# Query-driven compaction in LSM-trees

Ye Tian, Peiying Ye, Li Xi

# Log-Structured Merge-tree and Range Query

## LSM Tree

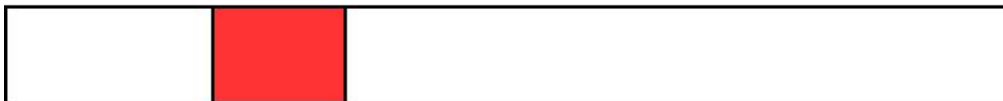**Fence pointers**

Buffer



Level 0

Level 1

Level 2

# Range Query and its Sorted View

Range query (min = 7, max = 40)

| $R_1$ | 6 | 7 | 17 | 29 | 73 |

| $R_2$ | 4 | 31 | 43 | 52 | 67 |

Access each run and
find the smallest key
>= 7 and <= 40.  ⟶

# Range Query and its Sorted View

Range query (min = 7, max = 40)

| $R_1$ | 6 | **7** | 17 | 29 | 73 |
|-------|---|-------|----|----|----|

| $R_2$ | 4 | **31** | 43 | 52 | 67 |
|-------|---|--------|----|----|----|

Access each run and find the smallest key >= 7 and <= 40.

→

Compare, build heap, and output the smallest.

# Range Query and its Sorted View

Range query (min = 7, max = 40)

| | | | | | |
|---|---|---|---|---|---|
| $R_1$ | 6 | **7** | **17** | **29** | 73 |

| | | | | | |
|---|---|---|---|---|---|
| $R_2$ | 4 | **31** | 43 | 52 | 67 |

| Access each run and find the smallest key >= 7 and <= 40. | → | Compare, build heap, and output the smallest. | → | Runtime Sort-merge completed. |
|---|---|---|---|---|

| **Duplicate Sort-merge next time!** | ← | **Discarded** after range query completed. | ← | **7, 17, 29, 31** |
|---|---|---|---|---|

# Research Motivation

Range query (min = 7, max = 40)

| R₁ | 6 | **7** | 17 | 29 | 73 |
|----|---|-------|----|----|----|

| R₂ | 4 | **31** | 43 | 52 | 67 |
|----|---|--------|----|----|----|

**7, 17, 29, 31** → **Duplicate Sort-merge next time!** → **How can we optimize it?**

# How does range-query-driven compaction work?
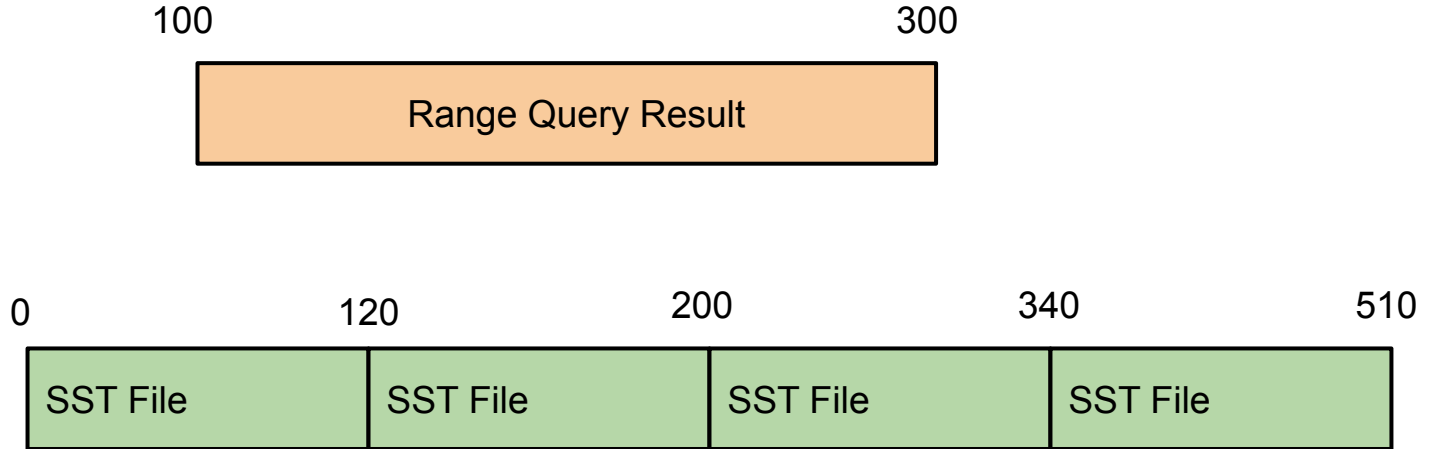


Data Boundary: 1-1000  Range Query: 40-500

Both count towards the beginning high write cost of range query-driven compaction.

# Algorithm

100                                    300

| Range Query Result |

0            120            200            340            510

| SST File | SST File | SST File | SST File |

# Algorithm

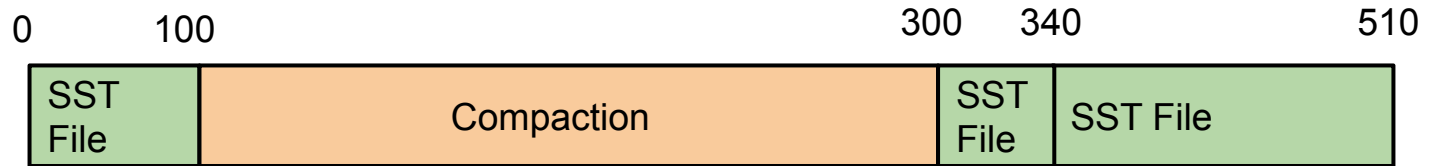# Algorithm

# Algorithm

# Algorithm

# Algorithm

# Algorithm

L1

L2

L3

L4

# Algorithm

# Algorithm

# Algorithm

L1

L2

L3

L4

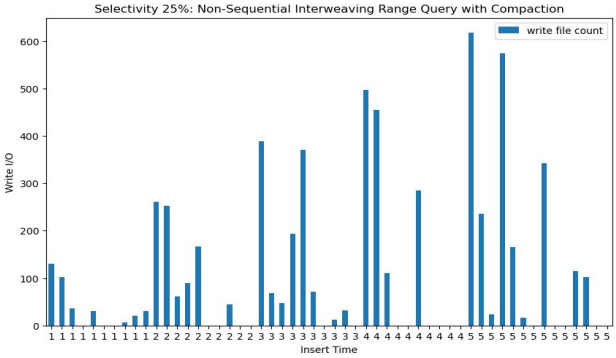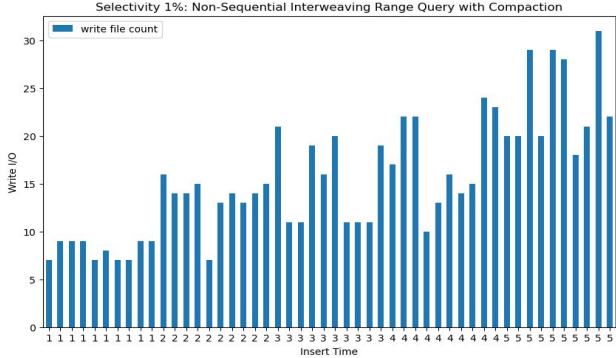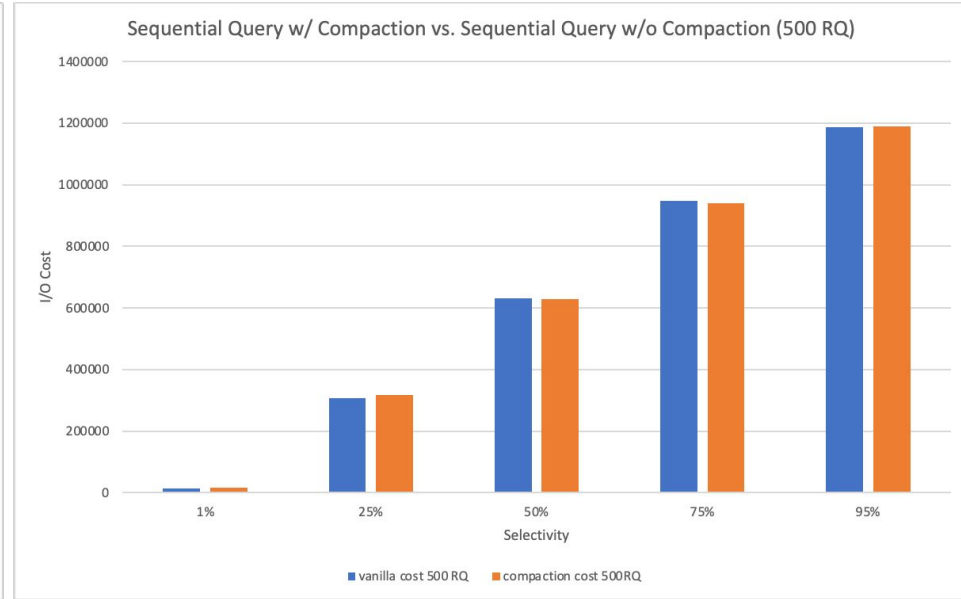# Benchmark Experiments  (100k Inserts)

## Sequential Workload (50RQ)

# Non-sequential Workload (5 times inserts, 100RQ each time)

# Comparison



Sequential Query w/ Compaction vs. Sequential Query w/o Compaction (50 RQ)

Sequential Query w/ Compaction vs. Sequential Query w/o Compaction (500 RQ)

# Challenges and Future Evaluations

- Further experiments on comparing non-sequential workload I/O for vanilla implementation and query-driven-compaction implementation
- Further experiments on 'QueryDrivenCompactionSelectivity' to figure out the trend and see if we can find a sweet spot

