

Department of Informatics, University of Zürich

BSc Thesis

An Adaptive Index for Hierarchical Database Systems

Rafael Kallis

Matrikelnummer: 14-708-887

Email: rk@rafaelkallis.com

February 1, 2018

supervised by Prof. Dr. Michael Böhlen and Kevin Wellenzohn



University of
Zurich^{UZH}

Department of Informatics



Abstract

The workload-aware property index is a hierarchical index which adapts to the database's recent transactional workload by not pruning volatile index nodes, i.e. nodes which are frequently inserted or deleted, in order to increase update performance. When the workload changes, these nodes cease to be volatile and become unproductive if they and their descendants, neither contribute to a query match, nor are volatile.

Unproductive nodes in hierarchical indexes waste space and slow down queries. We propose periodic Garbage Collection and Query-Time Pruning in order to clean unproductive nodes in the index. We implement the techniques in Apache Jackrabbit Oak and provide an extensive experimental evaluation to stress test the algorithms and show that the database throughput increases considerably when periodic Garbage Collection or Query-Time Pruning are applied.

Zusammenfassung

Der “workload-aware property index” ist ein hierarchischer Index, der sich an die jüngste Transaktionslast der Datenbank anpasst, indem er volatile Indexknoten, d.h. Knoten die häufig eingefügt oder gelöscht werden, nicht löscht um die Schreibleistung zu erhöhen. Wenn sich die Arbeitslast ändert, sind diese Knoten nicht mehr volatil und werden unproduktiv, falls sie und ihre Nachkommen weder zu einer Abfrage beitragen noch volatil sind.

Unproduktive Knoten in hierarchischen Indizes verschwenden Speicherplatz und verlangsamen die Abfragen. Wir schlagen periodische Indexreinigung und Abfragezeitbereinigung vor, um den Index von unproduktiven Knoten zu bereinigen. Wir implementieren unsere Techniken in Apache Jackrabbit Oak und bieten eine umfangreiche experimentelle Auswertung, um die Algorithmen unter hoher Last zu testen und zeigen, dass der Datenbankdurchsatz erheblich zunimmt, wenn periodische Indexreinigung oder Abfragezeitbereinigung benutzt wird.

Contents

1	Introduction	6
2	Background	7
2.1	CMS Workload	7
2.2	Workload Aware Property Index (WAPI)	7
3	Unproductive Nodes	11
3.1	Introduction	11
3.2	Impact on Query Runtime	12
4	Cleaning Unproductive Nodes	15
4.1	Periodic Garbage Collection (GC)	15
4.2	Query-Time Pruning (QTP)	19
5	Experimental Evaluation	22
5.1	Goals	22
5.2	Preliminaries	22
5.2.1	Setup	22
5.2.2	Datasets	22
5.2.3	Workload	23
5.3	Unproductive Nodes	25
5.3.1	Volatility Threshold τ	25
5.3.2	Sliding Window of Length L	27
5.3.3	Workload Skew s	28
5.3.4	Update Operations per Second	30
5.4	Garbage Collection	31
5.4.1	GC Period T	33
5.5	Query-Time Pruning	35
5.6	Comparison	37
5.7	Summary	39
6	Conclusion	40
7	Appendix	42

List of Figures

1	An instance of a hierarchical database	8
2	Volatile nodes becoming unproductive	10
3	Query Runtime over the accumulated number of update operations . . .	12
4	Index composition during query execution	13
5	Node Ratio during query execution	14
6	Worst case scenario during naïve GC	16
7	Postorder tree walk	17
8	GC applied on Oak	17
9	Periodic GC implemented in Java	18
10	QTP applied on Oak	20
11	QTP implemented in Java	21
12	Excerpts of both datasets visualized	22
13	CDF for the Zipf distribution for both datasets	24
14	CMS Workloads visualized	24
15	Unproductive Nodes over accumulated update operations with $\tau \in \{1, 5, 10\}$	25
16	Unproductive Nodes over volatility threshold τ	25
17	Unproductive Nodes over update operations with lengths $L \in \{10s, 20s, 30s\}$	27
18	Unproductive Nodes over sliding window length L	28
19	Hotspots affected by different skew values	28
20	Unproductive Nodes over update operations with skew $s \in \{1, \frac{3}{2}, 2\}$. . .	29
21	Unproductive Nodes over skew s	29
22	Unproductive Nodes over time for 38, 46 and 57 updates per second . . .	30
23	Unproductive Nodes over updates per second	31
24	Query Runtime over accumulated update operations, periodic GC enabled	31
25	Index Nodes over accumulated update operations, periodic GC enabled .	32
26	Avg. Query Runtime over GC Period T	33
27	Unproductive Nodes over GC period T	34
28	Unproductive Nodes over update operations with GC period $T \in \{5s, 50s\}$	34
29	Query Runtime over accumulated update operations, QTP enabled . . .	35
30	Index Nodes over accumulated update operations, QTP enabled	36
31	QTP overhead	36
32	Periodic GC vs. QTP	37
33	<code>postOrder()</code> implementation in Java	42
34	<code>map()</code> implementation in Java	43
35	<code>filter()</code> implementation in Java	44

1 Introduction

Frequently adding and removing data from hierarchical indexes causes them to repeatedly grow and shrink. A single insertion or deletion can trigger a sequence of structural index modifications (node insertions/deletions) in a hierarchical index. Skewed and update-heavy workloads trigger repeated structural index updates over a small subset of nodes to the index. Informally, a frequently added or removed node is called *volatile*. Volatile nodes decrease index update performance due to two reasons. First, frequent structural index modifications are expensive since they cause many disk accesses. Second, frequent structural index modifications also increase the likelihood of conflicting index updates by concurrent transactions. Conflicting index updates further decrease update performance since concurrency control protocols need to resolve the conflict.

Wellenzohn et al. [6] propose the Workload-Aware Property Index (WAPI). The WAPI exploits the workload’s skewness by identifying and not removing volatile nodes from the index, thus significantly reducing the number of expensive structural index modifications. Since fewer nodes are inserted/deleted, the likelihood of conflicting index updates by concurrent transactions is reduced.

Some Content Management Systems (CMS) make use of hierarchical databases. The Adobe Experience Manager,¹ Adobe’s enterprise CMS, works with the hierarchical database system Apache Jackrabbit Oak (Oak). CMSs yield skewed, update-heavy and changing workloads. They frequently update a small changing subset of index nodes. Such workloads decrease WAPI’s query performance.

When the workload characteristics change, new index nodes can become volatile while others cease to be volatile and become *unproductive*. Unproductive index nodes slow down queries, as traversing an unproductive node is useless because unproductive nodes do not contain any data and thus cannot yield a query match. Additionally, unproductive nodes occupy storage space that could otherwise be reclaimed. If the workload changes frequently, unproductive nodes accumulate in the index and the query performance deteriorates over time. Therefore, unproductive nodes must be cleaned to keep query performance stable over time and reclaim disk space as the workload changes.

Wellenzohn et al. [6] propose periodic Garbage Collection (GC), which traverses the entire index subtree and prunes all unproductive index nodes at once. Additionally we propose Query-Time Pruning (QTP), an incremental approach to cleaning unproductive nodes in the index. The idea is to turn queries into updates. Since Oak already traverses unproductive nodes as part of query processing, these nodes could be pruned at the same time. With QTP, only one query has to traverse an unproductive node, while subsequent queries can skip this overhead and thus perform better. The goal of this BSc thesis is to study, implement and empirically compare GC and QTP as proposed by [6] in the open-source hierarchical database Apache Jackrabbit Oak.

¹<http://www.adobe.com/marketing-cloud/experience-manager.html>

2 Background

2.1 CMS Workload

Depending on the data model, applications might use a hierarchical [4] database system such as Apache Jackrabbit Oak, called Oak. A content management system (CMS) might choose Oak because the database reflects the hierarchical structure of a webpage. Content management systems have specific workloads. These workloads have distinct properties: they are *skewed*, *update-heavy* and *changing* [6]. CMSs frequently use a job-queuing system that has the noted characteristics.

Consider a social media feed as a running example. Only a few posts are popular. These posts have many comments or likes. Since most of the interactions (comments, likes) are on a small subset of the posts, we have a skewed workload. Users submit new posts or interact with existing posts by writing comments for example, creating an update-heavy workload. As time passes, new posts are created. Users are more likely to interact with recent posts than older ones, hence the workload changes over time.

When a user submits a new post, a job is sent to the CMS for processing. A background thread is periodically checking for pending jobs. A pending job is signaled using node properties in Oak. The CMS adds a property to the respective node in order to signal the background thread that the specific node has a pending job requiring processing. A node’s “pub” property signals when the job has to be processed. Setting the value of “pub” to “now” indicates that the job must be processed immediately. Once the background thread detects the node, the job is processed and the “pub” property is removed.

From now on, we shall refer to a workload with the properties mentioned above as a *CMS workload*. When a hierarchical database operates under a CMS workload, we can increase its update performance using the workload aware property index.

2.2 Workload Aware Property Index (WAPI)

Oak mostly executes content-and-structure (CAS) queries [5]. We denote node n ’s property k as $n[k]$ and node n ’s descendants as $desc(n)$.

Definition 1 (CAS Query). Given node m , property k and value v , a CAS query $Q(k, v, m)$ returns all descendants of m which have k set to v , i.e.,

$$Q(k, v, m) = \{n | n \in desc(m) \wedge n[k] = v\}$$

Example 1 (CAS Query). Consider Figure 1. CAS-Query $Q(\text{pub}, \text{now}, /a)$, which queries for every descendant of $/a$ with “pub” set to “now”, would evaluate to $Q(\text{pub}, \text{now}, /a) = \{/a/b/d\}$, since $/a/b/d$ is the only descendant of $/a$ with “pub” set to “now”.

Figure 1 depicts a database instance with the workload aware property index (WAPI). The WAPI is a hierarchical index which indexes the properties of nodes in order to answer

should be close to zero. We infer that, an index subtree which resembles a job-queuing system under a CMS workload, should have almost no matching nodes.

The CMS workload is also skewed and update-heavy, therefore causing repeated structural index updates (insertions/deletions) over a small subset of nodes to the index. The WAPI takes into account if an index node is frequently added and removed, i.e. *volatile* (see Definition 4), before performing structural index modifications. If a node is considered volatile, we do not remove it from the index.

Volatility is the measure which is used by the WAPI in order to distinguish whether to remove a node or not from the index. Wellenzohn et al. [6] propose to look at the recent transactional workload to check whether a node n is volatile. The workload on Oak instance O_i is represented by a sequence $H_i = \langle \dots, G^a, G^b, G^c \rangle$ of snapshots [4], called a *history*. A snapshot represents an immutable committed tree of the database. Let t_n be the current time and $t(G^b)$ be the point in time snapshot G^b was committed, $N(G^a)$ is the set of nodes which are members of snapshot G^a . We use a superscript a to emphasize that a node n^a belongs to tree G^a . $pre(G^b)$ is the predecessor of snapshot G^b in H_i .

Given two snapshots G^a and G^b , we write n^a and n^b to emphasize that nodes n^a and n^b are two versions of the same node n , i.e., they have the same absolute path from the root node.

Node n is volatile iff n 's volatility count is at least τ , called volatility threshold. The volatility count of n is defined as the number of times n was added or removed from snapshots in a sliding window of length L over history H_i .

Definition 3 (Volatility Count). The volatility count $vol(n)$ of index node n on Oak instance O_i , is the number of times node n was added or removed from snapshots contained in a Sliding Window of Length L over history H_i .

$$vol(n) = |\{G^b | G^b \in H_i \wedge t(G^b) \in [t_n - L + 1, t_n] \wedge \exists G^a [\\ G^a = pre(G^b) \wedge ([n^a \notin N(G^a) \wedge n^b \in N(G^b)] \vee \\ [n^a \in N(G^a) \wedge n^b \notin N(G^b)])]\}|$$

Definition 4 (Volatile Node). Index node n is volatile iff n 's volatility count (see Definition 3) is greater or equal than the volatility threshold τ , i.e.,

$$volatile(n) \iff vol(n) \geq \tau$$

Example 3 (Volatile Node). Consider the snapshots depicted in Figure 2. Assume volatility threshold $\tau = 1$, sliding window of length $L = 2$ and history $H_h = \langle G^0, G^1, G^2, G^3 \rangle$. Oak instance O_h executes transactions T_1, \dots, T_3 . Note that volatile index nodes are color-coded blue in Figure 2. Snapshot G^0 was committed at time $t(G^0) = t$. Given snapshot G^0 , transaction T_1 adds the property-value pair “pub” = “now” to /a/b/d and commits snapshot G^1 at time $t(G^1) = t+1$. Next, transaction T_2 removes property “pub”

from $/a/b/d$ given snapshot G^1 and commits snapshot G^2 at time $t(G^2) = t + 2$. During T_2 we have $t_n = t + 2$ and the sliding window is $[t_n - L + 1, t_n] = [t + 2 - 2 + 1, t + 2] = [t + 1, t + 2]$. Snapshot G^1 is in the sliding window since $t(G^1) = t + 1$ and $t + 1 \in [t + 1, t + 2]$. Snapshot G^0 is not in the sliding window because $t(G^0) = t$ and $t \notin [t + 1, t + 2]$. Index node $n = /i/pub/now/a/b/d$ exists in the set of nodes which are members of snapshot G^1 , $n^1 \in N(G^1)$ but not in its predecessor $pre(G^1) = G^0$, $n^0 \notin N(G^0)$ and therefore has a volatility count of $vol(n) = 1$. Since the threshold is $\tau = 1$, the node is volatile. The same holds for n 's ancestors within the index subtree.

