

PIM-tree: A Skew-resistant Index for Processing-in-Memory

Hongbo Kang Tsinghua University

Yiwei Zhao Carnegie Mellon University khb20@mails.tsinghua.edu.cn yiweiz3@andrew.cmu.edu

Guy E. Blelloch Carnegie Mellon University guyb@cs.cmu.edu

Laxman Dhulipala University of Maryland laxman@umd.edu

Yan Gu UC Riverside ygu@cs.ucr.edu Charles McGuffey Reed College cmcguffey@reed.edu

Phillip B. Gibbons Carnegie Mellon University gibbons@cs.cmu.edu

ABSTRACT

The performance of today's in-memory indexes is bottlenecked by the memory latency/bandwidth wall. Processing-in-memory (PIM) is an emerging approach that potentially mitigates this bottleneck, by enabling low-latency memory access whose aggregate memory bandwidth scales with the number of PIM nodes. There is an inherent tension, however, between minimizing inter-node communication and achieving load balance in PIM systems, in the presence of workload skew. This paper presents PIM-tree, an ordered index for PIM systems that achieves both low communication and high load balance, regardless of the degree of skew in data and queries. Our skew-resistant index is based on a novel division of labor between the host CPU and PIM nodes, which leverages the strengths of each. We introduce push-pull search, which dynamically decides whether to push queries to a PIM-tree node or pull the node's keys back to the CPU based on workload skew. Combined with other PIM-friendly optimizations (shadow subtrees and chunked skip lists), our PIM-tree provides high-throughput, (guaranteed) low communication, and (guaranteed) high load balance, for batches of point queries, updates, and range scans. We implement PIM-tree, in addition to prior proposed PIM indexes, on the latest PIM system from UPMEM, with 32 CPU cores and 2048 PIM nodes. On workloads with 500 million keys and batches of 1 million queries, the throughput using PIM-trees is up to 69.7× and 59.1× higher than the two best prior PIM-based methods. As far as we know these are the first implementations of an ordered index on a real PIM system.

PVLDB Reference Format:

Hongbo Kang, Yiwei Zhao, Guy E. Blelloch, Laxman Dhulipala, Yan Gu, Charles McGuffey, and Phillip B. Gibbons. PIM-tree: A Skew-resistant Index for Processing-in-Memory. PVLDB, 16(4): 946 - 958, 2022. doi:10 14778/3574245 3574275

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at https://github.com/cmuparlay/PIM-tree.

INTRODUCTION

The mismatch between CPU speed and memory speed (a.k.a. the "memory wall") makes memory accesses the dominant cost in today's data-intensive applications. Traditional architectures use

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit https://creativecommons.org/licenses/by-nc-nd/4.0/ to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment

Proceedings of the VLDB Endowment, Vol. 16, No. 4 ISSN 2150-8097. doi:10.14778/3574245.3574275

multi-level caches to reduce data movement between CPUs and memory, but if the application exhibits limited locality, most data accesses are still serviced by the memory. This excessive data movement incurs significant energy cost, and performance is bottlenecked by high memory latency and/or limited memory bandwidth.

Processing-in-Memory (PIM) [25, 30], a.k.a. near-data-processing, is emerging as a key technique for reducing costly data movement. By integrating computing resources in memory modules, PIM enables data-intensive computation to be executed in the PIMenabled memory modules, rather than moving all data to the CPU to process. Recent studies have shown that, for programs with high data-intensity and low cache-locality, PIM provides significant advantages in increasing performance and reducing power consumption by reducing data movement [14, 15]. Although proposals for processing-in-memory/processing-in-storage date back to at least 1970 [29], including forays by the database community in active disks [26], PIM is emerging today as a key technology thanks to advances in 3D-stacked memories [18, 22] and the recent availability of commercial PIM system prototypes [30]. Typical applications exploiting state-of-the-art PIM architectures include neural networks [3, 21, 23, 32], graph processing [1, 17, 35], databases [6, 7], sparse matrix multiplication [13, 33], and genome analysis [2, 34].

PIM systems are typically organized as a host (multicore) CPU that pushes compute tasks to a set of P PIM modules (computeenhanced memory modules), and collects the results. Thus, cost is incurred for moving both task descriptors and data-the sum of these costs is the communication cost between the CPU and the PIM modules. The host CPU can be any commodity multicore processor, and is typically more powerful than the wimpy CPUs within the PIM modules. Thus, an interesting feature of a PIM system is the potential to utilize both sets of resources (CPU side and PIM side).

In this paper, we focus on designing a PIM-friendly ordered index for in-memory data. Ordered indexes (e.g., B-trees [11]) are one of the backbone components of databases/data stores, supporting efficient search queries, range scans, insertions, and deletions. Prior works targeting PIM [10, 24] proposed ordered indexes based on range partitioning: the key space is partitioned into P subranges of equal numbers of keys, and each of the P PIM modules stores one subrange. Each PIM module maintains a local index over the keys in its subrange, and the host CPU maintains the top portion of the index down to the P roots of the local indexes. This approach works well for data and queries with uniformly random keys-the setting studied by these works-but it suffers from load imbalance under data or query skew. In more realistic workloads, batches of queries/updates may concentrate on the data in a small subset of the partitions, overwhelming those PIM modules, while the rest are idle. In the extreme, only one PIM module is active processing queries and the rest are idle, fully serializing an entire batch of queries on a single (wimpy) processor. The approach also suffers the cost of all data movements required to keep partitions (roughly) balanced in size. In a recent paper [19], we designed a PIM-friendly skip list that asymptotically achieves load-balance (details in Section 2.4), but the solution is not practical (as discussed below).

To address the above challenges with query and data skew, we present the *PIM-tree*, a practical ordered index for PIM that achieves both low communication cost and high load balance, regardless of the degree of skew in data and queries. Our skew-resistant index is based on a novel division of labor between the host CPU and PIM nodes, which leverages the strengths of each. Moreover, it combines aspects of both a B+-tree and a skip list to achieve its goals. We focus on achieving *high-throughput*, processing *batches* of queries at a time in a bulk-synchronous fashion. The PIM-tree supports a wide range of batch-parallel operations, including point queries (Get, Predecessor), updates (Insert, Delete, and range Scan.

We introduce *push-pull search*, which dynamically decides, based on workload skew, whether (i) to *push* queries from the CPU to a PIM-tree node residing on a PIM module or (ii) to *pull* the tree-node's keys back to the CPU. Combined with other PIM-friendly optimizations—*shadow subtrees* and *chunked skip lists*—our PIM-tree provides high-throughput, (guaranteed) low communication costs, and (guaranteed) high load balance, for batches of point queries, updates, and range scans. For example, each point query and update is answered using only $O(\log_B \log_B P)$ expected communication cost, where B is the expected fanout of a PIM-tree node and P is the number of PIM modules, independent of the number of keys n in the data structure, or the data skew. Note that it would take over 10^{19} PIM modules for $\log_B \log_B P$ to exceed 1, under our selection of B = 16; hence, the communication cost is constant in practice.

We implement the PIM-tree on the latest PIM system from UP-MEM [30], with 32 CPU cores and 2048 PIM modules. We choose four state-of-the-art ordered indexes as competitors, including two PIM-friendly approaches [19, 24] implemented by ourselves, and two traditional approaches[8, 9] implemented in SetBench [4]. On workloads with 500 million keys and batches of 1 million queries, the PIM-tree achieves (i) up to 59.1× higher throughput than the range-partitioned solution [24], (ii) up to 69.7× higher throughput than the prior skew-resistant solution [19], and (iii) comparable throughput in all cases regardless of skew and as low as 0.3× less communication than two state-of-the-art non-PIM indexes [8, 9].

The main contributions of the paper are:

- We design the PIM-tree, a high-throughput skew-resistant PIM
 data structure that efficiently supports a wide range of batchparallel point queries, updates, and scans, even under highlyskewed workloads. It causes nearly constant data movement
 (communication cost) for a point query or update, and linear
 data movement for scans. Key ideas include push-pull search and
 shadow subtrees.
- We implement and evaluate the PIM-tree on a commercial PIM system prototype, demonstrating significant performance improvements at modest skew, and performance gains that increase linearly with larger skew. As far as we know these are the first implementations of an ordered index on a real PIM system.

2 BACKGROUND

2.1 PIM System Architecture and Model

The Processing-in-Memory Model. We use the *Processing-in-Memory Model (PIM Model)* (first described in [19]) as an abstraction of generic PIM systems. It is comprised of a host CPU front end (CPU side) and a collection of P PIM modules (PIM side). The CPU side is a standard multicore processor, with an on-chip cache of M words. Each PIM module is comprised of a DRAM memory bank (local PIM memory) with an on-bank processor (PIM processor) and a local memory of $\Theta(n/P)$ words (where n denotes the problem size). The PIM processor is simple but general-purpose (e.g., a single in-order core capable of running C code). The CPU host can send code to the PIM modules, launch the code, and detect when the code completes. It can also send data to and receive data from PIM memory. The model assumes there is no direct PIM-to-PIM communication, although we could take advantage of such communication on PIM systems supporting it.

As the PIM model combines a shared-memory side (CPU and its cache) and a distributed side (PIM modules), algorithms are analyzed using both shared-memory metrics (work, depth) and distributed metrics (local work, communication time). On the CPU side, the model accounts for *CPU work* (total work summed over all cores) and *CPU depth* (all work on the critical path). On the PIM side, the model accounts for *PIM time*, which is the maximum local work on any one PIM core, and *IO time*, which is the maximum number of messages to/from any one PIM module. Programs execute in bulk-synchronous rounds [31], and the overall complexity metrics of an algorithm is the sum of the complexity metrics of each round. We focus on IO time and IO rounds in this paper.

