

A Simple Abstraction for Complex Concurrent Indexes

Extended Version

Pedro da Rocha Pinto

Imperial College London
pmd09@doc.ic.ac.uk

Thomas Dinsdale-Young

Imperial College London
td202@doc.ic.ac.uk

Mike Dodds

University of Cambridge
mike.dodds@cl.cam.ac.uk

Philippa Gardner

Imperial College London
pg@doc.ic.ac.uk

Mark Wheelhouse

Imperial College London
mjw03@doc.ic.ac.uk

Abstract

Indexes are ubiquitous. Examples include associative arrays, dictionaries, maps and hashes used in applications such as databases, file systems and dynamic languages. Abstractly, a sequential index can be viewed as a partial function from keys to values. Values can be queried by their keys, and the index can be mutated by adding or removing mappings. Whilst appealingly simple, this abstract specification is insufficient for reasoning about indexes that are accessed concurrently.

We present an abstract specification for concurrent indexes. We verify several representative concurrent client applications using our specification, demonstrating that clients can reason abstractly without having to consider specific underlying implementations. Our specification would, however, mean nothing if it were not satisfied by standard implementations of concurrent indexes. We verify that our specification is satisfied by algorithms based on linked lists, hash tables and B^{Link} trees. The complexity of these algorithms, in particular the B^{Link} tree algorithm, can be completely hidden from the client's view by our abstract specification.

General Terms Algorithms, Concurrency, Theory, Verification.

Keywords B-Trees, Concurrent Abstract Predicates, Separation Logic.

1. Introduction

An *index* is a data structure where data is associated with identifying *keys*, through which the data can be efficiently retrieved. Indexes are ubiquitous in computer systems: they are integral to databases, caches, file systems, and even the objects of dynamic languages such as JavaScript. Concurrent systems use indexes for: *database sanitation* – to concurrently remove patients who have been cured or transferred; *graphics rendering* – to clip all objects outside horizontal or

vertical bounds; *garbage collection* – to concurrently mark reachable objects; and *web applications* – to allow multiple clients to add and remove pictures and comments, for instance. A variety of implementations of indexes exist, such as skip lists, hash tables and B-trees. Different implementations offer different performance characteristics, but all exhibit the same abstract behaviour.

To a sequential client, an index can be viewed abstractly as a partial function from keys to values. A client can query or mutate the index without having to take into account the complexities of its underlying implementation. This simple, yet powerful, abstract specification largely accounts for the popularity of indexes. However, this abstraction breaks down if an index is accessed concurrently. When several threads insert, remove and query keys, clients can no longer model the whole index by a single partial function. Each client must take account of potential interference from other threads.

In this paper, we present a novel abstract specification for *concurrent indexes*, and use it to verify a number of client programs. Crucially, clients can reason abstractly using our specification without having to consider specific underlying implementations. However, we can also verify our specification against complex concurrent index implementations.

Our approach is based on concurrent abstract predicates [6], recently introduced to reason about concurrent modules. With this technology, we can view the index as *divisible*: keys are a resource which can be divided between the threads. When threads operate on disjoint keys, they can do so independently of each other. When threads operate on shared keys, concurrent abstract predicates can account for the interference caused by other threads.

Intuitive description of the approach. First, consider the *disjoint* case, where each key is manipulated by a single thread. In this case, we can verify each thread in terms of the keys it uses, and combine the results to understand the composed system. In our specification, we have the predi-

cates $\text{in}(h, k, v)$ and $\text{out}(h, k)$: $\text{in}(h, k, v)$ declares that the key k is mapped to value v in index h ; $\text{out}(h, k)$ declares that there is no mapping of k . A thread must hold one of these predicates in order to modify k . A disjointness axiom enforces that only one thread can hold such a predicate on k at any one time. We describe these predicates as *abstract*, because they do not reveal how they are implemented.

Given these predicates, we can give the following specification to `remove`:

$$\{\text{in}(h, k, v)\} \quad \text{remove}(h, k) \quad \{\text{out}(h, k)\}$$

With this specification, we can prove the following property of a simple client program performing parallel removes (our proof assumes that $k_1 \neq k_2$):

$$\begin{array}{c} \{\text{in}(h, k_1, v_1) * \text{in}(h, k_2, v_2)\} \\ \{\text{in}(h, k_1, v_1)\} \quad \parallel \quad \{\text{in}(h, k_2, v_2)\} \\ \text{remove}(h, k_1) \quad \parallel \quad \text{remove}(h, k_2) \\ \{\text{out}(h, k_1)\} \quad \parallel \quad \{\text{out}(h, k_2)\} \\ \{\text{out}(h, k_1) * \text{out}(h, k_2)\} \end{array}$$

In this proof, we reason about the parallel threads individually. We then join the disjoint pre- and postconditions to form the overall proof. Disjointness is expressed by the separating conjunction, $*$, of concurrent separation logic [15]. The disjointness axiom requires that $k_1 \neq k_2$.

Now consider the shared case, where threads can interfere with each other: for example, when $k_1 = k_2$ in the parallel removes. We introduce the more refined predicates $\text{in}_{\text{def}}(h, k, v)_i$, $\text{out}_{\text{def}}(h, k)_i$, $\text{in}_{\text{rem}}(h, k, v)_i$ and $\text{out}_{\text{rem}}(h, k)_i$. These predicates are extended in two ways:

1. `def` and `rem` are *restrictions* on the type of interference that is allowed on the key: `def` prohibits any interference, while `rem` only permits removal of the key. All threads must agree on the type of interference for a given key.
2. The interference *permissions* $i \in (0, 1]$ determine whether a thread has shared ($0 < i < 1$) or exclusive ($i = 1$) access to a key. If a thread holds shared permission, it can only perform operations that respect the interference restrictions.

Using the `rem` predicates, we can give the following specification for `remove`:

$$\{\text{in}_{\text{rem}}(h, k, v)_i\} \quad \text{remove}(h, k) \quad \{\text{out}_{\text{rem}}(h, k)_i\}$$

Predicates can be split and joined by permission, so for example we have the axiom:

$$\text{in}_{\text{rem}}(h, k, v)_{i+j} \iff \text{in}_{\text{rem}}(h, k, v)_i * \text{in}_{\text{rem}}(h, k, v)_j,$$

where the sum of permissions held by all threads cannot exceed 1. In addition, if the current thread holds exclusive permission, we have axioms to change the type of the interference restriction without violating the expectations of other threads, such as:

$$\text{in}_{\text{def}}(h, k, v)_1 \iff \text{in}_{\text{rem}}(h, k, v)_1.$$

Using our specification and these axioms, we can prove a natural specification for parallel remove on a shared key:

$$\begin{array}{c} \{\text{in}_{\text{def}}(h, k, v)_1\} \\ \{\text{in}_{\text{rem}}(h, k, v)_1\} \\ \{\text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}} * \text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}}\} \\ \{\text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}}\} \quad \parallel \quad \{\text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}}\} \\ \text{remove}(h, k) \quad \parallel \quad \text{remove}(h, k) \\ \{\text{out}_{\text{rem}}(h, k)_{\frac{1}{2}}\} \quad \parallel \quad \{\text{out}_{\text{rem}}(h, k)_{\frac{1}{2}}\} \\ \{\text{out}_{\text{rem}}(h, k)_{\frac{1}{2}} * \text{out}_{\text{rem}}(h, k)_{\frac{1}{2}}\} \\ \{\text{out}_{\text{def}}(h, k)_1\} \end{array}$$

This specifies the strong property that, if we definitely know that key k has a value then, after the parallel remove, we definitely know that the value has been removed. We do not know *which* thread has performed the remove, but this fact is irrelevant to correctness.

Verifying clients and index implementations. Our concurrent index specification allows us to present a single abstract interface to clients, irrespective of the choice of underlying implementation. We demonstrate that our specification is useful by verifying several representative client programs such as function memoization, a prime number sieve and a mapping of a function onto an index.

We also verify that several concurrent index algorithms satisfy our specification: in particular, a naïve linked list algorithm with coarse-grained locking for expository purposes; a simple algorithm using a hash table linked to a set of (abstract) secondary indexes, to demonstrate the verification of a more complex implementation; and Sagiv’s substantial B^{Link} tree algorithm [19] to demonstrate the scalability of our techniques to a real-world algorithm. During verification, we found a subtle bug in the B^{Link} tree algorithm.

We use the *concurrent abstract predicate* methodology [6] to hide low-level sharing in the implementations from clients. In particular, the underlying sharing mechanism used by the B^{Link} tree algorithm to permit non-blocking reads is exceedingly complex. This complexity is completely hidden from the client’s view by our abstract specification.

Related work. We build directly on concurrent abstract predicates (CAP) [6], which provides a logic for verifying concurrent modules based on separation logic. CAP developed from three lines of work: racy concurrent variants of separation logic such as RGSep [8, 9, 22]; sequential modular reasoning based on abstract predicates [16]; and fine-grained modular reasoning based on context logic [3, 7]. We originally used RGSep [22] to verify concurrent B-trees [4]. However, RGSep and similar approaches depend on global conditions; consequently, they cannot verify abstract specifications such as our index specification. This observation formed part of our original motivation for CAP.

Our concurrent index specification descends from the set specification verified in [6]. In that paper, we focussed on building a sound logic, and verified only simple, disjoint

specifications against small implementations. As far as we are aware, our specification is the first in separation logic to allow thread-local reasoning combined with races over elements of a shared structure. We have verified our index specification against Sagiv’s real-world concurrent B^{Link} tree algorithm [19]¹, a substantial jump in the complexity of the verification compared with [6]. Our work is beginning to develop the idioms necessary to scale to large examples.

Others have worked on reasoning abstractly about index-like data structures for sequential clients. For example, Dillig *et al.* propose a static analysis for C-like programs which represents the abstract content of containers [5]. Kuncak *et al.* propose an analysis that represents various kinds of data by abstract sets, while proving these abstractions [13].

One of the most challenging parts of our work was verifying that the concurrent B^{Link} tree implementation satisfies our specification. Some prior work exists on verifying *sequential* B-trees. In [20], B-tree search and insert operations are verified as fault-free in a simplified sequential setting. In [14], a sequential B-tree implementation is verified in Coq as part of a relational database management system. The authors comment that the proof was difficult and in need of abstraction. They go on to state that ‘*verifying the correctness of high-performance, concurrent B+ trees will be a particularly challenging problem*’.

The only prior verification of a *concurrent* B-tree we are aware of is a highly-abstracted version of the algorithm modelled in process algebra [17]. It verifies a global specification, rather than allowing elements to be divided between threads. We believe that our work provides the first direct, formal verification of Sagiv’s widely-used algorithm [19].

Paper structure. §2 gives technical background. §3 give the disjoint index specification, and §4 extends it to sharing. §5 discusses iteration over indexes. §6 describes verifying our specification against index implementations. §3-5 can be understood from the simple summary in §2. A complete understanding of the technicalities in §6 requires knowledge of the original CAP paper [6].

2. Separation Logic & Abstraction

This paper is based on separation logic [18], a Hoare-style program logic for reasoning *locally* about programs that manipulate resource: for example, C programs that manipulate the heap. Local reasoning focusses on the specific part of the resource that is relevant at each point in the program. This supports scalable and compositional reasoning, since disjoint resource neither impinges upon nor is affected by the behaviour of the program at that point.

Separation logic specifications have a fault-avoiding partial-correctness interpretation. Consider the following specification for a command \mathbb{C} (here P, Q are assertions):

$$\{P\} \mathbb{C} \{Q\}$$

¹Without compression, which is beyond the scope of this paper.

The interpretation of this specification is that (1) executing \mathbb{C} in a state satisfying assertion P will result in a state satisfying assertion Q , if the command terminates; and (2) the resources represented by P are the only resources needed for \mathbb{C} to execute successfully.

Other resources can be conjoined with such a specification without affecting its validity. This is expressed by the following proof rule:

$$\text{FRAME} \quad \frac{\{P\} \mathbb{C} \{Q\}}{\{P * F\} \mathbb{C} \{Q * F\}} \quad \langle \text{side-condition} \rangle$$

This rule allows us to extend a specification on a small resource with an unmodified *frame assertion* F , giving a larger resource. Here, ‘*’ is the so-called *separating conjunction*. Combining two assertions P and F into a separating conjunction $P * F$ asserts that both resources are independent of each other. The side-condition simply states that no variable occurring free in F is modified by the program \mathbb{C} .

Separation logic provides straightforward reasoning about sequential programs. It also handles concurrency [15], using the following rule:

$$\text{PAR} \quad \frac{\{P_1\} \mathbb{C}_1 \{Q_1\} \quad \{P_2\} \mathbb{C}_2 \{Q_2\}}{\{P_1 * P_2\} \mathbb{C}_1 \parallel \mathbb{C}_2 \{Q_1 * Q_2\}}$$

In a concurrent setting, the precondition and postcondition are interpreted as resources owned exclusively by the thread. Reasoning using PAR is *thread-local*. We reason about each thread purely using the resources that are mentioned in its precondition, without requiring global reasoning about interleaving. As with sequential reasoning, locality is the key to compositional reasoning about threads.

Abstraction. Abstract specifications are a mechanism for specifying the external behaviour of a module’s functions, while hiding their implementation details from clients. Resources are represented by *abstract predicates* [16]. Clients do not need to know the concrete definitions of these predicates; they can reason purely in terms of the module’s operations. For example, `insert` in a set module might be specified as:

$$\{\text{set}(x, S)\} \quad \text{insert}(x, v) \quad \{\text{set}(x, S \cup \{v\})\}$$

`insert` updates the abstract contents of the set at address x from S to $S \cup \{v\}$. A client can reason about the high-level behaviour of `insert` without knowing about the concrete definition of the set predicate.

Abstract predicates, however, can only represent the set as a single entity, because implementation details disrupt finer-grained abstractions. *Concurrent abstract predicates* [6], on the other hand, can achieve finer abstractions. We can break the set down into predicates representing individual elements: $\text{in}(x, v)$ if v belongs to the set x ; $\text{out}(x, v)$ if it does not. Different threads can hold access to different set elements. When element v is not in the set, the command

insert can be specified by:

$$\{\text{out}(x, v)\} \quad \text{insert}(x, v) \quad \{\text{in}(x, v)\}$$

Concurrent abstract predicates provide a finer granularity of local reasoning, whilst still hiding implementation details from clients. We follow the concurrent abstract predicate approach in our reasoning about concurrent indexes.

3. Index Specification: Disjointness

We start by giving a simple specification which divides an index up into its constituent keys. Our specification ensures that each key is accessed by at most one thread (in §4 we discuss a refined specification that supports sharing). Our specification hides the fact that each key is part of an underlying shared data structure, allowing straightforward high-level reasoning about keys and values.

Abstractly, the state of an index can be seen as a partial function mapping keys to values²:

$$H : \text{Keys} \rightarrow \text{Vals}$$

There are three basic operations on an index – `search`, `insert` and `remove` – which operate on index h (with current state H) as follows:

- `search(h, k)` looks for the key k in the index. It returns $H(k)$ if it is defined, and `nil` otherwise.
- `insert(h, k, v)` tries to modify H to associate the key k with value v . If $k \in \text{dom}(H)$ then `insert` does nothing. Otherwise it modifies the shared index to $H \uplus \{k \mapsto v\}$.
- `remove(h, k)` tries to remove the value of the key k from the index. If $k \notin \text{dom}(H)$ then `remove` does nothing. Otherwise it rewrites the index to $H \setminus \{k\}$.

This view of operations on the index is appealingly simple, but cannot be used for practical concurrent reasoning. This is because it depends on *global* knowledge of the underlying index H . To reason in this way, a thread would require perfect knowledge of the behaviour of other threads.

To avoid this, we give a specification that breaks the index up by key value. Our specification allows threads to hold the exclusive ownership of an individual key. Each key in the index is represented by a predicate, either `in` or `out` depending on whether the key is associated with a value or not. The predicates have this intuitive interpretation:

$\text{in}(h, k, v)$: there is a mapping in the index h from k to v , and only the thread holding the predicate can modify k .

$\text{out}(h, k)$: there is no mapping in the index h from k , and only the thread holding the predicate can modify k .

² Where possible, we treat the key and value sets abstractly. Implementations require certain properties of these sets, however: all require keys to be comparable for equality, hash tables require the ability to compute hashes of keys, and B-trees require a linear ordering on keys.

These predicates combine knowledge about state – whether a key is in the index – with knowledge about ownership – whether the thread is allowed to alter that key. A thread holding the predicate for a given key knows the value of the key, and can be sure that no other thread will modify it. This entangling of state with ownership is essential to our approach: each predicate is invariant under the behaviour of other threads, meaning its implementation can be abstracted.

The index operations have the following specifications with respect to these predicates:

$$\begin{array}{ll} \{\text{in}(h, k, v)\} & r := \text{search}(h, k) \quad \{\text{in}(h, k, v) \wedge r = v\} \\ \{\text{out}(h, k)\} & r := \text{search}(h, k) \quad \{\text{out}(h, k) \wedge r = \text{nil}\} \\ \{\text{in}(h, k, v')\} & \text{insert}(h, k, v) \quad \{\text{in}(h, k, v')\} \\ \{\text{out}(h, k)\} & \text{insert}(h, k, v) \quad \{\text{in}(h, k, v)\} \\ \{\text{in}(h, k, v)\} & \text{remove}(h, k) \quad \{\text{out}(h, k)\} \\ \{\text{out}(h, k)\} & \text{remove}(h, k) \quad \{\text{out}(h, k)\} \end{array}$$

Predicates can be composed using the separating conjunction $*$, indicating that they hold independently of each other. Note that our specification allows us to reason about an index as a collection of disjoint, independent elements, despite the fact that indexes are generally implemented as a single shared data structure.

Each predicate represents exclusive ownership of a particular key. Our specification represents this fact by exposing the following axiom:

$$\left(\begin{array}{l} (\text{in}(h, k, v) \vee \text{out}(h, k)) * \\ (\text{in}(h, k, v') \vee \text{out}(h, k)) \end{array} \right) \implies \text{false}$$

Given the above specifications, we can reason locally about programs that use concurrent indexes. Consider for example the following simple program:

```
r := search(h, k2);
insert(h, k1, r) || remove(h, k2)
```

This program retrieves the value v associated with the key k_2 . It then concurrently associates v with the key k_1 and removes the key k_2 . When the program completes, k_1 will be associated with v , and k_2 will have been removed from the index. This specification can be expressed as:

$$\{\text{out}(h, k_1) * \text{in}(h, k_2, v)\} - \{\text{in}(h, k_1, v) * \text{out}(h, k_2)\}$$

We can prove this specification as follows:

$$\begin{array}{l} \{\text{out}(h, k_1) * \text{in}(h, k_2, v)\} \\ r := \text{search}(h, k_2); \\ \{\text{out}(h, k_1) * \text{in}(h, k_2, v) \wedge r = v\} \\ \{\text{out}(h, k_1) \wedge r = v\} \quad \parallel \quad \{\text{in}(h, k_2, v)\} \\ \text{insert}(h, k_1, r) \quad \parallel \quad \text{remove}(h, k_2) \\ \{\text{in}(h, k_1, v)\} \quad \parallel \quad \{\text{out}(h, k_2)\} \\ \{\text{in}(h, k_1, v) * \text{out}(h, k_2)\} \end{array}$$

In this proof, the search operation first uses the predicate $\text{in}(h, k_2, v)$ to retrieve the value v . Then, the parallel rule hands `insert` and `remove` the $\text{out}(h, k_1)$ and $\text{in}(h, k_2, v)$ predicates respectively. The postcondition of the program consists of the separating conjunction of the two thread postconditions.

3.1 Example: Map

A common operation on a concurrent index is applying a particular function to every value held in the index: *mapping* the function onto the index. We consider a simple algorithm `rangeMap` that maps function f (implemented by `f`) onto keys within a specified range. We implement `rangeMap` with a divide-and-conquer approach, which splits the key range into sub-intervals on which the map operation is recursively applied in parallel.

```

rangeMap(h, k1, k2) {
  if (k1 = k2) {
    r := search(h, k1);
    if (r ≠ nil) {
      remove(h, k1);
      r := f(r);
      insert(h, k1, r);
    }
  } else {
    rangeMap(h, k1, k1 + ((k2 - k1) / 2))
    || rangeMap(h, k1 + ((k2 - k1) / 2) + 1, k2)
  }
}

```

We specify `rangeMap` as follows, where S is a set of key-value pairs:

$$\left\{ \begin{array}{l} \bigotimes_{k_1 \leq i \leq k_2} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \\ \quad \left. (\exists v. \text{in}(h, i, v) \wedge (i, v) \in S) \right) \\ \text{rangeMap}(h, k_1, k_2) \\ \bigotimes_{k_1 \leq i \leq k_2} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \\ \quad \left. (\exists v. \text{in}(h, i, f(v)) \wedge (i, v) \in S) \right) \end{array} \right\}$$

(Here, \bigotimes is the iterated separating conjunction. That is, $\bigotimes_{x \in \{1,2,3\}} \cdot P$ is equivalent to $P[1/x] * P[2/x] * P[3/x]$. The set $\text{keys}(S)$ is the set of keys associated with values in S .)

In the specification, the logical variable S describes the initial state of the index (in the key range $[k_1, k_2]$). Assuming that S contains at most one key-value pair for each key, the key i (for $k_1 \leq i \leq k_2$) initially has value v if and only if $(i, v) \in S$. After execution of `rangeMap`, the postcondition ensures that if the key i had an initial value v , then it now has value $f(v)$, and if it had no value then it still has no value. A proof that `rangeMap` conforms to this specification is given in Figure 1.

`rangeMap` might not be considered truly typical of map operations, as it maps over a range of keys rather than the entire index. In §5, we introduce a specification for iterators, allowing all keys in an index to be enumerated. Using an iterator, we implement and verify a map function over all values in the index.

$$\left\{ \begin{array}{l} \bigotimes_{k_1 \leq i \leq k_2} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \\ \quad \left. (\exists v. \text{in}(h, i, v) \wedge (i, v) \in S) \right) \\ \text{rangeMap}(h, k_1, k_2) \{ \\ \quad \text{if } (k_1 = k_2) \{ \\ \quad \quad \left\{ k_1 = k_2 \wedge ((\text{out}(h, k_1) \wedge k_1 \notin \text{keys}(S)) \vee \right. \\ \quad \quad \left. (\exists v. \text{in}(h, k_1, v) \wedge (k_1, v) \in S)) \right\} \\ \quad \quad r := \text{search}(h, k_1); \\ \quad \quad \left\{ ((\text{out}(h, k_1) \wedge k_1 \notin \text{keys}(S) \wedge r = \text{nil}) \vee \right. \\ \quad \quad \left. (\text{in}(h, k_1, r) \wedge (k_1, r) \in S)) \wedge k_1 = k_2 \right\} \\ \quad \quad \text{if } (r \neq \text{nil}) \{ \\ \quad \quad \quad \left\{ \text{in}(h, k_1, r) \wedge (k_1, r) \in S \wedge k_1 = k_2 \right\} \\ \quad \quad \quad \text{remove}(h, k_1); \\ \quad \quad \quad \left\{ \text{out}(h, k_1) \wedge (k_1, r) \in S \wedge k_1 = k_2 \right\} \\ \quad \quad \quad r := f(r); \\ \quad \quad \quad \left\{ \exists v. \text{out}(h, k_1) \wedge (k_1, v) \in S \wedge r = f(v) \wedge k_1 = k_2 \right\} \\ \quad \quad \quad \left\{ \exists v. \text{in}(h, k_1, f(v)) \wedge (k_1, v) \in S \wedge k_1 = k_2 \right\} \\ \quad \quad \quad \} \\ \quad \quad \left\{ k_1 = k_2 \wedge ((\text{out}(h, k_1) \wedge k_1 \notin \text{keys}(S)) \vee \right. \\ \quad \quad \left. (\exists v. \text{in}(h, k_1, f(v)) \wedge (k_1, v) \in S)) \right\} \\ \quad \} \\ \} \text{ else } \{ \\ \quad \left\{ \bigotimes_{k_1 \leq i \leq \lfloor \frac{k_1+k_2}{2} \rfloor} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \right. \\ \quad \quad \left. \left. (\exists v. \text{in}(h, i, v) \wedge (i, v) \in S) \right) * \right\} \\ \quad \left\{ \bigotimes_{\lfloor \frac{k_1+k_2}{2} \rfloor < i \leq k_2} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \right. \\ \quad \quad \left. \left. (\exists v. \text{in}(h, i, v) \wedge (i, v) \in S) \right) \right\} \\ \quad // \text{Apply the PAR rule.} \\ \quad \text{rangeMap}(h, k_1, k_1 + ((k_2 - k_1) / 2)) \\ \quad || \text{rangeMap}(h, k_1 + ((k_2 - k_1) / 2) + 1, k_2) \\ \quad \left\{ \bigotimes_{k_1 \leq i \leq \lfloor \frac{k_1+k_2}{2} \rfloor} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \right. \\ \quad \quad \left. \left. (\exists v. \text{in}(h, i, f(v)) \wedge (i, v) \in S) \right) * \right\} \\ \quad \left\{ \bigotimes_{\lfloor \frac{k_1+k_2}{2} \rfloor < i \leq k_2} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \right. \\ \quad \quad \left. \left. (\exists v. \text{in}(h, i, f(v)) \wedge (i, v) \in S) \right) \right\} \\ \} \} \\ \left\{ \bigotimes_{k_1 \leq i \leq k_2} \cdot \left(\text{out}(h, i) \wedge i \notin \text{keys}(S) \vee \right. \right. \\ \quad \left. \left. (\exists v. \text{in}(h, i, f(v)) \wedge (i, v) \in S) \right) \right\} \end{array} \right\}$$

Figure 1. Proof for `rangeMap`.

