

MULTICORE PROGRAMMING

Optimizing lock-based KCAS with HTM, OpenMP, debugging/perf

Lecture 16

Trevor Brown

LAST TIME

- Try-locks
 - Naïve (incorrect) KCAS using these
- **Versioned** (try-)locks
 - KCAS using these

THIS TIME

- Accelerating version-lock-based KCAS with HTM
 - Purpose: to demonstrate accelerating a lock-based algorithm using HTM!
- OpenMP
- Debugging and performance
 - Checksums/validation

USING HTM TO ACCELERATE TRY-LOCK BASED KCAS

- Similar to accelerating the lock-free algorithm (but a bit simpler)
- Goal
 - HTM-based KCAS that uses **try-lock based KCAS as the fallback path**
- Approach
 - Fast path algorithm:
 - Wrap try-lock based KCAS in a transaction: `xstart ; KCAS ; xend`
 - Now each KCAS on the fast path is atomic, just because it is in a transaction
 - Get rid of parts of the algorithm that are unnecessary

ADDING TRANSACTIONS

```
template <int K>
bool KCAS(addr_t ** addr, val_t * expv, val_t * newv) {
    int retries = 5;
    int status;
retry:
    if ((status = _xbegin()) == _XBEGIN_STARTED) {
        bool result = KCAS_txn<K>(addr, expv, newv);
        _xend();
        return result;
    } else { // transaction aborted
        if (--retries >= 0) goto retry;
        bool result = KCAS_locks<K>(addr, expv, newv);
        return result;
    }
}
```

Fast path: try-lock based KCAS code inside a transaction

Fallback path: try-lock based KCAS code

(We use the same **version-lock-based KCASRead** for both code paths)

KCAS_txn is initially **the same as KCAS_locks**. We will optimize it.

Before retrying, should wait a bit...

Idea: if we aborted because a lock was held, could we wait until that lock is released?

But that information was lost when we aborted...

Maybe we could sneak it out using an explicit `_xabort()` status code?

Could pass an index < 127 into `addr[]` to `_xabort!`

TURNING LOCK ACQUISITION INTO READING

```
1  template <int K>
2  bool KCAS_txn(addr_t ** addr, val_t * expv, val_t * newv) {
3      retry:
4      TryLockSet<K> tls;          // auto-unlock when this goes out of scope
5      for (int i=0;i<K;++i) {    // try to acquire locks
6          if (!tls.tryLock(addr[i]->addrOfLock()))
7              goto retry;       // (release locks and) retry
8          if (addr[i]->value != expv[i])
9              return false;     // (release locks and) fail
10     }
11     for (int i=0;i<K;++i) {
12         *addr[i] = newv[i];
13     }
14     return true;               // (release locks and) succeed
15 }
```

No need to lock!
HTM ensures atomicity.
However, we must **respect others' locks**.
This becomes a **read**.

But, be careful! Still need to increment **version numbers** so fallback path can tell we changed addresses!

With no locking, we don't need TryLockSet

We modify the versioned try-lock implementation to **expose the "lock" word...**

TURNING RETRIES INTO ABORTS

```
1  template <int K>
2  bool KCAS_txn(addr_t ** addr, val_t * expv, val_t * newv) {
3      retry:
4
5      for (int i=0;i<K;++i) {
6          if (addr[i]->addrOfLock()->lock & 1)
7              goto retry;
8          if (addr[i]->value != expv[i])
9              return false;
10     }
11     for (int i=0;i<K;++i) {
12         *addr[i] = newv[i];
13         addr[i]->addrOfLock()->lock += 2; // update version numbers
14     }
15     return true;
16 }
```

No point retrying here, since we can only make progress if the lock is released, **which will abort us**. So, just **xabort**.

So, no need for the **“retry:”** label anymore.