Programming Interface. For concreteness, we assume the following programming interface for our generic PIM system, although our techniques would also work with other interfaces. Programs consist of two parts: a *host program* executed on the host CPU, and a *PIM program* executed on PIM modules. The host program has additional functions (discussed below) to communicate with the PIM side, including functions to invoke PIM programs on PIM modules and to transfer data to/from PIM modules. The PIM program is a traditional program that is invoked in all PIM processors when launched by the host program. It executes using the module's local memory, with no visibility into the CPU side or other PIM modules. The specific functions are (named MPI-style [16]):

- PIM_Load(PIM_Program_Binary): loads a binary file to the PIM modules.
- PIM_Launch(): launches the loaded PIM program on all PIMs.
- PIM_Status(): checks whether the PIM program has finished on all PIMs.
- PIM_Broadcast(src, length, PIM_Local_Address): copies a fixed length buffer to the same local memory address on each PIM.
- PIM_Scatter(srcs[], length[], PIM_Local_Address): similar to PIM_BROADCAST, but with a distinct buffer for each PIM module.
- PIM_Gather(dsts[], length[], PIM_Local_Address): the reverse
 of PIM SCATTER, reading into the buffer array dsts[].

¹There is no separate accounting needed for messages to/from the CPU side because any well-balanced system should provide bandwidth out of the host CPU that matches bandwidth into the PIM modules (and vice-versa).

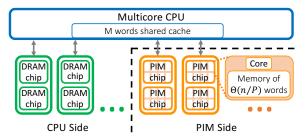


Figure 1: The architecture for the UPMEM PIM system, a specific example of our generic PIM system architecture. PIM modules are packed into memory DIMMs connected to the host CPU via normal memory channels. The CPU side also includes traditional DRAM modules, which are not part of the PIM model.

Algorithm 1. Batch-parallel Execution(O: batch of operations)

Repeat the following steps until done processing *O*:

- (1) Prepare a buffer of tasks for each PIM module.
- (2) Scatter the task buffers to the local memory of each PIM module using either PIM_SCATTER or PIM_BROADCAST.
- (3) Launch PIM programs using PIM_LAUNCH, to run their tasks and fill their reply buffers. Wait until all tasks finish.
- (4) Gather reply buffers using PIM_GATHER.

Based on this interface, our PIM-friendly ordered index processes batches of operations in *bulk-synchronous rounds*, like in [28], using the steps in Algorithm 1. As discussed in Section 5, when implementing our PIM-friendly programs, we use pipelining to overlap step 1 at the CPU and step 3 at the PIM modules.

A Concrete Example: UPMEM. We evaluate our techniques on the latest PIM system from UPMEM [30]. UPMEM's architecture (Figure 1) is one way to instantiate the PIM model. Its PIM modules are plug-and-play DRAM DIMM replacements, and therefore can be configured with various ratios of traditional DRAM memory to PIM-equipped ones (current maximum available configuration has 2560 PIM modules). The CPU has access to both the main memory (traditional DRAM) and all the PIM memory, but each PIM processor only has access to its local memory. Each PIM module has up to 628 MB/s local DRAM bandwidth, so a machine with 2560 PIM modules can provide up to 1.6 TB/s aggregate bandwidth [15]. To move data between PIM modules, the CPU reads from the origin and writes to the target. UPMEM's SDK supports the programming interface functions listed above, but with the restriction that the scatter/gather functions must transmit same length buffers to/from all PIM modules (i.e., the buffers are padded out to equal lengths).

UPMEM's main memory (a component not in the PIM model) enables running programs with CPU-side memory footprints over M words, but these additional memory accesses bring another type of communication not existing in the PIM model: CPU-DRAM communication. Thus the cache efficiency is important for host programs. Our solution in the PIM-tree is to use only a small amount of CPU-side memory: $\Theta(S) < M$ words for a batch of S operations.

2.2 Load Balance Preliminaries

A key challenge for PIM systems is to keep load balance among the PIM modules, which we define as follows:

Definition 2.1. A program achieves **load balance** if the work (unit-time instructions) performed by each PIM program is O(W/P) and the *communication* (data sent/received) by each PIM module is O(C/P), where W and C are the sums of the work and communication, respectively, across all P PIM modules. For programs with multiple bulk-synchronous rounds, the program achieves load balance if each round achieves load balance.

The challenge in achieving load balance is that the PIM module with the maximum work or communication must be bounded. Note that randomization does not directly lead to load balance, e.g., randomly scattering P tasks of equal work and communication to P PIM modules fails to achieve load balance. This is because one of the PIM modules receives $\Theta(\log P/\log\log P)$ tasks with high probability $(whp)^2$ in P [5], causing the work and communication at that module to be a factor of $\Theta(\log P/\log\log P)$ higher than balanced.

We use balls-into-bins lemmas to prove load balance, where a bin is a PIM module and a ball with weight w corresponds to a task with w work or w communication. We will use the following:

LEMMA 2.2 ([19, 27]). Placing weighted balls with total weight $W = \sum w_i$ and each $w_i < W/(P \log P)$ into P bins uniformly randomly yields O(W/P) weight in each bin whp.

2.3 Prior Work on Indexes for PIM

There are several prior works for indexes on PIM systems. Two prior works [10, 24] proposed PIM-friendly skip lists. Their skip lists are based on *range partitioning*: they partition the skip list by disjoint key ranges and maintain each part locally on one PIM module. As discussed in Section 1, such range partitioning can suffer from severe load imbalance under data and query skew.

The load imbalance problem of range-partitioned ordered indexes is also studied in traditional distributed settings. Ziegler et al. [36] discussed other choices for tree-based ordered indexes in order to avoid load imbalance, including: (i) partitioning by the hash value of keys, (ii) fine-grained partitioning that randomly distributes all index nodes, and (iii) a hybrid method that does fine-grained partitioning in leaves, and range partitioning for internal nodes. They also experimentally evaluated their approaches on an 8 machine cluster. However, each of these choices has its own problem in the case of a PIM system with thousands of PIM modules: (i) partitioning by hash makes range operations costly, because they must be processed by all PIM modules, (ii) fine-grained partitioning causes too much communication because all accesses will be non-local, and (iii) the hybrid method suffers from the load balance problem in its range partitioned part.

2.4 Prior Work: PIM-balanced Skip List

In a recent paper [19], we presented the first provably load-balanced batch-parallel skip list index, the *PIM-balanced skip list*, under adversary-controlled workloads on the PIM model. A key insight was to leverage the CPU side to solve the load balance problem.

The **PIM-balanced skip list** horizontally splits the skip list into two parts, an *upper part* and a *lower part*, replicating the upper part in all PIM modules and distributing lower part nodes randomly to PIMs. This is shown in Figure 2, where nodes in different PIM modules have different colors, and the replicated upper part is explicitly

²We use O(f(n)) with high probability (whp) (in n) to mean O(cf(n)) with probability at least $1-n^{-c}$ for $c\geq 1$.

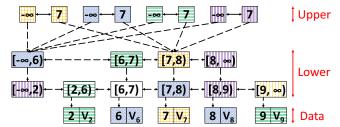


Figure 2: PIM-balanced skip list [19] with the upper part replicated on a 4-PIM system. Nodes on different PIM modules are different colors. PIM pointers are dashed lines. The lower part is $\log P$ levels.

drawn as multiple copies. For a system with P PIM-modules, the lower part is $\log P$ levels. We can afford to replicate (only) the top part because (i) it is small relative to the rest of the skip list and (ii) it is updated relatively infrequently (recall that an inserted key reaches a height h in a skip list with probability $1/2^h$).

Queries are executed by pointer chasing in the "tree" of skip list nodes. The batched queries are first evenly divided and sent to all PIM modules, each progressing through the upper part locally. Then the skip list goes through the lower part by sending the query to the host PIM module of each lower part node on the search path one-by-one, until reaching a leaf. We call this the **Push** method, because queries are sent ("pushed") to PIM modules to execute.

Executing a batch of parallel queries using only *Push* can cause severe imbalance, despite the lower part nodes being randomly distributed. For skewed workloads, many queries may share a common node on their search path, meaning that they are all sent to the host PIM module of that node, causing a load imbalance. These nodes are called *contention points*. An example is when multiple non-duplicate Predecessor queries return the same key, with *all* nodes on the search path to that key being contention points.

The **PIM-balanced skip list** [19] solves this problem by avoiding contention points, based on a key observation: once the search paths of keys l and r share a lower part node v, searching any key $u \in [l,r]$ will also reach node v. Thus the search for u can start directly from the LCA (lowest common ancestor) of these two paths. We call this the Jump-Push method. Jump-Push search has a preprocessing stage to record search paths. It is a multi-round sample search starting with one sample: In each round, it doubles the sample size and uses the search paths recorded in previous rounds to decide start nodes of sample queries in this round. This approach limits the contention on each node, avoiding load imbalance.

However, the preprocessing cost is high. For P PIM modules and a batch of $P\log^2 P$ operations, it takes $O(\log P)$ sampling rounds, each of which takes $O(\log P)$ steps of inter-module pointer chasing to search the lower part. The main stage, in contrast, takes only $O(\log P)$ steps. Moreover, the preprocessing stage requires recording entire search paths—another overhead for the CPU side.