4. Index Specification: Sharing

The specification we defined in the previous section requires that each key in the index is accessed by at most one thread. However, often threads read and write to keys at the same time. In this section, we define a refined specification that allows for concurrent access to keys. As before, our specification hides implementation details and allows threads to reason locally.

Consider the following program:

$$\text{remove}(h, k) \parallel r := \text{search}(h, k) \quad (1)$$

If we know at the start of the program that key k maps to some value v , we should be able to establish that there will not be a mapping from the key k at the end. However, we will not know the value of r , because we do not know at which point during the `remove` operation that the `search` operation will read the value associated with k .

Implementations have many different ways of handling the sharing of keys (for example using mutual exclusion locks or transactions), but at the abstract level they all behave in the same way. If a thread reads a key multiple times, the reads all return the same result, unless another thread also writes to that key.

Our refined specification is based on abstract predicates that express three facts about a given key:

1. whether there is a mapping from the key to some value in a set;
2. whether the thread holding the predicate can add or remove the value of the key in the index;
3. whether any other concurrently running threads (the *environment*) can add or remove the value of the key in the index.

These facts are related. If a key maps to a value in the index, but other threads are allowed to remove the value of the key, the current thread cannot assume the value will remain in the index. Our predicates therefore reflect the uncertainty generated by sharing in a local way.

We define the following set of predicates, parametric on key k and index h :

$\text{in}_{\text{def}}(h, k, v)_i$: there is a mapping from key k to value v and a thread can only modify this key if it has exclusive permission ($i = 1$).

$\text{out}_{\text{def}}(h, k)_i$: there is no mapping from key k and a thread can only modify this key if it has exclusive permission ($i = 1$).

$\text{in}_{\text{ins}}(h, k, S)_i$: there is a mapping from key k to a value in set S and threads can only insert values in set S at this key.

$\text{out}_{\text{ins}}(h, k, S)_i$: there may be a mapping from key k to a value in set S , threads can only insert values in set S at this key, and the current thread has not made such an insertion so far.

$\text{in}_{\text{rem}}(h, k, v)_i$: there may be a mapping from key k to value v , threads can only remove the value at this key, and the current thread has not done this so far.

$\text{out}_{\text{rem}}(h, k)_i$: there is no mapping from key k and threads can only remove the value at this key.

$\text{unk}(h, k, S)_i$: there may be a mapping from key k to a value v in set S and threads can search, remove and insert any value in set S at this key.

$\text{read}(h, k)$: there may be a mapping from key k to some value, the current thread may not change it, but other threads can make any modification.

The subscripts def , ins and rem and the fractional components $i \in (0, 1]$ record the behaviours allowed by the current thread and its environment on key k .

Access to keys can be shared between threads. We represent this in our specification by splitting predicates. Our specification includes axioms which define the ways that predicates can be split and joined. For example:

$$\text{in}_{\text{rem}}(h, k, v)_{i+j} \iff \text{in}_{\text{rem}}(h, k, v)_i * \text{in}_{\text{rem}}(h, k, v)_j \\ \text{if } i + j \leq 1$$

As in Boyland [2], fractional *permissions* are used to record splittings. A permission $i \in (0, 1)$ records that a key is shared with other threads, while $i = 1$ records it is held exclusively by the current thread.

When a thread holds exclusive access to a key ($i = 1$), the thread can add or remove the key freely. When a thread shares access to the key ($i \in (0, 1)$), the subscripts def , ins and rem *restrict* what the thread and its environment are able to do. Subscript def specifies that no thread is able to modify the key. Subscript ins specifies that both thread and environment can insert on the key, but not remove the key, while subscript rem specifies the converse.

Modifying keys concurrently can result in different threads holding different predicates for the same key. For example, suppose a thread holds the $\text{in}_{\text{rem}}(h, k, v)_1$ predicate, which denotes that the key k has value v in the index. Since the permission is 1, this knowledge is assured. However, we can split this predicate into two halves, $\text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}}$ and $\text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}}$, and give each half to two sub-threads. Assume the first thread does not modify the key, but the second calls $\text{remove}(h, k)$, which has the following specification:

$$\{\text{in}_{\text{rem}}(h, k, v)_i\} \quad \text{remove}(h, k) \quad \{\text{out}_{\text{rem}}(h, k)_i\}$$

The result is uncertainty: one thread holds the $\text{out}_{\text{rem}}(h, k)_{\frac{1}{2}}$ predicate, stating that k is not in the index, while the other holds the $\text{in}_{\text{rem}}(h, k, v)_{\frac{1}{2}}$ predicate, stating that k may have associated value v . We define joining axioms that resolve this uncertainty. Since rem allows removal but not insertion, we know that once the key has been removed from the index, it stays removed. So out_{rem} dominates in_{rem} , which is reflected in the following axiom:

$$\text{in}_{\text{rem}}(h, k, v)_i * \text{out}_{\text{rem}}(h, k)_j \implies \text{out}_{\text{rem}}(h, k)_{i+j} \\ \text{if } i + j \leq 1$$

Some predicates take sets of value arguments, while others take singleton values. We use singleton values when we know a key has that value. We use a set of values when concurrent inserts are possible (that is, in the ins and unk cases), because we cannot know which thread will be the first to insert. However, if a value is inserted, it will be one of the values in the set S .

Our full specification is given in Figure 3. The choice of predicates is not arbitrary; each represents a stable combination of facts about the key k and the behaviours permitted by the thread and environment. Figure 2 shows how various

		Thread		Env.	
Predicate	Perm.	Ins.	Rem.	Ins.	Rem.
$\text{in}_{\text{def}} / \text{out}_{\text{def}}$	1	Yes	Yes	No	No
$\text{in}_{\text{def}} / \text{out}_{\text{def}}$	i	No	No	No	No
$\text{in}_{\text{ins}} / \text{out}_{\text{ins}}$	1	Yes	No	No	No
$\text{in}_{\text{ins}} / \text{out}_{\text{ins}}$	i	Yes	No	Yes	No
$\text{in}_{\text{rem}} / \text{out}_{\text{rem}}$	1	No	Yes	No	No
$\text{in}_{\text{rem}} / \text{out}_{\text{rem}}$	i	No	Yes	No	Yes
unk	i	Yes	Yes	Yes	Yes
read	-	No	No	Yes	Yes

Figure 2. Predicates and their interference.

combinations of fractional permissions and subscripts correspond to various behaviours. Our predicates give almost complete coverage of all possible combinations. The missing combinations are either cases where the current thread has no access to a key, or where it is only safe to conclude that a key has an unknown value, in which case we can use one of the read or unk predicates. We do not claim that our specification is definitive, just one natural choice. We expect to adapt our specification when looking at real-world applications such as the POSIX file system, the concurrent database algorithm ARIES, and `java.util.concurrent`. We believe that our specification is robust enough to be able to support such applications with minor modification.

4.1 Proving Simple Examples

Recall the program labelled (1) with which we began this section. This program satisfies the following specifications:

$$\begin{aligned} \{ \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 \} &- \{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} \\ \{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} &- \{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} \end{aligned}$$

Using our abstract specifications, we can prove the first of these specifications as follows:

$$\begin{aligned} &\{ \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 \} \\ &\{ \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 * \text{read}(\mathbf{h}, \mathbf{k}) \} \\ \{ \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 \} &\parallel \{ \text{read}(\mathbf{h}, \mathbf{k}) \} \\ \text{remove}(\mathbf{h}, \mathbf{k}) &\parallel \text{r} := \text{search}(\mathbf{h}, \mathbf{k}) \\ \{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} &\parallel \{ \text{read}(\mathbf{h}, \mathbf{k}) \} \\ &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 * \text{read}(\mathbf{h}, \mathbf{k}) \} \\ &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} \end{aligned}$$

The proof starts with the predicate $\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1$, which specifies that there is a mapping from key \mathbf{k} to a value v in the index. The def subscript asserts that no other thread can modify the value mapped by this key. We use the following axiom to create a `read`(\mathbf{h}, \mathbf{k}) predicate:

$$X_i \iff X_i * \text{read}(\mathbf{h}, \mathbf{k})$$

This allows the right-hand thread to perform a simple search operation, although the postcondition establishes

nothing about the result. This captures the fact that we do not know at which point during the `remove` operation the `search` operation will read the key's value. The $\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1$ predicate allows the left-hand thread to remove the value successfully, as we know that it is the only thread changing the shared state for the key \mathbf{k} . When both threads finish their execution we use the same axiom to merge `read`(\mathbf{h}, \mathbf{k}) back into the $\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1$. We can prove the second specification in a similar fashion.

We can establish natural specifications for all the various combinations of `insert`, `remove` and `search`. For example, consider the parallel composition of two removes on the same key \mathbf{k} :

$$\text{remove}(\mathbf{h}, \mathbf{k}) \parallel \text{remove}(\mathbf{h}, \mathbf{k})$$

Regardless of whether \mathbf{k} is in the index, we definitely know that there will be no mapping from key \mathbf{k} afterwards. By splitting the predicates, we can share this knowledge between the threads.

$$\begin{aligned} &\{ \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 \} \\ &\{ \text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_1 \} \\ &\{ \text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_{\frac{1}{2}} * \text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_{\frac{1}{2}} \} \\ \{ \text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_{\frac{1}{2}} \} &\parallel \{ \text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_{\frac{1}{2}} \} \\ \text{remove}(\mathbf{h}, \mathbf{k}) &\parallel \text{remove}(\mathbf{h}, \mathbf{k}) \\ \{ \text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_{\frac{1}{2}} \} &\parallel \{ \text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_{\frac{1}{2}} \} \\ \{ \text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_{\frac{1}{2}} * \text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_{\frac{1}{2}} \} & \\ &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} \end{aligned}$$

We cannot always establish the exact state of an index at all points during a program, but our specification will always allow us to be as precise as possible. For example, consider the following program:

$$\text{remove}(\mathbf{h}, \mathbf{k}) \parallel \begin{array}{l} \text{insert}(\mathbf{h}, \mathbf{k}, v) \\ \text{remove}(\mathbf{h}, \mathbf{k}) \end{array}$$

When run in a state where key \mathbf{k} is initially unassigned, we will not know if there is a mapping from key \mathbf{k} in the index at the end of the parallel call. However, after the final `remove` operation we know that the key \mathbf{k} will be unassigned.

$$\begin{aligned} &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} \\ &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \vee \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 \} \\ &\{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_1 \} \\ &\{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} * \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} \} \\ \{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} \} &\parallel \{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} \} \\ \text{remove}(\mathbf{h}, \mathbf{k}) &\parallel \text{insert}(\mathbf{h}, \mathbf{k}, v) \\ \{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} \} &\parallel \{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} \} \\ \{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} * \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_{\frac{1}{2}} \} & \\ &\{ \text{unk}(\mathbf{h}, \mathbf{k}, \{v\})_1 \} \\ &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \vee \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1 \} \\ &\text{remove}(\mathbf{h}, \mathbf{k}) \\ &\{ \text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1 \} \end{aligned}$$

SPECIFICATIONS:

$\{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_i \wedge \mathbf{r} = v\}$
$\{\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_i \wedge \mathbf{r} = \text{nil}\}$
$\{\text{in}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{\text{in}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i \wedge \mathbf{r} \in S\}$
$\{\text{out}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{(\text{out}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i \wedge \mathbf{r} = \text{nil}) \vee (\text{in}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i \wedge \mathbf{r} \in S)\}$
$\{\text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{(\text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_i \wedge \mathbf{r} = v) \vee (\text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_i \wedge \mathbf{r} = \text{nil})\}$
$\{\text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{\text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_i \wedge \mathbf{r} = \text{nil}\}$
$\{\text{unk}(\mathbf{h}, \mathbf{k}, S)_i\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{\text{unk}(\mathbf{h}, \mathbf{k}, S)_i \wedge (\mathbf{r} \in S \vee \mathbf{r} = \text{nil})\}$
$\{\text{read}(\mathbf{h}, \mathbf{k})\}$	$\mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k})$	$\{\text{read}(\mathbf{h}, \mathbf{k})\}$
$\{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_i\}$	$\text{insert}(\mathbf{h}, \mathbf{k}, v')$	$\{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_i\}$
$\{\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1\}$	$\text{insert}(\mathbf{h}, \mathbf{k}, v)$	$\{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1\}$
$\{(\text{in}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i \vee \text{out}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i) \wedge v \in S\}$	$\text{insert}(\mathbf{h}, \mathbf{k}, v)$	$\{\text{in}_{\text{ins}}(\mathbf{h}, \mathbf{k}, S)_i\}$
$\{\text{unk}(\mathbf{h}, \mathbf{k}, S)_i \wedge v \in S\}$	$\text{insert}(\mathbf{h}, \mathbf{k}, v)$	$\{\text{unk}(\mathbf{h}, \mathbf{k}, S)_i\}$
$\{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, v)_1\}$	$\text{remove}(\mathbf{h}, \mathbf{k})$	$\{\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_1\}$
$\{\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_i\}$	$\text{remove}(\mathbf{h}, \mathbf{k})$	$\{\text{out}_{\text{def}}(\mathbf{h}, \mathbf{k})_i\}$
$\{\text{in}_{\text{rem}}(\mathbf{h}, \mathbf{k}, v)_i \vee \text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_i\}$	$\text{remove}(\mathbf{h}, \mathbf{k})$	$\{\text{out}_{\text{rem}}(\mathbf{h}, \mathbf{k})_i\}$
$\{\text{unk}(\mathbf{h}, \mathbf{k}, S)_i\}$	$\text{remove}(\mathbf{h}, \mathbf{k})$	$\{\text{unk}(\mathbf{h}, \mathbf{k}, S)_i\}$

AXIOMS:

$X_i * X_j$	$\Leftrightarrow X_{i+j}$	if $i + j \leq 1$
$\text{in}_{\text{ins}}(h, k, S)_i * \text{out}_{\text{ins}}(h, k, S)_j$	$\Rightarrow \text{in}_{\text{ins}}(h, k, S)_{i+j}$	if $i + j \leq 1$
$\text{in}_{\text{rem}}(h, k, v)_i * \text{out}_{\text{rem}}(h, k)_j$	$\Rightarrow \text{out}_{\text{rem}}(h, k)_{i+j}$	if $i + j \leq 1$
$\text{in}_{\text{def}}(h, k, v)_1$	$\Leftrightarrow \text{in}_{\text{rem}}(h, k, v)_1$	
$\exists v \in S. \text{in}_{\text{def}}(h, k, v)_1$	$\Leftrightarrow \text{in}_{\text{ins}}(h, k, S)_1$	
$\text{out}_{\text{def}}(h, k)_1$	$\Leftrightarrow \text{out}_{\text{rem}}(h, k)_1 \Leftrightarrow \text{out}_{\text{ins}}(h, k, S)_1$	
X_i	$\Leftrightarrow X_i * \text{read}(h, k)$	
$\text{read}(h, k)$	$\Leftrightarrow \text{read}(h, k) * \text{read}(h, k)$	
$\text{unk}(h, k, S)_1$	$\Leftrightarrow \text{out}_{\text{def}}(h, k)_1 \vee \exists v \in S. \text{in}_{\text{def}}(h, k, v)_1$	

CONTRADICTION AXIOMS:

$X_i * X_j$	$\Rightarrow \text{false}$	if $i + j > 1$
$\text{in}_{\text{def}}(h, k, v)_i * X_j$	$\Rightarrow \text{false}$	if $X \neq \text{in}_{\text{def}}(h, k, v)$
$\text{out}_{\text{def}}(h, k)_i * X_j$	$\Rightarrow \text{false}$	if $X \neq \text{out}_{\text{def}}(h, k)$
$(\text{in}_{\text{ins}}(h, k, S)_i \vee \text{out}_{\text{ins}}(h, k, S)_i) * X_j$	$\Rightarrow \text{false}$	if $X \neq \text{in}_{\text{ins}}(h, k, S) \wedge X \neq \text{out}_{\text{ins}}(h, k, S)$
$(\text{in}_{\text{rem}}(h, k, v)_i \vee \text{out}_{\text{rem}}(h, k)_i) * X_j$	$\Rightarrow \text{false}$	if $X \neq \text{in}_{\text{rem}}(h, k, v) \wedge X \neq \text{out}_{\text{rem}}(h, k)$
$(\text{in}_{\text{ins}}(h, k, S)_i * \text{in}_{\text{ins}}(h, k, S')_j) \vee (\text{out}_{\text{ins}}(h, k, S)_i * \text{out}_{\text{ins}}(h, k, S')_j)$	$\Rightarrow \text{false}$	if $S \neq S'$
$\text{unk}(h, k, S)_i * X_j$	$\Rightarrow \text{false}$	if $X \neq \text{unk}(h, k, S)$

Figure 3. Full specification for concurrent indexes. X denotes $\text{in}_{\text{def}}(h, k, v)$, $\text{out}_{\text{def}}(h, k)$, $\text{in}_{\text{ins}}(h, k, S)$, $\text{out}_{\text{ins}}(h, k, S)$, $\text{in}_{\text{rem}}(h, k, v)$, $\text{out}_{\text{rem}}(h, k)$ or $\text{unk}(h, k, S)$ in the axioms.


```

 $\{\exists i \in (0, 1]. \otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 
memoized_f(v) {
   $\{\otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 
  // frame the irrelevant values off
   $\{\text{unk}(\text{memo}, v, \{f(v)\})_i\}$ 
  r := search(memo, v);
   $\{\text{unk}(\text{memo}, v, \{f(v)\})_i \wedge (r = f(v) \vee r = \text{nil})\}$ 
  if (r = nil) {
     $\{\text{unk}(\text{memo}, v, \{f(v)\})_i\}$ 
    r := f(v);
     $\{\text{unk}(\text{memo}, v, \{f(v)\})_i \wedge r = f(v)\}$ 
    insert(memo, v, r);
     $\{\text{unk}(\text{memo}, v, \{f(v)\})_i \wedge r = f(v)\}$ 
  }
   $\{\text{unk}(\text{memo}, v, \{f(v)\})_i \wedge r = f(v)\}$ 
  // frame the values back on
   $\{r = f(v) \wedge \otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 
  return r;
}
 $\{\text{ret} = f(v) \wedge \exists i \in (0, 1]. \otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 

```

Figure 4. Proof outline for memoized_f.

The key step in this proof is the use of the final axiom from Figure 3 to convert a complete unk predicate into the disjunction of an in and out predicate. In both cases, the remove operation results in an index where the key k is definitely unassigned.

4.2 Example: Memoization

A common client application of indexes is memoization: storing the results of expensive computations to avoid having to recompute them. Our specification can verify that a memoized function gives the same result as the original function.

Suppose that f is a side-effect free procedure implementing the (mathematical) function f. A memoized version of f, memoized_f, can be implemented using the index memo as follows:

```

memoized_f(v) {
  r := search(memo, v);
  if (r = nil) {
    r := f(v);
    insert(memo, v, r);
  }
  return r;
}

```

We give memoized_f the following specification:

$\{\text{memo}\} r := \text{memoized_f}(v) \{r = f(v) \wedge \text{memo}\}$

where the abstract predicate memo is

$\text{memo} \triangleq \exists i \in (0, 1]. \otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i$

The definition of memo states that, for each value v', we do not know if v' is in the index. The predicate is splittable: that

```

 $\{\exists i \in (0, 1]. \otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 
evict_f() {
  while (...) {
     $\{\otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 
    k := nondet();
    // frame the irrelevant values off
     $\{\text{unk}(\text{memo}, k, \{f(k)\})_i\}$ 
    remove(memo, k);
     $\{\text{unk}(\text{memo}, k, \{f(k)\})_i\}$ 
  }
}
 $\{\exists i \in (0, 1]. \otimes_{v'} \text{unk}(\text{memo}, v', \{f(v')\})_i\}$ 

```

Figure 5. Proof outline for evict_f.

is, $\text{memo} \Leftrightarrow \text{memo} * \text{memo}$. The memoized_f specification therefore allows calls to f to be replaced with memoized_f, even in parallel. A proof of the specification for memoized_f is shown in Figure 4.

Evicting memoised values. We may want to periodically evict memoised values from the index, for example to ensure that the number of stored values does not grow indefinitely. Using our index specification, we can show that values can be evicted in parallel with memoized_f().

We model eviction by the function evict_f, which non-deterministically removes keys from the index:

```

evict_f() {
  while (...) {
    k := nondet();
    remove(memo, k);
  }
}

```

where nondet() returns an arbitrary key value and the Boolean assertion for while is not given. (A more nuanced eviction function might store timestamps along with the memoised values, and evict only old values. For simplicity, we choose not to model this.)

We give evict_f the following specification:

$\{\text{memo}\} \text{evict_f}() \{\text{memo}\}$

A proof of this specification is given in Figure 5. Because we can split and join memo arbitrarily, we can reason as follows:

$$\begin{array}{c}
\{\text{memo}\} \\
\{\text{memo} * \text{memo}\} \\
\{\text{memo}\} \quad \parallel \quad \{\text{memo}\} \\
\text{evict_f}() \quad \parallel \quad r := \text{memoized_f}(v) \\
\{\text{memo}\} \quad \parallel \quad \{\text{memo} \wedge r = f(v)\} \\
\{\text{memo} * (\text{memo} \wedge r = f(v))\} \\
\{\text{memo} \wedge r = f(v)\}
\end{array}$$

Consequently, it is safe to run the memoised version of f in parallel with eviction from the index.

```

sieve(max) {
  idx := idxrange(2, max);
  parwork(2, max, idx);
  return idx;
}

parwork(v, max, idx) {
  if (v ≤ sqrt(max)) {
    worker(v, max, idx)
    ||
    parwork(v+1, max, idx)
  }
}

worker(v, max, idx) {
  c := v + v;
  while (c ≤ max)
    remove(idx, c);
  c := c + v;
}

```

Figure 6. Prime sieve functions.

4.3 Example: The Sieve of Eratosthenes

Let us consider an example where many threads require write access to the same shared value in a concurrent index. We choose the Sieve of Eratosthenes [1, 12], an algorithm for generating all of the prime numbers up to a given maximum value \max . The sieve is a simple algorithm, but it is representative of a class of algorithms where threads cooperatively race to delete elements of shared data. Similar behaviour occurs in databases when deleting stale records, and in rendering when removing objects outside of a clipped region.

The algorithm starts by constructing a set of integers from 2 (since 1 is not a prime number) to \max . We use an index to represent the set of (candidate) prime numbers. A set can be viewed as an instance of an index where the set of values is a singleton (in this example, we use $\{0\}$). A key is either present, representing that it is in the set, or not: the value itself conveys no information. We assume a function `idxrange` that creates an index with mappings for keys in a specified range.

For each integer in the range $2 \dots \lfloor \sqrt{\max} \rfloor$, a thread is created that removes multiples of that integer from the set. Once all threads have completed, the remaining elements of the set are exactly those with no factors in the range $2 \dots \lfloor \sqrt{\max} \rfloor$ (excluding themselves), and hence exactly the prime numbers less than or equal to \max .

The code for the implementation is given in Figure 6. The procedure `sieve` is the main sieve function, which uses the recursive `parwork` procedure to run each worker thread in parallel. The procedure `worker` is the implementation of the worker threads.

The specification for `sieve` is

$$\begin{aligned}
& \{ \text{emp} \wedge \max > 1 \} \\
& x := \text{sieve}(\max) \\
& \left\{ \bigotimes_{i \in [2.. \max]}. \text{isPrime}(i) \Rightarrow \text{in}_{\text{def}}(x, i, 0)_1 \right. \\
& \quad \left. \wedge \neg \text{isPrime}(i) \Rightarrow \text{out}_{\text{def}}(x, i)_1 \right\}
\end{aligned}$$

where the predicate ‘`emp`’ denotes no resource at all, and the predicate ‘`isPrime(i)`’ holds exactly when i is prime. We also define the predicate ‘`fac(i, v, v')`’, which holds when i

has a factor (distinct from itself) in the range $[v \dots v']$:

$$\text{fac}(i, v, v') \triangleq \exists j. v \leq j \leq v' \wedge j \neq i \wedge (i \bmod j) = 0$$

The proof that `sieve` meets its specification is given in Figure 7. This proof requires we establish the following specification for `worker`:

$$\begin{aligned}
& \{ 2 \leq v \wedge \bigotimes_{i \in [2.. \max]}. \text{in}_{\text{rem}}(\text{idx}, i, 0)_t \} \\
& \text{worker}(v, \max, \text{idx}) \\
& \left\{ \bigotimes_{i \in [2.. \max]}. \text{fac}(i, v, v) \Rightarrow \text{out}_{\text{rem}}(\text{idx}, i)_t \wedge \right. \\
& \quad \left. \neg \text{fac}(i, v, v) \Rightarrow \text{in}_{\text{rem}}(\text{idx}, i, 0)_t \right\}
\end{aligned}$$

This specification expresses that the worker removes all multiples of v from the set; any other elements will still be present unless they are removed by another thread. The fact that (for $v \leq v'$)

$$\text{fac}(i, v, v) \vee \text{fac}(i, v+1, v') \iff \text{fac}(i, v, v')$$

allows us to conclude that the `parwork` procedure eliminates exactly the set elements with factors different from themselves in the range $v \dots \max$. Since $p > 1$ is prime if and only if it has no factor in the range $2 \dots \lfloor \sqrt{p} \rfloor$, for $i \in [2 \dots \max]$

$$\neg \text{fac}(i, 2, \lfloor \sqrt{\max} \rfloor) \iff \text{isPrime}(i).$$

Together with the index axioms that allow `rem` predicates to be switched to `def` predicates when full permission is held, this lets us establish the postcondition of `sieve`.

