

Huffman Coding in Kotlin: A Complete Implementation

Olin DSA-style Project Report (Compressor, Visualizer, and CLI)

Kotlin Implementation

HuffmanCompressor + HuffmanVisualizer (single-file)

October 26, 2025

Abstract

This report documents a complete, single-file implementation of *Huffman coding* in Kotlin. The project provides two primary classes: a `HuffmanCompressor` for the core logic and a `HuffmanVisualizer` for generating ASCII-art trees and frequency analyses. We present the API, the core algorithms (frequency analysis, heap-based tree construction, code generation, encoding, and decoding), and a robust I/O strategy using a text-based header for binary data. Following Olin DSA conventions, we provide correctness arguments (prefix-free property) and a full asymptotic runtime analysis. The report concludes with reflections on Kotlin-specific implementation choices and a discussion of the interactive CLI.

Contents

1	Overview and Goals	1
2	Interface and Code Architecture	2
3	Algorithms & Data Structures	3
3.1	Frequency Analysis	3
3.2	Huffman Tree Construction	3
3.3	Code Generation	4
3.4	Encoding	4
3.5	Decoding	5
4	Correctness (DSA-style)	6
5	Asymptotic Runtime and Space	6
6	Interactive Interface & Visualization	6
7	Implementation Reflection (Kotlin)	7
8	Acknowledgments	7

1 Overview and Goals

Huffman coding is a fundamental lossless data compression algorithm. It constructs an optimal, prefix-free set of binary codes based on the frequencies of input symbols, minimizing the expected code length.

The goals for this project were:

- Implement a **full round-trip** Huffman encoder and decoder in Kotlin.

- Provide a **clean API** separating compression logic from visualization.
- Build an **interactive CLI visualizer** to inspect frequencies, codes, and the tree structure.
- Ensure **correctness** for edge cases (empty files, single-symbol files).
- Analyze the **asymptotic runtime** and memory usage.

2 Interface and Code Architecture

The project is built around three core classes, all contained in a single file for portability.

Listing 1: Public API surface

```

1 // Basic node structure for the huffman tree
2 data class Node(
3     val symbol: Char? = null,
4     val frequency: Int,
5     val left: Node? = null,
6     val right: Node? = null
7 ) : Comparable<Node> {
8     override fun compareTo(other: Node): Int
9     fun isLeaf(): Boolean
10 }
11
12 class HuffmanCompressor {
13     // count how many times each character appears
14     fun analyzeFrequency(text: String): Map<Char, Int>
15
16     // build the huffman tree using a min-heap
17     fun buildHuffmanTree(frequencies: Map<Char, Int>): Node
18
19     // traverse tree to generate binary codes
20     fun generateCodes(root: Node): Map<Char, String>
21
22     // encode text using the code table and write to file
23     fun encode(text: String, codeTable: Map<Char, String>, outputFile: File)
24
25     // decode a compressed file back to original text
26     fun decode(inputFile: File): String
27
28     // full compression pipeline
29     fun compress(inputFile: File, outputFile: File)
30
31     // full decompression pipeline
32     fun decompress(inputFile: File): String
33 }
34
35 // visualization tool for understanding the huffman tree
36 class HuffmanVisualizer(private val compressor: HuffmanCompressor) {
37     fun visualizeFromText(text: String)
38     private fun printTree(node: Node?, prefix: String, isTail: Boolean)
39     fun interactiveMode()
40 }
41
42 // main interactive menu
43 fun interactiveMain()
44 fun main()

```

Key design choices.

- **ADT:** The `Node` is a simple `data class`. Leaf nodes are identified by a non-null `symbol` and null children (via `isLeaf()`).
- **Heap:** We use `java.util.PriorityQueue<Node>` for the min-heap, relying on the `Node`'s `Comparable` implementation.
- **Header Format:** To decompress, the decoder needs the code tree. We store this as a human-readable text header:
 1. Line 1: Number of symbols k .
 2. Lines 2 to $k + 1$: The code table, one entry per line, as `[char.code]:[bitstring]`. Using `char.code` (its integer value) ensures newlines and control characters are serialized safely.
 3. Line $k + 2$: A sentinel string "END_HEADER".
 4. Line $k + 3$: The exact `bitLength` of the data, used to trim padding from the final byte.
- **Payload:** The header is followed immediately by the raw binary payload (the packed bitstream).

3 Algorithms & Data Structures

3.1 Frequency Analysis

Goal: Build a histogram of character frequencies.

Listing 2: Frequency analysis

```
1 fun analyzeFrequency(text: String): Map<Char, Int> {
2     val freqMap = mutableMapOf<Char, Int>()
3     for (char in text) {
4         freqMap[char] = freqMap.getOrDefault(char, 0) + 1
5     }
6     // ... sanity check ...
7     require(freqMap.values.sum() == text.length) { /* ... */ }
8     return freqMap
9 }
```

Runtime: $\Theta(n)$ where n is the number of characters in the text.

Space: $O(k)$ where k is the number of distinct symbols (the alphabet size).

