

# Programming Workbench

## Hand-Over-Hand Locking for Highly Concurrent Collections

TERENCE KELLY

Terence Kelly studied computer science at Princeton and the University of Michigan, followed by a long stint at Hewlett-Packard Laboratories. Kelly now writes code and documentation promoting persistent memory programming and other programming techniques. He usually avoids falling off the monkey bars on the playground by remembering to grab the next bar before letting go of the previous one. His publications are listed at <http://ai.eecs.umich.edu/~tpkelly/papers/> and he welcomes feedback at [tpkelly@eecs.umich.edu](mailto:tpkelly@eecs.umich.edu).

Welcome to “Programming Workbench,” a new column that will delve into interesting programming problems and solve them with working software. All code is available in machine-readable form at [7]. I welcome feedback from readers, the best of which I may discuss in future columns.

This first installment of “Programming Workbench” reviews a concurrent programming pattern that every developer should know: hand-over-hand locking. Over the past year, I’ve been surprised more than once to meet well-educated, experienced, proficient programmers who aren’t familiar with this versatile and powerful technique. After a bit of digging I began to understand why it’s underappreciated: hand-over-hand locking isn’t mentioned at all in numerous places where I’d expect a detailed treatment: for example, several Pthreads books and several other books on systems programming in my personal library. A few books mention it without going into great detail [1, 8]. One magazine article discusses the technique at length without providing code [10]. I found only one source with both a detailed discussion and an implementation (in Java) [2].

Why should programmers care about concurrency control in general and hand-over-hand locking in particular? In a word, performance. Even in the bygone age of uniprocessors, multithreaded code made servers more efficient and made interactive software more responsive by overlapping computation with I/O. Today, well-designed *concurrent* software enjoys genuine *parallel* execution on ubiquitous multicore and multiprocessor hardware. Embarrassing parallelism, in which different threads don’t interact at all, remains “good work if you can find it”; most multithreaded software, however, isn’t so lucky and must orchestrate orderly access to shared memory. Mutex-based concurrency control is the most conventional way to do so, and hand-over-hand locking is a primordial pattern that embodies timeless principles—and sometimes outperforms the alternatives.

So let’s brush up on hand-over-hand locking. We’ll start with the simplest dynamic data structure, the singly linked list, and review several ways to arrange safe access to linked lists in multithreaded programs. We’ll consider hand-over-hand locking in detail, describing its advantages over the alternatives. We’ll walk through a working C program whose threads employ the hand-over-hand protocol to access a linked list. Finally, we’ll conclude with generalizations and extensions of the basic techniques that we’ve covered.

### Concurrent Lists

A linked list is an easy way to implement the abstraction of an unordered, unindexed, dynamic collection of items. Lists support all of the operations that make sense for such collections: traversing the contents of a collection and inserting, reading, writing, and deleting items along the way. I’d use the word “set” rather than “collection” but in some contexts, e.g., the C++ Standard Template Library, `<set>` confusingly refers to an *ordered* container. Lists are useful in themselves and also as building blocks in more elaborate data structures, e.g., hash tables.

## Programming Workbench: Hand-Over-Hand Locking for Highly Concurrent Collections

If multiple threads access a collection concurrently, they must avoid data races, which lead to undefined behavior according to the C and C++ language standards. There are several ways to implement a concurrent list safely.

### Transactional Memory

Arguably the easiest concurrency control mechanism from the programmer's point of view is transactional memory (TM). TM allows a thread to execute a sequence of instructions atomically and in isolation, preventing other threads from observing intermediate states of the data that the instructions manipulate. A concurrent linked list based on TM avoids data races, and some TM research prototypes would allow genuine parallel access to a linked list, but mainstream TM implementations would effectively *serialize* access to the list. In other words, for the present purpose, off-the-shelf industrial-strength TM-based concurrency control would combine the safety and simplicity of single-threaded code with the performance of single-threaded code, defeating one of the main motives for multithreading.

