

Adv. Computer Systems - Project 4

Running Instructions	2
Grading Criteria	2
Software Implementation	2
Encoding Functionality	2
Query Operations	2
Multi-Threading Implementation	2
SIMD Utilization	2
Performance and Analysis	2
Vanilla Column Scan	2
Dictionary Encoding	2
Query	2
Documentation	2
Readme Clarity	2
Experimental Setup and Analysis	2
Conclusion & Analysis	3

Running Instructions:

This assignment is split into 4 program files. DictionaryCodec.h and DictionaryCodec.cpp handle implementing the dictionary encoding and query logic, as well as setting up the dictionary encoding scheme. The encoding utilizes the variable num_threads, to define the number of threads used in encoding. The query operations include VanillaScan, Item Search (SIMD/no SIMD), and Prefix Scan (SIMD/No SIMD). Item Search(SIMD/No SIMD) and Prefix Scan (SIMD/No SIMD). The last two files are testbenches that handle encoding and query tests. test_encoding.cpp handles encoding benchmarks with a varying num_thread value (Testing 1,2,4 and 8 thread encoding), and test_query.cpp measures query speeds for the five query operations.

To run this assignment, download the 4 program files, as well as the given raw column data file (Column.txt, provided by the professor at

https://drive.google.com/file/d/195XTg8HWDILc1JlsGX6_jJ5PUhi9KDG/view?usp=drive_link).

To compile the program, use the following two commands:

```
g++ -std=c++17 -O2 -msse4.2 -mavx2 -march=native -o test_encoding test_encoding.cpp DictionaryCodec.cpp
g++ -std=c++17 -O2 -msse4.2 -mavx2 -march=native -o test_query test_query.cpp DictionaryCodec.cpp
```

The program testbenches can now be run with ./test_encoding and ./test_query respectively.

To modify the query string, adjust line 61 in test_query.cpp. The default implementation is [**std::string query = "spin";**], meaning the search query is the string "spin".

Grading Criteria:

Software Implementation

Encoding Functionality

Dictionary encoding is handled in DictionaryCodec.cpp and manages the encoding of the given Column.txt file. The encoded file is outputted as EncodedColumn.txt (as handled in test_query.cpp) and displays the dictionary and encoded data column.

```

17 void DictionaryCodec::encodeColumn(const std::string inputfile, const std::string outputfile, int numThreads) {
18     std::ifstream input(inputfile);
19     if (!input.is_open()) {
20         std::cerr << "Error: could not open input file: " << inputfile << std::endl;
21         return;
22     }
23     std::vector<string> lines;
24     std::string line;
25     while (std::getline(input, line)) {
26         lines.push_back(line);
27     }
28     input.close();
29     std::cout << "Loaded " << lines.size() << " lines from input file: " << inputfile << std::endl;
30     // initialize local dictionaries for each thread
31     std::vector<std::vector<string>> localDicts(numThreads);
32     std::vector<int> localDictSizes(numThreads);
33     try {
34         for (int t = 0; t < numThreads; ++t) {
35             std::vector<string> dict;
36             size_t start = 0;
37             size_t end = (lines.size() - 1) / numThreads;
38             if (t == numThreads - 1) end = lines.size();
39             std::cout << "Thread " << t << " processing lines " << start << " to " << end << std::endl;
40             for (size_t i = start; i < end; ++i) {
41                 localDicts[t].push_back(lines[i]);
42             }
43             localDictSizes[t] = localDicts[t].size();
44         }
45     } catch (std::exception & e) {
46         std::cerr << "Exception during encoding: " << e.what() << std::endl;
47         return;
48     }
49     // merge dictionaries
50     try {
51         std::vector<string> mergedDict;
52         for (int t = 0; t < numThreads; ++t) {
53             mergedDict.insert(mergedDict.end(), localDicts[t].begin(), localDicts[t].end());
54         }
55         std::cout << "Merging complete. Writing to file: " << outputfile << std::endl;
56         // write dictionary and encoded column
57         std::ofstream output(outputfile);
58         if (!output.is_open()) {
59             std::cerr << "Error: could not open output file: " << outputfile << std::endl;
60             return;
61         }
62         for (const auto & pair : dictionary) {
63             output << pair.first << " " << pair.second << "\n";
64         }
65         for (const auto & line : lines) {
66             output << dictionary[line] << "\n";
67         }
68         output.close();
69         std::cout << "Encoded file written successfully." << std::endl;
70     }
71 }
72
73 void DictionaryCodec::buildLocalDictionaries(const std::vector<string> lines, const std::vector<int> localDictSizes, int numThreads) {
74     for (size_t i = 0; i < lines.size(); ++i) {
75         size_t thread = i % numThreads;
76         localDicts[thread].push_back(lines[i]);
77     }
78     localDictSizes[thread] = localDicts[thread].size();
79 }
80
81 void DictionaryCodec::mergeLocalDictionaries(const std::vector<std::vector<string>> localDicts, const std::vector<int> localDictSizes) {
82     for (const auto & dict : localDicts) {
83         for (const auto & word : dict) {
84             localDict[word] = localDict[word] + 1;
85         }
86     }
87     dictionary = localDict;
88 }
89
90 int main() {
91     DictionaryCodec codec;
92     std::string inputfile = "column.txt"; // input raw column file
93     std::string outputfile = "transformed.txt"; // encoded output file
94     std::cout << "Measuring encoding performance..." << std::endl;
95     measureEncodingPerformance(codec, inputfile, outputfile);
96     return 0;
97 }