5. Iterating an Index

The high-level specification discussed so far does not allow us to explore the contents of an arbitrary index. To use `search`, we must know which keys we seek. If we do not (and the set of keys is infinite), we cannot write a program that examines all the values stored in the index. To handle this case, we add imperative iterators, based loosely on those in Java. Iterators have three operations:

- `it := createIter(h)` creates a new iterator for index h .
- `(k, v) := next(it)` returns some key-value pair in the index for which `it` is an iterator. The returned pair will be one that has not been returned by a previous call to `next` on `it`. When all key-value pairs have been returned, the call returns `(nil, nil)`.
- `destroyIter(it)` frees the iterator `it`.

To iterate an index, one creates a new iterator, calls `next` until it returns `(nil, nil)`, then frees the iterator. Notice that the `next` procedure just returns *some* key-value pair, placing no order on the iteration. This keeps the iterator specification general, as many underlying implementations have no natural ordering.

As in Java, we do not allow full mutability of an index being iterated. We allow partial mutability: keys can be safely modified once they have been returned by `next`.

```

{emp ∧ max > 1}
sieve(max) {
  idx := idxrange(2, max);
  {⊗i∈[2..max]. inrem(idx, i, 0)1}
  parwork(2, max, idx);
  {⊗i∈[2..max]. fac(i, 2, ⌊√max⌋) ⇒ outrem(idx, i)1 ∧
   ¬fac(i, 2, ⌊√max⌋) ⇒ inrem(idx, i, 0)1}
  // By properties of prime numbers and
  // index axioms
  {⊗i∈[2..max]. isPrime(i) ⇒ indef(idx, i, 0)1
   ∧ ¬isPrime(i) ⇒ outdef(idx, i)1}
  return idx;
}
{ret = idx ∧ ⊗i∈[2..max]. isPrime(i) ⇒ indef(idx, i, 0)1
 ∧ ¬isPrime(i) ⇒ outdef(idx, i)1}

{2 ≤ v ∧ ⊗i∈[2..max]. inrem(idx, i, 0)t}
parwork(v, max, idx) {
  if (v ≤ sqrt(max)) {
    { (2 ≤ v ∧ ⊗i∈[2..max]. inrem(idx, i, 0)t/2) *
      (2 ≤ v + 1 ∧ ⊗i∈[2..max]. inrem(idx, i, 0)t/2) }
    worker(v, max, idx) || parwork(v+1, max, idx)
    { (⊗i∈[2..max]. fac(i, v, v) ⇒ outrem(idx, i)t/2 ∧
      ¬fac(i, v, v) ⇒ inrem(idx, i, 0)t/2) *
      (⊗i∈[2..max]. fac(i, v+1, ⌊√max⌋) ⇒ outrem(idx, i)t/2
      ∧ ¬fac(i, v+1, ⌊√max⌋) ⇒ inrem(idx, i, 0)t/2) }
    // Using permission combination axioms
    {⊗i∈[2..max]. fac(i, v, ⌊√max⌋) ⇒ outrem(idx, i)t ∧
     ¬fac(i, v, ⌊√max⌋) ⇒ inrem(idx, i, 0)t}
  } }
{⊗i∈[2..max]. fac(i, v, ⌊√max⌋) ⇒ outrem(idx, i)t ∧
 ¬fac(i, v, ⌊√max⌋) ⇒ inrem(idx, i, 0)t}

{2 ≤ v ∧ ⊗i∈[2..max]. inrem(idx, i, 0)t}
worker(v, max, idx) {
  c := v + v;
  while (c ≤ max) {
    {⊗i∈[2..(c-1)]. fac(i, v, v) ⇒ outrem(idx, i)t ∧
     ¬fac(i, v, v) ⇒ inrem(idx, i, 0)t
     * ⊗j∈[c..max]. inrem(idx, j, 0)t}
    remove(idx, c);
    c := c + v;
  }
}
{⊗i∈[2..max]. fac(i, v, v) ⇒ outrem(idx, i)t ∧
 ¬fac(i, v, v) ⇒ inrem(idx, i, 0)t}

```

Figure 7. Proofs for the sieve and worker programs.

Iterator specification. An iterator is represented by the abstract predicate $\text{iter}(it, h, S, K, i)$, which describes an iterator it , iterating over index h . The set S contains the key-value pairs that are in the index and have not yet been returned by next , while K is the set of keys that are not assigned in the index. The iterator has definite permission i for every key in $\text{keys}(S) \cup K$.

Our specification for the three iterator operations is shown in Figure 8. Creating an iterator for an index requires definite information about the state of each key in that index, in the form of in_{def} and out_{def} predicates for all keys. It is not sensible for two threads to share the same iterator, as each thread will iterate over an unknown subset of the underlying index. As such, the iter predicate cannot be split for sharing between threads. However, notice that we can create multiple iterators for a single index, as createIter requires only fractional permission for each key.

The two specifications for next handle the case where the client has not yet seen all key-value pairs in the iterator (in which case, a pair is returned non-deterministically), and when it has (in which case, nil is returned for both the key and value). Destroying an iterator liberates all of the index predicates that have not been returned by next , including the out_{def} predicates.

5.1 Example: a more powerful map.

In §3.1, we verified rangeMap , an algorithm that mapped all values in an index from a given key range through a function, replacing the values with the result. Using an iterator, we can define a concurrent map that does not require a key range, and works over all entries in an index. To avoid having to reason about function pointers, we assume the particular function f is baked into the algorithm source.

```

map_f(h) {
  it := createIter(h);
  map_worker(it, h);
  destroyIter(it);
}
map_worker(it, h) {
  (k,v) := next(it);
  if (k ≠ nil) {
    remove(h, k);
    insert(h, k, f(v));
  }
  || map_worker(it, h);
}

```

A proof of correctness for map_f is given in Figure 9.

5.2 Example: counting distinct values.

We can use an index to store discovered information, and then use iteration to summarise what has been discovered. To illustrate this, we give an algorithm which counts the number of distinct values stored in a tree. Both the tree and the secondary store used for recording distinct values are implemented using our index specification.

Our algorithm is defined as follows:

$$\begin{aligned}
& \{ \bigotimes_{(k,v) \in S} \text{in}_{\text{def}}(\mathbf{h}, k, v)_i * \bigotimes_{k \notin \text{keys}(S)} \text{out}_{\text{def}}(\mathbf{h}, k)_i \} \text{it} := \text{createIter}(\mathbf{h}) \{ \text{iter}(\text{it}, \mathbf{h}, S, \overline{\text{keys}(S)}, i) \} \\
& \quad \{ \text{iter}(\text{it}, \mathbf{h}, S, K, i) \wedge S \neq \emptyset \} (k, v) := \text{next}(\text{it}) \quad \left\{ (k, v) \in S \wedge \text{iter}(\text{it}, \mathbf{h}, S \setminus \{(k, v)\}, K, i) * \right. \\
& \quad \quad \left. \text{in}_{\text{def}}(\mathbf{h}, k, v)_i \right\} \\
& \quad \{ \text{iter}(\text{it}, \mathbf{h}, \emptyset, K, i) \} (k, v) := \text{next}(\text{it}) \quad \{ \text{iter}(\text{it}, \mathbf{h}, \emptyset, K, i) \wedge k = \text{nil} \wedge v = \text{nil} \} \\
& \quad \{ \text{iter}(\text{it}, \mathbf{h}, S, K, i) \} \text{destroyIter}(\text{it}) \quad \{ \bigotimes_{(k,v) \in S} \text{in}_{\text{def}}(\mathbf{h}, k, v)_i * \bigotimes_{k \in K} \text{out}_{\text{def}}(\mathbf{h}, k)_i \}
\end{aligned}$$

Figure 8. Specification for iterators. For `createIter`, set S denotes the key-value pairs of h , $\text{keys}(S)$ denotes the assigned keys of h , and $\overline{\text{keys}(S)}$ denotes the unassigned keys.

$$\begin{aligned}
& \{ \bigotimes_{(k,v) \in S} \text{in}_{\text{def}}(\mathbf{h}, k, v)_1 * \bigotimes_{k \notin \text{keys}(S)} \text{out}_{\text{def}}(\mathbf{h}, k)_1 \} \\
& \text{map_f}(\mathbf{h}) \{ \\
& \quad \text{it} := \text{createIter}(\mathbf{h}); \\
& \quad \{ \text{iter}(\text{it}, \mathbf{h}, S, \overline{\text{keys}(S)}, 1) \} \\
& \quad \text{map_worker}(\text{it}, \mathbf{h}); \\
& \quad \{ \text{iter}(\text{it}, \mathbf{h}, \emptyset, \overline{\text{keys}(S)}, 1) * \bigotimes_{(k,v) \in S} \text{in}_{\text{def}}(\mathbf{h}, k, f(v))_1 \} \\
& \quad \text{destroyIter}(\text{it}); \\
& \} \\
& \{ \bigotimes_{(k,v) \in S} \text{in}_{\text{def}}(\mathbf{h}, k, f(v))_1 * \bigotimes_{k \notin \text{keys}(S)} \text{out}_{\text{def}}(\mathbf{h}, k)_1 \} \\
& \{ \text{iter}(\text{it}, \mathbf{h}, S, K, 1) \} \\
& \text{map_worker}(\text{it}, \mathbf{h}) \{ \\
& \quad (k, v) := \text{next}(\text{it}); \\
& \quad \left\{ (k, v) \in S \wedge \text{iter}(\text{it}, \mathbf{h}, S \setminus \{(k, v)\}, K, 1) * \text{in}_{\text{def}}(\mathbf{h}, k, v)_1 \right. \\
& \quad \quad \left. \vee (\text{iter}(\text{it}, \mathbf{h}, \emptyset, K, 1) \wedge k = \text{nil} \wedge v = \text{nil}) \right\} \\
& \quad \text{if } (k \neq \text{nil}) \{ \\
& \quad \quad \left\{ (k, v) \in S \wedge \text{iter}(\text{it}, \mathbf{h}, S \setminus \{(k, v)\}, K, 1) * \text{in}_{\text{def}}(\mathbf{h}, k, v)_1 \right\} \\
& \quad \quad (\\
& \quad \quad \quad \{ (k, v) \in S \wedge \text{in}_{\text{def}}(\mathbf{h}, k, v)_1 \} \\
& \quad \quad \quad \text{remove}(\mathbf{h}, k); \text{insert}(\mathbf{h}, k, f(v)); \\
& \quad \quad \quad \{ (k, v) \in S \wedge \text{in}_{\text{def}}(\mathbf{h}, k, f(v))_1 \} \\
& \quad \quad) \parallel \\
& \quad \quad \{ \text{iter}(\text{it}, \mathbf{h}, S \setminus \{(k, v)\}, K, 1) \} \\
& \quad \quad \text{map_worker}(\text{it}, \mathbf{h}); \\
& \quad \quad \{ \text{iter}(\text{it}, \mathbf{h}, \emptyset, K, 1) * \bigotimes_{(k',v') \in S \setminus \{(k,v)\}} \text{in}_{\text{def}}(\mathbf{h}, k', f(v'))_1 \} \\
& \quad \} \\
& \} \\
& \{ \text{iter}(\text{it}, \mathbf{h}, \emptyset, K, 1) * \bigotimes_{(k,v) \in S} \text{in}_{\text{def}}(\mathbf{h}, k, f(v))_1 \}
\end{aligned}$$

Figure 9. Proof outline for `map_f`.

$$\begin{aligned}
& \text{count}(\text{it}, k, \text{is}) \{ \\
& \quad \text{fetch}(\text{it}, k, \text{is}); \\
& \quad \text{itr} := \text{createIter}(\text{is}); \\
& \quad \text{num} := 0; \\
& \quad (k, v) := \text{next}(\text{itr}); \\
& \quad \text{while } (k \neq \text{nil}) \{ \\
& \quad \quad \text{num} := \text{num} + 1; \\
& \quad \quad \text{remove}(\text{is}, k); \\
& \quad \quad (k, v) := \text{next}(\text{itr}); \\
& \quad \} \\
& \quad \text{destroyIter}(\text{itr}); \\
& \quad \text{return num}; \\
& \} \\
& \text{fetch}(\text{it}, k, \text{is}) \{ \\
& \quad \text{if } (k \neq \text{nil}) \{ \\
& \quad \quad (k1, k2, v) := \\
& \quad \quad \quad \text{search}(\text{it}, k); \\
& \quad \quad \text{insert}(\text{is}, v, k); \\
& \quad \quad (\text{fetch}(\text{it}, k1, \text{is}) \parallel \\
& \quad \quad \quad \text{fetch}(\text{it}, k2, \text{is})); \\
& \quad \} \\
& \}
\end{aligned}$$

The function `count` calls `fetch` to construct the index `is` from the values of the tree in the index `it`. Values can appear at more than one tree node, but are only recorded once in the `is` index. `count` then iterates the index `is`, counting the number of distinct values discovered. We define a tree predicate annotated with the set of values stored in the tree:

$$\begin{aligned}
& \text{tree}(h, k, vs) \triangleq \exists k_1, k_2, vs_1, vs_2, v. \\
& \quad (k = \text{nil} \wedge vs = \emptyset \wedge \text{emp}) \vee \\
& \quad \left(\text{tree}(h, k_1, vs_1) * \text{tree}(h, k_2, vs_2) \right) \\
& \quad \quad * \text{in}_{\text{def}}(h, k, \langle k_1, k_2, v \rangle)_1 \\
& \quad \quad \wedge vs = vs_1 \cup vs_2 \cup \{v\}
\end{aligned}$$

`fetch` and `count` satisfy the following specifications:

$$\begin{aligned}
& \{ \text{tree}(\text{it}, k, vs) * \bigotimes_{k'} \text{out}_{\text{ins}}(\text{is}, k', \text{Keys})_i \} \\
& \quad \text{fetch}(\text{it}, k, \text{is}) \\
& \{ \text{tree}(\text{it}, k, vs) * \bigotimes_{k' \notin vs} \text{out}_{\text{ins}}(\text{is}, k', \text{Keys})_i \} \\
& \quad * \bigotimes_{k' \in vs} \text{in}_{\text{ins}}(\text{is}, k', \text{Keys})_i \\
& \{ \text{tree}(\text{it}, k, vs) * \bigotimes_{k'} \text{out}_{\text{def}}(\text{is}, k')_1 \} \\
& \quad \text{count}(\text{it}, k, \text{is}) \\
& \{ \text{tree}(\text{it}, k, vs) * \bigotimes_{k'} \text{out}_{\text{def}}(\text{is}, k')_1 \wedge \text{ret} = |vs| \}
\end{aligned}$$

Figure 10 shows an outline proof of these specifications. The part of the proof associated with searching the tree is similar in structure to O'Hearn *et al*'s proof of tree disposal using concurrent separation logic [15]. The difference is that we are able to reason abstractly about concurrently inserting into the `is` index.

6. Verifying Index Implementations

In this section, we verify three quite different concurrent index implementations against our abstract specification. Note that proving implementations is an obligation on the writer of the module – clients can reason using our specification without any knowledge of such proofs. We first introduce a simple list-based implementation and show that it satisfies the disjoint specification of §3. This example is given to develop our technical approach. We then prove that a hash table implementation satisfies the sharing specification of §4. Finally, we show that our approach scales to quite complex implementations, by outlining our proof that the B^{Link} tree algorithm satisfies the sharing specification.

```

{tree(it, k, vs) *  $\bigotimes_{k'} \text{out}_{\text{ins}}(\text{is}, k', \text{ltKeys})_i$ }
fetch(it, k, is) {
  if (k  $\neq$  nil) {
    { $\exists k_1, k_2, vs_1, vs_2, v. \text{tree}(\text{it}, k_1, vs_1) * \text{tree}(\text{it}, k_2, vs_2)$ 
      *  $\text{in}_{\text{def}}(\text{it}, k, \langle k_1, k_2, v \rangle)_1 * \bigotimes_{k'} \text{out}_{\text{ins}}(\text{is}, k', \text{ltKeys})_i$ 
       $\wedge vs = vs_1 \cup vs_2 \cup \{v\}$ }
    (k1, k2, v) := search(it, k);
    { $\exists vs_1, vs_2. \text{tree}(\text{it}, k_1, vs_1) * \text{tree}(\text{it}, k_2, vs_2)$ 
      *  $\text{in}_{\text{def}}(\text{it}, k, \langle k_1, k_2, v \rangle)_1 * \bigotimes_{k'} \text{out}_{\text{ins}}(\text{is}, k', \text{ltKeys})_i$ 
       $\wedge vs = vs_1 \cup vs_2 \cup \{v\}$ }
    insert(is, v, k);
    { $\exists vs_1, vs_2. \text{tree}(\text{it}, k_1, vs_1) * \text{tree}(\text{it}, k_2, vs_2)$ 
      *  $\text{in}_{\text{def}}(\text{it}, k, \langle k_1, k_2, v \rangle)_1 * \text{in}_{\text{ins}}(\text{is}, v, \text{ltKeys})_{\frac{i}{2}}$ 
      *  $\bigotimes_{k' \neq v} \text{out}_{\text{ins}}(\text{is}, k', \text{ltKeys})_{\frac{i}{2}}$ 
      *  $\bigotimes_{k'} \text{out}_{\text{ins}}(\text{is}, k', \text{ltKeys})_{\frac{i}{2}} \wedge vs = vs_1 \cup vs_2 \cup \{v\}$ }
    ( fetch(it, k1, is) || fetch(it, k2, is) );
  }
}
{tree(it, k, vs) *  $\bigotimes_{k' \notin vs} \text{out}_{\text{ins}}(\text{is}, k', \text{ltKeys})_i$ 
  *  $\bigotimes_{k' \in vs} \text{in}_{\text{ins}}(\text{is}, k', \text{ltKeys})_i$ }

{tree(it, k, vs) *  $\bigotimes_{k'} \text{out}_{\text{def}}(\text{is}, k')_1$ }
count(it, k, is) {
  fetch(it, k, is);
  {tree(it, k, vs) *  $\bigotimes_{k' \notin vs} \text{out}_{\text{def}}(\text{is}, k')_1$  *
    { $\bigotimes_{k' \in vs} \exists v' \in \text{ltKeys}. (k', v') \in S \wedge \text{in}_{\text{def}}(\text{is}, k', v')_1$ 
       $\wedge vs = \text{keys}(S)$ }
  }
  itr := createIter(is);
  num := 0;
  {tree(it, k, vs) * iter(itr, is, S,  $\overline{vs}$ , 1)  $\wedge vs = \text{keys}(S) \wedge \text{num} = 0$ }
  (k, v) := next(itr);
  while(k  $\neq$  nil) {
    num := num + 1;
    remove(is, k);
    { $\exists vs', S'. \text{tree}(\text{it}, k, vs) * \text{iter}(\text{itr}, \text{is}, S', \overline{vs}, 1) *$ 
      { $\bigotimes_{k' \in vs \setminus vs'} \text{out}_{\text{def}}(\text{is}, k')_1 \wedge |vs'| + \text{num} = |vs|$ 
       $\wedge vs' = \text{keys}(S')$ }
    }
    (k, v) := next(itr);
  }
  {tree(it, k, vs) * iter(itr, is,  $\emptyset$ ,  $\overline{vs}$ , 1) *
    { $\bigotimes_{k' \in vs} \text{out}_{\text{def}}(\text{is}, k')_1 \wedge \text{num} = |vs|$ }
  }
  destroyIter(itr);
  return num;
}
{tree(it, k, vs) *  $\bigotimes_{k'} \text{out}_{\text{def}}(\text{is}, k')_1 \wedge \text{ret} = |vs|$ }

```

Figure 10. Outline proofs of count and fetch.

Approach: Concurrent Abstract Predicates. We use the techniques developed in the work on concurrent abstract predicates (CAP) [6] to prove that index implementations satisfy our specification. This approach extends separation logic with both explicit reasoning about sharing within modules, and a powerful abstraction mechanism that can hide sharing from clients.

Sharing between threads is represented in CAP by shared regions, denoted by boxed assertions of the form \boxed{P}_I^r . The assertion P describes the contents of the region, r is the name of the region, and I is an interface environment specifying type of mutations threads can perform on P . Assertions on shared regions behave additively under $*$, that is,

$$\boxed{P}_I^r * \boxed{Q}_I^r \triangleq \boxed{P \wedge Q}_I^r$$

A shared region can be mutated by the environment threads. This means that assertions about shared regions must be *stable*: that is, invariant under other threads' interference.

Often, different threads can perform different operations over a shared resource: for example, they may be able to mutate different keys in a shared index. To represent this behaviour, CAP introduces *capabilities*. These are resources giving a thread the ability to perform particular operations. Threads can hold both non-exclusive and exclusive capabilities. When an exclusive capability is held, no other thread can perform the associated operation.

Shared regions and capabilities can be abstracted using predicates in the manner described in §2. Each predicate represents both some information about a shared region, and some ability held by the thread to modify the shared region. If the combination of capabilities held ensures that the shared assertion is invariant, then stability need not be considered by clients, and the predicate can be treated abstractly.

In the discussion below, we assume the proof system and semantics given in [6], and only give details necessary for understanding the proof structure. The interested reader is referred to [6] for other technical details, including a proof of soundness for the CAP logic.

6.1 Linked List Implementation

To illustrate our approach, we consider a very simple index implementation which uses a linked list with a single lock protecting the entire list³. The code for this implementation is given in Figure 11. In order to simplify the presentation, we only consider the disjoint specification of §3 in this section. Additional technicalities are required to handle the full sharing specification of §4. We give these technicalities in §6.3, when we verify the B^{Link} tree implementation against the sharing specification.

Before performing any operation on the list, the thread first acquires the lock. The `search` operation traverses the list checking if an element matches the key; if so, it returns the corresponding value. The `insert` operation is similar to `search`. However, if it cannot find the key, it creates a new node and adds it to the head of the list. The `remove` operation searches for the key to be removed. If it finds the key, it updates the previous node in the list to point to the following node. The node, having been thus removed from the list, is then deleted.

³This example is quite similar to the coarse-grained set example from [6].

```

search(h, k) {
  lock(h.lk);
  e := h.nxt;
  while (e ≠ nil) {
    if (e.key = k) {
      unlock(h.lk);
      return e.val;
    }
    e := e.nxt;
  }
  unlock(h.lk);
  return nil;
}

insert(h, k, v) {
  lock(h.lk);
  e := h.nxt;
  while (e ≠ nil) {
    if (e.key = k) {
      unlock(h.lk);
      return;
    }
    e := e.nxt;
  }
  e := makeNode(k, v, h.nxt);
  h.nxt := e;
  unlock(h.lk);
}

remove(h, k) {
  lock(h.lk);
  e := h.nxt;
  prev := h;
  while (e ≠ nil) {
    if (e.key = k) {
      prev.nxt := e.nxt;
      disposeNode(e);
      unlock(h.lk);
      return;
    }
    prev := e;
    e := e.nxt;
  }
  unlock(h.lk);
}

```

Figure 11. Linked list operations.

Interpretation of abstract predicates. In order to prove that the operations of the implementation are correct with respect to our specification, we first give concrete interpretations to the abstract predicates.

