

NOBtree: A NUMA-Optimized Tree Index for Nonvolatile Memory

Zhaole Chu^{1b}, Peiquan Jin^{1b}, *Member, IEEE*, Yongping Luo^{1b}, Xiaoliang Wang^{1b}, *Member, IEEE*, and Shouhong Wan

Abstract—Nonvolatile memory (NVM) suffers from more serious nonuniform memory access (NUMA) effects than DRAM because of the lower bandwidth and higher latency. While numerous works have aimed at optimizing NVM indexes, only a few of them tried to address the NUMA impact. Existing approaches mainly rely on local NVM write buffers or DRAM-based read buffers to mitigate the cost of remote NVM access, which introduces memory overhead and causes performance degradation for lookup and scan operations. In this article, we present *NOBtree*, a new NUMA-optimized persistent tree index. The novelty of *NOBtree* is two-fold. First, *NOBtree* presents per-NUMA replication and an efficient node-migration mechanism to reduce remote NVM access. Second, *NOBtree* proposes a NUMA-aware NVM allocator to improve the insert performance and scalability. We conducted experiments on six workloads to evaluate the performance of *NOBtree*. The results show that *NOBtree* can effectively reduce the number of remote NVM accesses. Moreover, *NOBtree* outperforms existing persistent indexes, including TLBtree, Fast&Fair, ROART, and PACtree, by up to 3.23× in throughput and 4.07× in latency.

Index Terms—Nonuniform memory access (NUMA) effect, nonvolatile memory (NVM), tree index.

I. INTRODUCTION

THE NON-UNIFORM memory access (NUMA) architecture is a prevalent design in modern multicore systems [1], [2]. In a NUMA system, CPU cores and memory DIMMs are organized into clusters known as NUMA nodes, connected by internode links like the Intel Ultra Path Interconnect (UPI). Each processor can access either the memory on its own NUMA node or that of another, resulting in local/remote memory access. Local memory access is inherently faster than remote memory access. Such an asymmetry in accessing cost is known as the NUMA effect [1], [2], which impacts application performance substantially.

Nonvolatile memory (NVM) [3] is an emerging memory technology characterized by byte-addressability and persistent storage capabilities. It challenges the traditional storage hierarchy by bridging the gap between DRAM and SSD,

Manuscript received 30 July 2024; accepted 30 July 2024. This work was supported by the National Science Foundation of China under Grant 62072419. This article was presented at the International Conference on Compilers, Architectures, and Synthesis for Embedded Systems (CASES) 2024 and appeared as part of the ESWEEK-TCAD Special Issue. This article was recommended by Associate Editor S. Dailey. (*Corresponding author: Peiquan Jin.*)

The authors are with the School of Computer Science and Technology, University of Science and Technology of China, Hefei 230027, China (e-mail: czle@mail.ustc.edu.cn; jpq@ustc.edu.cn; ypluo@mail.ustc.edu.cn; wxl147@mail.ustc.edu.cn; wansh@ustc.edu.cn).

Digital Object Identifier 10.1109/TCAD.2024.3438111

prompting a fundamental reevaluation of storage system design principles. In response to this paradigm shift, significant efforts have been devoted to optimizing conventional data management techniques for NVM, including index structures, file systems, and key-value stores. Current NVM devices, e.g., Intel Optane Persistent Memory (“Optane” for short in this article) [4], are integrated into the system much like DRAM, thus also experiencing the NUMA effect. Due to the higher latency and lower bandwidth than DRAM, NVM is more susceptible to the NUMA effect. With NUMA becoming ubiquitous, mitigating its impact is crucial when devising NVM indexes.

So far, many NVM-aware indexes have been proposed [5], [6], [7]. These indexes have introduced various techniques to optimize performance on NVM, but the consideration of the NUMA effect is notably absent in many of these designs. Several works have focused on mitigating the NUMA effect of NVM indexes, such as Nap [1] and PACtree [8]. However, Nap needs additional DRAM buffers and is tailored for specific workloads, and PACtree has huge extra space overhead caused by the log. These overheads will become more influential in machines with more NUMA nodes. In addition, Nap can not work well under uniform workloads and always undergoes a degraded scan performance because of the existence of DRAM buffers.

In this article, we introduce *NOBtree*, a novel NUMA-optimized tree index tailored for NVM environments. *NOBtree* adopts a decoupled tree structure comprising a static read-optimized upper layer and a write-optimized bottom layer to mitigate the NUMA effect. Briefly, this article makes the following contributions.

- 1) We present *NOBtree*, a new NUMA-optimized NVM-aware tree index with a decoupled structure. *NOBtree* proposes two new designs, including per-NUMA replication and an efficient node-migration mechanism, to reduce remote NVM access.
- 2) We propose a dedicated NUMA-aware NVM allocator that supports round-robin, local, and specific NUMA node allocations. It incorporates a post-crash garbage collection (GC) mechanism to reduce the overhead of persistence and adopts a two-layer architecture to reduce the contention of threads, thus improving the performance of NVM allocation.
- 3) We conduct comprehensive experiments on a two-socket server equipped with two CPUs and real NVM devices to compare *NOBtree* with state-of-the-art NVM indexes. The results show that *NOBtree* can effectively reduce the number of remote NVM accesses. Moreover,

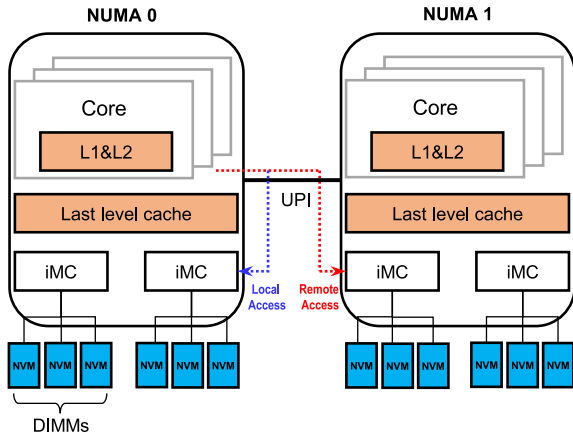


Fig. 1. Architecture of a typical NUMA system.

NOBtree outperforms existing persistent indexes, including TLBtree, Fast&Fair, ROART, and PACtree, by up to $3.23\times$ in throughput and $4.07\times$ in latency.

The remainder of this article is structured as follows. Section II summarizes related work. Section III details the design of NOBtree. Section V reports the experimental results. Section VI discusses the impact of Optane’s discontinuation. Finally, in Section VII, we conclude this article.

II. RELATED WORK

In this section, we first discuss the NUMA effect on NVM. Then, we summarize the recent advances in NVM-aware indexes and NUMA optimizations. Finally, we introduce techniques for NVM management.

A. NUMA Effect on NVM

NVM integrates into a system much like DRAM. Fig. 1 depicts a typical NUMA system. Given its higher latency and lower bandwidth than DRAM, NVM experiences more pronounced NUMA effects [9] than DRAM. Prior research [1], [8], [9] has underscored a noteworthy decrease in remote bandwidth for NVM. This effect is primarily attributed to the cache coherency protocol, particularly evident when threads on different NUMA nodes access the same NVM address. More specifically, today’s Intel processor architectures rely on the directory coherence protocol among NUMA domains. For NVM devices, the directory information is also stored in the 3D-XPoint media. When accessing the NVM address on a remote NUMA node, it will incur the coherence state change, causing a write operation to the directory. That means even a read operation will result in small write operations to NVM, causing the performance drop of remote NVM accesses. Consequently, it is necessary to avoid or reduce the NUMA impact in designing data management systems or storage systems tailored for NVM.

B. NVM-Aware Indexes

There have been a lot of indexes designed for NVM, either for NVM-only or DRAM-NVM hybrid architecture [5], [6], [10]. The main challenge of the NVM index is how to ensure crash consistency without too much

performance loss. The NVM index’s operation must guarantee crash consistency because of NVM’s feature of persistent storage. A commonly used approach is to use cacheline flush and fence instructions (e.g., *CLWB* and *SFENCE*) to guarantee the memory write order. However, using flush and fence instructions will lead to performance degradation. Previous works have introduced various techniques to optimize the performance of NVM indexes, such as fingerprints, indirect order arrays, selective persistence, and log-free node splitting mechanisms. However, the consideration of the NUMA impact is notably absent in many of these designs.