Our new ordered index (PIM-tree) uses some of the same ideas as this work, but includes key new ideas to make it simpler, and more efficient both theoretically and practically.

3 PIM-TREE DESIGN

Overview. The PIM-tree is a batch-parallel skew-resistant ordered index designed for PIM systems. It supports fundamental key-value

operations, including Get(key), Update(key, value), Predecessor(key), Insert(key, value), Delete(key), and Scan(Lkey, Rkey). It executes operations in same-type atomic batches in parallel, similar to [28]. As such, the data structure avoids conflicts caused by operations of different types. We design it starting from the structure discussed in §2.4.

In this section, we describe PIM-tree's design by studying the impact of our techniques/optimizations on the Predecessor operation. We review the basic data structure design discussed in §2.4 in detail, then introduce our three key techniques/optimizations, and finally analyze the resulting communication cost and load balance. Later in §4, we describe the design of PIM-tree's other operations.

Notations. We refer to the number of PIM modules as P, the total number of elements in the index as n, the batch size as S, and the expected fanout of PIM-tree nodes as B.

Key Ideas. We observe that the two components of the PIM architecture, the CPU side and the PIM side, prefer different workloads. The distributed PIM side prefers uniformly random workloads and suffers from the load imbalance caused by highly skewed ones. Meanwhile, the shared CPU side prefers skewed workloads since we can explore spatial and temporal locality in these workloads that lead to better cache efficiency. This difference shows the complementary nature of shared-memory and distributed computing, and their coexistence in the PIM architecture motivates our hybrid method: we design a dynamic labor division strategy that automatically switches between the CPU side and the PIM side to use the more ideal platform. This strategy, called *Push-Pull search*, is the core technique of the PIM-tree.

With load balance achieved by Push-Pull search, we further propose two other optimizations, called *shadow subtrees* and *chunking*, to reduce communication. These optimizations are motivated by two basic ideas respectively: caching remote accesses at PIM modules to build local shortcuts (thereby eliminating communication), and blocking nodes into chunks (for better locality). We will show how these "traditional" techniques combine with the Push-Pull search optimization to bring an asymptotic reduction of communication from $O(\log P)$ to $O(\log_B \log_B P)$ (B is the expected fanout of chunked nodes), and a throughput increase of up to $69.7\times$, compared with the PIM-friendly skip list of §2.4.

3.1 Basic Structure

The PIM-balanced skip list is a distributed skip list horizontally divided into three parts, the upper part, the lower part, and the data nodes (Figure 2). Data nodes are key-value pairs randomly distributed to PIM modules to support hash-based lookup in one round and O(1) communication. Every upper part node is replicated across all PIM modules, and every lower part node is stored in a single random PIM module. The ID of the PIM module hosting a lower part node is called the node's PIM ID. Remote pointers, called PIM pointers, are comprised of (PIM ID, address) pairs. In Figure 2, PIM pointers are represented by dashed arrows, while traditional (i.e., intra-PIM) pointers are represented by solid arrows. To save communication during search, each lower part node stores the key of the next skip list node at the same level (called a right-key). There is a node at the lowest level for each key, and the probability of a node joining the next higher level in the skip list is set to 1/2.

The upper part is replicated to enable local executions for queries on PIM modules, but the replication brings an overhead of P to both space complexity and update costs. To mitigate this overhead, the lower part height is set to be $H_{\rm low} = \log P$, so that only a $1/2^{\log P} = 1/P$ fraction of the keys reach the upper part. By replicating the upper part, the number of remote accesses needed for a PREDECESSOR query is reduced from $O(\log n)$ to $O(\log P)$.

In the following sections, we call the upper part *L3* and the lower part *L2*. After applying the Shadow Subtree optimization (§3.3), we will further divide the lower part horizontally into two parts, called *L2* and *L1*.

3.2 Push-Pull Search

Push-Pull search is our proposed search method that guarantees load balance even under skewed workloads. In the **Push** method, the CPU sends the query to the host PIM module of the next node along the search path, the PIM module runs the query, then the CPU fetches the result; in the **Pull** method, the CPU retrieves the next node along the path back to the CPU side, running the query itself. **Push-Pull** search chooses between Push and Pull by counting the number of queries to each node: when the number of queries to a node exceeds a specific threshold, denoted as *K* in the following sections, we Pull that node, otherwise we Push the query.

In further detail, Push-Pull search performs multi-round pointer chasing over the basic structure mentioned in §3.1 in three stages, where the CPU records the next pointer for each query as an array of PIM pointers throughout the process.

- Traverse L3 using the replicated upper parts. The CPU evenly distributes queries to PIM modules. Each PIM module runs its queries using its local copy of L3, until reaching a pointer to an L2 node. The CPU retrieves these pointers (using PIM_GATHER).
- (2) Traverse L2 using contention-aware Push-Pull. The CPU performs multiple Push-Pull rounds. In each round, the CPU counts the number of queries to each L2 node. If there are more than *K* queries to a node, choose Pull by sending a task to the PIM-side to retrieve that node to the CPU-side, then partition the queries (in parallel) based on the PIM IDs in the retrieved node's pointers on the CPU-side. Otherwise, choose Push to send a Query task to the PIM and retrieve the next pointer for the query.
- (3) When the search reaches a data node, return the data.

We can record the addresses of all nodes on the pointer-chasing path for a query on the CPU side to get the *search trace* for each query. Note that these traces are used when performing updates (in §4). For the basic structure mentioned in §3.1, we choose K = 1, as it minimizes communication for constant size nodes.

Discussion. The most interesting part of Push-Pull search is that it is based on integrating two fundamental methods from distributed and shared-memory computing to achieve provable load balance with low cost (see §3.5 for analysis). We observe that the Push method is a distributed computing technique, as it uses the CPU as a router and always runs queries on PIM modules. Meanwhile, the Pull method is a shared memory technique, treating the PIM modules as standard memory modules and running the queries on the CPU. As discussed in §3, combining such fundamental methods works because of the complementary nature of the CPU side and the PIM side in the load balance issue: contention-causing (thus PIM-unfriendly) workloads are meanwhile CPU-friendly workloads.

As a solution only to the load balance issue, Push-Pull search provides no asymptotic improvements in worst-case bounds compared with Push-Only or Pull-Only methods. Such improvements are provided by our optimizations, *shadow subtrees* and *chunking*, which we describe next.

3.3 Shadow Subtrees

Shadow subtrees are auxiliary data structures in L2 that act as shortcuts to reduce communication from $O(\log P)$ to $O(\log \log P)$ for each query, while ensuring that the space complexity is still O(n). The shadow subtree optimization is based on the idea of the search tree defined by a skip list, which is an imaginary tree generated by merging all possible search paths of a skip list. It contains all nodes and all edges of the skip list, except some horizontal edges. The shadow subtree of each node is a shadow copy of its search subtree stored together with this node. By using shadow subtrees, a PIM module can run queries locally through L2. Although shadow subtrees and replicating the top of the tree both involve copying nodes across different PIM modules with the purpose of reducing communication, they are actually quite different. When replicating the top, a single tree is copied *P* times across the modules. In the shadow subtree, every ancestor of a node has a copy of that node as part of its shadow subtree (in our case just the ancestors in L2).

Building shadow subtrees on all (O(n)) L2 nodes would require $O(n \log P)$ space. Instead, to maintain O(n) space, we build them only on a small proportion of L2 nodes. In particular, we divide L2 into two layers, denoting the upper levels to be the new L2 and the lower levels to be L1. We build shadow subtrees only on the new L2. We set the height of L1 to be $H_{L1} = \log \log P$, so only $(1/\log P)$ -fraction of nodes $(O(n/\log P))$ nodes) are in the new L2, and the space complexity summing over all shadow subtrees is O(n). Thus, the PIM-tree now has three layers: L3 under full replication, L2 with shadow subtrees, and L1 under random distribution without any replication. Each layer requires O(n) space, so the total space complexity is O(n). This is shown in Figure 3. We refer to original tree nodes and pointers to them as *physical* nodes (pointers), and mark them in black. Shadow-tree nodes and pointers to them are referred to as *shadow* nodes and pointers, and are marked in red.

Accelerating PREDECESSOR using Shadow Subtrees. Shadow subtrees strengthen the Push side of Push-Pull search: a single Push round can send a query through the whole L2, rather than going forward by just one level, by running the query on the shadow subtree. Therefore the search process takes only $O(\log\log P)$ rounds for uniform-random workloads: one Push through L3, one Push through L2, and $O(H_{\rm L1}) = O(\log\log P)$ Push-Pull rounds for L1. However, for skewed workloads, we cannot simply perform a single Push round through L2, because multiple queries may still be pushed to the contention points in L2 and cause load imbalance. We solve this problem again by Pull, by introducing a multi-round Pull process to eliminate contention points.

In further detail, Push-Pull search in L2 has two stages: we first perform up to $O(H_{L2})$ Pull rounds for nodes with $\geq K$ queries until no such node exists, where $H_{L2} = \log P - \log \log P$ denotes the new L2's height, then execute one "Push" round to send all queries through L2. We set the threshold $K = H_{L2}$ instead of 1 since "Push" is now more powerful and we tend to use it more. Both stages take

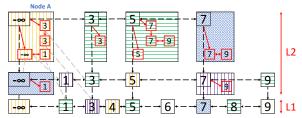


Figure 3: The structure of L2 and L1 after introducing shadow subtrees. Shadow nodes and shadow pointers are marked in red. Note that blue 1 does not have a shadow tree node for 3 because node 3 is not in its search subtree. Right-keys are omitted. We also omit pointers from shadow nodes to physical nodes except for node A. The L1 part (L2 part) is $\log \log P$ levels ($\log P - \log \log P$ levels, respectively).

