

**Course: MsDS**

**Data Structures and Algorithms**

**Module: 6**

## **Learning Objectives:**

1. Understand how Tries are structured and used for effective string manipulation and retrieval.
2. Understand the fundamentals of Huffman Coding and how it is used to compress data using Binary Trees.
3. Understand the value of graphs in describing relationships and investigate different methods of storing them.
4. Understand the fundamental graph operations, such as identifying Paths, performing Traversals, and evaluating connectivity.

## **Structure:**

- 6.1 Tries
- 6.2 Huffman Coding using Binary Trees
- 6.3 Graphs for Capturing Relations; Storing Graphs
- 6.4 Basic Operations on Graphs
- 6.5 Summary
- 6.6 Keywords
- 6.7 Self-Assessment Questions
- 6.8 Case Study
- 6.9 References

## 6.1 Tries

Within the domain of data structures, the Trie also referred to as a prefix tree, distinguishes itself as a very efficient tree-based structure meticulously devised to handle a dynamic collection of strings. In contrast to binary search trees, which establish pathways based on numerical sequencing, Tries is specifically constructed to accommodate the characters comprising the strings. This differentiation enables them to provide a distinctive combination of benefits in the context of managing and interrogating large collections of strings.

In a Trie data structure, the nodes are often connected to the letters of the string in a manner that each route traversed in the tree corresponds to either a whole word or a prefix. As an example, it can be seen that the strings "bat", "batman", and "bath" exhibit a common prefix of "bat" inside the Trie data structure. However, subsequent to this shared prefix, each string diverges along distinct paths to form its own full strings. The use of this shared structure not only serves to reduce the amount of storage space needed for common prefixes but also enhances the efficiency of queries pertaining to string retrieval and existence checks.

One of the pragmatic uses of Tries lies in their use for the development of dictionaries and phonebooks. The use of a Trie data structure enables rapid traversal over its nodes to provide auto-completions when a user begins inputting the beginning letters of a word or name. This functionality enhances user interactions by providing a more intuitive and efficient experience. Moreover, the height of a Trie is precisely proportional to the length of the longest string it contains, resulting in speedier search operations compared to standard balanced trees.

Nevertheless, similar to other data structures, Tries possess their own distinct set of trade-offs. Although these data structures provide improved efficiency in terms of quicker insertions, deletions, and search times for strings, it is important to note that they may have a higher memory consumption, particularly when the size of the alphabet is big. However, in scenarios where there is a need for efficient string operations, the benefits of using a Trie data structure generally outweigh its drawbacks.

## **6.2 Huffman Coding using Binary Trees**

### **6.2.1 Huffman Coding**

Within the field of data compression, Huffman Coding is well-recognised as a leading method for achieving lossless compression. The approach, developed by David Huffman during his PhD research, aims to maximise the efficiency of code lengths for symbols by taking into account their individual frequencies. The fundamental idea of this approach is to allocate shorter codes to symbols that occur more often, whereas symbols that occur less frequently are assigned longer codes. Huffman Coding distinguishes itself by the use of a variable-length prefix code mechanism, which effectively prevents any code assigned to a particular symbol from serving as a prefix for another symbol. Consequently, this property ensures the unambiguous decodability of the encoded data.

### **6.2.2 Binary Trees: The Core of Huffman Coding**

Binary trees play a crucial part in the implementation of the Huffman Coding method. The procedure commences by generating distinct nodes for every symbol and its corresponding frequency. The assemblage of nodes constitutes an initial "forest", whereby each individual node exists as an independent tree. During the algorithm's progression, the process involves selecting two trees with the lowest frequencies of their root nodes and then merging them into a singular tree. The aggregate occurrence of these trees becomes established, with the initial two trees extending as the left and right offspring. The aforementioned process is repeated repeatedly until a single overarching Huffman Tree exists. In order to derive codes, the tree is traversed, with a leftward movement indicating the value '0' and a rightward movement indicating the value '1'. Upon reaching a leaf node, the route travelled provides the Huffman code for the corresponding symbol.

### **6.2.3 The Practical Implementation of Huffman Coding**

The usefulness of Huffman Coding is wide-ranging, especially in the context of file compression programs. Huffman Coding is widely used in several data compression technologies to provide lossless compression across different data formats, including

text, pictures, and audio. When combined with other methodologies, it is possible to achieve substantial reductions in file sizes without compromising the integrity of the original data. However, the use of Huffman Coding presents a series of problems. The inclusion of the Huffman Tree alongside the compressed data as an extra requirement may sometimes undermine the advantages of compression, particularly when dealing with small datasets. In dynamic environments characterised by fluctuating symbol frequencies, it may be necessary to regularly update the Huffman Tree. However, Huffman Coding remains a prominent example of algorithmic innovation, offering an optimised approach to address challenges related to data storage and transmission. It guarantees performance in applications that handle large amounts of data, even in the presence of constraints such as limited bandwidth or storage capacity.