We begin by defining a predicate $ls(a, H)$, corresponding to list with address a and representing the index state $H : \text{Keys} \rightarrow \text{Vals}$. This is defined in terms of the inductive predicate $lseg(a, b, H)$, which represents a list segment with address a and final pointer b , having key-value elements given by H . A list segment is either empty, in which case $a = b$ and $H = \emptyset$, or it consists of a node at address a whose key and value are taken from H , and whose `nxt` field points to a list segment of the rest of the keys and values. The definition of $lseg$ is, in turn, defined in terms of the predicate $node(a, k, v, n)$, which simply represents a node at address a whose key, val and `nxt` fields are k, v , and n respectively. The formal definitions of these predicates are as follows:

$$\begin{aligned}
node(a, k, v, n) &\triangleq a.\text{key} \mapsto k * a.\text{val} \mapsto v * a.\text{nxt} \mapsto n \\
lseg(a, b, H) &\triangleq (a = b \wedge H = \emptyset) \vee \\
&\quad \exists k, v, n, H'. H = H' \uplus \{k \mapsto v\} \wedge \\
&\quad \quad node(a, k, v, n) * lseg(n, b, H') \\
ls(a, H) &\triangleq lseg(a, nil, H)
\end{aligned}$$

Using the ls predicate, we can give a concrete interpretation to our index predicates for the linked list implementation of an index, as follows:

$$\begin{aligned}
in(h, k, v) &\triangleq \exists r, l, H. H(k) = v \wedge [\text{LOCK}(k)]_1^r * \\
&\quad \boxed{\text{lock}(h.lk, r, k) * h.nxt \mapsto l * ls(l, H)}_{I(r, h)}^r \\
out(h, k) &\triangleq \exists r, l, H. k \notin \text{dom}(H) \wedge [\text{LOCK}(k)]_1^r * \\
&\quad \boxed{\text{lock}(h.lk, r, k) * h.nxt \mapsto l * ls(l, H)}_{I(r, h)}^r
\end{aligned}$$

Here, the boxed assertion describes the region r shared between all the threads that can access the list. This boxed assertion says that region r contains a lock at $h.lk$ (we define the predicate `lock` below) and a pointer $h.nxt$ to a list representing the contents of the index. The index state H is existentially quantified; the assertions only specify whether the key k is in the index, and its value, if any.

Both predicates also include the (unshared) capability resource $[\text{LOCK}(k)]_1^r$. A thread with such a capability in its local state is able to update the contents of the corresponding region r by performing the `LOCK(k)` action that is defined in the interference environment $I(r, h)$ associated with the region. We will give the formal definition of $I(r, h)$ presently; intuitively, the `LOCK(k)` action allows a thread to acquire the lock in order to subsequently add or remove the key k . The subscript 1 in the capability denotes that it is an *exclusive* capability: no other thread can perform the action. The exclusivity of the permission ensures that the predicates are stable, since the state of key k in the index cannot be changed by any other thread.

We define the predicate $\text{lock}(x, r, k)$ as follows:

$$\begin{aligned}
unlocked(x, r) &\triangleq x \mapsto 0 * \bigotimes_{i \in \text{Keys}} [\text{MOD}(i)]_1^r \\
locked(x, r, j) &\triangleq x \mapsto 1 * \bigotimes_{i \in \text{Keys} \setminus \{j\}} [\text{MOD}(i)]_1^r \\
\text{lock}(x, r, k) &\triangleq unlocked(x, r) \vee \exists j \neq k. \text{locked}(x, r, j)
\end{aligned}$$

This lock predicate contains a shared lock bit and a collection of capabilities. Each capability $[\text{MOD}(k)]_1^r$ controls the ability to add or remove a particular key k from the shared list in region r . When these capabilities are in the shared region, no thread is able to modify the list; such is the case when the lock is unlocked. When the lock is locked, a single $[\text{MOD}(j)]_1^r$ capability is held by some thread, allowing it to perform the necessary update, but only to the key j . The $\text{lock}(h.lk, r, k)$ predicate ensures, that no other thread can have the $[\text{MOD}(k)]_1^r$ capability, and hence update key k .

Describing Interference. The meaning of the capabilities $[\text{LOCK}(k)]_1^r$ and $[\text{MOD}(k)]_1^r$ is determined by the *interference environment* associated with the region r : $I(r, h)$. This defines the possible state mutations that can occur over a given shared region. The environment defines the meaning of capabilities in terms of actions, written $P \rightsquigarrow Q$. When a

thread holds a capability mapped to an action $P \rightsquigarrow Q$, it is permitted to replace a part of the region matching P with a part matching Q . To perform the action, a thread may transfer resource between the region and its own local state, and may mutate it in an atomic operation.

For the linked list implementation, the interference environment $I(r, h)$ is defined as follows:

$$\text{MOD}(k): \begin{cases} h.\text{next} \mapsto l * \text{ls}(l, H) \\ \rightsquigarrow h.\text{next} \mapsto l' * \text{ls}(l', H \uplus \{k \mapsto v\}) \\ h.\text{next} \mapsto l * \text{ls}(l, H) \\ \rightsquigarrow h.\text{next} \mapsto l' * \text{ls}(l', H \setminus \{k\}) \end{cases}$$

$$\text{LOCK}(k): \begin{cases} h.\text{lk} \mapsto 0 * [\text{MOD}(k)]_1^r \rightsquigarrow h.\text{lk} \mapsto 1 \\ h.\text{lk} \mapsto 1 \rightsquigarrow h.\text{lk} \mapsto 0 * [\text{MOD}(k)]_1^r \end{cases}$$

The definition of $\text{MOD}(k)$ says that a thread holding a capability $[\text{MOD}(k)]_1^r$ is allowed to update the list by adding or removing the key k . The definition of $\text{LOCK}(k)$ says that the thread is allowed to set or unset the lock bit. Recall that actions replace part of the shared state, so the definition of $\text{LOCK}(k)$ implies that a thread acquiring the lock also acquires the capability $[\text{MOD}(k)]_1^r$, which leaves the shared state. Similarly, when releasing the lock it must give up the capability $[\text{MOD}(k)]_1^r$. In this way, acquiring the lock gives a thread the ability to modify the contents of the list.

Verifying the operations. Having given concrete definitions to the index predicates, we can verify that the module's implementations of `add`, `remove` and `search` match our high-level specification. Figure 12 shows one such proof, establishing that the implementation of `insert` matches the following abstract specification:

$$\{\text{out}(h, k)\} \text{ insert}(h, k, v) \{\text{in}(h, k, v)\}$$

In the proof, mutations of the shared state require that the thread holds a capability permitting the mutation. These points in `insert` are annotated by program comments. For example, towards the end of `insert`, the assignment `h.next := e` assigns to the shared location `h.next`. This mutation corresponds to performing the first of the actions associated with the $[\text{MOD}(k)]_1^r$ capability, held in the local state. The action requires that initially `h.next` should point to a list representing some index state H , and that after the assignment it should point to a list representing the state $H \uplus \{k \mapsto v\}$ for some v . By considering the predicate definitions, this is clearly the case.

It is necessary to check that the all assertions in the proof are stable. In fact, once the lock has been acquired, the only actions which can affect the shared state are $\text{MOD}(k)$ and $\text{LOCK}(k)$. Since full permission to both is held in local state, no interference can happen, and so the assertions are stable.

For the pre- and postconditions, the list may be locked or unlocked, but it can only be modified with respect to

```

{out(h, k)}
insert(h, k, v) {
  {∃r, l, H. k ∉ dom(H) ∧ [LOCK(k)]_1^r *
   {lock(h.lk, r, k) * h.next ↦ l * ls(l, H)}_{I(r, h)}^r}
  lock(h.lk); // use the capability [LOCK(k)]_1^r.
  {∃r, l, H. k ∉ dom(H) ∧ [MOD(k)]_1^r * [LOCK(k)]_1^r *
   {locked(h.lk, r, k) * h.next ↦ l * ls(l, H)}_{I(r, h)}^r}
  e := h.next;
  while (e ≠ nil) {
    {∃r, l, H, H_1, H_2, k', v', n. k ∉ dom(H) ∧
     [MOD(k)]_1^r * [LOCK(k)]_1^r *
     {locked(h.lk, r, k) * h.next ↦ l *
      lseg(l, e, H_1) * node(e, k', v', n) * ls(n, H_2)
      ∧ H_1 ∪ H_2 ∪ {k' ↦ v'} = H}_{I(r, h)}^r}
    if (e.key = k) {
      {false} // this branch is for k in the set
      unlock(h.lk);
      return;
    }
    e := e.next;
    {∃r, l, H, H_1, H_2. k ∉ dom(H) ∧
     [MOD(k)]_1^r * [LOCK(k)]_1^r *
     {locked(h.lk, r, k) * h.next ↦ l *
      lseg(l, e, H_1) * ls(e, H_2) ∧ H_1 ∪ H_2 = H}_{I(r, h)}^r}
  }
  {∃r, l, H. k ∉ dom(H) ∧ [MOD(k)]_1^r * [LOCK(k)]_1^r *
   {locked(h.lk, r, k) * h.next ↦ l * ls(l, H)}_{I(r, h)}^r}
  e := makeNode(k, v, h.next);
  {∃r, l, H. k ∉ dom(H) ∧ node(e, k, v, l) *
   [MOD(k)]_1^r * [LOCK(k)]_1^r *
   {locked(h.lk, r, k) * h.next ↦ l * ls(l, H)}_{I(r, h)}^r}
  h.next := e; // use the capability [MOD(k)]_1^r.
  {∃r, l, H. k ∉ dom(H) ∧ [MOD(k)]_1^r * [LOCK(k)]_1^r *
   {locked(h.lk, r, k) * h.next ↦ e *
    node(e, k, v, l) * ls(l, H)}_{I(r, h)}^r}
  unlock(h.lk); // use the capability [LOCK(k)]_1^r.
  {∃r, l, H. H(k) = v ∧ [LOCK(k)]_1^r *
   {lock(h.lk, r, k) * h.next ↦ l * ls(l, H)}_{I(r, h)}^r}
}
{in(h, k, v)}

```

Figure 12. Proof outline for linked list `insert`.

keys other than k . Since no information about such keys is contained in these assertions, they are also stable.

Figure 13 shows a proof that `remove` satisfies the specification:

$$\{\text{in}(h, k, v)\} \text{ remove}(h, k) \{\text{out}(h, k)\}$$

This proof makes use of the second MOD action to remove the node with key k from the linked list. The stability of assertions in this proof follows by the same argument as for the `insert` proof.

```

{in(h, k, v)}
remove(h, k) {
  {
     $\exists r, l, H. H(k) = v \wedge [\text{LOCK}(k)]_1^r * \text{lock}(h.lk, r, k) * h.nxt \mapsto l * \text{ls}(l, H)$ 
  }I(r,h)
  lock(h.lk); // use the capability  $[\text{LOCK}(k)]_1^r$ 
  {
     $\exists r, l, H. H(k) = v \wedge [\text{MOD}(k)]_1^r * [\text{LOCK}(k)]_1^r * \text{locked}(h.lk, r, k) * h.nxt \mapsto l * \text{ls}(l, S)$ 
  }I(r,h)
  e := h.nxt;
  prev := h;
  // loop invariant:
  {
     $\exists r, l, H. H(k) = v \wedge [\text{MOD}(k)]_1^r * [\text{LOCK}(k)]_1^r * \text{locked}(h.lk, r, k) * h.nxt \mapsto l * \exists H'. (\text{prev} = h \wedge H = H' \wedge e = l \wedge \text{emp}) \vee (\exists H'', k_p, v_p. H = H' \uplus H'' \uplus \{k_p \mapsto v_p\} \wedge \text{lseg}(l, \text{prev}, H'') * \text{node}(\text{prev}, k_p, v_p, e) * \exists k_e, v_e, a. \text{node}(e, k_e, v_e, a) * \text{lseg}(a, \text{nil}, H' \setminus \{k_e\}) \wedge H'(k_e) = v_e \wedge H'(k) = v)$ 
  }I(r,h)
  while (e ≠ nil) {
    {
       $\exists r, l, H. H(k) = v \wedge [\text{MOD}(k)]_1^r * [\text{LOCK}(k)]_1^r * \text{locked}(h.lk, r, k) * h.nxt \mapsto l * \exists H'. (\text{prev} = h \wedge H = H' \wedge e = l \wedge \text{emp}) \vee (\exists H'', k_p, v_p. H = H' \uplus H'' \uplus \{k_p \mapsto v_p\} \wedge \text{lseg}(l, \text{prev}, H'') * \text{node}(\text{prev}, k_p, v_p, e) * \exists k_e, v_e, a. \text{node}(e, k_e, v_e, a) * \text{lseg}(a, \text{nil}, H' \setminus \{k_e\}) \wedge H'(k_e) = v_e \wedge H'(k) = v)$ 
    }I(r,h)
  }
}

if (e.key = k) {
  {
     $\exists r, l, H. H(k) = v \wedge [\text{MOD}(k)]_1^r * [\text{LOCK}(k)]_1^r * \text{locked}(h.lk, r, k) * h.nxt \mapsto l * \exists H'. (\text{prev} = h \wedge H = H' \wedge e = l \wedge \text{emp}) \vee (\exists H'', k_p, v_p. H = H' \uplus H'' \uplus \{k_p \mapsto v_p\} \wedge \text{lseg}(l, \text{prev}, H'') * \text{node}(\text{prev}, k_p, v_p, e) * \exists a. \text{node}(e, k, v, a) * \text{lseg}(a, \text{nil}, H' \setminus \{k\})$ 
  }I(r,h)
  prev.nxt := e.nxt;
  // use the capability  $[\text{MOD}(k)]_1^r$ 
  {
     $\exists r, l, H, a. H(k) = v \wedge [\text{MOD}(k)]_1^r * [\text{LOCK}(k)]_1^r * \text{node}(e, k, v, a) * \text{locked}(h.lk, r, k) * h.nxt \mapsto l * \exists H'. (\text{prev} = h \wedge H = H' \wedge a = l \wedge \text{emp}) \vee (\exists H'', k_p, v_p. H = H' \uplus H'' \uplus \{k_p \mapsto v_p\} \wedge \text{lseg}(l, \text{prev}, H'') * \text{node}(\text{prev}, k_p, v_p, a) * \text{lseg}(a, \text{nil}, H' \setminus \{k\})$ 
  }I(r,h)
  disposeNode(e);
  {
     $\exists r, l, H. k \notin \text{dom}(H) \wedge [\text{MOD}(k)]_1^r * [\text{LOCK}(k)]_1^r * \text{locked}(h.lk, r, k) * h.nxt \mapsto l * \text{ls}(l, S)$ 
  }I(r,h)
  unlock(h.lk); // use the capability  $[\text{LOCK}(k)]_1^r$ 
  {
     $\exists r, l, H. k \notin \text{dom}(H) \wedge [\text{LOCK}(k)]_1^r * \text{lock}(h.lk, r, k) * h.nxt \mapsto l * \text{ls}(l, H)$ 
  }I(r,h)
  return;
}
prev := e;
e := e.nxt;
}
// since  $\text{lseg}(\text{nil}, \text{nil}, H') \wedge H'(k) = v \implies \text{false}$ 
{false}
unlock(h.lk);
}
{out(h, k, v)}

```

Figure 13. Proof outline for linked list remove.

Verifying the axioms. As well as proving the specifications for the operations, our other obligation is establishing that implementation satisfies the axioms of the abstract specification. To do this, we use the concrete definitions for the abstract predicates. For example, we prove the following axiom from the disjoint specification:

$$\text{in}(h, k, v) * \text{out}(h, k) \implies \text{false}$$

If we expand the predicate definitions on the left-hand side of this implication, we end up with the following assertion:

$$\begin{aligned} & \exists r, l, H. H(k) = v \wedge [\text{LOCK}(k)]_1^r * \\ & \quad \text{lock}(h.lk, r, k) * h.nxt \mapsto l * \text{ls}(l, H) \Big|_{I(r,h)}^r * \\ & \exists r, l, H. k \notin \text{dom}(H) \wedge [\text{LOCK}(k)]_1^r * \\ & \quad \text{lock}(h.lk, r, k) * h.nxt \mapsto l * \text{ls}(l, H) \Big|_{I(r,h)}^r \end{aligned}$$

The memory location $h.nxt$ cannot belong to more than one region at once, so we can infer that both existentially-quantified rs must refer to the same shared region. The capability $[\text{LOCK}(k)]_1^r$ is exclusive, denoted by the 1 subscript.

Now

$$[\text{LOCK}(k)]_1^r * [\text{LOCK}(k)]_1^r \implies \text{false}$$

and so the axiom holds.

6.2 Hash Table Implementation

We now consider a second index implementation which uses a hash table. The hash table algorithm consists of a fixed-size array and a hashing function mapping from keys to offsets in the array. Each element of the array is a pointer to a secondary index storing the key-value pairs that hash to the associated array offset.

Secondary indexes are often implemented as linked lists, but in fact any kind of index implementation can be used. In this section, we assume that secondary indexes are implemented by *some* module matching our abstract specification, but do not specify which. (To avoid naming conflicts, we rename the methods and predicates of the secondary index to search' , insert' , remove' , in'_{def} , in'_{rem} , etc.) We then show that the resulting hash table module also matches our


```

search(h, k) {
  w := hash(k);
  a := [h+w];
  return (search'(a, k));
}

insert(h, k, v) {
  w := hash(k);
  a := [h+w];
  insert'(a, k, v);
}

remove(h, k) {
  w := hash(k);
  a := [h+w];
  remove'(a, k);
}

```

Figure 14. Hash table operations.

abstract specification. That is, we show that we can build a concurrent index using a (different) index module.

The hash table implementations of `search`, `insert` and `remove` are given in Figure 14. This code assumes a pure hashing function `hash` which takes a key `k` and returns an integer `hash(k)` between 0 and $max - 1$, where max is the size of the hash table array.

Although the implementation we consider here is very simple, it captures the essence of more complicated implementation's such as Java's `ConcurrentHashMap`, which uses resizable hash tables as a secondary index. It would be invaluable to consider such real-world implementations in detail, but this is beyond the scope of the present work.

Interpretation of abstract predicates. All of our index predicates – in_{ins} , out_{ins} , in_{rem} , and so on – consist of a shared region containing a hash table pointer, and a local predicate representing the associated secondary index. Picking an arbitrary example, we define $in_{rem}(h, k, v)_i$ as follows:

$$in_{rem}(h, k, v)_i \triangleq \exists r, h'. \boxed{h + \text{hash}(k) \mapsto h' * \text{true}}^r * in'_{rem}(h', k, v)_i$$

(The predicates have exactly the same form. Only the predicate pertaining to the secondary index changes.)

The shared region contains a pointer from $h + \text{hash}(k)$ to the address of the secondary index, h' . The rest of the hash table array also belongs to the shared region; it is represented in the assertion by `true`. The array of pointers representing the hash table is read only, so the interference environment for the shared region is empty.

The secondary index is represented by the predicate $in'_{rem}(h', k, v)_i$. Note that this definition hides completely the implementation of the secondary index. The hash table simply knows that this element of the index can be queried according to the abstract index specification. State mutations on the secondary index are already captured by the predicate representing it, meaning that they need not be considered when verifying the hash table implementation.

Verifying the operations. A sketch-proof for the hash table implementation of `search` is given in Figure 15. Notice that this proof appeals to the specification of `search'` when

```

{in_def(h, k, v)_i}
search(h, k) {
  {∃r, h'. (h + hash(k) ↦ h' * true)}^r * in'_def(h', k, v)_i}
  w := hash(k);
  a := [h+w];
  {∃r. (h + hash(k) ↦ a * true)}^r * in'_def(a, k, v)_i}
  return (search'(a, k)); // search' specification
  {∃r, h'. (h + hash(k) ↦ h' * true)}^r * in'_def(h', k, v)_i}
  ∧ ret = v
}
{in_def(h, k, v)_i ∧ ret = v}

```

Figure 15. Proof outline for hash table search.

retrieving a value from the appropriate secondary index. Since there are no actions defined for the shared region, stability of our assertions is trivial.

Verifying the axioms. The axioms follow from the axioms of the secondary index. In particular, two predicates involving the same key will be defined in terms of predicates which must be on the same secondary index.

6.3 B^{Link} Tree Implementation

Our final index implementation is Sagiv's B^{Link} tree algorithm [19]. (Note that we only consider the algorithm without compression here.) A B^{Link} tree is a balanced search tree. An example is shown in Figure 16. The leaves of the tree contain they key-value pairs stored in the index in order. Non-leaf (or inner) nodes associate keys with pointers to nodes at the next level down, which direct the traversal of the tree. In addition, the final pointer in each node's list, the link pointer, points to the next node at that level (if it exists). The tree is accessed through a prime block which holds pointers to the first node at each level in the tree.

During inserts, nodes of the tree that are at full capacity may be split by creating a new right sibling and transferring half of the keys to the new node. This new node must then be attached to the level above, which might require further splittings. However, other operations may still need to traverse the tree before this operation is completed. A traversal in progress may therefore have to use link pointers to find the correct leaf. Since the minimum values of nodes are always preserved, and every leaf with a minimum value no less than that of a given node is reachable from that node, such traversals are always possible.

Search operations on a B^{Link} tree are lock-free, and insert and remove operations lock only one node (or two if they are modifying the root node) at a time, making this a highly concurrent implementation of an index. This index algorithm is much more complex than the list or hash table, and is therefore considerably more challenging to verify.

The code the B^{Link} implementation is given in Figure 17.

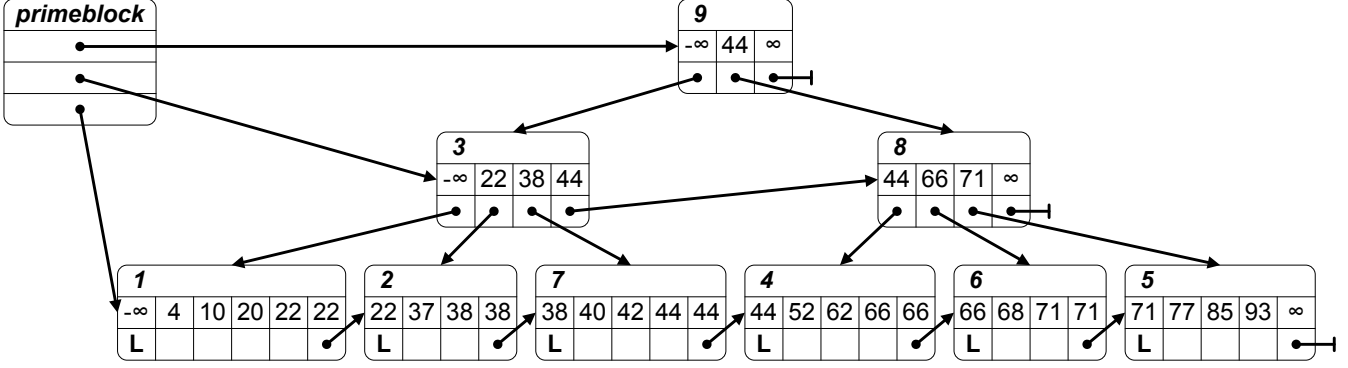


Figure 16. A B^{Link} tree.

Node notation. We use the notation $node(l, k, p, D, k', p')$ to denote the contents of a node, where:

- l determines whether the node is locked ($l = 1$) or not ($l = 0$),
- k is the minimum key for the node (less than or equal to all keys that are reachable from it),
- p is the pointer to the left-most child for an inner node, or nil for a leaf node,
- D is a list of pairs of keys and child pointers (for inner nodes), or keys and stored values (for leaf nodes),
- k' is the maximum key for the node,
- p' is the pointer to the next sibling of the node (or nil if it is the last).

We also use $inner(l, k, p, D, k', p')$ to denote an inner node with the given contents (requiring that $p \neq nil$) and $leaf(l, k, D, k', p')$ to denote a leaf node. In this notation, the contents of node 3 in Figure 16 would be represented as

$$inner(0, -\infty, 1, [(22, 2), (38, 7)], 44, 8).$$

Interpretation of abstract predicates. All of our index predicates are defined as a shared region containing a B^{Link} tree and a collection of shared capabilities, as well as some thread-local capabilities. For example, the predicate $in_{def}(h, k, v)$ is defined as follows:

$$in_{def}(h, k, v)_i \triangleq \exists r. [B_{\in}(h, k, v)]_{I(r, h)}^r * dcaps(k, r, i)$$

The shared assertion $B_{\in}(h, k, v)$ denotes a B^{Link} tree at address h containing the key-value pair (k, v) . The formal definition of this predicate can be found in Appendix A. The predicate $dcaps(k, r, i)$, which is defined in Figure 18, consists of capabilities associated with the current thread.

The permission subscripts of the capabilities are more complex than those we have seen so far: they are deny-guarantee permissions [8]. A guarantee permission, indicated by the subscript (g, i) for $0 < i \leq 1$ (or simply g when

we do not care about the exact value of i), allows a thread to perform the associated action. If the permission is less than 1, other threads may have guarantee permissions to perform the same action – it is a non-exclusive permission. A deny permission, indicated by the subscript (d, i) (or, again, simply d), does not allow the thread to perform the associated action, but precludes the possibility that any other thread will either. Fractions of the same type may be combined by addition, and $(g, 1) = 1 = (d, 1)$ represents exclusive permission; however, a deny permission and a guarantee permission cannot be combined, since they are conflicting. (For further details, see Dodds *et al.* [8].)

The intuitive meaning of the capabilities in the $dcaps$ predicate is as follows. The $[LOCK]_g^r$ capability says that the current thread is allowed to lock nodes in the region r . The $[SWAP]_g^r$ capability allows the in_{def} predicate to be modified to represent different behaviour (for example, by converting it to in_{rem} or unk) provided $i = 1$. The $[REM(0, k)]_{(d, i)}^r$ capability says that neither the current thread, nor any other thread, is allowed to remove the key k from the B^{Link} tree in region r . However, if $i = 1$, then the current thread has the exclusive capability to remove key k from the tree. The $[INS(0, k, v')]_{(d, i)}^r$ capabilities similarly restrict the ability to insert value v' at key k .