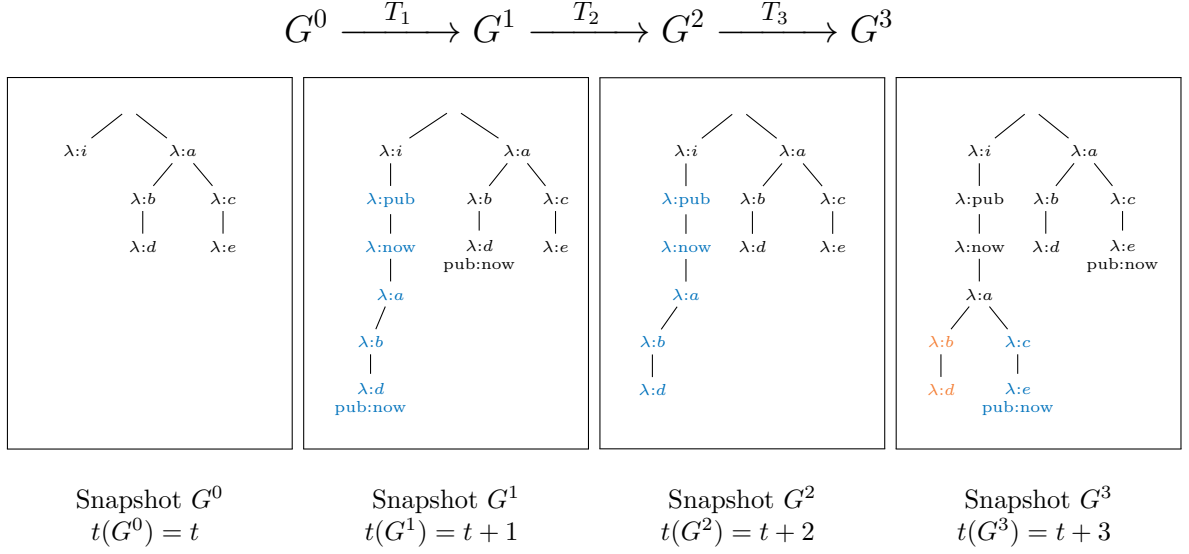


Figure 2: Volatile nodes becoming unproductive. Given $\tau = 1$, $L = 2$, nodes $/i/pub/now/a/b/d$ and $/i/pub/now/a/b$ are unproductive in snapshot G^3 . They are not volatile and don't match either. Note that volatile and unproductive index nodes are color-coded blue and orange, respectively.

3 Unproductive Nodes

3.1 Introduction

When time passes and the database workload changes, volatile nodes cease to be volatile and they become unproductive. When nodes are volatile, their volatility count has to be at least τ . When time passes, insertions and deletions that increased the volatility count drop out of the sliding window, causing the volatility count to decrease. If the volatility count drops below threshold τ , the node ceases to be volatile. If the now non-volatile node is also non-matching, and the same holds for its descendants, we call the node and its descendants unproductive.

Unproductive index nodes slow down queries as traversing an unproductive node is useless, because neither the node itself nor any of its descendants are matching and thus cannot yield a query match. Additionally, unproductive nodes occupy storage space that could otherwise be reclaimed. If the workload changes frequently, unproductive nodes accumulate in the index and the query performance deteriorates over time [6].

Definition 5 (Unproductive Node). Index node n is unproductive iff n , and any descendant of n , is neither matching (see Definition 2) nor volatile (see Definition 4), i.e.,

$$\text{unproductive}(n) \iff \forall m(m \in (\{n\} \cup \text{desc}(n)) \wedge \neg \text{matching}(m) \wedge \neg \text{volatile}(m))$$

Example 4 (Unproductive Node). In our running example (cf. Figure 2), transaction T_3 adds the property-value pair “pub” = “now” to $/a/c/e$ to G^2 and commits G^3 at time $t(G^3) = t+3$. We assume the same parametrization as in the last example ($\tau = 1, L = 2$). Volatile and unproductive index nodes are color-coded blue and orange, respectively. In snapshot G^3 , index nodes $/i/pub/now/a/b/d$ and $/i/pub/now/a/b$ cease to be volatile because their volatility counts drop below the threshold. The sliding window has length $L = 2$, so we only consider snapshots G^2, G^3 towards the volatility count. The two nodes were not inserted or deleted in any of the mentioned snapshots and therefore the volatility count is $\text{vol}(/i/pub/now/a/b/d) = \text{vol}(/i/pub/now/a/b) = 0$. Since the threshold is $\tau = 1$, the nodes are not volatile. In addition, they are not matching either, therefore they are unproductive (cf. Definition 5). Index node $/i/pub/now/a$ in snapshot G^3 has two volatile descendants ($/i/pub/now/a/c/e, /i/pub/now/a/c$) one of which is matching and therefore is not unproductive.

In the example above, we saw how index nodes become unproductive after being volatile. Index nodes do not necessarily have to be volatile before they become unproductive. If a non-volatile and non-matching index node with no matching descendants has a volatile descendant, and that descendant ceases to be volatile, the index node becomes unproductive, even though it was not volatile.

Example 5 (CAS Query with Unproductive Nodes). Consider CAS-Query $Q(\text{pub}, \text{now}, /a)$ from Example 1 again. We apply the query to snapshot G^3 in Figure 2. The query executor has to traverse the four descendants of $/i/pub/now/a$. Traversing nodes $/i/pub/now/a/b/d$ and $/i/pub/now/a/b$ is useless and slows down the query because these nodes are unproductive. The query evaluates to $Q(\text{pub}, \text{now}, /a) = \{/a/c/e\}$.

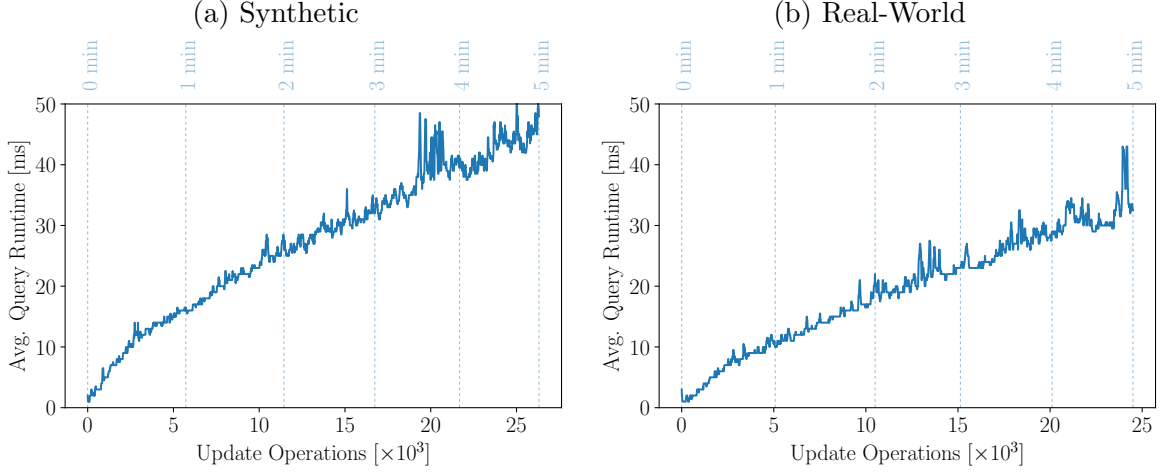


Figure 3: Query Runtime over the accumulated number of update operations.

3.2 Impact on Query Runtime

In this section we study and quantify the impact of unproductive nodes on query runtime. During query execution, traversing an unproductive node is useless, because neither the node itself nor any of its descendants are matching and therefore cannot contribute a query match. We hypothesize that unproductive nodes significantly slow down queries under a CMS workload. An index under a CMS workload (see Section 2.1) is dominated by unproductive nodes, that is unproductive nodes constitute a large percentage of all index nodes.

In order to find supporting evidence for the hypotheses above, we conduct a series of experiments on Oak. The setup of the experimental evaluation and datasets will be described in detail in Section 5. We record the query runtime throughout the experiment and present the data below.

Figures 3a and 3b show the query runtime of the same query over the accumulated number of update operations for the synthetic and real-world dataset, respectively. To provide a better sense of time to the reader, we add a vertical line for each passed minute during the simulation. Each point corresponds to the running median [3] over 10 time points. We observe a sublinear increase of the runtime from $2ms$ to $50ms$ after running the simulation for 5 minutes on the synthetic dataset and an increase from $2ms$ to $35ms$ on the real-world dataset.

Next, we present data regarding the type of index nodes traversed during query execution. Figures 4a and 4b depict the total number of traversed nodes in addition to the number of traversed volatile and unproductive nodes during query execution for each dataset.

The total number of traversed nodes is increasing sublinearly over time. This explains the increase in query runtime in Figures 3a and 3b. We believe the index becomes *static* over time. As time passes, more and more content nodes are randomly drawn by the CMS workload and their corresponding index nodes become volatile. After some time

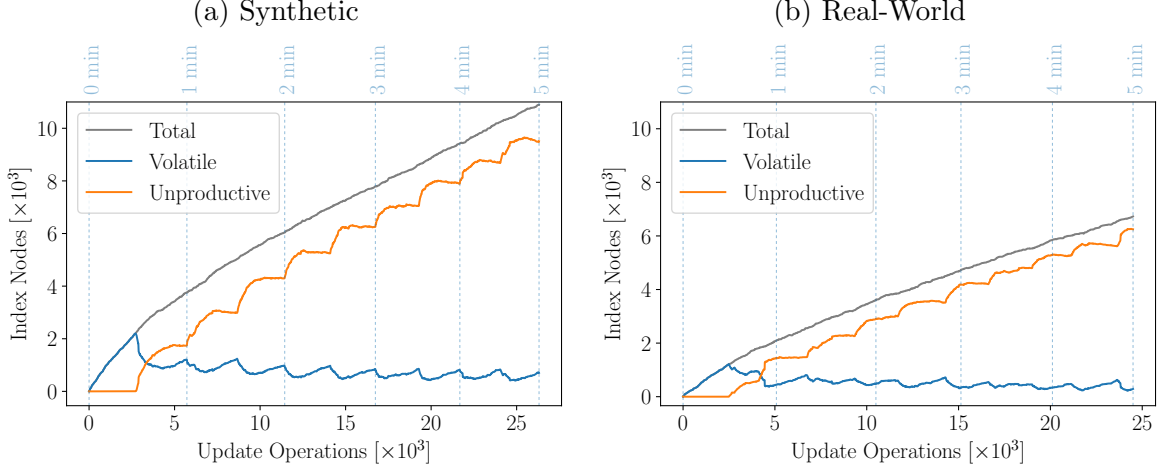


Figure 4: Index composition during query execution.

these nodes, most likely become unproductive and are not being pruned from the index anymore. Therefore, the probability of picking a content node with no corresponding index node (non-indexed) decreases over time. Since it becomes less and less likely for a non-indexed content node to be randomly picked by the CMS workload, the rate of growth of total traversed index nodes decreases over time. We expect that if we change the workload sufficiently often, the number of total index nodes converges to the number of content nodes.

Furthermore, we see the number of volatile nodes have a downward trend from the 30 second mark till the end of the experiment. We believe it becomes less likely for nodes to become volatile if unproductive nodes are not pruned, as time passes. When a workload randomly picks a content node whose index node is unproductive, the index node becomes matching but was not physically added, thus the volatility count of the index node does not increment. In comparison, if the workload randomly picks a non-indexed content node, the corresponding index node’s volatility count increments. We infer that it is less likely for an unproductive node to become volatile than a non-indexed one. Our experimental evaluation suggests that the number of unproductive nodes increases over time. Therefore it becomes less likely for any node to become volatile over time. This explains the downward trend of volatile nodes over time.

The sliding window of length L is set to 30 seconds, therefore we encounter no unproductive nodes during the first 30 seconds of the simulation. Once we reach the 30 second mark, the query executor encounters unproductive nodes. From that point, we observe a steep increase in traversed unproductive nodes. After 1 minute, we observe that the total traversed nodes are dominated by unproductive nodes. The rate of growth of traversed unproductive nodes seems to decrease over time. When content nodes are picked by the workload, their corresponding index nodes are inserted into the WAPI, where some become volatile. When the workload changes, these volatile index nodes most likely become unproductive and accumulate in the index, hence the number of non-indexed content nodes decreases. Since fewer non-indexed content nodes are available to become

volatile and thereafter unproductive, the rate of growth of traversed unproductive nodes decreases.

Additionally, we observe the functions of the unproductive and volatile index nodes having the shape of a cycloid. In Figures 4a and 4b, each cusp [1] ($30s, 60s, \dots, 300s$) represents a point in which the workload changes during the experiment. When the workload changes, we initially see a steep decrease in volatile nodes. During that phase, more nodes cease to be volatile than become volatile. Nodes need to reach the threshold in order to become volatile and few do, since the skew in the workload only picks a subset of nodes frequently. Before a new workload kicks in, we observe the opposite phenomenon. More nodes become volatile than cease to be volatile, because many nodes are on the verge of becoming volatile and therefore need to be picked only a few times more to become volatile.

Lastly, we also observe the real-world dataset having a more gentle slope over the synthetic dataset. Since the real-world dataset has more content nodes, it is less likely for each content node to be picked by the workload. Having a smaller chance to be picked by the workload implies fewer volatile and consequently fewer unproductive nodes in the index.

