrated on adapting SHA to use NEON intrinsics and seamlessly integrating it into the entire system along with the initial version of 12-bit packing, while Ruturaj Nanoti tackled the testbench, optimized portions of the LZW kernel, had to come up with optimized 13 bit packing to support bigger binaries files. Additionally, we collaborated on debugging the final stage performance and issue and on

further optimizing the LZW collectively. Our combined efforts ensured a well-rounded approach to achieving the project's complete functionality.

We, **Ruturaj A. Nanoti**, **Akhil Gunda**, and **Siddhant Mathur**, certify that we have complied with the University of Pennsylvania's Code of Academic Integrity in completing this final exercise.

# 8  Appendix

## 8.1  Vivado Results

- Control interface of the accelerator wrapper lzw_1 screenshot:



Figure 12: Vivado Ports

For lzw_1 the s_axi_control is mapped to 0x00_B000_0000 to 0x00_B000_FFFF.

- Vivado Block design under IP integrator:



Figure 13: Vivado Design Block Diagram

Figure 14: Vivado Design Overall Block Diagram

- Device (Design Mapping):



Figure 15: Vivado Design Mapping

## 8.2  Vitis HLS and Vitis Analyzer Results

- Co-Simulation Screenshot:



Figure 16: Vitis Co-Sim Result

- System Diagram:



Figure 17: Vitis System Diagram

- Profile Summary:



Figure 18: Vitis Profile Summary

Figure 19: Kernel and Compute Units

## 8.3 Code

### 8.3.1 CDC

```
1 #include <iostream>
2 #include <math.h>
3 #include <stdint.h>
4 #include <stdio.h>
5 #include <stdlib.h>
6 #include <vector>
7
8 using namespace std;
9
10 #define CHUNK_SIZE 4096
11 #define MODULUS_MASK (CHUNK_SIZE - 1)
12 #define MODULUS_MASK_S ((CHUNK_SIZE * 2) - 1)
13 #define MODULUS_MASK_L ((CHUNK_SIZE / 2) - 1)
14 #define TARGET 0
15 #define BITS 12
16 #define MINSIZE 1024
17 #define BUFFER_LEN 16384
18
19 uint64_t GEAR[256] = {
20     1553318008, 574654857,  759734804,  310648967,  1393527547,
    1195718329,
21     694400241,  1154184075, 1319583805, 1298164590, 122602963,
    989043992,
22     1918895050, 933636724,  1369634190, 1963341198, 1565176104,
    1296753019,
23     1105746212, 1191982839, 1195494369, 29065008,   1635524067,
    722221599,
24     1355059059, 564669751,  1620421856, 1100048288, 1018120624,
    1087284781,
25     1723604070, 1415454125, 737834957,  1854265892, 1605418437,
    1697446953,
26     973791659,  674750707,  1669838606, 320299026,  1130545851,
    1725494449,
```

28

```
27    939321396,   748475270,   554975894,   1651665064,  1695413559,
      671470969,
28    992078781,   1935142196,  1062778243,  1901125066,  1935811166,
      1644847216,
29    744420649,   2068980838,  1988851904,  1263854878,  1979320293,
      111370182,
30    817303588,   478553825,   694867320,   685227566,   345022554,
      2095989693,
31    1770739427,  165413158,   1322704750,  46251975,    710520147,
      700507188,
32    2104251000,  1350123687,  1593227923,  1756802846,  1179873910,
      1629210470,
33    358373501,   807118919,   751426983,   172199468,   174707988,
      1951167187,
34    1328704411,  2129871494,  1242495143,  1793093310,  1721521010,
      306195915,
35    1609230749,  1992815783,  1790818204,  234528824,   551692332,
      1930351755,
36    110996527,   378457918,   638641695,   743517326,   368806918,
      1583529078,
37    1767199029,  182158924,   1114175764,  882553770,   552467890,
      1366456705,
38    934589400,   1574008098,  1798094820,  1548210079,  821697741,
      601807702,
39    332526858,   1693310695,  136360183,   1189114632,  506273277,
      397438002,
40    620771032,   676183860,   1747529440,  909035644,   142389739,
      1991534368,
41    272707803,   1905681287,  1210958911,  596176677,   1380009185,
      1153270606,
42    1150188963,  1067903737,  1020928348,  978324723,   962376754,
      1368724127,
43    1133797255,  1367747748,  1458212849,  537933020,   1295159285,
      2104731913,
44    1647629177,  1691336604,  922114202,   170715530,   1608833393,
      62657989,
45    1140989235,  381784875,   928003604,   449509021,   1057208185,
      1239816707,
46    525522922,   476962140,   102897870,   132620570,   419788154,
      2095057491,
47    1240747817,  1271689397,  973007445,   1380110056,  1021668229,
      12064370,
48    1186917580,  1017163094,  597085928,   2018803520,  1795688603,
      1722115921,
49    2015264326,  506263638,   1002517905,  1229603330,  1376031959,
      763839898,
50    1970623926,  1109937345,  524780807,   1976131071,  905940439,
      1313298413,
51    772929676,   1578848328,  1108240025,  577439381,   1293318580,
      1512203375,
52    371003697,   308046041,   320070446,   1252546340,  568098497,
      1341794814,
53    1922466690,  480833267,   1060838440,  969079660,   1836468543,
      2049091118,
54    2023431210,  383830867,   2112679659,  231203270,   1551220541,
      1377927987,
55    275637462,   2110145570,  1700335604,  738389040,   1688841319,
      1506456297,
```

```
56     1243730675,  258043479,   599084776,   41093802,    792486733,
       1897397356,
57      28077829,   1520357900,  361516586,   1119263216,  209458355,
       45979201,
58      363681532,   477245280,   2107748241,  601938891,   244572459,
       1689418013,
59      1141711990,  1485744349,  1181066840,  1950794776,  410494836,
       1445347454,
60      2137242950,  852679640,   1014566730,  1999335993,  1871390758,
       1736439305,
61      231222289,   603972436,   783045542,   370384393,   184356284,
       709706295,
62      1453549767,  591603172,   768512391,   854125182,
63 };
64
65 void fast_cdc(unsigned char *buff, unsigned int buff_size,
66               unsigned int chunk_size, vector<uint32_t> &vect) {
67     unsigned int hash = 0;
68     unsigned int i = MINSIZE;
69     unsigned int prev = 0;
70
71     vect.push_back(0);
72     unsigned int modulus_mask_s = (chunk_size * 2) - 1;
73
74     while (i < buff_size) {
75         hash = (hash >> 1) + GEAR[buff[i]];
76         if (((hash & modulus_mask_s) == TARGET) ||
77             (((i - prev)) == chunk_size)) {
78             vect.push_back(i);
79             prev = i;
80             i += MINSIZE;
81             hash = 0;
82         }
83         i += 1;
84     }
85
86     if (vect[vect.size() - 1] != buff_size)
87         vect.push_back(buff_size);
88 }
```

### 8.3.2   SHA

```
1 #include "common.h"
2 #include "sha256/sha256.h"
3 #include <iomanip>
4 #include <iostream>
5 #include <sstream>
6 #include <stdio.h>
7 #include <string.h>
8 #include <vector>
9
10 using namespace std;
11
12 string bytearray_hex_string(unsigned char *bytes, int size) {
13     stringstream ss;
14     ss << hex << setfill('0');
15     for (int i = 0; i < size; i++) {
```

```
16          ss << setw (2) << static_cast < int >( bytes [i ]);
17      }
18      return ss.str ();
19 }
20
21 string sha_256(unsigned char *chunked_data, uint32_t chunk_start_idx,
22                 uint32_t chunk_end_idx) {
23      SHA256_CTX ctx;
24      BYTE hash_val [32];
25
26      sha256_hash(&ctx, (const BYTE *)(chunked_data + chunk_start_idx),
27                  chunk_start_idx, chunk_end_idx, hash_val, 1);
28
29      string sha_fingerprint = bytearray_hex_string(hash_val, 32);
30
31      return sha_fingerprint;
32 }
```