The other index predicates are defined in a similar way to in_{def} . For example, the definition of the $in_{rem}(h, k, v)_i$ predicate will include a REM capability for k with permission (g, i) , so that any thread may remove the key from the tree, as well as all INS capabilities for k with permission (d, i) , so that no thread may insert values for the key into the tree. We give the full definitions of the predicates in Appendix A.

Describing Interference. The interference environment, $I(r, h)$, for the B^{Link} tree implementation is markedly more complex than for the list or hash table. It involves a substantial amount of capability swapping to track changes to the shared state and to thread behaviour. Figure 19 gives a few examples of definitions in the interference environment. These definitions can be read as follows:

```

search(h, k) {
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  while(isLeaf(N) = false) {
    cur := next(N, k);
    N := get(cur);
  }
  while(k > highValue(N)) {
    cur := next(N, k);
    N := get(cur);
  }
  if(isIn(N, k)) {
    return(lookup(N, k));
  } else {
    return nil;
  }
}

insertIntoSafe {
  addPair(N, m, w);
  put(N, cur);
  unlock(cur);
}

insertIntoUnsafeRoot {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, x);
  lock(x);
  put(N, cur);
  y := lowValue(N);
  t := highValue(N);
  u := highValue(M);
  r := new();
  R := newNode(y, cur, t, x, u);
  PB := getPrimeBlock(h);
  put(R, r);
  addRoot(PB, r);
  putPrimeBlock(h, PB);
  unlock(cur);
  unlock(x);
}

insert(h, k, v) {
  stack := newStack();
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  while (isLeaf(N) = false) {
    if (k < highValue(N)) {
      push(stack, cur);
    }
    cur := next(N, k);
    N := get(cur);
  }
  level := 1;
  m := k;
  w := v;
  while (true) {
    found := false;
    while (found = false) {
      found := true;
      lock(cur);
      N := get(cur);
      if (isIn(N, m)) {
        unlock(cur);
        return;
      }
      if (m > highValue(N)) {
        unlock(cur);
        found := false;
        while (m > highValue(N)) {
          cur := next(N, m);
          N := get(cur);
        }
      }
    }
  }
  if (isSafe(N)) {
    insertIntoSafe;
    return;
  } else {
    PB := getPrimeBlock(h);
    if (isRoot(PB, cur)) {
      insertIntoUnsafeRoot;
      return;
    } else {
      insertIntoUnsafe;
    }
  }
}

remove(h, k) {
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  while (isLeaf(N) = false) {
    cur := next(N, k);
    N := get(cur);
  }
  while (true) {
    lock(cur);
    N := get(cur);
    if (isIn(N, k)) {
      removePair(N, k);
      put(A, cur);
      unlock(cur);
      return;
    } else {
      unlock(cur);
      if (k > highValue(N)) {
        while (k > highValue(N)) {
          cur := next(N, k);
          N := get(cur);
        }
      } else { // value is not in the tree
        return;
      }
    }
  }
}

insertIntoUnsafe {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, q);
  put(N, cur);
  unlock(cur);
  w := x;
  m := highValue(N);
  level := level + 1;
  if (isEmpty(stack)) {
    PB := getPrimeBlock(h);
    cur := getNodeLevel(PB, level);
  } else {
    cur := pop(stack);
  }
}

```

Figure 17. B^{Link} tree operations.

$$\begin{aligned}
\text{dcaps}(k, r, i) &\triangleq [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(0, k)]_{(d,i)}^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(0, k, v)]_{(d,i)}^r \\
\text{niceNode}(N, k, v, r, h) &\triangleq \exists k_0, p_0, D, k', p'. \left(k' = +\infty \vee \boxed{p' \mapsto \text{node}(-, k', -, -, -, -) * \text{true}}_{I(r,h)}^r \right) \wedge \\
&\quad \left(\left(\left(N = \text{inner}(-, k_0, p_0, D, k', p') \wedge \forall (k, p) \in D. \right. \right. \right. \\
&\quad \left. \left. \left. \boxed{p \mapsto \text{node}(-, k, -, -, -, -) * \text{true}}_{I(r,h)}^r \right) \vee \left(N = \text{leaf}(-, k_0, D, k', p') \wedge \right. \right. \right. \\
&\quad \left. \left. \left. \left(\wedge \boxed{p_0 \mapsto \text{node}(-, k_0, -, -, -, -) * \text{true}}_{I(r,h)}^r \right) \vee \left(k_0 < k \leq k' \Rightarrow (k, v) \in D \right) \right) \right) \\
\text{present}(n, k, p_0) &\triangleq \boxed{n \mapsto \text{node}(-, k, p_0, -, -, -) * \text{true}}_{I(r,h)}^r \\
\text{stLf}(n, N, k, v, r, h) &\triangleq \exists k', k'', D, p'. \left(k'' = +\infty \vee \boxed{p' \mapsto \text{leaf}(-, k'', -, -, -) * \text{true}}_{I(r,h)}^r \right) \\
&\quad \wedge N = \text{leaf}(1, k', D, k'', p') \wedge (k' < k \leq k'' \Rightarrow (k, v) \in D) \wedge \boxed{n \mapsto \text{leaf}(1, k', D, k'', p') * \text{true}}_{I(r,h)}^r
\end{aligned}$$

Figure 18. Predicates used in the B^{Link} tree proofs.

$$\begin{aligned}
\text{LOCK} : \quad &x \mapsto \text{node}(0, k_0, p, D, k', p') * [\text{UNLOCK}(x)]_1^r \rightsquigarrow x \mapsto \text{node}(1, k_0, p, D, k', p') \\
\text{REM}(t, k) : \quad &[\text{MODLR}(t, x, k, i)]_1^r \rightsquigarrow [\text{REM}(t, k)]_{(g,i)}^r * [\text{UNLOCK}(x)]_1^r \\
\text{MODLR}(t, x, k, i) : \quad &\left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D, k', p') * [\text{UNLOCK}(x)]_1^r \\ * \left([\text{REM}(t, k)]_{(g,i)}^r \wedge t = 0 \vee \text{emp} \wedge t = 1 \right) \\ \wedge (k, -) \in D \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D', k', p') \\ * [\text{MODLR}(t, x, k, i)]_1^r \\ \wedge D = D' \uplus (k, -) \end{array} \right)
\end{aligned}$$

Figure 19. Example actions from the B^{Link} tree interference environment.

- **LOCK** allows a thread to lock a node in the B^{Link} tree. When locking, the thread acquires the exclusive capability $[\text{UNLOCK}(x)]_1^r$, allowing it to unlock the node again.
- **REM**(t, k) allows a thread to give up $[\text{REM}(t, k)]_{(g,i)}^r$ and $[\text{UNLOCK}(x)]_1^r$ and acquire the exclusive capability $[\text{MODLR}(t, x, k, i)]_1^r$. This means that a thread which is allowed to remove the key k from the tree and holds the lock on a node x can acquire the right to remove the key k from the leaf node x (the value t is used to track capability transfer in some environments).
- **MODLR**(t, x, k, i) allows a thread to remove a key-value pair $(k, -)$ from a leaf node. In doing so, the thread gives up the capability $[\text{MODLR}(t, x, k, i)]_1^r$ and reacquires the capability $[\text{UNLOCK}(x)]_1^r$, and, if $t = 0$, the capability $[\text{REM}(k)]_{(g,i)}^r$. (We write “ $-$ ” to indicate an unspecified, existentially quantified value.)

We give the full interference environment for the B^{Link} tree implementation in Appendix A.

Note that both $[\text{REM}(0, k)]_{(g,i)}^r$ and $[\text{REM}(1, k)]_{(g,i)}^r$ capabilities allow a thread to remove the key k ; however, the latter requires the thread to leave a $[\text{REM}(1, k)]_{(g,i)}^r$ capability behind in the shared state when it does so. This is used to implement the in_{rem} predicate: if none of the threads

with $\text{in}_{\text{rem}}(k, v)$ predicates remove k then between them they must still be able to produce the full $[\text{REM}(1, k)]_1^r$ capability, proving that none of them did so. Thus the $\text{in}_{\text{rem}}(k, v)_1$ can be converted to $\text{in}_{\text{def}}(k, v)_1$.

Verifying the operations. We give a sketch proof in Figure 20, showing that the B^{Link} tree implementation of `search` matches the following specification:

$$\{\text{in}_{\text{def}}(\mathbf{k}, \mathbf{v}, \mathbf{v})\} \mathbf{r} := \text{search}(\mathbf{h}, \mathbf{k}) \{\text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, \mathbf{v})_i \wedge \mathbf{r} = \mathbf{v}\}$$

The `search` operation only mutates thread-local state, so the thread does not require capabilities to perform actions. However, by owning deny permissions (d, i) on all the **REM** and **INS** capabilities for key \mathbf{k} , the thread can establish that no other thread can modify the value associated with \mathbf{k} . Thus, the assertion that the key-value pair (\mathbf{k}, \mathbf{v}) is contained in the B^{Link} tree is stable.

The proof uses the predicate $\text{niceNode}(N, k, v, r, h)$, defined in Figure 18. The definition of `niceNode` asserts that the node descriptor N contains legitimate information about the tree. If N is an inner node, then the children and link pointers of N must all point to extant nodes in the tree, which have the minimum values specified by N – this ensures that following a pointer reaches an appropriate node. If N is a leaf node into whose range the key k falls, then the key-value

```

{indef(h, k, v)i}
search(h, k) {
  {B∈(h, k, v)I(r,h)r * dcaps(k, r, i)}
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  {
    {B∈(h, k, v)I(r,h)r * dcaps(k, r, i)
    * niceNode(N, k, v, r, h)
    ∧ N = node(-, k0, p, D, k', p') ∧ k0 = -∞}
    while(isLeaf(N) = false) {
      cur := next(N, k);
      N := get(cur);
    }
    {
    {B∈(h, k, v)I(r,h)r * dcaps(k, r, i)
    * niceNode(N, k, v, r, h)
    ∧ N = leaf(-, k0, D, k', p') ∧ k0 < k}
    while(k > highValue(N)) {
      cur := next(N, k);
      N := get(cur);
    }
    {
    {B∈(h, k, v)I(r,h)r * dcaps(k, r, i)
    * niceNode(N, k, v, r, h)
    ∧ N = leaf(-, k', D, k'', -) ∧ k' < k ≤ k''}
    if(isIn(N, k)) {
      {
        {B∈(h, k, v)I(r,h)r * dcaps(k, r, i)
        * niceNode(N, k, v, r, h)
        ∧ N = leaf(-, k', D, k'', -) ∧ (k, v) ∈ D}
        return(lookup(N, k));
      }
    } else {
      {false}
      return nil;
    }
  }
  {B∈(h, k, v)I(r,h)r * dcaps(k, r, i) ∧ ret = v}
}
{indef(h, k, v)i ∧ ret = v}

```

Figure 20. Proof outline for the B^{Link} tree search.

pair (k, v) must be stored in N – this ensures that the search will return the correct value.

Assertions in the proof must be stable – that is, invariant under interference from other threads. The stability of `niceNode` is ensured by the fact that the capabilities held by the thread do not allow nodes to be removed, the minimum values of nodes to change, or key k to be changed.

In Figures 21 and 22, we give a sketch proof of the following specification for the implementation of `remove`:

$$\{indef(h, k, v)_1\} \text{remove}(h, k) \{out_{def}(h, k)_1\}$$

The proof uses the predicates defined in Figure 18. The definition of `stLf` asserts that the node descriptor N contains legitimate information about the tree and contains the same information as the node n in the tree. If a leaf node into

```

{indef(h, k, v)1}
remove(h, k) {
  {B∈(h, k, v)I(r,h)r * dcaps(k, r, 1)}
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  {
    {B∈(h, k, v)I(r,h)r * dcaps(k, r, 1) * niceNode(N, k, v, r, h)
    * present(cur, k0, p) ∧ N = node(-, k0, p, D, k', p')
    ∧ k0 = -∞}
    while (isLeaf(N) = false) {
      cur := next(N, k);
      N := get(cur);
    }
    {
    {B∈(h, k, v)I(r,h)r * dcaps(k, r, 1) * niceNode(N, k, v, r, h)
    * present(cur, k', nil) ∧ N = leaf(-, k', D, k'', p')
    ∧ k' < k}
    while (true) {
      // see Figure 22
    }
  }
}
{out_{def}(h, k)1}

```

Figure 21. Proof outline for B^{Link} tree remove (excluding loop body).

whose range the key k falls, then the key-value pair (k, v) must be stored in.

A bug in the B^{Link} tree algorithm. While verifying the algorithm, we discovered a subtle bug in the original presentation [19]. The bug can occur during an `insert`, when a thread splits a tree node which itself was the result of another thread splitting the tree root. In order to insert the new node into the tree, the first thread will look in the prime block for the node’s parent. However, the second thread might not yet have written a pointer to the new root, resulting in an invalid dereference. Our solution was to require that a thread splitting the current the root locks the new node. A thread trying to insert must wait until the creation of the root is complete. A detailed trace exhibiting this bug can be found in appendix B

7. Conclusions

We have proposed a simple, abstract specification for reasoning about concurrent indexes. We have demonstrated the versatility of our specification, verifying a representative range of client applications ranging from common programming patterns such as memoization and `map`, to algorithms such as a prime number sieve. We have demonstrated that our particular choice of index specification is satisfied by three radically different concurrent implementations, based on simply linked lists, hash tables, and Sagiv’s complex and highly concurrent B^{Link} trees respectively.

```


$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{niceNode}(N, k, v, r, h) \\ * \text{present}(\text{cur}, k', \text{nil}) \wedge N = \text{leaf}(-, k', D, k'', p') \wedge k' < k \end{array} \right\}$$

lock(cur); // use LOCK
N := get(cur);

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{stLf}(\text{cur}, N, k, v, r, h) \\ * [\text{UNLOCK}(\text{cur})]_1^r \wedge N = \text{leaf}(1, k', D, k'', p') \wedge k' < k \end{array} \right\}$$

if (isIn(N, k)) {

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{stLf}(\text{cur}, N, k, v, r, h) \\ * [\text{UNLOCK}(\text{cur})]_1^r \wedge N = \text{leaf}(1, k', D, k'', p') \\ \wedge k' < k \leq k'' \wedge (k, -) \in D \end{array} \right\}$$

// use REM

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r \\ * [\text{MODLR}(0, \text{cur}, k, 1)]_1^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(0, k, v)]_{(d, 1)}^r \\ * \text{stLf}(\text{cur}, N, k, v, r, h) \wedge N = \text{leaf}(1, k', D, k'', p') \\ \wedge k' < k \leq k'' \wedge (k, -) \in D \end{array} \right\}$$

removePair(N, k);
put(A, cur); // use MODLR

$$\left\{ \begin{array}{l} \boxed{B_{\notin}(h, k)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{stLf}(\text{cur}, N, k, v, r, h) \\ * [\text{UNLOCK}(\text{cur})]_1^r \wedge N = \text{leaf}(1, k', D, k'', p') \\ \wedge k' < k \leq k'' \wedge D = D' \cup (k, -) \end{array} \right\}$$

unlock(cur); // use UNLOCK

$$\left\{ \begin{array}{l} \boxed{B_{\notin}(h, k)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) \end{array} \right\}$$

outdef(h, k)1
return;
} else {

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{stLf}(\text{cur}, N, k, v, r, h) \\ * [\text{UNLOCK}(\text{cur})]_1^r \wedge N = \text{leaf}(1, k', D, k'', p') \wedge k'' < k \end{array} \right\}$$

unlock(cur); // use UNLOCK

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{niceNode}(N, k, v, r, h) \\ * \text{present}(\text{cur}, k', \text{nil}) \wedge N = \text{leaf}(-, k', D, k'', p') \\ \wedge k'' < k \end{array} \right\}$$

if (k > highValue(N)) {

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{niceNode}(N, k, v, r, h) \\ * \text{present}(\text{cur}, k', \text{nil}) \wedge N = \text{leaf}(-, k', D, k'', p') \\ \wedge k'' < k \end{array} \right\}$$

while (k > highValue(N)) {
cur := next(N, k);
N := get(cur);
}

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{niceNode}(N, k, v, r, h) \\ * \text{present}(\text{cur}, k', \text{nil}) \wedge N = \text{leaf}(-, k', D, k'', p') \\ \wedge k' < k \end{array} \right\}$$

} else { // value is not in the tree
{false}
outdef(h, k)1
return;
}

$$\left\{ \begin{array}{l} \boxed{B_{\in}(h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) * \text{niceNode}(N, k, v, r, h) \\ * \text{present}(\text{cur}, k', \text{nil}) \wedge N = \text{leaf}(-, k', D, k'', p') \\ \wedge k' < k \end{array} \right\}$$

}

```

Figure 22. Proof outline for B^{Link} tree remove (main loop body).

Relationship to linearizability. Linearizability [11] is the current de-facto correctness criterion for concurrent algorithms. It requires that the methods of concurrent objects behave as atomic operations, thus providing a proof technique for observational refinement [10]. We could employ linearizability, or other atomicity refinement techniques such as [21], as a proof technique for verifying that implementations meet our abstract specification: an implementation that meets the sequential specification of an index and whose operations behave atomically can easily be shown to meet the concurrent specification. However, this simply shifts the proof burden; our approach is able to verify clients and implementations in a single coherent proof system.

While linearizability assures that index operations behave atomically, our abstract specification makes no such guarantee. Instead, our client proofs enforce abstract constraints on the possible interactions between threads, such as only allowing removals on a certain key. Consequently, while all linearizable indexes can be shown to implement our specification, our specification also admits implementations that are not linearizable. For instance, an index that implemented removal by performing the operation twice in succession could meet our specification, but would not be linearizable. As a more realistic example, consider the program:

$$\begin{array}{l} \text{insert}(k, 0) \parallel \text{insert}(k, 1) \\ x := \text{search}(k) \parallel y := \text{search}(k) \end{array}$$

Linearizability ensures that, after executing the program, the variables x and y will be equal (one or the other insert must come first, and a second insert has no effect). However, our specification does not ensure this. An implementation in which writes are cached, for instance, may satisfy our specification, but fail to provide this stronger guarantee. In practice, the strength of specifications is often traded against performance. We have shown how our approach can provide weak (§3) and strong (§4) specifications of concurrent behaviour. Our approach could therefore be seen as a flexible alternative to linearizability as a correctness criterion for concurrent programs.

Acknowledgments. Special thanks to Moshe Vardi for challenging us to verify Sagiv’s concurrent B^{Link} tree algorithm, and to Adam Wright for substantial contributions to the section on iteration (§5), and invaluable discussions and feedback overall. Thanks also to Bornat, Jones, Shapiro, and many researchers at Cambridge, Imperial and Queen Mary working on separation logic, for discussions and feedback. We acknowledge funding from an EPSRC DTA (da Rocha Pinto), EPSRC programme grant EP/H008373/1 (da Rocha Pinto, Dinsdale-Young, Gardner and Wheelhouse) and EPSRC grant EP/H010815/1 (Dodds).

References

- [1] BLELLOCH, G. E. Programming parallel algorithms. *Commun. ACM* 39 (March 1996), 85–97.

- [2] BOYLAND, J. Checking interference with fractional permissions. In *Static Analysis* (2003).
- [3] CALCAGNO, C., GARDNER, P., AND ZARFATY, U. Context Logic and tree update. In *POPL* (2005), ACM.
- [4] DA ROCHA PINTO, P. Reasoning about Concurrent Indexes. Master’s thesis, Imperial College London, Sept. 2010.
- [5] DILLIG, I., DILLIG, T., AND AIKEN, A. Precise reasoning for programs using containers. *SIGPLAN Not.* 46 (2011).
- [6] DINSDALE-YOUNG, T., DODDS, M., GARDNER, P., PARKINSON, M., AND VAPEIADIS, V. Concurrent abstract predicates. In *ECOOP* (2010).
- [7] DINSDALE-YOUNG, T., GARDNER, P., AND WHEELHOUSE, M. Abstraction and Refinement for Local Reasoning. In *VSTTE* (2010).
- [8] DODDS, M., FENG, X., PARKINSON, M., AND VAPEIADIS, V. Deny-guarantee reasoning. In *ESOP* (2009).
- [9] FENG, X., FERREIRA, R., AND SHAO, Z. On the relationship between concurrent separation logic and assume-guarantee reasoning. In *ESOP* (2007).
- [10] FILIPOVIC, I., O’HEARN, P., RINETZKY, N., AND YANG, H. Abstraction for concurrent objects. In *ESOP* (2010).
- [11] HERLIHY, M. P., AND WING, J. M. Linearizability: a correctness condition for concurrent objects. *ACM Trans. Program. Lang. Syst.* 12 (July 1990), 463–492.
- [12] HOARE, C. A. R. Proof of a structured program: ‘The sieve of Eratosthenes’. *The Computer Journal* 15, 4 (1972), 321–325.
- [13] KUNCAK, V., LAM, P., ZEE, K., AND RINARD, M. C. Modular pluggable analyses for data structure consistency. *IEEE Trans. Softw. Eng.* 32 (December 2006), 988–1005.
- [14] MALECHA, G., MORRISETT, G., SHINNAR, A., AND WISNESKY, R. Toward a verified relational database management system. In *POPL* (2010).
- [15] O’HEARN, P. W. Resources, concurrency, and local reasoning. *Theor. Comput. Sci.* 375 (April 2007), 271–307.
- [16] PARKINSON, M., AND BIERMAN, G. Separation logic and abstraction. In *POPL* (2005).
- [17] PHILIPPOU, A., AND WALKER, D. A process-calculus analysis of concurrent operations on b-trees. *J. Comput. Syst. Sci.* 62, 1 (2001), 73–122.
- [18] REYNOLDS, J. Separation logic: a logic for shared mutable data structures. In *LICS* (2002).
- [19] SAGIV, Y. Concurrent operations on B*-trees with overtaking. *Journal of Computer and System Sciences* 33 (October 1986).
- [20] SEXTON, A., AND THIELECKE, H. Reasoning about B+ trees with operational semantics and separation logic. *ENTCS* 218 (2008).
- [21] TURON, A. J., AND WAND, M. A separation logic for refining concurrent objects. In *POPL* (2011).
- [22] VAPEIADIS, V., AND PARKINSON, M. A marriage of rely/guarantee and separation logic. *CONCUR* (2007).

A. B^{Link} Tree Implementation Details

In this appendix we provide an in depth discussion of our B^{Link} tree index implementation, based on Sagiv’s BTree algorithms, and how this implementation satisfies our abstract specification. We give concrete interpretations to each of our abstract predicates and define the interference environment for the B^{Link} tree. Together these allow us to prove the correctness of our implementation.

B^{Link} Tree Data Structure

To begin the verification of our B^{Link} tree implementation, we first define a series of predicates representing the concrete B^{Link} tree data structure. There are two types of node in a B^{Link} tree: *leaf* nodes and *inner* nodes. Leaf nodes are at the fringe of the structure and contain the key-value pairs from the abstract interface. Inner nodes make up the rest of the tree and contain key-pointer pairs that provide the search structure of the tree. We assume two basic predicates for representing these nodes in the tree: a leaf predicate, and an inner predicate.

$$x \mapsto \text{leaf}(l, k_0, D, k', p') \quad x \mapsto \text{inner}(l, k_0, p, D, k', p')$$

Here, x is the address of the node. The value l is the node’s lock. If the node is unlocked then $l = 0$ and if the node is locked then $l = 1$. The ordered list D contains the key-value pairs (k, v) represented by the node. In each node the list D may contain up to $2K$ key-value pairs for some fixed constant K given by the implementation (K is often chosen so that a node fills a single page in memory). The values k_0 and k' are the lower and upper bound, respectively, on the keys contained in this list. So, for every key-value pair $(k, v) \in D$ we have $k_0 < k \leq k'$. The pointer p (only present in an inner node) points to the subtree which contains all of the keys which are greater than the minimum value of this node. The pointer p' , known as the link pointer, points to the node’s right sibling, if it exists.

We define some additional notation for handling lists. We require a notion of iterated concatenation which we denote

$$\bigcup_{i=1}^n D_i = D_1 :: D_2 :: \dots D_n.$$

We also require an insertion operation $D \uplus (k, v)$ which adds the key-value pair (k, v) to the ordered list D in the correct place,

$$D \uplus (k, v) = D_1 :: (k, v) :: D_2$$

where $D = D_1 :: D_2$ and $D_1 = D'_1 :: (k_1, v_1)$ and $D_2 = (k_2, v_2) :: D'_2$ and $k_1 < k < k_2$ (undefined otherwise).