```

Figure 1: DictionaryCodec.cpp encoding

```

1 void measureEncodingPerformance(const std::string inputfile, const std::string outputfile, const std::string outputfile) {
2     for (int numThreads : {1, 2, 4, 8}) {
3         auto start = std::chrono::high_resolution_clock::now();
4
5         // use the correct method name
6         codec.encodeColumn(inputfile, outputfile, numThreads);
7
8         auto end = std::chrono::high_resolution_clock::now();
9         auto duration = std::chrono::duration_cast<std::chrono::milliseconds>(end - start).count();
10
11         std::cout << "Encoding with " << numThreads << " threads took: " << duration << " ms" << std::endl;
12     }
13 }
14
15 int main() {
16     DictionaryCodec codec;
17     std::string inputfile = "column.txt"; // input raw column file
18     std::string outputfile = "transformed.txt"; // encoded output file
19     std::cout << "Measuring encoding performance..." << std::endl;
20     measureEncodingPerformance(codec, inputfile, outputfile);
21     return 0;
22 }

```

Figure 2: test_encoding.cpp handling EncodedColumn output file

Query Operations

There are five essential query operations, vanillaScan, queryItem (non SIMD single scan), queryPrefix (non SIMD prefix scan), simdQueryItem (SIMD single scan), simdQueryPrefix (SIMD prefix scan). The test_query.cpp testbench file handles measuring the timing for all of these operations, given a string query.

```

110 std::vector<int> DictionaryCodec::queryItem(const std::string encodedFile, const std::string item) {
111     std::ifstream infile(encodedFile);
112     std::string line;
113     std::vector<int> indices;
114
115     size_t index = 0;
116     while (std::getline(infile, line)) {
117         if (line == item) {
118             indices.push_back(index);
119         }
120         index++;
121     }
122     infile.close();
123     return indices;
124 }
125
126 std::vector<int> DictionaryCodec::findQueryItem(const std::vector<std::string> column, const std::string item) {
127     std::vector<int> indices;
128     size_t itemLength = item.length();
129
130     // Loop through the column data in chunks
131     for (size_t i = 0; i < column.size(); i += 8) {
132         // Process multiple strings (8 strings at a time)
133         for (size_t j = 0; j < 8 * 8; j += 8) {
134             if (column[i + j] != item) continue;
135             std::string currentString = column[i + j];
136             indices.push_back(i + j);
137         }
138     }
139     return indices;
140 }
141
142 std::vector<std::pair<std::string, std::vector<int>>> DictionaryCodec::queryPrefix(const std::string encodedFile, const std::string prefix) {
143     std::ifstream infile(encodedFile);
144     std::string line;
145     std::vector<std::pair<std::string, std::vector<int>>> result;
146
147     size_t index = 0;
148     while (std::getline(infile, line)) {
149         if (line.find(prefix) == 0) {
150             result.push_back({line, indices});
151         }
152         index++;
153     }
154     infile.close();
155     return result;
156 }
157
158 std::vector<int> DictionaryCodec::findAllItems(const std::string readFile, const std::string query) {
159     std::ifstream infile(readFile);
160     std::vector<int> indices;
161     std::string line;
162
163     size_t index = 0;
164     while (std::getline(infile, line)) {
165         if (line == query) {
166             indices.push_back(index);
167         }
168         index++;
169     }
170     infile.close();
171     return indices;
172 }
173
174 std::vector<std::pair<std::string, std::vector<int>>> DictionaryCodec::findQueryPrefix(const std::vector<std::string> column, const std::string prefix) {
175     std::vector<std::pair<std::string, std::vector<int>>> result;
176     size_t prefixLength = prefix.length();
177
178     // Loop through the column data in chunks
179     for (size_t i = 0; i < column.size(); i += 8) {
180         // Process multiple strings (8 strings at a time)
181         for (size_t j = 0; j < 8 * 8; j += 8) {
182             const std::string currentString = column[i + j];
183
184             // Check if the string has at least the prefix length
185             if (currentString.size() < prefixLength) continue;
186             std::string currentStringSubstr = currentString.substr(0, prefixLength);
187             if (currentStringSubstr == prefix) {
188                 result.push_back({currentString, {i + j}});
189             }
190         }
191     }
192     return result;
193 }