3.2 Huffman Tree Construction

Goal: Use a min-priority queue to implement the standard greedy Huffman algorithm.

Listing 3: Greedy tree construction with a min-heap

```
1 fun buildHuffmanTree(frequencies: Map<Char, Int>): Node {
2     require(frequencies.isNotEmpty()) { "can't build tree from nothing" }
3
4     val pq = PriorityQueue<Node>()
5     // start with all symbols as leaf nodes
6     frequencies.forEach { (char, freq) ->
7         pq.offer(Node(symbol = char, frequency = freq))
8     }
9
10    // merge until we have a single root
11    while (pq.size > 1) {
12        val left = pq.poll()
```

```

13     val right = pq.poll()
14     val merged = Node(
15         frequency = left.frequency + right.frequency,
16         left = left,
17         right = right
18     )
19     pq.offer(merged)
20 }
21 return pq.poll()
22 }

```

Runtime: The loop runs $k - 1$ times. Each iteration involves two `poll` and one `offer`, each costing $O(\log k)$. The initial heap build is $O(k \log k)$ (or $O(k)$ with ‘heapify’, but k pushes is equivalent here). Total time: $O(k \log k)$.

Space: $O(k)$ for the priority queue and the $2k - 1$ nodes in the final tree.

3.3 Code Generation

Goal: Perform a DFS traversal on the tree to build the bitstring codes.

Listing 4: DFS code generation

```

1 fun generateCodes(root: Node): Map<Char, String> {
2     val codeTable = mutableMapOf<Char, String>()
3
4     // recursive DFS to build codes
5     fun traverse(node: Node, code: String) {
6         if (node.isLeaf()) {
7             // edge case: single symbol gets code "0"
8             codeTable[node.symbol!!] = if (code.isEmpty()) "0" else code
9         } else {
10            node.left?.let { traverse(it, code + "0") }
11            node.right?.let { traverse(it, code + "1") }
12        }
13    }
14
15    traverse(root, "")
16    // ... prefix-free verification ...
17    return codeTable
18 }

```

Runtime: The traversal visits every node and edge in the tree exactly once. The tree has k leaves and $k - 1$ internal nodes, so $2k - 1$ nodes in total. Runtime is $\Theta(k)$.

Space: $O(k)$ to store the code table. The recursion depth is at most k (in a skewed tree).

3.4 Encoding

Goal: Write the text header, then map the input text to a bitstream and pack it into bytes.

Listing 5: Bitstring building and byte packing

```

1 fun encode(text: String, codeTable: Map<Char, String>, outputFile: File) {
2     // convert text to bit string
3     val bitString = StringBuilder()
4     for (char in text) {
5         bitString.append(codeTable[char] ?: error("no code for '$char'"))
6     }
7

```

```

8      // build header with code table... (omitted for brevity)
9      val header = /* ... build header ... */
10
11     // pack bits into bytes (pad last byte with 0s if needed)
12     val bytes = mutableListOf<Byte>()
13     for (i in bitString.indices step 8) {
14         val chunk = bitString.substring(i, minOf(i + 8, bitString.length))
15             .padEnd(8, '0')
16         bytes.add(chunk.toInt(2).toByte())
17     }
18
19     // write header as text, then binary payload
20     outputFile.writeText(header.toString(), Charsets.UTF_8)
21     outputFile.appendBytes(bytes.toByteArray())
22 }

```

Runtime: Let B be the total number of bits in the encoded message. Building the ‘bitString’ takes $\Theta(B)$ time (assuming n lookups in $O(n)$ total). Packing into bytes takes $\Theta(B/8)$ time. Total: $\Theta(n + B)$.

Space: $\Theta(B)$ for the intermediate ‘bitString’. This is a memory bottleneck for large files and could be replaced by a streaming bit-writer.

3.5 Decoding

Goal: Read the file, safely parse the text header, unpack the binary payload, and decode.

Listing 6: Binary-safe header parse + greedy decode

```

1 fun decode(inputFile: File): String {
2     val bytes = inputFile.readBytes()
3     var i = 0
4
5     // helper to read lines from byte array (avoids text/binary mixing issues)
6     fun readLine(): String {
7         val start = i
8         while (i < bytes.size && bytes[i] != '\n'.code.toByte()) i++
9         val s = String(bytes, start, i - start, Charsets.UTF_8)
10        i++ // skip newline
11        return s
12    }
13
14    // parse header (omitted, see Listing 2) ...
15    val codeTable = mutableMapOf<String, Char>() // Note: inverse map
16    // ...
17    val bitLength = readLine().toInt()
18    // ...
19    val payload = bytes.copyOfRange(i, bytes.size)
20
21    // convert bytes back to bit string
22    val bitString = buildString(payload.size * 8) { /* ... */ }
23        .substring(0, bitLength) // trim to actual bit length
24
25    // decode by matching prefixes
26    val decoded = StringBuilder()
27    var cur = ""
28    for (bit in bitString) {
29        cur += bit
30        codeTable[cur]?.let { ch ->

```

```

31         decoded.append(ch)
32         cur = "" // reset for next code
33     }
34 }
35 return decoded.toString()
36 }

```

Runtime: $\Theta(B)$ to rebuild the ‘bitString’ and $\Theta(B)$ to scan it (since each bit is appended and checked once). Total: $\Theta(B)$.