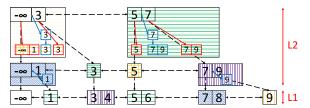


Figure 4: The intermediate state of the chunking transformation. We merge non-pivot nodes (nodes whose keys do not go to upper levels) to their left-side neighbors. Redundant shadow subtrees after merging are marked in blue, and will be removed in Figure 5. All physical pointers from shadow nodes are omitted.

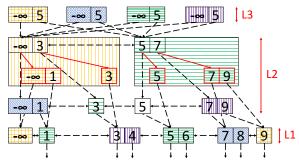


Figure 5: The actual structure of the PIM-tree with redundant shadow subtrees removed from Figure 4. We no longer need right-keys after chunking. Data nodes are omitted.

O(1) balanced communication per query. This is partly proved in Lemma 3.2, and fully proved in the full paper [20].

In practice, we use another optimization to reduce the number of Pull rounds. Note that although contention points are the only source of load imbalance, we may reach a reasonable level of load balance before eliminating all contention points. Therefore, to avoid unnecessary Pull rounds, before starting a Pull round, we measure the load balance across PIMs by counting the number of queries that will be sent to each PIM module; if the one with the most is below $3\times$ the average load, we stop the Pull round and start to Push.

Replication and Space/Imbalance Tradeoffs. Compared with (i) full replication used for L3 and (ii) range partitioning (which performs no replication) used in related works, shadow subtree is a novel scheme that supports queries with O(1) communication by improving the locality of a distributed ordered index. Specifically, shadow subtree is a selective replication approach that lies between these two prior schemes. If we replicate nodes not only to their L2

Table 1: Comparison between three types of replication schemes that run queries with O(1) communication. The larger the overhead factor, the more space it takes and the slower updates will be. The larger the maximum query number is, the more imbalanced the execution will be under skewed workloads.

Scheme	Overhead factor	Maximum query number
Full Replication	P	Perfect Balance
Range Partitioned	1	P
Shadow Subtrees	$O(\log P)$	$\log P$

ancestors, but to all PIM modules, we obtain full replication. On the other hand, if we only keep the shadow subtrees of the L2 roots, we obtain the range partitioned scheme.

The cost and the skew-resistance of shadow subtrees also lies between those of the other two schemes. Table 1 shows the bounds when applying different schemes to a skip list with height $\log P$ and size P in expectation. In the full replication scheme, we can run queries with perfect load balance, but it brings an overhead factor *P* to both space complexity and update costs. On the other hand, for the range partitioned scheme, each query can only be executed by a single PIM module. We can still do Push-Pull to avoid contention with a threshold of K = P: we will choose to simply Pull the whole tree when the number of queries exceeds the size of the whole part. There can be load imbalance, as some part gets up to Pqueries and others get none. For this approach there is no overhead on space complexity or for updates. Lastly, using shadow subtrees, the overhead factor is $O(\log P)$ as each node is replicated in all its L2 search tree ancestors, and the maximum query number is log P according to our choice of Push-Pull threshold $K = H_{1,2}$.

Shadow subtrees therefore yield a balanced compromise between the two schemes, providing a sweet spot for both overhead and skew-resistance.

3.4 Chunked Skip List

Chunking or "blocking" is a classic idea widely used in locality-aware data structures, e.g., B-trees and B+-trees. To improve locality, we apply a similar chunking approach to improve the access granularity of the PIM computation, while decreasing the tree height. As chunking increases the access granulariy, each PIM processor obtains larger local memory bandwidth, therefore better performance. The effect of access granularity in PIM is discussed in detail in [15].

We apply chunking to all layers of the PIM-tree. In L3, we replace the multi-thread skip list with a batch-parallel multi-threaded B+tree [28]. In L2 and L1, we chunk the nodes in our skip list to obtain a *chunked skip list*. We first merge horizontal non-pivot nodes (whose keys do not go to upper levels) into a single chunk, then remove redundant shadow subtrees. Applying this two step process on Figure 3 first gives Figure 4 as an intermediate state, and finally the PIM-tree in Figure 5.

The result with shadow subtrees looks similar to a B+-tree. The difference is that while the B+-tree sends nodes to upper levels on overflow of lower level nodes, the chunked skip list uses random heights generated during Insert, so the fanout holds in expectation. We decrease the probability of reaching the next level in the skip list from 1/2 to 1/B, so that the expected fanout is B. We choose the same chunking factor B in L3, L2 and L1 for simplicity, but different factors could be used in each part. As discussed in §4.2,

We use a chunked skip list instead of a classical B+-tree in L2 to make batch-parallel distributed INSERT and DELETE simpler and more efficient. We use a B+-tree in L3 because the structure is not distributed, making batch-parallel INSERT and DELETE easier.

Chunking reduces tree height at all levels, which improves multiple aspects of our design. We denote the new L2 (L1) height as H'_{L2} (H'_{L1} , respectively). The L2 part of the search path to each node is reduced from $O(H_{L2}) = O(\log P)$ to $H'_{L2} = \log_B P - \log_B \log_B P$, as there are no longer horizontal pointer-chasing processes. Therefore, the space and replication overhead of shadow subtree reduces from $O(H_{L2})$ to H'_{L2} , as each node is only replicated in its L2 ancestors. Furthermore, lower overhead enables us to reduce H'_{L1} to $\log_B \log_B P$. H'_{L1} is effectively 1 in practice, because with our choice of B=16 it will take over 10^{19} PIM modules for H'_{L1} to exceed 1.

Therefore, in practice, the height of L2 is reduced to 2 levels, and L1 reduced to 1 level. The probability for a key to reach L3 is 1/4096 < 1/P, and the probability of reaching L2 is $1/16 < \log_{16} P$. **Implementing Predecessor after Chunking.** Chunking brings only one modification to the search process: changing the Push-Pull threshold K from $H_{\rm L2}$ to $B \cdot H'_{\rm L2}$, because we now "Pull" chunks with expected size O(B) instead of O(1). The detailed algorithm is explained in §3.5.

Chunking also improves the communication costs of Predecessor. First, as the height of L1 is reduced from $\log \log P$ to $\log_B \log_B P$, each query now causes only $O(\log_B \log_B P)$ communication in L1. Second, chunking reduces the maximum possible number of Pullonly rounds in L2 from $O(H_{\rm L2}) = O(\log P)$ to exactly $H'_{\rm L2}$, which is $\log_B P - \log_B \log_B P$. This helps reduce the number of communication rounds under skewed workloads.

3.5 Predecessor Algorithm and Bounds

Next, we describe the complete algorithm for PREDECESSOR, and discuss its cost complexity. We provide proofs for the communication cost and load balance of PREDECESSOR queries. For simplicity throughout the paper, our cost analyses assume that hash functions provide uniform random maps to PIM modules, so that the lemma in §2.2 can be applied. Algorithm 2 summarizes the search process.

Algorithm 2. Predecessor (Q: batch of query keys)

- (1) Push queries from Q evenly to PIM modules, and traverse L3
- (2) While the number of queries that will be sent to each PIM module for L2 is not balanced (i.e., the busiest PIM module gets more than 3× the average load), do the following:
 - (a) Pull all nodes with more than K = B ⋅ H'_{L2} queries back to the CPU.
 - (b) Use these nodes to progress the pointer-chasing process of these queries by one step.
- (3) Push each query to the PIM module holding its search node, and traverse L2 using the shadow subtrees.
- (4) Perform H'_{L1} Push-Pull rounds with K = B to traverse L1, and retrieve the data nodes.

We demonstrate here a mini step-by-step example of a PREDECES-SOR batch with four queries on the PIM-tree in Figure 5 (note that real batches should have more queries on this tree to achieve loadbalance). The queries request the Predecessors of keys 1, 3, 4 and 7. PIM-tree first evenly distributes one query for each of the four PIM modules to search through L3, returning three queries falling onto the L2 node $[-\infty, 3]$ and one falling onto node [5, 7]. The context of node $[-\infty, 3]$ will be pulled to the CPU from the yellow-masked PIM module due to its large contention, and the pointer-chasing searching of keys 1, 3 and 4 over L2 will be executed on the CPU side. After that, query 1, query (3, 4), and query 7 will be pushed to the PIM module containing the blue-masked node $[-\infty, 1]$, greenmasked node [3] and green-masked node [5, 7] respectively on the local shadow subtrees to search through L2. Finally, all queries will be carried out in a similar Push-Pull way to return the results from L1 and data nodes.

Theorem 3.1. A batch of Predecessor queries can be executed in $O(\log_B P)$ communication rounds, with a cost of $O(\log_B \log_B P)$ communication for each operation in total whp. The execution is load balanced if the batch size $S = \Omega(P \log P \cdot B \cdot H'_{L2}) = \Omega(P \log P \cdot B \cdot \log_B P)$. The CPU-side memory footprint is O(S).

We provide part of the proof here, and give the full proof with all details in the full paper [20]. The key challenge is to prove the communication bounds and load balance, and we do this by proving separately for each stage of Algorithm 2. We take Lemma 3.2, the proof for the L2 Push stage (stage 3) as an example.