```

Figure 3: DictionaryCodec.cpp query

Multi-Threading Implementation

Multithreading is used within the Column File encoding to speed up the process. The number of lines to encode is equally divided amongst the number of threads specified. In this case, encoding was measured with 1,2,4, and 8 threads.

```

128 void measureEncodingPerformance(DictionaryCodec& codec, const std::string& inputFile, const std::string& outputFile) {
129     for (int numThreads : {1, 2, 4, 8}) {
130         auto start = std::chrono::high_resolution_clock::now();
131
132         // Use the correct method name
133         codec.encodeColumnFile(inputFile, outputFile, numThreads);
134
135         auto end = std::chrono::high_resolution_clock::now();
136         auto duration = std::chrono::duration_cast<std::chrono::milliseconds>(end - start).count();
137
138         std::cout << "Encoding with " << numThreads << " threads took: " << duration << " ms" << std::endl;
139     }
140 }

```

Figure 4: Multithreading test in test_encoding.cpp

SIMD Utilization

As mentioned in the previous section, SIMD accelerates single-query and prefix-query searches. This program includes `<immintrin.h>` to handle SIMD operations.

```
std::vector<size_t> DictionaryCodec::simdQueryItem(const std::vector<std::string>& column, const std::string& item) {
    std::vector<size_t> indices;
    size_t itemLength = item.length();

    // Loop through the column data in chunks
    for (size_t i = 0; i < column.size(); i += 8) {
        // Process multiple strings (8 strings at a time)
        for (size_t j = 0; j < 8 && i + j < column.size(); ++j) {
            if (column[i + j].size() == itemLength &&
                std::memcmp(column[i + j].c_str(), item.c_str(), itemLength) == 0) {
                indices.push_back(i + j);
            }
        }
    }

    return indices;
}
```

Figure 5: simdQueryItem

```
209 std::vector<std::pair<std::string, std::vector<size_t>>> DictionaryCodec::simdQueryPrefix(const std::vector<std::string>& column, const std::string& prefix) {
210     std::vector<std::pair<std::string, std::vector<size_t>>> result;
211     size_t prefixLength = prefix.length();
212
213     // Loop through the column data in chunks
214     for (size_t i = 0; i < column.size(); i += 8) {
215         // Process multiple strings (8 strings at a time)
216         for (size_t j = 0; j < 8 && i + j < column.size(); ++j) {
217             const std::string& currentString = column[i + j];
218
219             // Check if the string has at least the prefix length
220             if (currentString.size() >= prefixLength &&
221                 std::memcmp(currentString.c_str(), prefix.c_str(), prefixLength) == 0) {
222                 // Prefix matches, add to result
223                 result.push_back({currentString, {i + j}});
224             }
225         }
226     }
227
228     return result;
229 }
230
```

Figure 6: simdQueryPrefix

Performance and Analysis:

Vanilla Column Scan

The program uses a vanilla column scan (aka scanning without using the dictionary encoding scheme) as a baseline to compare the effect of encoded queries.

