

# DHASH: Non-Blocking Dynamic Hash Tables

Junchang Wang, Dunwei Liu, Xiong Fu, Fu Xiao, Chen Tian

**Abstract**—Once started, existing hash tables cannot change their pre-defined hash functions, even if the incoming data cannot be evenly distributed to the hash table buckets. In this paper, we present DHASH, a type of hash table for shared memory systems, that can change its hash function and rebuild the hash table on the fly, without noticeably degrading its service. The major technical novelty of DHASH stems from an efficient distributing mechanism that can atomically distribute every node when rebuilding, without locking the corresponding hash table buckets. This not only enables non-blocking lookup, insert, and delete operations, but more importantly, makes DHASH independent of the implementation of hash table buckets, such that DHASH allows programmers to select the set algorithms that meet their requirements best from a variety of existing lock-free and wait-free set algorithms. Evaluations show that DHASH can efficiently change its hash function on the fly. Moreover, when rebuilding, DHASH consistently outperforms the state-of-the-art hash tables in terms of throughput and response time of concurrent operations, at different concurrency levels, and with different operation mixes and average load factors.

**Index Terms**—Concurrent programming, parallelism and concurrency, dynamic hash tables.



## 1 INTRODUCTION

TYPICALLY, a hash table consists of an array of *buckets*. Nodes to be inserted into the hash table are first reduced to 32-bit or 64-bit *hash value* by using a pre-defined *hash function*. The hash table takes the hash value modulo the size of its bucket array to determine the buckets where the corresponding nodes will be stored. Once two or more distinct nodes are mapped to the same bucket, a *hash collision* has happened. One typical approach to dealing with hash collisions is *separate chaining*, in which each hash bucket contains a list of nodes mapped to the bucket. The average length of the lists, referred to as *average load factor*, is defined as the total number of nodes in the hash table divided by the size of the bucket array. In general cases, given a specified average load factor, hash tables offer the advantage of constant-time lookup operations, such that they have been widely used throughout computer systems.

Most existing hash tables, however, lacks robustness; the pre-defined hash functions cannot be changed even if they cannot evenly distribute incoming data to the hash table buckets. For example, an adversary can launch algorithmic complexity attacks [1] by first spying out the details of the hash function and then creating malicious nodes that will be hashed to a few buckets. This causes a few buckets (henceforth referred to as *target buckets*) to contain many more nodes than the average load factor, effectively making the hash table unavailable for legitimate requests accessing these buckets. This robustness issue has affected the hash table implementations in a long list of operating systems and programming languages, including the Linux kernel [2], Perl [1], PHP [3], and .NET [4].

Researchers have proposed various techniques to ad-

dress this robustness issue. The most well-studied solution is *universal hashing* [5], which, by randomly selecting a hash function from a set of well-designed hash functions, can, in theory, provide good average-case hashing performance and address the above robustness issue. Unfortunately, empirical evaluation shows that *universal hashing* is not fit for applications such as Perl [6], which is heavily used to process strings that are often selected from restricted character sets, and thus have an unexpected distribution property. Moreover, emerging research shows that it is possible for the adversary to observe or guess the random choice of the pre-defined hash function, rendering universal hashing useless [1], [7], [8], [9].

In this paper, we first demonstrate that rebuilding a hash table by dynamically changing its hash function on the fly is a promising approach to build robust hash tables. Specifically, once the pre-defined hash function cannot evenly distribute incoming data to the buckets, a new instance of the hash table with a distinct, well-designed hash function is created. Then, the existing nodes in the old hash table are distributed to the new one. During the time period of distributing nodes, both hash tables are used to serve concurrent *insert*, *delete*, and *lookup* operations (henceforth simply *common* operations). After all nodes have been distributed to the new hash table, the old hash table can be deleted. By changing the hash function, the new hash table can distribute existing nodes evenly over all buckets, and, therefore, provide the expected constant-time lookup operations. We use the term *dynamic* to describe a hash table algorithm that can provide this flexibility feature, and refer to the function that dynamically changes its hash function as *rebuild*. Researchers have proposed hash tables that can only enlarge or shrink their bucket sizes by constant factors but cannot change hash functions [10], [11], [12], [13], which we refer to as *resizable* hash tables. In Section 6.3, we demonstrate that for each run, resizable hash tables can only cut the list length of the target buckets by half, and hence cannot resolve the robustness issue. We thus mainly focus on dynamic hash tables. Note that in this paper, for ease

- J. Wang, D. Liu, X. Fu, and F. Xiao are with the School of Computer Science, Nanjing University of Posts and Telecommunications, Nanjing 210023, China. E-mail: {wangjc, 1020041113, fux, xiaoff}@njupt.edu.cn.
- C. Tian is with the State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210008, China. E-mail: tianchen@nju.edu.cn.

Manuscript received September 2020.  
(Corresponding author: Fu Xiao.)

of presentation, we only say that a dynamic hash table can change its hash function, even though it actually can change its hash function and/or hash seed. In contrast, a resizable hash table cannot change any of them, and the only aspect a resizable hash table can change is the size of its bucket array. Moreover, in some sense, a *dynamic* hash table is the hash table that is *resizable* and can *rebuild* by using a different function.

The core problem in designing a dynamic hash table is how to atomically and efficiently distribute every node from the old hash table to the new one. This poses the following unique challenge: Even though the delete and insert operations on a single hash table can be atomic and non-blocking [14], [15], [16], [17], there are no non-blocking approaches that can atomically distribute every node from the old hash table to the new one. Consequently, prior dynamic hash tables [18], [19] all depend on locking mechanisms and need to first lock the corresponding buckets in the old and the new hash tables, before distributing each node (see Related Work for the pseudocode). However, it is well known that locking can prevent common operations from executing simultaneously, leading to unexpected delays or even suspension in a busy system. In Section 6.2, we demonstrate that for a typical hash-table usage pattern, when a rebuild operation is in progress, the 99.9%-percentile response time of these lock-based dynamic hash tables' lookup operations reaches 5 microseconds, 3.6 times worse than the normal case, adding significant delays to normal operations from legitimate users.

This paper presents DHASH, a novel non-blocking Dynamic Hash table. The major technical novelty of DHASH stems from a new *distributing mechanism* that can atomically and efficiently distribute each node without needing to acquire any per-bucket mutex locks. (To serialize concurrent rebuilding attempts, the rebuild operation uses a global lock, which, however, can never block any common operations. We discuss this in Section 4.4.) Experimental results show that the distributing mechanism is lightweight and does not noticeably degrade the performance of concurrent common operations. Furthermore, the distributing mechanism is independent of the implementation of hash table buckets, such that DHASH can utilize a variety of existing lock-free and wait-free set algorithms as the implementation of its buckets. This modular design allows programmers to make trade-offs among their own DHASH's progress guarantee, performance, and engineering efforts, and to select the set algorithms that meet their requirements best, bridging the gap between building robust hash tables and utilizing existing non-blocking set algorithms from the parallel-processing community.

We prove the correctness of our key techniques and implement two versions of DHASH that respectively provide lock-free and wait-free lookup operations. Experimental results show that DHASH can effectively change its hash function and seed. Moreover, our distributing mechanism is lightweight; when rebuilding, DHASH consistently outperforms the alternatives in terms of throughput and response time of concurrent common operations, at different concurrency levels, and with different operation mixes and average load factors. Furthermore, we demonstrate that DHASH's rebuild operation is scalable; the parallelized rebuild op-

eration of DHASH with eight or more rebuilding threads outperforms the state-of-the-art alternatives.

The rest of the paper is organized as follows. We first discuss related work in Section 2. Section 3 gives an overview of DHASH and the distributing mechanism. Section 4 presents the details of the algorithm. We prove the correctness of DHASH in Section 5, present evaluations in Section 6, and conclude in Section 7.

## 2 RELATED WORK

This section sketches a high-level overview of the existing dynamic and resizable hash tables.

**Dynamic hash tables in Linux kernel:** Xu et al. designed a dynamic hash table [18] for the management of IGMP packets in the Linux kernel in 2010. The key idea behind Xu's algorithm is to manage two sets of pointers in each node, so that common operations traverse one set of pointers while the rebuild operation is updating the other set. The two sets are exchanged upon the completion of every rebuild operation. Xu's algorithm is straightforward and easy to be implemented, but it has two major drawbacks in practice. (1) Locking mechanisms are used to serialize concurrent update and rebuild operations. (2) A linked list algorithm must be customized by adding an extra set of pointers before it can be used by Xu's hash table. In contrast, DHASH overcomes these drawbacks in its design.

Based on Triplett et al.'s ATC'11 paper [11], Graf introduced a dynamic hash table into the Linux kernel in 2014 [19], and this algorithm has been widely used in the kernel. Graf's hash table maintains a single pointer in each node, and utilizes per-bucket mutex locks to synchronize concurrent update and rebuild operations on each bucket. The skeleton of the rebuild operation of Graf's dynamic hash table is as follows.

```
for each bucket htbp in the old hash table
  mutex_lock(htbp);
  for each node htnp in htbp
    remove htnp from htbp;
    htbp_new = new_hash(htnp);
    mutex_lock(htbp_new);
    insert htnp into htbp_new;
    mutex_unlock(htbp_new);
  mutex_unlock(htbp);
```

Correspondingly, the insert and delete operations of Graf's dynamic hash table must first acquire the corresponding per-bucket mutex lock before they can operate on a bucket. Graf's algorithm is a practical design. However, this algorithm has the following drawbacks. (1) It uses locks to serialize updates to the same bucket. (2) The rebuild operation must traverse and reach the tail of a list before distributing each node. (3) It maintains unordered lists as its buckets and temporarily concatenates two lists when rebuilding, which dramatically increases the overhead of lookup operations. In contrast, DHASH overcomes these drawbacks in its design.

**Dynamic hash table in Perl:** Perl (versions between 5.8.2 and 5.16.2) uses a much more coarse-grained locking mechanism when rebuilding; the rebuild operation needs to acquire a global lock before distributing every node, preventing concurrent insert and delete operations from

making any progress, even though they work on different buckets.

**Lock-free resizable hash tables:** The resizable hash tables [10], [11], [13] do not change their hash functions; all they can do is enlarge or shrink their bucket sizes by a constant factor. One classic resizable hash table [10] was presented by Ori Shalev and Nir Shavit in 2006. To atomically distribute every node from the old hash table to the new one, Shalev and Shavit’s algorithm introduces a novel data structure called “recursive split-ordered list” which is fundamentally a lock-free linked list. Each hash table maintains a single “recursive split-ordered list,” and all nodes in the hash table are chained in this list. This hash table gradually assigns the bucket pointers to the places in the list where a sublist containing the requested node can be found. In solving the atomic-distribution problem when resizing, Shalev and Shavit’s algorithm does not move nodes among the buckets. Instead, it moves the buckets among the nodes by referencing buckets to the proper nodes in the list.

Even though the split-ordered list is lock-free, it has drawbacks in practice. (1) The algorithm must first bit-reverse the hash results of the keys before performing any common operations. Unfortunately, the bit-reverse operation is inefficient on X86 and Power9, where hardware does not provide native instructions. (2) Evaluations show that when the pre-defined hash function cannot evenly distribute incoming data, resizing cannot effectively help a hash table recover.

**Wait-free resizable hash tables:** Liu et al. [12] presented a wait-free, resizable hash table based on *freezable set abstraction* [20]. When the hash table array needs to be resized, the algorithm “freezes” and then distributes buckets. To be wait-free, Liu’s algorithm utilizes the fast-path-slow-path strategy [21]. Fatourou et al. [13] presented a wait-free, scalable resizable hash table based on extendible hashing [22]. To be wait-free, Fatourou’s algorithm utilizes the *PSim* universal construction [23]. Even though the wait-free hash tables provide the strongest progress guarantee, their implementation is much more complex compared with Shalev’s split-ordered list based hash table. Combined with the fact that they cannot change their hash functions and hence face the same challenges when there are attacks, for ease of presentation, in this paper, we choose Shalev’s algorithm as representative of resizable hash tables.

**Open-addressing based hash tables:** Researchers have proposed several resizable hash tables [24], [25] based on *open addressing* [26]. In particular, Maier et al. presented a resizable hash table [25] based on *linear probing* [26]—a form of open addressing—demonstrating that it is possible to build fast, open-addressing based resizable hash tables. The main technical contributions of Maier et al.’s algorithm and that of DHASH are orthogonal. We focus on separate-chaining based hash tables in this paper and leave applying our distributing mechanism to open-addressing based hash tables (e.g., Maier et al.’s algorithm) as future work.

### 3 DHASH ALGORITHM OVERVIEW

This section first presents the key challenge in designing dynamic hash tables, and then sketches a high-level overview

of the distributing mechanism of DHASH, leaving technical details to Section 4.