Numerous works have endeavored to optimize indexes for NVM, such as Fast&Fair [5], TLBtree [6], and ROART [11]. Fast&Fair introduced a novel mechanism to maintain the order of records inside a node by shifting slots and splitting/merging nodes in a failure-tolerable manner. TLBtree proposed a decoupled two-layer structure to accelerate index performance, which also inspires the structure of NOBtree. ROART presented a series of techniques to optimize ART on NVM, including entry compression and selective persistence. Overall, these indexes represent the recent advances in NVM-aware indexes.

However, within the extensive body of work dedicated to optimizing indexes for NVM, only a small portion has specifically focused on mitigating the NUMA effect. Nap, a black-box approach, transforms any NVM index into a NUMA-aware counterpart [1]. This transformation involves introducing an in-DRAM NUMA-aware layer (NAL) to expedite hot lookups and utilizing a local NVM buffer to absorb insertions, updates, and deletions. While Nap is a versatile method that demonstrates commendable performance on skewed workloads, the presence of read/write buffers complicates the operation process, leading to a decline in scan performance. Additionally, Nap is tailored for skewed workloads and may not handle uniform workloads as efficiently.

Node replication (NR) [2] is a widely recognized technique for addressing the NUMA effect in DRAM. It involves replicating data structures across NUMA nodes and utilizing a NUMA-aware shared log for synchronization. However, full replication consumes substantial memory, and the synchronization overhead is considerable, making it unsuitable for NVM-oriented indexes. Therefore, PACtree [8] divides the tree index into two parts: 1) the search layer and 2) the data layer. PACtree selectively replicates the search layer and employs per-NUMA logs to synchronize the replications. The per-NUMA log records the structural modifications of the data layer, and background updating threads will replay the log to update the search layer asynchronously. However, a drawback of the logging method is the additional NVM consumption for the log, as PACtree maintains a log file in every NUMA node. In Section V, where we assess each index’s NVM usage, we observe that PACtree requires 500 MB of NVM space for the per-NUMA log. This space overhead becomes significant in machines with multiple NUMA nodes.

C. NVM Allocator

Dynamic NVM allocation is crucial to building high-performance and scalable index structures. The memory

181 allocators are well-defined for DRAM to achieve high
 182 scalability and low latency. However, NVM allocators must
 183 also consider crash consistency, necessitating a reevaluation
 184 of design principles and implementations. Several NVM
 185 allocators have been developed, primarily falling into two
 186 categories based on crash consistency mechanisms: 1) log-
 187 based and 2) GC-based. Log-based allocators leverage logging
 188 to record all changes in metadata and memory addresses.
 189 Upon a crash, replaying the log rebuilds the allocator’s
 190 correct metadata, ensuring crash consistency. PAllocator [12],
 191 Poseidon [13], and NValloc [14] adopt the log approach.

192 Another mechanism for ensuring crash consistency is the
 193 GC mechanism, which provides an alternative to the write-
 194 ahead log approach. Unlike log-based allocators, GC-based
 195 allocators do not rely on persistent logs for maintaining crash
 196 consistency. Instead, they rebuild metadata by traversing the
 197 heap from a predefined persistent root pointer. GC-based allo-
 198 cators offer faster allocation and deallocation than log-based
 199 allocators because they avoid the overhead of persisting logs
 200 for every allocation or deallocation operation. However, they
 201 may incur longer recovery times after a crash, as they need
 202 to traverse the entire memory space to reconstruct metadata.
 203 Ralloc [15] and DCMM [11] employ the GC mechanism.

204 Previous studies have emphasized the importance of
 205 NUMA-aware NVM allocation to minimize costs [1], [8].
 206 However, existing approaches typically rely on PMDK to
 207 build customized NVM managers, which may suffer from
 208 performance and scalability limitations. To address the short-
 209 coming, we propose to develop a dedicated high-performance
 210 NVM allocator with NUMA-awareness. This allocator aims
 211 to optimize NVM allocation while considering the NUMA
 212 effect, thereby enhancing the performance and scalability of
 213 NVM-based indexes.

214 III. DESIGN OF NOBTREE

215 In this section, we first introduce the motivation and overall
 216 structure of NOBtree. Then, we describe the optimizations of
 217 NOBtree for the NUMA architecture. Finally, we present the
 218 architecture and operation of the NVM allocator.

219 A. Motivation of NOBtree

220 To explore the impact of the NUMA effect on NVM
 221 indexes, we first perform a comparative experiment involving
 222 TLBtree [6] and Fast&Fair [5], along with a variant of TLBtree
 223 named TLBtree-NR. In TLBtree-NR, the tree undergoes repli-
 224 cation across all NUMA nodes, ensuring that all operations
 225 exclusively interact with the tree within the local NVM,
 226 eliminating all the remote NVM access. Note that TLBtree-
 227 NR does not incorporate the NR algorithm [2]; thus, it does
 228 not support crash consistency. We utilize TLBtree-NR solely
 229 to gauge the highest-achievable read performance by TLBtree
 230 without leveraging DRAM.

231 Fig. 2 depicts the performance gap between TLBtree and
 232 TLBtree-NR. The experiment took place on a two-socket
 233 server equipped with two Intel Xeon Gold 6242R CPUs,
 234 each having 20 physical cores. Thread assignment follows a
 235 configuration where if the number of threads is below 20, they

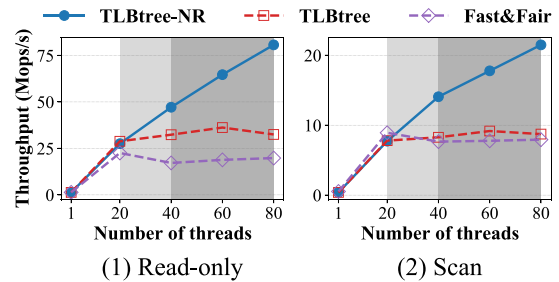


Fig. 2. Performance gain of reducing the NUMA effect: TLBtree-NR uses replications across NUMA nodes to reduce the NUMA effect, showing much higher performance than TLBtree and Fast&Fair, which incur NUMA effects with the increase of threads.

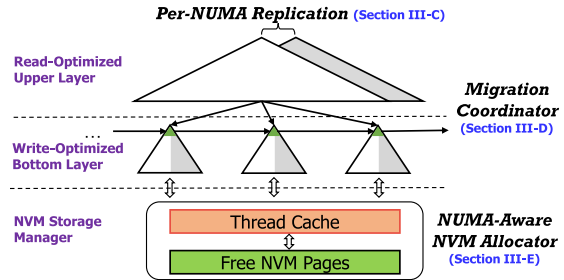


Fig. 3. Structure of NOBtree.

are allocated to a single node. However, if the thread count
 236 surpasses 20, threads are distributed across two nodes, and
 237 hyper-threading is enabled when exceeding 40. We employ
 238 Yahoo! cloud serving benchmark (YCSB)-like workloads to
 239 assess the lookup and scan performance, with keys following
 240 a Zipfian distribution parameterized at 0.99.
 241

242 As shown in Fig. 2, TLBtree and Fast&Fair struggle to scale
 243 beyond a single node as thread count increases beyond 20.
 244 Conversely, TLBtree-NR consistently demonstrates robust
 245 scalability even though threads are distributed across the two
 246 sockets. The main reason for the poor scalability of TLBtree
 247 and Fast&Fair is the limited bandwidth of UPI and the
 248 cache coherency mapping cost across sockets. These results
 249 underscore the substantial NUMA impact on NVM index
 250 performance, revealing significant potential for enhancement.
 251 Therefore, there is a need to design a NUMA-optimized NVM
 252 index.

253 B. Overall Structure of NOBtree

254 Fig. 3 illustrates the overall structure of NOBtree. NOBtree
 255 includes two layers: 1) a read-optimized upper layer and
 256 2) a write-optimized bottom layer. The upper layer is a read-
 257 optimized static structure used to rapidly locate the target
 258 subindex in the bottom layer during a read/write operation.
 259 The bottom layer is a series of write-optimized subindexes
 260 linked by pointers. The rationale behind this design is rooted
 261 in the observation that B+ tree nodes closer to the root
 262 are read-dominated, while write operations (insertion/deletion)
 263 primarily affect the bottom levels. The structure of these two
 264 layers has been designed to adapt to the features of the NVM.
 265 Moreover, both the read-optimized upper layer and the write-
 266 optimized bottom layer are designed to reduce the NUMA
 267 effect (the details will be described below).