For the query string “spin”, the vanilla scan took 3510 milliseconds. This value is used in the query section later in the report.

Dictionary Encoding

Dictionary encoding was handled across four different thread counts- those being 1, 2, 4 and 8. A table and graph showing the thread count speed difference is shown below.

Table 1: Encoding Time (ms) vs. Thread Count

Thread Count	Encoding Time (ms)
1	94226
2	80589
4	75451
8	76140

Encoding Time (ms) vs. Thread Count

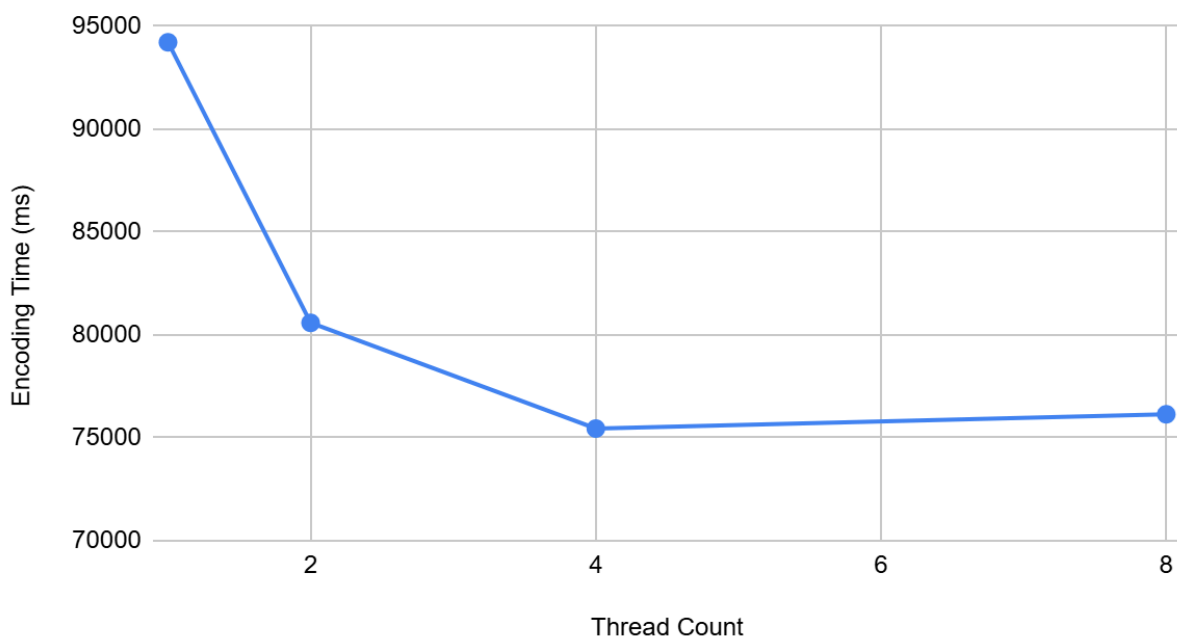


Figure 7: Encoding Time (ms) vs. Thread Count

The graph shows that as thread count increases, the time to encode decreases. This functionality works as expected. It is worth noting that there seems to be a bottleneck or diminishing return after 4 threads. 8 threads do not provide any additional performance benefit and seem to be the same as 4 threads (within the margin of error).

Query

For the query test, the query string “spin” was used to run 5 different operations, vanilla scan, item search (No SIMD), item search (SIMD), query scan (No SIMD), query scan (SIMD). A table and graph showing the time to complete these operations are included below.