Lemma 3.2 (Push round for L2). Push using the shadow subtrees (stage 3) takes 1 round, O(1) communication whp for each query, and is load balanced.

PROOF. In this stage we send each query as a task to the corresponding PIM module, incurring O(1) communication per query and 1 round overall. For load balance, we analyze it as a weighted balls-into-bins game, where we take the target nodes as balls, the numbers of queries on the target nodes as weights, and PIMs as bins. The weight limit is $K = B \cdot H'_{L2}$ by assumption, as each node gets at most K queries, and the weight sum is at most S. Applying Lemma 2.2, each PIM module incurs O(S/P) communication. \square

4 PIM-TREE: OTHER OPERATIONS

Having described the design of the PIM-tree data structure in §3, using the Predecessor operation as the running example, we now briefly introduce how other PIM-tree operations are implemented. Please refer to the full paper [20] for more detail.

4.1 GET and UPDATE using Hashing

Get and Update are operations with a given key. These operations also do not modify the structure of the data structure. Therefore, we solve these in one round and O(1) communication per operation through a hash-based approach by first (i) using a fixed hash function to map keys to PIM modules, and (ii) building a local hash table on each PIM module to map keys to the local memory addresses of their data nodes.

Because the data nodes are distributed by a hash function, we achieve good load balance even for skewed workloads, assuming that there are not duplicate operations to the same key. If such redundant operations exist, we can solve this by preprocessing to

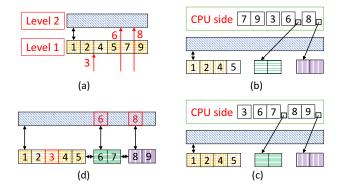


Figure 6: The process to insert keys 3, 6, and 8 into L2 of the PIM-tree. Insertion 3 has height 1, and insertion 6 and 8 both have height 2. These heights are generated beforehand by coin tossing with probability 1/B. The height of the yellow node is 1, and that of the blue node is 2. As a result, key 3 is inserted into the yellow node, and keys 6 and 8 split the yellow node.

combine operations on the CPU-side, using a user-defined combinining mechanism. In practice, we use a linear-probing hash table on the PIM-side, but other hash-table variants could also be used.

4.2 INSERT

An Insert(key,value) operation inserts the key into nodes on its search path, and the primary challenge is to avoid contention and conflicts when multiple Inserts in a batch go to the same node(s). We solve this by preprocessing: we perform searches in parallel and record the trace of each search, and use the traces to detect and handle contention points. Our algorithm has three stages: (1) perform searches to record the search trace, (2) modify the physical skip list based on the search trace, and (3) update the shadow subtrees.

Update Physical Skip List. After getting the search traces of each INSERT, we INSERT into these nodes according to random heights we generate prior to updating the PIM-tree. Figure 6 is an example of our contention solving strategy. According to their pre-generated heights, the insertion of 3 is into the yellow node, and the insertions of 6 and 8 will split that node. The insertion takes three steps: in (b) we fetch the right-side part of the node to the CPU side and generate empty new nodes in random PIM modules; in (c) we derive the correct element to be inserted to each node in CPU; finally in (d) we insert them.

Choosing skip lists instead of B+-trees as the basis of L1 and L2 helps reduce the number of rounds, since we can insert to all nodes in parallel, rather than level-by-level bottom up from the leaves like the B+-tree. Insert to L3 is also executed in parallel with that of L2, performed by broadcasting INSERTS reaching L3 to all PIMs.

Update Shadow Subtrees. To maintain the invariant that shadow subtrees are copies of the search subtrees, we update the shadow subtrees after updating the physical skip list. There are three types of updates: (1) **build** the new shadow subtree for a new node, (2) **insert** a new node into the shadow subtrees of its ancestors, and (3) **trim** a shadow subtree after a node split.

Our shadow subtree updating technique is straightforward. For build, we pull the L2 search tree and send it to the new node. For insert and trim, we observe that only shadow subtrees of nodes on the search trace need updating, so we send the newly-inserted node to all these nodes.

Discussion: Load Balance in Insert. There is a load balance issue in our shadow subtree update algorithm: To keep shadow subtrees up to date, an L2 node may need updates of size $O(P/\log_B P)$. For example, a new L2 root needs to build its shadow subtree of expected $O(P/\log_B P)$ nodes (given $H'_{L2} = \log_B P - \log_B \log_B P$). This contention factor $O(P/\log_B P)$ will grow faster than the factor $K = B \cdot \log_B P$ of Predecessor as P grows. This contention has minor effect at present, but we propose an algorithm to address this problem (not implemented at present).

The solution is not to keep all shadow subtrees up to date, but instead to mark some nodes as *unfinished*, update their shadow subtrees gradually in future rounds, and avoid using a shadow subtree until it is up to date. This helps smooth out the imbalance. Predecessor bounds still hold when the number of such nodes is below a threshold, and we achieve this by additional update rounds, which are load balanced when the number of such nodes is high.

Theorem 4.1. A batch of Insert operations can be executed in $O(\log_B P)$ IO rounds, incurring $O(\log_B \log_B P)$ communication for each operation. The execution is load balanced if the batch size $S = \Omega(P\log P \cdot B \cdot H'_{L2}) = \Omega(P\log P \cdot B \cdot \log_B P)$. The CPU-side memory footprint is O(S).

Implementing Delete. We handle deletions similarly to insertions: first obtain the search trace, then delete keys from nodes on the trace, finally apply updates to shadow subtrees. While insertion causes node split, deletion causes nodes to merge when removing the pivot key from a node. See the full paper [20] for details.

4.3 SCAN

The Scan(LKey,Rkey) operation (a.k.a. range query) returns all the (key, value) pairs whose keys fall into the range of [Lkey, Rkey]. Its algorithm is similar to Predecessor.

When running a batch of Scan queries, we first on the CPU merge all batched overlapping ranges into groups of disjoint ranges. Then PIM-tree labels two boundary nodes (the predecessors of Lkey and Rkey on the current search level of the tree) for each range as SEARCHREQUIRED, and all the intermediate nodes as FETCHALL.

FETCHALL nodes are required to return all their leaf data nodes. Note that all ranges are disjoint, so no contention exists in FETCHALL nodes. All FETCHALL nodes are pushed to PIMs and recursively return all child nodes. Meanwhile, SEARCHREQUIRED nodes of a range are maintained similar to PREDECESSOR on each level, using push-pull search to deal with potential contention. It might generate new FETCHALL nodes. Please see the full paper [20] for details.

5 IMPLEMENTATION

CPU-PIM Pipelining. Thus far, we have introduced algorithms where tasks on the CPU and PIM run in a synchronized, tick-tock manner in each round as depicted in Algorithm 1. The total execution time of this approach consists of three non-overlapping components: CPU-only time, PIM-only time, and communication time. Communication requires both CPU and PIM, but the other two components only utilize one part of the system, which presents an opportunity to reduce execution time by pipelining the CPU-only and PIM-only components.

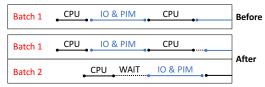


Figure 7: Program traces before/after CPU-PIM pipelining

For pipelining, we consider executions that run multiple batches in parallel in the PIM-tree. This is shown in Figure 7, where "CPU" represents time spent in CPU-only execution, while "IO & PIM" represents time spent in CPU-PIM communication and the PIM program. On our UPMEM system, CPU-PIM communication requires exclusive control of the PIM side, and any concurrent use of the PIM side will cause a hardware fault. Hence, one batch needs to wait for the PIM side to finish the current execution tasks. We only pipeline queries in our experiments, since update batches cannot be carried out concurrently. For mixed operations, we protect the PIM-tree by a read-write lock to prevent update batches from running concurrently with other batches.

PIM Program. PIM-tree's PIM program is a parallel executor of the tasks in the buffer sent from the CPU. It is designed to address two features of UPMEM's current PIM processors. First, the PIM processor is a fine-grained multi-threaded computing unit [15], and requires at least 11 threads to fill the pipeline, so we write PIM programs in the form of 12 threads. Second, UPMEM's system only supports PIM programs with fewer than 4K instructions, but the implementation of PIM-tree exceeds this bound. To bypass this restriction, we write the PIM program as multiple separate modules, and load each module when needed. Only INSERT and DELETE operations require swapping modules; program loading currently takes around 25% of the execution time. The remaining operations on PIM-tree fit within the 4K instruction limit.

6 EVALUATION

In this section, we evaluate our new PIM-optimized indexes on a PIM-equipped machine provided by UPMEM, and two traditional state-of-the-art indexes on a machine with similar performance. We summarize our experimental results from this section as follows:

- (1) The PIM-tree performs better than the range-partitioned skip list under skewed workloads in terms of throughput, memory-bus communication, and energy consumption.
- (2) The PIM-tree causes lower communication on the memory bus compared with traditional indexes without PIM.
- (3) All optimizations mentioned, including Push-Pull search, shadow subtrees, chunked skip list and CPU-PIM pipelining, yield performance increases to (some) PIM-tree operations.

6.1 Experiment Setup

UPMEM's PIM Platform. We evaluate PIM-tree on a PIM-equipped server provided by UPMEM(R). The server has two Intel(R) Xeon(R) Silver 4126 CPUs, each CPU with 16 cores at 2.10 GHz and 22 MB cache. Each socket has 6 memory channels: 4 DIMMs of conventional DRAM are implemented on 2 channels, while 8 UPMEM DIMMs are on the other 4 channels. Each of the 16 UPMEM DIMMs has 2 ranks, each rank has 8 chips, and each chip has 8 PIM modules. There are 2048 PIM modules in total.