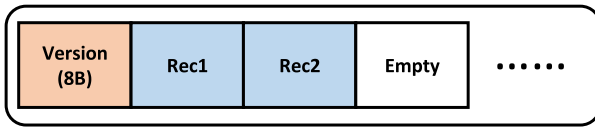


Fig. 4. Leaf-node structure in the upper layer.

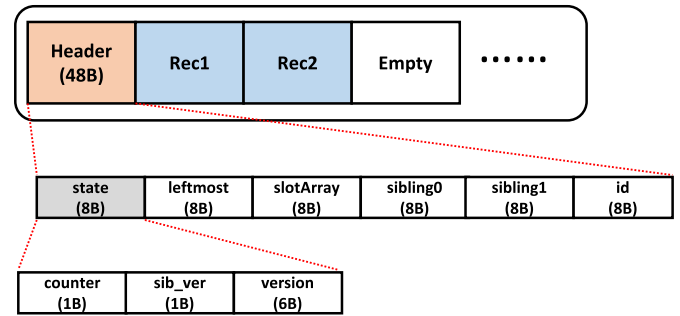


Fig. 5. Node structure in the bottom layer.

268 1) *Read-Optimized Upper Layer*: The upper layer of
 269 NOBtree can take the form of any data structure capable of
 270 delivering high-search performance. In our current implemen-
 271 tation, we opt for a customized k -ary tree-like index to serve
 272 as the upper layer. This layer is designed as a static structure
 273 with no structural modifications, utilizing a contiguous node
 274 array and eliminating all pointers. The main features of the
 275 upper layer are *static structure*, *read-optimized node layout*,
 276 and *per-NUMA replication*.

277 *Static Structure*: The static structure means there is no
 278 structural modification in the upper layer. We employ this
 279 design mainly for two reasons. First, the static structure is
 280 beneficial for read operations. Using the static structure means
 281 that we can neglect the write operations and employ a fully
 282 read-optimized index layout. Second, according to the previous
 283 work [6], write operations are rare in this layer. That is because
 284 the write operations in a B+-tree's inner node are triggered by
 285 the node split occurring in the bottom layer. The number of
 286 write operations decreases exponentially from bottom to top.
 287 Thus, it is reasonable to sacrifice the write performance for
 288 the read performance in the upper layer.

289 Insertions to the upper layer may fail if the target leaf node
 290 is full because there will not be any structural modifications.
 291 Such failures may lead to performance decline by lengthening
 292 the traversal path if the target subindex is not inserted into
 293 the upper layer. To address the performance issues stemming
 294 from a full and imbalanced upper layer, NOBtree will rebuild
 295 the upper layer once a predefined threshold is reached. The
 296 threshold will influence the frequency of reconstructions. It's
 297 crucial to strike a balance with this threshold because if
 298 rebuilding occurs too frequently, it can degrade NOBtree's
 299 performance due to the high cost of reconstruction. Therefore,
 300 we choose a well-tuned threshold based on the experimental
 301 result to achieve the highest performance. We will give a
 302 detailed description of the rebuilding process in Section IV-C.

303 *Read-Optimized Node Layout*: The upper layer of NOBtree
 304 plays a crucial role in swiftly locating the subindex, influ-
 305 encing the performance of both read and write operations.
 306 We utilize an array-based k -ary tree as the upper layer,
 307 categorizing nodes into immutable inner nodes and gapped
 308 leaf nodes. Inner nodes are ordered arrays and are always
 309 100% full with no pointers. This design choice enables them to
 310 accommodate a maximum number of keys, thus lowering the
 311 tree height. The inner nodes are determined at the rebuilding
 312 time and will not change until the next reconstruction. Leaf
 313 nodes contain records with a key and a pointer to subindexes
 314 and employ a gapped structure. Fig. 4 shows the leaf-node
 315 structure. We reserve a few empty slots in each leaf node
 316 (preallocating some gaps when rebuilding the top layer). The
 317 empty slots in a leaf node can absorb new records generated

by a subindex splitting of the bottom layer. The version is used
 for optimistic concurrency control like previous works [16].
 In our implementation, the slot number is set to 15, and we
 reserve three empty slots when building the upper layer. Note
 that the size of inner nodes and leaf nodes is both 256B,
 aligning with Optane's internal page size, which can reduce the
 number of NVM accesses when traversing the tree. Moreover,
 all the nodes are allocated in contiguous memory space. Such
 a design can reduce the number of TLB misses during the
 operation process, leading to higher performance.

Per-NUMA Replication: While the read-optimized static
 structure of the upper layer can offer the ability to rapidly
 locate the subindex. Threads from other sockets still experi-
 ence reduced search performance due to remote NVM access if
 the upper layer resides in only one NUMA node. To tackle this
 issue, we propose replicating the upper layer across all NUMA
 nodes. During an index operation, each thread accesses the
 upper layer on its local socket, effectively reducing the costly
 remote NVM access. We will give a detailed description in
 Section III-C.

2) *Write-Optimized Bottom Layer*: The bottom layer
 employs a group of subindexes, as illustrated in Fig. 3, which
 are indexed by the upper layer. The subindex is a write-
 optimized persistent B+-tree. All the roots of the subindexes
 are linked via pointers. This linked structure enables the asyn-
 chronous update of the upper layer while ensuring correctness
 simultaneously. The subindex incorporates various techniques
 to adapt to NVM, including a node size optimized for NVM,
 physically unsorted but logically sorted keys, and a log-free
 splitting/merging mechanism. We will describe the design of
 the bottom layer next.

NVM-Friendly Structure: The bottom layer of NOBtree is
 tailored to handle the majority of insert operations, neces-
 sitating the creation of an NVM-friendly structure. Fig. 5
 illustrates the node structure in the bottom layer. Each node
 contains a header and several slots to store records. The header
 contains a *state* field, a *leftmost* child pointer, a *slotArray*,
 two sibling pointers, and a node *id*. The *state* field contains
 a *counter* (1 byte), a *sib_ver* (1 byte) and a node *version* (6
 bytes). The *leftmost* pointer points to the child node whose
 keys are all less than the current node's minimal key. It is
 valid only for the inner node. The *slotArray* records the keys'
 order information in the node. The two sibling pointers and the
sib_ver are used to implement the log-free splitting/merging

mechanism. The *counter* counts the total number of records in this node. The node *id* is used for node migration. It is assigned when the node is created. The node *version* is used for concurrency control. We adopt optimistic concurrency control in our implementation. The header's size is less than 64 bytes, and the size of *state* is 8 bytes. NOBtree's bottom layer employs several techniques to improve the efficiency of write operations.

The node size is set to 256B, aligned with Optane's internal page size, ensuring at most one NVM access when traversing a node. In addition, the records within a node are physically unsorted but logically sorted. Such a design eliminates the need to shift records while inserting new keys. However, a completely unordered key arrangement harms searching and node-splitting operations. To cope with the problem, a *slotArray* is introduced to maintain the order information of the keys. For instance, *slotArray*[3] = 4 indicates that the third smallest key is stored in the fourth slot. In our implementation, each node comprises only 12 records, allowing us to represent the position information using just 4 bits. Therefore, we can embed the *slotArray* within an 8-byte field and update it atomically. This approach significantly improves search performance without incurring additional costly NVM writes.

Also, we implement a log-free splitting/merging mechanism using shadow sibling pointers, eliminating the logging overhead typically employed for ensuring crash consistency. This approach utilizes two pointers in one node to reference the sibling node. A parameter, denoted as *sib_ver*, indicates the functional pointer, i.e., the pointer pointing to the actual sibling node. During a node splitting, we initially allocate a new node and copy half of the records from the splitting node to this new node. Subsequently, we install the new node into the nonfunctional pointer. Finally, we visualize the new nodes by atomically updating the *sib_ver*. This strategy ensures crash consistency without resorting to the costly logging mechanism.

NUMA Optimization: The NVM-friendly structure enhances the insert performance of NOBtree. However, a notable performance degradation will occur when all subindexes are placed in only one NUMA node while access threads are distributed across multiple nodes. This is attributed to two main factors. First, half of the threads consistently undergo remote NVM access, experiencing a costly operation process. Second, previous work [9] reveals that reading the same address from multiple sockets (denoted as the *near-far* access pattern) achieves very low-NVM bandwidth due to cache coherency protocol. A naive solution to this problem is replicating the entire bottom layer across all NUMA nodes [2], enabling each thread to access the local bottom layer and thus reducing remote NVM access. However, this approach will introduce huge extra NVM consumption and synchronization costs between the replications.