Table 2: Time to complete Query Operations

Operation	Time to Complete (ms)
Vanilla Scan	3510
Item Search	2650
Item Search (SIMD)	289
Prefix Scan	3531
Prefix Scan (SIMD)	489

Time to Complete (ms) vs. Operation

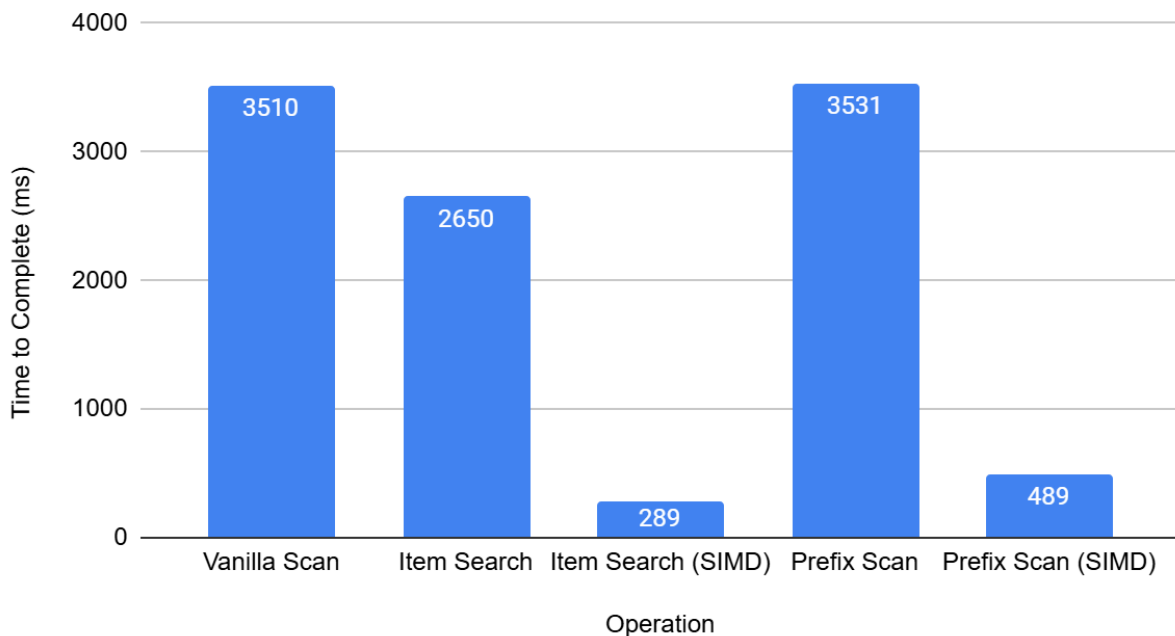


Figure 8: Time to complete Query Operations

As expected, dictionary encoding provided far faster results than the vanilla scan. Both item searches (SIMD and No-SIMD), provided much faster execution times than the non-encoded vanilla scan did, especially the SIMD implementation. Looking at the prefix scan, we see an expected rise in execution time, due to the new task of searching for a prefix through the entire data file (rather than stopping at the first find). SIMD once again provided a much faster execution time than the No-SIMD version.

Documentation:

Readme Clarity

Readme.txt provided in the repository.

Experimental Setup and Analysis

Experimental setup and analysis are provided in the “Performance and Analysis” portion of this report.

Conclusion & Analysis

The dictionary codec is based on compressing repeated data (encoding) and using indexing to make querying more efficient. In the context of large datasets, such as the one tested here, many values tend to repeat (e.g., string values of common phrases). A dictionary codec exploits this redundancy by doing two things. One, storing unique values in a dictionary, and two, encoding the original dataset as a list of indices pointing to the corresponding unique value in the dictionary.

Benefits of dictionary codec include sparse/query efficiency, scalability to larger datasets, faster prefix and range queries, and parallelization by utilizing multiple threads. It is clear in this assignment that multithreaded encoding can improve performance but can diminish returns after a certain number of threads. Dictionary encoding (combined with SIMD optimizations) can greatly reduce execution times compared to a vanilla/raw scan.

In summary, the philosophy of dictionary codecs focuses on reducing data redundancy, improving memory efficiency, and speeding up data access, making it an excellent choice for large-scale data storage and retrieval, especially when combined with parallel and SIMD processing for further acceleration.