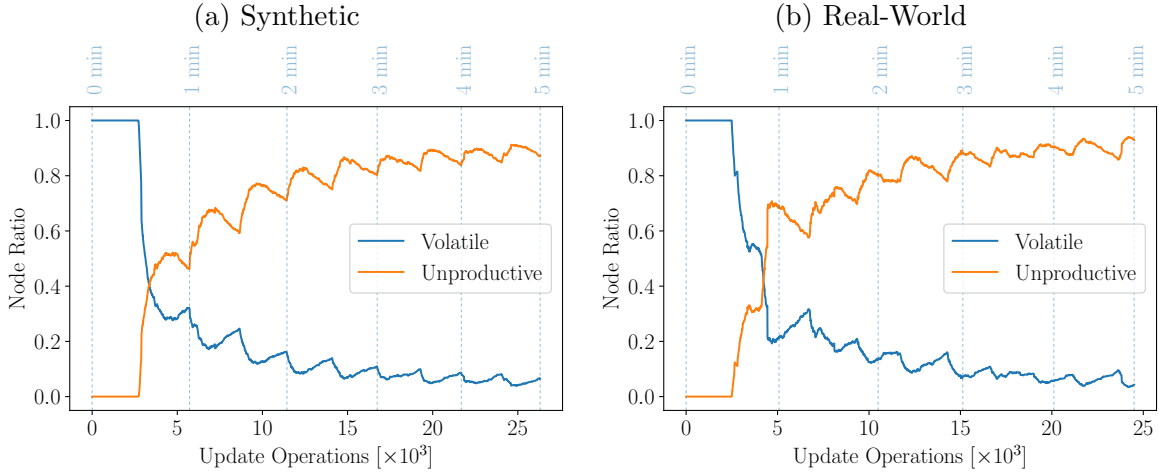


Figure 5: Node Ratio during query execution.

Figures 5a to 5b show the ratio of 1) volatile over total and 2) unproductive over total index nodes, over the accumulated number of update operations for each dataset. These Figures quantify how strongly unproductive nodes dominate the total traversed nodes. After 5 minutes, unproductive nodes account for over 80% of the traversed nodes whilst less than 20% are volatile on the synthetic dataset. Similarly, on the real-world dataset, we observe 90% traversed unproductive and 10% traversed volatile index nodes.

Concluding, the experiments strongly support our hypothesis. The query runtime increases by an order of magnitude after 5 minutes. Also, unproductive nodes, which dominate the index, are accountable for the increase in query runtime.

4 Cleaning Unproductive Nodes

In the previous section, we saw how unproductive index nodes slow down query execution. To prevent unproductive nodes from accumulating in the index, we need strategies to clean up the index. In the following two subsections, we suggest two different approaches for dealing with unproductive nodes. We will empirically investigate their query performance in Sections 5.4 and 5.5.

4.1 Periodic Garbage Collection (GC)

First, we propose to clean the index periodically with a garbage collector. We add a background process that periodically executes garbage collection of unproductive nodes wrapped inside a single Oak transaction.

The naïve approach for garbage collection is to traverse the entire index subtree, apply Definition 5 to each visited node n , and delete n if it is unproductive. Deciding if n is unproductive, requires us to check that no descendant $desc(n)$ of n is matching or volatile. Checking the descendants of each index node n causes naïve GC to have a quadratic time complexity.

Example 6 (Naïve GC). Assume we apply naïve GC to the index subtree depicted in Figure 6. The subtree resembles the worst-case scenario for GC. It contains n nodes, no node is unproductive because the leaf node is matching and each internal node has a single child. In order to determine if root node a_1 is unproductive, we have to traverse each descendant of a_1 and check if it is matching or volatile. The same holds for determining if any of a_1 's descendants is unproductive. In this example, a node a_k has exactly $n - k$ descendants and to check if a_k is unproductive, these $n - k$ descendants need to be traversed.

To determine if the n nodes in the subtree rooted at a_1 are unproductive, we need to traverse the following number of nodes:

$$\begin{aligned}
 & (n - 1) + (n - 2) + \cdots + (n - n - 1) + (n - n) \\
 &= \sum_{k=1}^n n - k \\
 &= \sum_{k=1}^n n - \sum_{k=1}^n k \\
 &= n^2 - \frac{n(n + 1)}{2} \\
 &= n^2 - \frac{n^2 + n}{2} \\
 &= \frac{n^2 - n}{2} \\
 &= \Theta(n^2)
 \end{aligned}$$

We showed that applying naïve GC on an index subtree of n nodes is bound by $\Theta(n^2)$ in the worst case.

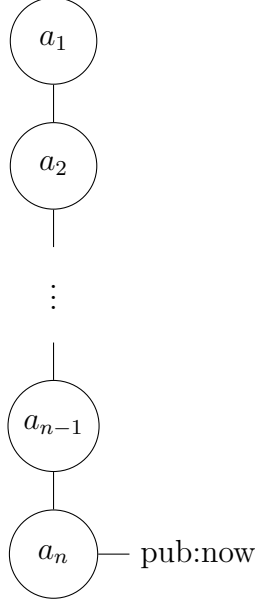


Figure 6: Worst case scenario during naïve Garbage Collection. All internal nodes have a single child. The leaf node is matching, therefore no node is unproductive.

By applying a *postorder tree walk* [2], we can garbage collect the index subtree in linear time complexity. The postorder tree walk allows us to process a node n 's descendants before n . When visiting node n during a postorder tree walk, any unproductive descendant of n was already pruned, hence no child of n can be unproductive. Thus, if n has at least one child, n cannot be unproductive and therefore checking if n has at least one child, in addition to checking if n is matching or volatile, is sufficient to determine if n is unproductive, when applying a postorder tree walk. Figure 7 depicts a postorder tree walk on snapshot G^3 from Figure 2. The numbers represent the order in which the nodes are checked and pruned if unproductive. A node's descendants are always evaluated before the node itself.

Algorithm 1 prunes all unproductive descendants of the index subtree rooted at $/i$. The algorithm traverses the subtree rooted at $/i$ using a postorder tree walk. If a descendant $n \in \text{desc}(/i)$ has no children and is neither matching, nor volatile, it is unproductive and therefore pruned from the index. If n has at least one child, we infer that n has at least one matching or volatile descendant, thus n cannot be unproductive. The postorder tree walk ensures that the algorithm prunes a child before its parent node. This guarantees that all unproductive nodes are pruned.

Algorithm 1: GarbageCollect

```

for node  $n \in \text{desc}(/i)$  in postorder tree walk do
  if  $\text{chd}(n) = \emptyset \wedge \neg \text{matching}(n) \wedge \neg \text{volatile}(n)$  then
    delete node  $n$ 

```

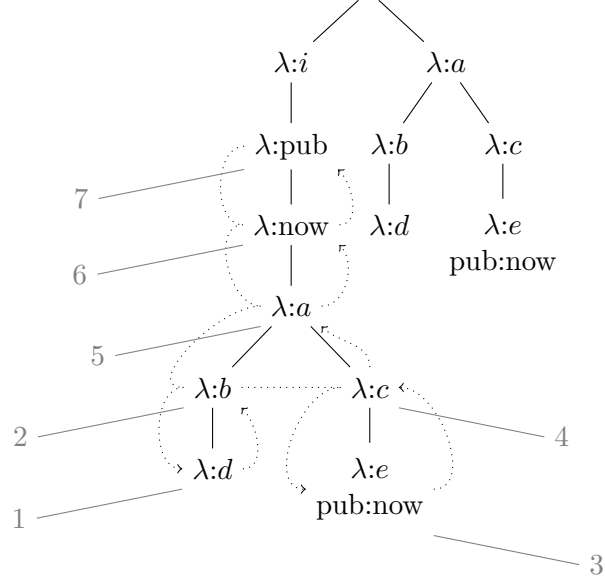


Figure 7: Postorder tree walk on the index subtree rooted at `/i` of G^3 . The numbers represent the order the corresponding nodes were visited, e.g. `/i/pub/now/a/b/d` was visited first, `/i/pub/now/a/b` second, etc.



Figure 8: Garbage Collection applied on Oak. Assume nodes `/i/pub/now/a/b/d` and `i/pub/now/a/b` are unproductive (color-coded orange) in snapshot G^3 . Garbage collection is executed during transaction T_4 . G^4 is the committed snapshot.

Example 7 (Periodic GC). Figure 8 shows snapshot G^3 from Figure 2. Index nodes `/i/pub/now/a/b/d` and `/i/pub/now/a/b` are unproductive in snapshot G^3 , as explained in Example 4. Oak’s background process executes garbage collection during T_4 . While executing GC, the index subtree of G^3 is traversed in postorder. The first unproductive node we visit and prune is `/i/pub/now/a/b/d`. Next, the garbage collector visits and prunes `/i/pub/now/a/b`. No further unproductive node is left for pruning. We see the cleaned index after garbage collection in snapshot G^4 .

```
/**
 * Removes any unproductive descendant from the index subtree.
 *
 * @param root: latest Oak snapshot
 */
void garbageCollect(Root root) {

    /* index subtree root */
    Tree indexRoot = root.getTree(OAK_INDEX_PATH);

    /* filter nodes which have children, are matching or volatile */
    for (Tree unproductiveNode : filter(
        (Tree n) -> n.getChildrenCount(1) == 0 &&
                    !isMatching(n) &&
                    !isVolatile(n),

        /* postorder tree walk iterable */
        postOrder(indexRoot)
    ) {
        unproductiveNode.remove();
    }
}
```

Figure 9: Periodic Garbage Collection implemented in Java.

Figure 9 shows the implementation of the periodic GC in Java inside Apache Jackrabbit Oak. Calling `postOrder()` returns a lazy sequence of nodes which correspond to a postorder tree walk of the index subtree (cf. Figure 33 in the Appendix). Next, using `filter()` we remove any node that has children, is matching or volatile from the sequence. All other nodes are unproductive and therefore pruned from the index.

By applying periodic GC on Oak, we introduce a new parameter. *GC Period T* defines how often garbage collection is run by the background process on Oak. If we pick a smaller period T , garbage collection is run more often and this reduces the number of unproductive nodes in the index. As a result, we increase the query performance of Oak, since the query executor has to traverse fewer unproductive nodes on average.

It is also worth mentioning that GC uses system resources in order to clean the index. If run too often, GC might degrade query performance, because the system is busy cleaning the index, instead of processing queries.

We suggest running GC when the system is not busy executing queries, e.g, every day during the early morning hours. Like so, garbage collection minimally interferes with other transactions. We will revisit GC period in Section 5.4.1, where we conduct an experiment in order to determine the impact of period T on unproductive nodes.

4.2 Query-Time Pruning (QTP)

While periodic garbage collection is explicitly executed by Oak in order to clean unproductive nodes, Query-Time Pruning is an approach which cleans unproductive nodes whilst Oak answers queries. Doing so, we benefit by avoiding the cost of explicitly traversing the index in comparison to GC.

Frequent queries also benefit from QTP. If a query is executed for the first time, QTP cleans all unproductive nodes traversed during query execution. If the query is executed a second time shortly thereafter and no new nodes become unproductive in the meantime, the query executor encounters no unproductive nodes anymore and therefore the query is answered faster.

If the path filter \mathbf{m} of CAS query $Q(\mathbf{k}, \mathbf{v}, \mathbf{m})$ (cf. Definition 1) changes often, the query executor starts traversing new subtrees in the index and consequently, some subtrees are not queried anymore. If an unproductive node is a member of a subtree that is not queried anymore, it is not pruned and remains in the index. Such unproductive nodes waste space because they contain no data but do not impact query performance since they are not traversed during query execution.

When using QTP to query for all descendants of a content node with path \mathbf{m} , we apply a postorder tree walk. With the postorder tree walk, we traverse and clean all unproductive nodes in the subtree of $/\mathbf{i}/\mathbf{k}/\mathbf{v}/\mathbf{m} = /\mathbf{i}/\mathbf{k}/\mathbf{v}/\lambda_1/\dots/\lambda_d$ in linear time complexity, as explained in Section 4.1. In the same Section, we mention that during a postorder tree walk, checking if node n has at least one child, in addition to checking if n is matching or volatile, is sufficient to determine if n is unproductive. The postorder tree walk ensures that the algorithm prunes a child before its parent node. This guarantees that all unproductive nodes are pruned in the subtree rooted at $/\mathbf{i}/\mathbf{k}/\mathbf{v}/\mathbf{m}$.

Algorithm 2 takes a CAS query $Q(\mathbf{k}, \mathbf{v}, \mathbf{m})$ as an argument, where \mathbf{k} is a property, \mathbf{v} a value and \mathbf{m} a content node's path. We initialize set r as the empty set. r will hold all content nodes satisfying the CAS query (cf. Definition 1). We traverse any descendant n of $/\mathbf{i}/\mathbf{k}/\mathbf{v}/\mathbf{m} = /\mathbf{i}/\mathbf{k}/\mathbf{v}/\lambda_1/\dots/\lambda_d$ in a postorder tree walk. If node n is matching, we add its corresponding content node to r and proceed to the next descendant. If node n has no children and is neither matching, nor volatile, it is unproductive and therefore pruned from the index. After we finish traversing all descendants n , we return the result set r .

Algorithm 2: QueryQTP

Data: Query $Q(k, v, m)$, where k is a property, v a value and $m (= / \lambda_1 / \dots / \lambda_d)$ a content node's path.

Result: A set of nodes satisfying $Q(k, v, m)$

$r \leftarrow \emptyset$

for node $n \in \text{desc}(/ i / k / v / \lambda_1 / \dots / \lambda_d)$ **in postorder tree walk** **do**

if $\text{matching}(n)$ **then**

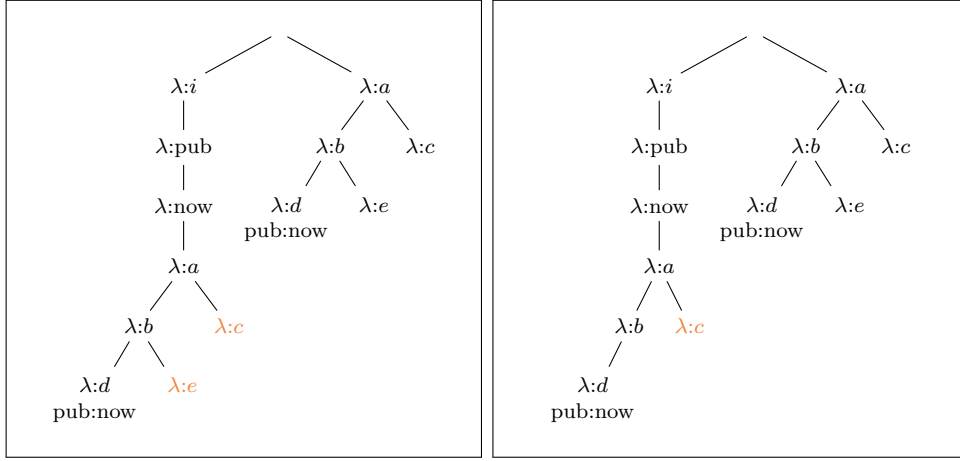
$r \leftarrow r \cup \{ *n \}$

else if $\text{chd}(n) = \emptyset \wedge \neg \text{volatile}(n)$ **then**

 delete node n

return r

$$G^5 \xrightarrow{T_6} G^6$$



Snapshot G^5

Snapshot G^6

Figure 10: Query-Time Pruning applied on Oak. Assume nodes $/i/pub/now/a/b/e$ and $/i/pub/now/a/c$ are unproductive in snapshot G^5 . Transaction T_6 executes CAS query $Q(\text{pub}, \text{now}, /a/b)$ which queries for all descendants of $/a/b$ with “pub” set to “now” and commits the resulting snapshot G^6 . QTP is used during query execution.

