# Benchmark Compression with Near-Sortedness

By -
Shivangi
Vani Singhal
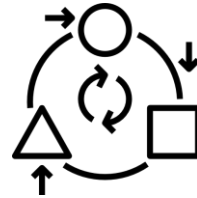
# Background

# Why need data compression?

Large amount of data generated

Run out of resources soon

# Motivation



Degree of Sortedness



Completely Sorted or Unsorted Data

## Problem Statement



Explore the Performance

# Workload

Different size of Workload
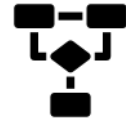
Different values of K-L
(Varying Sortedness)

# 5 Compression Algorithms
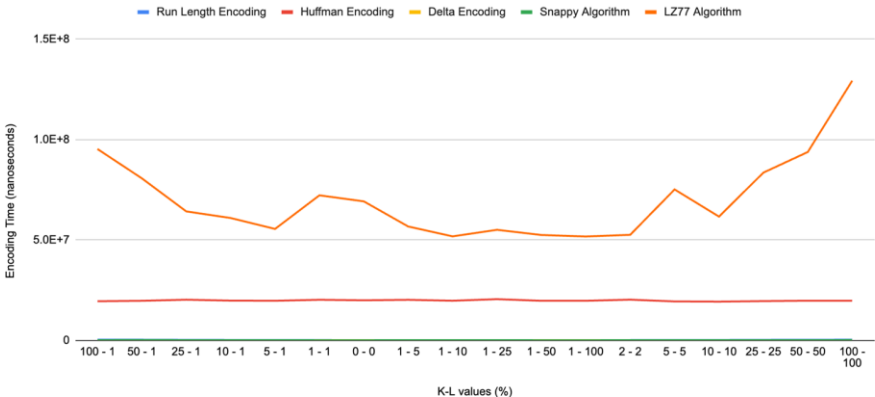
Run Length Encoding
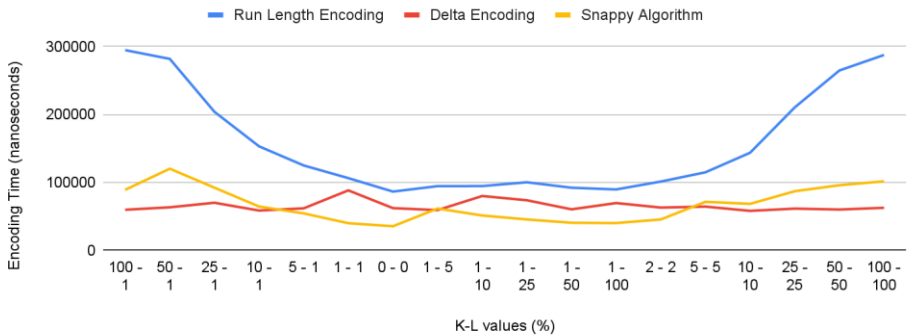
Delta Encoding

Snappy

LZ77

Huffman

# Results

# Workload 40KB



Workload 40KB - Encoding Time vs K-L values

Run Length Encoding — Huffman Encoding — Delta Encoding — Snappy Algorithm — LZ77 Algorithm



Workload 40KB - Encoding Time vs K-L values

Run Length Encoding — Delta Encoding — Snappy Algorithm

Workload 40KB -Compression Ratio vs K- L values

# Workload 400KB



Encoding Time vs K-L values

Run Length Encoding — Huffman Encoding — Delta Encoding — Snappy Algorithm — LZ77 Algorithm