### 8.3.3  DEDUP

```
1 #include "common.h"
2 #include < cstdlib >
3 #include < iostream >
4 #include < stdbool .h >
5 #include < unordered_map >
6
7 using namespace std;
8
9 // Return -1 on failure.
10 int64_t dedup(string sha_fingerprint) {
11
12      static unordered_map < string , int64_t > sha_chunk_id_map;
13      static int64_t chunk_id = 0;
14
15      bool found =
16          (sha_chunk_id_map.find(sha_fingerprint) == sha_chunk_id_map.end
    ()
17                ? false
18                : true);
19
20      // Perform lookup in map here.
21      if (!found) {
22          // Insert into map before calling LZW.
23          sha_chunk_id_map[sha_fingerprint] = chunk_id;
24          ++chunk_id;
25          return -1;
26      } else
27          return sha_chunk_id_map[sha_fingerprint];
28 }
```

### 8.3.4  LZW

```
1 #include < hls_stream .h >
2 #include < stdint .h >
3 #include < stdio .h >
4 #include < stdlib .h >
```

31

```c
#include <string.h>

#define CAPACITY                                                        \
    16384 // hash output is 15 bits, and we have 1 entry per bucket, so
          // capacity
          // is 2^15
#define SEED 524057
#define ASSOCIATIVE_MEM_STORE 64
#define CHUNK_SIZE 4096
#define MAX_CHUNK_SIZE 8192
#define INFO_LEN 4
#define BUFFER_LEN 16384
#define MAX_ITERATIONS 20
#define MAX_OUTPUT_CODE_SIZE 40960

#define SEED_MURMUR 524057
#define SEED_CRAP 75768593

#define FNV_PRIME 16777619
#define FNV_OFFSET_BASIS 2166136261U

typedef enum InfoParams {
    INFO_START_IDX,
    INFO_END_IDX,
    INFO_OUT_PACKET_LENGTH,
    INFO_FAILURE,
} InfoParams;

//********************************************************************
typedef struct {
    // Each key_mem has a 9 bit address (so capacity = 2^9 = 512)
    // and the key is 20 bits, so we need to use 3 key_mems to cover all
    the key
    // bits. The output width of each of these memories is 64 bits, so
    we can
    // only store 64 key value pairs in our associative memory map.

    unsigned long upper_key_mem[512]; // the output of these  will be 64
    bits
                                      // wide (size of unsigned long).
    unsigned long middle_key_mem[512];
    unsigned long lower_key_mem[512];
    unsigned long
        value[ASSOCIATIVE_MEM_STORE]; // value store is 64 deep, because
    the
                                      // lookup mems are 64 bits wide
    unsigned int fill; // tells us how many entries we've currently
    stored
} assoc_mem;

// cast to struct and use ap types to pull out various feilds.
//********************************************************************

static inline uint32_t murmur_32_scramble(uint32_t k) {
    k *= 0xcc9e2d51;
    k = (k << 15) | (k >> 17);
    k *= 0x1b873593;
```

```
56      return k;
57  }
58
59  unsigned int inline murmur_hash(unsigned long key) {
60      uint32_t h = SEED;
61      uint32_t k = key;
62
63      h ^= murmur_32_scramble(k);
64      h = (h << 13) | (h >> 19);
65      h = h * 5 + 0xe6546b64;
66
67      h ^= murmur_32_scramble(k);
68      /* Finalize. */
69      h ^= h >> 16;
70      h *= 0x85ebca6b;
71      h ^= h >> 13;
72      h *= 0xc2b2ae35;
73      h ^= h >> 16;
74      return h;
75  }
76
77  uint32_t inline Crap8(unsigned int key) {
78  #define c8fold(a, b, y, z)                      \
79      {                                           \
80          p = (uint32_t)(a) * (uint64_t)(b);      \
81          y ^= (uint32_t)p;                       \
82          z ^= (uint32_t)(p >> 32);               \
83      }
84  #define c8mix(in)                               \
85      {                                           \
86          h *= m;                                 \
87          c8fold(in, m, k, h);                    \
88      }
89
90      const uint32_t m = 0x83d2e73b, n = 0x97e1cc59;
91      const uint32_t key4[4] = {((key)&0xFF), ((key >> 8) & 0xFF),
92                                ((key >> 16) & 0xFF), ((key >> 24) & 0xFF)
93      };
93      uint32_t h = SEED_CRAP, k = n;
94      uint64_t p;
95
96      c8mix(key4[0]) c8mix(key4[1]) c8mix(key4[2])
97          c8mix(key4[3] & ((1 << (2 * 8)) - 1)) c8fold(h ^ k, n, k, k)
        return k;
98  }
99
100 unsigned int fnv1a_hash(unsigned int key) {
101     uint32_t hash = FNV_OFFSET_BASIS;
102
```

```c
103      hash ^= ((key >> 24) & 0xFF);
104      hash *= FNV_PRIME;
105
106      hash ^= ((key >> 16) & 0xFF);
107      hash *= FNV_PRIME;
108
109      hash ^= ((key >> 8) & 0xFF);
110      hash *= FNV_PRIME;
111
112      hash ^= ((key)&0xFF);
113      hash *= FNV_PRIME;
114
115      return hash;
116 }
117
118 unsigned int djb2_hash(unsigned int key) {
119      unsigned int hash = 5381;
120      int c;
121
122      hash = ((hash << 5) + hash) + ((key >> 24) & 0xFF); // hash * 33 + c
123      hash = ((hash << 5) + hash) + ((key)&0xFF);         // hash * 33 + c
124      hash = ((hash << 5) + hash) + ((key >> 8) & 0xFF);  // hash * 33 + c
125      hash = ((hash << 5) + hash) + ((key >> 24) & 0xFF); // hash * 33 + c
126
127      return hash;
128 }
129
130 unsigned int my_hash(unsigned long key) {
131      unsigned int hash_1 = murmur_hash(~key);
132      // unsigned int hash_2 = Crap8(hash_1);
133      // unsigned int hash_1 = djb2_hash(~(key));
134      // unsigned int hash_2 = fnv1a_hash(key);
135      // return (((((hash_2 >> 3) & 0x3F) << 7) | (((hash_3 >> 5) & 0x7) <<
        4) |
136      // ((hash_1 >> 7) & 0xF)); return (((((hash_2 >> 3) & 0x7F) << 7) |
        ((hash_1
137      // >> 7) & 0x7F)); return ((hash_1 + (SEED_CRAP ^ hash_2)) >> 7) &
138      // (CAPACITY_MOD);
139      return hash_1 & (CAPACITY - 1);
140 }
141
142 void inline hash_lookup(unsigned long (*hash_table)[2], unsigned int key
        ,
143                          bool *hit, unsigned int *result) {
144
145      key &= 0x1FFFFF; // make sure key is only 21 bits
146
147      unsigned int hash_val = my_hash(key);
148
149      unsigned long lookup = hash_table[hash_val][0];
150
151      // [valid][value][key]
152      unsigned long stored_key = lookup & 0x1FFFFF;      // stored key is
        21 bits
153      unsigned long value = (lookup >> 21) & 0x1FFF;     // value is 13
        bits
154      unsigned long valid = (lookup >> (21 + 13)) & 0x1; // valid is 1 bit
155
```

```
156     if (valid && (key == stored_key)) {
157         *hit = 1;
158         *result = value;
159         //*is_exists = 1;
160         return;
161     }
162
163     lookup = hash_table[hash_val][1];
164
165     // [valid][value][key]
166     stored_key = lookup & 0x1FFFFF;      // stored key is 21 bits
167     value = (lookup >> 21) & 0x1FFF;     // value is 13 bits
168     valid = (lookup >> (21 + 13)) & 0x1; // valid is 1 bit
169     if (valid && (key == stored_key)) {
170         *hit = 1;
171         *result = value;
172         return;
173     }
174     *hit = 0;
175     *result = 0;
176 }
177
178 void inline hash_insert(unsigned long (*hash_table)[2], unsigned int key
    ,
179                         unsigned int value, bool *collision) {
180
181     key &= 0x1FFFFF; // make sure key is only 21 bits
182     value &= 0x1FFF; // value is only 13 bits
183
184     unsigned int hash_val = my_hash(key);
185
186     unsigned long lookup = hash_table[hash_val][0];
187     unsigned long valid = (lookup >> (21 + 13)) & 0x1;
188
189     if (!valid) {
190         hash_table[hash_val][0] =
191             (1UL << (21 + 13)) | ((unsigned long)(value) << 21) | key;
192         *collision = 0;
193         return;
194     }
195
196     lookup = hash_table[hash_val][1];
197     valid = (lookup >> (21 + 13)) & 0x1;
198     if (valid) {
199         *collision = 1;
200         return;
201     }
202     hash_table[hash_val][1] =
203         (1UL << (21 + 13)) | ((unsigned long)(value) << 21) | key;
204     *collision = 0;
205 }
206
207 void inline assoc_insert(assoc_mem *mem, unsigned int key, unsigned int
    value,
208                         bool *collision) {
209     key &= 0x1FFFFF; // make sure key is only 21 bits
210     value &= 0x1FFF; // value is only 13 bits
211
```

```
212    unsigned int mem_fill = mem->fill;
213
214    if (mem_fill < ASSOCIATIVE_MEM_STORE) {
215        unsigned int key_high = (key >> 18) & 0x1FF;
216        unsigned int key_middle = (key >> 9) & 0x1FF;
217        unsigned int key_low = (key)&0x1FF;
218        mem->upper_key_mem[key_high] |=
219            (1UL << mem_fill); // set the fill'th bit to 1, while
    preserving
220                                // everything else
221        mem->middle_key_mem[key_middle] |=
222            (1UL << mem_fill); // set the fill'th bit to 1, while
    preserving
223                                // everything else
224        mem->lower_key_mem[key_low] |=
225            (1UL << mem_fill); // set the fill'th bit to 1, while