Traditional Machine w/o PIM. We evaluate traditional indexes on a machine with two Intel(R) Xeon(R) CPU E5-2630 v4 CPUs, each CPU with 10 cores at 2.20 GHz and 25 MB cache. Each socket has 4 memory channels. There are no PIM-equipped DIMMs. We cannot evaluate traditional indexes on the server of UPMEM because 2/3 of its memory channels are used by PIM-equipped DIMMs, which cannot be used as the main memory. Directly running traditional indexes on the server would cause unfairness in main memory bandwidth for the traditional indexes. In our experiments we choose the state-of-the-art binary search tree [8] and (a,b)-tree [9] as competitors. Both implementations are obtained from SetBench [4].

Range-Partitioned and Jump-Push Baselines. We implement a range-partitioned-based ordered PIM index as our primary baseline, where both data nodes and index nodes are distributed to PIM modules based on the ranges of the key [10, 24]. We record the range splits in the CPU side, and use these splits to find the targeted PIM module of each operation. Point operations are sent to and executed on the corresponding PIM module. Running a batch of SCAN operations is similar, except that it runs an additional splitting in queries according to the range splits before tasks are sent to the PIMs. We also build a local hash table on all PIM modules for GET.

We also implement the PIM-balanced skip list [19] described in §2.4 as another baseline. We experimentally evaluate this approach when discussing the impact of the optimizations proposed in this paper, in Figure 11, with the algorithm called "Jump-Push based".

Test Framework. We run multiple types of operations on the PIM-tree, range-partitioning skip list we implemented, and the state-of-the-art traditional indexes. In all experiments, we first warm up the index by running the **initialize set** that INSERT key-value pairs, then evaluate the index by the **evaluation set** of multiple operations. All operations are loaded from pre-generated test files. PIM algorithms (the PIM-tree and range-partitioning) run operations in batches, and traditional indexes run them directly with multi-threaded parallelism. In all experiments, the sizes of both keys and values are set to 8 bytes.

To study the algorithms, we measure both the time spent, and the memory bus traffic. Memory bus traffic is measured by adding CPU-PIM and CPU-DRAM communication, the prior one measured by a counter increased whenever a PIM function(e.g., PIM_Broadcast) is called, and the later one measured as cache misses by PAPI. We bind the program to a single NUMA node and disable the CPU-PIM pipeline when measuring cache misses for an accurate traffic measurement. As each CPU of the PIM-equipped machine has two NUMA nodes, the effective cache of the PIM algorithms is reduced to 11 MB, half of the full cache, under this setup. Time is measured with full interleave over all NUMA nodes.

Instability in performance exists in the current generation of PIM hardwares. We oberserved an approximately $\pm 15\%$ fluctuation in the measured metrics mentioned above in our experiments.

6.2 Microbenchmarks

Workload Setup. Each test first warms up the index by inserting 500 million uniform random key-value pairs;³ then for testing it executes (i) 100 million point operations or (ii) 1 million SCAN

³This is a favorable setting for the range-partitioned baseline, because the range boundaries are stable. The performance of the PIM-tree is not impacted by the distribution of key-value pairs over time.

operations that each retrieve 100 elements in expectation. Point operations use batch size S=1 million, and Scan operations use batch size S=10 thousand.

We generate skewed workloads with Zipfian distribution [37]. However, workloads generated by Zipf-skew over elements is not ideal for evaluating batch-parallel ordered indexes, because this skew can be easily handled by a deduplication in preprocessing on the CPU side, by merging operations of the same key into one. To better represent the spatial bias, where keys in some ranges are more likely to be accessed in the same batch, we slightly modify the way to generate our Zipfian workload, as follows: (i) we divide the key space evenly into P = 2048 parts; (ii) for each operation, we first choose a part according to the Zipfian distribution, then choose a uniformly random element in that part. For operations for existing keys (Get, Delete), we divide and choose among the keys currently in the index; for operations on arbitrary keys (Predecessor, Insert, Scan), the key space consists of all 64-bit integers. We periodically shuffle the probability of each part in Zipfian distribution. This helps alleviate, but not eliminate the PIM memory overflow problem of the range-partitioned baseline caused by INSERTS accumulating in high-probability parts. PIM-tree gains no benefit from this shuffle.

To show results on different amounts of skew, we evaluate the algorithms on different α values in the Zipfian distribution, ranging from 0 (uniformly random) to 1.2. With this skewness generation approach, less than 10% operations are eliminated because of duplicated keys, under the most skewed case ($\alpha = 1.2$).

Performance. Figure 8 illustrates the throughput of the range-partitioned skip list and the PIM-tree on microbenchmarks. The performance of the range-partitioned baseline drops drastically as the query skew increases, while the PIM-tree shows robust resistance to query skew. In fact, across all operations, it is observed that PIM-tree is essentially unaffected by data skew, obtaining similar running times for $\alpha=0$ and $\alpha=1.2$. For $\alpha=1.2$, PIM-tree outperforms the range-partitioned baseline by $3.87-59.1\times$.

It is observed that GeT operations are significantly simpler and achieve higher throughput since a hash table is used as a shortcut (the same holds for UPDATE operations), whereas PREDECESSOR and SCAN operations must go through the entire ordered index. In Figure 8(c), INSERT on the range-partitioned baseline crashes when $\alpha = 1.2$, because skewed INSERT causes imbalanced data placement over PIM modules, then causes overflow of local memory on some PIM modules. Although this problem could be solved by a rebalancing scheme, the rebalancing process itself will cause load imbalance as it requires sending data from the overflowing PIM modules to other less-loaded PIM modules. It is observed that even if this improvement to the baseline (with which the existing rangepartitioning solutions in the literature are not equipped) were to be made, the throughput of PIM-tree would still be significantly larger than range-partitioning, as this would be extra work that the baseline must perform during the execution.

Figure 9 shows the performance of PIM-tree compared with state-of-the-art binary search tree [8] and (a,b)-tree [9] under our workload. PIM-tree outperforms traditional indexes in all test cases, except throughput of PREDECESSOR compared with (a,b)-tree.

Execution breakdown. Figure 10 shows the percentage of time spent in each component mentioned in Section 5. These results

are derived with our pipelining optimization turned off, because pipelining would cause an overlapping of different components. We select as typical examples the throughput of Predecessor and Insert, on range-partitioned skip list and the PIM-tree, for uniform random workload and Zipfian-skewed workload with $\alpha=0.6$. Similar results also exist in the cases of other α values.

For range-partitioned skip list, *PIM Execution* and *CPU-PIM Communication* dominates the time cost of skewed workloads, mainly because the bottleneck of the entire execution—the busiest PIM modules— are receiving a growing number of tasks.

For uniform random workloads, *PIM Execution* only takes a small proportion of the total time cost, though almost all comparisons are executed in PIM modules. It is inferred that parallelism is fully exploited when a large number of PIM modules are involved in this case. We believe that this implies that PIM-based systems are an ideal platform for parallel index structures.

PIM-tree INSERT spends time loading PIM program modules during execution, as the full program size exceeds the current size of instruction memory on PIMs intended to store the PIM program. This limit is discussed in Section 5.

Effect of Optimizations. Figure 11 shows the impact of different optimizations on our ordered index. Here, we start with our Jump-Push baseline (the PIM-balance skip list in [19]). Replacing Jump-Push with Push-Pull provides up to 6.8× throughput improvements across all test cases. Adding Chunking provides the biggest improvement jump, up to 9.0×, across all test cases, while adding shadow subtrees mostly benefits Predecessor under no skew. (Insert get minor benefits because it needs to maintain this supplementary data structure.) Finally, adding pipelining—thereby implementing the complete PIM-tree algorithm—provides additional benefit for Predecessor. (Pipelining is not implemented for Insert because it would require interleaved Insert batches, which is not supported in our implementation.) Compared to the Jump-Push baseline, PIM-trees are up to 69.7× higher throughput for the settings studied.

Memory bus communication. Figure 12 shows the average amount of communication for PIM-tree, range-partitioned skip list, and traditional non-PIM indexes. PIM-tree needs less communication than all traditional indexes. Range partitioned skip lists outperform all competitors by much under uniform random workload, but perform much worse in skewed workloads.

Another observation is that, while the PIM-tree stores all the data and does most comparisons in PIM modules, most memory bus traffic is between CPU and the DRAM. This is because though PIM-tree algorithms requires O(S) CPU-side memory for a batch of S operations, the available setup with 11 MB cache is too small for batches of one million operations. As the result, CPU side data overflow to DRAM and cause significant CPU-DRAM communication. To show the effect of this overflow, in Figure 13 we study the CPU side communication as we run the 100 million uniform random predecessor operations with different batch sizes. Results show that the CPU-DRAM communication is reduced by 67% as we reduce batch size from 1M to 50K. We cannot directly use smaller batch size because of the load balance requirements, but this result hints that we can get much less CPU-DRAM communication when running the PIM-tree on a machine with larger cache size.