A B^{Link} tree is a superimposed structure made up of both a tree and several layers of linked lists. At the leaf level the linked list contains pointers to data entries, while at other levels the linked list contains pointers to nodes deeper in the tree structure. These linked lists always have at least one element, the first node has minimum value $-\infty$ and the last node has maximum value $+\infty$. Each node in these linked

lists is disjoint, so we can use a separation logic predicate to define this structure precisely. Given ordered key-value list T and D , which contain the key-value pairs that point into the current level of the tree and the key-value pairs contained in the current level of the tree respectively, we can define the linked list structure for a layer of the B^{Link} tree. Let $T = [(k_1, v_1), \dots, (k_n, v_n)]$ then,

$$\begin{aligned} \text{leafList}(T, D) &\triangleq \exists D_1, \dots, D_n. \\ &\quad \bigotimes_{i=1}^{n-1} v_i \mapsto \text{leaf}(-, k_i, D_i, k_{i+1}, v_{i+1}) \\ &\quad * v_n \mapsto \text{leaf}(-, k_n, D_n, +\infty, \text{nil}) \\ &\quad \wedge k_1 = -\infty \wedge n > 0 \wedge D = \bigcup_{i=1}^n D_i \end{aligned}$$

$$\begin{aligned} \text{innerList}(T, D) &\triangleq \exists D_1, \dots, D_n. \\ &\quad \bigotimes_{i=1}^{n-1} v_i \mapsto \text{inner}(-, k_i, v'_i, D_i, k_{i+1}, v_{i+1}) \\ &\quad * v_n \mapsto \text{inner}(-, k_n, v'_n, D_n, +\infty, \text{nil}) \\ &\quad \wedge k_1 = -\infty \wedge n > 0 \\ &\quad \wedge D = \bigcup_{i=1}^n (k_i, v'_i) \cup D_i \end{aligned}$$

We choose not to define the tree structure directly, as at some points in time the tree structure of the B^{Link} tree can actually be broken by the `insert` operation. When the `insert` operation creates a new node in the tree, it is added to the linked list structure before it is given a reference in the layer above. If the `search` operation did not use the link pointers as well as the tree pointers, it would not be able to find this new node at this point in time. To capture this behaviour we instead choose to build up our tree predicate by layering our lists on top of one another. Using our linked list predicates, we can build up a predicate for the tree-like structure of the B^{Link} tree.

$$\begin{aligned} \text{Btree}_1(PB, x :: T, D) &\triangleq \exists p. \text{leafList}(x :: T, D) \\ &\quad \wedge x = (-\infty, p) \wedge PB = [p] \end{aligned}$$

$$\begin{aligned} \text{Btree}_{n+1}(p :: PB, x :: T, D) &\triangleq \exists L, L'. \text{innerList}(x :: T, L) \\ &\quad * \text{Btree}_n(PB, L', D) \\ &\quad \wedge x = (-\infty, p) \wedge L \subseteq L' \end{aligned}$$

The prime block PB contains a list of pointers to the leftmost node at each level of the tree. The key-value list D is the concatenation of all key-value pairs at the fringe of the tree and corresponds to our abstract index view of the B^{Link} tree structure.

Finally, using these predicates, we can now define a predicate for the complete B^{Link} tree structure.

$$\text{BLTree}(h, D) \triangleq \exists PB, n, T. \text{Btree}_n(PB, T, D) * h \mapsto PB$$

This describes a B^{Link} tree whose prime block is stored at address h and contains a set of key-value pairs D . Figure 16 shows an example of a B^{Link} tree. The fringe of the tree forms a `leafList` that contains all of the key-value pairs mapped to by the index. Each of the other layers of the tree forms an `innerList` that makes up the search structure of the tree. Each list has minimum value $-\infty$ and maximum value ∞ and the primeblock points to the head of each layer's list.

Interpretation of Abstract Predicates

Now that we have a predicate describing a B^{Link} tree, we can turn our attention to providing concrete interpretations of our abstract predicates. In § 6.3 we introduced the interpretation of the $\text{in}_{\text{def}}(h, k, v)$ predicate. Here we go into more detail about the auxiliary predicates we used in our interpretations, and then provide the concrete interpretations of the full abstract specification.

First we define a number of predicates which will come in useful for our later definitions:

$$\begin{aligned} \diamond P &\triangleq \text{true} * P \\ \text{isNode}(x, l) &\triangleq \exists k_0, p, D, k', p'. \\ &\quad x \mapsto \text{node}(l, k_0, p, D, k', p') \\ \text{locked}(x) &\triangleq \text{isNode}(x, 1) \\ \text{unlocked}(x) &\triangleq \text{isNode}(x, 0) \\ \text{child}(h, x) &\triangleq \exists l. \text{isNode}(x, l) \\ &\quad \wedge \exists p, ps. \diamond h \mapsto p : ps \\ &\quad \wedge x = p \\ &\quad \vee p \mapsto \text{node}(-, -, -, -, -, x) \\ &\quad \vee \\ &\quad \exists y, k_0, v_0, D, k. \\ &\quad \diamond y \mapsto \text{inner}(-, k_0, v_0, D, -, -) \\ &\quad \wedge (k, x) \in (k_0, v_0) :: D \\ \text{orphan}(h, x) &\triangleq \neg \text{child}(h, x) \\ \text{dualRoot}(h, x, y) &\triangleq \exists p, D, k, p', D'. h \mapsto x : xs \\ &\quad * x \mapsto \text{node}(1, -\infty, p, D, k, y) \\ &\quad * y \mapsto \text{node}(1, k, p', D', \infty, \text{nil}) \\ \text{allMods}(x) &\triangleq \forall l, t, k, v, y, i. \text{isNode}(x, l) \\ &\quad \wedge \diamond [\text{MODLR}(t, x, k, i)]_1^r \\ &\quad \wedge \diamond [\text{MODLI}(t, x, k, v, i)]_1^r \\ &\quad \wedge \diamond [\text{FIX}(k, x)]_1^r \\ &\quad \wedge \diamond [\text{MODII}(x, k, y)]_1^r \\ &\quad \wedge \diamond [\text{NEW}(x, k, y)]_1^r \end{aligned}$$

Informally, these predicates have the following meanings:

- $\diamond P$ describes a heap where P is satisfied somewhere in the heap.
- $\text{isNode}(x, l)$ describes a node x in the B^{Link} tree with lock value l .
- $\text{locked}(x)$ describes a locked node x in the B^{Link} tree.
- $\text{unlocked}(x)$ describes an unlocked node x in the B^{Link} tree.
- $\text{child}(h, x)$ describes a node x in the B^{Link} tree at address h which is either at the root level, or has a parent in the tree's search structures; some node in the tree contains a key-value pair $(-, x)$.
- $\text{orphan}(h, x)$ describes a node x in the B^{Link} tree at address h which does not have a parent in the tree's search structure; it is not at the root and no node contains a key-value pair $(-, x)$.

- $\text{dualRoot}(h, x, y)$ describes a B^{Link} tree at address h that currently has two nodes at its root level (so an insert operation has just split the root and is about to create a new one).
- $\text{allMods}(x)$ describes the set of all modification capabilities, with exclusive permission, for node x .

As we saw in § 6.3, our concrete interpretations describe the shared state with one of the following assertions:

$$B_{\in}(h, k, v) \triangleq \exists D. \text{BLTree}(h, D) * \neg \exists x, l. \diamond \text{isNode}(x, l) \wedge (k, v) \in D \wedge \text{Tokens}(h)$$

$$B_{\notin}(h, k) \triangleq \exists D. \text{BLTree}(h, D) * \neg \exists x, l. \diamond \text{isNode}(x, l) \wedge k \notin \text{keys}(D) \wedge \text{Tokens}(h)$$

The assertion $B_{\in}(h, k, v)$ describes a B^{Link} tree at address h that contains the key-value pair (k, v) . Similarly the assertion $B_{\notin}(h, k, v)$ describes a B^{Link} tree at address h where the key k is unassigned. However, both assertions also describe an additional part of the shared state. The assertion $\neg \exists x, l. \diamond \text{isNode}(x, l)$ ensures that there are no nodes in this additional state; it consists only of capabilities. The assertion $\text{Tokens}(h)$ ensures that these capabilities are consistent with the current state of the B^{Link} tree at address h .

The $\text{Tokens}(h)$ predicate is quite complex and is defined in Figure 23. The predicate describes the capabilities that are in the shared state on a capability by capability basis dependent on the current state of the B^{Link} tree. The predicate is built up of the conjunction of a number of disjuncts.

The first disjunct describes if a node's $[\text{UNLOCK}(x)]_1^r$ capability is present in the shared state. If x is not a node, or if x is an unlocked node then the UNLOCK capability must be present in the shared state. If x is a locked node then this capability may be missing from the shared state. However, it is also possible that the thread that has locked the node may have acquired a MOD capability for that node, that is it is about to make some change to the node. In this case the UNLOCK capability will be present in the shared state, but so will some REM or INS capability. This may appear to allow some other thread to acquire the UNLOCK capability for this node, but recall that the node is still locked. We shall see later, when we define the interference environment, that a thread may only acquire a nodes UNLOCK capability if that node is unlocked and in doing so the thread locks the node.

The second and third disjuncts describe if we are in an action tracking state or not. If $t = 0$, then we are not tracking the actions on this key (we are in a def or unk environment) and all of the REM and INS capabilities for $t = 1$ must be in the shared state. If $t = 1$, then we are tracking the actions on this key (we are in an ins or rem environment) and all of the REM and INS capabilities for $t = 0$ must be in the shared state. We shall see later, when we define the interference environment, that when we are tracking the actions on a key, threads leave behind some fraction of their REM or INS capabilities after performing a modification action. This

allows us to track if a value has been inserted or removed from a given key value and return to a def state.

The fourth and final disjunct describes which of the modification capabilities are present in the shared state for each node in the B^{Link} tree. It is always the case that either all of the modification capabilities are in the shared state, or one such capability is missing. If one of the modification capabilities is missing then the node must be locked and the locking thread must have placed the UNLOCK capability and some other capability, describing the action it is about to perform on that node (e.g. REM or INS). This represents a thread that has locked the node and is about to make some update to that node. Due to the locking, it is only ever possible for at most one thread to be in this state, hence why at most one modification capability is ever missing for any given node.

We define the concrete interpretations of our abstract predicates in Figure 24. Each case describes the current state of the shared B^{Link} tree, as well as which capabilities are known to be in the shared and thread local state. For example, the definition of $\text{in}_{\text{def}}(h, k, v)_i$ states that the key-value pair (k, v) must be stored in the tree. Notice that this definition also gives the thread deny permission (d, i) on all REM and INS capabilities for k . When $i \in (0..1)$ no thread is able to modify the value of k in the tree, and when $i = 1$ only the current thread may modify the value of k in the tree, so this assertion is self-stable. The thread also has the $[\text{LOCK}]_g^r$ capability, which allows it to lock nodes in the tree, and the $[\text{SWAP}]_g^r$ capability, which allows it to change between tracking actions or not (by swapping $t = 0$ and $t = 1$ capabilities).

Some of the other definitions make more complicated assertions about the shared state. Take, for example, the definition of the $\text{in}_{\text{rem}}(h, k, v)_i$ predicate. Recall from our abstract specification that this predicate states that key k was assigned value v , but that any thread can remove this value. We track which actions have occurred so far by using the $t = 1$ capabilities. If a thread removes the value for the key, then it must leave some $[\text{REM}(1, k)]_{(g, i')}^r$ capability in the shared state. The uncertainty about the current assignment of k is represented by the disjunction in the shared state. In the first case no thread has yet removed the key from the tree, since there is no REM capability for that k in the shared state. In the second case some thread has just acquired the modification capability $[\text{MODLR}(1, x, k, i')]_1^r$ allowing it to remove the key from the tree, but it has yet to perform this action, so the key is still currently assigned. In the last case some thread has removed the key from the tree and left part of its REM capability in the shared state to signify this.

The other predicates are defined in similar ways.

We can now verify that our interpretations satisfy the axioms from §4 for our abstract specification. For example we can verify,

$$\text{in}_{\text{rem}}(h, k, v)_i * \text{out}_{\text{rem}}(h, k)_j \implies \text{out}_{\text{rem}}(h, k)_{i+j}$$

$$\begin{aligned}
\text{Tokens}(h) &\triangleq \\
&\forall x. \left(\begin{array}{l} \neg \exists l. \Diamond \text{isNode}(x, l) \\ \wedge \Diamond [\text{UNLOCK}(x)]_1^r \end{array} \right) \vee \left(\begin{array}{l} \Diamond \text{unlocked}(x) \\ \wedge \Diamond [\text{UNLOCK}(x)]_1^r \end{array} \right) \vee \left(\begin{array}{l} \Diamond \text{locked}(x) \\ \wedge \neg \Diamond [\text{UNLOCK}(x)]_1^r \end{array} \right) \\
&\quad \vee \left(\begin{array}{l} \Diamond \text{locked}(x) \wedge \Diamond [\text{UNLOCK}(x)]_1^r \\ \wedge \exists k, v, i, t. \left(\begin{array}{l} [\text{REM}(t, k)]_i^r \\ \wedge \neg \Diamond [\text{MODLR}(t, x, k, i)]_1^r \end{array} \right) \vee \left(\begin{array}{l} [\text{INS}(t, k, v)]_i^r \\ \wedge \neg \Diamond [\text{MODLI}(t, x, k, v, i)]_1^r \end{array} \right) \end{array} \right) \\
&\quad \wedge \\
&\quad \forall k. \Diamond [\text{REM}(0, k)]_1^r \vee \Diamond [\text{REM}(1, k)]_1^r \\
&\quad \wedge \\
&\quad \forall k, v. \Diamond [\text{INS}(0, k, v)]_1^r \vee \Diamond [\text{INS}(1, k, v)]_1^r \\
&\quad \wedge \\
&\quad \forall x. \text{allMods}(x) \vee \exists t, k, i. \Diamond [\text{REM}(t, k)]_i^r \wedge \Diamond [\text{UNLOCK}(x)]_1^r \wedge ([\text{MODLR}(t, x, k, i)]_1^r \neg * \text{allMods}(x)) \\
&\quad \vee \exists t, k, v, i. \Diamond [\text{INS}(t, k, v)]_i^r \wedge \Diamond [\text{UNLOCK}(x)]_1^r \wedge ([\text{MODLI}(t, x, k, v, i)]_1^r \neg * \text{allMods}(x)) \\
&\quad \vee \exists k, y. \text{orphan}(h, x) \wedge \Diamond [\text{MODII}(y, k, x)]_1^r \wedge ([\text{FIX}(k, x)]_1^r \neg * \text{allMods}(x)) \\
&\quad \vee \exists k, y. \text{orphan}(h, y) \wedge \Diamond [\text{FIX}(k, y)]_1^r \wedge [\text{UNLOCK}(x)]_1^r \wedge ([\text{MODII}(x, k, y)]_1^r \neg * \text{allMods}(x)) \\
&\quad \vee \exists k, y. \text{dualRoot}(h, x, y) \wedge \Diamond [\text{UNLOCK}(x)]_1^r \wedge \Diamond [\text{UNLOCK}(y)]_1^r \wedge ([\text{NEW R}(x, k, y)]_1^r \neg * \text{allMods}(x))
\end{aligned}$$

Figure 23. Definition of the $\text{Tokens}(h)$ predicate.

$$\begin{aligned}
\text{in}_{\text{def}}(h, k, v)_i &\triangleq \exists r. \boxed{\text{B}_{\in}(h, k, v)}_{I(r, h)}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(0, k)]_{(d, i)}^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(0, k, v)]_{(d, i)}^r \\
\text{out}_{\text{def}}(h, k)_i &\triangleq \exists r. \boxed{\text{B}_{\notin}(h, k)}_{I(r, h)}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(0, k)]_{(d, i)}^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(0, k, v)]_{(d, i)}^r \\
\text{in}_{\text{ins}}(h, k, S)_i &\triangleq \exists v \in S, r, i', i''. \boxed{\text{B}_{\in}(h, k, v) \wedge \Diamond [\text{INS}(1, k, v)]_{i''}^r}_{I(r, h)}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(1, k)]_{(d, i)}^r \\
&\quad * \bigotimes_{v \in S} [\text{INS}(1, k, v)]_{(g, i')}^r * \bigotimes_{v \notin S} [\text{INS}(1, k, v)]_{(d, i)}^r \wedge (i' = i \vee i' + i'' = i) \\
\text{out}_{\text{ins}}(h, k, S)_i &\triangleq \exists v \in S, r, i'. \boxed{\begin{array}{l} \text{B}_{\notin}(h, k) \wedge \neg \Diamond [\text{INS}(1, k, v)]_{(g, i')}^r \\ \vee \text{B}_{\notin}(h, k) \wedge \Diamond [\text{INS}(1, k, v)]_{(g, i')}^r \wedge \Diamond [\text{UNLOCK}(x)]_1^r \wedge \neg \Diamond [\text{MODLI}(1, x, k, v, i')]_1^r \\ \vee \text{B}_{\in}(h, k, v) \wedge \Diamond [\text{INS}(1, k, v)]_{(g, i')}^r \wedge \Diamond [\text{MODLI}(1, x, k, v, i')]_1^r \end{array}}_{I(r, h)}^r \\
&\quad * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(1, k)]_{(d, i)}^r * \bigotimes_{v \in S} [\text{INS}(1, k, v)]_{(g, i)}^r \\
&\quad * \bigotimes_{v \notin S} [\text{INS}(1, k, v)]_{(d, i)}^r \wedge i' > 0 \\
\text{in}_{\text{rem}}(h, k, v)_i &\triangleq \exists r, i'. \boxed{\begin{array}{l} \text{B}_{\in}(h, k, v) \wedge \neg \Diamond [\text{REM}(1, k)]_{(g, i')}^r \\ \vee \text{B}_{\in}(h, k, v) \wedge \Diamond [\text{REM}(1, k)]_{(g, i')}^r \wedge \Diamond [\text{UNLOCK}(x)]_1^r \wedge \neg \Diamond [\text{MODLR}(1, x, k, i')]_1^r \\ \vee \text{B}_{\notin}(h, k) \wedge \Diamond [\text{REM}(1, k)]_{(g, i')}^r \wedge \Diamond [\text{MODLR}(1, x, k, i')]_1^r \end{array}}_{I(r, h)}^r \\
&\quad * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(1, k)]_{(g, i)}^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(1, k, v)]_{(d, i)}^r \wedge i' > 0 \\
\text{out}_{\text{rem}}(h, k)_i &\triangleq \exists r, i', i''. \boxed{\text{B}_{\notin}(h, k) \wedge \Diamond [\text{REM}(1, k)]_{(g, i'')}^r}_{I(r, h)}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(1, k)]_{(g, i')}^r \\
&\quad * \bigotimes_{v \in \text{Vals}} [\text{INS}(1, k, v)]_{(d, i)}^r \wedge (i' = i \vee i' + i'' = i) \\
\text{unk}(h, k, S)_i &\triangleq \exists v \in S, r. \boxed{\text{B}_{\in}(h, k, v) \vee \text{B}_{\notin}(h, k)}_{I(r, h)}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(0, k)]_{(g, i)}^r \\
&\quad * \bigotimes_{v \in S} [\text{INS}(0, k, v)]_{(g, i)}^r * \bigotimes_{v \notin S} [\text{INS}(0, k, v)]_{(d, i)}^r \\
\text{read}(h, k) &\triangleq \exists v, r. \boxed{\text{B}_{\in}(h, k, v) \vee \text{B}_{\notin}(h, k)}_{I(r, h)}^r
\end{aligned}$$

Figure 24. Concrete predicate interpretations for the B^{Link} tree implementation.

since the assertion on the shared state from the out_{rem} predicate collapses the disjunction in the shared state from the in_{rem} predicate into just one matching case, and the thread local capabilities sum together as expected.

Describing Interference

We model the possible interference on the shared state by an interference environment $I(r, h)$. The interference environment is made up of a set of actions that can be performed by the current thread, and other threads, so long as they possess sufficient resources and capabilities for the actions.

First, we introduce some additional predicates which will help us describe the actions in our interference environment. We have a node predicate for when we want to talk about a node of arbitrary type (a leaf node or an inner node).

$$x \mapsto \text{node}(l, k_0, p, D, k', p') \triangleq \begin{aligned} &(p = \text{nil} \wedge x \mapsto \text{leaf}(l, k_0, D, k', p')) \\ &\vee (p \neq \text{nil} \wedge x \mapsto \text{inner}(l, k_0, p, D, k', p')) \end{aligned}$$

We also have a root predicate which describes if a node is the root of the B^{Link} tree or not.

$$\text{root}(h, x) \triangleq \exists xs. \diamond h \mapsto PB \wedge PB = x :: xs$$

When the insertion operations tries to split a node (when adding a pair to full node) it is important to know if that node is the root or not. If the root is split, then a new root needs to be created and the prime block updated accordingly.

Finally, when describing the insertion action for inner nodes we require a notion of a list of nodes up to some point.

$$\text{nodeList}(p, N, p') \triangleq (N = [] \wedge p = p' \wedge \text{emp}) \vee \left(\begin{aligned} &\exists l, k_0, p_0, D, k_1, p_1, N'. \\ &N = (l, k_0, p_0, D, k_1, p_1) :: N' \\ &\wedge p \mapsto \text{node}(l, k_0, p_0, D, k_1, p_1) \\ &* \text{nodeList}(p_1, N', p') \end{aligned} \right)$$

The actions that make up the interference environment for the B^{Link} tree implementation are given in Figure 25. The LOCK and $\text{UNLOCK}(x)$ actions control the locking and unlocking of nodes in the tree. The $\text{INS}(t, k, v)$, $\text{REM}(t, k)$ and $\text{FIX}(k, y)$ actions allow a thread to gain the modification tokens for a node that they have locked. The SWAP action allows a thread with full permission for some key change if we are tracking actions for that key. The $\text{MODLI}(t, x, k, v, i)$ action allows a thread to insert a key-value pair (k, v) into some leaf node x . If this node was full, then the thread is given the $[\text{FIX}(k, y)]_1^r$ capability so that it may repair the search structure of the tree. The $\text{MODLR}(t, x, k, i)$ action allows a thread to remove a key-value pair $(k, -)$ from some leaf node x . Notice that there is no way for a thread to remove key-value pairs from inner nodes. The $\text{MODII}(x, k, y)$ action allows a thread to insert a key-value pair (k, y) into some inner node x . This action is used to repair the search structure of the tree after a node has been

split. The $\text{NEWR}(x, k, y)$ action allows a thread to create a new root, and update the prime block accordingly, after the old root has been split. This action can only be used if the thread has previously split the old root and thus acquired the $[\text{NEW}(x, k, y)]_1^r$ capability.

Verifying the Operations

Our B^{Link} tree implementation uses a language which includes a set of heap update commands, which directly modify nodes in the shared heap, and a set of store update commands, which work with nodes but do not manipulate the shared state. We assume that variables in a thread's local store can contain integer, pointer, Boolean, stack and node content information.

The heap update commands are:

```
lock(x)
unlock(x)
x := new()
N := get(x)
put(N, x)
PB := getPrimeBlock(h)
putPrimeBlock(h, PB)
```

Since these commands update the shared state, it is necessary that they each behave atomically so that they do not interfere with one another.

The store update commands are:

```
k := lowValue(N)
k := highValue(N)
x := next(N, k)
x := lookup(N, k)
addPair(N, k, v)
removePair(N, k)
M := rearrange(N, k, v, x)
x := root(PB)
N := newRoot(k', p, k, v, k'')
addRoot(PB, x)
x := getNodeLevel(PB, i)

b := isSafe(N)
b := isIn(N, k)
b := isLeaf(N)
b := isRoot(PB, x)

stack := newStack()
push(stack, x)
x := pop(stack)
b := isEmpty(stack)
```

The store update commands only modify the local store of a thread, so it is not necessary for these commands to be atomic.