## **6.3 Graphs for Capturing Relations; Storing Graphs**

### **6.3.1 Graphs**

Graphs are fundamental mathematical structures that play a crucial role in reflecting the interactions between different elements. Graphs are composed of vertices, also known as nodes, and edges, sometimes known as links, which serve to connect these vertices. Graphs have many uses in several disciplines, including computer science, network theory, sociology, and biology. One example of representing social networks involves the use of graphs, whereby people are shown as nodes, and their connections or interactions are depicted as edges. Within the realm of computer networks, the term "nodes" refers to entities that may take the form of either computers or servers. These nodes are interconnected by means of "edges," which symbolise the communication pathways facilitating the exchange of information between them. Graphs are very effective tools for representing and analysing linked and relational data, rendering them essential for visualising and processing interactions across many areas.

### 6.3.2 Storing Graphs: Representation Techniques

The representation of graphs plays a key role in computational applications. In the field of graph storage, two prevailing approaches have a significant influence: adjacency matrices and adjacency lists. The adjacency matrix is a square matrix that represents the connectivity between vertices in a graph. Each element in the matrix, located at row  $(i)$  and column  $(j)$ , denotes the existence of an edge between vertex  $(i)$  and vertex  $(j)$ , often denoted by a '1' for presence and '0' for absence. Although this approach is quite simple, it may suffer from inefficiency in terms of space utilisation when dealing with graphs that have a low density of edges. In contrast, adjacency lists use an array of lists. The list located at index  $(i)$  consists of the adjacent vertices of vertex  $(i)$ . The proposed technique demonstrates a high degree of space efficiency, particularly when used for graphs with a low density of edges. Additionally, it provides expedited traversal capabilities for certain processes. The selection between these methodologies often hinges upon the particular application scenario and the level of density shown by the graph.

### 6.3.3 From Theory to Real-World Applications

Graphs possess not just theoretical foundations but also significant practical implications. Graphs play a fundamental role in several applications, ranging from routing algorithms used in navigation systems to recommendation engines utilised in online platforms for suggesting friends or items. The efficiency of operations done on storage structures, such as adjacency matrices or lists, is determined by the chosen storage strategies. As the complexity of data increases, the significance of graphs in collecting, expressing, and analysing connections becomes more evident. The flexibility and complexity of graph topologies make them a prominent tool for comprehending and manipulating relational data.

## 6.4 Basic Operations on Graphs

### 6.4.1 Graph Connectivity

Graphs, being entities that depict associations and linkages, derive a significant portion of their usefulness from discerning the characteristics of these linkages. One

fundamental term that is pertinent in this context is the notion of "connectivity." Fundamentally, connectivity refers to the determination of the existence of a route connecting two vertices inside a network. A graph is considered linked if there is a route between every pair of vertices in the undirected graph. On the other hand, if there are vertices in a network that cannot be reached from other vertices, we classify the graph as disconnected. Furthermore, a more robust kind of connection, known as "strongly connected," is specifically applicable to directed graphs, in which there must be a directed route connecting every pair of vertices. The examination of connection enables the exploration of several applications, such as network architecture, where the preservation of optimal connectivity is crucial in order to avert disruptions in communication.

#### **6.4.2 Graph Traversals**

Traversal operations play a crucial role in the field of graph theory since they provide systematic approaches to visiting every vertex of a graph. There are two main techniques often used for graph traversal: Depth-First Search (DFS) and Breadth-First Search (BFS).

The Depth-First Search (DFS) algorithm starts at a chosen vertex and proceeds to examine each branch as much as feasible before retracing its steps. The comprehensive investigation conducted in DFS renders it especially advantageous in situations such as identifying linked components in a graph or performing topological sorting of vertices in a directed acyclic graph.

In contrast, the Breadth-First Search (BFS) algorithm starts by selecting a certain vertex and then examines all of its neighbouring vertices. Consequently, for each of the closest vertices, it then examines its neighbouring vertices, and this iterative procedure persists. The Breadth-First Search (BFS) algorithm is particularly useful in several applications, such as the determination of the shortest route in an unweighted network or the identification of all nodes inside a single linked component.

Both traversal approaches provide unique viewpoints and utility, with their suitability being decided by the individual activity or intended outcome.

### 6.4.3 Paths

In graph theory, a route is defined as an ordered series of vertices in which each consecutive pair of vertices is linked by an edge. Paths play a crucial role in several activities and difficulties pertaining to graphs. The most efficient route connecting two vertices is sometimes referred to as the shortest path, which is determined by the path with the fewest number of edges. In the case of weighted graphs, the shortest path is determined by the minimal total weight. Algorithms such as Dijkstra's and Floyd-Warshall have been devised to effectively ascertain the shortest routes inside graphs.

Nevertheless, it should be noted that pathways are not always confined to the shortest possible routes. Certain issues may pertain to the existence of pathways that traverse each vertex exactly once, known as Hamiltonian paths, or paths that traverse each edge exactly once, known as Eulerian paths. The comprehension and identification of such routes are fundamental to the complexities presented by problems such as the Travelling Salesman Problem or the Seven Bridges of Königsberg.