Sequential Query w/ Compaction vs. Sequential Query w/o Compaction - 50 RQ

Sequential Query w/ Compaction vs. Sequential Query w/o Compaction - 500 RQ

# Non-sequential Query: 500 RQ, 5 times insert



Non-sequential Query w/ Compaction vs. w/o Compaction 500 RQ with 5 Insert Times

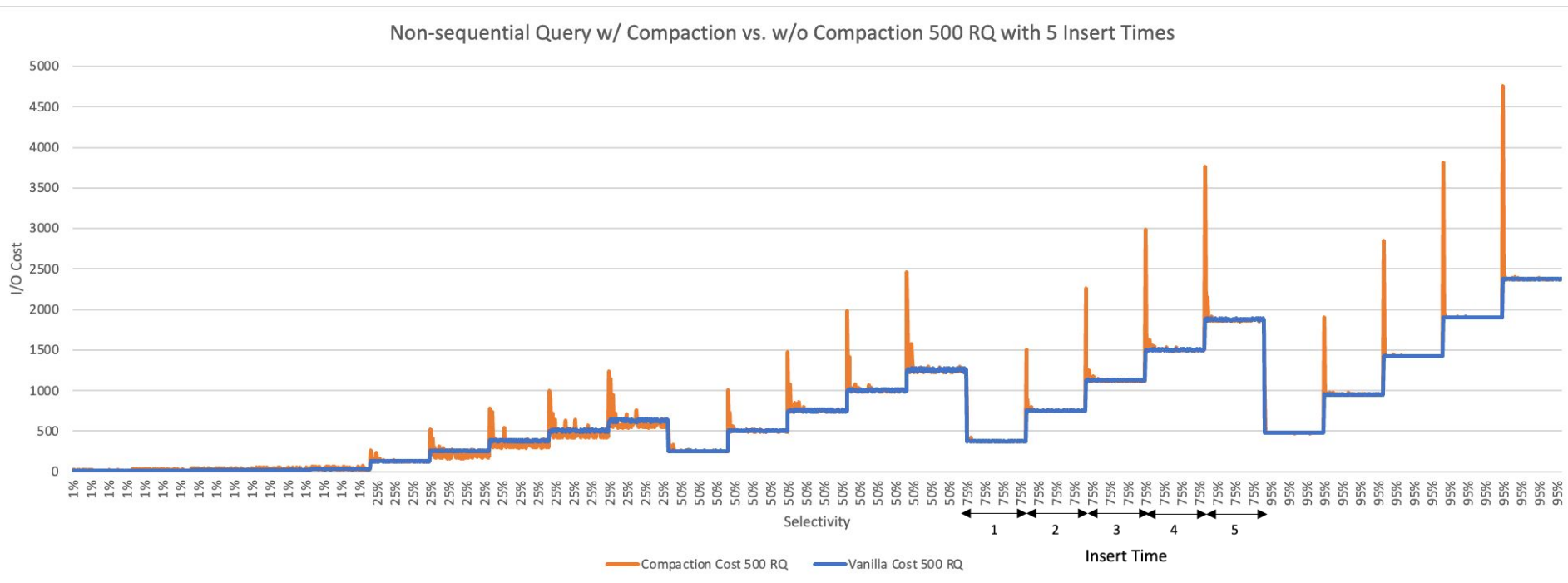# Non-sequential Query: 500 RQ, 5 times insert -1% Selectivity



Non-sequential Query 1% Selectivity Compaction vs. Vanilla: 500 RQ with 5 Insert Times

# Non-sequential Query: 500 RQ, 5 times insert -25% Selectivity



Non-sequential Query 25% Selectivity Compaction vs. Vanilla: 500 RQ with 5 Insert Times

# Non-sequential Query: 500 RQ, 5 times insert -50% Selectivity



Non-sequential Query 50% Selectivity Compaction vs. Vanilla: 500 RQ with 5 Insert Times