To cope with the problem without introducing extra overhead, we propose randomly distributing the nodes in the subindex across all NUMA nodes, as depicted in Fig. 3. This approach helps avoid *near-far* accesses, thereby improving index performance. We also introduce a node migration mechanism for the bottom layer. NOBtree migrates nodes within a subindex to the socket that accesses them most

frequently, thereby reducing remote NVM access. We will give a detailed description in Section III-D. It's worth noting that these designs necessitate specialized NVM allocation methods, e.g., round-robin allocation and allocation to specific NUMA nodes. However, existing NVM allocators do not support such allocation patterns. Therefore, we propose a new NVM storage manager to fulfil these requirements.

3) *NVM Storage Manager:* The NVM storage manager (NSM) is responsible for allocating NVM. NOBtree relies on a new NUMA-aware NVM allocator, as shown in Fig. 3. To optimize allocation performance, we adopt a two-layer structure to organize NVM storage: 1) a thread-cache layer and an 2) NVM-free-page layer. The thread-cache layer is used to minimize contention among concurrent NVM allocations, and the NVM-free-page layer manages all free NVM pages. The details of NSM will be discussed in Section III-E.

C. Per-NUMA Replication

The NUMA optimization in the upper layer is per-NUMA replication, where the upper layer of NOBtree is replicated in each NUMA node. Leveraging NOBtree's decoupled structure, it can tolerate inconsistency between the upper and bottom layers, allowing it to replicate the upper layer in each NUMA node. This design enables threads to access the upper layer on the local socket, thus mitigating the costly remote NVM access.

While the idea of NUMA NR is straightforward, it is challenging to synchronize the replicas across NUMA nodes. A common approach is to use a shared log for synchronization [2], [8]. However, relying on a log to replay operations to the upper layer introduces vulnerabilities to NOBtree. Since we use only one thread to process the log entry, the insertion of a newly split node into the upper layer might be delayed, especially during periods of high-splitting activity. While delayed updates to the upper layer do not compromise the correctness of subsequent operations, they may affect performance stability. If subsequent operations need to access the split node before inserting it into the upper layer, they must traverse the subindex list to locate the target subindex. Also, the log-based method requires extra NVM space for logs; for example, PACTree needs 500 MB of NVM space to store the log on each node. Considering the potential drawbacks, it is more practical to adopt a synchronous approach for insertions into the upper layer.

The second reason that we do not delay the insertion to the upper layer is that structural modifications in the upper layer are less costly than anticipated, making immediate insertion feasible. Hence, there's no need to employ the logging approach for delayed insertion. In summary, we utilize a synchronous approach to update the upper layer for the aforementioned reasons. Upon a subindex root node split, the new root node is inserted into the upper layers immediately.

D. Migration Coordinator

The NUMA optimization in the bottom layer mainly relies on the node migration mechanism, which can reduce the remote NVM access in the bottom layer. We propose a

475 migration coordinator for the bottom layer to handle the
 476 node migration. This optimization is based on an observation:
 477 *each thread has different data-accessing localities*. Traditional
 478 evaluations often distribute the workload equally among all
 479 threads executing the task, assuming the same hotspot data
 480 across them. However, systems like MySQL and Microsoft
 481 SQL Server typically assign a single core to service requests
 482 from a single session, resulting in different threads having
 483 distinct hotspots. This distribution characteristic presents an
 484 opportunity to design a node migration mechanism for the
 485 bottom layer in NOBtree. Leveraging the diverse hotspots of
 486 each thread, we can migrate nodes in the bottom layer to the
 487 socket that accesses them most frequently, thereby reducing
 488 remote NVM access. The node migration in the bottom
 489 layer involves two primary steps: 1) identifying the node to
 490 be migrated and putting it into a global queue and 2) the
 491 migration coordinator periodically retrieves node information
 492 from the global queue to execute migrations in the background.
 493 Inspired by the producer-consumer model, the process of node
 494 migration utilizes a global queue to record information about
 495 nodes slated for migration. When an operation thread accesses
 496 a specific node, it checks whether the node requires migration.
 497 If so, its information will be put into the global queue. A
 498 dedicated thread (migration coordinator) periodically polls the
 499 global queue and executes node migrations. Note that the
 500 migration coordinator runs in the background, minimizing
 501 disruptions to normal operations. To implement the node
 502 migration mechanism correctly and efficiently, we have to
 503 address the following challenges.

504 *Identifying the Node to Be Migrated:* NOBtree utilizes
 505 historical access statistics to determine nodes eligible for
 506 migration in the bottom layer. It monitors access to all
 507 leaf nodes, identifying nodes meeting two criteria: 1) they
 508 are considered hot nodes and 2) most accesses are remote
 509 NVM accesses. A node qualifies as a hot node if its
 510 total accesses surpass a threshold within a monitoring win-
 511 dows. Subsequently, we calculate the percentage of accesses
 512 from remote nodes. If remote access constitutes the major-
 513 ity, we add the node's information to the global queue.
 514 This information includes the node's address, its parent and
 515 previous sibling addresses, and the target migrating NUMA
 516 node. The migration coordinator fetches this information
 517 periodically from the global queue and performs the migra-
 518 tion accordingly. Other methods of identifying hot data can
 519 also be employed, maintaining the core concept of node
 520 migration.

521 However, maintaining access statistics can pose a scalability
 522 bottleneck due to the need for frequent modifications, which
 523 incur additional writes [17]. Given the expensive cost of NVM
 524 writes, we utilize an in-DRAM hash table to record the access
 525 information rather than maintaining statistics within the nodes
 526 themselves. The key in the hash table corresponds to the node
 527 *id*, as discussed in Section III-B2. To further alleviate the
 528 impact on the performance, we implement a sampling method,
 529 recording access statistics every ten operations. These statistics
 530 include the total number of local and remote accesses, enabling
 531 us to calculate the node's 'hotness' and determine whether
 532 migration is necessary.

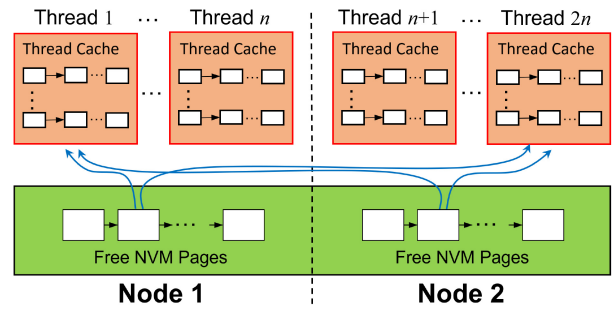


Fig. 6. Structure of the NSM.

533 *Coordinating Migration With Other Operations:*
 534 Coordinating migration with other operations is a key
 535 responsibility of the migration coordinator. Since node
 536 migration occurs in the background, structural modifications
 537 may occur in the migrating node, its parent, or its previous
 538 sibling node. If any structural modifications occur during
 539 migration, the migration process is aborted to maintain data
 540 integrity. During a migration, the node is first locked to block
 541 subsequent insertions. Then, we copy the node to the target
 542 NUMA node. Next, we update the node's parent and previous
 543 node to point to the new node. Note that we need to guarantee
 544 crash consistency during the two update operations. Because
 545 the two operations can not be done atomically, we use logs to
 546 ensure crash consistency. Finally, we mark the original node
 547 as obsolete and free it once no more readers access it. The
 548 migration does not block any reader during this process. We
 549 adopt an epoch-based reclamation strategy [18] to reclaim the
 550 obsolete node. Note that the migration relies on the dedicated
 551 NUMA-aware NVM allocator, which allows NVM allocation
 552 on specific NUMA nodes. In the subsequent section, we will
 553 present the details of the allocator.

E. NUMA-Aware NVM Allocator

554
 555 NVM allocation profoundly affects the insert performance
 556 of indexes. Prior works [8], [11], [14] have demonstrated that
 557 different NVM allocators can result in a performance gap of
 558 25% ~ 50% for indexes. Thus, NOBtree presents a NUMA-
 559 aware NVM allocator to enhance the insert performance of
 560 the index.