We assume that these commands satisfy the specifications given in Figure 26 and Figure 27. Our proof of the search operation given in §6.3 Figure 20 then follows. The other

$$\begin{aligned}
\text{LOCK} &: x \mapsto \text{node}(0, k_0, p, D, k', p') * [\text{UNLOCK}(x)]_1^r \rightsquigarrow x \mapsto \text{node}(1, k_0, p, D, k', p') \\
\text{UNLOCK}(x) &: x \mapsto \text{node}(1, k_0, p, D, k', p') \rightsquigarrow x \mapsto \text{node}(0, k_0, p, D, k', p') * [\text{UNLOCK}(x)]_1^r \\
\text{INS}(t, k, v) &: [\text{MODLI}(t, x, k, v, i)]_1^r \rightsquigarrow [\text{INS}(t, k, v)]_{(g,i)}^r * [\text{UNLOCK}(x)]_1^r \\
\text{REM}(t, k) &: [\text{MODLR}(t, x, k, i)]_1^r \rightsquigarrow [\text{REM}(t, k)]_{(g,i)}^r * [\text{UNLOCK}(x)]_1^r \\
\text{FIX}(k, y) &: [\text{MODII}(x, k, y)]_1^r \rightsquigarrow [\text{FIX}(k, y)]_1^r * [\text{UNLOCK}(x)]_1^r \\
\text{SWAP} &: \left\{ \begin{array}{l} [\text{REM}(1, k)]_1^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(1, k, v)]_1^r \rightsquigarrow [\text{REM}(0, k)]_1^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(0, k, v)]_1^r \\ \left(\begin{array}{l} [\text{REM}(0, k)]_1^r * [\text{REM}(1, k)]_i^r \\ * \bigotimes_{v \in \text{Vals}} [\text{INS}(0, k, v)]_1^r \\ * \bigotimes_{v \in \text{Vals}} [\text{INS}(1, k, v)]_{i_v}^r \end{array} \right) \rightsquigarrow [\text{REM}(1, k)]_1^r * \bigotimes_{v \in \text{Vals}} [\text{INS}(1, k, v)]_1^r \end{array} \right. \\
\text{MODLI}(t, x, k, v, i) &: \left\{ \begin{array}{l} \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D, k', p') * [\text{UNLOCK}(x)]_1^r \\ * \left([\text{INS}(t, k, v)]_{(g,i)}^r \wedge t = 0 \vee \text{emp} \wedge t = 1 \right) \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D', k', p') \\ * [\text{MODLI}(t, x, k, v, i)]_1^r \\ \wedge D' = D \uplus (k, v) \end{array} \right) \\ \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D, k', p') * [\text{UNLOCK}(x)]_1^r \\ * \left([\text{INS}(t, k, v)]_{(g,i)}^r \wedge t = 0 \vee \text{emp} \wedge t = 1 \right) \\ * [\text{FIX}(k_1, y)]_1^r \wedge |D| = 2K \wedge \neg \text{root}(h, x) \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D_1, k_1, y) \\ * y \mapsto \text{leaf}(0, k_1, D_2, k', p') \\ * [\text{MODLI}(t, x, k, v, i)]_1^r \\ \wedge D_1 :: D_2 = D \uplus (k, v) \end{array} \right) \\ \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D, k', p') \\ * [\text{NEWR}(x, k_1, y)]_1^r \\ * \left([\text{INS}(t, k, v)]_{(g,i)}^r \wedge t = 0 \vee \text{emp} \wedge t = 1 \right) \\ \wedge |D| = 2K \wedge \text{root}(h, x) \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D_1, k_1, y) \\ * y \mapsto \text{leaf}(1, k_1, D_2, k', p') \\ * [\text{MODLI}(t, x, k, v, i)]_1^r \\ \wedge D_1 :: D_2 = D \uplus (k, v) \end{array} \right) \end{array} \right. \\
\text{MODLR}(t, x, k, i) &: \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D, k', p') * [\text{UNLOCK}(x)]_1^r \\ * \left([\text{REM}(t, k)]_{(g,i)}^r \wedge t = 0 \vee \text{emp} \wedge t = 1 \right) \\ \wedge (k, -) \in D \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{leaf}(1, k_0, D', k', p') \\ * [\text{MODLR}(t, x, k, i)]_1^r \\ \wedge D = D' \uplus (k, -) \end{array} \right) \\
\text{MODII}(x, k, y) &: \left\{ \begin{array}{l} \left(\begin{array}{l} x \mapsto \text{inner}(1, k_0, p, D, k', p') * [\text{UNLOCK}(x)]_1^r \\ * y \mapsto \text{node}(l, k, p_y, D_y, k'_y, p'_y) \\ * \text{nodeList}(p, N, y) \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{inner}(1, k_0, p, D', k', p') \\ * y \mapsto \text{node}(l, k, p_y, D_y, k'_y, p'_y) \\ * \text{nodeList}(p, N, y) \\ * [\text{MODII}(x, k, y)]_1^r \\ \wedge D' = D \uplus (k, y) \end{array} \right) \\ \left(\begin{array}{l} x \mapsto \text{inner}(1, k_0, p, D, k', p') * [\text{UNLOCK}(x)]_1^r \\ * y \mapsto \text{node}(l, k, p_y, D_y, k'_y, p'_y) \\ * \text{nodeList}(p, N, y) * [\text{FIX}(k_z, z)]_1^r \\ \wedge |D| = 2K \wedge \neg \text{root}(h, x) \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{inner}(1, k_0, p, D'_1, k_z, z) \\ * z \mapsto \text{inner}(0, k_z, p_z, D'_2, k', p') \\ * y \mapsto \text{node}(l, k, p_y, D_y, k'_y, p'_y) \\ * \text{nodeList}(p, N, y) \\ * [\text{MODII}(x, k, y)]_1^r \\ \wedge D'_1 :: (k_z, p_z) :: D'_2 = D \uplus (k, v) \end{array} \right) \\ \left(\begin{array}{l} x \mapsto \text{inner}(1, k_0, p, D, k', p') \\ * y \mapsto \text{node}(l, k, p_y, D_y, k'_y, p'_y) \\ * \text{nodeList}(p, N, y) * [\text{NEWR}(x, k_z, z)]_1^r \\ \wedge |D| = 2K \wedge \text{root}(h, x) \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{inner}(1, k_0, p, D'_1, k_z, z) \\ * z \mapsto \text{inner}(1, k_z, p_z, D'_2, k', p') \\ * y \mapsto \text{node}(l, k, p_y, D_y, k'_y, p'_y) \\ * \text{nodeList}(p, N, y) \\ * [\text{MODII}(x, k, y)]_1^r \\ \wedge D'_1 :: (k_z, p_z) :: D'_2 = D \uplus (k, v) \end{array} \right) \end{array} \right. \\
\text{NEWR}(x, k, y) &: \left(\begin{array}{l} x \mapsto \text{node}(1, -\infty, p_0, D_1, k, y) \\ * y \mapsto \text{node}(1, k, p, D_2, \infty, \text{nil}) \\ * [\text{UNLOCK}(x)]_1^r * [\text{UNLOCK}(y)]_1^r \\ * h \mapsto PB \wedge PB = x :: xs \end{array} \right) \rightsquigarrow \left(\begin{array}{l} x \mapsto \text{node}(1, -\infty, p_0, D_1, k, y) \\ * y \mapsto \text{node}(1, k, p, D_2, \infty, \text{nil}) \\ * z \mapsto \text{inner}(0, -\infty, x, [(k, y)], \infty, \text{nil}) \\ * [\text{NEWR}(x, k, y)]_1^r \\ * h \mapsto z :: PB \end{array} \right)
\end{aligned}$$

Figure 25. The interference environment for the B^{Link} tree implementation.

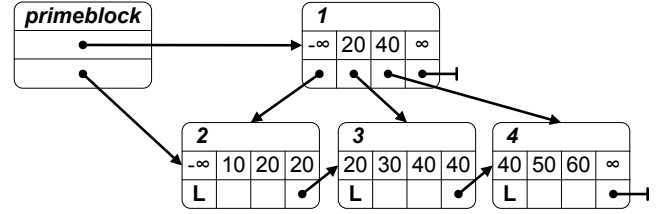
$\{x \mapsto \text{node}(0, k_0, p, D, k', p')\}$	<code>lock(x)</code>	$\{x \mapsto \text{node}(1, k_0, p, D, k', p')\}$
$\{x \mapsto \text{node}(1, k_0, p, D, k', p')\}$	<code>unlock(x)</code>	$\{x \mapsto \text{node}(0, k_0, p, D, k', p')\}$
$\{\text{emp}\}$	<code>x := new()</code>	$\{x \mapsto \text{node}(0, 0, \text{nil}, [], 0, \text{nil})\}$
$\{x \mapsto \text{node}(l, k_0, p, D, k', p')\}$	<code>N := get(x)</code>	$\left. \begin{array}{l} \{x \mapsto \text{node}(l, k_0, p, D, k', p')\} \\ \wedge N = \text{node}(l, k_0, p, D, k', p') \end{array} \right\}$
$\left. \begin{array}{l} \{x \mapsto \text{node}(-, -, -, -, -, -)\} \\ \wedge N = \text{node}(l, k_0, p, D, k', p') \end{array} \right\}$	<code>put(N, x)</code>	$\left. \begin{array}{l} \{x \mapsto \text{node}(l, k_0, p, D, k', p')\} \\ \wedge N = \text{node}(l, k_0, p, D, k', p') \end{array} \right\}$
$\{h \mapsto \text{stack}\}$	<code>PB := getPrimeBlock(h)</code>	$\{h \mapsto \text{stack} \wedge \text{PB} = \text{stack}\}$
$\{h \mapsto - \wedge \text{PB} = \text{stack}\}$	<code>putPrimeBlock(h, PB)</code>	$\{h \mapsto \text{stack} \wedge \text{PB} = \text{stack}\}$

Figure 26. Specification of the heap update commands.

implementations and cases can all be proven in a similar style.

B. B^{Link} Tree Bug

Consider the following B^{Link} tree, where nodes have capacity 2:



To illustrate the bug in Sagiv's proposed algorithm we consider two threads. Thread 1 performs an insert operation, for the purpose of this example, it will insert on key 15. Concurrently, thread 2 performs an insert operation on key 35 followed by an insert on key 55.

Initially the scheduler allows thread 1 to run and it traverses the tree and reaches the leaf node, where the new key value pair will be inserted, in the tree and locks it in order to insert. We have the following code trace and tree state visible to both threads:

```

insert(h, k, v) {
  stack := newStack();
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to true
  // (if) k < highValue(N) evaluates to true
  push(stack, cur);
  cur := next(N, k);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to false
  level := 1;
  m := k;
  w := v;
  // (while) true evaluates to true
  found := false;
  // (while) found = false evaluates to true
  found := true;
  lock(cur);
  N := get(cur);
}

```

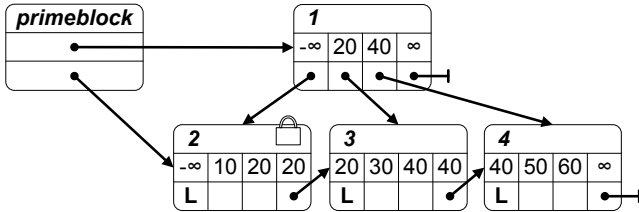
$\{\text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p')\}$	$k := \text{lowValue}(N)$	$\{\text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \wedge k = k_0\}$
$\{\text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p')\}$	$k := \text{highValue}(N)$	$\{\text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \wedge k = k'\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{inner}(l, k_0, v_0, D, k_{n+1}, p') \\ \wedge D = [(k_1, v_1), \dots, (k_n, v_n)] \\ \wedge k_i < k \leq k_{i+1} \end{array} \right\}$	$x := \text{next}(N, k)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{inner}(l, k_0, v_0, D, k_{n+1}, p') \\ \wedge x = v_i \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge k > k' \end{array} \right\}$	$x := \text{next}(N, k)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge x = p' \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{leaf}(l, k_0, D, k', p') \\ \wedge (k, v) \in D \end{array} \right\}$	$x := \text{lookup}(N, k)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{leaf}(l, k_0, D, k', p') \\ \wedge x = v \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge D < 2K \wedge k \notin \text{keys}(D) \end{array} \right\}$	$\text{addPair}(N, k, v)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D', k', p') \\ \wedge D' = D \uplus (k, v) \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', v') \\ \wedge (k, -) \in D \end{array} \right\}$	$\text{removePair}(N, k)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D', k', p') \\ \wedge D = D' \uplus (k, -) \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{leaf}(l, k_0, D, k', p') \\ \wedge k_0 < k \leq k' \wedge D = 2K \end{array} \right\}$	$M := \text{rearrange}(N, k, v, x)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{leaf}(l, k_0, D_1, k'', x) \\ \wedge M = \text{leaf}(0, k'', D_2, k', p') \\ \wedge D_1 :: D_2 = D \uplus (k, v) \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{inner}(l, k_0, p, D, k', p') \\ \wedge k_0 < k < k' \wedge D = 2K \end{array} \right\}$	$M := \text{rearrange}(N, k, v, x)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{inner}(l, k_0, p, D_1, k'', x) \\ \wedge M = \text{inner}(0, k'', p'', D_2, k', p') \\ \wedge D_1 :: (k'', p'') :: D_2 = D \uplus (k, v) \end{array} \right\}$
$\{\text{emp} \wedge \text{PB} = p : ps\}$	$x := \text{root}(\text{PB})$	$\{\text{emp} \wedge \text{PB} = p : ps \wedge x = p\}$
$\{\text{emp}\}$	$N := \text{newRoot}(k', p, k, v, k'')$	$\{\text{emp} \wedge N = \text{inner}(0, k', p, [(k, v)], k'', \text{nil})\}$
$\{\text{emp} \wedge \text{PB} = xs\}$	$\text{addRoot}(\text{PB}, x)$	$\{\text{emp} \wedge \text{PB} = x : xs\}$
$\{\text{emp} \wedge \text{PB} = [x_n, \dots, x_1] \wedge 1 \leq i \leq n\}$	$x := \text{getNodeLevel}(\text{PB}, i)$	$\{\text{emp} \wedge \text{PB} = [x_n, \dots, x_1] \wedge x = x_i\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge D = [(k_1, v_1), \dots, (k_n, v_n)] \\ \wedge n < 2K \end{array} \right\}$	$b := \text{isSafe}(N)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge D = [(k_1, v_1), \dots, (k_n, v_n)] \\ \wedge n < 2K \wedge b = \text{tt} \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge D = [(k_1, v_1), \dots, (k_n, v_n)] \\ \wedge n = 2K \end{array} \right\}$	$b := \text{isSafe}(N)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge D = [(k_1, v_1), \dots, (k_n, v_n)] \\ \wedge n = 2K \wedge b = \text{ff} \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge (k, v) \in D \end{array} \right\}$	$b := \text{isIn}(N, k)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge b = \text{tt} \end{array} \right\}$
$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge k \notin \text{keys}(D) \end{array} \right\}$	$b := \text{isIn}(N, k)$	$\left\{ \begin{array}{l} \text{emp} \wedge N = \text{node}(l, k_0, p, D, k', p') \\ \wedge b = \text{ff} \end{array} \right\}$
$\{\text{emp} \wedge N = \text{leaf}(l, k_0, D, k', p')\}$	$b := \text{isLeaf}(N)$	$\{\text{emp} \wedge N = \text{leaf}(l, k_0, D, k', p') \wedge b = \text{tt}\}$
$\{\text{emp} \wedge N = \text{inner}(l, k_0, p, D, k', p')\}$	$b := \text{isLeaf}(N)$	$\{\text{emp} \wedge N = \text{inner}(l, k_0, p, D, k', p') \wedge b = \text{ff}\}$
$\{\text{emp} \wedge \text{PB} = x : xs\}$	$b := \text{isRoot}(\text{PB}, x)$	$\{\text{emp} \wedge \text{PB} = x : xs \wedge b = \text{tt}\}$
$\{\text{emp} \wedge \text{PB} = y : ys \wedge x \neq y\}$	$b := \text{isRoot}(\text{PB}, x)$	$\{\text{emp} \wedge \text{PB} = y : ys \wedge b = \text{ff}\}$
$\{\text{emp}\}$	$\text{stack} := \text{newStack}()$	$\{\text{emp} \wedge \text{stack} = []\}$
$\{\text{emp} \wedge \text{stack} = xs\}$	$\text{push}(\text{stack}, x)$	$\{\text{emp} \wedge \text{stack} = x : xs\}$
$\{\text{emp} \wedge \text{stack} = y : ys\}$	$x := \text{pop}(\text{stack})$	$\{\text{emp} \wedge \text{stack} = ys \wedge x = y\}$
$\{\text{emp} \wedge \text{stack} = []\}$	$b := \text{isEmpty}(\text{stack})$	$\{\text{emp} \wedge \text{stack} = [] \wedge b = \text{tt}\}$
$\{\text{emp} \wedge \text{stack} = x : xs\}$	$b := \text{isEmpty}(\text{stack})$	$\{\text{emp} \wedge \text{stack} = x : xs \wedge b = \text{ff}\}$

Figure 27. Specification of the store update commands.

```

// (if) isIn(N, m) evaluates to false
// (if) m > highValue(N) evaluates to false
// (while) found = false evaluates to false
// (if) isSafe(N) evaluates to false
PB := getPrimeBlock(h);
// (if) isRoot(PB, cur) evaluates to false
insertIntoUnsafe;

```

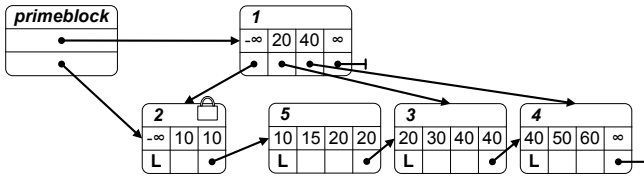


While having the lock on thread 1, it checks that the node is full and splits it into two by performing insertIntoUnsafe:

```

insertIntoUnsafe {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, x);
  put(N, cur);
}

```



```

unlock(cur);
w := x;
m := highValue(N);
level := level + 1;
// (if) isEmpty(stack) evaluates to false
cur := pop(stack);

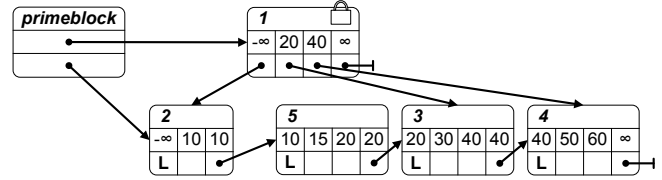
```

After that the thread unlocks the leaf node. Since the thread 1 inserted into a full node which was not the root, it is required to add a reference to the newly created node in the next level up the tree. It reads from the stack the next level up node and continues the insert algorithm.

```

// (while) true evaluates to true
found := false;
// (while) found = false evaluates to true
found := true;
lock(cur);
N := get(cur);
// (if) isIn(N, m) evaluates to false
// (if) m > highValue(N) evaluates to false
// (while) found = false evaluates to false
// (if) isSafe(N) evaluates to false
PB := getPrimeBlock(h);
// (if) isRoot(PB, cur) evaluates to true
insertIntoUnsafeRoot;

```



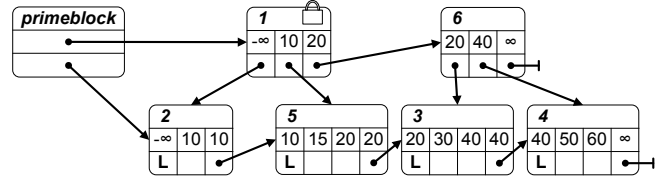
While having the lock on thread 1, it checks that the node is full and splits it into two by performing the first part of insertIntoUnsafeRoot up to the statement put(N, cur):

```

insertIntoUnsafeRoot {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, x);
  put(N, cur);
}

```

adding the new key-value pair to the node and updating the tree. We now have the following shared state:



Because the thread has created a new node it is required to create a new root, but before it is able to do so, the scheduler allows thread 2 to run. Thread 2 starts by traversing the tree until it reaches the node which is the right one to insert 35 into. It locks it as follows:

```

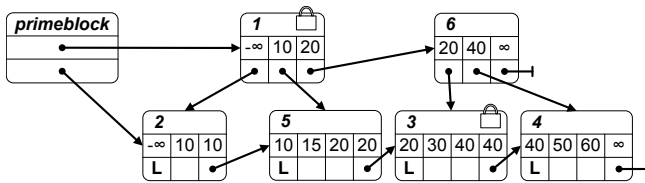
insert(h, v, p) {
  stack := newStack();
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to true
  // (if) k < highValue(N) evaluates to false
  cur := next(N, k);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to true
  // (if) k < highValue(N) evaluates to true
  push(stack, cur);
  cur := next(N, k);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to false
  level := 1;
  m := k;
  w := v;
  // (while) true evaluates to true
  found := false;
  // (while) found = false evaluates to true
  found := true;
  lock(cur);
  N := get(cur);
  // (if) isIn(N, m) evaluates to false
  // (if) m > highValue(N) evaluates to false
  // (while) found = false evaluates to false
}

```

```

// (if) isSafe(N) evaluates to false
PB := getPrimeBlock(h);
// (if) isRoot(PB, cur) evaluates to false
insertIntoUnsafe;

```

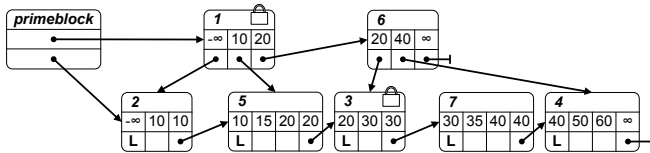


Since the node is full, it performs the insertIntoUnsafe section of the code, it splits the node by creating a new node and adding the new node to the tree as follows:

```

insertIntoUnsafe {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, x);
  put(N, cur);
}

```



```

unlock(cur);
w := x;
m := highValue(N);
level := level + 1;
// (if) isEmpty(stack) evaluates to false
cur := pop(stack);

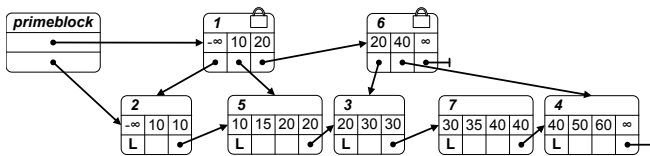
```

After that the thread unlocks the leaf node. Since the thread 2 inserted into a full node which was not the root, it is required to add a reference to the newly created node in the next level up the tree. It reads from the stack the next level up node and continues the insert algorithm.

```

// (while) true evaluates to true
found := false;
// (while) found = false evaluates to true
found := true;
lock(cur);
N := get(cur);
// (if) isIn(N, m) evaluates to false
// (if) m > highValue(N) evaluates to false
// (while) found = false evaluates to false
// (if) isSafe(N) evaluates to true
insertIntoSafe;

```

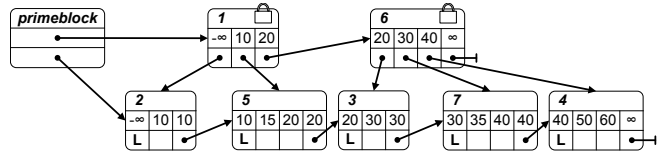


Since the node is not full, it performs the insertIntoSafe section of the code as follows:

```

insertIntoSafe {
  addPair(N, m, w);
  put(N, cur);
}

```



```

unlock(cur);
return;

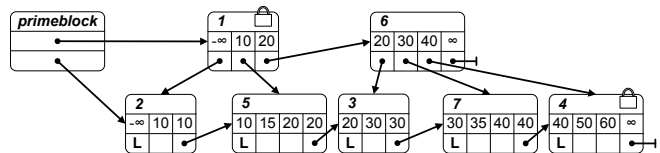
```

After that the thread unlocks the node which it inserted and returns. Thread 2 now starts the second insert operation for key 55. It traverses the tree and reaches the leaf node, where the new key value pair will be inserted, in the tree and locks it in order to insert. We have the following code trace and tree state:

```

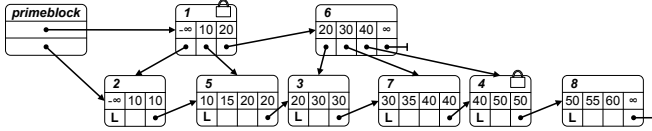
insert(h, k, v) {
  stack := newStack();
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to true
  // (if) k < highValue(N) evaluates to false
  cur := next(N, k);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to true
  // (if) k < highValue(N) evaluates to true
  push(stack, cur);
  cur := next(N, k);
  N := get(cur);
  // (while) isLeaf(N) = false evaluates to false
  level := 1;
  m := k;
  w := v;
  // (while) true evaluates to true
  found := false;
  // (while) found = false evaluates to true
  found := true;
  lock(cur);
  N := get(cur);
  // (if) isIn(N, m) evaluates to false
  // (if) m > highValue(N) evaluates to false
  // (while) found = false evaluates to false
  // (if) isSafe(N) evaluates to false
  PB := getPrimeBlock(h);
  // (if) isRoot(PB, cur) evaluates to false
  insertIntoUnsafe;
}

```



Since the node is full, it performs the insertIntoUnsafe section of the code, it splits the node by creating a new node and adding the new node to the tree as follows:

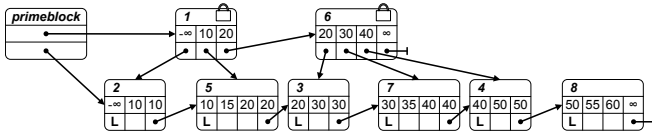

```
insertIntoUnsafe {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, x);
  put(N, cur);
}
```



```
unlock(cur);
w := x;
m := highValue(N);
level := level + 1;
// (if) isEmpty(stack) evaluates to false
cur := pop(stack);
```

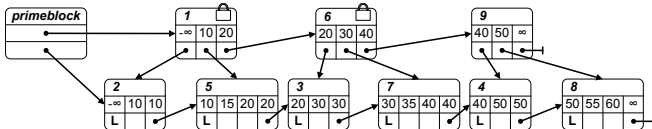
After that the thread unlocks the leaf node. Since the thread 2 inserted into a full node which was not the root, it is required to add a reference to the newly created node in the next level up the tree. It reads from the stack the next level up node and continues the insert algorithm by locking the node at the level up in order to insert.

```
// (while) true evaluates to true
found := false;
// (while) found = false evaluates to true
found := true;
lock(cur);
N := get(cur);
// (if) isIn(N, m) evaluates to false
// (if) m > highValue(N) evaluates to false
// (while) found = false evaluates to false
// (if) isSafe(N) evaluates to false
PB := getPrimeBlock(h);
// (if) isRoot(PB, cur) evaluates to false
insertIntoUnsafe;
```



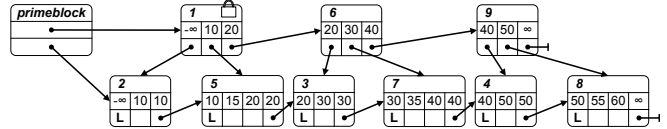
Since the node is full, it starts to perform the insertIntoUnsafe section of the code, it splits the node by creating a new node and adding the new node to the tree (by performing put(N, cur)) as follows:

```
insertIntoUnsafe {
  x := new();
  M := rearrange(N, m, w, x);
  put(M, x);
  put(N, cur);
}
```



After that the thread unlocks the node. Since thread 2 inserted into a full node which was not the root, it is required to add a reference to the newly created node in the next level up the tree.

```
unlock(cur);
w := x;
m := highValue(N);
level := level + 1;
if (isEmpty(stack)) {
  PB := getPrimeBlock(h);
  cur := getNodeLevel(PB, level);
}
```



When it tries to read the next level up from the prime block, by calling cur := getNodeLevel(PB, level) it steps out of bounds of the prime block array and faults. It is important to notice that real implementations have been corrected to avoid this bug. When thread 1 tries to insert on the full root it locks the newly created node before creating a new root node for the tree. This prevents other threads from inserting in the new node, or the old root node, until a new root has been created and both nodes are unlocked.