## 6.5 Summary

- ❖ Tries are a very effective data structure for performing string operations, particularly for activities such as word look-up. They provide notable advantages in terms of both performance and space optimisation when compared to conventional ways of storing strings.
- ❖ Huffman Coding, which utilises Binary Trees, is a key approach used in data compression. It achieves lower data sizes by assigning shorter binary codes to frequently occurring characters. The implementation of Huffman Coding relies on the organisational structure of binary trees.
- ❖ Graphs play a fundamental role in representing and analysing interactions and connections across many areas. These systems not only serve as representations of intricate networks but also provide efficient techniques for the storage of relational data.



- ❖ The fundamental operations performed on graphs, including the determination of connectedness between nodes, the use of systematic traversal methods like Depth-First Search (DFS) and Breadth-First Search (BFS), and the implementation of pathfinding strategies, serve as the foundational components of several algorithms and real-world applications. These operations provide a more comprehensive understanding of linked systems, allowing for the extraction of valuable insights.

## 6.6 Keywords

- **Tries:** Efficient tree-like data structures for string storage and retrieval.
- **Huffman Coding:** The use of binary trees in a lossless data compression technique.
- **Graphs:** The structures that depict items and their relationships.
- **Connectivity:** The metric that quantifies the level of interconnectivity among the various components of a graph.
- **Traversals:** Methods such as Depth-First Search (DFS) and Breadth-First Search (BFS) are used in a systematic manner to traverse all vertices inside a graph.
- **Paths:** Routes connecting two vertices in a graph without revisiting any vertex.

## 6.7 Self-Assessment Questions

1. Describe the advantages of Tries over conventional tree structures for the storing and retrieval of strings. What are the main uses of them?
2. How does Huffman coding work? How does it guarantee more effective compression without data loss, and how does this approach use binary trees?
3. Explain the differences between a directed and an undirected graph. Give instances of each's practical applicability.
4. Draw a flowchart outlining the Depth-First Search (DFS) traversal stages. What distinguishes it from Breadth-First Search (BFS)?
5. Describe the importance of finding the shortest pathways in graphs. How does this problem get solved by algorithms like Dijkstra's or Floyd-Warshall?

## 6.8 Case Study

### **Title: Optimising Search and Storage for a Web-Based Dictionary**

#### **Introduction:**

In the current epoch characterised by rapid digitalisation, online dictionaries have supplanted their tangible equivalents for a multitude of users, providing convenient access to word searches and definitions. The effectiveness of these platforms is greatly dependent on the foundational data structures, particularly in the context of string manipulation.

#### **Case Study:**

WebsterNet, a widely used digital dictionary platform, has garnered a substantial user population as a result of its comprehensive lexicon and functionalities such as word recommendations, translations, and related word suggestions. Nevertheless, with the expansion of the user base, a discernible delay in word retrieval becomes evident, leading to a sense of discontentment among users.

#### **Background:**

The early architecture of WebsterNet had a rudimentary approach to searching, relying on basic string-matching techniques rather than on more advanced data structures. As the platform's word volume and user base expanded, it became more apparent that there were constraints in place. It was observed that some processes, particularly those involving string retrievals and recommendations, exhibited a higher level of time consumption than expected.

#### **Your Task:**

You've been hired as a Data Structures Specialist to revamp WebsterNet's string storage and retrieval mechanisms. The objective is to provide expedited word searches, optimised storage, and prompt word recommendations during user input.

#### **Questions to Consider:**

1. What are the potential methods for incorporating Tries into the system to enhance the efficiency of word recommendations and retrievals?
2. In light of the dictionary's ability to support several languages, how might Huffman Coding be used to enhance storage efficiency for a wide range of character sets?

3. What kind of graph-based format would be appropriate for representing the relationships between words and their synonyms or translations?
4. When taking into account the need for prompt response times, which style of traversal would be most appropriate for traversing word connections inside the system?

**Recommendations:**

The use of a Trie data structure might provide substantial improvements in terms of both word suggestion and retrieval performance. The use of Huffman Coding may provide efficient storage solutions, particularly when dealing with a wide range of characters found in many languages. The use of a graph-based methodology enables the representation of word associations, whereby adjacency lists serve as a storage mechanism that exhibits potential efficiency. In the context of efficient information retrieval, it is possible that Depth-First Search (DFS) might provide faster outcomes when applied to a densely interconnected word network.

**Conclusion:**

Within the dynamic landscape of digital platforms, the pursuit of optimisation is an ongoing and perpetual endeavour. By incorporating sophisticated data structures and algorithms, WebsterNet has the potential to regain its status as a prominent online dictionary, providing users with a streamlined and effective user interface.

## 6.9 References

1. Goodrich, M.T., Tamassia, R. and Goldwasser, M.H., 2014. *Data structures and algorithms in Java*. John Wiley & sons.
2. Mehlhorn, K., 2013. *Data structures and algorithms 1: Sorting and searching* (Vol. 1). Springer Science & Business Media.
3. Storer, J.A., 2003. *An introduction to data structures and algorithms*. Springer Science & Business Media.
4. Drozdek, A., 2013. *Data Structures and algorithms in C++*. Cengage Learning.
5. Lafore, R., 2017. *Data structures and algorithms in Java*. Sams publishing.