    preserving
226                                // everything else
227        mem->value[mem_fill] = value;
228        mem->fill = mem_fill + 1;
229        *collision = 0;
230        return;
231    }
232    *collision = 1;
233 }
234
235 void inline assoc_lookup(assoc_mem *mem, unsigned int key, bool *hit,
236                         unsigned int *result) {
237    key &= 0x1FFFFF; // make sure key is only 21 bits
238    unsigned int key_high = (key >> 18) & 0x1FF;
239    unsigned int key_middle = (key >> 9) & 0x1FF;
240    unsigned int key_low = (key)&0x1FF;
241
242    unsigned long match_high = mem->upper_key_mem[key_high];
243    unsigned long match_middle = mem->middle_key_mem[key_middle];
244    unsigned long match_low = mem->lower_key_mem[key_low];
245
246    unsigned long match = match_high & match_middle & match_low;
247
248    if (match == 0) {
249        *hit = 0;
250        return;
251    }
252
253    unsigned int address = 0; //(unsigned int)(log2(match & -match) + 1)
    ;
254
255    // Right shift until the rightmost set bit is found
256    address += ((match & 0x00000000FFFFFFFFUL) == 0) ? 32 : 0;
257    match >>= ((match & 0x00000000FFFFFFFFUL) == 0) ? 32 : 0;
258
259    address += ((match & 0x000000000000FFFFUL) == 0) ? 16 : 0;
260    match >>= ((match & 0x000000000000FFFFUL) == 0) ? 16 : 0;
261
262    address += ((match & 0x00000000000000FFUL) == 0) ? 8 : 0;
263    match >>= ((match & 0x00000000000000FFUL) == 0) ? 8 : 0;
264
265    address += ((match & 0x000000000000000FUL) == 0) ? 4 : 0;
```

```
266      match >>= ((match & 0x000000000000000FUL) == 0) ? 4 : 0;
267
268      address += ((match & 0x0000000000000003UL) == 0) ? 2 : 0;
269      match >>= ((match & 0x0000000000000003UL) == 0) ? 2 : 0;
270
271      address += ((match & 0x0000000000000001UL) == 0) ? 1 : 0;
272      match >>= ((match & 0x0000000000000001UL) == 0) ? 1 : 0;
273
274      if (address != ASSOCIATIVE_MEM_STORE) {
275          *result = mem->value[address];
276          *hit = 1;
277          return;
278      }
279      *hit = 0;
280 }
281
282 void inline insert(unsigned long (*hash_table)[2], assoc_mem *mem,
283                    unsigned int key, unsigned int value, bool *collision
    ) {
284      hash_insert(hash_table, key, value, collision);
285      if (*collision) {
286          assoc_insert(mem, key, value, collision);
287      }
288 }
289
290 void inline lookup(unsigned long (*hash_table)[2], assoc_mem *mem,
291                    unsigned int key, bool *hit, unsigned int *result) {
292      hash_lookup(hash_table, key, hit, result);
293      if (!*hit) {
294          assoc_lookup(mem, key, hit, result);
295      }
296 }
297
298 static void compute_lzw(hls::stream<unsigned char> &input_stream,
299                         hls::stream<uint32_t> &out_stream,
300                         uint32_t generic_info[4]) {
301
302      // create hash table and assoc mem
303      unsigned long hash_table[CAPACITY][2];
304      assoc_mem my_assoc_mem;
305
306 #pragma HLS array_partition variable = hash_table complete dim = 2
307
308 // make sure the memories are clear
309 LOOP1:
310      for (int i = 0; i < CAPACITY; i++) {
311 #pragma HLS UNROLL factor = 512
312          hash_table[i][0] = 0;
313          hash_table[i][1] = 0;
314      }
315      my_assoc_mem.fill = 0;
316
317 LOOP2:
318      for (int i = 0; i < 512; i++) {
319 #pragma HLS UNROLL
320          my_assoc_mem.upper_key_mem[i] = 0;
321          my_assoc_mem.middle_key_mem[i] = 0;
322          my_assoc_mem.lower_key_mem[i] = 0;
```

```
323          }
324
325      unsigned int next_code = 256;
326      uint8_t failure = 0;
327      uint32_t start_idx = generic_info[INFO_START_IDX];
328      uint32_t end_idx = generic_info[INFO_END_IDX];
329      unsigned int prefix_code = input_stream.read();
330      unsigned int code = 0;
331      unsigned char next_char = 0;
332      uint64_t j = 0;
333
334 LOOP3:
335      for (int i = start_idx; i < end_idx - 1; i++) {
336          next_char = input_stream.read();
337          bool hit = 0;
338          lookup(hash_table, &my_assoc_mem, ((prefix_code << 8) +
     next_char),
339                  &hit, &code);
340          if (!hit) {
341              out_stream.write(prefix_code);
342              bool collision = 0;
343              insert(hash_table, &my_assoc_mem, ((prefix_code << 8) +
     next_char),
344                      next_code, &collision);
345              if (collision) {
346                  failure = 1;
347              }
348              next_code += 1;
349              ++j;
350              prefix_code = next_char;
351          } else {
352              prefix_code = code;
353          }
354      }
355
356      out_stream.write(prefix_code);
357      generic_info[INFO_OUT_PACKET_LENGTH] = j + 1;
358      if (failure)
359          generic_info[INFO_FAILURE] = failure;
360 }
361
362 static void store_data(hls::stream<uint32_t> &out_stream, uint32_t *
     lzw_codes) {
363      unsigned int i = 0;
364
365      while (!out_stream.empty()) {
366          if (i < MAX_OUTPUT_CODE_SIZE) {
367              lzw_codes[i] = out_stream.read();
368              i++;
369          }
370      }
371 }
372
373 static void load_data(unsigned char input[BUFFER_LEN],
374                      hls::stream<unsigned char> &input_stream) {
375      for (int i = 0; i < BUFFER_LEN; i++) {
376          input_stream.write(input[i]);
377      }
```

```
378  }
379
380  static void compute_data(hls::stream<unsigned char> &input_stream,
381                          hls::stream<uint32_t> &out_stream,
382                          uint32_t out_packet_lengths[MAX_ITERATIONS],
383                          uint32_t temp_chunk_indices[MAX_ITERATIONS]) {
384
385      // This is an array that packs generic information for the LZW
386      // function. The elements in the this array are as follows:
387      // 0 - start_idx
388      // 1 - end_idx
389      // 2 - out_packet_length
390      // 3 - failure
391      uint32_t generic_info[INFO_LEN];
392      generic_info[INFO_FAILURE] = 0;
393      generic_info[INFO_OUT_PACKET_LENGTH] = 0;
394
395      // The first element in the chunk_indices array contains the size
396      // of the chunk_indices array.
397      uint32_t chunk_indices_len = temp_chunk_indices[0];
398
399
400      const int bound = chunk_indices_len - 1;
401
402      for (int i = 1; i <= MAX_ITERATIONS; i++) {
403  #pragma HLS UNROLL factor = 2
404          if (i <= bound) {
405              generic_info[INFO_START_IDX] = temp_chunk_indices[i];
406              generic_info[INFO_END_IDX] = temp_chunk_indices[i + 1];
407              compute_lzw(input_stream, out_stream, generic_info);
408              out_packet_lengths[i] = generic_info[INFO_OUT_PACKET_LENGTH
409      ];
410          }
410      }
411
412      while (!input_stream.empty()) {
413          input_stream.read();
414      }
415
416      // The first element of the out_packet_lengths is going to signify
417      // failure to insert into the associative memory.
418      out_packet_lengths[0] = generic_info[INFO_FAILURE];
419  }
420
421  static void perform_lzw(unsigned char input[BUFFER_LEN],
422                          uint32_t temp_chunk_indices[MAX_ITERATIONS],
423                          uint32_t *lzw_codes,
424                          uint32_t out_packet_lengths[MAX_ITERATIONS]) {
425
426      hls::stream<unsigned char> in_stream("chunk_in");
427      hls::stream<uint32_t> out_stream("lzw_out");
428
429  #pragma HLS STREAM variable = in_stream depth = 16384
430  #pragma HLS STREAM variable = out_stream depth = 40960
431
432      // #pragma HLS DATAFLOW
433      load_data(input, in_stream);
434      compute_data(in_stream, out_stream, out_packet_lengths,
```

```
             temp_chunk_indices);
435        store_data(out_stream, lzw_codes);
436 }
437
438 void lzw(unsigned char input[BUFFER_LEN], uint32_t *lzw_codes,
439            uint32_t chunk_indices[MAX_ITERATIONS],
440            uint32_t out_packet_lengths[MAX_ITERATIONS]) {
441
442 #pragma HLS INTERFACE m_axi port = input depth = 16384 bundle = p0
443 #pragma HLS INTERFACE m_axi port = lzw_codes depth = 40960 bundle = p1
444 #pragma HLS INTERFACE m_axi port = chunk_indices depth = 20 bundle = p0
445 #pragma HLS INTERFACE m_axi port = out_packet_lengths depth = 20 bundle
       = p0
446
447     uint32_t temp_chunk_indices[MAX_ITERATIONS] = {0};
448
449 #pragma HLS array_partition variable = temp_chunk_indices block factor =
             \
450     10 dim = 1
451
452 LOOP4:
453     for (int i = 0; i < MAX_ITERATIONS; i++) {
454 #pragma HLS UNROLL
455          temp_chunk_indices[i] = chunk_indices[i];
456     }
457
458     perform_lzw(input, temp_chunk_indices, lzw_codes, out_packet_lengths
    );
459 };
```