Example 8 (QTP). Consider Figure 10. Transaction T_6 executes CAS query $Q(\text{pub}, \text{now}, /a/b)$ which queries for all descendants of $/a/b$ with “pub” set to “now” in snapshot G^5 . Assume the query executor uses QTP and nodes $/i/pub/now/a/b/e$ and $/i/pub/now/a/c$ are unproductive. The query executor traverses all descendants of $/i/pub/now/a/b$ and therefore prunes the unproductive index node $/i/pub/now/a/b/e$. Since the other unproductive index node ($/i/pub/now/a/c$) is not traversed during query execution, it is not pruned and remains in the index unproductive. The resulting snapshot G^6 is committed by T_6 after finishing query execution.

```

/**
 * Answers a CAS Query and prunes traversed unproductive index nodes.
 *
 * @param k: Property we query for
 * @param v: Value we query for
 * @param m: Path of content node which the descendants we query for
 * @param root: Latest Oak snapshot
 * @returns An iterable with content nodes satisfying the CAS query
 */
Iterable<Tree> QueryQTP(String k, String v, String m, Root root) {

    /* e.g.: /oak:index/pub/:index/now */
    String indexRootPath = concat(OAK_INDEX_PATH, k, v);

    /* e.g.: /oak:index/pub/:index/now/a */
    String queryNodePath = concat(indexRootPath, m);

    /* map index nodes to corresponding content nodes */
    return map((Tree n) -> {
        return root.getTree(relativize(indexRootPath, n.getPath()));
    },

    /* filter non-matching index nodes */
    filter((Tree n) -> {
        boolean isMatchingMemo = isMatching(n);

        /* prune if no children, not matching and not volatile */
        if (n.getChildrenCount(1) == 0 &&
            !isMatchingMemo &&
            !isVolatile(n)
        ) {
            n.remove();
        }
        return isMatchingMemo;
    },

    /* postorder tree walk of descendants of /i/k/v/m */
    postOrder(root.getTree(queryNodePath))
    );
}

```

Figure 11: Query-Time Pruning implemented in Java.

Figure 11 shows the Java implementation of QTP in Oak. The algorithm takes a property k , value v , path m , a `Root` and returns a lazy sequence of content nodes satisfying the CAS query $Q(k, v, m)$. By Calling `postOrder()` we get a sequence of nodes representing a postorder tree walk in the subtree rooted at $/i/k/v/m$. We remove any node that is not matching from the sequence using `filter()`. If a node has no children and is neither matching, nor volatile, it is unproductive and we prune it from the index.

5 Experimental Evaluation

5.1 Goals

The goal of our experimental evaluation is to evaluate and compare GC and QTP under different parametrizations and workloads. In Section 5.3 we see how the parameters volatility threshold τ , sliding window length L , workload skew s and update operations per second, impact the number of unproductive nodes in the index subtree. These parameters directly impact the number of volatile nodes and therefore can affect the production of unproductive nodes in the index. In Section 5.4, we evaluate the query performance under GC and in Section 5.4.1 we investigate how GC period T affects GC’s performance. The experimental methodology used on GC is also applied on QTP in Section 5.5. Lastly, Section 5.6 compares GC and QTP directly and suggests under which circumstances one might use periodic GC over QTP, and vice versa.

5.2 Preliminaries

5.2.1 Setup

We use a volatility threshold $\tau = 5$ and sliding window of length $L = 30s$ as defaults, unless otherwise noted. All experiments are conducted on a 15” Macbook Pro 2015 inside a virtual machine running Linux Arch.² We allocate 4 out of 8 virtual cores (Intel i7-4980HQ 2.7 - 4.0 GHz) to the virtual machine and 12 GB of RAM. We allocated half the available virtual cores of the machine because we wanted to ensure the machine’s limited CPU cooling performance will not affect the simulations.

5.2.2 Datasets

We use two datasets in our experiments. Each dataset resembles the content subtree of an Oak instance. The *synthetic* dataset is a complete binary tree of height 19, that is a binary tree in which all leaf nodes have depth 19 [2]. It contains $2^{20} \approx 10^6$ nodes, of which 50% are leaf nodes. Each node has a mean depth and fanout of 18 and 2, respectively. The *real-world* dataset is based on DELL’s website (<http://dell.com>) which is built on top of Adobe’s AEM that uses Oak.³ It is sparser than the synthetic dataset and contains

²<https://www.archlinux.org/about/>

³<https://www.images2.adobe.com/content/dam/acom/en/customer-success/pdfs/dell-case-study.pdf>

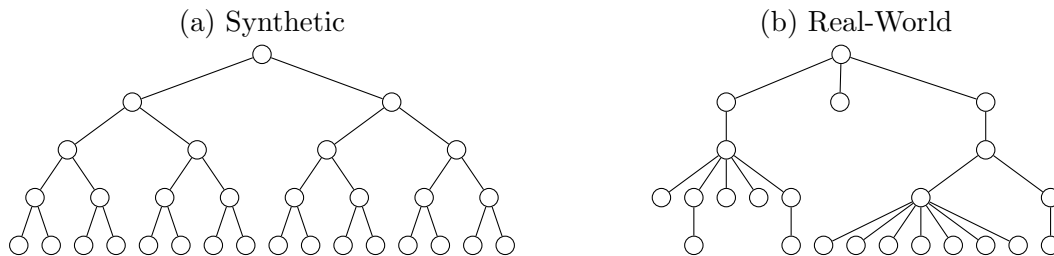


Figure 12: Excerpts of both datasets visualized.

$13 \cdot 10^6$ nodes, 65% of which are leaf nodes. Each node has a mean depth and fanout of 13.68 and approximately 1, respectively. The tree’s nodes have a max depth and fanout of 24 and 1729, respectively. Since the dataset is so sparse, there exist many nodes with a single child and few with many children, forming linked lists (called ”chains”) of nodes sometimes. Figure 12 shows a sample subtree for each dataset.

5.2.3 Workload

In Section 2.1, we introduced the characteristics of a CMS workload. Using a running example we explained that a CMS workload is skewed, update-heavy, changing and CMSs frequently make use of a job-queuing system. For our experiments, we designed a workload that has the noted characteristics.

The workload randomly draws content nodes and executes update operations on them. We denote the set of nodes of the content subtree with depth greater than the mean depth as the *lower content subtree*. The workload is only allowed to draw nodes from the lower content subtree. Doing so, we increase the number of nodes in the index subtree and make query execution more expensive. Assume the workload repeatedly draws a leaf n and executes update operations on it until its corresponding index node becomes volatile and therefore remains in the index. If the workload then randomly draws an ancestor of n , no nodes are added to the index. Since n is volatile, n and its ancestors already exist in the index, hence no nodes have to be added. By only drawing nodes from the lower content subtree, we decrease the likelihood of drawing nodes whose corresponding index nodes already exist in the index and therefore more index nodes are added. With more nodes in the index subtree, query execution becomes more expensive.

In order to have skew, we incorporate the Zipf distribution [7]. The Zipf distribution picks a small subset of nodes more frequent than others. Let N be the number of nodes inside the lower content subtree. We randomly and uniquely assign an integer $k \in [1, N]$ to each node in the lower content subtree, creating a $k \rightarrow \text{node}$ mapping. The probability of randomly drawing the k -th node from the lower content subtree according to the Zipf distribution is:

$$\text{Zipf}(k, N, s) = \left(k^s \cdot \sum_{i=1}^N \frac{1}{i^s} \right)^{-1}$$

The k -th node corresponds to the node which is assigned integer k in the mapping. Skewness s of the Zipf distribution is parametrizable and we use $s = 1$ by default in our experiments, unless otherwise noted. By having a smaller skew, the subset of nodes which is frequently drawn, becomes bigger. With $s = 0$, the Zipf distribution corresponds to the uniform distribution and every node has equal probability to be drawn. In Section 5.3.3, we make an experiment to determine skew’s impact on unproductive nodes.

Figures 13a and 13b show Zipf’s Cumulative Distribution Function (CDF) for both datasets and skew values $s \in \{0, \frac{1}{2}, 1, \frac{3}{2}, 2\}$. If we randomly draw a large sample of nodes according to the distribution with $s = 1$, 80% of these nodes are amongst the same 30k nodes (6% of all nodes) on the synthetic dataset and 40k nodes (5.9% of all nodes) on the real-world dataset.

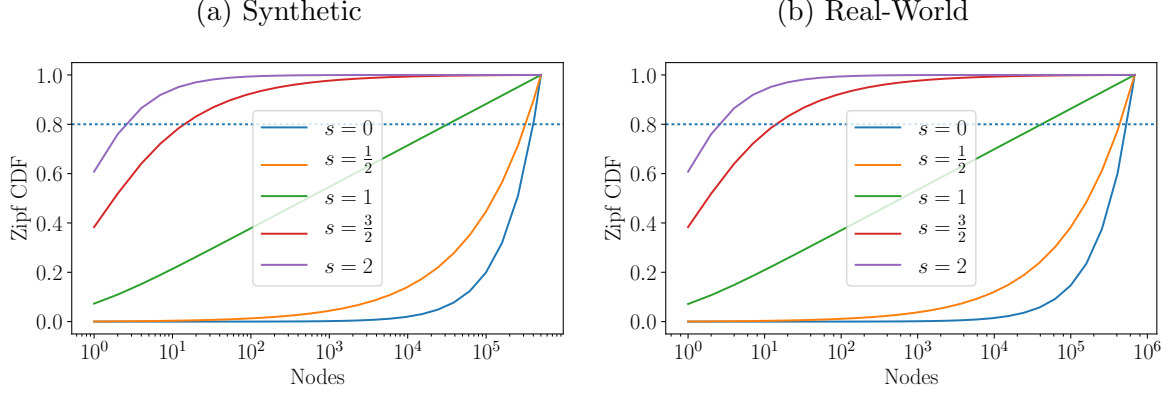


Figure 13: CDF for the Zipf distribution for both datasets.

A CMS workload is also update-heavy. To simulate this, we execute many update operations before executing a query operation. We fix the default update to query ratio to 10:1, i.e., 10 update operations are executed for each query operation. In Section 5.3.4, we investigate the effect of update operations per second.

Additionally, the workload we design has to be changing. We periodically permute the node mappings by reassigning a new random and unique integer $k \in [1, N]$ to each node from the lower content subtree in order to change the hotspots of the simulation, that is the subset of nodes which are frequently updated. We change the hotspot by default every 30 seconds. During a 5 minute experiment we expect 10 different workloads.

Lastly, the workload has to simulate a job-queuing system. An update operation executed during the experiment is composed from two actions. We first set the “pub” = “now” property-value pair to the node and then consecutively remove it. The actions simulate a node being inserted into the queue and then being processed and removed. Using our designed update operation, there can be at most one matching index node at any instance of time.

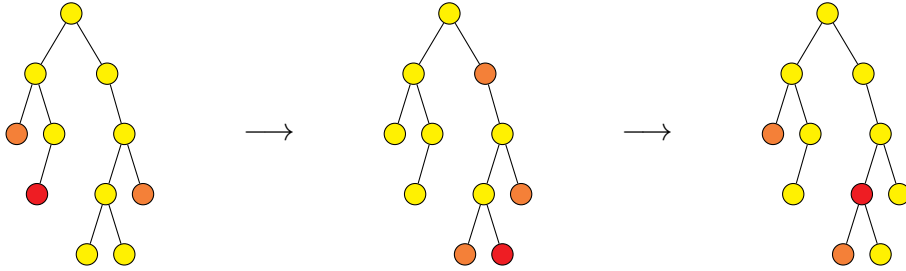


Figure 14: CMS Workloads visualized. Red shaded nodes are the most frequently drawn nodes.

Figure 14 depicts heatmaps of a content subtree of Oak over time. These heatmaps show how often the designed workload selects specific nodes inside the content subtree. Red shaded nodes are the most frequently drawn nodes.

5.3 Unproductive Nodes

5.3.1 Volatility Threshold τ

Volatility threshold τ determines after how many insertions/deletions of an index node it becomes volatile (see Definition 4). In this section, we study the impact of volatility threshold τ on unproductive nodes, which directly affect query runtime.

We hypothesize that an increase in τ yields a decrease to the number of traversed unproductive nodes during query execution under a CMS workload. If τ increases, it is less likely for a node to become volatile. Having fewer volatile nodes should cause a decrease to the number of unproductive nodes and consequently also query runtime in the CMS workload.