Encoding Time vs K-L values

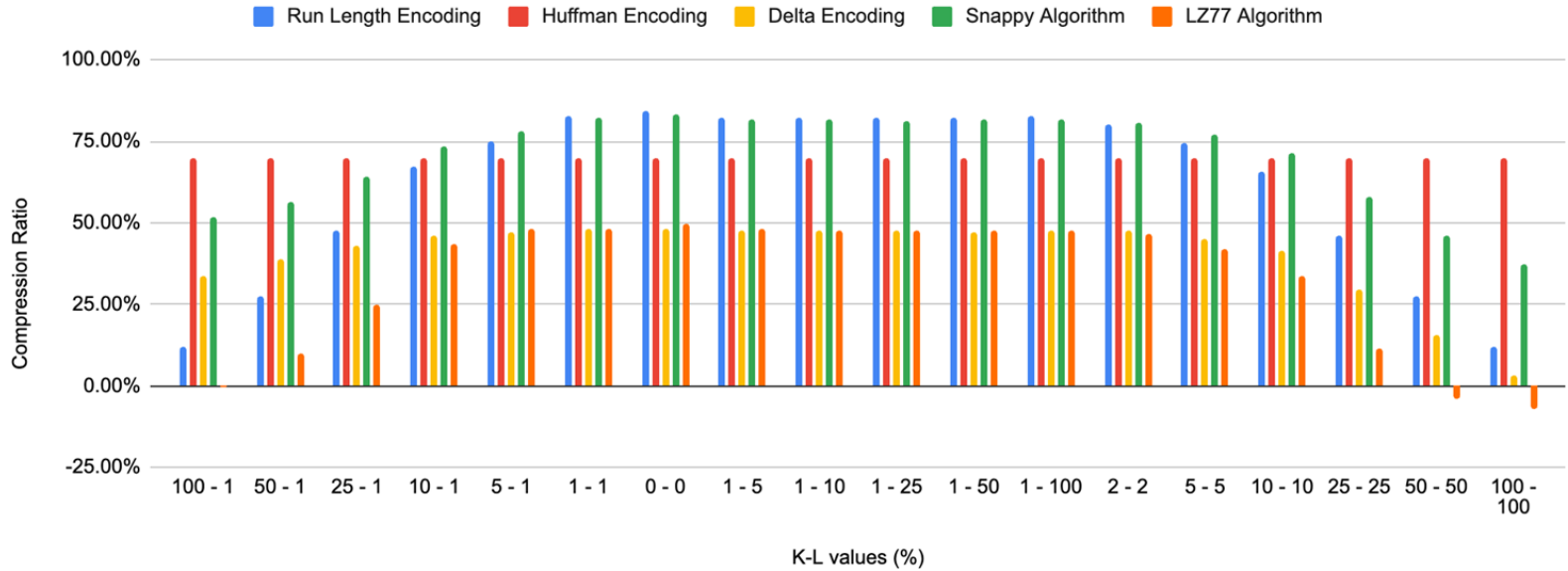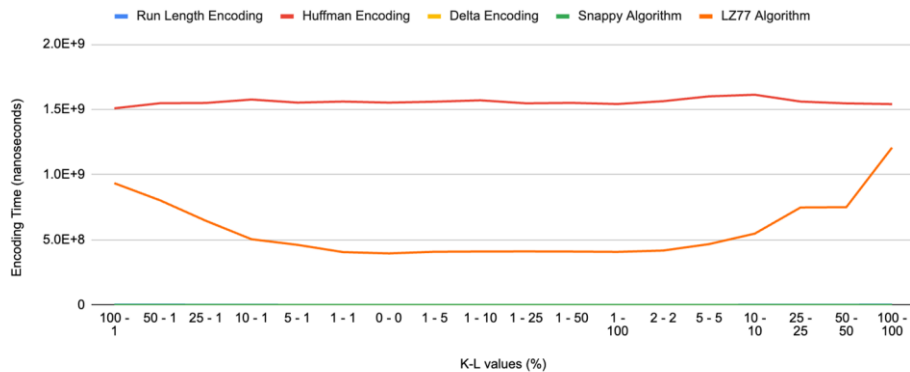Run Length Encoding — Delta Encoding — Snappy Algorithm