561 NVM allocators are generally categorized into two types:
 562 1) log-based [12], [13], [14] and 2) GC-based [11], [14], [15].
 563 The former relies on logging to persist memory changes and
 564 heap metadata, ensuring the atomicity of operations in case
 565 of failure. However, this approach introduces additional per-
 566 sistence overhead. To mitigate this overhead, recent allocators
 567 utilize GC to rebuild heap metadata post-crash by scanning the
 568 entire memory space [15]. The GC-based allocators generally
 569 offer faster allocation speeds than the log-based allocators.
 570 Therefore, we adopt the GC-based NVM allocator.

571 Fig. 6 depicts the overall structure of the NSM, which
 572 consists of two main parts. The bottom part operates as an
 573 NVM space manager, responsible for managing the entire
 574 NVM space. Meanwhile, the upper part serves as the thread
 575 cache for NVM allocation and deallocation operations.
 576

Algorithm 1: NVM Allocation

```

1 Function nvm_numa_malloc(size, node):
2   if node < number_of_numa_node then
3     // allocate to the given NUMA node
4     return thread_cache.slab[node].slab_alloc(size);
5   end
6   // allocate NVM according to the pre-set mode
7   return thread_cache.slab[mode].slab_alloc(size);
8 Function slab_alloc(size):
9   size_class = get_size_class(size);
10  address = freelist[size_class].pop();
11  if !address then
12    // get a NVM page according to the pre-set
13    mode
14    page = get_nvm_page(mode);
15    fill the freelist with the requested NVM page;
16    address = freelist[size_class].pop();
17  end
18  return address;
19 Function get_nvm_page(i):
20  if i < number_of_numa_node then
21    return NVM page on NUMA node i;
22  else if i == LOCAL then
23    numa_node = get_numa_node();
24    return NVM page on NUMA node numa_node;
25  else
26    return NVM page in a round-robin way;

```

576 The bottom part manages the per-NUMA NVM pool,
577 treating the NVM space as a free-page list with a default page
578 size of 256KB. The thread cache comprises segregated lists,
579 each containing blocks of the same size. Each block is used to
580 serve a single allocation request. This two-layer structure helps
581 alleviate contention for concurrent NVM allocations, resulting
582 in high performance and scalability.

583 NSM offers three allocation modes: 1) *local* (allocating
584 pages from the local node); 2) *round-robin* (allocating pages
585 among NUMA nodes in a round-robin manner); and 3) *node-*
586 *specific* (allocating pages from a specified node). Users initiate
587 the mode during the creation of NSM. There are a certain
588 number of slab managers inside a thread cache. Each contains
589 segregated free lists and serves different modes of NVM allo-
590 cation. Algorithm 1 shows the process of the NVM allocation.
591 The allocation needs to provide two parameters, *size* and
592 *node*. The slab manager responsible for the specified node is
593 selected to allocate NVM. The slab manager first gets the size
594 class of the request. Then, it retrieves an address from the
595 corresponding free list based on the size class. If the free list
596 returns an address successfully, the allocation is completed.
597 Otherwise, the slab manager requests a free NVM page from
598 the bottom NVM space manager, fills the free list with the
599 page, and allocates memory from it.

600 In the NVM allocation process, thread contention arises
601 primarily when threads request a page from the bottom
602 NVM space manager. However, this contention is infrequent
603 compared to the overall NVM allocations, allowing the NSM
604 to deliver high-allocation performance. Additionally, only the
605 page allocation requires persistence, while small allocations
606 from the thread cache do not need to persist any metadata.
607 These two features contribute to the high performance and

Algorithm 2: Lookup Operation

```

1 Function Lookup(key):
2   num_node = get_numa();
3   // identify current thread's NUMA node
4   bottom_root = uptree[numa_node].Up_Lookup(key);
5   while bottom_root.max_key < key do
6     | bottom_root = bottom_root.get_sib();
7   end
8   return Bottom_Lookup(bottom_root, key);
9 Function Bottom_Lookup(node, key):
10  cur = node;
11  while cur.leftmost != NULL do
12    | cur = cur.get_child(key);
13  end
14  return cur.get_child(key);
15 Function get_child(key):
16  retry: old_version = node.version;
17  if key > node.get_sib(key).min_key then
18    | return node.get_sib(key).get_child();
19  end
20  if leftmost != NULL then
21    // search in the inner node
22    pos = get_pos(key);
23    if pos == 0 then
24      | ret = leftmost;
25    else
26      | ret = records[slotArray[pos - 1]].val;
27  else
28    // search in the leaf node
29    pos = get_pos(key);
30    if records[pos].key != key then
31      | ret = NULL;
32    else
33      | ret = records[pos].val;
34  if node.version != old_version then
35    | goto retry;
36  end
37  return ret;

```

scalability of NSM. Moreover, the three allocation modes 608
provided by the NSM offer users greater flexibility in utilizing 609
NVM. 610

IV. OPERATIONS OF NOBTREE 611

In this section, we describe NOBtree's operations, including 612
lookup, insert, and the upper layer's rebuild. Due to the new 613
structure and NUMA-aware designs of NOBtree, all these 614
operations need to be redesigned. 615

A. Lookup 616

Algorithm 2 outlines the lookup process in NOBtree, which 617
comprises two main steps: 1) searching in the upper layer 618
to locate the subindex and 2) searching in the subindex to 619
retrieve the desired result. We initiate the lookup process by 620
selecting the local upper tree and then searching for the key 621
within this upper tree to obtain the root of the subindex (lines 622
2 and 3). This search process in the upper layer is similar to 623
the procedure employed in the traditional B+-tree. Starting 624
from the root node, we recursively traverse the inner nodes 625
until reaching the leaf node. Subsequently, we get the root 626
of the target subindex. However, the obtained result may be 627

628 incorrect due to subindex splits. To correct this, we traverse the
629 linked list to get the true root of the target subindex (lines 4–
630 6). Then we execute the lookup operation within this subindex
631 and return the result (line 7).

632 The lookup process of the upper layer is similar to that in a
633 traditional B+-tree. However, there are some distinctions due
634 to the specific characteristics of the upper layer’s structure.
635 The inner node of NOBtree’s upper layer is pointless to
636 accommodate more keys. So when searching in the inner node,
637 we will utilize a formula to calculate the position of the child
638 node. This approach is feasible due to the static structure of
639 the upper layer, where all nodes are allocated in a contiguous
640 space. Consequently, the offset of each node remains fixed,
641 and we can determine it by calculating the offset. Furthermore,
642 the structure of the upper layer remains unchanged until we
643 rebuild it, ensuring the reliability of this calculation method.

644 Similarly, the lookup process of the bottom layer, outlined in
645 Algorithm 2 (lines 8–34), comprises two main steps: 1) locat-
646 ing the leaf node and 2) conducting a local search within the
647 leaf. When probing a node, the process begins by snapshotting
648 the version (line 15). Subsequently, we verify whether the node
649 has undergone a split. If a split has occurred before accessing
650 the node, we switch to its next sibling node (lines 16–18). If
651 the node is an inner node, we find the first position where the
652 key exceeds the target key. This search is facilitated using the
653 *slotArray* (lines 19–24). Conversely, if the node is a leaf node,
654 we identify the first position where the key is greater than or
655 equal to the target key. We then check if the key matches the
656 target key (lines 25–30). Finally, we verify the node’s version.
657 If it has changed during the search, the process is retried.
658 Otherwise, the result is returned.

659 B. Insertion

660 The main steps of the insert operation are similar to those of
661 the lookup operation. We begin by searching the local upper
662 layer to locate the target subindex, followed by inserting the
663 key into the subindex. Subsequently, we check whether we
664 need to rebuild the upper layer. If reconstruction is necessary, it
665 is executed in the background (lines 11–13). During insertion
666 to the subindex, node splitting may occur in the subindex.
667 If the root of the subindex splits and we need to insert the
668 newly generated root into the upper layer, we will try to
669 insert the new root into all upper layers (lines 14–17). If the
670 insertion to the upper layer fails, we will store the new root
671 in a temporary list (lines 18–20). We use the temporary list
672 to accelerate the rebuilding process. Further details will be
673 provided in Section IV-C.