### Non-Blocking Approaches

At the opposite ends of the ergonomic and performance spectra lie non-blocking (lock-free, wait-free) techniques based on the careful use of atomic CPU instructions. The main attraction of non-blocking techniques is that the untimely suspension or death of one thread (due, for example, to a software bug or an unfortunate CPU scheduling decision) doesn't prevent other threads from doing useful work. That's a major advantage compared with mutex-based isolation, which offers no similar guarantee. The main downsides of non-blocking techniques are that they're rather esoteric, to put it mildly—every new contribution is a tour de force by experts—and sometimes they work best with automatic garbage collection. See Michael [6] for a good example of a non-blocking list and Herlihy and Shavit [2] for a broad discussion of non-blocking techniques.

### Mutex-Based Isolation

Mutex-based isolation is well understood, and good implementations of POSIX-standardized mutexes have been available for decades. Protecting an entire linked list with a single mutex is easy, but such *coarse-grained* locking serializes access to the list and creates a potential performance bottleneck.

*Fine-grained* locking for a linked list means associating a mutex with each list node. Per-node locks allow multiple threads to access different parts of the list simultaneously, potentially improving performance. Fine-grained locking, however, isn't guaranteed to be faster, and indeed it can be slower than coarse-grained locking, depending on myriad details beyond the scope of this column. A more worrisome downside of fine-grained locking is that it's just plain trickier than coarse-grained locking; opportunities abound for errors that can cause data races or deadlocks.

It pays to study carefully the correct access discipline, hand-over-hand locking. We'll walk through an implementation, and then we'll reflect on the protocol's properties and benefits.

### The Code

The C99/C11 program listed in this section is available at [7]. We'll pore over everything but boilerplate like `#includes`. The purpose of the example program is to emphasize the locking protocol, so it avoids frills for the sake of clarity.

The following `struct` is the building block of our linked list. Each node on the list contains the mutex that protects it, a simple data field, and a pointer to the next node on the list.

```
typedef struct node { pthread_mutex_t m;
                    int data;
                    struct node *next; } node_t;
```

For brevity and simplicity we'll just hard-wire a short list into the program. The list consists of a dummy head node followed by five "real" nodes, A through E, whose data fields are respectively initialized to values 1 through 5:

```
#define PMI PTHREAD_MUTEX_INITIALIZER
static node_t E = {PMI, 5, NULL},
                D = {PMI, 4, &E},
                C = {PMI, 3, &D},
                B = {PMI, 2, &C},
                A = {PMI, 1, &B},
                head = {PMI, 0, &A}; // dummy node
```

For diagnostic printouts, it's convenient to derive a human-readable name from a pointer to a node. Since our quick-and-dirty program uses a short hardwired list, we can get away with a static mapping of node pointers to name strings. Compared with the alternative of an `if/else` statement cascade, the ternary operator (`?:`) saves keystrokes and yields a pure expression:

```
#define NAME(p) ( &head == (p) ? "head" \
                 : &A == (p) ? "A" \
                 : &B == (p) ? "B" \
                 : &C == (p) ? "C" \
                 : &D == (p) ? "D" \
                 : &E == (p) ? "E" \
                 : NULL == (p) ? "NULL" : (assert(0), "?") )
```

A simple program isn't well served by elaborate, verbose runtime checks, so we use a handful of succinct function-like macros to consolidate error checking. All of our function-like macros expand to *expressions* rather than statements because expressions may appear in a wider range of contexts; later we'll see one in the initialization part of a `for` loop.

If anything unexpected happens, the program falls on its sword via the `DIE()` macro below, which expands to a parenthesized expression that uses the comma operator to evaluate `perror()` and `assert()` for their side effects: `perror()` prints an interpre-

## Programming Workbench: Hand-Over-Hand Locking for Highly Concurrent Collections

tation of `errno`; `assert()` prints the filename and line number where things went wrong and dumps a core file that we may autopsy with a debugger. `DIE()` appears in contexts like `func() && DIE("func")`, where `func()` returns nonzero to indicate failure. The short-circuit property of the `&&` operator ensures that `DIE()` is evaluated if and only if `func()` fails.

```
#define DIE(s)      ( perror(s), assert(0), 1 )
#define PT(f, ...) ( ( errno = pthread_ ## f (__VA_ARGS__) ) \
                    && DIE(#f) )
```