### 8.3.5  Host

```
1 #include "../Server/encoder.h"
2 #include "../Server/server.h"
3 #include "../Server/stopwatch.h"
4 #include "Utilities.h"
5 #include "common.h"
6 #include <chrono>
7 #include <condition_variable>
8 #include <errno.h>
9 #include <fcntl.h>
10 #include <iostream>
11 #include <pthread.h>
12 #include <stdint.h>
13 #include <stdio.h>
14 #include <stdlib.h>
15 #include <string.h>
16 #include <sys/mman.h>
17 #include <thread>
18 #include <unistd.h>
19 #include <unordered_map>
20 #include <vector>
21
22 #define NUM_PACKETS 2
23 #define pipe_depth 4
24 #define DONE_BIT_L (1 << 7)
25 #define DONE_BIT_H (1 << 15)
```

```
26
27  #define CL_HPP_CL_1_2_DEFAULT_BUILD
28  #define CL_HPP_TARGET_OPENCL_VERSION 120
29  #define CL_HPP_MINIMUM_OPENCL_VERSION 120
30  #define CL_HPP_ENABLE_PROGRAM_CONSTRUCTION_FROM_ARRAY_COMPATIBILITY 1
31  #define CL_USE_DEPRECATED_OPENCL_1_2_APIS
32
33  using namespace std;
34
35  uint64_t dedup_bytes = 0;
36  uint64_t lzw_bytes = 0;
37
38  stopwatch time_cdc;
39  stopwatch time_lzw;
40  stopwatch time_sha;
41  stopwatch time_dedup;
42  stopwatch total_time;
43
44  typedef struct __attribute__((packed)) RawData {
45      int length_sum;                      /// The total length of the input
        buffer.
46      unsigned char *pipeline_buffer; /// Input data that needs to be
        processed.
47      FILE *fptr_write;                    /// File pointer to write the output
        .
48      unsigned char *host_input;
49      uint32_t *output_codes;
50      uint32_t *chunk_indices;
51      uint32_t *output_code_lengths;
52  } RawData;
53
54  typedef struct CLDevice {
55      cl::Kernel kernel;
56      cl::CommandQueue queue;
57  } CLDevice;
58
59  void handle_input(int argc, char *argv[], int *blocksize, char **
        filename,
60                    char **kernel_name, unsigned int *chunk_size) {
61      int x;
62      extern char *optarg;
63
64      while ((x = getopt(argc, argv, ":b:f:k:c:")) != -1) {
65          switch (x) {
66          case 'k':
67              *kernel_name = optarg;
68              printf("Kernel name is set to %s optarg\n", *kernel_name);
69              break;
70          case 'b':
71              *blocksize = atoi(optarg);
72              printf("blocksize is set to %d optarg\n", *blocksize);
73              break;
74          case 'c':
75              *chunk_size = atoi(optarg);
76              printf("Chunk size is set to %d optarg\n", *chunk_size);
77              break;
78          case 'f':
79              *filename = optarg;
```

```
 80              printf("filename is %s optarg\n", *filename);
 81              break;
 82          case ':':
 83              printf("-%c without parameter\n", optopt);
 84              break;
 85          }
 86      }
 87  }
 88
 89  static unsigned char *create_packet(int32_t chunk_idx,
 90                                       uint32_t out_packet_length,
 91                                       uint32_t *out_packet, uint32_t
     packet_len) {
 92      unsigned char *data =
 93          (unsigned char *)calloc(packet_len, sizeof(unsigned char));
 94      CHECK_MALLOC(data, "Unable to allocate memory for new data packet");
 95
 96      uint32_t data_idx = 0;
 97      uint16_t current_val = 0;
 98      int bits_left = 0;
 99      int current_val_bits_left = 0;
100
101      for (uint32_t i = 0; i < out_packet_length; i++) {
102          current_val = out_packet[i];
103          current_val_bits_left = CODE_LENGTH;
104
105          if (bits_left == 0 && current_val_bits_left == CODE_LENGTH) {
106              data[data_idx] = (current_val >> 5) & 0xFF;
107              bits_left = 0;
108              current_val_bits_left = 5;
109              data_idx += 1;
110          }
111
112          if (bits_left == 0 && current_val_bits_left == 5) {
113              if (data_idx < packet_len) {
114                  data[data_idx] = (current_val & 0x1F) << 3;
115                  bits_left = 3;
116                  current_val_bits_left = 0;
117                  continue;
118              } else
119                  break;
120          }
121
122          if (bits_left == 3 && current_val_bits_left == CODE_LENGTH) {
123              data[data_idx] |= ((current_val >> 10) & 0x07);
124              bits_left = 0;
125              data_idx += 1;
126              current_val_bits_left = 10;
127          }
128
129          if (bits_left == 0 && current_val_bits_left == 10) {
130              if (data_idx < packet_len) {
131                  data[data_idx] = ((current_val >> 2) & 0xFF);
132                  bits_left = 0;
133                  data_idx += 1;
134                  current_val_bits_left = 2;
135              } else
136                  break;
```

```
137              }
138
139          if (bits_left == 0 && current_val_bits_left == 2) {
140              if (data_idx < packet_len) {
141                  data[data_idx] = (current_val & 0x03) << 6;
142                  bits_left = 6;
143                  current_val_bits_left = 0;
144                  continue;
145              } else
146                  break;
147          }
148
149          if (bits_left == 6 && current_val_bits_left == CODE_LENGTH) {
150              data[data_idx] |= ((current_val >> 7) & 0x3F);
151              bits_left = 0;
152              data_idx += 1;
153              current_val_bits_left = 7;
154          }
155
156          if (bits_left == 0 && current_val_bits_left == 7) {
157              if (data_idx < packet_len) {
158                  data[data_idx] = (current_val & 0x7F) << 1;
159                  bits_left = 1;
160                  current_val_bits_left = 0;
161                  continue;
162              } else
163                  break;
164          }
165
166          if (bits_left == 1 && current_val_bits_left == CODE_LENGTH) {
167              data[data_idx] |= ((current_val >> 12) & 0x1);
168              bits_left = 0;
169              data_idx += 1;
170              current_val_bits_left = 12;
171          }
172
173          if (bits_left == 0 && current_val_bits_left == 12) {
174              if (data_idx < packet_len) {
175                  data[data_idx] = ((current_val >> 4) & 0xFF);
176                  bits_left = 0;
177                  data_idx += 1;
178                  current_val_bits_left = 4;
179              } else
180                  break;
181          }
182
183          if (bits_left == 0 && current_val_bits_left == 4) {
184              if (data_idx < packet_len) {
185                  data[data_idx] = (current_val & 0x0F) << 4;
186                  bits_left = 4;
187                  current_val_bits_left = 0;
188                  continue;
189              } else
190                  break;
191          }
192
193          if (bits_left == 4 && current_val_bits_left == CODE_LENGTH) {
194              data[data_idx] |= ((current_val >> 9) & 0x0F);
```

```
195             data_idx += 1;
196             bits_left = 0;
197             current_val_bits_left = 9;
198         }
199
200         if (bits_left == 0 && current_val_bits_left == 9) {
201             if (data_idx < packet_len) {
202                 data[data_idx] = ((current_val >> 1) & 0xFF);
203                 bits_left = 0;
204                 data_idx += 1;
205                 current_val_bits_left = 1;
206             } else
207                 break;
208         }
209
210         if (bits_left == 0 && current_val_bits_left == 1) {
211             data[data_idx] = (current_val & 0x01) << 7;
212             bits_left = 7;
213             current_val_bits_left = 0;
214             continue;
215         }
216
217         if (bits_left == 7 && current_val_bits_left == CODE_LENGTH) {
218             data[data_idx] |= ((current_val >> 6) & 0x7F);
219             bits_left = 0;
220             current_val_bits_left = 6;
221             data_idx += 1;
222         }
223
224         if (bits_left == 0 && current_val_bits_left == 6) {
225             if (data_idx < packet_len) {
226                 data[data_idx] = ((current_val)&0x3F) << 2;
227                 bits_left = 2;
228                 current_val_bits_left = 0;
229                 continue;
230             } else
231                 break;
232         }
233
234         if (bits_left == 2 && current_val_bits_left == CODE_LENGTH) {
235             if (data_idx < packet_len) {
236                 data[data_idx] |= ((current_val >> 11) & 0x03);
237                 bits_left = 0;
238                 data_idx += 1;
239                 current_val_bits_left = 11;
240             } else
241                 break;
242         }
243
244         if (bits_left == 0 && current_val_bits_left == 11) {
245             if (data_idx < packet_len) {
246                 data[data_idx] = ((current_val >> 3) & 0xFF);
247                 bits_left = 0;
248                 data_idx += 1;
249                 current_val_bits_left = 3;
250             } else
251                 break;
252         }
```

```cpp
253
254            if (bits_left == 0 && current_val_bits_left == 3) {
255                if (data_idx < packet_len) {
256                    data[data_idx] = ((current_val)&0x07) << 5;
257                    bits_left = 5;
258                    current_val_bits_left = 0;
259                    continue;
260                } else
261                    break;
262            }
263
264            if (bits_left == 5 && current_val_bits_left == CODE_LENGTH) {
265                data[data_idx] |= ((current_val >> 8) & 0x1F);
266                bits_left = 0;
267                current_val_bits_left = 8;
268                data_idx += 1;
269            }
270
271            if (bits_left == 0 && current_val_bits_left == 8) {