Workload 400KB -Compression Ratio vs K- L values

# Workload 4MB



Encoding Time vs K-L values
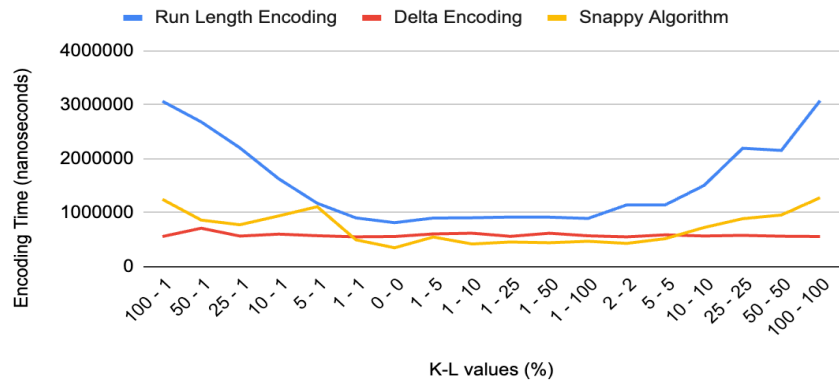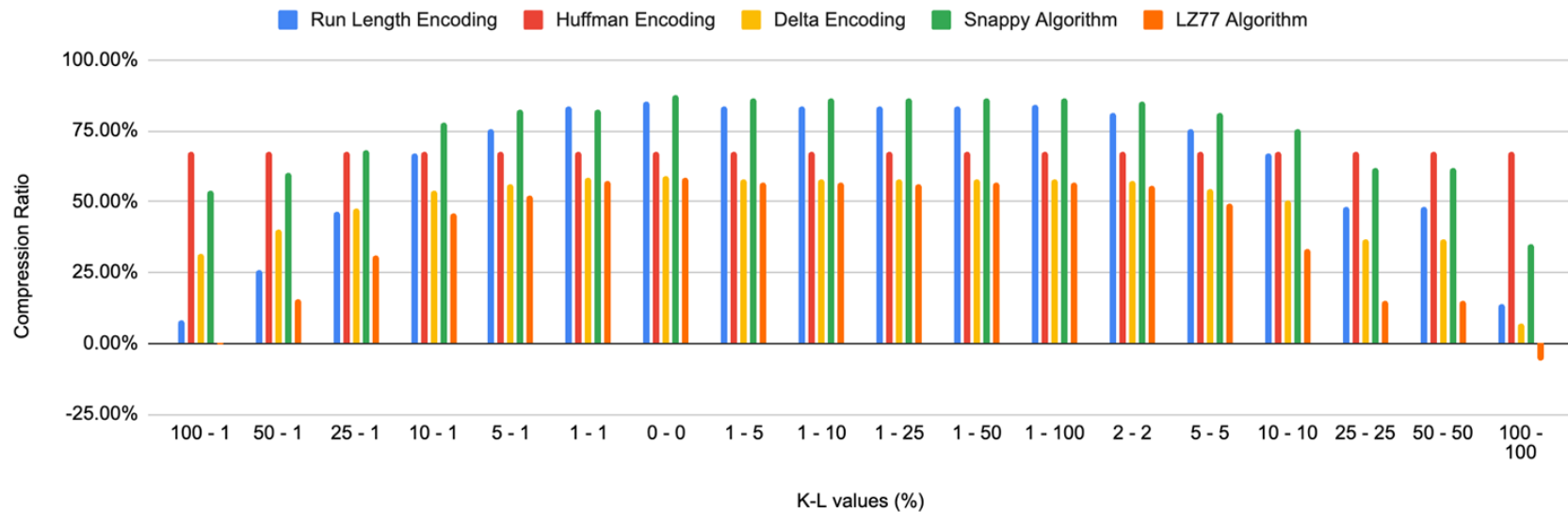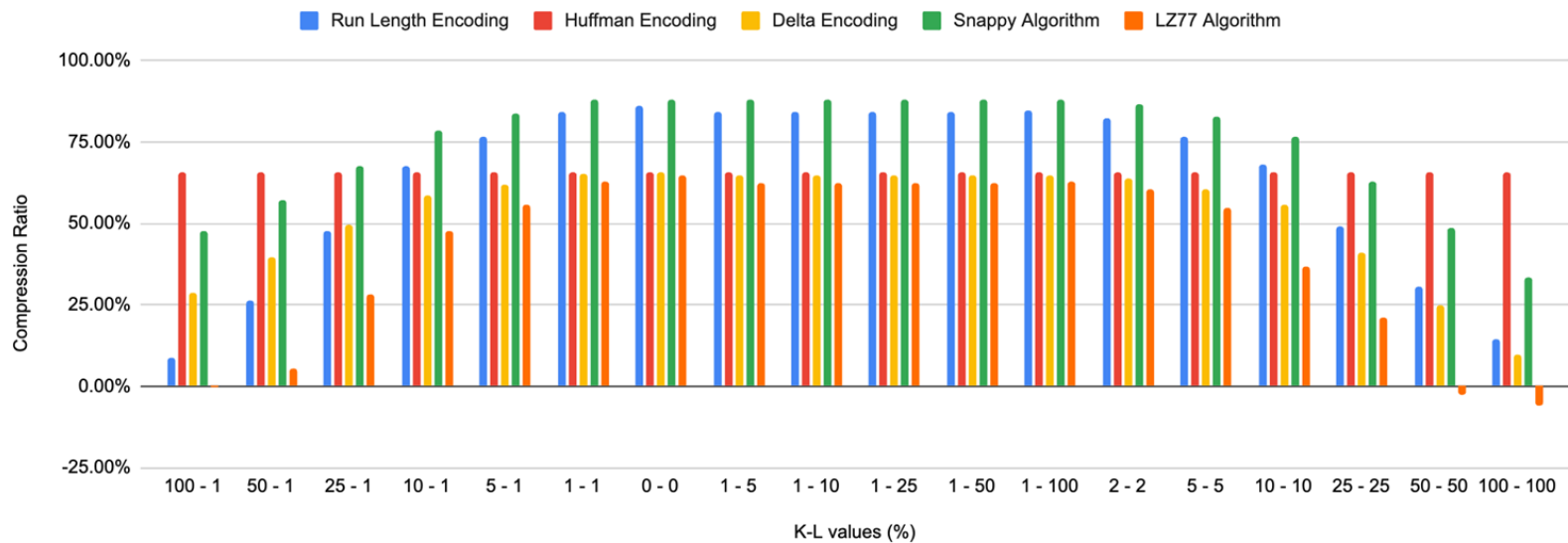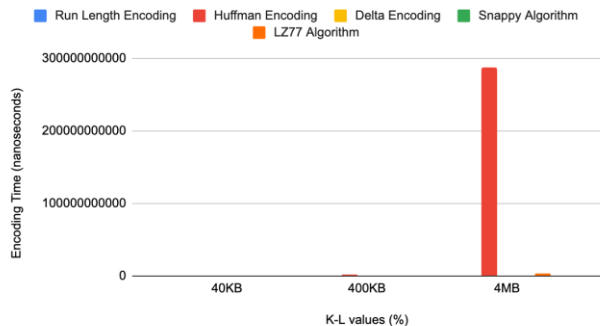
Workload 4MB -Compression Ratio vs K- L values