The key challenge in designing a dynamic hash table is to atomically and efficiently distribute every node from the old hash table to the new one. To distribute each node, the rebuild operation must update two hash tables. Even though the delete and insert operations on a single hash table can be atomic and non-blocking [14], [15], [16], [17], there are no non-blocking approaches that can atomically move a node from the old hash table to the new one. Prior research addressed this challenge by (1) maintaining two sets of linked lists in each bucket [18], (2) synchronizing the rebuild and other common operations by using locks [18], [19], (3) maintaining unordered linked lists [11], [19], and/or (4) degenerating dynamic hash tables to resizable hash tables [10]. These approaches, however, sacrifice the algorithms’ generality and/or performance.

**Distributing mechanism:** This paper presents a novel *distributing mechanism* that can atomically and efficiently distribute every node from the old hash table to the new one, without acquiring any mutex locks. The basic idea behind our distributing mechanism is that we can synchronize the rebuild operation and concurrent common operations by managing them to access hash tables and shared variables in specified orders, rather than serializing them by using expensive synchronization mechanisms such as locking. Specifically, to distribute a node, the rebuild operation of DHASH first deletes the node from the old hash table and then inserts it into the new hash table, without acquiring any mutex locks. The process of distributing the node leads to a short time period during which in neither the hash tables can this node be found. We call this a node’s *hazard period*. To allow other concurrent common operations to be able to access the node, DHASH employs a global pointer that always points to the node that is in the *hazard period*. Consequently, when a rebuild operation is in progress, concurrent common operations need to check different locations because a node may reside in either the new or the old hash table, or is referenced by the global pointer. In Section 4, we prove that if common operations search both hash tables and check the node currently in the *hazard period* in the specified orders, they can always find the node with the matching key and successfully apply the specified operations on it.

**Example of rebuilding:** To illustrate how DHASH’s distributing mechanism works, we use the example hash table shown in Figure 1. In this hash table, DHASH consists of two buckets: *Bkt 0* and *Bkt 1*. *Bkt 0* contains three nodes (a, b, and c), and *Bkt 1* contains two nodes (d and e). We assume that the new hash table contains three buckets, and that its user-provided hash function maps all of the keys to the new three-buckets array. The rebuild operation then performs a hash table traversal and distributes every node in the old hash table to the new one. Figure 1 illustrates the process of distributing the node a, with the initial state shown in Figure 1a and with the time advancing from figure to figure.

Specifically, the node a is first pointed to by the global pointer *rebuild\_cur*, resulting in the state shown in Figure 1b. Then, a is removed from the old hash table and enters its *hazard period*, shown in Figure 1c. When a is in the *hazard period*, other concurrent common operations can access it via the global pointer *rebuild\_cur*. Without loss of generality,

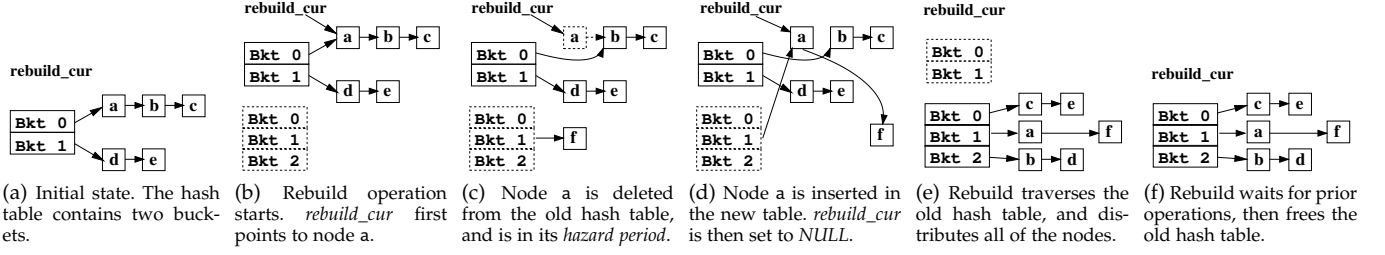


Fig. 1: The workflow of DHASH’s rebuild operation. *Bkt* is short for *Bucket*. *rebuild\_cur* points to the node in *hazard period*.

we assume that when the rebuild operation is in progress, other common operations concurrently insert a new node *f* into the new hash table, shown in Figure 1c. Then, the node *a* is inserted into the new hash table, shown in Figure 1d. After it has been successfully inserted into the new hash table, *rebuild\_cur* is set to NULL. The rebuild operation traverses the old hash table and distributes every node to the new hash table, shown in Figure 1e. After that, the rebuild operation waits for all prior unfinished operations to complete before safely reclaiming the old hash table, shown in Figure 1f.

Figure 1 illustrates the following key reasons that DHASH consistently outperforms prior dynamic and resizable hash tables. (1) By changing the hash function, DHASH effectively cuts the list length of the target buckets that contain many more nodes than the average load factor. Therefore, concurrent common operations can avoid traversing the lists that are unacceptably long. (2) Allowing other threads to access the node that is in the *hazard period* is the key reason that distributing a node is not necessary to involve expensive locking mechanisms in distributing the node. (3) In sharp contrast to prior algorithms [11], [19] that can only distribute the tail nodes of unordered lists, DHASH utilizes efficient set algorithms (e.g., lock-free skip lists [16]) and always distributes the head nodes. We demonstrate these advantages in Section 6.

## 4 DHASH IMPLEMENTATION

The design of DHASH presented in Section 3 leads to a relatively straightforward implementation, which is the subject of this section. Section 4.1 introduces concepts that provide a foundation for our design. We then discuss in detail how DHASH synchronizes rebuild operations and concurrent common operations.

Note that even though we use the RCU mechanism in this paper, it can be replaced by other synchronization mechanisms such as reference counters [15] and hazard pointers [27]. There is no direct impact on the distributing mechanism in DHASH by using the alternative techniques. However, they can lead to higher performance penalties. While using reference counters and hazard pointers, additional memory fence instructions are required each time a lookup operation moves forward to a new node. Updating the counter within each node and referencing a hazard pointer to the node may cause cache misses. Moreover, by using RCU, the rebuild operation of DHASH takes the benefit of the “waiting” mechanism provided by RCU [28] (in particular, the *synchronize\_rcu()* and *call\_rcu()* primitives). In contrast, while using reference counters or hazard pointers,

programmers need to implement the “waiting” mechanism by hand, which requires much more engineering effort.

### 4.1 Preliminaries

**Read-Copy Update:** RCU was originally developed for operating systems such as IBM’s K42 [29], and nowadays has been widely used in the Linux kernel [30] and user-space applications [31]. RCU works as a synchronization mechanism to address read-write conflicts. In particular, it distinguishes between the read-side code and the write-side code and has the following primitives:

- *rcu\_read\_lock()* / *rcu\_read\_unlock()* defines the read-side critical section. Each time a thread wants to access shared variables, it accesses them in a read-side critical section, which begins with the primitive *rcu\_read\_lock()* and ends with the primitive *rcu\_read\_unlock()*. Within a read-side critical section, the thread is safe to access the shared resources; it does not need to worry if these resources could be freed simultaneously by other threads.
- The hash table common operations of DHASH are protected by RCU read-side critical sections. That is, before invoking the function *ht\_lookup()*, *ht\_insert()*, or *ht\_delete()*, which will be discussed later, a caller must have entered an RCU read-side critical section by invoking *rcu\_read\_lock()*.
- *synchronize\_rcu()* works as a wait-for-readers barrier. This function blocks until all pre-existing RCU read-side critical sections have completed. It can be used in the rebuild operation of DHASH, for example, to ensure that any lookup operations that might have references to the old hash table complete before freeing the old hash table.
- *call\_rcu()* is an asynchronous version of *synchronize\_rcu()*. Instead of blocking, *call\_rcu()* register a callback function and argument, which can be invoked by a separate thread after all ongoing RCU read-side critical sections have completed. The caller to *call\_rcu()* can thus continue without blocking. The function *call\_rcu()* is particularly useful in situations where it is illegal to block or where update-side performance is critically important. For example, *call\_rcu()* is used in DHASH to safely reclaim nodes memory without blocking the delete operation.

RCU synchronizes readers with writers by using constrained access order, instead of shared variables [32]. Any RCU-protected node accessed by a reader is guaranteed to remain unreclaimed until the reader completes its access and calls *rcu\_read\_unlock()*. The production-quality implementations of *rcu\_read\_lock()* and *rcu\_read\_unlock()* are extremely lightweight; they have exactly zero overhead in the Linux kernels built for production use with *CONFIG\_PREEMPT=n* [30] and have extremely close to zero

overhead in user-space applications when the *QSBR flavor* model is used [31], such that readers of RCU-based data structures can execute as fast as single-threaded programs.

**API of DHASH’s bucket:** DHASH is modular, and its buckets can be implemented by using any lock-free or wait-free set algorithm that implements the Application Programming Interface (API) shown in Algorithm 1.

**Algorithm 1:** Structures and API of DHASH’s bucket.

---

```

1 struct node {long key; <node *ptr, flag> next};
  /* Has been logically removed by delete operation.          */
2 define LOGIC_RM (IUL < < 0)
  /* Is being distributed by rebuild operation.                */
3 define IN_HAZARD (IUL < < 1)
4 struct list {node *head};
5 struct snapshot {node **prev, *cur, *next};
  /* Search the node in list htpb.                             */
6 list_find(list *htpb, key, snapshot *sp)
  /* Insert node htpb into list htpb.                          */
7 list_insert(list *htpb, node *htnp)
  /* Search the node in list htpb, set the node's flag, and try
   to physically delete the node.                               */
8 list_delete(list *htpb, long key, long flag)

```

---

The data structures and API of DHASH’s bucket are shown in Algorithm 1. The structure *list*, which will be used as the implementation of DHASH’s buckets, is fundamentally a chain of *nodes*. For each *node*, the *key* field holds the key value, and the *next* field, which is a pointer, consists of two parts called *ptr* and *flag*. The *ptr* field points to the following node in the linked list if any, or has a *NULL* value otherwise. Since pointers are at least word aligned on all currently available architectures, the two least significant bits of *next* are used as the *flag* field indicating if the node is in a special state. The least significant bit, denoted as *LOGIC\_RM*, indicates that a node has been logically removed by a delete operation. The second least significant bit, denoted as *IN\_HAZARD*, indicates that a rebuild operation is distributing this node from the old hash table to the new one, and hence this node is in its *hazard period*. The difference between these two states is whether the node memory will be reclaimed when the node is physically removed from the list, which we will discuss in detail in the following paragraphs.

To return the search result of function *list\_find()* to the function invoking it, we borrowed the definition of the structure *snapshot* from Michael’s classic lock-free linked list [14]. Specifically, each time we want to search a node, an instance of *snapshot* is passed to *list\_find()*. Upon the completion of *list\_find()*, it is guaranteed that the *cur* field of the snapshot points to the list node containing the value that is greater than or equal to the specified search key.

The set algorithm is required to implement three functions: *list\_find()*, *list\_insert()*, and *list\_delete()*. The function *list\_delete()* takes the third parameter *flag*. If *flag* is set to *LOGIC\_RM*, *list\_delete()* deletes the matching node from the list and reclaims the node memory. In contrast, if *flag* is set to *IN\_HAZARD*, the node memory will not be reclaimed because the node will be inserted into the new hash table. The function *list\_delete()* does not block because it uses *call\_rcu()* to asynchronously reclaim a node. Note that *call\_rcu()* is safe to be invoked within an RCU read-side critical section [32].

**Choice of hash functions:** When rebuilding, DHASH heavily relies on the choice of the hash function for the

new hash table to resolve hash collisions. Although DHASH allows programmers to choose arbitrary hash functions, in practice, we suggest that they first utilize some widely-used non-cryptographic hash functions (e.g., Jenkins’ hash function [33]) that are fast and have been empirically proved to be “approximately” universal by changing their hash seeds [34].

If the target buckets are still unacceptably long after rebuilding, programmers are suggested to adopt a *universal hashing* approach [5]. In universal hashing, the system first creates a set of hash functions that are carefully designed, and programmers randomly select one hash function from the set for each hash table. Universal hashing guarantees that for the same group of nodes, two randomly-selected hash functions (one for the old hash table and another for the new one) are extremely unlikely to both cause bad hashing performance [5].

Overall, the proposed strategy for selecting new hash functions can provide good average-case hashing performance for DHASH, even if the existing hash function cannot evenly distribute the incoming data, or if an adversary has compromised the existing hash function. Specifically, for any existing nodes (including those in the target buckets), the expected length of the list in the new hash table that the node hashes to is at most the load factor  $\alpha = n/m$ , where  $n$  and  $m$  are the total number of the existing nodes in DHASH and the number of buckets in the new hash table, respectively (Theorem 11.3 of [35]). Moreover, for a rebuilding scenario where the  $m$  and  $q$  (the number of nodes in the target buckets) are 1,024 and 32, respectively, the probability of at least 4 and 16 out of these  $q$  nodes being hashed to the same bucket are  $3.34 \times 10^{-5}$  and  $4.21 \times 10^{-37}$ , respectively (Theorem 4 in [36]), which shows that by using the universal hashing approach, it is extremely unlikely that the hash collisions in DHASH can remain after rebuilding.