There's a lot to unpack in the `PT()` macro above, so we'll walk through it slowly to see how it leverages several C preprocessor features. The problem `PT()` solves is that several Pthreads functions we use don't *set* the standard `errno` variable but instead *return* an error number; they return zero to indicate success. `PT()` allows us to call any of these functions, arranging for `errno` to be set and `DIE()` to be called if the function returns nonzero. The easiest way to understand how `PT()` does its job is to expand a typical use with the compiler's preprocessor (`gcc -E`). For example, expanding `PT(join, t[i], &tr)`; and formatting for clarity yields:

```
(
  ( errno = pthread_join (t[i], &tr) )
  &&
  ( perror("join"), assert(0), 1 )
);
```

The token-pasting operator `##` glues `PT()`'s first argument, `join` in the example above, to `pthread_`. All subsequent arguments to `PT()` correspond to the ellipsis parameter ("`...`") in the macro definition, so they get dropped in where `__VA_ARGS__` appears in the macro replacement list. In the example above, the last two `PT()` arguments `t[i]` and `&tr` end up as arguments to `pthread_join()`. The return value of `pthread_join()` is assigned to the standard `errno` variable; POSIX defines `errno` to be a per-thread variable, so there's no data race if two threads call `PT()` concurrently. The `&&` operator ensures that control reaches the expanded `DIE()` macro if and only if `pthread_join()` returns nonzero to indicate failure. Finally, notice in `PT()`'s definition that parameter `f` appears a second time in its replacement list, in "`DIE(#f)`". A single `#` is the "stringification" operator: in the example above, `PT()` argument `join` corresponds to `PT()` parameter `f`, so `DIE(#f)` in `PT()`'s replacement list expands to `DIE("join")`, whose expansion places the "join" in `perror("join")`.

Given a pointer to a list node, the `LOCK()` and `UNLK()` macros below lock and unlock the mutex embedded in the node. `UNLK()` also sets the pointer to `NULL`, which helps to catch a common and insidious bug: dereferencing a pointer to a node after unlocking the node. That would have been a silent data race, but we've turned it into a loud `SIGSEGV`.

```
#define LOCK(p)      PT(mutex_lock, (&((p)->m)))
#define UNLK(p)     ((void)PT(mutex_unlock, (&((p)->m))), (p)=NULL)
```

The next macro isn't strictly necessary, but it facilitates testing on my computer. The standard `printf()` function is thread-safe, but two race detectors that I use, Helgrind and DRD from the Valgrind family of tools, falsely attribute data races to `printf()`. Protecting `printf()` with a mutex squelches these false positives. The print mutex can't cause a deadlock because we never try to lock any other mutex while holding it.

```
static pthread_mutex_t pm = PMI; // print mutex
#define printf(...) \
do { PT(mutex_lock, &pm); printf(__VA_ARGS__); \
    PT(mutex_unlock, &pm); } while (0)
```

Now we're ready for the interesting part: function `hoh()` below traverses our linked list, observing the hand-over-hand locking protocol. `hoh()` will be the start routine passed to `pthread_create()`. Its lone argument will be an identifier string that prefixes each thread's diagnostic printouts. These prefixes make it easy to separate out per-thread reports to see what each thread saw as it traversed our linked list.

```
static void * hoh(void * ID) {
  char *id = (char *)ID;
  node_t *p, *n; // "previous" follows "next" down the list
  printf("%s: begin\n", id);
  for (p = &head, LOCK(p); NULL != (n = p->next); p = n) {
    // A: *p locked & might be dummy head
    // *n not yet locked & can't be head
    LOCK(n);
    // B: we may remove *n here
    UNLK(p);
    // C: best place to inspect *n or insert node after *n
    printf("%s: node %s @ %p data %d\n",
           id, NAME(n), (void *)n, n->data);
    n->data++;
  }
  // D
  sleep(1) && DIE("sleep"); // stall for "convoy" interleaving
  UNLK(p);
  printf("%s: end\n", id);
  return id;
}
```

The for loop of `hoh()` walks two pointers down the linked list: `n` ("next") goes first, followed by `p` ("previous"). The loop initialization locks the dummy head node, and the loop body iterates once per non-dummy node. At comment A, node `*p` is a locked node whose successor node `*n` exists; `*p` might point to the dummy head node—it does on the first iteration—but `n` never points to the head.