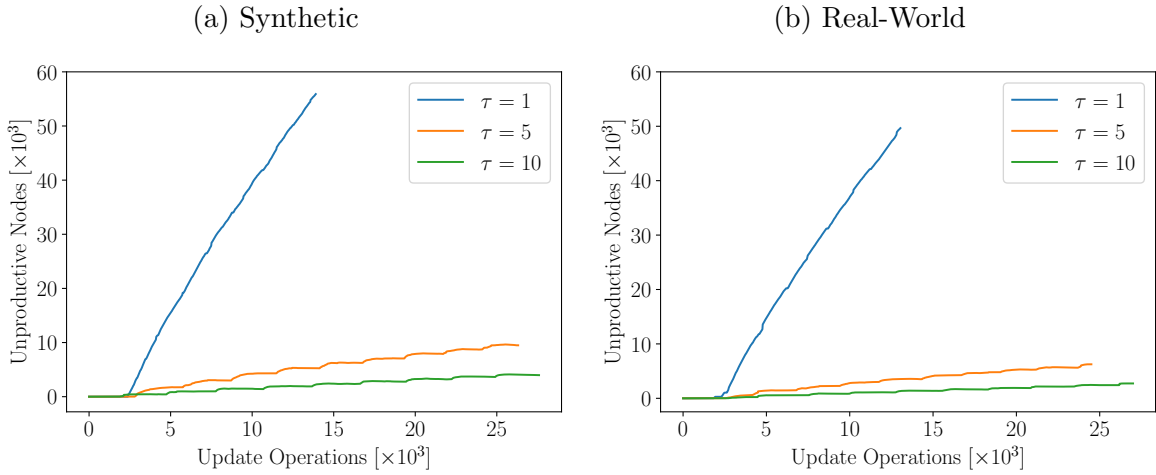


Figure 15: Unproductive Nodes over the accumulated number of update operations with thresholds $\tau \in \{1, 5, 10\}$.

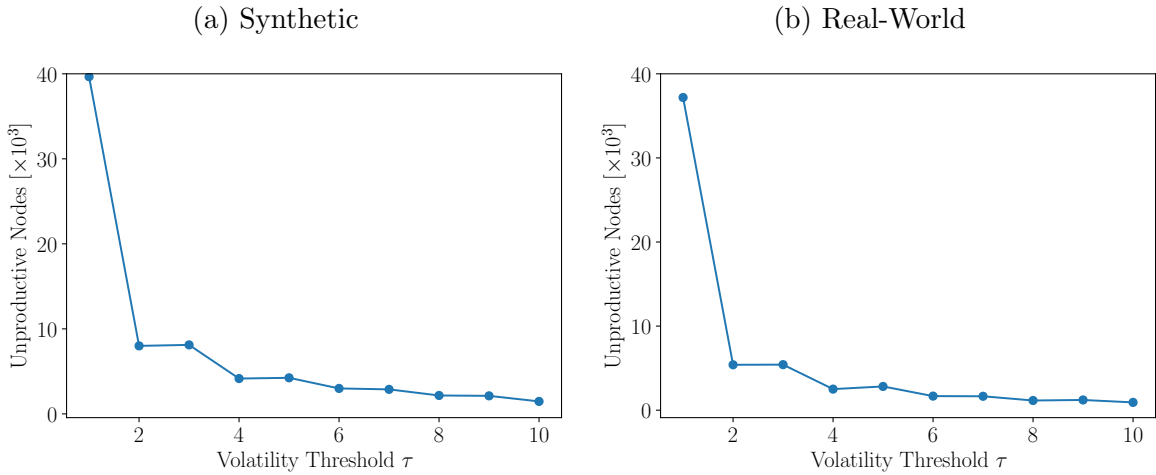


Figure 16: Unproductive Nodes over volatility threshold τ .

Figures 15a and 15b show the increase in unproductive index nodes for different thresholds $\tau \in \{1, 5, 10\}$ over the course of a five minute experiment. We observe lower thresholds τ yielding a steeper slope. A lower volatility threshold increases the likelihood of a node becoming volatile. The amount of volatile nodes also affect the number of unproductive nodes, since volatile nodes eventually stop being frequently updated when the workload changes and become unproductive. The increase in unproductive nodes in the index also affects query runtime because Oak has to traverse these nodes during query execution.

The Figures show that unproductive nodes increase linearly as update operations increase. We expect the increase to be sublinear for longer experiments. As time passes by, the index becomes more static (see Section 3.2), because more volatile/unproductive index nodes are retained and fewer nodes are pruned. As a result, fewer nodes become volatile as the experiment progresses. The number of traversed unproductive nodes is upper-bounded by the number of content nodes. This is the case, because the index subtree for a given property-value pair (e.g. “pub” = “now”) can not have more nodes than the content subtree, because the index subtree is a subset of the content subtree. Eventually, we will reach the upper bound.

Figures 16a and 16b compare the number of traversed unproductive nodes during query execution over a range of thresholds. We observe a decrease in unproductive nodes while increasing threshold τ . As suggested earlier, a lower volatility threshold increases the amount of volatile nodes in the index and consequently also increases the number of unproductive nodes.

We see the two variables sharing a power law relationship. We believe the workload’s skew to be responsible for the power law relationship. At $\tau = 0$, any node drawn by the workload becomes volatile. That is why we see so many unproductive nodes. At $\tau = 2$, we observe a sudden drop in unproductive nodes. That is because content nodes have to be drawn twice in order to become unproductive, and since the workload is skewed, only a few nodes are drawn twice. For even larger values, the function flattens out. Since the workload is skewed, only a handful of nodes are drawn enough by the workload to become volatile.

The query executor has to traverse $5k$ unproductive index nodes when Oak has a threshold of $\tau = 5$. By increasing the threshold to $\tau = 6$, the average number of traversed unproductive nodes does not decrease significantly because of the Zipf distribution’s skewness. The majority of update operations are executed amongst a small subset of nodes. This explains the power law relationship between traversed unproductive nodes and τ .

Summarizing, all observations verify our hypotheses. Increasing volatility threshold τ decreases the number of unproductive nodes traversed which decreases query runtime. Increasing the volatility threshold causes fewer nodes to become volatile. Since we create fewer volatile nodes, we also reduce the number of unproductive nodes. Less unproductive nodes yield lower WAPI query runtimes.

5.3.2 Sliding Window of Length L

The Sliding Window Length L determines the length of the recent workload that WAPI considers to compute an index node’s volatility count (cf. Definition 4). Greater values of L allow WAPI to count more updates and therefore increase the chances of a node becoming volatile. In this section, we study the effect of the sliding window on the number unproductive nodes.

We hypothesize that an increase in L yields an increase to the number of traversed unproductive nodes during query execution. If L increases, it is more likely for a node to become volatile, since more updates are considered towards the volatility count. Having more volatile nodes should imply an increase in unproductive nodes and consequently also query runtime.

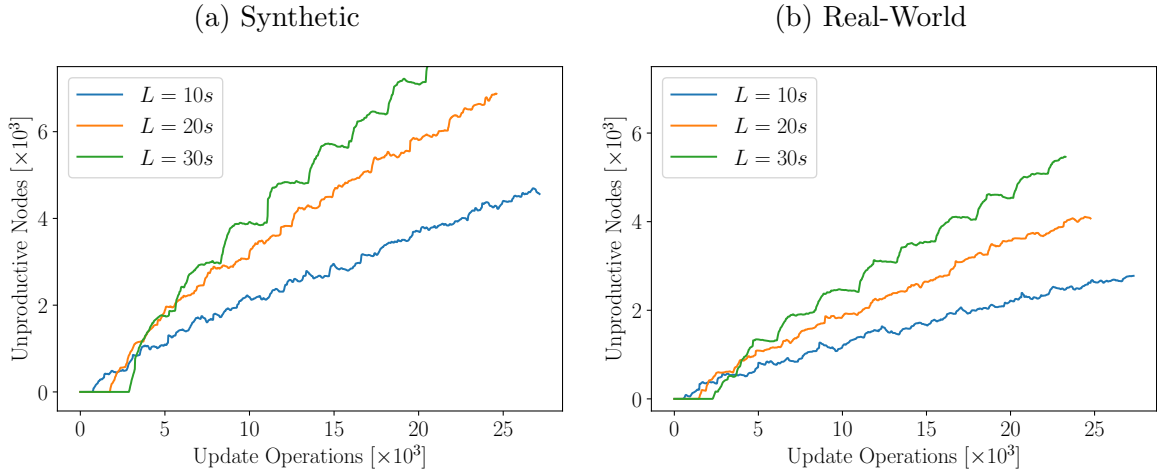


Figure 17: Unproductive Nodes over accumulated update operations with sliding window lengths $L \in \{10s, 20s, 30s\}$.

Figures 17a and 17b visualize the number of unproductive nodes WAPI traverses during query execution for three distinct lengths $L \in \{10s, 20s, 30s\}$ of the sliding window w.r.t. update operations. As expected, the number of unproductive nodes increases the fastest for the longest sliding window length ($L = 30s$). In the long term, the number of unproductive nodes increases sublinearly w.r.t. time. We explained earlier that as time passes, the index becomes more static and the number of unproductive nodes converges towards the upper bound set by the number of content nodes.

Figures 18a and 18b show the number of unproductive nodes the query executor has to traverse w.r.t. the sliding window length L . Greater sliding window lengths cause a sublinear increase to the number of unproductive nodes traversed by WAPI during query execution. Since the index becomes more static as explained earlier, the function converges towards the upper bound. The Figures seem to suggest a linear relationship between the two variables. That is because the number of content nodes is greater by two orders of magnitude compared to the number of unproductive nodes we can produce

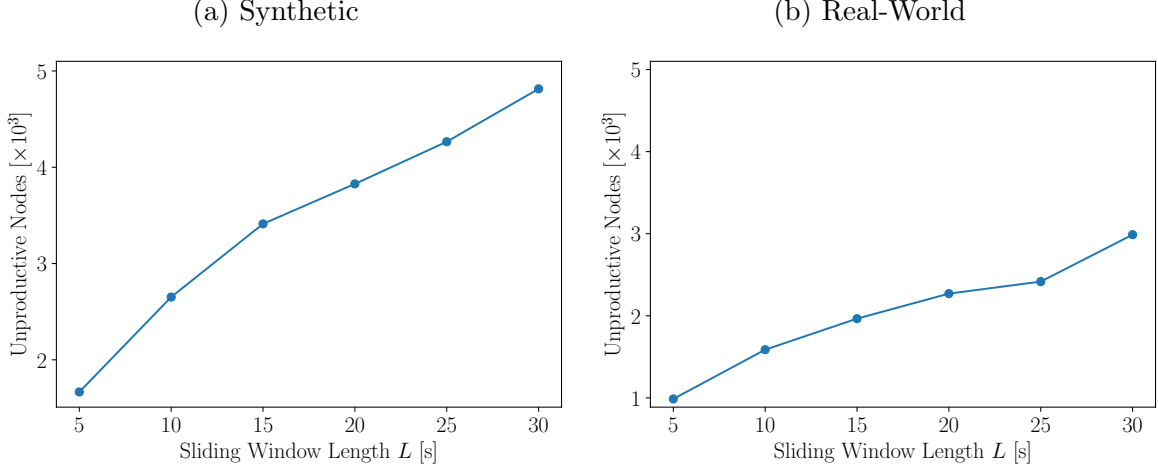


Figure 18: Unproductive Nodes over sliding window length L .

during the entire five-minute experiment. We expect to see the function converge if we run the experiment for a significantly longer duration of time.

Concluding, we see sliding window of length L to be affecting the rate of growth of unproductive nodes. Increasing L does increase the likelihood of an index node to become volatile and later unproductive.

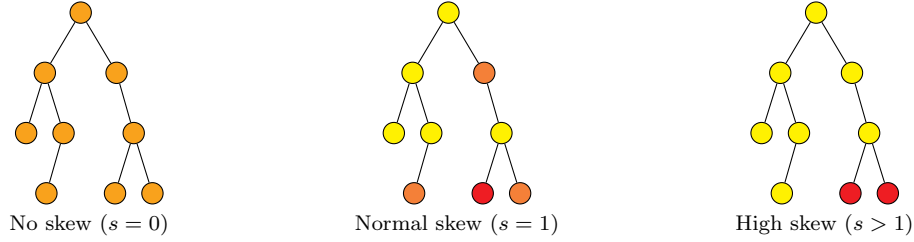


Figure 19: Hotspots affected by different skew values. Red shaded nodes denote frequently drawn nodes.

5.3.3 Workload Skew s

One of the key characteristics of the CMS Workload is skewness, as mentioned in Section 2.1. A small subset of nodes is frequently updated, whereas the rest of the nodes is not. We refer to the set of content nodes which are frequently drawn by the skewed workload as *hotspot*. Nodes inside the hotspot most likely become volatile and later on unproductive. We use the Zipf distribution to model a skewed workload (see Section 5.2.3). If skew s increases, the hotspot becomes smaller but the hotspot's member nodes are drawn more frequent.

We expect skew s to affect the number of unproductive nodes WAPI traverses during query execution. Since some nodes are updated more often, these nodes are more likely to become volatile and also unproductive, later on. On the other hand, increasing skew

s also reduces the size of the hotspot and therefore reduces the number of volatile nodes, hence fewer nodes become unproductive.

Figure 19 visualizes workloads with different skew values for the reader’s convenience. In the uniform workload, that is the workload with no skew ($s = 0$), each node is equally likely to be drawn. The workload with skew $s = 1$ has a hotspot and nodes in the hotspot are drawn more frequent than other nodes. The workload with high skew ($s > 1$) has an even smaller hotspot but member nodes are drawn more frequent.