Note that if a programmer really cares about the worst-case behaviour of a universal hashing, it is possible for her to create a *perfect hashing* [37] for the new hash table that allows  $O(1)$  memory accesses to perform a search in the worst case, because when DHASH rebuilds, all keys of nodes are known ahead of time. However, a perfect hashing, which notably uses two levels of hashing with universal hashing at each level, can lead to a more intricate algorithm design. We therefore leave it for future work.

**Atomic primitives:** We model the main memory as an array of 64-bit values, and use the notation  $m[addr]$  for the value stored in the address *addr* of the memory. The pseudocode in this paper use the following 64-bit atomic primitives:

- *cas(addr, old, new)*, which stores the value *new* into  $m[addr]$  and returns the old value *old* if  $m[addr] = old$ ; otherwise, it returns the current value of  $m[addr]$ .
- *atomic\_fetch\_and(addr, val)*, which performs a bitwise AND operation on the value of  $m[addr]$  and *val*, and stores the result back to  $m[addr]$ .

**Two implementations of DHASH:** In this paper, we built and evaluated two versions of DHASH: one by using Michael’s classic lock-free linked list [14] as the implementation of its buckets, and another by using a concurrent linked list algorithm [38] that provides wait-free lookup operations. To keep our presentation self contained, we take Michael’s

linked list [14] as an example and discuss how we modify it before using it in DHASH in Appendix A.

## 4.2 Data structures

---

### Algorithm 2: Structures and auxiliary functions.

---

```

9 struct ht {ht *ht_new; long (*hash)(long key); int nbuckets; list *bkets[]};
/* Global variables. */
10 struct node *rebuild_cur;
11 mutex rebuild_lock;
12 #define logically_removed(htnp) (htnp->next & LOGIC_RM)
13 clean_flag(node *htnp, long flag) {
14     atomic_fetch_and(&htnp->next, <NULL, ~flag>);
15 }
16 ht_alloc(ht *htp, int nbuckets, long (*hash)(long key)) {
17     htp->ht_new := NULL; htp->hash := hash;
18     htp->nbuckets := nbuckets; htp->bkets := allocate(nbuckets);
19 }

```

---

Algorithm 2 lists the data structures and auxiliary functions of DHASH. The main structure, *ht*, is an array of buckets (*bkets[]*), with the array length stored in the *nbuckets* field. Each element of *bkets* is fundamentally a pointer to an instance of *list*, an RCU-based lock-free or wait-free set algorithm that provides the API shown in Algorithm 1. The *hash* field is a function pointer to the user-specified hash function. The *ht\_new* field is set to *NULL* unless a rebuild operation is in progress, in which case it points to the new hash table that is going to replace the old one. The global variable *rebuild\_cur* points to the node that is currently in the *hazard period* or is equal to *NULL* if there is no such a node. The mutex lock *rebuild\_lock* is to serialize attempts to rebuild the hash table. Note that the workload of rebuilding a hash table can be parallelized by creating a group of rebuilding threads, which will be discussed in Section 4.8.

The helper function *logically\_removed(htnp)* checks whether the *LOGIC\_RM* bit of the *next* field of the node pointed to by *htnp* has been set. The helper function *ht\_alloc()* creates a new instance of the hash table and allows users to specify its hash function. Note that the hash function, once set, cannot be changed, and is shared among all rebuilding threads and worker threads accessing this hash table instance.

## 4.3 Solving read-write conflicts within each bucket

Before diving into DHASH’s distributing mechanism, we first briefly discuss how we use the RCU mechanism to solve the read-write conflicts within each bucket. Without proper protection, lookup operations may point to nodes that have been reclaimed by concurrent delete operations, resulting in segmentation fault errors.

Specifically, before invoking the function *ht\_lookup()*, which will be discussed in Section 4.5, a caller must have entered an RCU read-side critical section by invoking *rcu\_read\_lock()*. The functions *ht\_delete()* and *ht\_insert()* also need read-side protection because they need to first traverse one or two lists. As one example, the following code snippet shows how the function *ht\_lookup()* can be safely used.

```

rcu_read_lock();
node *cur = ht_lookup(htp, key);
/* Accessing *cur is safe here. */
rcu_read_unlock();
/* Accessing *cur becomes unsafe. */

```

## 4.4 Rebuild operation

The pseudocode for the rebuild operation is shown in Algorithm 3. Line 22 attempts to acquire the global lock *rebuild\_lock*, which serializes concurrent rebuild attempts. Note that the lock *rebuild\_lock* is transparent to concurrent common operations and thus can never block any of them. Once a rebuild operation has the lock, it checks again that the rebuild is still required on line 23. Line 24 allocates a new hash table with the user-specified bucket array size and hash function. Line 25 assigns the new hash table to the *ht\_new* field of the old hash table, allowing subsequent operations to access the new hash table. Line 26 performs an RCU synchronization barrier to wait for prior common operations, which may not be aware of the new hash table, to complete before the rebuild operation continues.

---

### Algorithm 3: Rebuild operation of DHASH.

---

**Parameters:** *nbuckets*: Size of the bucket array.

*hash*: User-specified hash function.

```

20 Local variables: struct node *htnp, *htbp, *htnp_new; struct ht *htp_new;
21 void ht_rebuild(ht *htp, nbuckets, hash) {
22     if ( trylock(rebuild_lock) != SUCCESS ) return -EBUSY;
23     if ( ! rebuild_is_required() ) return -EPERM;
24     htp_new := ht_alloc(htp, nbuckets, hash);
25     htp->ht_new := htp_new;
/* Wait for operations not aware of htp_new. */
26     synchronize_rcu();
27     for (each bucket htbp in htp) {
28         for (each node htnp in htbp) {
29             rcu_read_lock();
30             rebuild_cur := htnp;
31             smp_wmb0();
32             key := htnp->key;
33             if (list_delete(htbp, key, IN_HAZARD) != SUCCESS) / Hflag
34                 continue;
35             clean_flag(htnp, IN_HAZARD)/* unHflag */;
36             htp_new := htp_new->bkets[htp_new->hash(key)];
37             if (list_insert(htbp_new, htnp) != SUCCESS) {
38                 | call_rcu(htnp, free);
39             }
40             smp_wmb0();
41             rebuild_cur := NULL;
42             rcu_read_unlock();
43         }
44     }
/* Wait for operations accessing nodes via htp->bkets[] */
45     synchronize_rcu();
46     htp_tmp := htp; htp := htp_new;
/* Wait for operations accessing old hash table. */
47     synchronize_rcu();
48     unlock(rebuild_lock);
49     free(htp_tmp);
50     return SUCCESS;
51 }

```

---

The function *ht\_rebuild()* traverses the old hash table and distributes every node to the new hash table (lines 27–44). Before distributing each node, the function *ht\_rebuild()* first enter an RCU read-side critical section (line 29) to prevent concurrent delete operations from erroneously reclaiming the node, while *ht\_rebuild()* is still accessing it. Then, the global variable *rebuild\_cur* points to the node on line 30. The two write barriers on lines 31 and 40 pair with the read barriers in *ht\_lookup()* and *ht\_delete()*. They together guarantee that *ht\_rebuild()*’s updates to *rebuild\_cur* and the two hash tables are performed in program order with respect to concurrent common operations. Note that, for ease of presentation, we omit memory order specifications in the pseudocode. In practice, all accesses to the bucket pointers (e.g., *htbp*), the node pointers (e.g., *htnp*), *ht\_new*,

and *rebuild\_cur* must be made with the specifications of `std::memory_order_acquired` or `release` [39].

Line 33 deletes the node from the old hash table. The function *list\_delete()* receives a third argument *IN\_HAZARD*, indicating that the node with the matching key will be deleted from the old hash table, but its memory will not be reclaimed. If this delete operation fails, which implies that the node has been deleted by other concurrent delete operations since the reference to the node was fetched on line 28, the rebuild operation skips this node (line 34). Otherwise, the node has been physically removed from the old hash table. Line 35 then prepares the node for reuse by cleaning the *IN\_HAZARD* bit and the *ptr* part of its *next* field. Note that if the *LOGIC\_RM* bit has been set by another concurrent delete operation, it remains. Since the node can be logically removed by other concurrent delete operations, *clean\_flag()* uses an atomic primitive (line 13). Then, line 37 inserts the node into the proper bucket of the new hash table. This insertion generally succeeds. It fails only if the node referenced by the global pointer *rebuild\_cur* is logically deleted by a concurrent delete operation (discussed in Section 4.6), and then a new node with identical key is inserted into the new hash table before line 37 is performed. In this case, line 38 invokes *call\_rcu()*, which reclaims the node after currently unfinished operations pointing to this node have been completed. After the node has been inserted into the new hash table, the global pointer *rebuild\_cur* is set back to *NULL*.

After distributing the existing nodes in the old hash table, line 45 waits for prior unfinished operations, which may be holding references to the distributed nodes. Line 46 sets the new hash table as the current one, and line 47 waits for prior unfinished operations that may be holding references to the old hash table. Then, line 48 releases the global lock, line 49 frees the old hash table, and finally, line 50 returns success.

For each iteration, *ht\_rebuild()* deletes a node from the old hash table and then inserts it into the new hash table, reusing the node's memory. One potential issue with this distributing mechanism is that it may redirect lookup operations traversing the old hash table to the wrong lists. For example, suppose that a lookup operation is traversing a hash bucket in the old hash table and is pointing to the node  $\alpha$ . At the same time, the rebuild operation distributes  $\alpha$  by inserting it into the new hash table. This can redirect the lookup operation to the bucket in the new hash table, and result in a false negative if the node with the matching key is at the bottom of the linked list in the old hash table. There are two solutions to this problem.

(1) Each node contains an extra integer field called *bkt\_id*, which records the *id* of the bucket where the node is stored. Before traversing a list, a lookup operation first reads and remembers the *id* of the bucket. Each time the lookup operation moves to a new node, it compares the value stored in the *bkt\_id* field with the stored bucket *id*, and starts over if they are not equal. This approach introduces negligible performance overhead because accessing the *bkt\_id* field does not involve any atomic instructions. The shortcoming of this approach, however, is that lookup operations may start over when a rebuild operation is distributing nodes. Note that, even though lookup operations may start over, they will not

be blocked because when rebuilding, the length of the lists in the old hash table is limited, and therefore the lookup operations can complete in a finite number of steps.

(2) After a node has been deleted from the old hash table, and before this node is inserted into the new hash table, the rebuild operation invokes a *synchronize\_rcu()* barrier and waits until all of the prior common operations complete and leave their RCU read-side critical sections. Even though this approach can slow down the rebuild operation, it saves the *bkt\_id* field in each node and prevents lookup operations from starting over. Given that rebuild operations should be relatively infrequent, we believe this approach is a valuable option for use cases where extremely fast lookup operations are required.

Overall, for DHASH implementations that need to provide extremely fast lookup operations, the second solution best meets the requirements. Otherwise, the first one is the solution of choice. The two implementations of DHASH presented in this paper (discussed in Section 6.1) use the first solution.

#### 4.5 Lookup operation

The pseudocode for the lookup operation is presented in Algorithm 4. This function first searches for the specified *key* value in the old hash table (line 54). If a node with the matching key can be found, a pointer to the node is returned (line 54). Otherwise, line 55 checks whether a rebuild operation is in progress, and line 55 returns *-ENOENT* if rebuild operations are absent. The two read barriers on lines 56 and 61 pair with the two write barriers in *ht\_rebuild()*. Line 58 continues the lookup operation by checking the node pointed to by the global pointer *rebuild\_cur*. Recall that *rebuild\_cur* always points to the node that is currently in its *hazard period*. If the node pointed to by *rebuild\_cur* matches, and if the *LOGIC\_RM* bit of the *next* field of this node has not been set, which means that the node has not been logically deleted by concurrent delete operations, line 59 returns a pointer to the node. Otherwise, *ht\_lookup()* continues by searching the new hash table and returns a pointer to the node if the function *list\_find()* succeeds (line 64).

---

#### Algorithm 4: Lookup operation of DHASH.

---

```

Local variables: struct node *cur, *htbp, *htbp_new; struct ht *htp_new;
                  struct snapshot ss;
52 node *ht_lookup(ht *htp, long key) {
53     htp := htp->bks[htp->hash(key)];
54     if (list_find(htbp, key, &ss) == SUCCESS) { return ss.cur };
55     if (htp->ht_new == NULL) { return -ENOENT };
56     smp_rmb();
57     cur := rebuild_cur;
58     if (cur and (cur->key == key) and !logically_removed(cur)) {
59         | return cur;
60     }
61     smp_rmb();
62     htp_new := htp->ht_new;
63     htbp_new := htp_new->bks[htp_new->hash(key)];
64     if (list_find(htbp_new, key, &ss) == SUCCESS) { return ss.cur };
65     else return -ENOENT;
66 }

```

---

Algorithm 4 shows that a lookup operation first searches for the node with the matching key in the old hash table (line 54), then checks the node referenced by the global pointer *rebuild\_cur* (line 58), and finally searches in the new hash table (line 64). This manipulation order guarantees that

lookup operations can always find the node. That is, the following lemma holds:

**Lemma 1.** *If DHASH contains a node  $\alpha$  with the key value of  $k$ , the function  $ht\_lookup(k)$  can find  $\alpha$ , no matter if a rebuild operation is in progress.*

*Proof.* Obviously, if there are no rebuild operations, the node  $\alpha$  resides in the only hash table. The function  $ht\_lookup(K)$  can find the node in the only hash table (lines 53 - 54).