# Non-sequential Query: 500 RQ, 5 times insert -75% Selectivity

# Non-sequential Query: 500 RQ, 5 times insert -95% Selectivity
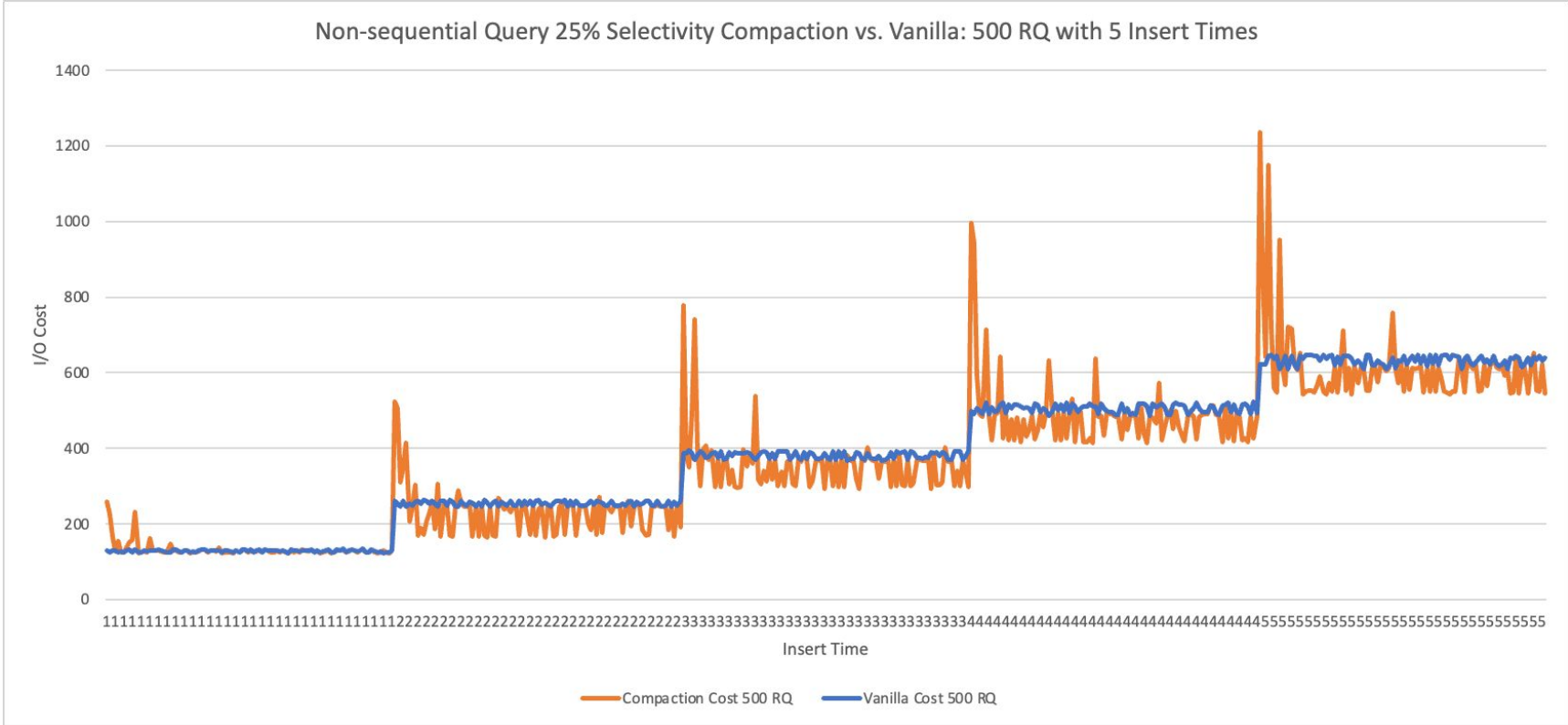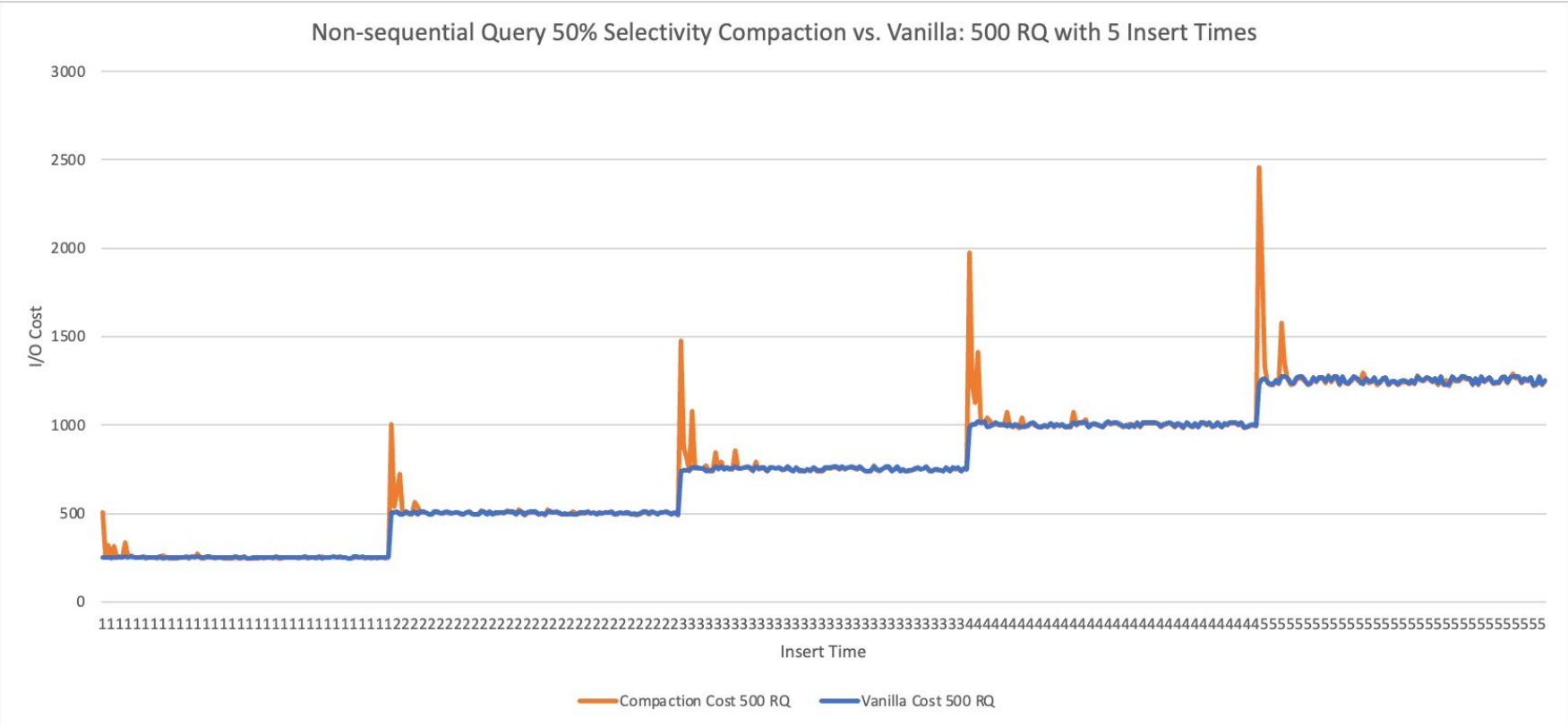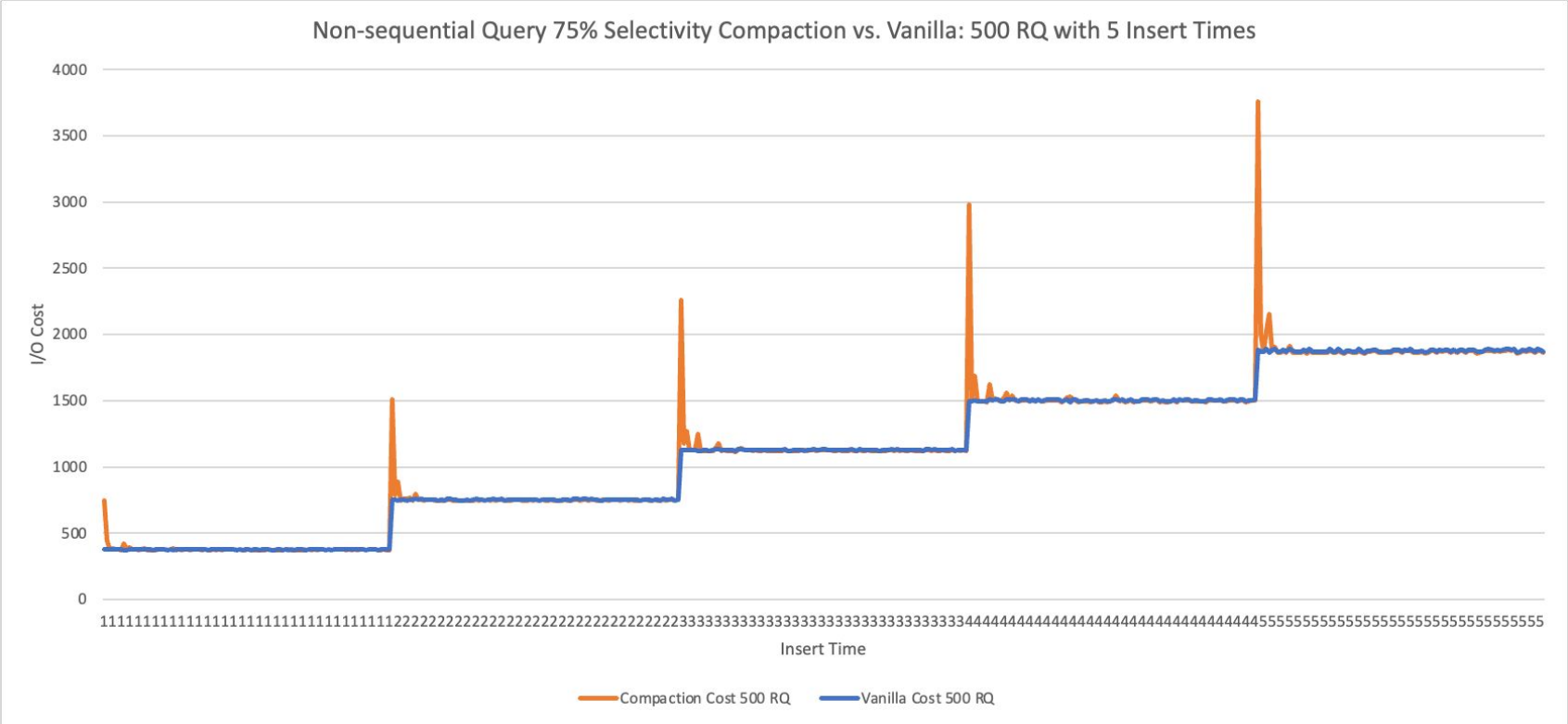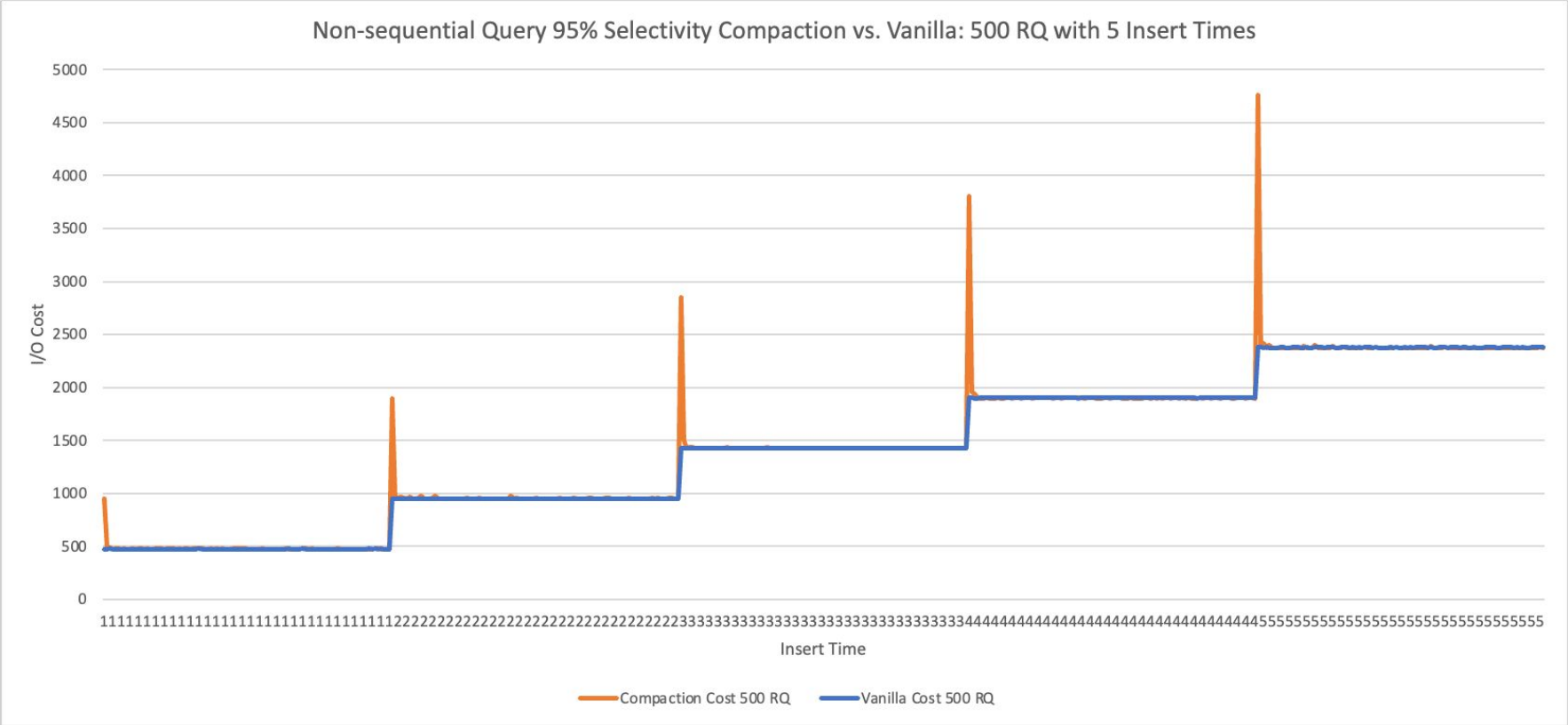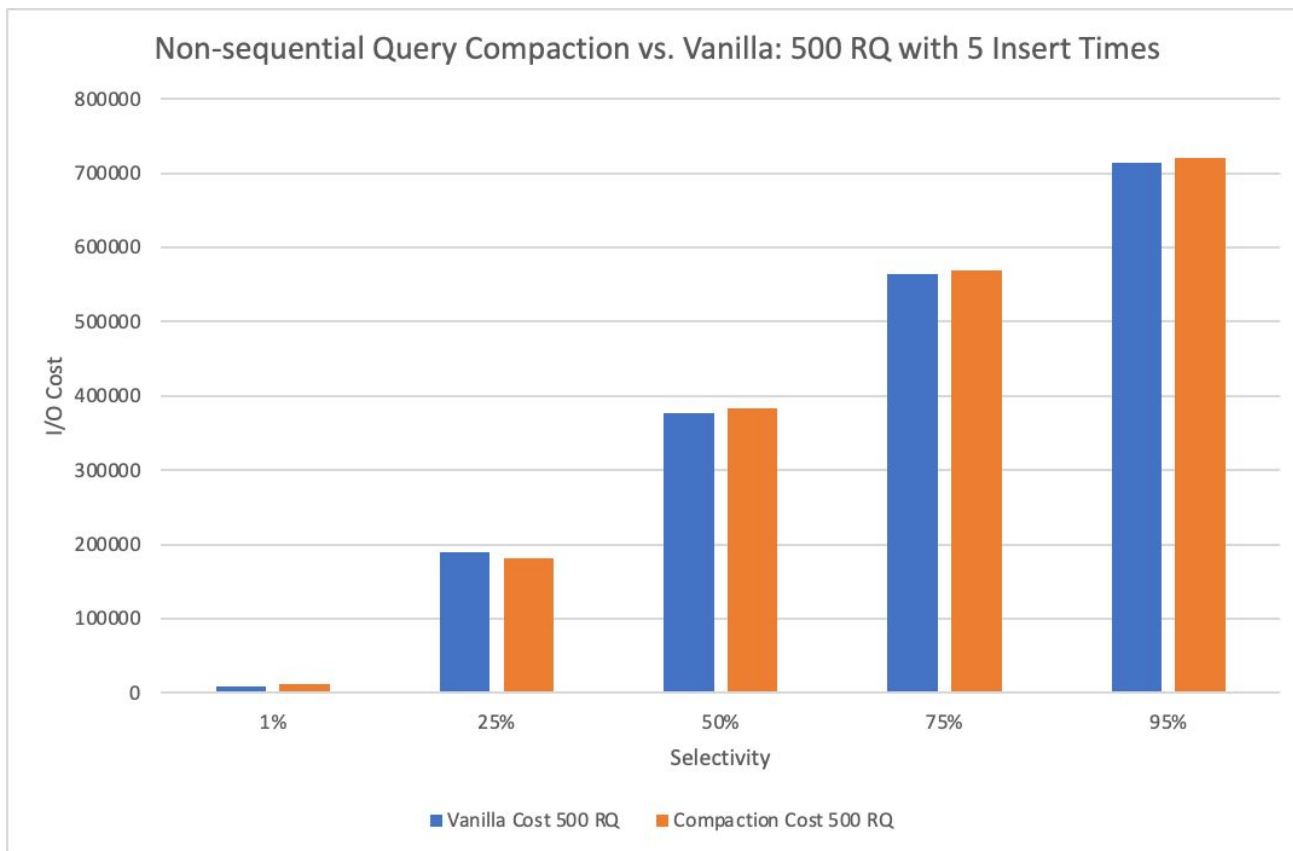