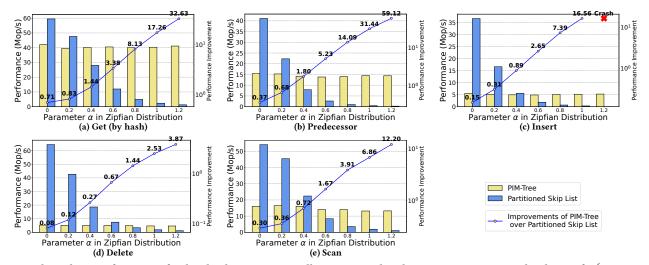


Figure 8: Throughput performance of ordered index operations. All operations other than SCAN are run using a batch size of 10^6 ; SCAN uses a batch size of 10^4 , with 100 elements retrieved by each SCAN operation in expectation.

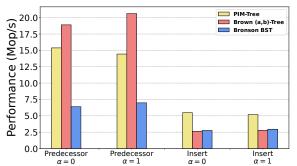


Figure 9: Throughput of PIM-tree versus SOTA traditional indexes.

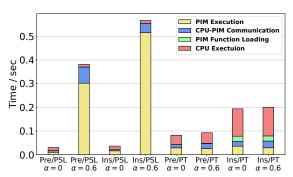


Figure 10: Time spent on each component of the program (without pipelining). "PSL" for partitioned skip list, "PT" for PIM-tree. "Pre" for predecessor, "Ins" for insert.

Push-Pull threshold choice We study the PIM-tree PREDECESSOR performance under different Push-Pull threshold. In our microbenchmark with $\alpha=1$, we find that choosing a lower threshold leads to about 10% throughput drop and up to 28% more CPU-PIM communication. A higher threshold brings minor performance increase. Please refer to the full paper [20] for more details.

Energy Evaluation. PIM-tree costs roughly $5 \times -10 \times$ less energy on skewed cases, compared to the range-partitioned baseline on PIM. Please refer to the full version of this paper [20] for details.

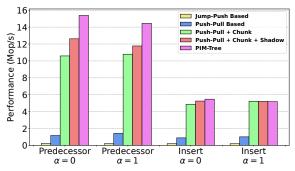


Figure 11: Impact of three optimizations on final performance.

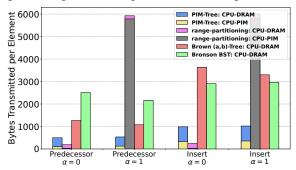


Figure 12: Average communication on the memory bus per operation in bytes.

YCSB workload. PIM-tree achieves roughly 9.5×-32× higher throughput on skewed cases, compared to the range-partitioned baseline on PIM. Please refer to the full version of this paper [20].

6.3 Workload of Real-world Skewness

In this section, we test the PIM-tree over a workload with real-world skewness using the publicly available wikipedia dataset [12], which is a collection of documents from wikipedia. To use this dataset in our test framework, we need to transform it into a collection of 8 byte key-value pairs, then run operations over them. To be specific, we first extract words from each document, lowercase them, then use (word, document id) pairs as keys, and a random

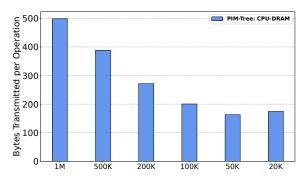


Figure 13: Average CPU-DRAM communication per predecessor operation for PIM-Tree with different batch size in bytes.



Figure 14: An example that convert word "ordered" in document id "1000000" to a 64-bit integer.

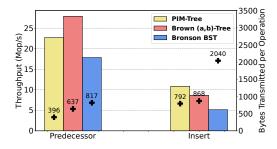


Figure 15: Throughput on the wikipedia workload.

8 byte integer as values. Because our indexes only support 8 byte integer keys, we need to transform the (word, document id) pairs into 8 byte integers. The transformation is shown in Figure 14. We use 40 bits to represent the word, and 23 bits to store the document id (as there's less than 2²³ documents). In the word part, we use 5 bits for each of the first 5 letters, then store the hash value of the whole word in the following 15 bits to avoid collision. After this transformation, the generated integers preserve the two skewness of english words: (i) word frequency skewness (some words are used more than others) and (ii) word distribution skewness in the dictionary order space (words with some prefix are used more).

In this test, we pick the first 1.2 billion keys: the first 1 billion words used for initialization, and the following 200 million used for evaluation. These keys covers the first 5.1 million documents, which is 63% of the whole dataset. There are 3.9 million unique words, and pairing the word and document id generates 529 million unique keys. We get duplicated keys only for the same word in the same document. Because the duplication rate of keys is about 2X, we also double the batch size of the PIM-tree to two million.

The result is shown in Figure 15, where the throughputs are shown as the bars, and communication (bytes transmitted peroperation) as labeled points. All indexes experience higher throughput and lower communication in this workload than in microbenchmarks because of the replicated keys. Comparing different indexes gives results similar to that of microbenchmarks: PIM-tree has

lower predecessor throughput than the (a,b)-tree, but outperforms traditional indexes in all other metrics.

7 DISCUSSION

PIM-tree outperforms conventional indexes in throughput in most cases, but very occasionally cannot win, e.g. only in PREDECESSOR compared with (a,b)-tree in our paper. We address here three hardware limits of the current PIM system by way of explanation, and to describe future changes to the hardware that would result in even better performance for PIM-optimized data structures.

The first factor is the limited CPU-PIM bandwidth on UPMEM's newly developed hardware. When carrying out a 50% read - 50% write task, the bandwidth obtained on UPMEM machine is 16GB/s, 1.9× slower than the shared-memory machine we use with a bandwidth of 31GB/s on the same workload. Even under such significant bandwidth limitations, PIM-tree still achieves better or comparable performance to DRAM-only indexes, primarily because it greatly reduces inter-module communication. Designing hardware to improve CPU-PIM bandwidth is thus an important direction, and one that we expect improvements for in the future. Therefore, we believe that PIM-tree will outperform conventional indexes in all cases in terms of throughput in the future.

Another issue is that the limited size of PIM programs prevents us from more complicated designs. Current workaround—dynamic program loading—is too costly. We believe this problem will be solved in future hardware by a larger instruction memory.

The last limit is that of inadequate CPU cache, as mentioned in Section 6.2. CPU-DRAM communication caused by cache overflow comprises most of memory bus communication, and this can be alleviated by a larger cache. We believe an adequate cache will be important in future PIM systems.

8 CONCLUSION

This paper presented PIM-tree, the first ordered index for PIM systems that achieves both low communication and high load balance in the presence of data and query skew. We presented the first experimental evaluation of ordered indexes on a real PIM system, demonstrating up to $69.7\times$ and $59.1\times$ higher throughput than the two best prior PIM-based methods and down to $0.3\times$ less communication than two state-of-the-art conventional indexes. Key ideas include push-pull search and shadow subtrees—techniques likely to be useful for other applications on PIM systems due to their effectiveness in reducing communication costs and managing skew. Our future work will explore such applications (e.g., radix-based indexes, graph analytics).

ACKNOWLEDGMENTS

We thank Rémy Cimadomo, Julien Legriel, Damien Lagneux, etc. at UPMEM for providing extensive access to their system and help whenever needed. This work would not have been possible without their support. This research was supported by NSF grants CCF-1910030, CCF-1919223, CCF-2028949, and CCF-2103483, VMware University Research Fund Award, Parallel Data Lab (PDL) Consortium (Alibaba, Amazon, Datrium, Facebook, Google, Hewlett-Packard Enterprise, Hitachi, IBM, Intel, Microsoft, NetApp, Oracle, Salesforce, Samsung, Seagate, and TwoSigma) and National Key Research & Development Program of China (2020YFC1522702).