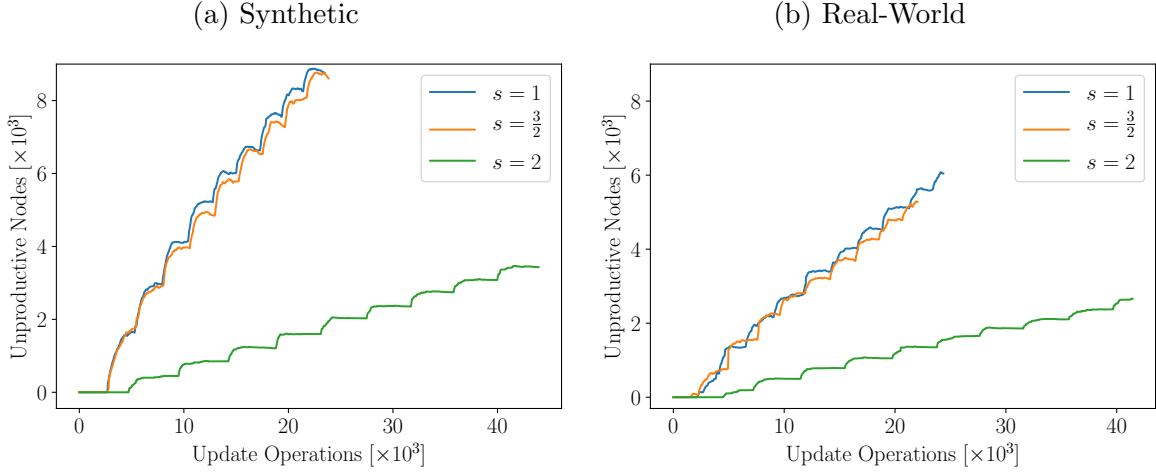


Figure 20: Unproductive Nodes over accumulated update operations with skew $s \in \{1, \frac{3}{2}, 2\}$.

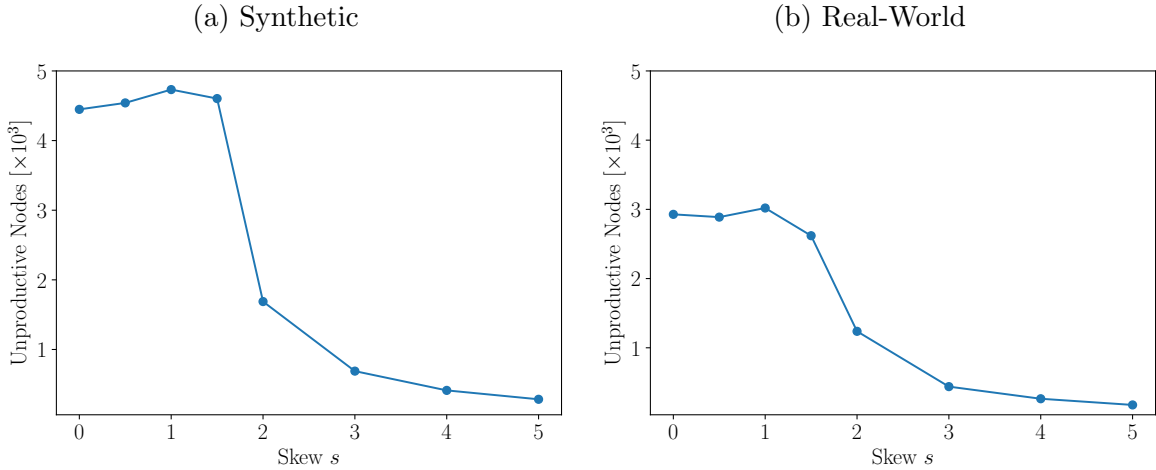


Figure 21: Unproductive Nodes over skew s .

Figures 20a and 20b show the number of unproductive nodes traversed by WAPI during query execution with workloads having skew $s \in \{1, \frac{3}{2}, 2\}$. We see a smaller skew s increases the number of unproductive nodes. We believe that when increasing s , the

size of the hotspot decreases and therefore the number of nodes becoming volatile, and eventually unproductive, decreases. We also see that a workload with $s = 1$ has barely more unproductive nodes over time than a workload with $s = \frac{3}{2}$.

Figures 21a and 21b offer a more detailed perspective by showing the traversed unproductive nodes w.r.t. skew s . We observe that the number of unproductive nodes maximizes when $s = 1$. When the workload draws nodes with little to no skew ($s < 1$), update operations are spread (almost) uniformly amongst nodes. As a result, fewer nodes become volatile, and later on unproductive. We see a sudden drop in unproductive nodes when skew is $s = 2$. The hotspot reaches a point where it loses a significant amount of nodes because the workload is so skewed. From there on, we see a sublinear rate of convergence.

5.3.4 Update Operations per Second

One of our initial assumptions was that Oak’s workload is update heavy, amongst others. To emphasize the implications of such an update heavy workload, we conduct the following experiment. Oak is benchmarked with different update operations per second, in order to show the effect of such a workload on unproductive nodes.

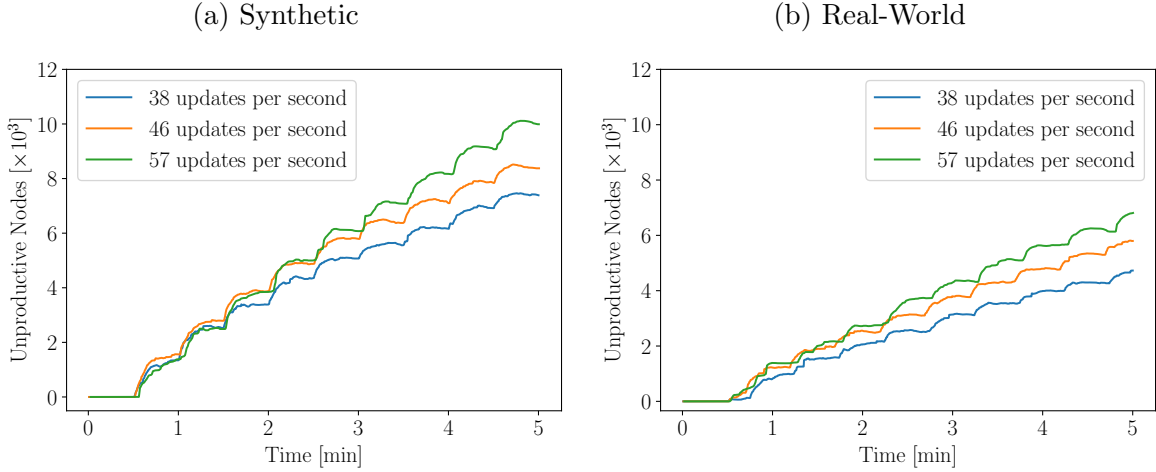


Figure 22: Unproductive Nodes over time for 38, 46 and 57 updates per second.

Figures 22a and 22b show the number of traversed unproductive nodes over time. We run the simulation for workloads having 38, 46 and 57 updates per second, respectively. An increase in time yields an increase in traversed unproductive nodes, an observation depicted in Figures 23a and 23b, too. The latter Figures show the number of traversed unproductive nodes over updates per second. A node becomes volatile if it is updated τ times during a sliding window of length L . If we increase the number of updates during a fixed interval, the node is more likely to become volatile. If the number of volatile nodes increases, there is an increase in unproductive nodes, too, as shown before.

In addition, the increase we observe is linear. Increasing the number of update operations, should yield a proportional increase in volatile nodes which later on become

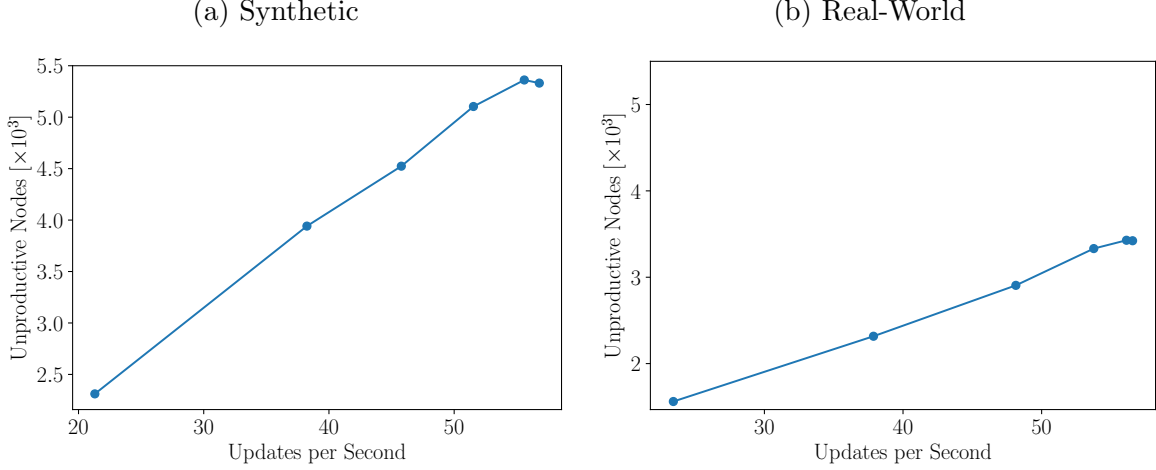


Figure 23: Unproductive Nodes over updates per second.

unproductive.

5.4 Garbage Collection

Our next experiment evaluates Oak’s query performance under periodic Garbage Collection. We record the average query runtime and the number of volatile and unproductive index nodes during query execution while having periodic GC enabled.

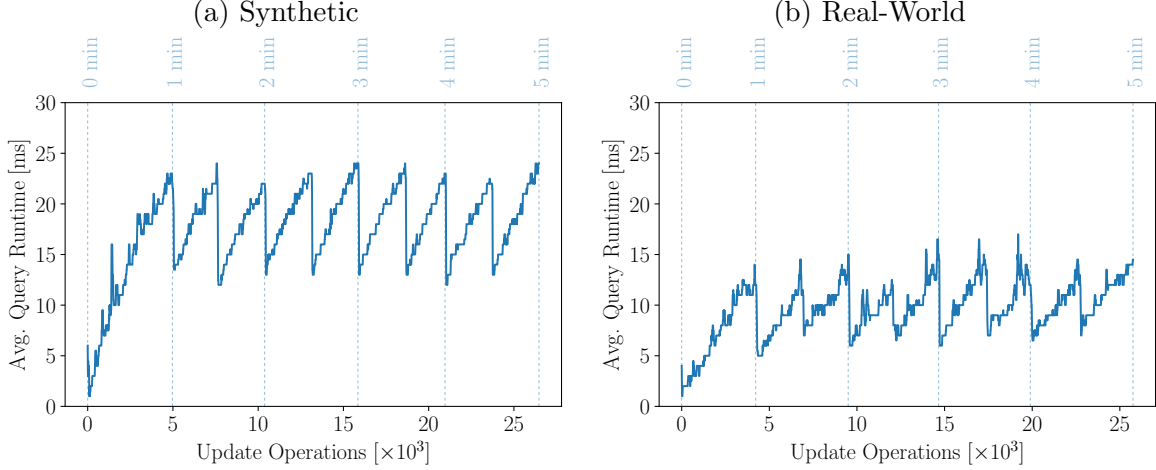


Figure 24: Query Runtime over accumulated update operations, periodic GC enabled.

Figures 24a and 24b show the resulting decrease in query runtime when applying periodic GC on Oak. The garbage collector is run every 30 seconds. We observe the query runtime increase during the first minute of the simulation. When GC is run for the first time (at 30 seconds), it does not encounter any unproductive nodes because with $L = 30s$ no volatile node has yet ceased to be volatile. Afterwards, every succeeding

GC encounters and prunes unproductive nodes. The pruned nodes are made visible by the sawtooth pattern depicted in the Figure. Each drop in query runtime reflects the speedup caused by cleaning unproductive nodes. The query runtime oscillates between 20 and 30 milliseconds with a mean of 25 milliseconds on the synthetic dataset. The real world’s query runtime oscillates between 8 and 18 milliseconds with a mean value of 13 milliseconds.

We also observe that queries executed on the real-world dataset are faster, on average. We will see later that the synthetic dataset has more index nodes compared to the real-world dataset. The query executor has to traverse more nodes under the synthetic dataset, hence the query runtime increases.

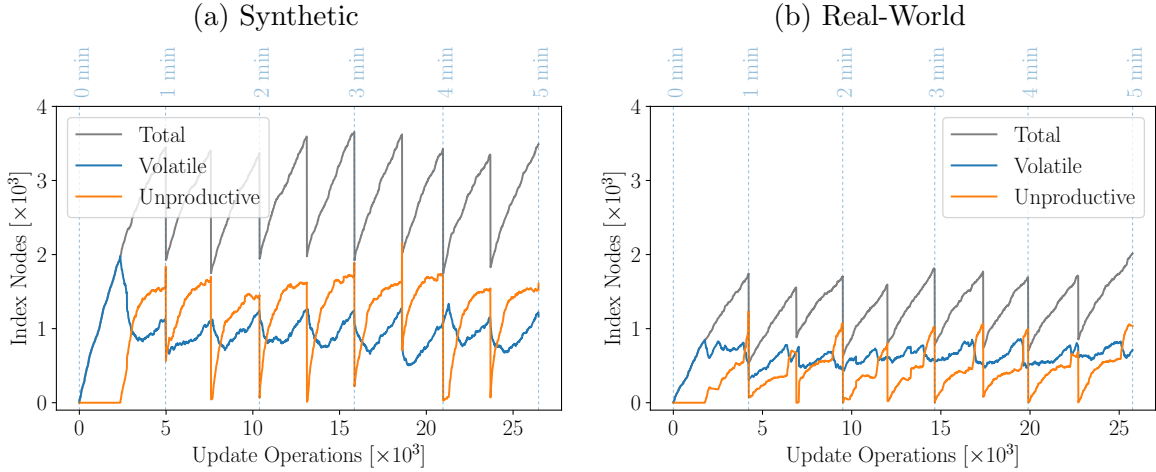


Figure 25: Index Nodes over accumulated update operations, periodic GC enabled.

Figures 25a and 25b visualize the index structure during query execution with periodic GC. Unproductive nodes increase until GC is run, after which they are completely removed again. Observing the Figures, the number of unproductive nodes does not always reach zero. That is because we use a running median. The running median removes outliers, and in our case also points which had zero unproductive nodes.

The volatile nodes also have a cycloid pattern. The cycloid pattern, as seen and explained in Figures 4a and 4b, appears because many nodes cease to be volatile immediately after a workload change and many nodes become volatile before a workload change. In contrast to Figures 4a and 4b, we do not see the number of volatile nodes having a downward trend anymore from the 30 second mark onward. Instead in Figures 25a and 25b, the number of unproductive nodes oscillates around a constant throughout the experiment. Since the likelihood of node becoming volatile remains constant, there is no increase or decrease in volatile nodes, on average.