Non-sequential Query 95% Selectivity Compaction vs. Vanilla: 500 RQ with 5 Insert Times

# Non-sequential Query: 500 RQ, 5 times insert - Sum i/o

# Non-sequential Query: 500 RQ, 5 times insert - Insertion



Non-sequential 100K Insertion with 1% - 95% Selectivity and 5 Insert Times

# Non-sequential Query: 500 RQ, 10 times insert -1% Selectivity



Non-sequential Query 1% Selectivity Compaction vs. Vanilla: 500 RQ with 10 Insert Times

Compaction Cost 500 RQ    Vanilla Cost 500 RQ

# Non-sequential Query: 500 RQ, 10 times insert -25% Selectivity



Non-sequential Query 25% Selectivity Compaction vs. Vanilla: 500 RQ with 10 Insert Times

# Non-sequential Query: 500 RQ, 10 times insert -50% Selectivity



Non-sequential Query 50% Selectivity Compaction vs. Vanilla: 500 RQ with 10 Insert Times

# Non-sequential Query: 500 RQ, 10 times insert -75% Selectivity



Non-sequential Query 75% Selectivity Compaction vs. Vanilla: 500 RQ with 10 Insert Times

# Non-sequential Query: 500 RQ, 10 times insert -95% Selectivity



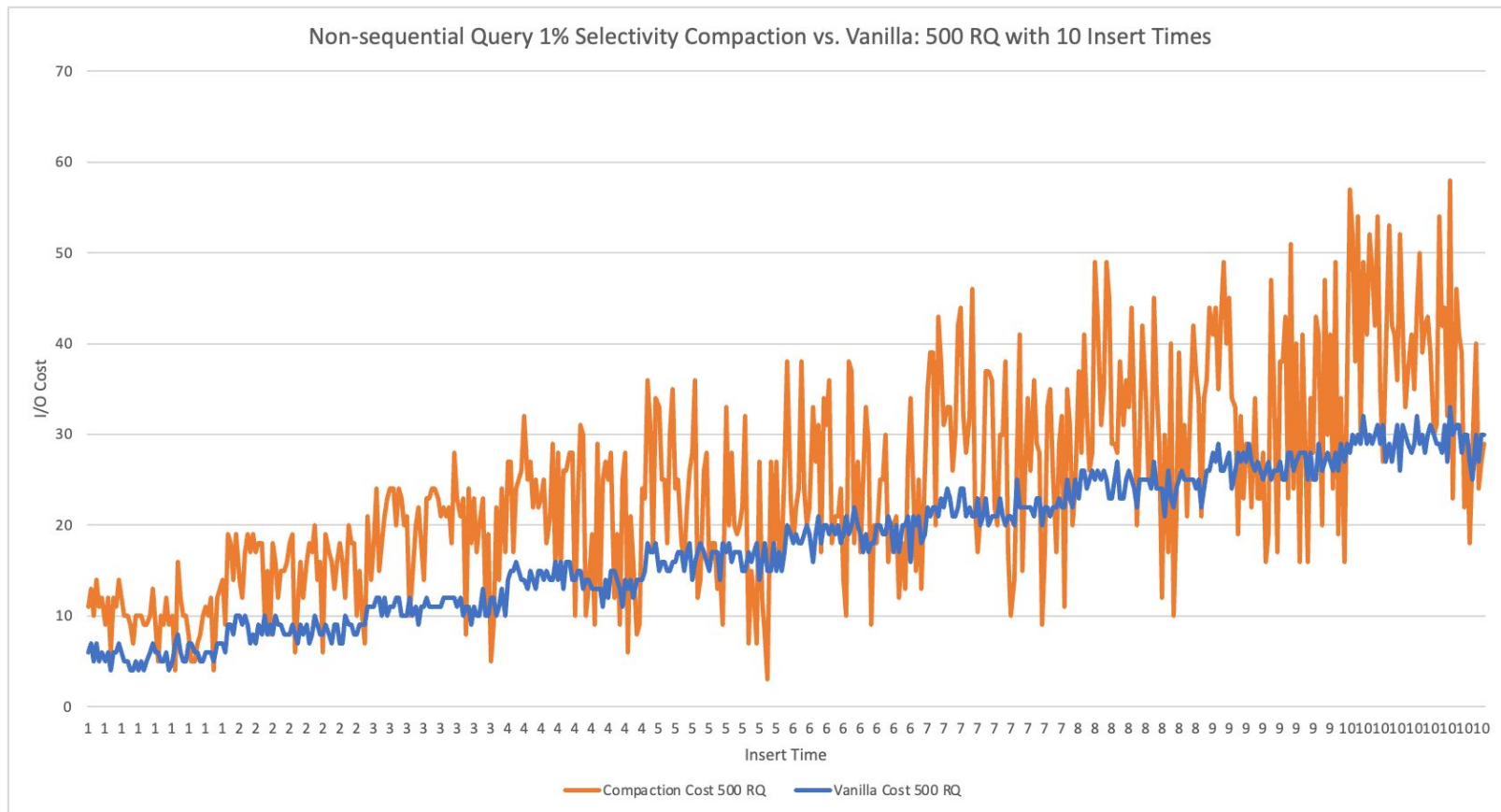Non-sequential Query 95% Selectivity Compaction vs. Vanilla: 500 RQ with 10 Insert Times