We then prove that  $ht\_lookup(K)$  can find the node when a rebuild operation is in progress. The code snippet to distribute a node is shown on lines 30 - 41 in Algorithm 3. We use  $write_{rebuild}(rebuild\_cur, \alpha)$  to denote the event in which the thread running the rebuild operation (henceforth *rebuild thread* for short) writes the address of the node  $\alpha$  to the global variable  $rebuild\_cur$  (line 30), and use  $delete_{rebuild}(old, \alpha)$  and  $insert_{rebuild}(new, \alpha)$  to denote the events in which  $\alpha$  is deleted from and inserted into the old and the new hash tables, respectively (lines 33 and 37). Similarly, we use  $lookup_{lookup}(old, \alpha)$  and  $lookup_{lookup}(new, \alpha)$  to denote the events in which the thread running lookup operations (henceforth *lookup thread* for short) searches for the node  $\alpha$  in the old and the new hash tables, respectively (lines 54 and 64). We use  $lookup_{lookup}(rebuild\_cur, \alpha)$  to denote the event in which the lookup thread checks the node pointed to by  $rebuild\_cur$  (line 58). In the following proof, since the rebuild thread is the only thread that performs write, delete, and insert operations, and the lookup thread is the only thread that performs lookup operations, we omit the thread symbol without introducing any ambiguity. For brevity, we use the acronym  $rbc$  to stand for  $rebuild\_cur$ . One event  $e_1$  precedes another event  $e_2$ , written  $e_1 \prec e_2$ , if  $e_1$  occurs at an earlier time.

By inspecting the code of  $ht\_rebuild()$  in Algorithm 3, we get the following event sequence:

$$write(rbc, \alpha) \prec delete(old, \alpha) \prec insert(new, \alpha) \prec write(rbc, \perp) \quad (1)$$

By inspecting the code of  $ht\_lookup()$  in Algorithm 4, we get the following event sequence:

$$lookup(old, \alpha) \prec lookup(rbc, \alpha) \prec lookup(new, \alpha) \quad (2)$$

When the rebuild and the lookup threads are simultaneously accessing the node  $\alpha$ , there are three types of interleaving between these two threads:

- $lookup(old, \alpha) \prec delete(old, \alpha)$ , which implies that the lookup thread searches for the node  $\alpha$  before the rebuild thread starts distributing the node. Thus,  $\alpha$  can be found in the old hash table, and the lookup operation can return a pointer to  $\alpha$  on line 54.
- $insert(new, \alpha) \prec lookup(new, \alpha)$ , which implies that the lookup thread searches for the node  $\alpha$  after it has been inserted into the new hash table by the rebuild thread. Thus,  $\alpha$  can be found in the new hash table, and the lookup operation can return a pointer to  $\alpha$  on line 64.
- $delete(old, \alpha) \prec lookup(old, \alpha) \prec \dots \prec lookup(new, \alpha) \prec insert(new, \alpha)$ , which implies that the lookup thread searches for the node  $\alpha$  which is in the *hazard period*. Combined with the event sequences 1 and 2, we get the following event sequence:

$$write(rbc, \alpha) \prec delete(old, \alpha) \prec lookup(old, \alpha) \prec lookup(rbc, \alpha) \prec lookup(new, \alpha) \prec insert(new, \alpha) \prec write(rbc, \perp)$$

It follows that:

$$write(rbc, \alpha) \prec lookup(rbc, \alpha) \prec write(rbc, \perp)$$

Once the global pointer  $rbc$  is set to point to the node  $\alpha$  it remains. Hence the lookup thread can find  $\alpha$  via  $rbc$  and can return a pointer to it on line 59.

Overall, if there is a node with the matching key in DHASH, it is guaranteed that the function  $ht\_lookup(k)$  can find the node and return a pointer to it, no matter if a rebuild operation is in progress.  $\square$

#### 4.6 Delete operation

Recall that when DHASH is rebuilding, a node can (1) resides in either the old or the new hash table, or (2) is in the *hazard period* and is referenced by the global pointer  $rebuild\_cur$ . Therefore, the challenges that the delete operation of DHASH must address are (1) how to find the node, and then (2) how to delete this node which is concurrently being distributed from the old hash table to the new one by the rebuild operation.

To address the first challenge, the delete operation uses the same manipulation order as  $ht\_lookup()$  (shown in Algorithm 4). That is, the delete operation first searches for the node with the matching key in the old hash table, then checks the node referenced by the global pointer  $rebuild\_cur$ , and finally searches in the new hash table. Lemma 1 guarantees that if a node with the matching key is in DHASH, the function  $ht\_delete()$  can always find this node.

To address the second challenge, DHASH adopts a classic, lightweight mechanism, by separating the deletion of a node into two stages: *logical* and *physical* deletions [14], [40]. The first stage is to mark a node to prevent subsequent lookup operations from returning this node, and to prevent subsequent insert/delete operations from inserting/deleting nodes after this node. Marking a node is notably by using a *cas* operation that sets the least significant bit (i.e., the *LOGIC\_RM* in this paper) of the *next* field of the node. The second stage, which is typically performed by subsequent lookup operations, is to physically remove the node from the list by swinging the next pointer of the previous node to the next node in the list and then reclaiming the node memory.

The pseudocode for the delete operation is shown in Algorithm 5. The function  $ht\_delete()$  first attempts to delete the node from the old hash table on line 70. Then, it continues by checking if a rebuild operation is in progress on line 73. The two read barriers on lines 74 and 83 pair with the two write barriers in  $ht\_rebuild()$ . If a rebuild operation is in progress,  $ht\_delete()$  checks the node referenced by  $rebuild\_cur$ , and attempts to delete the node by setting the *LOGIC\_RM* bit if the node has the expected key value (line 79). Since the rebuild operation can simultaneously clean the *IN\_HAZARD* bit (line 35) and/or update the *ptr* field (line 35 or 37) of this node,  $ht\_delete()$  uses a loop that repeatedly sets the *LOGIC\_RM* bit by using a *cas* instruction, which can fail at most two times. If other concurrent delete operations have successfully deleted this node,  $ht\_delete()$  skips this node (line 78). The delete operation continues by attempting to delete the node with the matching key from the new hash table (line 85). If the delete operation fails, line 87 returns *-ENOENT*, indicating that no node with the matching key can be found in DHASH.



---

**Algorithm 5: Delete operation of DHASH.**


---

```

67 Local variables: struct node *cur, *htbp, *htbp_new; struct ht *htp_new;
68 int ht_delete(ht *htp, long key) {
69     htbp := htp->bkts[htp->hash(key)];
70     if (list_delete(htbp, key, LOGIC_RM) = SUCCESS)/*Rflag*/
71         return SUCCESS;
72     htp_new := htp->ht_new;
73     if (htp_new = NULL) return -ENOENT;
74     smp_rmb();
75     cur := rebuild_cur;
76     while (cur and (cur->key = key)) {
77         cur_next := cur->next;
78         if (cur_next & LOGIC_RM) break;
79         if (cas(cur->next, cur_next, cur_next | LOGIC_RM))/*Rflag*/
80             return SUCCESS;
81         cur := rebuild_cur;
82     }
83     smp_rmb();
84     htp_new := htp_new->bkts[htp_new->hash(key)];
85     if (list_delete(htbp_new, key, LOGIC_RM) = SUCCESS)/*Rflag*/
86         return SUCCESS;
87     return -ENOENT;
88 }

```

---

To prove that the delete operation of DHASH can successfully delete the expected node, while a rebuild operation is distributing this node, we illustrate the possible values of the two least significant bits of the *next* field (the *flag* field henceforth) of a node and the sequence of changes the node's *flag* field can go through in Figure 2. The labels in the boxes show the values of the *flag* field, and the labels of the transitions show the corresponding *cas* instructions and line numbers in the pseudocode. For example, the right-most box shows that both the *IN\_HAZARD* and *LOGIC\_RM* bits have been set, indicating that while the node is being distributed by the rebuild operation (between lines 33 and 35), a delete operation finds it via the global pointer *rebuild\_cur* and successfully deletes it on line 79. The *flag* field of a node is initially clean (the left-most box) and can only be changed by the following *cas* operations. (1) To distribute the node, the rebuild operation sets the *IN\_HAZARD* bit of the node on line 33, and after physically removing this node out of the old hash table, cleans the *IN\_HAZARD* bit on line 35. We refer to these two *cas* steps as *Hflag* and *unHflag*, respectively. (2) To delete the node, a delete operation sets the *LOGIC\_RM* bit of the node on either line 70, 79, or 85. We refer to these *cas* steps as *Rflag*.

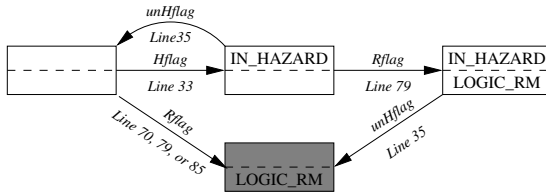


Fig. 2: Sequence of changes a node's *flag* field can go through. The labels in the boxes show the values of the *flag* field of the node. Each of the transitions corresponds to a successful *cas* instruction, and the label shows the specific *cas* instruction along with its line number in the pseudocode.

By examining the pseudocode in Algorithms 3 and 5, we can see that these *cas* operations proceed in the orderly way shown in Figure 2, and that the sequence of changes shown in Figure 2 reflects all possible changes to the *flag* field of the node. For example, each time a node is distributed from the

old hash table to the new one, its *flag* field walks through the top left circuit (transitions labeled *Hflag* and *unHflag*). Figure 2 illustrates that the following statements are true for each node.

- No matter if the *IN\_HAZARD* bit has been set, exactly one delete operation can succeed, by setting the *LOGIC\_RM* bit of the node (via one of the two *Rflag* CASes).
- If the *IN\_HAZARD* bit is set, it is eventually cleaned via one of the two *unHflag* CASes, and the *unHflag* CASes leave the *LOGIC\_RM* bit untouched.
- Once the *LOGIC\_RM* bit is set, it remains; a node that has been deleted eventually reached the state represented by the shaded box in Figure 2.

Moreover, by examining the code, we can see that once a node's *LOGIC\_RM* bit is set, no concurrent lookup operations can return this node. Instead, they help physically remove this node. With all these facts, we get the following lemma:

**Lemma 2.** *No matter where the node is (either in one of the hash tables, or is referenced by rebuild\_cur), a delete operation can successfully delete it by setting its LOGIC\_RM bit.*

Overall, we can prove that the following lemma holds:

**Lemma 3.** *If DHASH contains a node  $\alpha$  with the key  $k$ , the function  $ht\_delete(k)$  can delete the node  $\alpha$  by successfully setting its LOGIC\_RM bit, no matter if a rebuild operation is in progress.*

*Proof.* The function  $ht\_delete(k)$  fundamentally performs a lookup operation along with a *logical deletion* if the node can be found. Lemma 1 guarantees that if a node with the required key is in DHASH, the function  $ht\_delete()$  can always find this node. Moreover, Lemma 2 guarantees that once the node can be found, the delete operation can successfully delete it.  $\square$

Note that Figure 2 also allows us to naturally choose the linearization points for delete operations: each successful delete operation is linearized at the successful *Rflag* CAS.

#### 4.7 Insert operation

The basic idea of DHASH's insert operation is as follows. When there are no rebuild operations, which is the common case, DHASH inserts new nodes into the only (i.e., old) hash table. Otherwise, new nodes are inserted into the new hash table. (For ease of presentation, we refer to insert operations that add a node with key  $k$  into the old and the new hash table as  $IO(k)$  and  $IN(k)$ , respectively.)

However, simply applying this strategy in practice might lead to concurrency issues because when rebuild operations are in progress,  $IO(k)$  and  $IN(k)$  may exist simultaneously, such that some synchronization mechanisms are required to prevent  $IO(k)$  and  $IN(k)$  from inserting duplicate nodes in DHASH.

To that end, we adopt a lightweight and conceptually straightforward synchronization strategy that consists of two aspects. (1) Once a rebuild operation starts,  $IO(k)$  stops inserting nodes into the old hash table. Instead,  $IO(k)$  starts over and then attempts to insert the node into the new hash table (i.e., it becomes an  $IN(k)$ ). (2) An  $IN(k)$  first checks if

any node with key  $k$  already exists in the old hash table or is currently referenced by *rebuild\_cur*, before attempting to insert its own node into the new hash table.

Specifically, we implement the function *list\_insert\_dcsc()*, a variant of the standard insert operation of the underlying buckets discussed in Algorithm 1, that allows us to insert a node into the underlying bucket only if rebuild operations are absent. The function *list\_insert\_dcsc()* returns *-EADDR1* if a rebuild operation started before it linearizes; otherwise, it works as if it were the standard *list\_insert()* (i.e., it returns *SUCCESS* if the node is successfully inserted into the list, and *-EEXIST* if such a node already existed in the bucket.)