Comparing the synthetic dataset to the real-world dataset, we observe another phenomenon. The gap between traversed volatile and total traversed nodes is significantly bigger in the synthetic dataset. In other words, the percentage of volatile index nodes is higher in the real-world dataset. We believe that non-volatile ancestors of volatile nodes account for the gap because ancestors of volatile nodes are more likely to become volatile

too in sparser subtrees. Sparse subtrees form long chains of nodes. If a leaf node in such a chain becomes volatile, so do its ancestors, since they are also updated when the leaf node is added/removed. The more sparse a subtree is, the longer the chains become. Since the real-world dataset is sparser than the synthetic dataset (Section 5.2.2), the percentage of volatile nodes is higher compared to the synthetic dataset.

5.4.1 GC Period T

In this section we discuss the performance impact of GC period T . Oak can run GC arbitrary many times. Given our setup, we would like to find out what the optimal period of GC is. We run GC under varying period and compare the results.

We expect the number of unproductive nodes to decrease as we run GC more often, i.e., T decreases. GC uses system resources and if GC is run too often, we might decrease the query performance because the system is busy garbage collecting instead of executing queries. This negative impact of GC should be especially pronounced if the system is under high load.

For the following experiments, we allocate only a single virtual core to the virtual machine in order to disable parallel computation and simulate system resource shortage. By doing so, we can demonstrate that GC, run on a different thread from query execution, steals CPU time from the query executor.

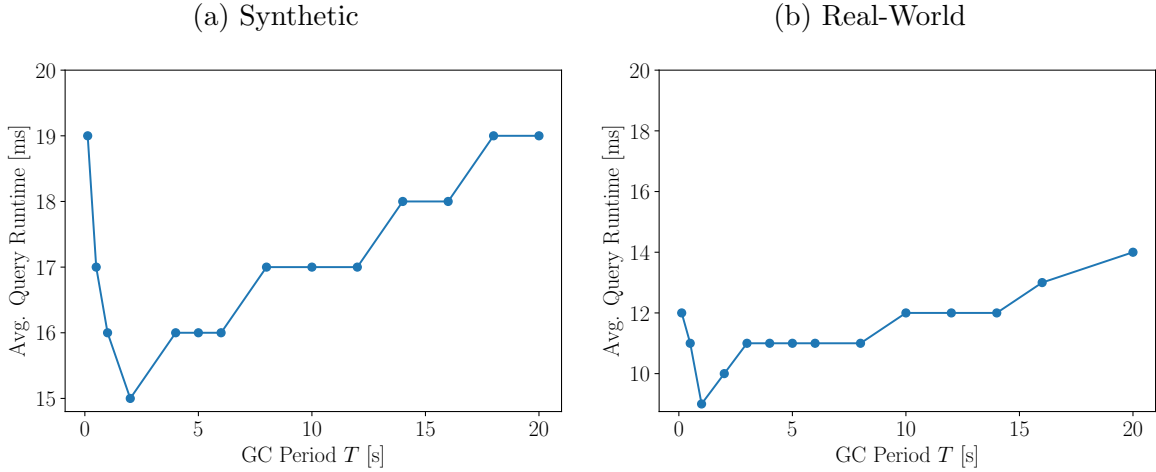


Figure 26: Avg. Query Runtime over GC Period T .

Figures 26a and 26b show the average query runtime w.r.t. period T . The optimal period T^* seems to be in the interval $T^* \in [500ms, 2000ms]$ for both datasets. If $T < T^*$, GC starts stealing CPU time from the query executor and queries take longer to process. On the other hand, when $T > T^*$, the query runtime increases sublinearly because the index fills up with unproductive nodes. We see a sublinear increase since the index becomes more static over time, as explained previously.

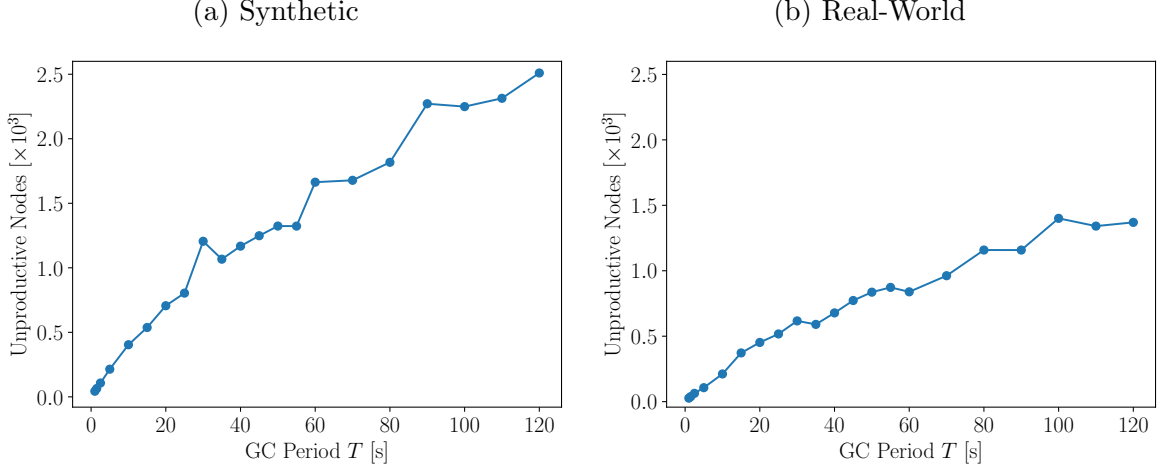


Figure 27: Unproductive Nodes over GC period T .

Figures 27a and 27b depict the number of unproductive nodes w.r.t. T . We see that the number of unproductive nodes increases sublinearly as T increases. The function converges because the index becomes more static over time, as explained previously.

When comparing Figures 27a and 27b to Figures 26a and 26b it becomes more clear that although we have fewer unproductive index nodes under frequent GC ($T < 2s$), the average query runtime increases because of the system resource shortage. For any period $T > 2s$, the average query runtime seems to correlate with the number of unproductive nodes.

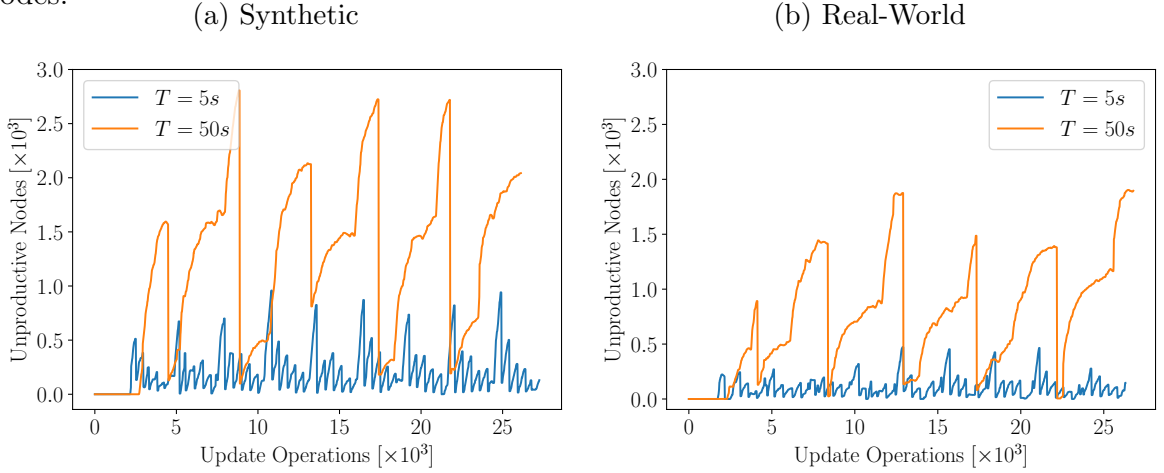


Figure 28: Unproductive Nodes over accumulated update operations with GC period $T \in \{5s, 50s\}$.

Let us consider Figures 28a and 28b. They show the number of traversed unproductive nodes for two different periods $T \in \{5s, 50s\}$, w.r.t. the accumulated number of update operations. We clearly see that GC with period $T = 5s$ has a smaller number of unproductive nodes over time, on average.

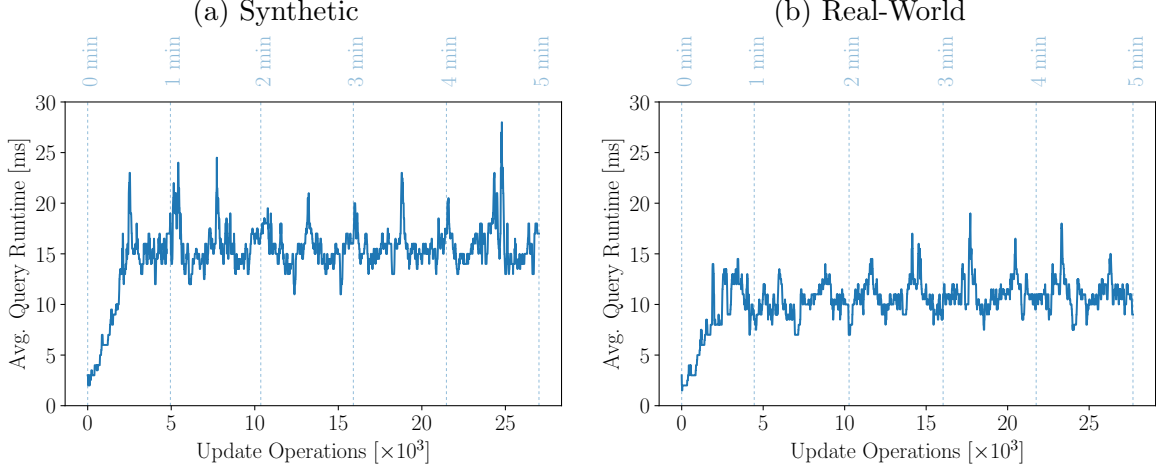


Figure 29: Query Runtime over accumulated update operations, QTP enabled.

5.5 Query-Time Pruning

In this section, we evaluate Oak’s query performance under Query-Time Pruning. We record the average query runtime and the number of volatile and unproductive index nodes, similar to Section 5.4.

Figures 29a and 29b depict the resulting decrease in query runtime when applying QTP on Oak. We see the runtime rise during the first 30 seconds and then see the runtime fluctuate around 15ms for the synthetic and 10ms for the real-world dataset.. We believe that traversing volatile nodes dominates query runtime.

Figures 30a and 30b show the composition of the index with QTP enabled during query execution. We observe the number of unproductive nodes to be close to 0 throughout the entire simulation. This is because the index is continuously cleaned by QTP. As long as query execution is continuous, so is our index cleaning. Since after ten updates we always query the root content node in our experiment, we clean the whole index subtree from unproductive nodes. QTP also cleans unproductive nodes the moment they emerge. In contrast, periodic GC cleans these unproductive nodes 30 seconds after they appear in our experiment. As a result, frequent queries benefit the most from QTP, since a subsequent query with the same filter as its preceding query, does not traverse any unproductive nodes under QTP.

As seen previously with GC (Section 5.4), we also observe a bigger gap between volatile and total index nodes in the synthetic dataset compared to the real-world one due to more non-volatile ancestors of volatile nodes.

Figure 31 depicts a cost breakdown of a CAS-query with QTP enabled. Assuming we use Algorithm 2, a 20ms query spends on average 10ms (50%) traversing the index subtree, 0.5ms (2.5%) checking the match property of visited nodes, 6ms (30%) checking if the visited nodes have children, 0.5ms (2.5%) checking if the visited nodes are volatile and 3ms (15%) writing, that is pruning unproductive nodes and committing the transaction.

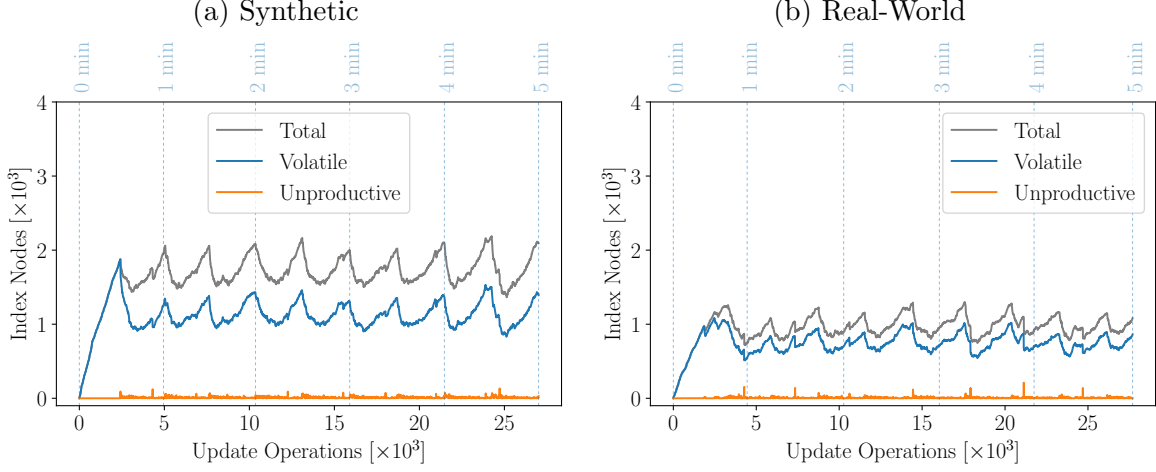


Figure 30: Index Nodes over update operations, QTP enabled.

QTP added a 9.5ms overhead to a 11.5ms query, which corresponds to an additional 90% runtime. The main overhead added by QTP in this experiment is checking if the visited nodes have children. Checking the match property and checking if a node is volatile, are cheap operations because Oak caches node properties internally. Writing changes becomes a significant overhead when there are many unproductive nodes, e.g. when a query has not been executed for a long time. Although QTP is adding an overhead, it is still faster than using no QTP in the long-term because of unproductive nodes (cf. Figure 3).