674 Algorithm 4 offers a detailed overview of the insert opera-
675 tion within a node of subindexes. Initially, the node is locked
676 to prevent subsequent insertions. If the node has been split
677 just before the insertion attempt, we switch to the next node
678 to perform the insertion (Algorithm 4, lines 3–6). If the node
679 is full, a new node is created, and the keys are rearranged.
680 Notably, the split operation is log-free due to the utilization
681 of shadow sibling pointers. To ensure crash consistency, any
682 updates to the node’s header are followed by *CLWB* and
683 *SFENCE* instructions. Node splits may propagate from the leaf

Algorithm 3: Insert Operation

```

1 Function Insert(key, val):
2   num_node = get_numa();
3   bottom_root = uptree[numa_node].Up_Lookup(key);
4   query_time = query_time + 1;
5   while bottom_root.max_key < key do
6     | bottom_root = bottom_root.get_sib();
7     | goes_step = goes_step + 1;
8   end
9   res = Bottom_Insert(bottom_root, key, val);
10  average_goes_step = goes_step/query_time;
11  if average_goes_step > threshold then
12    | rebuild_upper();
13  end
14  if res.flag = TRUE then
15    | // the root of the subindex splits
16    | for i = 0; i < num_numa; i ++ do
17      | succ = uptree[i].Up_Insert(res.key, res.val);
18    | end
19    | if succ = FALSE then
20      | // failed to insert to upper layer
21      | mutable.append(res.key, res.val);
22    | end
23  end
24  Function Bottom_Insert(node, key, val):
25    res = recursive_insert(node, key, val);
26    if res.is_split then
27      | // node split spreads to the root
28      | if res.level < threshold then
29        | create a new root and update the upper layer;
30      | else
31        | return {TRUE, res.split_k, res.split_node};
32    | end
33    return {FALSE, 0, NULL};
34  Function recursive_insert(node, key, val):
35    if node.leftmost != NULL then
36      | child = node.get_child(key);
37      | res = recursive_insert(child, key, val);
38      | if res.is_split then
39        | return node.store(res.split_k, res.split_node)
40      | return {node.level, FALSE, 0, NULL};
41    else
42      | return node.store(key, val)

```

to the root. If we find that the root of the subindex splits, 684
we will check whether the height of the subindex exceeds 685
the predefined threshold. If not, we will create a new root to 686
accommodate the splitting node and update the upper layer 687
to point to this new root (Algorithm 3, lines 25 and 26). 688
Otherwise, we will insert the newly created node to the upper 689
layer (Algorithm 3, lines 14–17). 690

691 C. Rebuilding the Upper Layer

692 Rebuilding the upper layer starts with collecting all the 693
roots of subindexes. Then, based on the number of records, 694
we calculate the number of leaf nodes and the tree height. 695
Subsequently, we fill the leaf nodes with the records and 696
recursively construct inner nodes until the root is created. 697

698 To mitigate the overhead of rebuilding the upper layer, we 699
implement several optimizations to reduce the frequency and 700
associated costs of this process, thus enhancing the index’s 701
stability and efficiency. 702

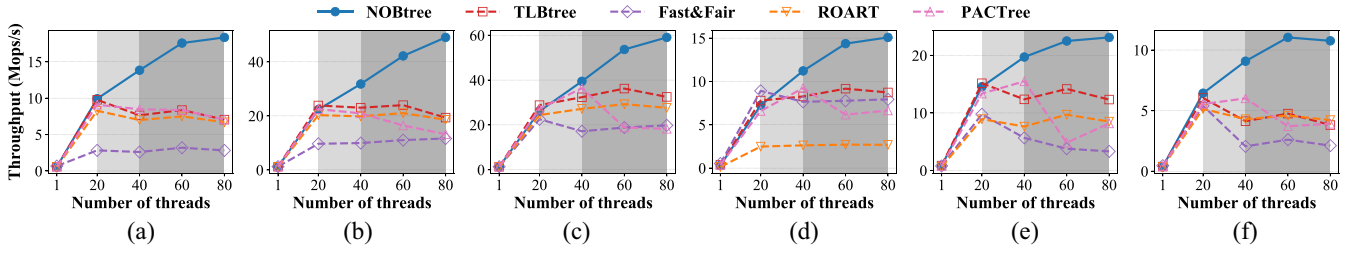


Fig. 7. Throughput under various workloads. (a) Write-heavy. (b) Read-heavy. (c) Read-only. (d) Scan. (e) Read-Modify-Write. (f) Write-only.

Algorithm 4: Node Insert Operation

```

1 Function store(key, val):
2   Lock();
3   if key > node.get_sib().min_key then
4     Unlock();
5     return node.get_sib().store(key, val);
6   end
7   if count == split_threshold then
8     // node needs to split
9     split_node ← create new node;
10    split_k ← the middle key;
11    rearrange the keys in the old node and insert key;
12    update the header of split_node;
13    clwb();sfence();
14    update the header of current node;
15    clwb();sfence();
16    Unlock();
17    return {level, TRUE, split_k, split_node};
18  else
19    insert the key;
20    Unlock();
21    return {level, FALSE, 0, NULL};

```

701 First, although the upper layer is static, its leaf node reserves
702 some empty slots to absorb limited insertions. Such a design
703 can effectively reduce the total number of rebuilds. Second,
704 instead of initiating a rebuilding process immediately after a
705 failed insertion, we defer this action until the average probing
706 length in the linked list exceeds a predefined threshold. As
707 shown in Algorithm 3 (lines 10–13), we use the parameter
708 *average_goes_step* to record the average probing length along
709 the linked list. Once this parameter exceeds the predefined
710 threshold, indicating a significant performance degradation, the
711 upper layer is rebuilt. This approach ensures that rebuilding
712 is only triggered when failed insertions noticeably impact
713 performance. As a result, the frequency of rebuilding occur-
714 rences is reduced. In our evaluation, we observed less than 20
715 reconstructions after inserting 200 million keys.

V. PERFORMANCE EVALUATION

717 We conduct experiments on a 2-socket server equipped
718 with two Intel Xeon Gold 6242R CPUs, each featuring 20
719 cores. The server contains 256-GB DRAM and 2048-GB Intel
720 Optane DC persistent memory, equally distributed over two
721 sockets. We configure all Optane modules to App-Direct
722 mode and create a DAX-aware ext4 file system. Then, we
723 mount the file system using the DAX option.

724 *Competitors:* We conduct a comparative evaluation
725 of NOBtree against four NVM-oriented indexes:
726 1) Fast&Fair [5]; 2) PACTree [8]; 3) ROART [11]; and
727 4) TLBtree [6]. We leverage their respective open-source
728 codes for our assessment. It is noteworthy that PACTree also
729 incorporates measures to mitigate the NUMA effect. Nap is
730 not considered in the comparison as it is a black-box approach
731 applicable to any index, which is orthogonal to our design,
732 and its use of buffers may result in performance degradation
733 for scan workloads.

734 *Workloads:* The dataset consists of 200M randomly gener-
735 ated integers (8 bytes long). There are mainly six kinds
736 of YCSB-like workloads, including *read-only* (100% lookup),
737 *read-heavy* (95% lookup and 5% insert), *write-heavy* (50%
738 lookup and 50% insert), *write-only* (100% insert), *read-*
739 *modify-write* (50% lookup and 50% read-modify-write), and
740 *scan* (100% scan). The YCSB is an open-source benchmarking
741 suite for evaluating the maintenance and retrieval capabilities
742 of computer programs [19]. It is widely used to compare the
743 performance of database systems and database indexes. We
744 made slight modifications to its core workload generator by
745 replacing all update operations with insert operations and then
746 generated the experimental workloads. In these workloads, all
747 the lookup and scan operations follow a Zipfian distribution
748 with a parameter of 0.99. The experiments scale from 1 thread
749 to 80 threads. When the number of threads is less than 20,
750 they are assigned to a single node. However, if the thread
751 number exceeds 20, threads are distributed across two nodes,
752 and hyper-threading is enabled when exceeding 40.