Now comes the "hand-over-hand" aspect: We lock the next node `*n` while still holding a lock on its predecessor `*p`. At no point in the for loop is it safe to access `*n`'s successor (`*n->next`), which is unlocked, but after locking `*n` we may access pointer `n->next`,

## Programming Workbench: Hand-Over-Hand Locking for Highly Concurrent Collections

for example, to see if we're at the end of the list by comparing it to NULL. Comment B, where both `*p` and `*n` are locked, is the right place to remove `*n`. We must lock two consecutive nodes to remove the one farther down the list, otherwise concurrent attempts by different threads to remove two adjacent nodes may interfere in such a way that only *one* node is removed [2].

After comment B we unlock `*p`. At comment C, node `*n` alone is locked; this is the best place for inspecting or modifying the contents of `*n` alone because other nodes may access `*p` simultaneously. We can insert a node after `*n` here too, but first we should ask why we care *where* to place a new node if the list represents an *unordered* collection—why not simply insert at the head of the list? At comment C we no longer hold a lock on `*p` so it's no longer safe to read or write `*p`; as noted above, the `UNLK()` macro sets `*p` to NULL to catch careless errors. Our example program prints the name and data field of node `n` and then increments the data field.

After the `for` loop terminates at the end of the list, at comment D we gratuitously `sleep()` while holding a lock on the last list node to produce an interesting “convoy” interleaving of threads.

Inserting a node into the list doesn't require anything like `hoh()`. Simply lock the head node and splice in the new node after it. If a list represents an *unordered* collection, there's seldom a good reason to insert anywhere else. It's possible to use a list to represent an *ordered* collection by inserting nodes into proper position according to some comparison criterion, but if the collection is large and we must frequently search it to find particular nodes, then a list will be inefficient compared with a search tree or skip list. If an ordered collection is *not* large it might be reasonable to store it as a list, but coarse-grained locking might outperform fine-grained locking.

The `main()` function below runs `hoh()` twice single-threaded then spawns several threads that concurrently traverse the list.

```
#define NTHREADS 4
int main(void) {
    pthread_t t[NTHREADS]; int i; void *tr;
    char m1[] = "1st (serial) traversal",
        m2[] = "2nd (serial) traversal",
        id[NTHREADS][3] = {"T0"}, {"T1"}, {"T2"}, {"T3"};
    hoh((void *)m1);
    hoh((void *)m2);
    printf("\nmain: going multi-threaded:\n\n");
    for (i = 0; i < NTHREADS; i++)
        PT(create, &t[i], NULL, hoh, (void *)id[i]);
    for (i = 0; i < NTHREADS; i++) {
        PT(join, t[i], &tr);
        printf("main: joined %s\n", (char *)tr);
    }
    printf("\nmain: all threads finished\n");
    return 0;
}
```

The example code tarball at [7] includes a README containing the commands that I use to compile and run the example program. When the program runs, the interleaving of individual thread outputs reflects the interleaving of the threads themselves as they walk down the list. In the typical output below, thread T1 zooms down the list, then stalls at the `sleep(1)` call while holding a lock on the last node, E. T2 then gets as far as D, T0 advances to C, and T3 makes it only to B. T1 wakes, releases its lock on E and exits, allowing T2, T0, and T3 to each take a step forward on the list in that order. When T2 exits, T0 and T3 each advance one hop forward. Thus the convoy of threads plods down the list in the manner of an inchworm.

```
T1: begin
T1: node A @ 0x557caa098120 data 3
T1: node B @ 0x557caa0980e0 data 4
T1: node C @ 0x557caa0980a0 data 5
T1: node D @ 0x557caa098060 data 6
T1: node E @ 0x557caa098020 data 7
T2: begin
T2: node A @ 0x557caa098120 data 4
T2: node B @ 0x557caa0980e0 data 5
T2: node C @ 0x557caa0980a0 data 6
T2: node D @ 0x557caa098060 data 7
T0: begin
T0: node A @ 0x557caa098120 data 5
T0: node B @ 0x557caa0980e0 data 6
T0: node C @ 0x557caa0980a0 data 7
T3: begin
T3: node A @ 0x557caa098120 data 6
T3: node B @ 0x557caa0980e0 data 7
T1: end
T2: node E @ 0x557caa098020 data 8
T0: node D @ 0x557caa098060 data 8
T3: node C @ 0x557caa0980a0 data 8
T2: end
T0: node E @ 0x557caa098020 data 9
T3: node D @ 0x557caa098060 data 9
T0: end
T3: node E @ 0x557caa098020 data 10
T3: end
```