# Non-sequential Query: 500 RQ, 10 times insert - Sum i/o

# Non-sequential Query: 500 RQ, 10 times insert - Insertion



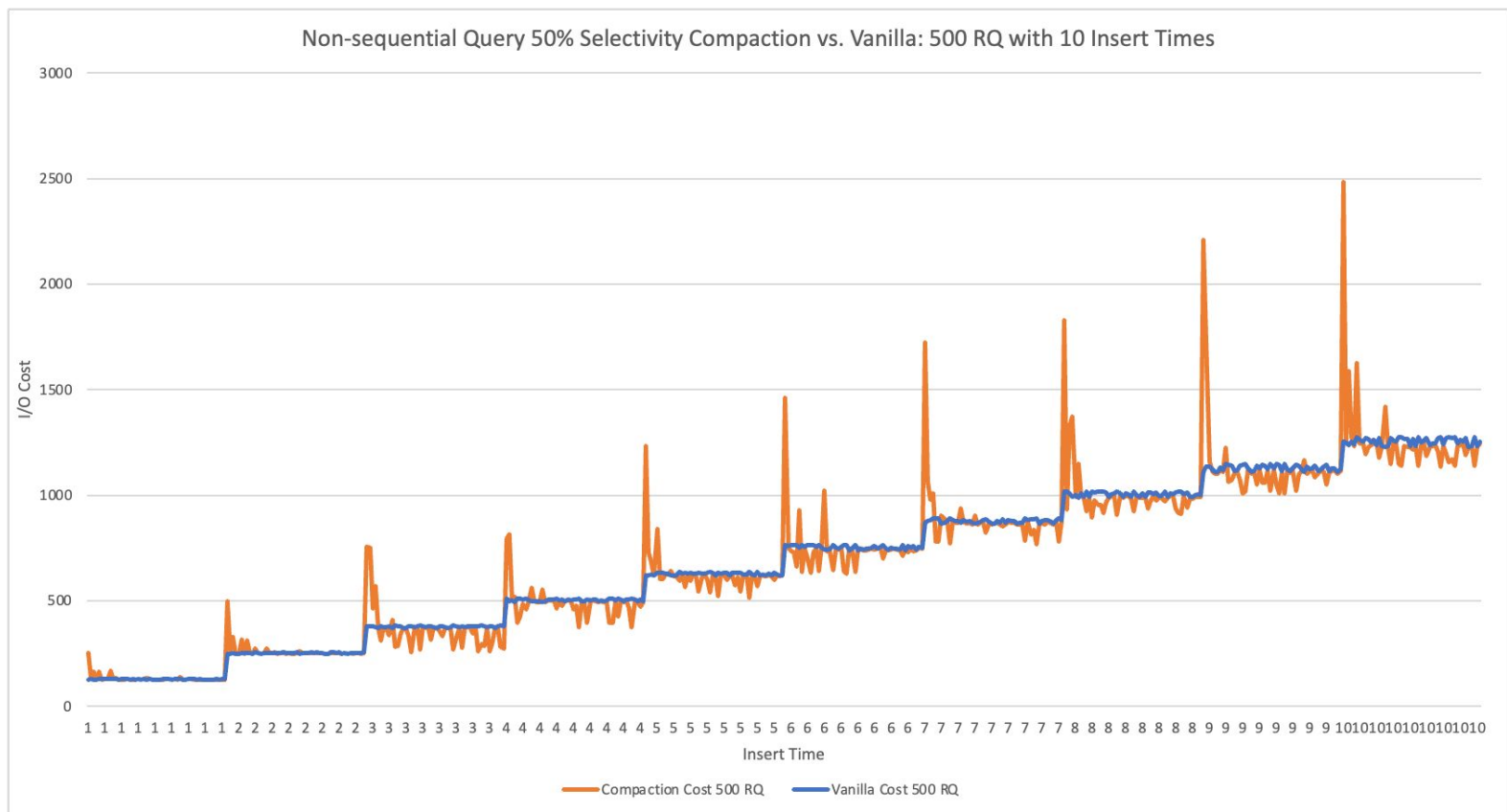Non-sequential 100K Insertion with 1% - 95% Selectivity and 10 Insert Times

# Compaction Selectivity: Sequential Query 1% - 500 RQ



Sequential Query 1% Selectivity Compaction: 500 RQ

# Compaction Selectivity: Sequential Query 25% - 500 RQ



Sequential Query 25% Selectivity Compaction: 500 RQ

# Compaction Selectivity: Sequential Query 50% - 500 RQ



Sequential Query 50% Selectivity Compaction: 500 RQ

# Compaction Selectivity: Sequential Query 75% - 500 RQ



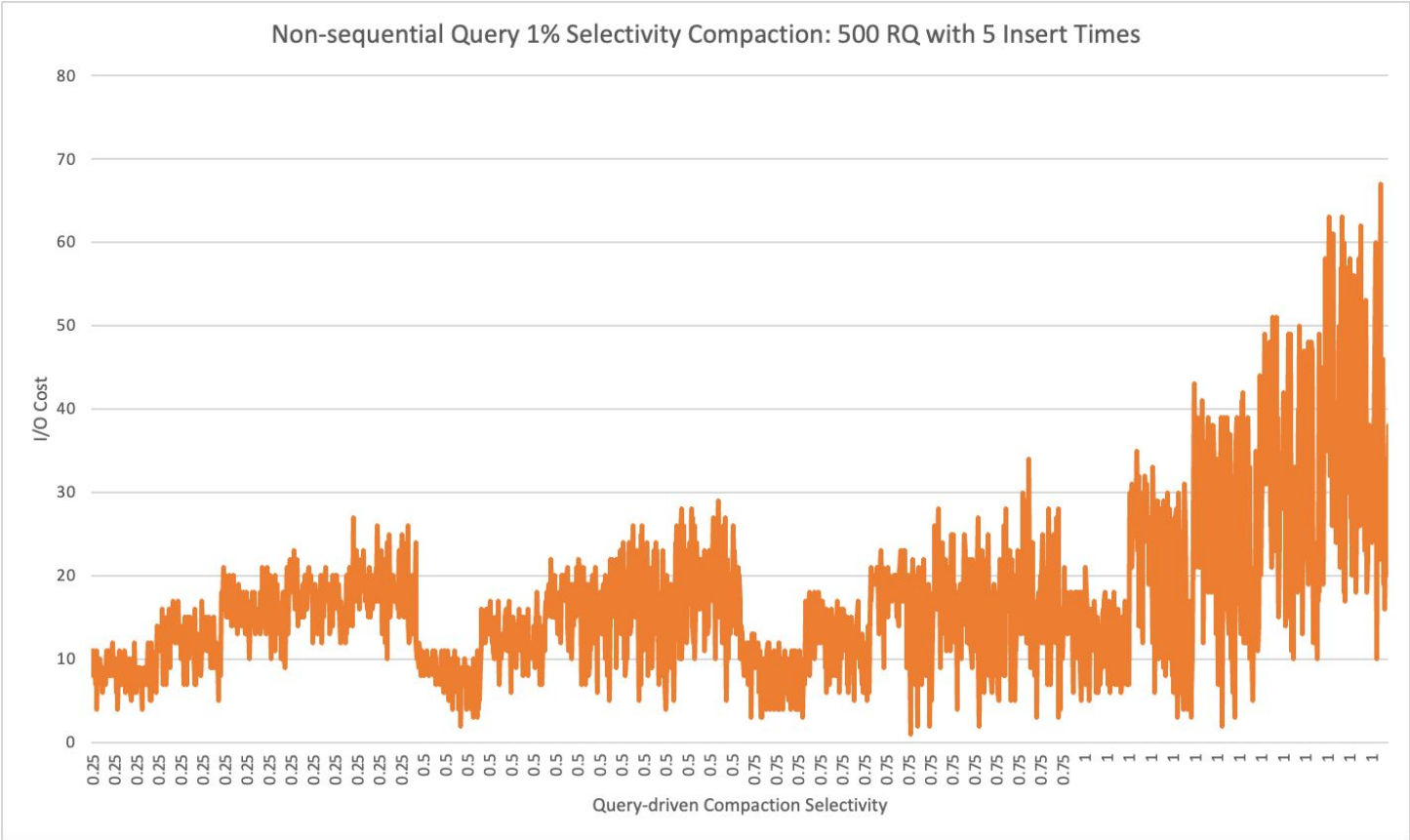Sequential Query 75% Selectivity Compaction: 500 RQ

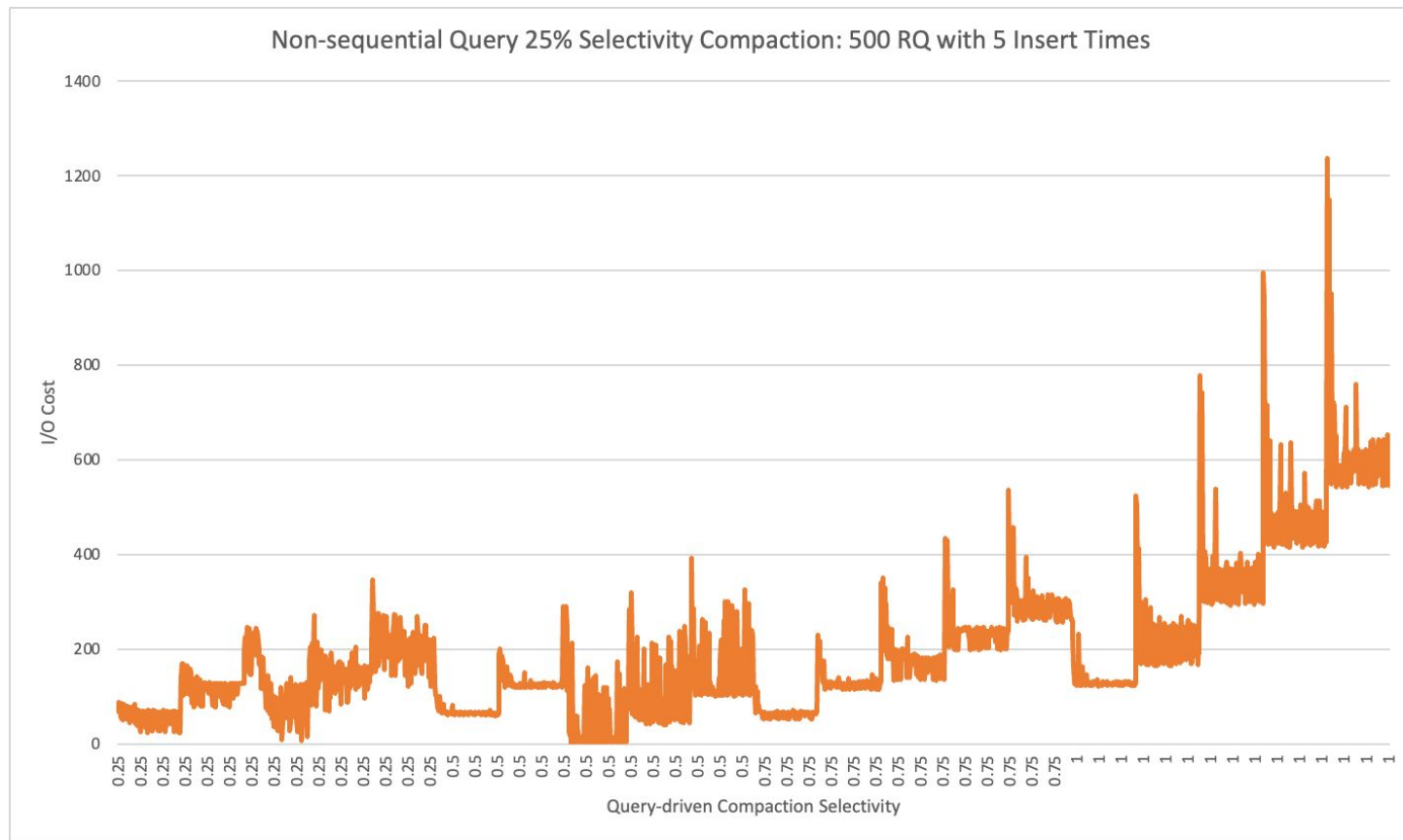# Compaction Selectivity: Sequential Query 95% - 500 RQ