The implementation of the function *list\_insert\_dcsc()* is similar to that of the standard *list\_insert()*, except that the statement where a successful *list\_insert()* operation linearizes, which is notably a *cas* or write instruction, is replaced with Harris et al.’s lock-free double-compare-and-single-swap (*dcsc*) primitive [41]. The *dcsc* primitive is implemented from normal *cas* instructions, and takes five arguments: two addresses, two expected values, and one new value, and can atomically (1) read the two memory addresses, (2) check if they contain the expected values, and (3) if so, write the new value into the second address. Therefore, by setting the first memory address and the expected value in this address to *htp->htp\_new* and *NULL*, respectively, the function *list\_insert\_dcsc()* can successfully insert the node into the list only if rebuild operations are absent. Once a rebuild operation has started, and the *htp->htp\_new* field has been set to point to the newly-allocated hash table on line 25, *list\_insert\_dcsc()* will fail and start over.

The use of the *dcsc* primitive is motivated by Arbel-Raviv and Brown’s non-blocking range query algorithms [42]. Appendix A shows the implementation of *list\_insert\_dcsc()* for Michael’s classic lock-free linked list algorithm [14].

---

**Algorithm 6: Insert operation of DHASH.**

---

```

89 Local variables: node *htnp, *htbp, *htbp_new; ht *htp_new;
90 int ht_insert(ht *htp, long key) {
91     htp := allocate_node(key);
92     retry:
93     htp_new := htp->htp_new;
94     if (htp_new = NULL) {
95         /* IO(k) */
96         htp := htp->bkt[hash(key)];
97         ret := list_insert_dcsc(&htp->htp_new, NULL, htp, htp);
98         if (ret = SUCCESS) return SUCCESS;
99         else if (ret = -EADDR1) goto retry (line 92);
100     }
101     else {
102         /* IN(k) */
103         if (ht_lookup(htp, key) = SUCCESS) return -EEXIST;
104         htp_new := htp->bkt[hash(key)];
105         if (list_insert(htbp_new, htp) = SUCCESS)
106             return SUCCESS;
107     }
108     free(htnp);
109     return -EEXIST;
110 }
```

---

Algorithm 6 shows the pseudocode for the insert operation of DHASH. The function first allocates a new node, initializes it (line 91), and then checks if a rebuild operation is in progress (line 94). For an *IO(k)*, the function *list\_insert\_dcsc()* will fail if a rebuild operation starts during the time gap between lines 93 and 96 of this insert operation. In this case, *list\_insert\_dcsc()* returns an error message *-EADDR1*. Upon receiving this message, *IO(k)* starts over

(line 98), and in the subsequent insertion attempt, it can see the new hash table (line 93) and will insert the node into it. That is, it becomes *IN(k)*. Note that the RCU mechanism (line 26) prevents another rebuild operation from starting before this *IO(k)* completes, such that an *IO(k)* restarts at most once. That is, the following lemma holds:

**Lemma 4.** *Once a rebuild operation has started,  $IO(k)$  will fail, start over, and then attempt to insert the node into the new hash table; that is, it becomes  $IN(k)$ .*

An *IN(k)* first performs a lookup operation on line 101, and then attempts to insert the node into the new hash table (line 103). By examining the code, we can prove that the following lemma holds:

**Lemma 5.** *When  $IN(k)$  attempts to insert a node  $\alpha$  with key  $k$  into the new hash table on line 103, it is guaranteed that (1) the old hash table does not contain any node with key  $k$ , and (2) the node referenced by *rebuild\_cur* does not have the key  $k$ .*

*Proof.* When *IN(k)* performs the lookup operation on line 101, which implies that a rebuild operation has started, Lemma 4 guarantees that no new nodes will be inserted into the old hash table. Moreover, if any existing node in DHASH has the key  $k$ , Lemma 1 guarantees that this lookup operation can find it, and return the error message *-EEXIST*. Since a new node with key  $k$  can be inserted into the new hash table during the time gap between lines 101 and 103 of this insert operation, we can only guarantee that the node with key  $k$  cannot appear in the old hash table and cannot be referenced by *rebuild\_cur* when *IN(k)* performs the line 103.  $\square$

With Lemma 4 and 5, now we can prove that the following lemma holds:

**Lemma 6.** *If DHASH does not contain a node with key  $k$ , the function *ht\_insert(k)* can successfully insert a node  $\alpha$  with key  $k$  in DHASH; otherwise, the error message *-EEXIST* is returned without changing the hash table, no matter if a rebuild operation is in progress.*

*Proof.* Obviously, if there are no rebuild operations, the value of *htp->htp\_new* is always equal to *NULL*, such that the function *ht\_insert(k)* always attempts to insert the node  $\alpha$  into the only hash table (lines 94 – 98).

When rebuild operations are in progress, *IO(k)* and *IN(k)* can exist simultaneously. However, Lemma 4 guarantees that *IO(k)* will start over and become *IN(k)*. Therefore, the function *ht\_insert(k)* always attempts to insert the node  $\alpha$  into the new hash table. Moreover, Lemma 5 guarantees that a duplicate node with key  $k$  can neither exist in the old hash table nor be referenced by *rebuild\_cur*.

Overall, DHASH always attempts to insert new nodes into either the old or the new hash table, depending on whether rebuild operation are absent, such that the insert operation of DHASH complies with the standard specification for the insert operation of a normal hash table.  $\square$

Note that to support the function *list\_insert\_dcsc()*, the insert operation of the set data structure used as the implementation of DHASH’s buckets needs to linearize at a write or a *cas* instruction, and we assume that programmers know where this write or *cas* occurs in the code. This assumption is

true for most non-blocking set algorithms that are linearizable, including the two used as DHASH's buckets in this paper [14], [38]. Moreover, the original implementation of the *dcss* primitive [41] uses a helping mechanism, which can lead to poor performance because each *dcss* operation needs to dynamically allocate (and later, free) an *Info* object that holds the required information for helping. To address this performance issue, we applied the following optimizations in our implementation. (1) Motivated by Arbel-Raviv and Brown's research results [43], the *Info* objects for the *dcss* operation are pre-allocated and reused. Additional fields such as sequence numbers and flags are used to prevent insert operations from erroneously reading improper information from these shared, global objects [43]. (2) The *Info* objects are carefully placed and properly aligned in memory to avoid cache thrashing [44], when multiple insert operations update them simultaneously. (3) Lookup operations that might access these objects utilize hardware prefetch unit as much as possible by prefetching the possible objects before accessing their contents. Experimental results show that our *dcss* implementation introduces negligible overhead to DHASH in terms of throughput and operation latency. We refer the interested readers to the source code of DHASH for our implementation of the *dcss* primitive.

Note that another approach to synchronizing  $IO(k)$  and  $IN(k)$  is by explicitly using the classic *helping mechanism* [45], [46]. This approach allows an  $IN(k)$  to first help  $IO(k)$  complete to avoid inserting duplicate nodes into DHASH, avoiding the use of the *dcss* primitive. However, the insert operation with this explicit helping mechanism is relatively intricate. Moreover, since each insert operation must go through at least four state transitions, each of which consists of a *cas* instruction, its performance is inferior to the *dcss*-based version. Therefore, we present the function *ht\_insert()* using the *dcss* primitive in the paper, and refer interested readers to the code, which will be open-sourced, for the version utilizing the helping mechanism.

#### 4.8 Parallelizing rebuild operation

Even though the rebuild operation of DHASH is lightweight (demonstrated in Section 6.3), it is expected to be completed as soon as possible. To that end, we parallelize the rebuild operation by leveraging the *thread-level parallelism* [44].

The basic idea behind our parallelizing strategy is that in distributing the nodes of the old hash table, different buckets can be processed in parallel, since they are mainly independent of one another. Specifically, when rebuilding, the rebuild operation of DHASH create a group of  $n$  threads (referred to as *rebuilding threads*), each of which is bind to a dedicated CPU core. We then divide the bucket range of the old hash table to  $n$  interleaved regions, and then assign different regions to different rebuilding threads. That is, the bucket  $b$  is assigned to the rebuilding thread  $i$  if  $b \% n = i$ . Each rebuilding thread takes the argument of a region of buckets, and distributes the nodes in these buckets to the new hash table by using the program logic shown on lines 27 – 44 in Algorithm 3. The only modification is that *rebuild\_cur* becomes a global, per-thread variable; each rebuilding thread has its own *rebuild\_cur*, which points to the node that is being distributed by this thread and is in the

*hazard period*. Rebuilding threads may simultaneously insert nodes into the same bucket in the new hash table, which is synchronized by the lock-free or wait-free set algorithm used as the implementation of the buckets.

Accordingly, each time a lookup, insert, or delete operation wants to check if a node with the matching key is in the *hazard period*, it checks the *rebuild\_cur* pointers of all rebuilding threads. For example, for the function *ht\_lookup()* (Algorithm 4), the only modification is that the code snippet on lines 57 – 60 will be included in a for-loop statement, and that each iteration checks the node that is being pointed to by one rebuilding thread's *rebuild\_cur* pointer. One benefit of our parallelizing strategy is that except for letting the lookup, insert, and delete operations to check a group of *rebuild\_cur*, the control flow of these operations remains, significantly simplifying the proof of the correctness of DHASH when a parallelized rebuild operation is running.

### 5 CORRECTNESS

In this section, we prove that DHASH is linearizable and implements an abstract set object in a lock-free or wait-free manner, depending on the implementation of the hash table buckets.

**System model:** We design DHASH for an *asynchronous shared memory system* [47] with multiprocessors. On such a system, a program is executed by  $p$  deterministic threads, where  $p$  may exceed the number of physical processors. A scheduler decides which threads to run and may suspend the execution of any thread at any time for arbitrarily long. We assume a TSO memory model [48] in this paper. However, today's programming languages and architectures provide weaker ordering guarantees and allow accesses to be reordered for performance. To prevent any undesired reordering, additional memory fences and/or compiler directives are required to guarantee the order in which multiprocessors access different shared memory locations. We omit these memory fences and compiler directives in the pseudocode and refer the interested readers to the source code of DHASH.

**Set semantics:** We first define the set  $H$  of nodes of DHASH in any given state as follows: if rebuild operations are absent,  $H$  is composed of the non-logically-removed nodes in the only (i.e., old) hash table; otherwise,  $H$  is the union of (1) the non-logically-removed nodes in both the old and the new hash tables, and (2) the non-logically-removed nodes that are referenced by the *rebuild\_cur* pointers.

We then claim that DHASH complies with the abstract set semantics, and we use the standard sequential specification defined in [35] that includes the following three functions:

- The *ht\_lookup(k)* operation returns a pointer to the node with key  $k$  in  $H$ , *-ENOENT* otherwise.
- The *ht\_insert(k)* operation returns *SUCCESS* if a node with key  $k$  is successfully inserted into  $H$ , and *-EEXIST* if such a node already existed in  $H$ .
- The *ht\_delete(k)* operation returns *SUCCESS* if the node with key  $k$  was successfully deleted from  $H$ , and *-ENOENT* if such a node cannot be found in  $H$ .

We first prove that DHASH is linearizable to this sequential specification. That is, each of DHASH's common

operations has specific linearization points, where the operation takes effect and maps the operation in our concurrent implementation to sequential operations so that the histories meet the specification.

Recall that DHASH is modular, and that any linearizable lock-free or wait-free set algorithm (e.g., linked lists) that implements the API listed in Algorithm 1 can be used as the implementation of DHASH’s bucket. The execution histories of DHASH include sequences of list operations (*list\_find*, *list\_insert*, and *list\_delete*) which can be treated as atomic operations because they are linearizable objects [14], [47]. Moreover, we prove in Appendix A that the function *ht\_insert\_dcscs()* is linearizable.

**Linearization points:** For each of DHASH’s common operations, linearization points exist within its execution interval, so that at this point, it atomically reads or changes the abstract state  $H$ .

(1) Every lookup operation that finds the node with the expected key via *rebuild\_cur*, or in the old or the new hash table is linearized on line 57, 54, or 64, respectively. When rebuild operations are absent, every lookup operation that returns *-ENOENT* is linearized on the first invocation of *list\_find()* (line 54). In contrast, when rebuild operations are in progress, defining the linearization point for a lookup operation that returns *-ENOENT* is a bit intricate, because there is a time gap between *ht\_lookup()* searches the old hash table (on line 54) and checks the existence of rebuild operations (on line 55), such that a new node with the desired key might be inserted into the old hash table after it has been searched. Therefore, in this case, we linearize the unsuccessful lookup operation within its execution interval at the earlier one of the following points: (1) the point immediately before a new node with the expected key is inserted into the old hash table, and (2) the point where the search on the new hash table returns *-ENOENT*. This linearization point definition is motivated by Heller et al.’s classic concurrent set implementation [38].