Space: $\Theta(B)$ for the ‘bitString’ and $\Theta(B)$ for the output ‘StringBuilder’.

4 Correctness (DSA-style)

Prefix property. The tree construction guarantees the prefix-free property. All symbols are at leaves. A code is a path from the root to a leaf. By definition, no path to a leaf can be a prefix of another path to a different leaf. The explicit check in `generateCodes` (omitted from listing) provides a redundant safeguard.

Greedy decoding. Because the code set is prefix-free, the first match found when scanning the bitstream is guaranteed to be the **only** valid match. There is no ambiguity. This allows for a simple, linear-time greedy scan (as seen in Listing 6). The loop invariant is: at the end of each iteration, the `cur` buffer is either empty or a valid prefix of one or more codes in the table.

Tree optimality. The standard Huffman exchange argument proves optimality. The greedy choice is to merge the two subtrees with the smallest frequencies. Any optimal tree can be transformed into the one generated by this algorithm without increasing the total weighted path length, proving the greedy choice is always part of an optimal solution.

5 Asymptotic Runtime and Space

Let n be the number of input characters, k the number of distinct symbols ($k \leq n$), and B the total number of bits in the compressed payload.

Phase	Time	Space
Frequency analysis	$\Theta(n)$	$O(k)$
Tree construction (heap)	$O(k \log k)$	$O(k)$
Code generation (DFS)	$\Theta(k)$	$O(k)$
Encoding (build+pack)	$\Theta(n + B)$	$O(B)$
Decoding (parse+scan)	$\Theta(B)$	$O(B)$

For typical text files, $k \ll n$. The dominant cost is linear in the input size n and compressed size B . The $O(k \log k)$ term is usually negligible. The $O(B)$ space for intermediate bitstrings is the main scalability limit.

6 Interactive Interface & Visualization

A key feature of this project is the `HuffmanVisualizer` class, which provides insight into the algorithm’s state. The `visualizeFromText` method prints:

- A **frequency table** with scaled ASCII bars.
- The **generated code table**, sorted by code length.
- The full **ASCII tree structure**, implemented with a recursive printer.

Listing 7: Recursive ASCII tree printer

```

1 private fun printTree(node: Node?, prefix: String = "", isTail: Boolean = true) {
2     if (node == null) return
3     val connector = if (isTail) " " else " "
4     val extension = if (isTail) " " else " "
5
6     if (node.isLeaf()) {
7         val display = when (node.symbol) {
8             ' ' -> "SPACE"
9             '\n' -> "NEWLINE"
10            else -> "'${node.symbol}'"
11        }
12        println("$prefix$connector$display (${node.frequency})")
13    } else {
14        println("$prefix$connector[${node.frequency}]")
15        node.left?.let { printTree(it, prefix + extension, false) }
16        node.right?.let { printTree(it, prefix + extension, true) }
17    }
18 }

```

The `interactiveMain()` function wraps all functionality in a simple menu-driven CLI, allowing users to visualize text, load from files, compress, and decompress.

7 Implementation Reflection (Kotlin)

Strengths. Kotlin was an excellent choice. `data class` makes the `Node` ADT trivial. Local functions (like `traverse` in Listing 4 and `readLine` in Listing 6) are perfect for helper routines that share parent scope. The `require(...)` function provides clean, expressive precondition checks. Java interop with `PriorityQueue` is seamless.

I/O Pitfalls & Fix. The most critical part of this implementation is mixing text and binary data. Reading the whole file as text will corrupt the binary payload. Writing the header as text and *appending* bytes ('appendBytes') is correct, but reading is harder. The solution (Listing 1st:decode) is to read the *entire file* as raw `ByteArray`, then parse the header portion manually with a byte-aware `readLine` helper, and finally slice the remaining `payload` bytes. This is robust and binary-safe.

Memory Considerations. As noted, building the entire `bitString` as a `StringBuilder` (or `String`) consumes $O(B)$ memory. For a 1GB file, this could require gigabytes of RAM for the intermediate string. A production-grade implementation would use a streaming `BitWriter` and `BitReader` to operate in small $O(1)$ or $O(\text{bufferSize})$ memory.

Alphabet Choice. This implementation uses `Char` as the symbol. This works well for text files but is not general. A more robust implementation would treat `Byte` (0-255) as the symbol, allowing it to compress *any* file (images, executables, etc.). The logic of the algorithm remains identical.

8 Acknowledgments

General Acknowledgments I express my gratitude to various resources that contributed to this work. Specifically, I utilized ChatGPT to assist with LaTeX writing and formatting. Additionally, I consulted several reputable online sources for coding guidance, including Wikipedia,

GeeksforGeeks, MIT's resources on Huffman coding, and Berkeley EECS materials. ChatGPT also provided valuable support in developing the command-line interface (CLI) tool.

End of report.