Filtering the output (e.g., `./hoh | grep '^T0:'` for thread T0) makes it easier to see what individual threads encountered as they traversed the list:

```
T0: begin
T0: node A @ 0x557caa098120 data 5
T0: node B @ 0x557caa0980e0 data 6
T0: node C @ 0x557caa0980a0 data 7
T0: node D @ 0x557caa098060 data 8
T0: node E @ 0x557caa098020 data 9
T0: end

T1: begin
T1: node A @ 0x557caa098120 data 3
T1: node B @ 0x557caa0980e0 data 4
T1: node C @ 0x557caa0980a0 data 5
T1: node D @ 0x557caa098060 data 6
T1: node E @ 0x557caa098020 data 7
T1: end
```

## Programming Workbench: Hand-Over-Hand Locking for Highly Concurrent Collections

```
T2: begin
T2: node A @ 0x557caa098120 data 4
T2: node B @ 0x557caa0980e0 data 5
T2: node C @ 0x557caa0980a0 data 6
T2: node D @ 0x557caa098060 data 7
T2: node E @ 0x557caa098020 data 8
T2: end

T3: begin
T3: node A @ 0x557caa098120 data 6
T3: node B @ 0x557caa0980e0 data 7
T3: node C @ 0x557caa0980a0 data 8
T3: node D @ 0x557caa098060 data 9
T3: node E @ 0x557caa098020 data 10
T3: end
```

Each thread saw the list as previous threads left it, *precisely as though the list were protected by a single mutex*.

### Properties and Benefits

That’s an important attraction of hand-over-hand locking: we get the *parallelism* of fine-grained locking with the simple, sane *semantics* of coarse-grained locking; the changes that one thread makes while traversing the list are, from the viewpoint of all other threads, *atomic*. As the list grows large, at some point fine-grained locking usually begins to improve performance compared with coarse locking, though exactly when depends on the details. Deadlock is impossible because all threads acquire locks in the same order, i.e., list order.

The major limitation of hand-over-hand locking is that threads must traverse the list in one direction only. One implication of this “don’t look back” rule is that a thread can’t atomically splice a node out of the middle of a long list and splice it back in at the head, which is a bummer, because move-to-front lists offer outstanding performance for some purposes [9]. More generally, hand-over-hand locking doesn’t allow us to arbitrarily rearrange a linked list. If we want to rearrange a list with per-node mutexes we can simply lock the head node *and hold that lock* while locking hand-over-hand to the end of the list, thus ensuring that no other threads are accessing any node; then we may alter the list arbitrarily, because effectively we’ll be holding a big lock on the entire list.

### Generalizations and Extensions

Linked lists are a natural way to implement unordered, unindexed collections. Hash tables implement unordered but *indexed* collections, and search trees implement *ordered and indexed* collections. The techniques we’ve discussed generalize beyond linked lists to hash tables and search trees: hash tables can represent hash buckets as linked lists, each of which may employ fine-grained locking, and search trees can employ hand-over-hand locking directly.

Unfortunately, the fine-grained locking story for hash tables and search trees isn’t as tidy and compelling as that for linked lists. Hash tables invite medium-grained locking—one mutex per hash bucket—which makes more sense than fine-grained locking in the typical case where each bucket contains only a handful of items. Implementing hand-over-hand locking for *balanced* search trees is quite tricky [10].

### Persistence

Making a linked list *persistent* is conceptually straightforward: we lay out the list in a file-backed memory mapping with help from a few simple persistent memory programming techniques [3, 4]. Supporting high *concurrency* in a persistent linked list using the techniques discussed above requires “persistence-friendly” mutexes suitable for embedding in persistent data structures, which ordinary `pthread_mutex_ts` aren’t. The design of persistence-friendly mutexes is beyond the scope of this column; the main difficulty involves mutex initialization when a program restarts.