272                if (data_idx < packet_len) {
273                    data[data_idx] = (current_val & 0xFF);
274                    bits_left = 0;
275                    current_val_bits_left = 0;
276                    data_idx += 1;
277                    continue;
278                } else
279                    break;
280            }
281        }
282        return data;
283 }
284
285 static void compression_pipeline(
286     RawData *r_data, unsigned char *host_input, uint32_t *output_codes,
287     uint32_t *chunk_indices, uint32_t *output_code_lengths, CLDevice dev
     ,
288     cl::Buffer lzw_input_buffer, cl::Buffer lzw_output_buffer,
289     cl::Buffer chunk_indices_buffer, cl::Buffer
     out_packet_lengths_buffer,
290     unsigned int chunk_size) {
291
292     vector<uint32_t> vect;
293     string sha_fingerprint;
294     int64_t chunk_idx = 0;
295     uint32_t packet_len = 0;
296     uint32_t header = 0;
297     vector<int64_t> dedup_out;
298     vector<pair<pair<int, int>, unsigned char *>> final_data;
299
300     //
     ----------------------------------------------------------------
301     // Step 3: Run the kernel
302     //
     ----------------------------------------------------------------
303
304     std::vector<cl::Event> write_event(1);
305     std::vector<cl::Event> compute_event(1);
306     std::vector<cl::Event> done_event(1);
```

```
307
308     // double total_time_2 = 0;
309
310     memcpy(host_input, r_data->pipeline_buffer,
311            sizeof(unsigned char) * r_data->length_sum);
312
313     total_time.start();
314
315     // RUN CDC
316     time_cdc.start();
317     fast_cdc(r_data->pipeline_buffer, r_data->length_sum, chunk_size,
       vect);
318     time_cdc.stop();
319
320     chunk_indices[0] = vect.size();
321
322     std::copy(vect.begin(), vect.end(), chunk_indices + 1);
323
324     for (int i = 0; i < (int)(vect.size() - 1); i++) {
325         // RUN SHA
326         time_sha.start();
327         sha_fingerprint =
328             sha_256(r_data->pipeline_buffer, vect[i], vect[i + 1]);
329         time_sha.stop();
330
331         // RUN DEDUP
332         time_dedup.start();
333         chunk_idx = dedup(sha_fingerprint);
334         dedup_out.push_back(chunk_idx);
335         time_dedup.stop();
336     }
337
338     // RUN LZW
339     time_lzw.start();
340
341     dev.kernel.setArg(0, lzw_input_buffer);
342     dev.kernel.setArg(1, lzw_output_buffer);
343     dev.kernel.setArg(2, chunk_indices_buffer);
344     dev.kernel.setArg(3, out_packet_lengths_buffer);
345
346     dev.queue.enqueueMigrateMemObjects({lzw_input_buffer,
       chunk_indices_buffer},
347                                        0 /* 0 means from host*/, NULL,
348                                        &write_event[0]);
349
350     dev.queue.enqueueTask(dev.kernel, &write_event, &compute_event[0]);
351
352     // Profiling the kernel.
353     /* compute_event[0].wait(); */
354     /* total_time_2 +=
355      * compute_event[0].getProfilingInfo<CL_PROFILING_COMMAND_END>() -
       */
356     /* compute_event[0].getProfilingInfo<CL_PROFILING_COMMAND_START>();
       */
357
358     dev.queue.enqueueMigrateMemObjects(
359         {lzw_output_buffer, out_packet_lengths_buffer},
360         CL_MIGRATE_MEM_OBJECT_HOST, &compute_event, &done_event[0]);
```

```
361    clWaitForEvents(1, (const cl_event *)&done_event[0]);
362    time_lzw.stop();
363
364    if (output_code_lengths[0] & 0x1) {
365        printf("FAILED TO INSERT INTO ASSOC MEM!!\n");
366        exit(EXIT_FAILURE);
367    }
368
369    uint32_t *output_codes_ptr = output_codes;
370
371    for (int i = 1; i < (int)chunk_indices[0]; i++) {
372        if (dedup_out[i - 1] == -1) {
373            packet_len = ((output_code_lengths[i] * 13) / 8);
374            packet_len = ((output_code_lengths[i] & 0x7) != 0) ?
    packet_len + 1
375                                                                :
    packet_len;
376
377            unsigned char *data_packet =
378                create_packet(chunk_idx, output_code_lengths[i],
379                              output_codes_ptr, packet_len);
380
381            header = packet_len << 1;
382            final_data.push_back({{header, packet_len}, data_packet});
383            lzw_bytes += 4;
384
385            lzw_bytes += packet_len;
386
387        } else {
388            header = (dedup_out[i - 1] << 1) | 1;
389            final_data.push_back({{header, -1}, NULL});
390            dedup_bytes += 4;
391        }
392
393        output_codes_ptr += output_code_lengths[i];
394    }
395    total_time.stop();
396
397    for (auto it : final_data) {
398        if (it.second != NULL) {
399            fwrite(&it.first.first, sizeof(uint32_t), 1, r_data->
    fptr_write);
400            fwrite(it.second, sizeof(unsigned char), it.first.second,
401                   r_data->fptr_write);
402            free(it.second);
403        } else {
404            fwrite(&it.first.first, sizeof(uint32_t), 1, r_data->
    fptr_write);
405        }
406    }
407
408    /* cout << "Total Kernel Execution Time using Profiling Info: " <<
409     * total_time_2 */
410    /*      << " ms." << endl; */
411 }
412
413 int main(int argc, char *argv[]) {
414
```

```
415    stopwatch ethernet_timer;
416    stopwatch compression_timer;
417    unsigned char *input[NUM_PACKETS];
418    int writer = 0;
419    int done = 0;
420    int length = -1;
421    uint64_t offset = 0;
422    ESE532_Server server;
423
424    int blocksize = BLOCKSIZE;
425    char *file = strdup("compressed_file.bin");
426    char *kernel_name = strdup("lzw.xclbin");
427    unsigned int chunk_size = CHUNK_SIZE;
428
429    // set blocksize if decalred through command line
430    handle_input(argc, argv, &blocksize, &file, &kernel_name, &
       chunk_size);
431
432    RawData *r_data = (RawData *)calloc(1, sizeof(RawData));
433    CHECK_MALLOC(r_data, "Unable to allocate memory for raw data");
434
435    r_data->fptr_write = fopen(file, "wb");
436    if (r_data->fptr_write == NULL) {
437        printf("Error creating file for compressed output!!\n");
438        exit(EXIT_FAILURE);
439    }
440
441    r_data->pipeline_buffer =
442        (unsigned char *)calloc(NUM_PACKETS * blocksize, sizeof(unsigned
       char));
443    CHECK_MALLOC(r_data->pipeline_buffer,
444                 "Unable to allocate memory for pipeline buffer");
445
446    for (int i = 0; i < (NUM_PACKETS); i++) {
447        input[i] = (unsigned char *)calloc((blocksize + HEADER),
448                                            sizeof(unsigned char));
449        CHECK_MALLOC(input, "Unable to allocate memory for input buffer"
       );
450    }
451
452    server.setup_server(blocksize);
453
454    writer = 0;
455
456    //
       ----------------------------------------------------------------------------
457    // Step 1: Initialize the OpenCL environment
458    //
       ----------------------------------------------------------------------------
459    cl_int err;
460    std::string binaryFile = kernel_name;
461    unsigned fileBufSize;
462    std::vector<cl::Device> devices = get_xilinx_devices();
463    devices.resize(1);
464    cl::Device device = devices[0];
465    cl::Context context(device, NULL, NULL, NULL, &err);
466    char *fileBuf = read_binary_file(binaryFile, fileBufSize);
467    cl::Program::Binaries bins{{fileBuf, fileBufSize}};
```

```cpp
468    cl::Program program(context, devices, bins, NULL, &err);
469    CLDevice dev;
470    dev.queue = cl::CommandQueue(context, device,
471                                 CL_QUEUE_OUT_OF_ORDER_EXEC_MODE_ENABLE,
       &err);
472    dev.kernel = cl::Kernel(program, "lzw", &err);
473
474    //
       ------------------------------------------------------------------------
475    // Step 2: Create buffers and initialize test values
476    //
       ------------------------------------------------------------------------
477
478    cl::Buffer lzw_input_buffer =
479        cl::Buffer(context, CL_MEM_READ_ONLY,
480                   sizeof(unsigned char) * NUM_PACKETS * blocksize, NULL
       , &err);
481    unsigned char *host_input = (unsigned char *)dev.queue.
       enqueueMapBuffer(
482        lzw_input_buffer, CL_TRUE, CL_MAP_WRITE, 0,
483        sizeof(unsigned char) * NUM_PACKETS * blocksize);
484
485    cl::Buffer lzw_output_buffer =
486        cl::Buffer(context, CL_MEM_WRITE_ONLY,
487                   sizeof(uint32_t) * MAX_OUTPUT_BUF_SIZE, NULL, &err);
488    uint32_t *output_codes = (uint32_t *)dev.queue.enqueueMapBuffer(
489        lzw_output_buffer, CL_TRUE, CL_MAP_READ, 0,
490        sizeof(uint32_t) * MAX_OUTPUT_BUF_SIZE);
491
492    cl::Buffer chunk_indices_buffer =
493        cl::Buffer(context, CL_MEM_READ_ONLY, sizeof(uint32_t) *
       MAX_LZW_CHUNKS,
494                   NULL, &err);
495    uint32_t *chunk_indices = (uint32_t *)dev.queue.enqueueMapBuffer(
496        chunk_indices_buffer, CL_TRUE, CL_MAP_WRITE, 0,
497        sizeof(uint32_t) * MAX_LZW_CHUNKS);
498
499    cl::Buffer out_packet_lengths_buffer =
500        cl::Buffer(context, CL_MEM_WRITE_ONLY,
501                   sizeof(uint32_t) * MAX_LZW_CHUNKS, NULL, &err);
502    uint32_t *output_code_lengths = (uint32_t *)dev.queue.
       enqueueMapBuffer(
503        out_packet_lengths_buffer, CL_TRUE, CL_MAP_READ, 0,
504        sizeof(uint32_t) * MAX_LZW_CHUNKS);
505
506    compression_timer.start();
507
508    // Loop until last message
509    while (!done) {
510
511        ethernet_timer.start();
512        server.get_packet(input[writer]);
513        ethernet_timer.stop();
514
515        // get packet
516        unsigned char *buffer = input[writer];
517
518        // decode
```