# Compaction Selectivity: Sequential Query: 500 RQ - Sum i/o

No vanilla data

# Compaction Selectivity: Non-sequential Query 1% - 500 RQ - 5 Times



Non-sequential Query 1% Selectivity Compaction: 500 RQ with 5 Insert Times

# Compaction Selectivity: Non-sequential Query 25% - 500 RQ - 5 Times



Non-sequential Query 25% Selectivity Compaction: 500 RQ with 5 Insert Times

# Compaction Selectivity: Non-sequential Query 50% - 500 RQ - 5 Times



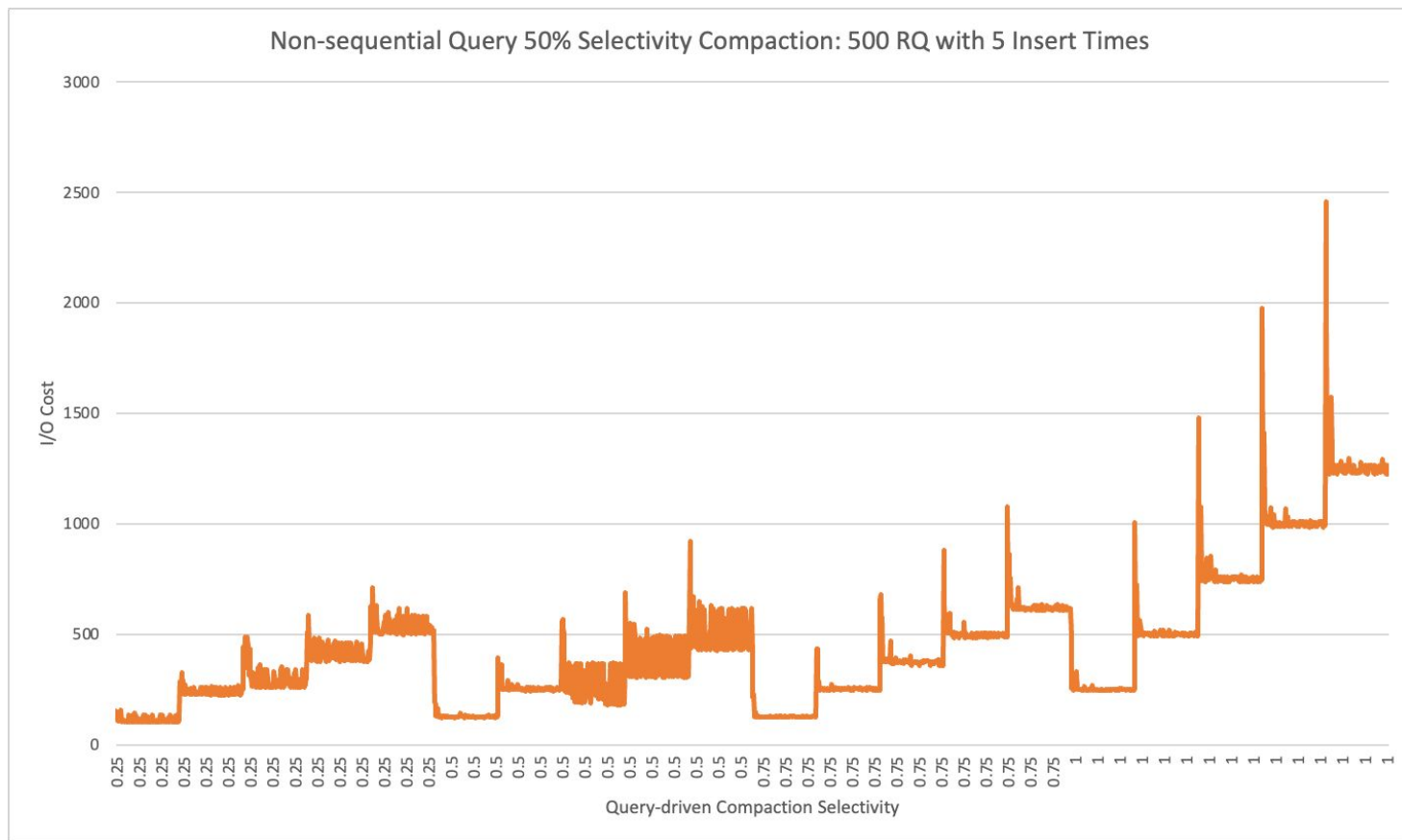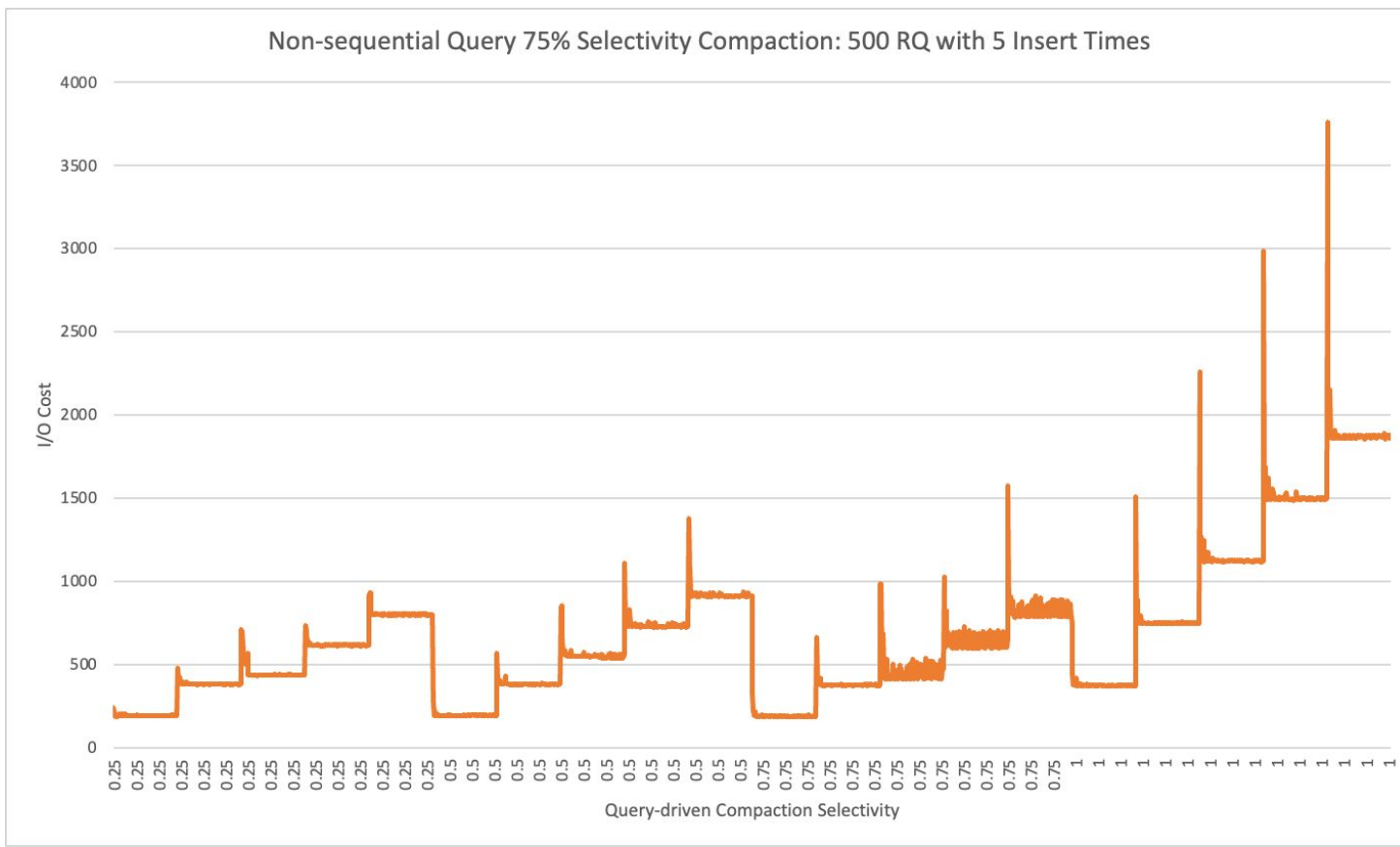Non-sequential Query 50% Selectivity Compaction: 500 RQ with 5 Insert Times