A. Throughput

754 In this experiment, we evaluate the throughput of each
755 index under various workloads. Fig. 7 shows the throughput
756 of NOBtree and other indexes under the six workloads. Across
757 the six workloads, NOBtree consistently outperforms other
758 indexes. Specifically, under 80 threads, NOBtree achieves
759 3.10 \times , 3.13 \times , 2.41 \times , 2.32 \times , 2.86 \times , and 3.06 \times higher
760 throughput compared to other indexes on average under the six
761 workloads, respectively. The high-read-throughput of NOBtree
762 is attributed to the per-NUMA upper-layer replication and
763 the node migration mechanism in the bottom layer, which
764 can reduce the costly remote NVM accesses effectively.
765 Additionally, the dedicated NUMA-aware NVM allocator con-
766 tributes to NOBtree’s superior insert performance. Although
767 rebuilding the upper layer is time-consuming, the rebuilding
768 frequency is rather low because we have reserved some empty

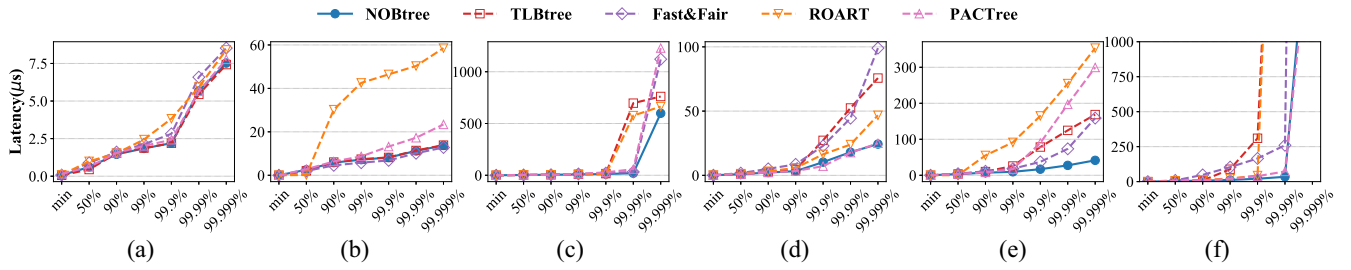


Fig. 8. Tail latency under skewed workloads (20t: 20 threads, 40t: 40 threads). (a) Lookup(20t). (b) Scan(20t). (c) Insert(20t). (d) Lookup(40t). (e) Scan(40t). (f) Insert(40t).

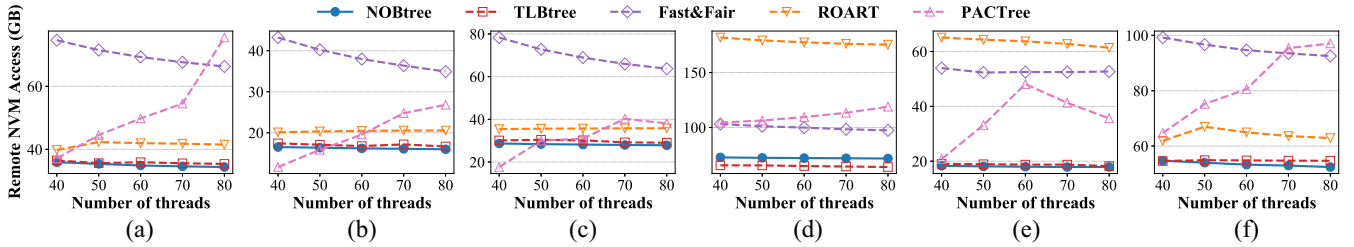


Fig. 9. Remote NVM access under various workloads. (a) Write-heavy. (b) Read-heavy. (c) Read-only. (d) Scan. (e) Read-Modify-Write. (f) Write-only.

769 slots to absorb insertions. We have tuned the threshold to
770 reduce the rebuilding time, which improves the performance.

771 Despite PACTree’s efforts to address the NUMA effect, its
772 performance struggles to scale effectively, especially beyond
773 40 threads. In particular, PACTree’s performance on the read-
774 only workload experiences a sharp decline, ultimately ranking
775 as the least efficient index. This decline could be attributed
776 to the influence of hyper-threading and cache coherence
777 protocols. We measure the data volume of the remote access in
778 Section V-C and find that PACTree experiences an increasing
779 volume of remote NVM access when the thread number
780 increases from 40 to 80. This is the main reason for PACTree’s
781 degraded performance. PACTree recommends using snooping
782 protocols for NVM, but most platforms use a directory-based
783 protocol by default because snooping may not scale well
784 to more CPU cores. The server in our experiments adopts
785 the directory-based protocol and does not allow changes to
786 coherence protocols, potentially contributing to the suboptimal
787 performance of PACTree.

788 B. Tail Latency

789 In this section, we assess the tail latency of the index’s
790 lookup, scan, and insert operations under 20 and 40 threads.
791 As shown in Fig. 8, NOBtree consistently exhibits the lowest-
792 tail latency across most cases and achieves $1.46\times$, $2.59\times$, and
793 $8.25\times$ lower latency than other indexes on average under the
794 three workloads, respectively. Notably, for the scan operation,
795 Fast&Fair outperforms others under 20 threads, attributed
796 to its larger leaf node size compared to NOBtree, which
797 is particularly advantageous for scan operations. However,
798 under 40 threads, the impact of the NUMA effect leads to
799 increased tail latency for all indexes across all workloads. In
800 this scenario, NOBtree maintains the lowest-tail latency across
801 all three workloads, providing further evidence of its superior
802 performance.

C. Remote NVM Access

803 In this section, we measure the data traffic caused by remote
804 NVM accesses. We utilized the Intel Performance Counter
805 Monitor (Intel PCM)¹ to quantify remote NVM access during
806 runtime. Fig. 9 presents the results for six workloads, where
807 we varied the thread number from 40 to 80 and distributed
808 threads across two sockets. NOBtree consistently achieves
809 less remote NVM access amount than TLBtree except the
810 *scan* workload. This can be owing to the replication of
811 the upper layer and the node migration. Under the *scan*
812 workload, NOBtree’s smaller leaf nodes incur more remote
813 NVM access compared to TLBtree. Since multiple nodes
814 need to be scanned to obtain results, smaller node sizes lead
815 to scanning more nodes. Nevertheless, as shown in Fig. 7,
816 NOBtree still outperforms TLBtree in throughput under the
817 *scan* workload because TLBtree places data on a single socket,
818 causing threads on the other socket to consistently access
819 remote NVM. In addition, Fast&Fair experiences a significant
820 amount of remote NVM access, which is primarily due to its
821 large node size.
822

D. Throughput Under Uniform Workloads

823 In this section, we will evaluate all indexes under uni-
824 form workloads. We omit the *write-heavy* and *write-only*
825 workloads as the insert operation is always randomly dis-
826 tributed. As shown in Fig. 10, the performance of all the
827 indexes degrades compared with those under skewed work-
828 loads. However, NOBtree’s performance still outperforms
829 others, indicating the efficiency of NOBtree under uniform
830 workloads, which is mainly owing to the per-NUMA upper-
831 layer replication as it can effectively reduce costly remote
832 NVM access.
833

¹<https://github.com/intel/pcm>

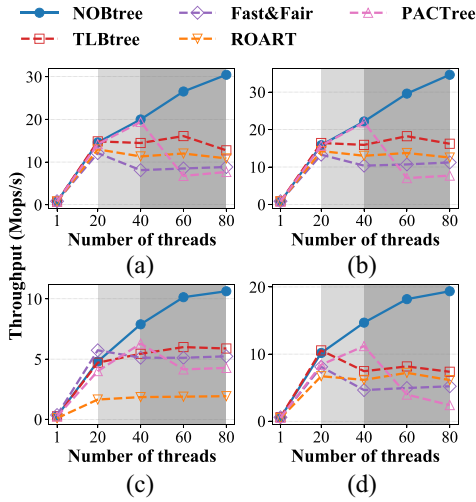


Fig. 10. Throughput under uniform workloads. (a) Read-heavy. (b) Read-only. (c) Scan. (d) Read-Modify-Write.

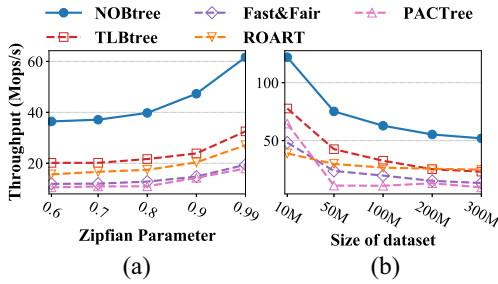


Fig. 11. Impact of workload skewness and data volume. (a) Impact of Skewness. (b) Impact of Data Volume.

E. Sensitivity Analysis

Impact of Workload Skewness: The skewness of the workloads reflects the deviation of query keys toward hotkeys, a characteristic often described by the Zipf parameter in the Zipfian distribution. In this experiment, we first load 100M keys and then perform 200M lookup operations. We vary the Zipf parameter from 0.6 to 0.99. Fig. 11(a) shows the throughput under 80 threads. One can see that all the indexes' performance improves with the increase of skewness. The performance of all indexes improves with increasing skewness. This improvement can be attributed mainly to the effect of the CPU cache. Higher skewness implies that most queries concentrate on a small set of keys, enabling the index to utilize the CPU cache more efficiently. NOBtree achieves the highest performance across all skewness, indicating its stable performance for different access patterns.