# Scalability for Sorted Data
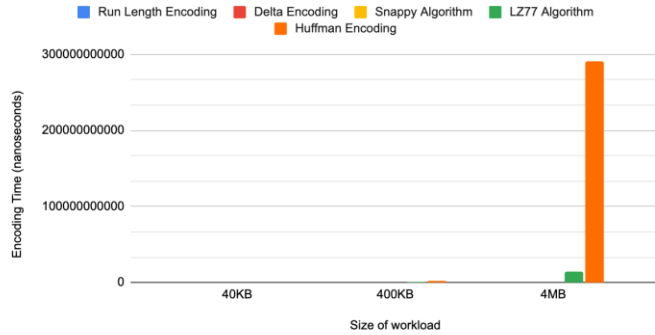


Sorted Data -varying Scalability

# Scalability for Unsorted Data



UnSorted Data -varying Scalability

Run Length Encoding · Delta Encoding · Snappy Algorithm · LZ77 Algorithm · Huffman Encoding



UnSorted Data -varying Scalability

Run Length Encoding · Delta Encoding · Snappy Algorithm · LZ77 Algorithm



UnSorted Data -varying Scalability

Run Length Encoding · Delta Encoding · Snappy Algorithm

# Scalability for Near Sorted Data



Near Sorted Data -varying Scalability

Run Length Encoding ▪ Delta Encoding ▪ Snappy Algorithm ▪ Huffman Encoding ▪ LZ77 Algorithm

Encoding Time (nanoseconds) vs Size of workload

Near Sorted Data -varying Scalability

Run Length Encoding ▪ Delta Encoding ▪ Snappy Algorithm ▪ LZ77 Algorithm

Encoding Time (nanoseconds) vs Size of workload

Near Sorted Data -varying Scalability

Run Length Encoding ▪ Delta Encoding ▪ Snappy Algorithm

Encoding Time (nanoseconds) vs Size of workload

# Next Steps

Unique Data

Larger Size of Workload - 40MB and 400MB

# Questions ?