C. B^{Link} Tree Insert Proof

The proof uses the predicate $\text{iniceNode}(N, r, h)$, defined in Figure 28. The definition of iniceNode asserts that the node descriptor N contains legitimate information about the tree. If N is an inner node, then the children and link pointers of N must all point to extant nodes in the tree, which have the minimum values specified by N – this ensures that following a pointer reaches an appropriate node. The definition of istLf asserts that the leaf node descriptor N contains legitimate information about the tree and contains the same information as the leaf node n in the tree. The definition of istIn is analogous but for inner nodes.

Assertions in the proof must be stable – that is, invariant under interference from other threads. The stability of iniceNode is ensured by the fact that the capabilities held by the thread do not allow nodes to be removed and the minimum values of nodes to change. The stability of istLf and istIn is ensured by the fact that the nodes are locked and the thread which makes use of the predicate holds the lock on the node, this makes sure that no other thread can insert (or remove on the leaf case) on the node.

```

{outdef(h, k)1}
insert(h, k, v) {
  {B∅(h, k)I(r,h)r * dcaps(k, r, 1)}
  stack := newStack();
  PB := getPrimeBlock(h);
  cur := root(PB);
  N := get(cur);
  {
    {B∅(h, k)I(r,h)r * dcaps(k, r, 1) ∧ stack = []}
    * iniceNode(N, r, h) * present(cur, k0, nil)
    ∧ N = node(-, k0, p, D, k', p') ∧ k0 = -∞
  }
  while (isLeaf(N) = false) {
    if (k < highValue(N)) {
      push(stack, cur);
    }
    cur := next(N, k);
    N := get(cur);
  }
  level := 1;
  m := k;
  w := v;
  {
    {B∅(h, k)I(r,h)r * dcaps(k, r, 1) ∧ stack = xs}
    * iniceNode(N, r, h) * present(cur, k', nil)
    ∧ N = leaf(-, k', D, k'', p') ∧ k' < k ∧ level = 1
    ∧ m = k ∧ w = v
  }
  while (true) {

```

```

    {
      {B∅(h, k)I(r,h)r * dcaps(k, r, 1) ∧ stack = xs}
      * iniceNode(N, r, h) * present(cur, k', nil)
      ∧ N = leaf(-, k', D, k'', p') ∧ k' < k ∧ level = 1
      ∧ m = k ∧ w = v
    }
    ∨
    {
      {B∈(h, k, v)I(r,h)r * dcaps(k, r, 1) ∧ stack = xs}
      * present(cur, k', p) * present(w, m, -)
      * nodeList(p, N', w)
      * [FIX(m, w)]1r ∧ k' < m ∧ level > 1
    }
    found := false;
    while (found = false) {
      found := true;
      lock(cur); // use LOCK
      N := get(cur);
      if (isIn(N, m)) {
        {false}
        unlock(cur);
        return;
      }
      if (m > highValue(N)) {
        unlock(cur); // use UNLOCK
        found := false;
        while (m > highValue(N)) {
          cur := next(N, m);
          N := get(cur);
        }
      }
    }
    {
      {B∅(h, k)I(r,h)r * dcaps(k, r, 1) ∧ stack = xs}
      * istLf(cur, N, r, h) * [UNLOCK(cur)]1r
      ∧ N = leaf(1, k', D, k'', p') ∧ k' < k ≤ k'' ∧ level = 1
      ∧ m = k ∧ w = v
    }
    ∨
    {
      {B∈(h, k, v)I(r,h)r * dcaps(k, 1) ∧ stack = xs}
      * istIn(cur, N, r, h) * present(w, m, -) * [FIX(m, w)]1r
      * [UNLOCK(cur)]1r ∧ N = inner(1, k', p', D, k'', p'')
      * nodeList(p', N', w) ∧ k' < m < k'' ∧ level > 1
    }
    if (isSafe(N)) {
      {
        {B∅(h, k)I(r,h)r * dcaps(k, r, 1) ∧ stack = xs}
        * istLf(cur, N, r, h) * [UNLOCK(cur)]1r
        ∧ N = leaf(1, k', D, k'', p') ∧ |D| < 2K
        ∧ k' < k ≤ k'' ∧ level = 1 ∧ m = k ∧ w = v
      }
      ∨
      {
        {B∈(h, k, v)I(r,h)r * dcaps(k, 1) ∧ stack = xs}
        * istIn(cur, N, r, h) * present(w, m, -) * [FIX(m, w)]1r
        * [UNLOCK(cur)]1r ∧ N = inner(1, k', p', D, k'', p'')
        * nodeList(p', N', w) ∧ |D| < 2K ∧ k' < m < k''
        ∧ level > 1
      }
      insertIntoSafe;
      {indef(h, k, v)1}
      return;
    } else {
      PB := getPrimeBlock(h);
      if (isRoot(PB, cur)) {

```

$$\begin{aligned}
\text{iniceNode}(N, r, h) &\triangleq \exists k_0, p_0, D, k', p'. \left(k' = +\infty \vee \boxed{p' \mapsto \text{node}(-, k', -, -, -) * \text{true}}_{I(r,h)}^r \right) \wedge \\
&\quad \left(\left(N = \text{inner}(-, k_0, p_0, D, k', p') \wedge \forall (k, p) \in D. \right. \right. \\
&\quad \quad \left. \left. \boxed{p \mapsto \text{node}(-, k, -, -, -) * \text{true}}_{I(r,h)}^r \vee N = \text{leaf}(-, k_0, D, k', p') \right) \right. \\
&\quad \quad \left. \wedge \boxed{p_0 \mapsto \text{node}(-, k_0, -, -, -) * \text{true}}_{I(r,h)}^r \right) \\
\text{istLf}(n, N, r, h) &\triangleq \exists k_0, D, k', p'. \left(k' = +\infty \vee \boxed{p' \mapsto \text{leaf}(-, k', -, -, -) * \text{true}}_{I(r,h)}^r \right) \\
&\quad \wedge N = \text{leaf}(-, k_0, D, k', p') \wedge \boxed{n \mapsto \text{leaf}(1, k_0, D, k', p') * \text{true}}_{I(r,h)}^r \\
\text{istIn}(n, N, r, h) &\triangleq \exists k_0, p_0, D, k', p'. \left(k' = +\infty \vee \boxed{p' \mapsto \text{inner}(-, k', -, -, -) * \text{true}}_{I(r,h)}^r \right) \\
&\quad \wedge N = \text{inner}(-, k_0, p_0, D, k', p') \wedge \boxed{n \mapsto \text{inner}(1, k_0, p_0, D, k', p') * \text{true}}_{I(r,h)}^r \\
&\quad \wedge \boxed{p_0 \mapsto \text{node}(-, k_0, -, -, -) * \text{true}}_{I(r,h)}^r \\
&\quad \wedge \forall (k, p) \in D. \boxed{p \mapsto \text{node}(-, k, -, -, -) * \text{true}}_{I(r,h)}^r
\end{aligned}$$

Figure 28. Predicates used in the B^{Link} tree insert proof.

```

{
  (
    (
       $\boxed{B_{\notin}(\mathbf{h}, \mathbf{k})}_{I(r,h)}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs$ 
      *  $\text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{UNLOCK}(\text{cur})]_1^r$ 
       $\wedge N = \text{leaf}(1, k', D, k'', p') \wedge |D| = 2K$ 
       $\wedge \text{root}(\mathbf{h}, \text{cur}) \wedge k' < \mathbf{k} \leq k'' \wedge \text{level} = 1$ 
       $\wedge \mathbf{m} = \mathbf{k} \wedge \mathbf{w} = \mathbf{v}$ 
    )
    \vee
    (
       $\boxed{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})}_{I(r,h)}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs$ 
      *  $\text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(\mathbf{w}, \mathbf{m}, -) * [\text{FIX}(\mathbf{m}, \mathbf{w})]_1^r$ 
      *  $[\text{UNLOCK}(\text{cur})]_1^r \wedge N = \text{inner}(1, k', p', D, k'', p'')$ 
      *  $\text{nodeList}(p', N', \mathbf{w}) \wedge |D| = 2K \wedge \text{root}(\mathbf{h}, \text{cur})$ 
       $\wedge k' < \mathbf{m} < k'' \wedge \text{level} > 1$ 
    )
  )
  insertIntoUnsafeRoot;
  {indef(h, k, v)}_1
  return;
} else {
  (
    (
       $\boxed{B_{\notin}(\mathbf{h}, \mathbf{k})}_{I(r,h)}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs$ 
      *  $\text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{UNLOCK}(\text{cur})]_1^r$ 
       $\wedge N = \text{leaf}(1, k', D, k'', p') \wedge |D| = 2K$ 
       $\wedge \neg \text{root}(\mathbf{h}, \text{cur}) \wedge k' < \mathbf{k} \leq k'' \wedge \text{level} = 1$ 
       $\wedge \mathbf{m} = \mathbf{k} \wedge \mathbf{w} = \mathbf{v}$ 
    )
    \vee
    (
       $\boxed{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})}_{I(r,h)}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs$ 
      *  $\text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(\mathbf{w}, \mathbf{m}, -) * [\text{FIX}(\mathbf{m}, \mathbf{w})]_1^r$ 
      *  $[\text{UNLOCK}(\text{cur})]_1^r \wedge N = \text{inner}(1, k', p', D, k'', p'')$ 
      *  $\text{nodeList}(p', N', \mathbf{w}) \wedge |D| = 2K \wedge \neg \text{root}(\mathbf{h}, \text{cur})$ 
       $\wedge k' < \mathbf{m} < k'' \wedge \text{level} > 1$ 
    )
  )
  insertIntoUnsafe;
  (
     $\boxed{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})}_{I(r,h)}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs$ 
    *  $\text{present}(\text{cur}, k', -) * \text{present}(\mathbf{w}, \mathbf{m}, -)$ 
    *  $[\text{FIX}(\mathbf{m}, \mathbf{w})]_1^r \wedge k' < \mathbf{m} \wedge \text{level} > 1$ 
  )
}
}

```

insertIntoSafe {

$$\left\{ \begin{array}{l} (\underline{B_{\notin}(\mathbf{h}, \mathbf{k})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{UNLOCK}(\text{cur})]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, k', D, k'', p') \wedge |D| < 2K \wedge k' < k \leq k'' \\ \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(w, m, -) * [\text{FIX}(m, w)]_1^r \\ * [\text{UNLOCK}(\text{cur})]_1^r \wedge \mathbf{N} = \text{inner}(1, k', p', D, k'', p'') \\ * \text{nodeList}(p', N', w) \wedge |D| < 2K \wedge k' < m < k'' \\ \wedge \text{level} > 1 \end{array} \right\}$$

// use INS or FIX

$$\left\{ \begin{array}{l} (\underline{B_{\notin}(\mathbf{h}, \mathbf{k})})_{I(r, \mathbf{h})}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(0, \mathbf{k})]_{(d, 1)}^r \\ * \bigotimes_{v \in \text{Vals} \setminus \{v\}} [\text{INS}(0, \mathbf{k}, v)]_{(d, 1)}^r \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{MODLI}(0, \text{cur}, \mathbf{k}, v, 1)]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, k', D, k'', p') \wedge |D| < 2K \wedge k' < k \leq k'' \\ \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(w, m, -) \\ * [\text{MODII}(\text{cur}, m, w)]_1^r \wedge \mathbf{N} = \text{inner}(1, k', p', D, k'', p'') \\ * \text{nodeList}(p', N', w) \wedge |D| < 2K \wedge k' < m < k'' \\ \wedge \text{level} > 1 \end{array} \right\}$$

addPair(N, m, w);

put(N, cur); // use MODLI or MODII

$$\left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{UNLOCK}(\text{cur})]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, k', D', k'', p') \wedge D' = D \uplus (\mathbf{k}, \mathbf{v}) \\ \wedge k' < k \leq k'' \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(w, m, -) \\ \wedge \mathbf{N} = \text{inner}(1, k', p', D', k'', p'') * \text{nodeList}(p', N', w) \\ \wedge D' = D \uplus (m, w) \wedge k' < m < k'' \wedge \text{level} > 1 \end{array} \right\}$$

unlock(cur); // use UNLOCK

$$\left\{ \underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})} \right\}_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, r, 1)$$

$$\left\{ \text{in}_{\text{def}}(\mathbf{h}, \mathbf{k}, \mathbf{v})_1 \right\}$$

}

insertIntoUnsafeRoot {

$$\left\{ \begin{array}{l} (\underline{B_{\notin}(\mathbf{h}, \mathbf{k})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{UNLOCK}(\text{cur})]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, k', D, k'', p') \wedge |D| = 2K \wedge \text{root}(\mathbf{h}, \text{cur}) \\ \wedge k' < k \leq k'' \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(w, m, -) * [\text{FIX}(m, w)]_1^r \\ * [\text{UNLOCK}(\text{cur})]_1^r \wedge \mathbf{N} = \text{inner}(1, k', p', D, k'', p'') \\ * \text{nodeList}(p', N', w) \wedge |D| = 2K \wedge \text{root}(\mathbf{h}, \text{cur}) \\ \wedge k' < m < k'' \wedge \text{level} > 1 \end{array} \right\}$$

// use INS or FIX

$$\left\{ \begin{array}{l} (\underline{B_{\notin}(\mathbf{h}, \mathbf{k})})_{I(r, \mathbf{h})}^r * [\text{LOCK}]_g^r * [\text{SWAP}]_g^r * [\text{REM}(0, \mathbf{k})]_{(d, 1)}^r \\ * \bigotimes_{v \in \text{Vals} \setminus \{v\}} [\text{INS}(0, \mathbf{k}, v)]_{(d, 1)}^r \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{MODLI}(0, \text{cur}, \mathbf{k}, v, 1)]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, k', D, k'', p') \wedge |D| = 2K \wedge \text{root}(\mathbf{h}, \text{cur}) \\ \wedge k' < k \leq k'' \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(w, m, -) \\ * [\text{MODII}(\text{cur}, m, w)]_1^r \wedge \mathbf{N} = \text{inner}(1, k', p', D, k'', p'') \\ * \text{nodeList}(p', N', w) \wedge |D| = 2K \wedge \text{root}(\mathbf{h}, \text{cur}) \\ \wedge k' < m < k'' \wedge \text{level} > 1 \end{array} \right\}$$

x := new();

M := rearrange(N, m, w, x);

put(M, x);

lock(x);

put(N, cur); // use MODLI or MODII

$$\left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{NEWR}(\text{cur}, k''', \mathbf{x})]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, k', D_1, k''', \mathbf{x}) * \text{istLf}(\mathbf{x}, \mathbf{M}, r, \mathbf{h}) \\ \wedge \mathbf{M} = \text{leaf}(-, k''', D_2, k'', p') \wedge D_1 :: D_2 = D \uplus (\mathbf{k}, \mathbf{v}) \\ \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * \text{present}(w, m, -) \\ * [\text{NEWR}(\text{cur}, k''', \mathbf{x})]_1^r \wedge \mathbf{N} = \text{inner}(1, k', p', D_1, k''', \mathbf{x}) \\ * \text{istIn}(\mathbf{x}, \mathbf{M}, r, \mathbf{h}) \wedge \mathbf{M} = \text{inner}(-, k''', p''', D_2, k'', p'') \\ * \text{nodeList}(p', N', w) \wedge D_1 :: (k''', p''') :: D_2 = D \uplus (m, w) \\ \wedge \text{level} > 1 \end{array} \right\}$$

y := lowValue(N);

t := highValue(N);

u := highValue(M);

$$\left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, r, 1) \wedge \text{stack} = xs \\ * \text{istLf}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{NEWR}(\text{cur}, \mathbf{t}, \mathbf{x})]_1^r \\ \wedge \mathbf{N} = \text{leaf}(1, y, D_1, \mathbf{t}, \mathbf{x}) * \text{istLf}(\mathbf{x}, \mathbf{M}, r, \mathbf{h}) \\ \wedge \mathbf{M} = \text{leaf}(-, \mathbf{t}, D_2, u, p') \wedge \text{level} = 1 \wedge m = k \wedge w = v \end{array} \right\} \\ \vee \\ \left\{ \begin{array}{l} (\underline{B_{\in}(\mathbf{h}, \mathbf{k}, \mathbf{v})})_{I(r, \mathbf{h})}^r * \text{dcaps}(\mathbf{k}, 1) \wedge \text{stack} = xs \\ * \text{istIn}(\text{cur}, \mathbf{N}, r, \mathbf{h}) * [\text{NEWR}(\text{cur}, \mathbf{t}, \mathbf{x})]_1^r \\ \wedge \mathbf{N} = \text{inner}(1, y, p', D_1, \mathbf{t}, \mathbf{x}) * \text{istIn}(\mathbf{x}, \mathbf{M}, r, \mathbf{h}) \\ \wedge \mathbf{M} = \text{inner}(-, \mathbf{t}, p''', D_2, u, p'') \wedge \text{level} > 1 \end{array} \right\}$$

r := new();

R := newNode(y, cur, t, x, u);

PB := getPrimeBlock(h);

put(R, r);

```

addRoot(PB, r);
putPrimeBlock(h, PB); // use NEWR
{
  ( [B∉(h, k, v)]I(r,h)r * dcaps(k, r, 1) ∧ stack = xs
  * istLf(cur, N, r, h) * [UNLOCK(cur)]1r * [UNLOCK(x)]1r
  ∧ N = leaf(1, y, D1, t, x) * istLf(x, M, r, h)
  ∧ M = leaf(-, t, D2, u, p') * iniceNode(R, r, h)
  * present(r, t, cur) ∧ R = inner(0, t, [(t, x)], u, nil)
  ∧ level = 1 ∧ m = k ∧ w = v )
  ∨
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, 1) ∧ stack = xs
  * istln(cur, N, r, h) * [UNLOCK(cur)]1r * [UNLOCK(x)]1r
  ∧ N = inner(1, y, p', D1, t, x) * istln(x, M, r, h)
  ∧ M = inner(-, t, p''', D2, u, p'') * iniceNode(R, r, h)
  * present(r, t, cur) ∧ R = inner(0, t, [(t, x)], u, nil)
  ∧ level > 1 )
  unlock(cur); // use UNLOCK
  unlock(x); // use UNLOCK
  { [B∈(h, k, v)]I(r,h)r * dcaps(k, r, 1) }
  { indef(h, k, v)1 }
}

```

```

insertIntoUnsafe {
  ( [B∉(h, k)]I(r,h)r * dcaps(k, r, 1) ∧ stack = xs
  * istLf(cur, N, r, h) * [UNLOCK(cur)]1r
  ∧ N = leaf(1, k', D, k'', p') ∧ |D| = 2K ∧ ¬root(h, cur)
  ∧ k' < k ≤ k'' ∧ level = 1 ∧ m = k ∧ w = v )
  ∨
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, 1) ∧ stack = xs
  * istln(cur, N, r, h) * present(w, m, -) * [FIX(m, w)]1r
  * [UNLOCK(cur)]1r ∧ N = inner(1, k', p', D, k'', p'')
  * nodeList(p', N', w) ∧ |D| = 2K ∧ ¬root(h, cur)
  ∧ k' < m < k'' ∧ level > 1 )
  // use INS or FIX
  ( [B∉(h, k)]I(r,h)r * [LOCK]gr * [SWAP]gr * [REM(0, k)](d,1)r
  * ⊗v∈Vals\{v} [INS(0, k, v)](d,1)r ∧ stack = xs
  * istLf(cur, N, r, h) * [MODLI(0, cur, k, v, 1)]1r
  ∧ N = leaf(1, k', D, k'', p') ∧ |D| = 2K ∧ ¬root(h, cur)
  ∧ k' < k ≤ k'' ∧ level = 1 ∧ m = k ∧ w = v )
  ∨
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, 1) ∧ stack = xs
  * istln(cur, N, r, h) * present(w, m, -)
  * [MODLI(cur, m, w)]1r ∧ N = inner(1, k', p', D, k'', p'')
  * nodeList(p', N', w) ∧ |D| = 2K ∧ ¬root(h, cur)
  ∧ k' < m < k'' ∧ level > 1 )
  x := new();
  M := rearrange(N, m, w, x);
  put(M, q);
  put(N, cur); // use MODLI or MODII
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, r, 1) ∧ stack = xs
  * istLf(cur, N, r, h) * [FIX(k''', x)]1r
  * [UNLOCK(cur)]1r * present(x, k''', -)
  ∧ N = leaf(1, k', D1, k''', x) * iniceNode(M, r, h)
  ∧ M = leaf(0, k''', D2, k'', p') ∧ D1 :: D2 = D ⊔ (k, v)
  ∧ ¬root(h, cur) ∧ level = 1 ∧ m = k ∧ w = v )
  ∨
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, 1) ∧ stack = xs
  * istln(cur, N, r, h) * present(w, m, -)
  * [FIX(k''', x)]1r * present(x, k''', -)
  ∧ N = inner(1, k', p', D1, k''', x) * iniceNode(M, r, h)
  ∧ M = inner(0, k''', p''', D2, k'', p'') * nodeList(p', N', w)
  ∧ D1 :: (k''', p''') :: D2 = D ⊔ (m, w)
  ∧ ¬root(h, cur) ∧ level > 1 )
  unlock(cur); // use UNLOCK
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, r, 1) ∧ stack = xs
  * iniceNode(N, r, h) * present(x, k''', -)
  * [FIX(k''', x)]1r
  ∧ N = leaf(1, k', D1, k''', x) * iniceNode(M, r, h)
  ∧ M = leaf(0, k''', D2, k'', p')
  ∧ level = 1 ∧ m = k ∧ w = v )
  ∨
  ( [B∈(h, k, v)]I(r,h)r * dcaps(k, 1) ∧ stack = xs
  * iniceNode(N, r, h) * present(x, k''', -)
  * [FIX(k''', x)]1r
  ∧ N = inner(1, k', p', D1, k''', x) * iniceNode(M, r, h)
  ∧ M = inner(0, k''', p''', D2, k'', p'')
  ∧ level > 1 )
  w := x;

```

```

m := highValue(N);
level := level + 1;
{
  (  $\boxed{B \in (h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) \wedge \text{stack} = xs$ 
  *  $\text{iniceNode}(N, r, h) * \text{present}(w, m, -)$ 
  *  $[\text{FIX}(m, w)]_1^r$ 
   $\wedge N = \text{leaf}(1, k', D_1, m, w) * \text{iniceNode}(M, r, h)$ 
   $\wedge M = \text{leaf}(0, m, D_2, k'', p') \wedge \text{level} = 2$ 
   $\vee$ 
  (  $\boxed{B \in (h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, 1) \wedge \text{stack} = xs$ 
  *  $\text{iniceNode}(N, r, h) * \text{present}(w, m, -)$ 
  *  $[\text{FIX}(m, w)]_1^r$ 
   $\wedge N = \text{inner}(1, k', p', D_1, m, w) * \text{iniceNode}(M, r, h)$ 
   $\wedge M = \text{inner}(0, m, p''', D_2, k'', p'')$ 
   $\wedge \text{level} > 1$ 
)
}
if (isEmpty(stack)) {
  PB := getPrimeBlock(h);
  cur := getNodeLevel(PB, level);
} else {
  cur := pop(stack);
}
{
  (  $\boxed{B \in (h, k, v)}_{I(r, h)}^r * \text{dcaps}(k, r, 1) \wedge \text{stack} = xs$ 
  *  $\text{present}(cur, k', p) * \text{present}(w, m, -)$ 
  *  $\text{nodeList}(p, N', w)$ 
  *  $[\text{FIX}(m, w)]_1^r \wedge k' < m \wedge \text{level} > 1$ 
)
}
}

```