49

```
519        done = buffer[1] & DONE_BIT_L;
520        length = buffer[0] | (buffer[1] << 8);
521        length &= ~DONE_BIT_H;
522
523        offset += length;
524        r_data->length_sum += length;
525
526        // Perform the actual computation here. The idea is to maintain
     a
527        // buffer, that will hold multiple packets, so that CDC can
     chunklength)
528        // at appropriate boundaries. Call the compression pipeline
     function
529        // after the buffer is completely filled.
530        if (length != 0)
531            memcpy(r_data->pipeline_buffer + (writer * blocksize),
532                    input[writer] + 2, length);
533
534        if (writer == (NUM_PACKETS - 1) || (length < blocksize && length
      > 0) ||
535            done == 1) {
536            compression_pipeline(
537                r_data, host_input, output_codes, chunk_indices,
538                output_code_lengths, dev, lzw_input_buffer,
     lzw_output_buffer,
539                chunk_indices_buffer, out_packet_lengths_buffer,
     chunk_size);
540            writer = 0;
541            r_data->length_sum = 0;
542        } else
543            writer += 1;
544    }
545
546    for (int i = 0; i < (NUM_PACKETS); i++)
547        free(input[i]);
548
549    fclose(r_data->fptr_write);
550    dev.queue.enqueueUnmapMemObject(lzw_input_buffer, host_input);
551    dev.queue.enqueueUnmapMemObject(lzw_output_buffer, output_codes);
552    dev.queue.enqueueUnmapMemObject(chunk_indices_buffer, chunk_indices)
     ;
553    dev.queue.enqueueUnmapMemObject(out_packet_lengths_buffer,
554                                    output_code_lengths);
555    dev.queue.finish();
556
557    free(r_data->pipeline_buffer);
558    free(r_data);
559
560    compression_timer.stop();
561
562    // Print Latencies
563    cout << "--------------- Total Latencies ---------------" << endl;
564    cout << "Total latency of CDC is: " << time_cdc.latency() << " ms."
     << endl;
565    cout << "Total latency of LZW is: " << time_lzw.latency() << " ms."
     << endl;
566    cout << "Total latency of SHA256 is: " << time_sha.latency() << " ms
     ."
```

```cpp
567         << endl;
568     cout << "Total latency of DeDup is: " << time_dedup.latency() << "
    ms."
569         << endl;
570     cout << "Total time taken: " << total_time.latency() << " ms." <<
    endl;
571     cout << "--------------- Average Latencies -----------" << endl;
572     cout << "Average latency of CDC per loop iteration is: "
573         << time_cdc.avg_latency() << " ms." << endl;
574     cout << "Average latency of LZW per loop iteration is: "
575         << time_lzw.avg_latency() << " ms." << endl;
576     cout << "Average latency of SHA256 per loop iteration is: "
577         << time_sha.avg_latency() << " ms." << endl;
578     cout << "Average latency of DeDup per loop iteration is: "
579         << time_dedup.avg_latency() << " ms." << endl;
580     cout << "Average latency: " << total_time.avg_latency() << " ms."
581         << std::endl;
582
583     std::cout << "\n\n";
584
585     std::cout << "--------------- Key Throughputs ---------------" <<
    std::endl;
586     float ethernet_latency = ethernet_timer.latency() / 1000.0;
587     float compression_latency = compression_timer.latency() / 1000.0;
588     float compression_latency_total_time = total_time.latency() /
    1000.0;
589
590     float cdc_latency_total_time = time_cdc.latency() / 1000.0;
591     float sha_latency_total_time = time_sha.latency() / 1000.0;
592     float lzw_latency_total_time = time_lzw.latency() / 1000.0;
593     float dedup_latency_total_time = time_dedup.latency() / 1000.0;
594
595     float compression_throughput =
596         (offset * 8 / 1000000.0) / compression_latency; // Mb/s
597     float compression_throughput_2 =
598         (offset * 8 / 1000000.0) / compression_latency_total_time; // Mb
    /s
599
600     float cdc_throughput =
601         (offset * 8 / 1000000.0) / cdc_latency_total_time; // Mb/s
602     float sha_throughput =
603         (offset * 8 / 1000000.0) / sha_latency_total_time; // Mb/s
604     float lzw_throughput =
605         (offset * 8 / 1000000.0) / lzw_latency_total_time; // Mb/s
606     float dedup_throughput =
607         (offset * 8 / 1000000.0) / dedup_latency_total_time; // Mb/s
608
609     float ethernet_throughput =
610         (offset * 8 / 1000000.0) / ethernet_latency; // Mb/s
611
612     cout << "Ethernet Latency: " << ethernet_latency << "s." << endl;
613     cout << "Bytes Received: " << offset << "B." << endl;
614     cout << "Latency for Compression: " << compression_latency << "s."
    << endl;
615     cout << "Latency for Compression (without fwrite): "
616         << compression_latency_total_time << "s." << endl;
617
618     std::cout << "\n";
```

```
619
620     cout << "Latency for CDC: " << cdc_latency_total_time << "s." <<
    endl;
621     cout << "Latency for SHA: " << sha_latency_total_time << "s." <<
    endl;
622     cout << "Latency for LZW: " << lzw_latency_total_time << "s." <<
    endl;
623     cout << "Latency for DEDUP: " << dedup_latency_total_time << "s." <<
    endl;
624
625     std::cout << "\n";
626
627     cout << "CDC Throughput: " << cdc_throughput << "Mb/s." << endl;
628     cout << "SHA Throughput: " << sha_throughput << "Mb/s." << endl;
629     cout << "LZW Throughput: " << lzw_throughput << "Mb/s." << endl;
630     cout << "DEDUP Throughput: " << dedup_throughput << "Mb/s." << endl;
631
632     std::cout << "\n";
633
634     cout << "Ethernet Throughput: " << ethernet_throughput << "Mb/s." <<
    endl;
635     cout << "Application Throughput: " << compression_throughput << "Mb/
    s."
636         << endl;
637     cout << "Application Throughput (without fwrite): "
638         << compression_throughput_2 << "Mb/s." << endl;
639     cout << "Bytes Contributed by Deduplication: " << dedup_bytes << "B.
    "
640         << endl;
641     cout << "Bytes Contributed by LZW: " << lzw_bytes << "B." << endl;
642
643     std::cout << "\n";
644
645     return 0;
646 }
```