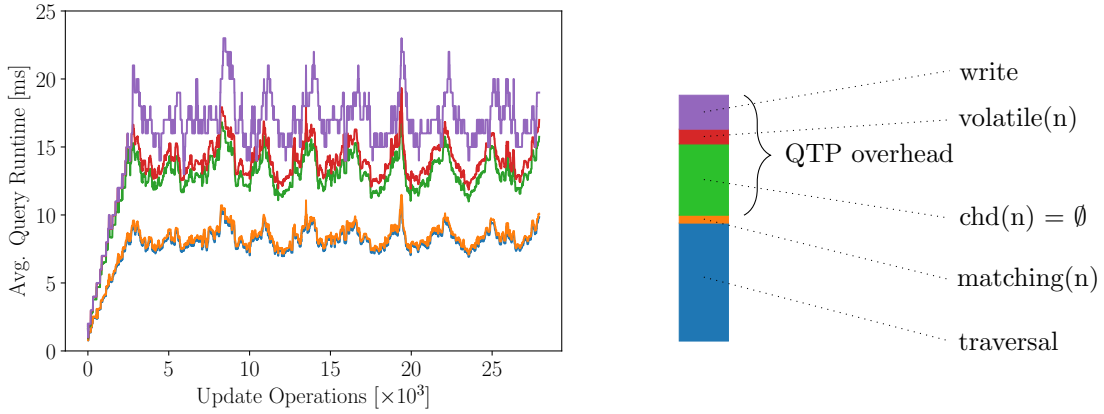


Figure 31: The overhead of various steps during QTP based on Algorithm 2 are highlighted with different colors.

5.6 Comparison

The previous Sections presented parameters affecting the rate of growth of unproductive nodes in the index subtree and also benchmarked the two proposed approaches, periodic GC and QTP. This Section directly compares the two solutions and provides general outlines on choosing the most suitable approach.

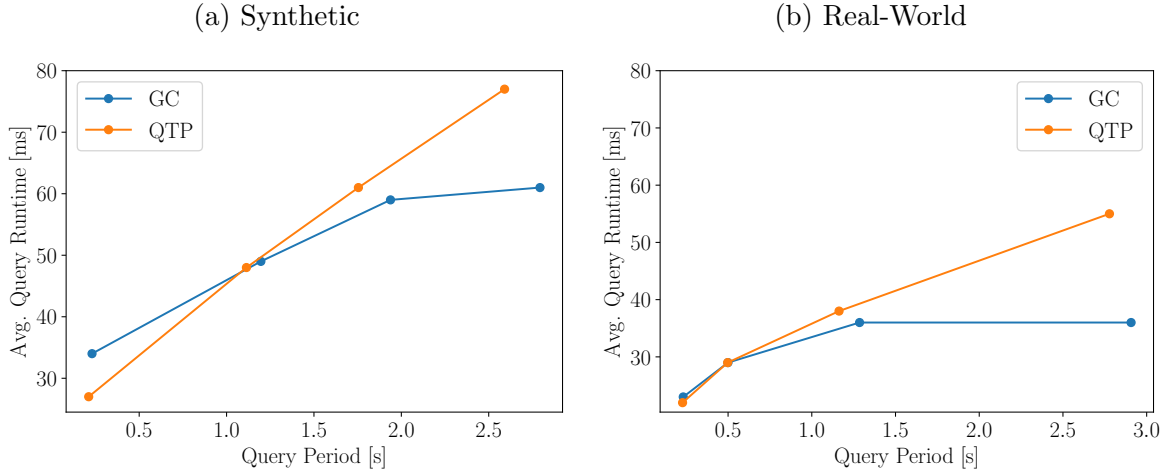


Figure 32: Periodic GC vs. QTP.

Figures 32a and 32b show the average query runtime over the query period, that is the average time in between two consecutive queries, under GC and QTP. We observe that QTP outperforms GC for a low query period. Queries occur often and the index is continuously cleaned by QTP. Since there are so many queries during a fixed interval, each query traverses a small number of unproductive nodes. The speedup gained by not continuously traversing many unproductive nodes is significantly greater than QTP's overhead (Section 5.5).

On the other hand, as queries become more unfrequent, periodic GC is preferred. When two consecutive queries under QTP are far apart, many unproductive nodes gather in the index subtree and when the latter query is run, it suffers from a hefty performance penalty from traversing many unproductive nodes. Say GC is executed between the two queries, then the latter query traverses considerably fewer unproductive nodes and is faster because unproductive nodes were already cleaned up by the GC run.

Figures 32a and 32b each have a data point on which the two functions (GC, QTP) intercept. For example, in Figure 32a, when $\tau = 5$, $L = 30s$, $s = 1$, $T = 30s$ and we execute a query every 1.1 second, we are indifferent choosing between GC and QTP. The same holds in Figure 32b for a query period of 0.5 seconds. We denote as *query runtime equilibrium point*, the point on which queries under GC take, on average, as long as queries under QTP, given τ , L , s , T and the query period.

The query performance under GC and QTP also depends on the production rate of unproductive nodes in the index subtree. Assume we produce an extremely low number of unproductive nodes over time. In that case, QTP should be more beneficial compared

to GC because no matter how far queries are apart, unproductive nodes are not produced fast enough so that queries take a notable performance hit.

The other extreme is if unproductive nodes spawn extremely frequently. In that case, running GC might be more beneficial, or at least less expensive. Assume we produce so many unproductive nodes over time, so that pruning an unproductive node with a query under QTP is useless and wastes time because it becomes unproductive again before the next query starts. In that scenario, GC is superior (depending on T) because queries only take a fraction of the performance penalty QTP imposes, since QTP has a large write overhead with many unproductive nodes.

The rate at which unproductive nodes are produced depends on the volatility threshold τ (Section 5.3.1), sliding window length L (Section 5.3.2), workload skew s (Section 5.3.3) and the number of updates per second (Section 5.3.4), as far as we know. We conjecture that the rate at which unproductive nodes are produced, controls the horizontal offset (x-coordinate) of the equilibrium point. If we produce more unproductive nodes during a fixed interval, the equilibrium point might move towards the left. Producing more unproductive nodes over time might impose a higher overhead on queries under QTP. If we produce fewer unproductive nodes, the equilibrium point might move towards the right because queries under QTP benefit from traversing fewer unproductive nodes.

5.7 Summary

Section 5 focused on the experimental evaluation. We first conducted a series of experiments to test parameters that impact the production rate of unproductive nodes. Our main findings with respect to the production rate of unproductive nodes are:

- A lower volatility threshold τ increases the number of unproductive nodes in index subtree. τ and unproductive nodes share a power law relationship (Section 5.3.1).
- A greater sliding window length L increases the number of unproductive nodes. The results of the experiments suggest a linear relationship between L and the number of unproductive nodes, but we expect a sublinear increase if the experiment is run longer (Section 5.3.2).
- The production rate of unproductive nodes is the highest with skew $s = 1$. The number of unproductive nodes drops rapidly for high values of s . The hotspot gets smaller and therefore we observe fewer unproductive nodes (Section 5.3.3).
- More update operations per second increase the production rate of unproductive nodes linearly. If the experiment is run longer, we expect to see a sublinear increase (Section 5.3.4).

We then evaluated and compared GC and QTP. Both approaches reduced the number of unproductive index nodes and increased Oak’s query performance significantly. We present our main findings below:

- Queries under GC are executed faster due to fewer unproductive index nodes. The number of unproductive nodes followed a sawtooth pattern over time. Running GC more often reduces the number of unproductive nodes, but running GC too often steals system resources from query execution and adds a performance penalty to the queries. The system is busy garbage collecting instead of answering queries (Section 5.4).
- Queries under QTP traverse near no unproductive nodes throughout the entire five-minute experiment. That was the case, because queries under QTP were cleaning the index continuously in comparison to the sawtooth pattern for GC. This results in more stable query runtimes. We later analyzed the cost structure of a query under QTP. The average query runtime increased from 11.5ms to 20ms. The main overhead of QTP in this experiment is checking if visited nodes have children, which is done to determine unproductive nodes. Although QTP is adding an overhead, it is still significantly faster than using no QTP in the long-term because of uproductive nodes (Section 5.5).
- QTP outperforms GC under workloads with frequent queries, depending on period T . Otherwise, GC is the method of preference, given the circumstances of our experiment. We conjecture that GC might outperform QTP under workloads with a high unproductive node production rate (Section 5.6).

6 Conclusion

We proposed two algorithms to detect and prune unproductive nodes in a workload-aware hierarchical index (WAPI). Unproductive nodes are an unwanted byproduct of WAPI’s decision to not prune volatile nodes. They do not contain any data and are therefore uselessly traversed during query execution. By pruning these unproductive nodes, we are able to improve the index query performance since we only encounter index nodes that do yield a query match, are volatile, or are an ancestor of such a node. Additionally, we free storage space that was previously allocated to unproductive nodes.

The first approach is to periodically apply garbage collection (GC) on WAPI. Doing so requires us to explicitly traverse the index subtree thus introducing a performance penalty, but we are able to prune all unproductive nodes from the index. The second approach is to prune unproductive nodes during query execution (QTP). Doing so we avoid traversing the index explicitly, but we add a small overhead on query execution because we have to check each visited node’s productiveness. Additionally, we only prune unproductive nodes from subtrees which are being queried. We extend WAPI by implementing periodic GC and QTP in Apache Jackrabbit Oak and our evaluation shows that both approaches significantly decrease the query runtime of WAPI.

Our experimental evaluation also showed that volatility threshold τ , sliding window of length L , skew s and the number of update operations per second, impact the number of volatile nodes in the index subtree, which later on most likely become unproductive and slow down queries. We then compare GC against QTP and conclude that deciding between one of the two depends on how fast unproductive nodes emerge and how often queries occur.

It is still open how concurrency control affects the performance of the two mentioned algorithms. We believe QTP will take a bigger performance hit in comparison to GC, since QTP continuously prunes and updates the index causing more index conflicts under concurrent transactions than GC with its periodic updates. We also do not know how much the query performance under QTP is affected when the query filter changes more often and frequent queries do not benefit from QTP anymore. An analysis of the impact of the unproductive node production rate on the query runtime equilibrium point would also be of importance. We conjecture that GC outperforms QTP under a high unproductive node production rate. Furthermore, we did not investigate what the benefits and caveats are of running both periodic GC and QTP simultaneously on the database system.

References

- [1] J. W. Bruce and P. J. Giblin. *Cusp*. Cambridge university press, 1992.
- [2] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms, Third Edition*. The MIT Press, 3rd edition, 2009.
- [3] J. R. Evans. Running median filters and a general despiker. *Bulletin of the Seismological Society of America*, 72(1):331, 1982.
- [4] L. LIU and M. T. ÖZSU, editors. *Hierarchical Data Model*. Springer US, Boston, MA, 2009.
- [5] C. Mathis, T. Härder, K. Schmidt, and S. Bächle. XML indexing and storage: fulfilling the wish list. *Computer Science - R&D*, 30(1):51–68, 2015.
- [6] K. Wellenzohn, M. Böhlen, S. Helmer, M. Reutegger, and S. Sakr. A Workload-Aware Index for Tree-Structured Data. To be published.
- [7] G. Zipf. *Human behavior and the principle of least effort: an introduction to human ecology*. Addison-Wesley Press, 1949.

7 Appendix

```
/**
 * Returns an iterable which lazily traverses a subtree rooted at root
 * in postorder.
 *
 * @param root: Root node of subtree we apply traversal on
 * @returns: Iterable for postorder tree walk
 */
Iterable<Tree> postOrder(Tree root) {
    return () -> {
        /* Stacks */
        Deque<Tree> s1 = new LinkedList();
        Deque<Tree> s2 = new LinkedList();
        s1.push(root);
        return new Iterator<Tree>() {
            @Override
            public boolean hasNext() {
                return s1.size() > 0 || s2.size() > 0;
            }
            @Override
            public Tree next() {
                while (s1.size() > 0 && (
                    s2.size() == 0 ||
                    isAncestor(
                        s2.peek().getPath(),
                        s1.peek().getPath()
                    )
                )) {
                    Tree n = s1.pop();
                    s2.push(n);
                    for (Tree child : n.getChildren()) {
                        s1.push(child);
                    }
                }
                return s2.pop();
            }
        };
    };
}
```

Figure 33: postOrder() implementation in Java.

```

/**
 * Higher order function, applies func to each element of iterable.
 *
 * @param func: The function to apply on each element of iterable
 * @param iterable: The iterable func is applied on
 * @returns: Resulting elements from applying func
 */
Iterable<R> map(Function<T,R> func, Iterable<T> iterable) {
    return () -> {
        Iterator<T> iterator = iterable.iterator();
        return new Iterator<R>() {
            @Override
            public boolean hasNext() {
                return iterator.hasNext();
            }
            @Override
            public R next() {
                return func.apply(iterator.next());
            }
        };
    };
}

```

Figure 34: map() implementation in Java.

```

/**
 * Higher order function that removes all elements from an iterable
 * not satisfying the predicate.
 *
 * @param predicate: The predicate that tests elements
 * @param iterable: The iterable whose members are tested against
 * @returns: An iterable with members satisfying the predicate
 */
Iterable<T> filter(Predicate<T> predicate, Iterable<T> iterable) {
    return () -> {
        Iterator<T> iterator = iterable.iterator();
        T n = null;
        return new Iterator<T>() {
            @Override
            public boolean hasNext() {
                nextIfNeeded();
                return n != null;
            }
            @Override
            public T next() {
                nextIfNeeded();
                T tmp = n;
                n = null;
                return tmp;
            }
            @Override
            private void nextIfNeeded() {
                while (n == null && iterator.hasNext()) {
                    T candidate = iterator.next();
                    if (predicate.test(candidate)) {
                        n = candidate;
                    }
                }
            }
        };
    };
}

```

Figure 35: filter() implementation in Java.