If a persistent and highly concurrent linked list must tolerate crashes, for example, because we can’t guarantee that the program accessing it will always enjoy an orderly shutdown, we’ll need a suitable crash tolerance mechanism. On conventional hardware the right crash tolerance mechanism for persistent memory programming is remarkably easy to implement by leveraging features present in certain file systems [4]. Crash tolerance imposes further requirements on persistence-friendly mutexes: post-crash recovery must quickly and conveniently restore all embedded mutexes to an *unlocked* as well as initialized state. The most onerous requirement on any program that purports to tolerate crashes is that it survive strenuous, realistic tests [5]. Documenting the design, implementation, and testing of persistent, crashproof, and highly concurrent data structures is future work, perhaps for a future installment of this column.

### Conclusion

Despite their well-known shortcomings, old-fashioned mutexes will be with us for a long time to come. Even today, conventional mutual exclusion sometimes outshines the alternatives, and fine-grained locking is sometimes the best foundation for high-performance concurrent data structures. Hand-over-hand locking is a conceptually simple protocol for safe multithreaded access to data structures protected by fine-grained locks. The simplest context where fine-grained locking and hand-over-hand traversal make sense is a linked list, and any serious student of concurrent programming should master this primordial pattern.

## Programming Workbench: Hand-Over-Hand Locking for Highly Concurrent Collections

Readers who want to go further might conduct experiments to explore the tradeoffs in different designs. For a concurrent linked list, when is it faster to use a single mutex on the entire list versus per-node mutexes? Are spinlocks faster than `pthread_mutex_ts`? If a single lock protects the entire list, how much does the move-to-front heuristic [9] help for realistic access patterns? Does it ever pay to maintain list items in sorted order? How do hand-crafted concurrent lists compare to off-the-shelf library implementations of unordered unindexed collections? Please share your results with me!

**References**

- [1] R. Arpaci-Dusseau and A. Arpaci-Dusseau, "Lock-Based Concurrent Data Structures," in *Operating Systems: Three Easy Pieces*, Chapter 29, p. 9: <http://pages.cs.wisc.edu/~remzi/OSTEP/threads-locks-usage.pdf>.
- [2] W. Herlihy and N. Shavit, *The Art of Multiprocessor Programming* (Morgan Kaufmann, 2008).
- [3] T. Kelly, "Persistent Memory Programming on Conventional Hardware," *ACM Queue*, vol. 17, no. 4 (July/August 2019): <https://queue.acm.org/detail.cfm?id=3358957>.
- [4] T. Kelly, "Good Old-Fashioned Persistent Memory," *login.*, vol. 44, no. 4 (Winter 2019): <https://www.usenix.org/publications/login/winter2019/kelly>.
- [5] T. Kelly, "Is Persistent Memory Persistent?" *ACM Queue*, vol. 18, no. 2 (March/April 2020): <https://queue.acm.org/detail.cfm?id=3400902>.
- [6] M. M. Michael, "High Performance Dynamic Lock-Free Hash Tables and List-Based Sets," in *Proceedings of the 14th ACM Symposium on Parallel Algorithms and Architectures (SPAA 2002)*: <https://docs.rs/crate/crossbeam/0.2.4/source/hash-and-skip.pdf>. DOI: <https://dl.acm.org/doi/10.1145/564870.564881>
- [7] T. Kelly, Example code to accompany this article: [https://www.usenix.org/sites/default/files/kelly0920\\_code.tgz](https://www.usenix.org/sites/default/files/kelly0920_code.tgz).
- [8] M. Scott, *Programming Language Pragmatics*, Third Edition (Morgan Kaufmann, 2009). See Exercise 12.14 on p. 642.
- [9] D. Sleator and R. Tarjan, "Amortized Efficiency of List Update and Paging Rules," *Communications of the ACM*, vol. 28, no. 2 (February 1985): <https://www.cs.cmu.edu/~sleator/papers/amortized-efficiency.pdf>. DOI: <https://dl.acm.org/doi/10.1145/564870.564881>
- [10] H. Sutter, "Choose Concurrency-Friendly Data Structures," *Dr. Dobbs's Journal*, June 27, 2008: <http://www.drdobbs.com/parallel/choose-concurrency-friendly-data-structu/208801371>.