### 8.3.6  Testbench

```
1  #include "../Encoder/common.h"
2  #include <iostream>
3  #include <stdlib.h>
4  #include <string.h>
5  #include <unordered_map>
6  #include <vector>
7
8  #define FILE_SIZE 16383
9  #define MAX_LZW_CODES (4096 * 20)
10
11 using namespace std;
12
13 // "Golden" functions to check correctness
14 std::vector<int> encoding(std::string s1) {
15     // std::cout << "Encoding\n";
16     std::unordered_map<std::string, int> table;
17     for (int i = 0; i <= 255; i++) {
18         std::string ch = "";
19         ch += char(i);
```

```cpp
20          table[ch] = i;
21      }
22
23      std::string p = "", c = "";
24      p += s1[0];
25      int code = 256;
26      std::vector<int> output_code;
27      // std::cout << "String\tOutput_Code\tAddition\n";
28      for (int i = 0; i < s1.length(); i++) {
29          if (i != s1.length() - 1)
30              c += s1[i + 1];
31          if (table.find(p + c) != table.end()) {
32              p = p + c;
33          } else {
34              //              std::cout << p << "\t" << table[p] << "\t\t"
35              //                        << p + c << "\t" << code << std::endl;
36              output_code.push_back(table[p]);
37              table[p + c] = code;
38              code++;
39              p = c;
40          }
41          c = "";
42      }
43      // std::cout << p << "\t" << table[p] << std::endl;
44      output_code.push_back(table[p]);
45      return output_code;
46  }
47
48  //**************************************************************************
49  int main() {
50      // FILE *fptr =
51      // fopen("/home1/r/ruturajn/ESE532/ese532_project/project/Text_Files
52      /Franklin.txt",
52      // "r"); FILE *fptr =
53      // fopen("/home1/r/ruturajn/ESE532/ese532_project/project/Text_Files
      /lotr.txt",
54      // "r"); FILE *fptr =
55      // fopen("/home1/r/ruturajn/ESE532/ese532_project/project/Text_Files
      /imaggeee.jpg",
56      // "r"); FILE *fptr =
57      // fopen("/home1/r/ruturajn/ESE532/ese532_project/project/Text_Files
      /IMAGE_390.jpg",
58      // "r"); FILE *fptr = fopen("./ruturajn.tgz", "r");
59      FILE *fptr = fopen("/home1/r/ruturajn/ESE532/ese532_project/project/
      "
60                          "Text_Files/LittlePrince.txt",
61                          "r");
62      // FILE *fptr =
63      // fopen("/home1/r/ruturajn/ESE532/ese532_project/project/Text_Files
      /small_prince.txt",
64      // "r"); FILE *fptr = fopen("/home1/r/ruturajn/Downloads/embedded_h5
      .JPG",
65      // "rb"); FILE *fptr =
66      // fopen("/home1/r/ruturajn/Downloads/ESE5070_Assignment3_ruturajn
      -1.pdf",
67      // "rb");
68      // FILE *fptr = fopen("/home1/r/ruturajn/Downloads/FiraCode.zip", "
      rb");
```

```cpp
      // FILE *fptr =
      // fopen("/home1/r/ruturajn/ESE532/ese532_project/project/encoder.xo
    ", "r");
      if (fptr == NULL) {
          printf("Unable to open file!\n");
          exit(EXIT_FAILURE);
      }

      fseek(fptr, 0, SEEK_END);      // seek to end of file
      int64_t file_sz = ftell(fptr); // get current file pointer
      fseek(fptr, 0, SEEK_SET);      // seek back to beginning of file

      unsigned char file_data[16384];
      uint32_t chunk_indices[20];
      uint32_t *lzw_codes = (uint32_t *)calloc(40960, sizeof(uint32_t));
      memset(lzw_codes, 23, 40960 * sizeof(uint32_t));
      if (lzw_codes == NULL) {
          cout << "Unable to allocate memory for lzw codes!" << endl;
          exit(EXIT_FAILURE);
      }
      uint32_t out_packet_lengths[20];
      unsigned int count = 0;
      bool fail_stat = false;

      while (file_sz > 0) {

          size_t bytes_read = fread(file_data, 1, 16384, fptr);

          if (file_sz >= 16384 && bytes_read != 16384)
              printf("Unable to read file contents");

          vector<uint32_t> vect;

          if (file_sz < 16384)
              fast_cdc(file_data, (unsigned int)file_sz, 4096, vect);
          else
              fast_cdc(file_data, (unsigned int)16384, 4096, vect);

          chunk_indices[0] = vect.size();

          std::copy(vect.begin(), vect.end(), chunk_indices + 1);

          lzw(file_data, lzw_codes, chunk_indices, out_packet_lengths);

          uint32_t *lzw_codes_ptr = lzw_codes;

          if (out_packet_lengths[0]) {
              cout << "TEST FAILED!!" << endl;
              cout << "FAILED TO INSERT INTO ASSOC MEM!!\n";
              exit(EXIT_FAILURE);
          }

          uint32_t packet_len = 0;
          uint32_t header = 0;
          vector<pair<pair<int, int>, unsigned char *>> final_data;

          for (int i = 1; i <= chunk_indices[0] - 1; i++) {
              std::string s;
```