(2) Similar to the lookup operation, every successful delete operation that finds the node with the required key via *rebuild\_cur* takes effect on line 79. For other successful cases, the delete operation linearizes in either of the two invocations of *list\_delete()* (line 70 or 85). When rebuild operations are absent, every delete operation that returns *-ENOENT* linearizes on the first invocation of *list\_delete()* (line 70). When rebuild operations are in progress, we linearize the unsuccessful delete operation within its execution interval at the earlier one of the following points: (1) the point immediately before a node with the expected key is inserted into the old hash table, and (2) the point where the invocations of *list\_delete()* on line 85 returns *-ENOENT*.

(3) Every successful *IN(k)* linearizes in the invocation of *list\_insert()* on line 103. Unsuccessful *IN(k)* linearizes on line 101 or 103, depending on where the existing node with the required key is found. *IO(k)* that fails because of the existence of concurrent rebuild operations (i.e., the return value of *list\_insert\_dcscs()* is equal to *-EADDR1*) will attempt to insert the node into the new hash table and linearizes as an *IN(k)* does. Otherwise, *IO(k)* linearizes in the invocation of *list\_insert\_dcscs()* (line 96).

Given the above definition, it is straightforward to deduce from Lemmas 1, 3, and 6 that for each of DHASH’s

common operations, within its execution interval, there is a linearization point at which the operation reads or modifies the abstract state  $H$ , according to the specified set semantics. That is, the following lemma hold:

**Theorem 1.** DHASH is a linearizable implementation of a sequential set object.

**Progress guarantee:** The lookup, insert, and delete operations of DHASH are non-blocking and can provide lock-free or wait-free progress guarantee, depending on the underlying set algorithm used.

(1) A lookup operation *ht\_lookup()* invokes the list operation *list\_find()* twice. Other statements are regular instructions, which can complete in a finite number of CPU cycles. One corner case is that for some list implementations, the function *list\_find()* may need to help other concurrent insert operations, which can be continuously inserting new nodes into the same linked list via the *dcscs* primitive (discussed in Appendix A). However, the progress guarantee of *list\_find()* does not degenerate, despite this helping workload, because each time it helps a pending insert operation, it moves forward one node, and this can happen only a limited number of times. Otherwise, a rebuild operation is incurred, preventing subsequent insert operations from inserting nodes into this list in the old hash table.

(2) A delete operation invokes the list operation *list\_delete()* twice, and utilizes *call\_rcu()* that is non-blocking. As discussed in Section 4.6, the *cas* operation on line 79 fails at most twice. Other statements are regular instructions.

(3) The function *list\_insert\_dcscs()* is a variation of the function *list\_insert()*, such that they have the same progress guarantee. An *IO(k)* jumps to the label *retry* and starts over at most once. Therefore, an insert operation invokes either the function *list\_insert\_dcscs()* or *list\_insert()* once for most cases, and when there are concurrent rebuild operations, it may invoke both of them once. Other statements are regular instructions.

In summary, DHASH provides a wait-free framework for common operations. For the implementation of DHASH utilizing Michael’s lock-free linked list, the lookup, insert, and delete operations are lock-free. However, since DHASH is modular, programmers can instead choose a wait-free linked list and build their own DHASH that provides a stronger progress guarantee.

Note that the rebuild operation of DHASH could block. Specifically, *ht\_rebuild()* serializes concurrent rebuild attempts by using a mutex lock and waits for prior hash table common operations by using the *synchronize\_rcu* barriers. However, this mutex lock is only used by the rebuild operations to serialize concurrent rebuild attempts, and hence can never block concurrent common operations. Similarly, the RCU barriers in the rebuild operation are used to force a build operation to wait for the existing common operations. However, they are transparent to common operations, and cannot block any of them. Moreover, rebuild operations are infrequent, and their speed is not the major concern if they do not noticeably degrade the performance of concurrent common operations. We thus leave making the helper function *ht\_rebuild* non-blocking as our future work.

## 6 EVALUATIONS

DHASH is the first dynamic hash table providing lock-free and wait-free common operations. Therefore, we need to prove that it is robust, efficient, and scalable. In this section, we demonstrate that (1) DHASH’s rebuild operation is effective, (2) DHASH’s rebuild operation is lightweight and does not noticeably degrade system performance in terms of throughput and response time of concurrent common operations, and (3) DHASH’s rebuild operation is fast and scalable; that is, users can obtain an approximately linear speedup by creating more rebuilding threads.

### 6.1 Evaluation Methodology

**Hash table implementations:** We chose Xu’s algorithm [18] as representative of dynamic hash tables that maintain two sets of list pointers in each node, and refer to this algorithm as HT-Dyn-Classic. Another dynamic hash table widely used is the *rhashtable* algorithm in the Linux kernel [19] (henceforth referred to as HT-Dyn-Linux). We chose *HT-Dyn-Linux* as representative of dynamic hash tables that maintain one unordered linked list for each hash bucket. Both *HT-Dyn-Classic* and *HT-Dyn-Linux* use locks to synchronize concurrent insert and delete operations on the same bucket. We chose Shalev and Shavit’s algorithm (henceforth referred to as HT-Resizable) as representative of resizable hash tables.

To demonstrate the benefits of DHASH’s modularity, we implemented and evaluated two versions of DHASH. *HT-DHash-lf* utilizes Michael’s classic lock-free linked list [14] as the implementation of its buckets. The lookup operation of *HT-DHash-lf* needs to start over from the beginning of the list once it encounters *logically* removed nodes. Moreover, we implemented *HT-DHash-wf*, another version of DHASH that provides wait-free lookup operations, by utilizing Heller et al.’s concurrent linked list algorithm [38]. Note that Heller et al.’s linked list provides wait-free lookup operations, but its insert and delete operations are serialized by an optimistic, fine-grained locking mechanism. We did not choose the linked list algorithms that provide wait-free insert and delete operations (e.g., [17]) for the following reasons. (1) Wait-free linked lists typically employ helping mechanisms that are inefficient in practice [49]. (2) They make the unrealistic assumption of a wait-free memory management scheme, which most programming environments (e.g., our benchmarking framework) do not support.

In the experiments, we chose Jenkins’ hash function [33]. The reasons are as follows. (1) Jenkins’ hash function is fast; it takes about 25 CPU cycles to hash a 64-byte integer and to produce a 32-bit hash result. (2) Jenkins’ hash function provides a good uniform distribution [33]. (3) Jenkins’ hash function takes a *hash\_seed* parameter, and for a benchmarking framework, the hash function with different *seed* values can be regarded as different hash functions. Specifically, in our benchmarking framework, when a hash table is created, an initial *seed* value is specified by users. For each common operation, Jenkins’ hash function is invoked to generate a 32-bit hash result for the key, and the hash result modulo the size of the bucket array is used as the index to the bucket array. When rebuilding, users change the *seed* value.

We implemented DHASH and the above-mentioned alternatives as user-space programs by using C. For all of

the hash tables, optimizations such as cache-line padding are applied. Our benchmarking framework utilizes Desnoyers’ user-space implementation of RCU [31]. Specifically, the *QSBR flavor* [31] model is used, such that the lookup operations can perform as fast as single-threaded programs. We compile the code with GCC 7.4.0 on all platforms where Ubuntu 18.04.4 is installed. We use -O3 as our optimization level without any special optimization flags.

**Benchmarking framework:** To compare DHASH with the alternatives, we built a benchmarking framework for evaluating dynamic hash tables. Some common features of benchmarking frameworks, such as creating worker threads and setting the CPU-affinity of the threads, were borrowed from Paul E. McKenney’s open-source *perfbok* project [32].

The framework consists of a specified number of *worker threads*, each of which performs the workload with the specified distribution of *lookup*, *insert*, and *delete* operations, denoted as  $m$ . The range of keys  $U$  is set to ten million to prevent the CPU caches from buffering the whole test set. When a test starts, each worker thread performs an infinite loop. In each iteration, the worker thread randomly selects an operation type according to the specified distribution  $m$ , randomly chooses a key ranging from 0 to the specified upper bound  $U$ , and then performs the specified operation. We set the initial size of the bucket array  $\beta$  to 1K. We controlled the average load factor  $\alpha$  indirectly by inserting  $\alpha * \beta$  nodes in a hash table before starting a test, and by selecting the ratio of insert operations to be equal to that of delete operations.

When a test starts, a *rebuild daemon* is created. Each bucket maintains a per-bucket atomic counter that records the number of nodes in this bucket. A successful insert operation increments this counter, and if it becomes larger than a pre-defined threshold, a signal is sent to the rebuild daemon. Once the rebuild daemon receives a rebuild signal, it first calculates the current average load factor of the hash table by adding all of the per-bucket counters up and dividing the result by the size of the hash table’s bucket array. If the current average load factor becomes 1.2 times larger than the user-specified threshold, the rebuild daemon (1) performs a rebuild operation by changing the *seed* value of Jenkins’ hash function, and at the same time, (2) enlarges the hash table by doubling the size of the bucket array. Otherwise, the rebuild thread only changes the *seed* value. The rebuild daemon then creates a group of rebuilding threads which, together, carry out the rebuilding task. For HT-Resizable, all its rebuild daemon can do is double the size of its bucket array, because HT-Resizable can change neither its hash function nor the *seed* value.

An *attacker thread* continuously inserts nodes with the keys mapped to a few target buckets of the evaluated hash table, simulating uneven distribution. The attacker thread inserts these nodes at the rate of 300 nodes per second, simulating malicious traffic that uses low bandwidth and is hard to be detected. The attacker thread is mapped to a dedicated CPU core. In mapping the worker threads and the rebuilding threads, we use a *performance-first* mapping; a new thread is mapped to the CPU core that has the smallest number of threads running on it. In this paper, experiments performed on a single CPU socket are marked with an \*, experiments performed on multiple CPU sockets

are marked with a #, and experiments in which all of the threads oversubscribe CPU cores are marked with a !.

TABLE 1: Summary of experimental platforms.

Processor Model	Speed	#Sockets	#Cores	LLC	Mem.
Intel Ivy Bridge	2.6 GHz	2	24	15 M	64 G
IBM Power9	2.9 GHz	1	16	80 M	16 G
Cavium ARMv8	2.0 GHz	2	96	16 M	32 G

**Hardware platforms:** It is worth noting that the performance of the evaluated hash tables varies at different architectures because of the unique hardware features. In particular, most modern x86 processors provide a hardware *fetch-and-add* instruction that can update the per-bucket counters in a few CPU cycles. However, Power9 and ARMv8 do not natively support this instruction, and the compiler needs to use the code snippet of a loop of *compare-and-swap* instructions to simulate one *fetch-and-add*, incurring overhead at high concurrency levels. Moreover, HT-Resizable needs to bit-reverse hash results before they can be used to index buckets. However, it takes the Intel Ivy Bridge and Power9 about 300 CPU cycles to reverse a 32-bit integer because there is no native instruction for this. In contrast, a bit-reversal operation can be completed on ARMv6 and above within a few CPU cycles [50].

For these reasons, we evaluated DHASH and the alternatives on three different servers. Table 1 lists the key characteristics of these servers. Because DHASH consistently outperforms the alternatives on all three servers, for brevity, we only present the results on Intel Ivy Bridge in this paper.

Each plotted data point, unless specified otherwise, represents the mean value of 20 trials. For performance tests, the standard deviation is denoted by vertical bars, which, however, may be too small to be visible in some figures.

## 6.2 Effects of Rebuilding

We first evaluate the effects of the rebuild operations. To that end, in this experiment, our benchmarking framework only measures the response time of lookup operations on the target buckets. The statistics of insert and delete operations and the statistics of lookup operations on the non-target buckets will be discussed in Section 6.3. The response time is measured by the worker threads as the time between an operation is invoked and that the hash table returns a result.

In this experiment, we evaluate the hash tables under a *typical hash-table usage pattern* that consists of eight worker threads, with the average load factor  $\alpha$  being 16, and with the operation distribution  $m$  being 5, 5, and 90 for insert, delete, and lookup operations, respectively. The initial and the maximum size of the bucket array are set to 1K and 64K, respectively. Once the length of any bucket becomes larger than 32, two times of  $\alpha$ , a rebuild request is sent to the rebuild daemon, which rebuilds the hash table by creating a single rebuilding thread.

To better evaluate the effectiveness of the rebuild operations, we further measure HT-DHash-1f under the same usage pattern but do not run the attacker thread; that is, we measure how concurrent operations of dynamic hash tables behave, when the length of the list of each bucket is roughly

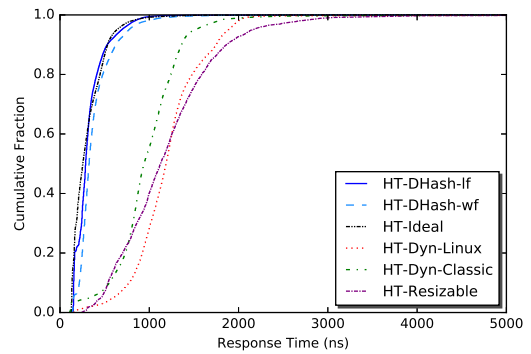


Fig. 3: Response time CDF of lookup operations when rebuilding.