Impact of Data Volume: In this experiment, we explore the impact of dataset volume on the index's performance. We load several keys to construct the index and then evaluate the lookup performance under 80 threads. The lookup operations follow the Zipfian distribution with a parameter of 0.99. The amount of keys varies from 10 M to 300 M. Fig. 11(b) illustrates the results. The performance of all indexes degrades as the data volume increases. This decline occurs because more keys will heighten the indexes, leading to longer traversal

TABLE I
COMPARISON OF NVM CONSUMPTION (GB)

NOBtree	TLBtree	Fast&Fair	PACTree	ROART
6.29	5.78	5.12	6.32	25.03

paths during operations. NOBtree consistently outperforms other indexes.

F. Space Cost

The NVM consumption of each index is shown in Table I. In this experiment, we measure the NVM usage of each index after loading 200 million records with 8-byte keys and values. Fast&Fair demonstrates the lowest-NVM consumption, mainly attributed to its larger node size. NOBtree consumes more NVM spaces than TLBtree, primarily due to two reasons. First, we reduce the number of keys that a single node can store to allocate space for maintaining migration statistics, which leads to more nodes after the key loading. Second, NOBtree replicates the upper layer across NUMA nodes, introducing additional NVM consumption. The size of the NOBtree's upper layer is only about 34 MB, which is relatively small compared to the total index size. Thus, it is possible to put the upper layer in a faster DRAM to further accelerate the performance of NOBtree.

VI. FURTHER DISCUSSIONS

As of July 2022, Intel decided to discontinue its Optane product [3]. It is a major setback for NVM-related research. However, the NVM technology has been regarded as an efficient solution to partially address the storage wall issue, and NVM-oriented structures and algorithms are needed in the storage architecture involving NVM. So far, except for the 3D-XPoint technology used by the Optane series, there are other technologies for implementing NVM, such as phase change memory (PCM), spin-transfer torque magnetoresistive RAM (STT-MRAM), and resistive RAM (RRAM) [20]. A recent report showed that the total sales of STT-MRAM might reach 98.3 billion dollars by 2033 [21]. Thus, there is still an urgent demand to study efficient structures optimized for NVM devices. There are similarities between Optane and other NVM devices, such as nonvolatility, read/write asymmetry, and byte addressability. Therefore, current research related to Optane can still inspire future work.

In addition, the growing demand for larger memory capacity has led to the development of new memory technologies like compute eXpress link (CXL). CXL provides a cache-coherent interface for connecting CPUs, memory, and accelerators. CXL-attached memory exhibits characteristics similar to NVM, such as byte addressability and near-DRAM performance with higher capacities [3], [22]. CXL-attached memory can be viewed as remote NUMA memory, and most current research on CXL memory is based on emulating memory in a remote NUMA node [22]. The proposed approach in this study can also be applied to indexes for CXL-attached memory in the future, given the similarity between remote NUMA memory and NVM.

VII. CONCLUSION

In this article, we presented NOBtree, a new index structure designed to mitigate the NUMA effect in NVM indexes. NOBtree employs a decoupled tree structure, which consists of a read-optimized upper layer and a write-optimized bottom layer to enhance both read and write performance. To improve the read performance in the NUMA architecture, we proposed per-NUMA replication for the upper layer and a node migration mechanism for the bottom layer. Additionally, we devised a dedicated NUMA-aware NVM allocator to optimize the insertion performance of NOBtree. The experimental results across diverse workloads suggested the effectiveness and efficiency of NOBtree.

REFERENCES

- [1] Q. Wang, Y. Lu, J. Li, and J. Shu, "Nap: A black-box approach to NUMA-aware persistent memory indexes," in *Proc. OSDI*, 2021, pp. 93–111.
- [2] I. Calciu, S. Sen, M. Balakrishnan, and M. K. Aguilera, "Black-box concurrent data structures for NUMA architectures," in *Proc. ASPLOS*, 2017, pp. 207–221.
- [3] D. Koutsoukos, R. Bhartia, M. Friedman, A. Klimovic, and G. Alonso, "NVM: Is it not very meaningful for databases?" *Proc. VLDB Endow.*, vol. 16, no. 10, pp. 2444–2457, 2023.
- [4] (Intel Corp., Santa Clara, CA, USA). *Intel Optane Technology*. 2020. [Online]. Available: <https://www.intel.com/content/www/us/en/architecture-and-technology/intel-optane-technology.html>
- [5] D. Hwang, W. Kim, Y. Won, and B. Nam, "Endurable transient inconsistency in byte-addressable persistent B+-tree," in *Proc. FAST*, 2018, pp. 187–200.
- [6] Y. Luo, P. Jin, Z. Zhang, J. Zhang, B. Cheng, and Q. Zhang, "Two birds with one stone: Boosting both search and write performance for tree indices on persistent memory," *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5s, pp. 1–25, 2021.
- [7] B. Zhang, S. Zheng, Z. Qi, and L. Huang, "NBTree: A lock-free PM-friendly persistent B+-tree for eADR-enabled PM systems," *Proc. VLDB Endow.*, vol. 15, no. 6, pp. 1187–1200, 2022.
- [8] W. Kim, M. K. Ramanathan, X. Fu, S. Kashyap, and C. Min, "PACTree: A high performance persistent range index using PAC guidelines," in *Proc. SOSP*, 2021, pp. 424–439.
- [9] B. Daase, L. J. Bollmeier, L. Benson, and T. Rabl, "Maximizing persistent memory bandwidth utilization for OLAP workloads," in *Proc. SIGMOD*, 2021, pp. 339–351.
- [10] Y. Luo, P. Jin, Q. Zhang, and B. Cheng, "TLBtree: A read/write-optimized tree index for non-volatile memory," in *Proc. ICDE*, 2021, pp. 1889–1894.
- [11] S. Ma et al., "ROART: Range-query optimized persistent ART," in *Proc. FAST*, 2021, pp. 1–16.
- [12] I. Oukid, D. Booss, A. Lespinasse, W. Lehner, T. Willhalm, and G. Gomes, "Memory management techniques for large-scale persistent-main-memory systems," *Proc. VLDB Endow.*, vol. 10, no. 11, pp. 1166–1177, 2017.
- [13] A. Demeri, W. Kim, M. K. Ramanathan, J. Kim, M. Ismail, and C. Min, "Poseidon: Safe, fast and scalable persistent memory allocator," in *Proc. 21st Int. Middlew. Conf.*, 2020, pp. 207–220.
- [14] Z. Dang et al., "NValloc: Rethinking heap metadata management in persistent memory allocators," in *Proc. ASPLOS*, 2022, pp. 115–127.
- [15] W. Cai, H. Wen, H. A. Beadle, C. Kjellqvist, M. Hedayati, and M. L. Scott, "Understanding and optimizing persistent memory allocation," in *Proc. ISMM*, 2020, pp. 60–73.
- [16] Y. Mao, E. Kohler, and R. T. Morris, "Cache craftiness for fast multicore key-value storage," in *Proc. 7th ACM Eur. Conf. Comput. Syst.*, 2012, pp. 183–196.
- [17] J. Ge et al., "SALI: A scalable adaptive learned index framework based on probability models," *Proc. ACM Manag. Data*, vol. 1, no. 4, pp. 1–25, 2023.
- [18] R. Binna, E. Zangerle, M. Pichl, G. Specht, and V. Leis, "HOT: A height optimized Trie index for main-memory database systems," in *Proc. SIGMOD*, 2018, pp. 521–534.
- [19] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears, "Benchmarking cloud serving systems with YCSB," in *Proc. 1st ACM Symp. Cloud Comput.*, 2010, pp. 143–154.
- [20] A. Chen, "A review of emerging non-volatile memory (NVM) technologies and applications," *Solid-State Electron.*, vol. 125, pp. 25–38, Nov. 2016.
- [21] C. Tom and H. Jim, *Emerging Memories Branch Out*, Santa Clara, CA, USA, Storage Netw. Ind. Assoc., 2023.
- [22] Y. Sun et al., "Demystifying CXL memory with genuine CXL-ready systems and devices," in *Proc. MICRO*, 2023, pp. 105–121.