# Compaction Selectivity: Non-sequential Query 75% - 500 RQ - 5 Times

# Compaction Selectivity: Non-sequential Query 95% - 500 RQ - 5 Times



Non-sequential Query 95% Selectivity Compaction: 500 RQ with 5 Insert Times

# Compaction Selectivity: Non-sequential Query: 500 RQ - 5 Insert - Sum i/o

# Compaction Selectivity: Non-sequential Query 1% - 500 RQ - 10 Times

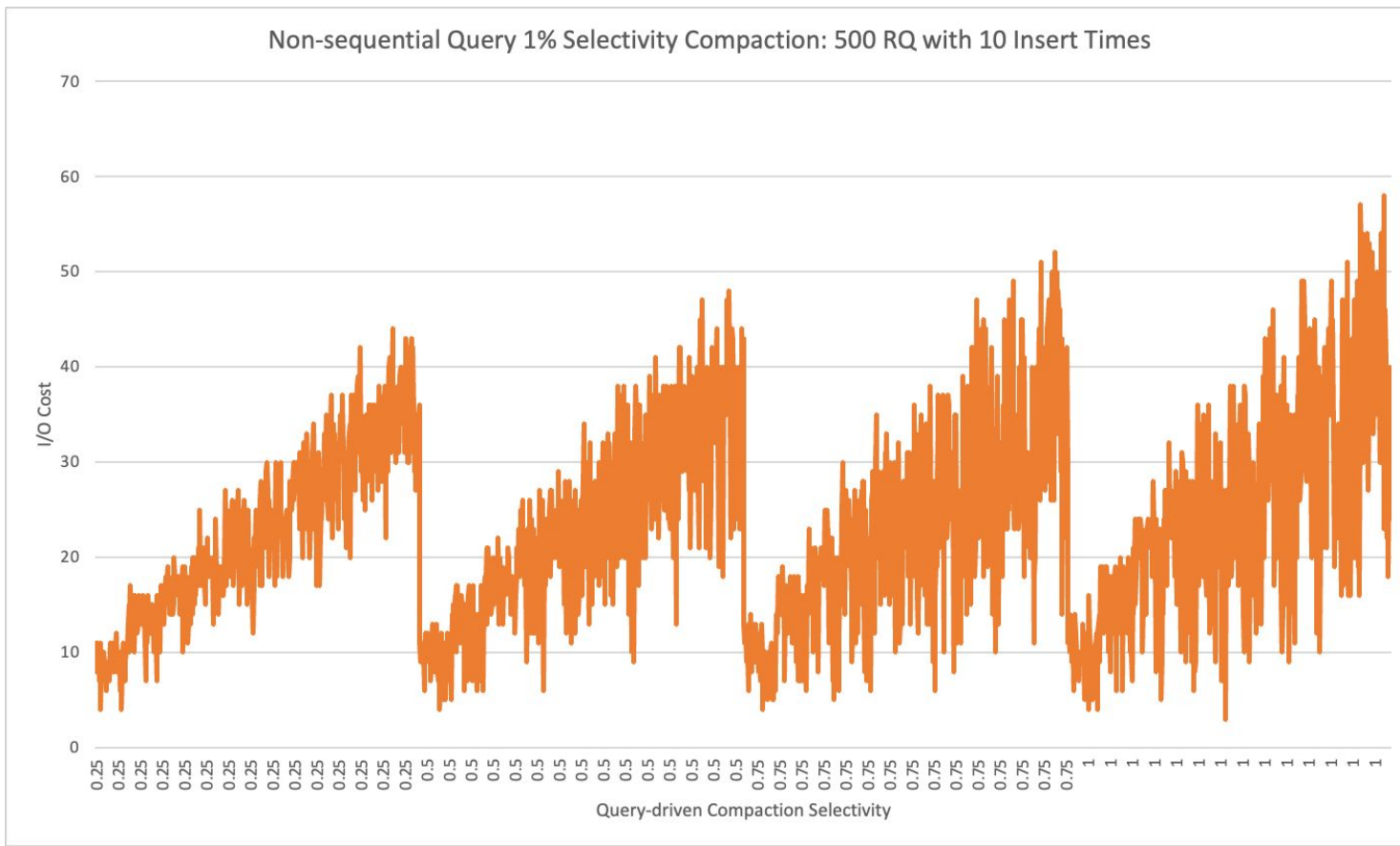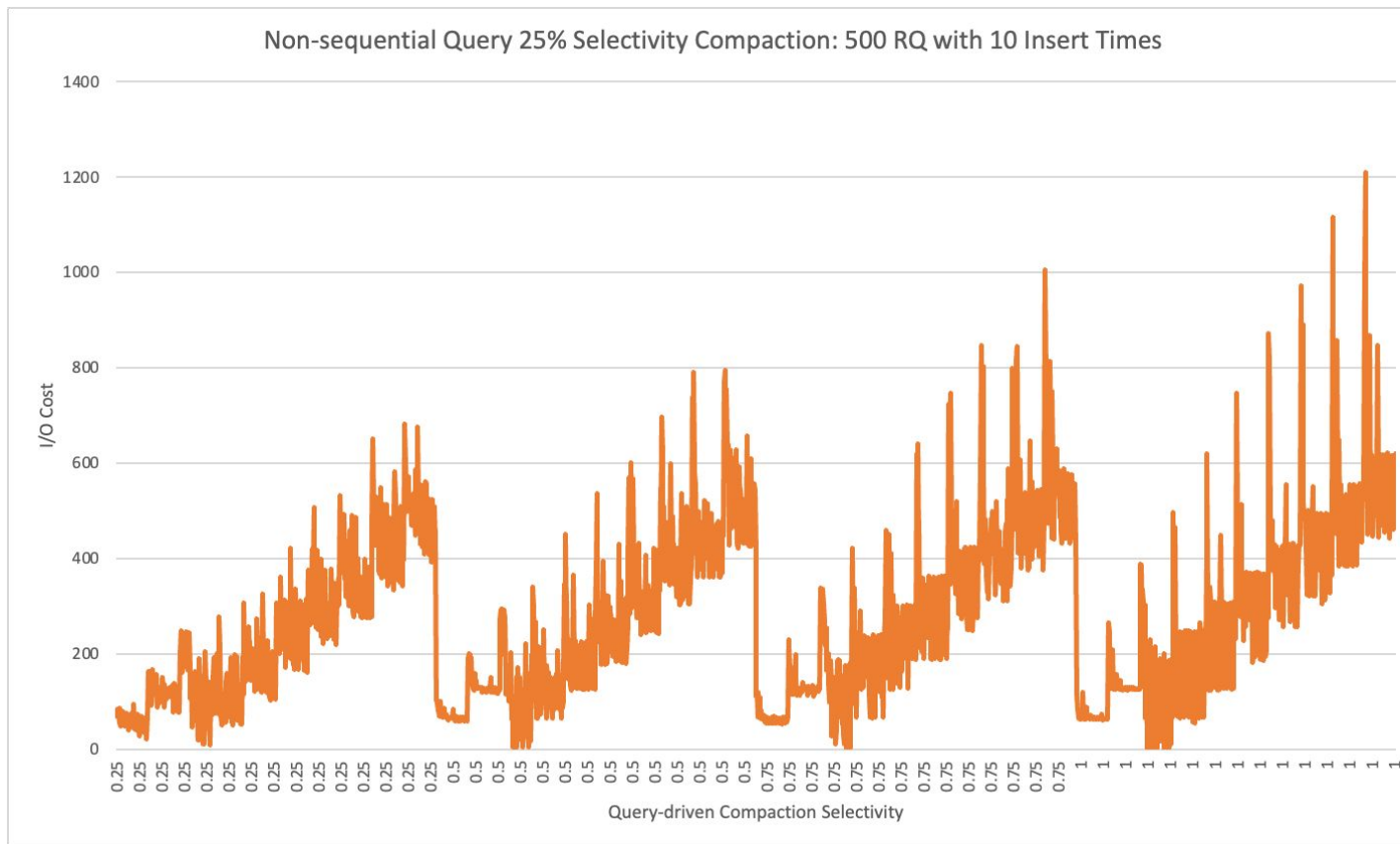# Compaction Selectivity: Non-sequential Query 25% - 500 RQ - 10 Times