TABLE 2: Response time statistics of lookup operations when rebuilding (nanoseconds).

Hash Table	Median	95%-percentile	99%-percentile	99.9%-percentile
HT-Dyn-Linux	1,211	1,888	2,060	5,492
HT-Dyn-Classic	1,076	1,512	2,038	5,393
HT-Resizable	1,196	2,152	3,003	7,665
HT-Ideal	256	636	944	1,571
HT-DHash-1f	288	692	920	1,546
HT-DHash-wf	328	813	1,361	2,012

equal to the average load factor. This evaluation is the ideal case for a hash table, which we refer to as HT-Ideal.

Figure 3 shows the Cumulative Distribution Function (CDF) of response time, and Table 2 lists the corresponding statistics. We make the following observations. First, once the pre-defined hash function cannot evenly distribute incoming data, DHASH is much more effective in cutting the length of the lists of the target buckets, such that concurrent lookup operations that need to traverse these lists do not experience increased delays. Figure 3 shows that both versions of DHASH outperform the alternatives in lookup operations' response time. In particular, Table 2 shows that the median response time of HT-DHash-1f and HT-DHash-wf is 288 and 328 nanoseconds, respectively. In contrast, the alternatives are more than one thousand nanoseconds, about 3 times slower than DHASH. Moreover, the 99.9%-percentile response time of DHASH is under 2 microseconds, but the alternatives reach 5 microseconds, introducing noticeable delays to legitimate users. The reasons for this are as follows. (1) Our distributing mechanism allows DHASH to use lock-free and wait-free lists as the implementation of the buckets. Therefore, the rebuild operation can run in parallel with concurrent common operations. In contrast, HT-Dyn-Classic and HT-Dyn-Linux employ a locking mechanism, such that the locks are contended, and the rebuild operation could be suspended, leading to a much longer time period during which concurrent lookup operations must traverse the lists of the target buckets. (2) DHASH is modular, and both HT-DHash-1f and HT-DHash-wf employ ordered linked lists. In contrast, HT-Dyn-Linux can only use unordered linked lists, which causes concurrent lookup operations to traverse many more nodes to get the results.

The second observation is that the performance of

HT-Resizable deteriorates sharply when its pre-defined hash function cannot evenly distribute incoming data. Figure 3 shows that HT-Resizable has the worst tail latencies compared with other hash tables. The reasons for this are as follows. (1) Upon receiving a rebuild request, all that HT-Resizable can do is double the size of its bucket array. However, this typically can only cut the length of the lists of the target buckets by half. (2) As we have discussed in Section 2, the original version of HT-Resizable ignores the basic statistical information, which is required when HT-Resizable is used in real-world applications. In particular, it is very costly to retrieve the information of list length. A practical implementation of HT-Resizable with this information becomes less efficient than the original version when resizing. Moreover, it is interesting that HT-DHash-1f even outperforms HT-Ideal in some cases. Table 2 shows that the 99.9%-percentile response time of HT-DHash-1f is 1,546 nanoseconds, smaller than HT-Ideal’s 1,571 nanoseconds.

### 6.3 DHASH performance

After evaluating the effectiveness of DHASH’s rebuild operation under a typical hash-table usage pattern, we now evaluate the performance of DHASH in a broad range of use cases.

We run each test for 3 minutes, which is long enough for our benchmarking framework to retrieve the characteristics of the evaluated hash tables but not to generate too many response-time log messages to analyze. In the experiments, HT-Resizable quickly exhausts its bucket array (64K in maximum) and requests more, and hence the benchmarking framework kills it in a few seconds. Combined with the fact that HT-Resizable cannot effectively recover from attacks (discussed in Section 6.2), we conclude that HT-Resizable is not the algorithm of choice for applications in the face of attacks and thus will omit HT-Resizable in the following experiments. The hash table configuration is the same as in Section 6.2, but we vary (1) the number of concurrent worker threads, (2) operation mixes, or (3) average load factors. We make the following observations.

**Throughput:** First, DHASH outperforms the alternatives at different concurrency levels. Figure 4a shows that all hash tables have similar throughput results when running with a single worker thread. As the number of worker threads increases, HT-DHash-1f outperforms the alternatives. The main reason for this is that to synchronize the rebuild operation and concurrent insert and delete operations, HT-Dyn-Linux and HT-Dyn-Classic employ per-bucket locks. Even though HT-DHash-wf uses an optimistic, fine-grained locking mechanism, it still suffers from the contention on the target buckets. When the server becomes overloaded (the number of worker threads exceeds that of the CPU cores), the throughput of HT-Dyn-Linux and HT-Dyn-Classic decreases sharply due to the increased contention on locks. HT-DHash-1f, in contrast, is much more robust, and there is no noticeable decrease in its throughput.

Second, Figure 4b shows that as the percentages of insert and delete operations increase, the throughput of all hash tables decreases. However, DHASH still outperforms the alternatives, except with the extreme operation mix consisting of 40% of insert, 40% of delete, and 20% of lookup

operations. The reason for this is that even for DHASH, insert and delete operations are much more costly compared with lookup operations because they must update shared variables and issue memory fence instructions.

Third, DHASH outperforms the alternatives with different average load factors. The main reason for this is that HT-DHash-1f employs a lock-free, ordered linked list.

**Response time:** When evaluating the throughput, we also measured the corresponding 99.9%-percentile response time of common operations. For a fair comparison, we separate the experimental results of lookup operations from that of insert and delete operations, and present them in Figures 5 and 6, respectively.

Recall that (1) DHASH synchronizes concurrent insert and delete operations by using lock-free or wait-free algorithms rather than locks, (2) its rebuild operation always distributes the head nodes, and (3) it utilizes ordered linked lists. Therefore, the rebuild operation of DHASH is lightweight and does not noticeably affect the response time of concurrent common operations. Figure 5 shows the results, and we make the following observations.

The response time of DHASH’s lookup operations increases mainly when the server is overloaded. Figure 5a shows that the worst 99.9%-percentile response time of HT-DHash-1f is around 3 microseconds, when there are 48 concurrent worker threads. In contrast, HT-Dyn-Linux and HT-Dyn-Classic reach as high as 8 microseconds. Second, the response time of HT-DHash-1f’s lookup operations is irrelevant to the operation mixes. Figure 5b shows that 99.9% of HT-DHash-1f’s lookup operations can complete in 2 microseconds, but HT-Dyn-Linux takes more than 8.3 microseconds, when the percentage of lookup operations is reduced to 20%. Third, Figure 5c shows that as the average load factor increases, the 99.9%-percentile response time of DHASH’s lookup operations increases slowly. In contrast, the lookup operations of both HT-Dyn-Linux and HT-Dyn-Classic suffer unexpected delays, mainly because of the unordered linked lists used.

We then evaluate the response time of insert and delete operations. Figure 6 shows that DHASH (especially, HT-DHash-1f) outperforms the alternatives for most cases. The main reason for this is that except HT-DHash-1f, other hash tables use locks to synchronize concurrent insert and delete operations. Note that the response time of DHASH also slowly increases, because the insert and delete operations have to update shared variables, and the overhead of these two types of operations mainly come from cache misses and memory fence instructions that are hard to be optimized.

### 6.4 Rebuilding efficiency

We next measure the efficiency of the rebuild operations; that is, how fast can DHASH and the alternatives change their hash functions and rebuild the hash tables. We use the same hash table usage pattern as in Section 6.2, but we vary (1) the total number of nodes in the hash tables, which relates to the workload of rebuild operations, and (2) the number of concurrent worker threads, which causes the rebuild operations to run at different concurrency levels. Since HT-DHash-1f and HT-DHash-wf show the same



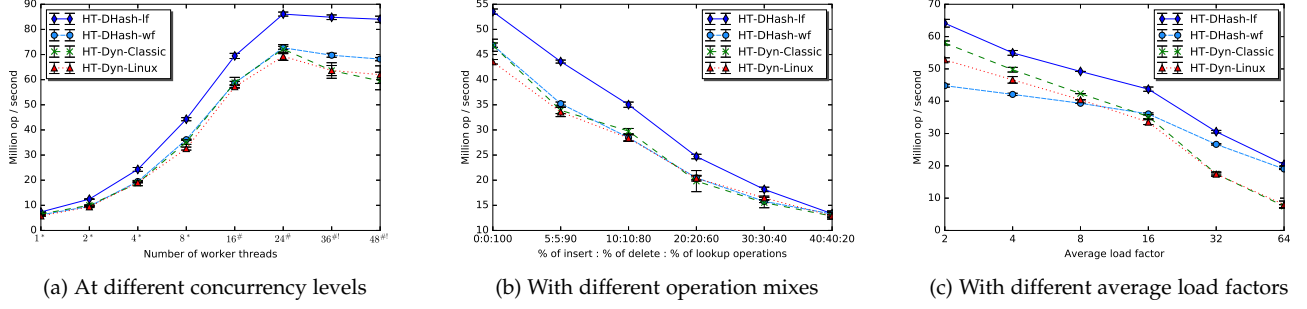


Fig. 4: Throughput (lookup, insert, and delete operations).

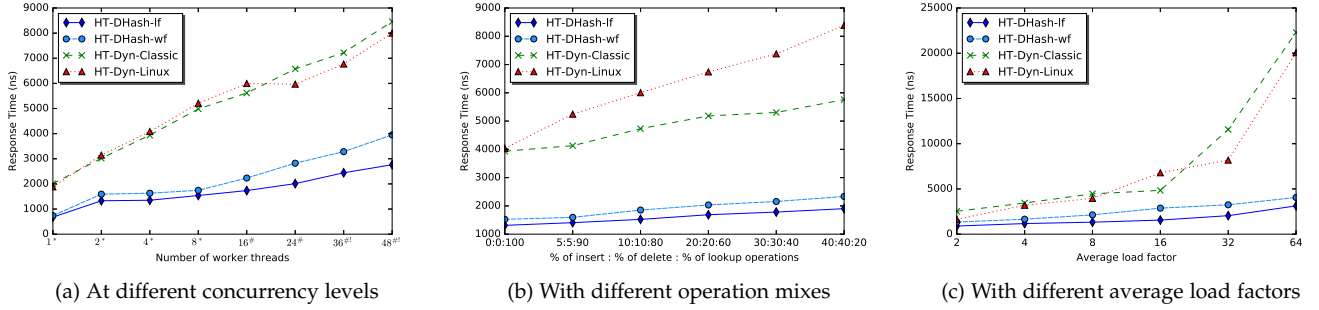


Fig. 5: 99.9%-percentile response time (lookup operations).

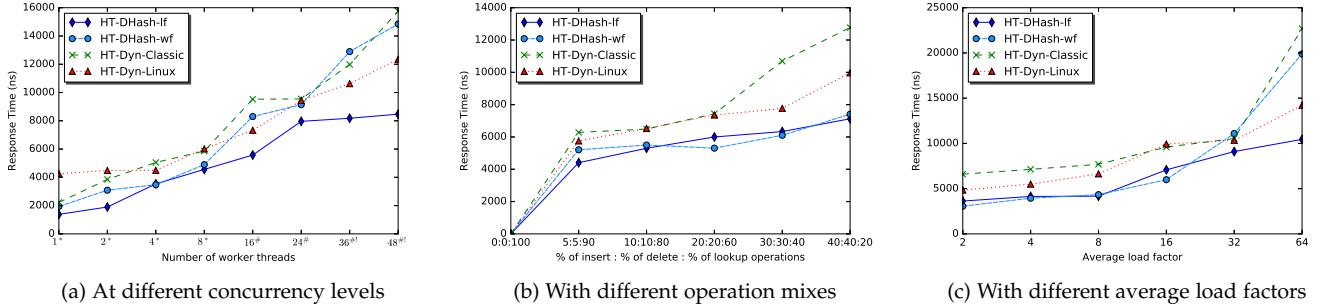


Fig. 6: 99.9%-percentile response time (insert and delete operations).

general trends, for clarity, we only present the results of HT-DHash-1f in this paper. To evaluate our parallelizing strategy, we measure HT-DHash-1f with different numbers of rebuilding threads, and the experimental results are marked with different suffixes. Figure 7 shows the results. We make the following observations.

linear speedup as a function of the number of rebuilding threads. For example, Figure 7a shows that it takes the sequential version (HT-DHash-1f-1) 40.3 ms to rebuild a hash table with 64K nodes, and that the time the 4-rebuilding-threads and the 8-rebuilding-threads versions need are 11.9 and 6.0 ms, achieving 3.4x and 6.7x performance speedups, respectively.

Second, for dynamic hash tables, which need to distribute nodes to the new hash tables, the time required to rebuild is almost linear to the number of nodes in the old hash tables. For example, Figure 7a shows that it takes HT-Dyn-Linux 13.7 and 27.1 ms to respectively rebuild hash tables containing 32K and 64K nodes. This implies that parallelizing rebuild operations is a must have for large dynamic hash tables.