54

```
126                 char *temp = (char *)file_data + chunk_indices[i];
127                 int count = chunk_indices[i];
128
129                 while (count++ < chunk_indices[i + 1]) {
130                     s += *temp;
131                     temp += 1;
132                 }
133
134                 std::vector<int> output_code = encoding(s);
135
136                 if (out_packet_lengths[i] != output_code.size()) {
137                     cout << "TEST FAILED!!" << endl;
138                     cout << "FAILURE MISMATCHED PACKET LENGTH!!" << endl;
139                     cout << out_packet_lengths[i] << "|" << output_code.size
    ()
140                         << "at i = " << i << endl;
141                     fail_stat = true;
142                     exit(EXIT_FAILURE);
143                 }
144
145                 for (int j = 0; j < output_code.size(); j++) {
146                     if (output_code[j] != lzw_codes_ptr[j]) {
147                         if (!fail_stat)
148                             fail_stat = true;
149                         cout << "FAILURE!!" << endl;
150                         cout << output_code[j] << "|" << lzw_codes_ptr[j]
151                             << " at j = " << j << " and i = " << i << endl;
152                     }
153                 }
154
155                 lzw_codes_ptr += out_packet_lengths[i];
156             }
157
158         file_sz -= 16384;
159
160         cout << "Iteration : " << count << endl;
161         count++;
162     }
163
164     fclose(fptr);
165
166     if (fail_stat)
167         cout << "TEST FAILED!!" << endl;
168     else
169         cout << "TEST PASSED!!" << endl;
170
171     return 0;
172 }
```