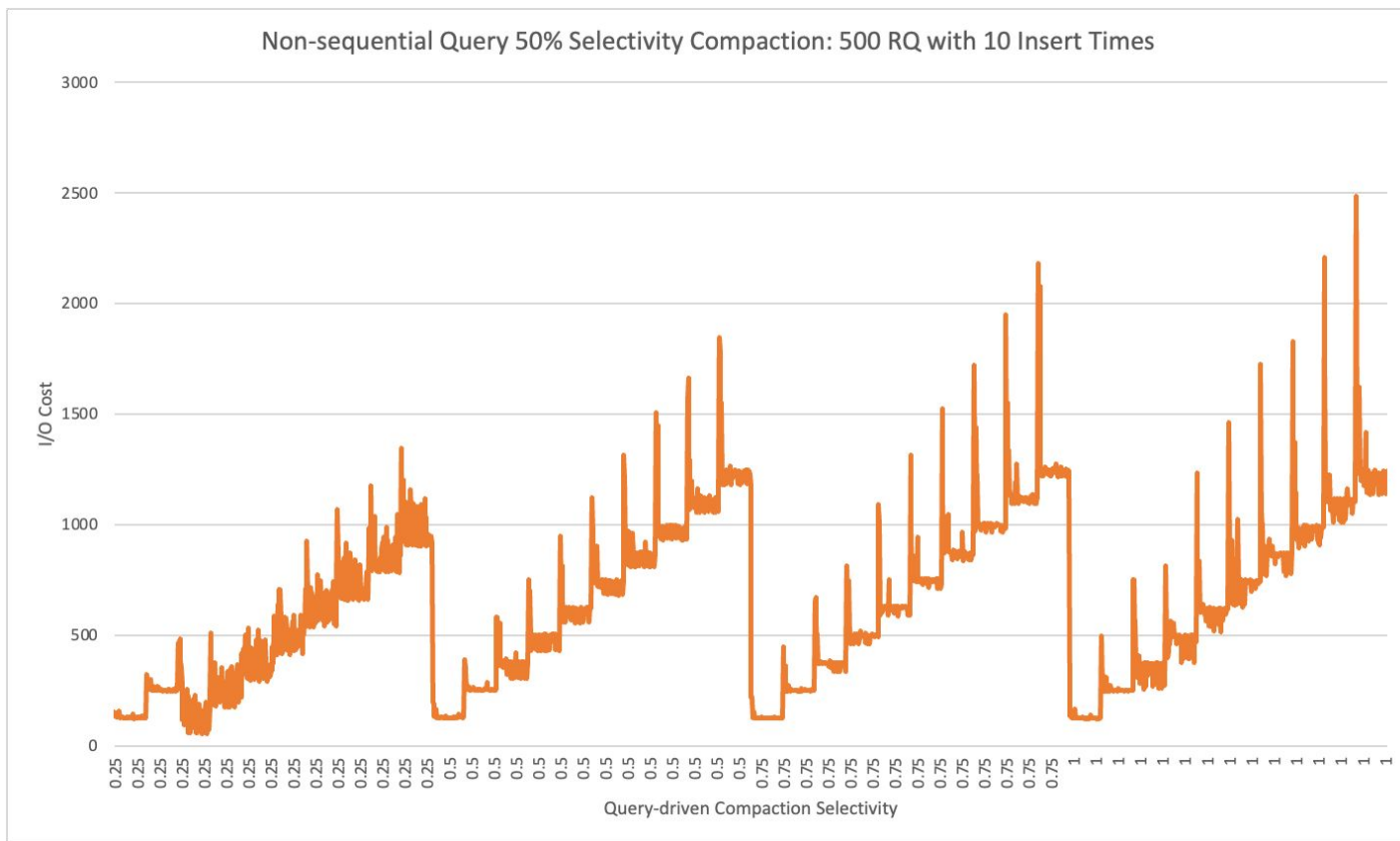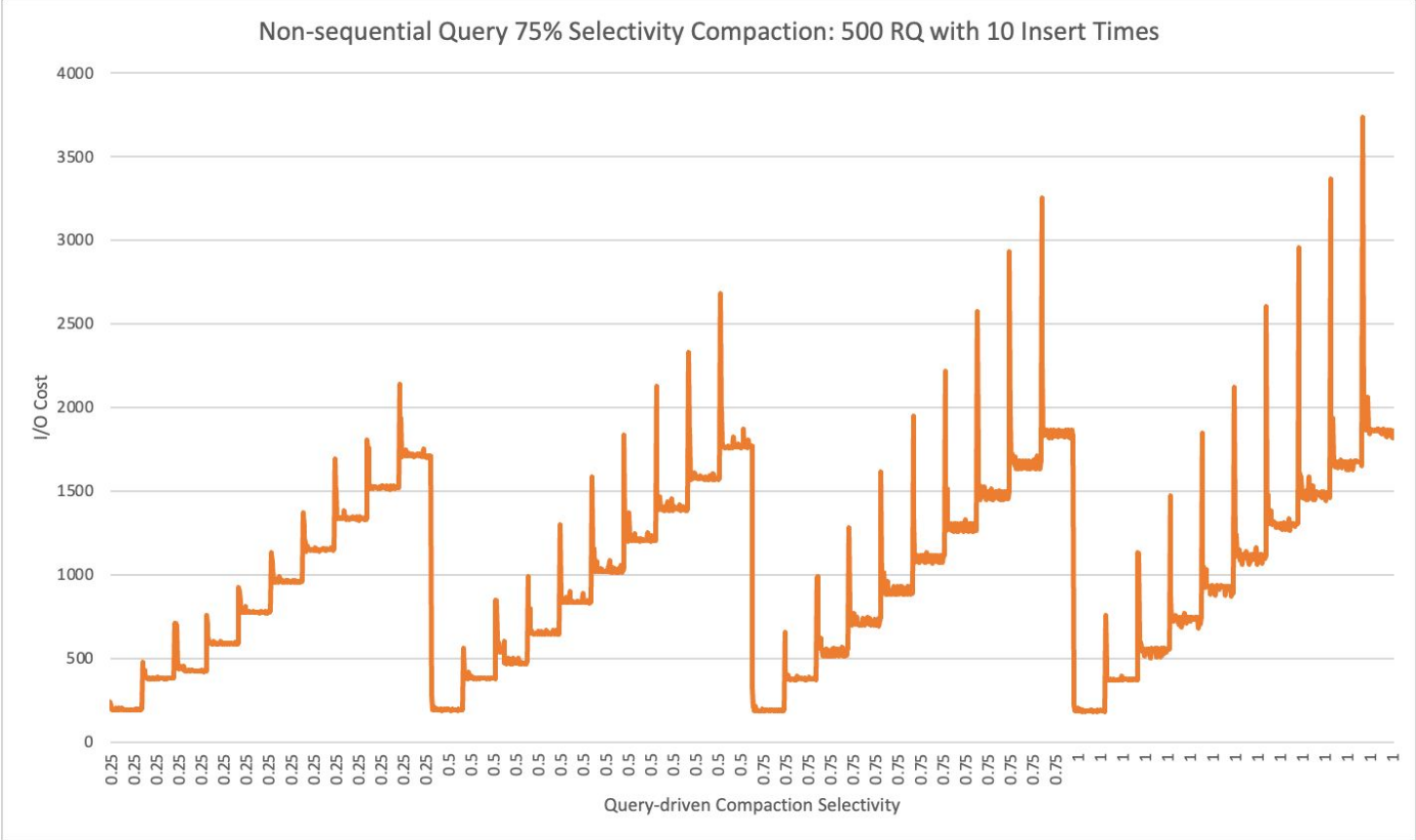Non-sequential Query 25% Selectivity Compaction: 500 RQ with 10 Insert Times

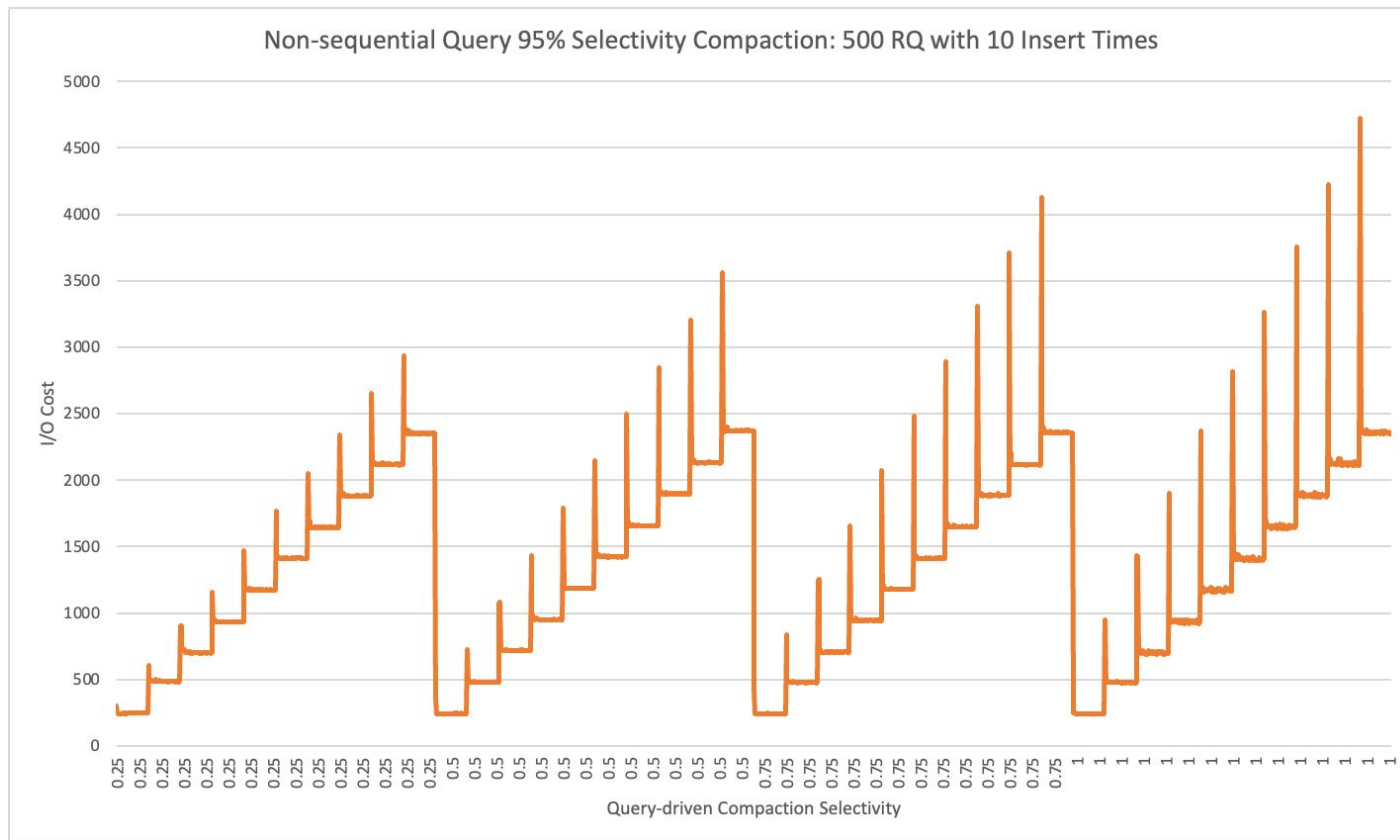# Compaction Selectivity: Non-sequential Query 50% - 500 RQ - 10 Times

# Compaction Selectivity: Non-sequential Query 75% - 500 RQ - 10 Times



Non-sequential Query 75% Selectivity Compaction: 500 RQ with 10 Insert Times

# Compaction Selectivity: Non-sequential Query 95% - 500 RQ - 10 Times



Non-sequential Query 95% Selectivity Compaction: 500 RQ with 10 Insert Times

# Compaction Selectivity: Non-sequential Query: 500 RQ - 10 Insert - Sum i/o

No vanilla data

# Sequential Query: 500 Point Query

No compaction data