Third, as the number of worker threads increases, the hash buckets are contended, resulting in performance degradation in the rebuild operations. This is the reason that, in general, the performance results shown in Figure 7b

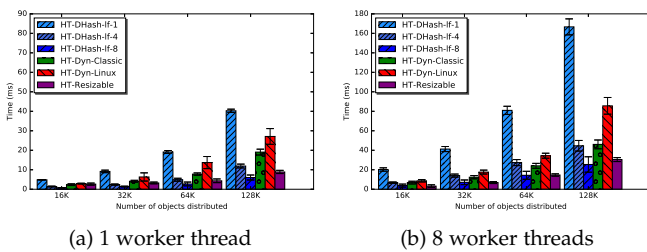


Fig. 7: Time required to rebuild hash tables.

First, our parallelizing strategy can achieve an almost



are inferior to those shown in Figure 7a. DHASH, however, by parallelizing its rebuild operation, can achieve a good speedup. For example, Figure 7b shows that the 8-rebuilding-threads version of DHASH (HT-DHash-1f-8) outperforms the alternatives in rebuilding efficiency. This demonstrates that our parallelization strategy is the right choice for accelerating the rebuild operations.

## 7 CONCLUSIONS

To recover from the emerging algorithmic complexity attacks, this paper presents DHASH, one type of hash table that can dynamically change its hash function and rebuild the hash table on the fly. The core of DHASH is a novel distributing mechanism that is atomic, non-blocking, and efficient in distributing every node from the old hash table to the new one. DHASH is modular and allows programmers to select from existing lock-free and wait-free set algorithms to create their own hash tables. We present the core techniques and demonstrate that DHASH is efficient and scalable. Experimental results indicate that DHASH is the algorithm of choice for operating systems and applications under service-level agreement (SLA) contracts.

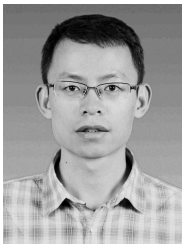
## ACKNOWLEDGMENTS

We thank Pedro Ramalhete and Paul E. McKenney, whose detailed suggestions greatly improved this paper. We thank OSU Open Source Lab and Packet Incorporation for allowing us to access the 16-core PowerPC and the 96-core ARM64 servers.

## REFERENCES

- [1] S. A. Crosby and D. S. Wallach, "Denial of service via algorithmic complexity attacks," in *USENIX Security Symposium*, 2003.
- [2] J. Edge, "Denial of service via hash collisions," <https://lwn.net/Articles/474912/>, 2012, [Online; accessed 9-Feb-2020].
- [3] P. Joye, "[PHP-DEV] 5.3.9 Release, Hash DoS," <https://lwn.net/Articles/474971/>, 2012, [Online; accessed 9-Feb-2020].
- [4] Microsoft, "Vulnerability in ASP.NET Could Allow Denial of Service," <https://docs.microsoft.com/en-us/security-updates/SecurityAdvisories/2011/2659883>, 2011, [Online; accessed 9-Sept-2021].
- [5] L. Carter and M. Wegman, "Universal classes of hash functions," *J. Comput. Syst. Sci.*, vol. 18, pp. 143–154, 1979.
- [6] Y. Orton, "Hardening Perl's Hash Function," <http://blog.booking.com/hardening-perls-hash-function.html>, 2013, [Online; accessed 9-Sept-2021].
- [7] R. J. Tobin and D. Malone, "Hash pile ups: Using collisions to identify unknown hash functions," *International Conference on Risks and Security of Internet and Systems*, pp. 1–6, 2012.
- [8] N. Bar-Yosef and A. Wool, "Remote algorithmic complexity attacks against randomized hash tables," in *SECURITY*, 2007.
- [9] C. Delimitrou and C. E. Kozyrakis, "Bolt: I know what you did last summer... in the cloud," *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*, 2017.
- [10] O. Shalev and N. Shavit, "Split-ordered lists: Lock-free extensible hash tables," *J. ACM*, vol. 53, pp. 379–405, 2006.
- [11] J. Triplett, P. E. McKenney, and J. Walpole, "Resizable, scalable, concurrent hash tables via relativistic programming," in *USENIX Annual Technical Conference*, 2011.
- [12] Y. Liu, K. Zhang, and M. Spear, "Dynamic-sized nonblocking hash tables," in *Proceedings of the 2014 ACM symposium on Principles of distributed computing*. ACM, 2014, pp. 242–251.
- [13] P. Fatourou, N. D. Kallimanis, and T. Ropars, "An efficient wait-free resizable hash table," in *SPAA*, 2018.
- [14] M. M. Michael, "High performance dynamic lock-free hash tables and list-based sets," in *SPAA*, 2002.
- [15] J. D. Valois, "Lock-free linked lists using compare-and-swap," in *PODC*, vol. 95. Citeseer, 1995, pp. 214–222.
- [16] M. Fomitchev and E. Ruppert, "Lock-free linked lists and skip lists," in *Proceedings of the twenty-third annual ACM symposium on Principles of distributed computing*. ACM, 2004, pp. 50–59.
- [17] S. Timnat, A. Braginsky, A. Kogan, and E. Petrank, "Wait-free linked-lists," in *PPoPP '12*, 2012.
- [18] H. XU, "Bridge: Add core igmp snooping support," Linux kernel git commit: eb1d16414339, February 2010.
- [19] T. Graf, "Resizable, scalable, concurrent hash table," Linux kernel git commit: 7e1e77636e36, August 2014.
- [20] K. Zhang, Y. Zhao, Y. Yang, Y. Liu, and M. F. Spear, "Practical non-blocking unordered lists," in *DISC*, 2013.
- [21] A. Kogan and E. Petrank, "A methodology for creating fast wait-free data structures," in *PPoPP '12*, 2012.
- [22] R. Fagin, J. Nievergelt, N. Pippenger, and H. R. Strong, "Extendible hashing? a fast access method for dynamic files," *ACM Trans. Database Syst.*, vol. 4, pp. 315–344, 1979.
- [23] P. Fatourou and N. D. Kallimanis, "Highly-efficient wait-free synchronization," *Theory of Computing Systems*, vol. 55, pp. 475–520, 2013.
- [24] H. Gao, J. F. Groote, and W. H. Hesselink, "Lock-free dynamic hash tables with open addressing," *Distributed Computing*, vol. 18, pp. 21–42, 2004.
- [25] T. Maier, P. Sanders, and R. Dementiev, "Concurrent hash tables: fast and general?(!)," *ACM Transactions on Parallel Computing (TOPC)*, vol. 5, pp. 1–32, 2019.
- [26] D. E. Knuth, *The art of computer programming*. Pearson Education, 1997, vol. 3.
- [27] M. M. Michael, "Hazard pointers: Safe memory reclamation for lock-free objects," *IEEE Transactions on Parallel and Distributed Systems*, vol. 15, no. 6, pp. 491–504, 2004.
- [28] J. Triplett and P. E. McKenney, "Scalable concurrent hash tables via relativistic programming," *Operating Systems Review*, vol. 44, pp. 102–109, 2010.
- [29] IBM, "K42 Operating System," <https://en.wikipedia.org/wiki/K42>, [Online; accessed 9-May-2019].
- [30] P. McKenney, "Introduction to RCU," <http://www.rdrop.com/~paulmck/RCU/>, 2020, [Online; accessed 9-May-2019].
- [31] M. Desnoyers, "Userspace RCU," <https://liburcu.org/>, 2012, [Online; accessed 9-May-2019].
- [32] P. E. McKenney et al., "Is parallel programming hard, and, if so, what can you do about it?" *Linux Technology Center, IBM Beaverton*, 2011.
- [33] B. Jenkins, "A hash function for hash table lookup," <http://www.burtleburtle.net/bob/hash/doobs.html>, 2013, [Online; accessed 9-Feb-2020].
- [34] N. Hua, E. Norige, S. S. Kumar, and B. Lynch, "Non-crypto hardware hash functions for high performance networking asics," *2011 ACM/IEEE Seventh Symposium on Architectures for Networking and Communications Systems*, pp. 156–166, 2011.
- [35] T. Cormen, C. Leiserson, R. Rivest, and C. Stein, "Introduction to algorithms, third edition," 2009.
- [36] K. Suzuki, D. Tonien, K. Kurosawa, and K. Toyota, "Birthday paradox for multi-collisions," in *ICISC*, 2006.
- [37] M. L. Fredman, J. H. Komlos, and E. Szemerédi, "Storing a sparse table with  $o(1)$  worst case access time," *23rd Annual Symposium on Foundations of Computer Science (sfcs 1982)*, pp. 165–169, 1982.
- [38] S. Heller, M. Herlihy, V. Luchangco, M. Moir, W. N. Scherer, and N. Shavit, "A lazy concurrent list-based set algorithm," in *OPODIS*, 2005.
- [39] C. Spec., "Memory order specification," [https://en.cppreference.com/w/cpp/atomic/memory\\_order](https://en.cppreference.com/w/cpp/atomic/memory_order), 2011, [Online; accessed 9-Feb-2020].
- [40] T. Harris, "A pragmatic implementation of non-blocking linked-lists," in *DISC*, 2001.
- [41] T. Harris, K. Fraser, and I. Pratt, "A practical multi-word compare-and-swap operation," in *DISC*, 2002.
- [42] M. Arbel-Raviv and T. Brown, "Harnessing epoch-based reclamation for efficient range queries," *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 2018.
- [43] M. Arbel-Raviv and T. Brown, "Reuse, don't recycle: Transforming lock-free algorithms that throw away descriptors," *ArXiv*, vol. abs/1708.01797, 2017.

- [44] J. Hennessy and D. Patterson, "Computer architecture - a quantitative approach, 5th edition," 1996.
- [45] F. Ellen, P. Fatourou, E. Ruppert, and F. Breugel, "Non-blocking binary search trees," *Proceedings of the 29th ACM SIGACT-SIGOPS symposium on Principles of distributed computing*, 2010.
- [46] P. Fatourou and E. Ruppert, "Persistent non-blocking binary search trees supporting wait-free range queries," *The 31st ACM Symposium on Parallelism in Algorithms and Architectures*, 2019.
- [47] M. P. Herlihy and J. M. Wing, "Linearizability: A correctness condition for concurrent objects," *ACM Transactions on Programming Languages and Systems (TOPLAS)*, vol. 12, no. 3, pp. 463–492, 1990.
- [48] P. Sewell, S. Sarkar, S. Owens, F. Z. Nardelli, and M. O. Myreen, "x86-tso: a rigorous and usable programmer's model for x86 multiprocessors," *Commun. ACM*, vol. 53, pp. 89–97, 2010.
- [49] M. Herlihy and N. Shavit, *The Art of Multiprocessor Programming*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2008.
- [50] ARM, "ARM Compiler armasm User Guide," <https://developer.arm.com/documentation/dui0473/m/arm-and-thumb-instructions/rbit>, 2019, [Online; accessed 9-May-2020].



**Junchang Wang** is currently an assistant professor at Nanjing University of Posts and Telecommunications. He received the Ph.D. degree in computer science from University of Science and Technology of China in 2014. From 2012 to 2013, he was a visiting student with the Department of Computer Science, Yale University. From 2014 to 2016, he was a research director in Jiangsu Future Networks Innovation Institute, China. From 2019 to 2020, he was a visiting scholar with the Department of Computer

Science, the Chinese University of Hong Kong. His research interest focuses on computer networks, and parallel and distributed computing.



**Dunwei Liu** received the BE degree from the Department of Electronic Information Engineering, Nanjing University of Posts and Telecommunications, in 2020. He is currently working toward the master's degree in the Department of Computer Science and Technology, Nanjing University of Posts and Telecommunications, China. His research interest includes parallel and distributed computing.



**Xiong Fu** received his BS and PhD in computer science from the University of Science and Technology, Hefei, in 2002 and 2007, respectively. He is currently a professor in School of Computer Science at the Nanjing University of Posts and Telecommunications. His research interest focuses on parallel and distributed computing, cloud computing.



**Fu Xiao** received the Ph.D. degree in computer science and technology from the Nanjing University of Science and Technology, Nanjing, China, in 2007. He is currently a Professor and Ph.D. supervisor with the School of Computer Science, Nanjing University of Posts and Telecommunications, Nanjing, China. His research papers have been published in many prestigious conferences and journals such as IEEE INFOCOM, IEEE ICC, IEEE IPCCC, IEEE/ACM ToN, IEEE JSAC, IEEE TMC, ACM TECS, IEEE TVT and so on.

His research interests are mainly in the areas of Internet of things and mobile computing.



**Chen Tian** is an associate professor at State Key Laboratory for Novel Software Technology, Nanjing University, China. He was previously an associate professor at School of Electronics Information and Communications, Huazhong University of Science and Technology, China. Dr. Tian received the BS (2000), MS (2003) and PhD (2008) degrees at Department of Electronics and Information Engineering from Huazhong University of Science and Technology, China. From 2012 to 2013, he was a postdoctoral re-

searcher with the Department of Computer Science, Yale University. His research interests include data center networks, network function virtualization, distributed systems, Internet streaming and urban computing.