REFERENCES

- [1] Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi. 2015. A scalable processing-in-memory accelerator for parallel graph processing. In 2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA). 105–117. https://doi.org/10.1145/2749469.2750386
- [2] Shaahin Angizi, Naima Ahmed Fahmi, Wei Zhang, and Deliang Fan. 2020. PIM-Assembler: A Processing-in-Memory Platform for Genome Assembly. In 2020 57th ACM/IEEE Design Automation Conference (DAC). 1–6. https://doi.org/10. 1109/DAC18072.2020.9218653
- [3] Shaahin Angizi, Zhezhi He, Adnan Siraj Rakin, and Deliang Fan. 2018. CMP-PIM: An Energy-Efficient Comparator-Based Processing-in-Memory Neural Network Accelerator. In Proceedings of the 55th Annual Design Automation Conference (San Francisco, California) (DAC '18). Association for Computing Machinery, New York, NY, USA, Article 105, 6 pages. https://doi.org/10.1145/3195970.3196009
- [4] Maya Arbel-Raviv, Trevor Brown, and Adam Morrison. 2018. Getting to the Root of Concurrent Binary Search Tree Performance. In 2018 USENIX Annual Technical Conference, USENIX ATC 2018, Boston, MA, USA, July 11-13, 2018, Haryadi S. Gunawi and Benjamin Reed (Eds.). USENIX Association, 295–306. https://www.usenix.org/conference/atc18/presentation/arbel-raviv
- [5] Petra Berenbrink, Tom Friedetzky, Zengjian Hu, and Russell Martin. 2008. On Weighted Balls-into-bins Games. Theoretical Computer Science 409, 3 (2008), 511–520.
- [6] Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Rachata Ausavarungnirun, Kevin Hsieh, Nastaran Hajinazar, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu. 2019. CoNDA: Efficient Cache Coherence Support for near-Data Accelerators. In Proceedings of the 46th International Symposium on Computer Architecture (Phoenix, Arizona) (ISCA '19). Association for Computing Machinery, New York, NY, USA, 629–642. https://doi.org/10. 1145/3307650.3322266
- [7] Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T Malladi, Hongzhong Zheng, and Onur Mutlu. 2016. LazyPIM: An efficient cache coherence mechanism for processing-in-memory. IEEE Computer Architecture Letters 16, 1 (2016), 46–50.
- [8] Nathan Grasso Bronson, Jared Casper, Hassan Chafi, and Kunle Olukotun. 2010. A practical concurrent binary search tree. In Proceedings of the 15th ACM SIG-PLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2010, Bangalore, India, January 9-14, 2010, R. Govindarajan, David A. Padua, and Mary W. Hall (Eds.). ACM, 257–268. https://doi.org/10.1145/1693453.1693488
- [9] Trevor Brown. 2017. Techniques for constructing efficient lock-free data structures. Ph.D. Dissertation. University of Toronto (Canada).
- [10] Jiwon Choe, Amy Huang, Tali Moreshet, Maurice Herlihy, and R. Iris Bahar. 2019. Concurrent Data Structures with Near-Data-Processing: an Architecture-Aware Implementation. In ACM Symposium on Parallelism in Algorithms and Architectures (SPAA). 297–308.
- [11] Douglas Comer. 1979. Ubiquitous B-tree. ACM Computing Surveys (CSUR) 11, 2 (1979), 121–137.
- [12] Wikimedia Foundation. 2016. Wikepedia:Database download. https://en. wikipedia.org/wiki/Wikipedia:Database_download. Accessed March 15, 2022.
- [13] Christina Giannoula, Ivan Fernandez, Juan Gómez Luna, Nectarios Koziris, Georgios Goumas, and Onur Mutlu. 2022. SparseP: Towards Efficient Sparse Matrix Vector Multiplication on Real Processing-In-Memory Architectures. Proc. ACM Meas. Anal. Comput. Syst. 6, 1, Article 21 (feb 2022), 49 pages. https://doi.org/10.1145/3508041
- [14] Christina Giannoula, Nandita Vijaykumar, Nikela Papadopoulou, Vasileios Karakostas, Ivan Fernandez, Juan Gómez-Luna, Lois Orosa, Nectarios Koziris, Georgios I. Goumas, and Onur Mutlu. 2021. SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures. In IEEE International Symposium on High-Performance Computer Architecture, HPCA 2021, Seoul, South Korea, February 27 March 3, 2021. IEEE, 263–276. https://doi.org/10.1109/HPCA51647.2021.00031
- [15] Juan Gómez-Luna, Izzat El Hajj, Ivan Fernandez, Christina Giannoula, Geraldo F Oliveira, and Onur Mutlu. 2021. Benchmarking a new paradigm: An experimental analysis of a real processing-in-memory architecture. arXiv preprint arXiv:2105.03814 (2021).
- [16] William Gropp, William D Gropp, Ewing Lusk, Anthony Skjellum, and Argonne Distinguished Fellow Emeritus Ewing Lusk. 1999. Using MPI: portable parallel programming with the message-passing interface. Vol. 1. MIT press.
- [17] Yu Huang, Long Zheng, Pengcheng Yao, Jieshan Zhao, Xiaofei Liao, Hai Jin, and Jingling Xue. 2020. A Heterogeneous PIM Hardware-Software Co-Design for Energy-Efficient Graph Processing. In 2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS). 684–695. https://doi.org/10.1109/IPDPS47924.2020.00076

- [18] Joe Jeddeloh and Brent Keeth. 2012. Hybrid memory cube new DRAM architecture increases density and performance. In 2012 Symposium on VLSI Technology (VLSIT). 87–88. https://doi.org/10.1109/VLSIT.2012.6242474
- [19] Hongbo Kang, Phillip B Gibbons, Guy E Blelloch, Laxman Dhulipala, Yan Gu, and Charles McGuffey. 2021. The Processing-in-Memory Model. In Proceedings of the 33rd ACM Symposium on Parallelism in Algorithms and Architectures. 295–306.
- [20] Hongbo Kang, Yiwei Zhao, Guy E. Blelloch, Laxman Dhulipala, Yan Gu, Charles McGuffey, and Phillip B. Gibbons. 2022. PIM-tree: A Skew-resistant Index for Processing-in-Memory. arXiv preprint arXiv:2211.10516 (2022).
- [21] Ji-Hoon Kim, Juhyoung Lee, Jinsu Lee, Jaehoon Heo, and Joo-Young Kim. 2021. Z-PIM: A Sparsity-Aware Processing-in-Memory Architecture With Fully Variable Weight Bit-Precision for Energy-Efficient Deep Neural Networks. *IEEE Journal of Solid-State Circuits* 56, 4 (2021), 1093–1104. https://doi.org/10.1109/JSSC.2020.3039206
- [22] Dong Uk Lee, Kyung Whan Kim, Kwan Weon Kim, Hongjung Kim, Ju Young Kim, Young Jun Park, Jae Hwan Kim, Dae Suk Kim, Heat Bit Park, Jin Wook Shin, Jang Hwan Cho, Ki Hun Kwon, Min Jeong Kim, Jaejin Lee, Kun Woo Park, Byongtae Chung, and Sungjoo Hong. 2014. 25.2 A 1.2V 8Gb 8-channel 128GB/s high-bandwidth memory (HBM) stacked DRAM with effective microbump I/O test methods using 29nm process and TSV. In 2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC). 432–433. https://doi.org/10.1109/ISSCC.2014.6757501
- [23] Wen Li, Ying Wang, Huawei Li, and Xiaowei Li. 2019. P³M: A PIM-Based Neural Network Model Protection Scheme for Deep Learning Accelerator (ASPDAC '19). Association for Computing Machinery, New York, NY, USA, 633–638. https://doi.org/10.1145/3287624.3287695
- [24] Zhiyu Liu, Irina Calciu, Maurice Herlihy, and Onur Mutlu. 2017. Concurrent Data Structures for Near-memory Computing. In ACM Symposium on Parallelism in Algorithms and Architectures (SPAA). 235–245.
- [25] Onur Mutlu, Saugata Ghose, Juan Gómez-Luna, and Rachata Ausavarungnirun. 2020. A Modern Primer on Processing in Memory. CoRR abs/2012.03112 (2020).
- [26] Erik Riedel, Garth A. Gibson, and Christos Faloutsos. 1998. Active Storage for Large-Scale Data Mining and Multimedia. In VLDB'98, Proceedings of 24rd International Conference on Very Large Data Bases, August 24-27, 1998, New York City, New York, USA, Ashish Gupta, Oded Shmueli, and Jennifer Widom (Eds.). Morgan Kaufmann, 62-73. http://www.vldb.org/conf/1998/p062.pdf
- [27] Peter Sanders. 1996. On the Competitive Analysis of Randomized Static Load Balancing. In Workshop on Randomized Parallel Algorithms (RANDOM).
- [28] Jason Sewall, Jatin Chhugani, Changkyu Kim, Nadathur Satish, and Pradeep Dubey. 2011. PALM: Parallel architecture-friendly latch-free modifications to B+ trees on many-core processors. Proceedings of the VLDB Endowment 4, 11 (2011), 795–806.
- [29] Harold S. Stone. 1970. A Logic-in-Memory Computer. IEEE Trans. Comput. C-19, 1 (1970), 73–78.
- [30] UPMEM. 2022. UPMEM Technology. https://www.upmem.com/technology/. Accessed March 15, 2022.
- [31] Leslie G Valiant. 1990. A bridging model for parallel computation. Commun. ACM 33, 8 (1990), 103–111.
- [32] Zhao Wang, Yijin Guan, Guangyu Sun, Dimin Niu, Yuhao Wang, Hongzhong Zheng, and Yinhe Han. 2020. GNN-PIM: A Processing-in-Memory Architecture for Graph Neural Networks. In Advanced Computer Architecture, Dezun Dong, Xiaoli Gong, Cunlu Li, Dongsheng Li, and Junjie Wu (Eds.). Springer Singapore, Singapore, 73–86.
- [33] Xinfeng Xie, Zheng Liang, Peng Gu, Abanti Basak, Lei Deng, Ling Liang, Xing Hu, and Yuan Xie. 2021. SpaceA: Sparse Matrix Vector Multiplication on Processing-in-Memory Accelerator. In 2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA). 570–583. https://doi.org/10. 1109/HPCA51647.2021.00055
- [34] Fan Zhang, Shaahin Angizi, Naima Ahmed Fahmi, Wei Zhang, and Deliang Fan. 2021. PIM-Quantifier: A Processing-in-Memory Platform for mRNA Quantification. In 2021 58th ACM/IEEE Design Automation Conference (DAC). 43–48. https://doi.org/10.1109/DAC18074.2021.9586144
- [35] Mingxing Zhang, Youwei Zhuo, Chao Wang, Mingyu Gao, Yongwei Wu, Kang Chen, Christos Kozyrakis, and Xuehai Qian. 2018. GraphP: Reducing Communication for PIM-Based Graph Processing with Efficient Data Partition. In 2018 IEEE International Symposium on High Performance Computer Architecture (HPCA). 544–557. https://doi.org/10.1109/HPCA.2018.00053
- [36] Tobias Ziegler, Sumukha Tumkur Vani, Carsten Binnig, Rodrigo Fonseca, and Tim Kraska. 2019. Designing Distributed Tree-based Index Structures for Fast RDMAcapable Networks. In ACM SIGMOD International Conference on Management of Data. 741–758.
- [37] George Kingsley Zipf. 2016. Human behavior and the principle of least effort: An introduction to human ecology. Ravenio Books.