OPTIMIZED CODE

```
1  template <int K>
2  bool KCAS_txn(addr_t ** addr, val_t * expv, val_t * newv) {
3
4
5     for (int i=0;i<K;++i) {
6         if (addr[i]->addrOfLock()->lock & 1)
7             _xabort(i);
8         if (addr[i]->value != expv[i])
9             return false;
10    }
11    for (int i=0;i<K;++i) {
12        *addr[i] = newv[i];
13        addr[i]->addrOfLock()->lock += 2; // update version numbers
14    }
15    return true;
16 }
```

User defined `_xabort` status
code (assuming `i < 127`)

CLEANED UP THE SPACING/COMMENTS

```
1  template <int K>
2  bool KCAS_txn(addr_t ** addr, val_t * expv, val_t * newv) {
3      for (int i=0;i<K;++i) {
4          if (addr[i]->addrOfLock()->lock & 1) _xabort(i); // read locks
5          if (addr[i]->value != expv[i]) return false; // read values
6      }
7      for (int i=0;i<K;++i) {
8          *addr[i] = newv[i]; // update values
9          addr[i]->addrOfLock()->lock += 2; // update version numbers
10     }
11     return true;
12 }
```

USING OUR NEW **_XABORT** STATUS CODE

```
1  template <int K>
2  bool KCAS(addr_t ** addr, val_t * expv, val_t * newv) {
3      int retries = 5;
4      int status;
5  retry:
6      if ((status = _xbegin()) == _XBEGIN_STARTED) {
7          bool result = KCAS_txn(addr, expv, newv);
8          _xend();
9          return result;
10     } else { // transaction aborted
11         if (status & _XABORT_EXPLICIT) {
12             int idx = _XABORT_CODE(status);
13             while (addr[idx]->addrOfLock()->read() & 1) { /* wait */ }
14         }
15         if (--retries >= 0) goto retry;
16         bool result = KCAS_locks<K>(addr, expv, newv);
17         return result;
18     }
19 }
```

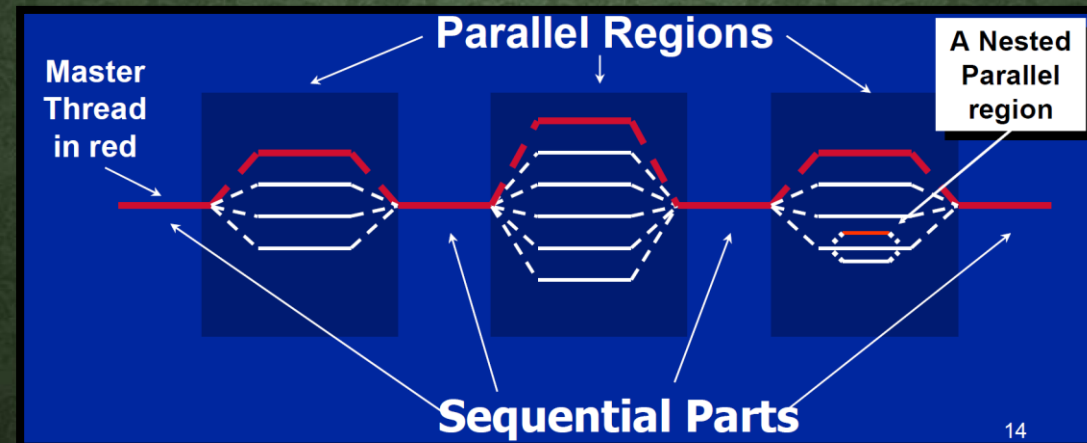
If the abort was caused
by a **user _xabort** call

Get the user-provided abort code

Wait until the returned
index is no longer locked

OPENMP (OPEN MULTI PROCESSING)

A library for fork-join parallelism



OVERVIEW

- Easy to use in your own C/C++ projects
- Tons of features; we just look at a couple of simple tools
 1. Parallel sections
 2. Parallel for loops
 3. Reductions
 4. OMP Single, OMP Tasks
- Usage:
 - `#include <omp.h>`
 - GCC: compile with **-fopenmp**

Warning: **Windows Subsystem for Linux v1** appears to use a **global lock** in its implementation of OpenMP.
It offered **no speedup** in my testing.
(Try running on a real Linux box or a virtual machine!)

1. PARALLEL SECTION

- Shortcut for: **spawning n threads**, where $n = \#$ of logical processors in the system, having them all execute the **same code block**, and then **joining them**

Sequential code

```
#include <stdio>
void main() {
    int id = 0;
    printf(" hello(%d)", id);
    printf(" world(%d)\n", id);
}
```

Output:
" hello(0) world(0)"

Concurrent/parallel code

```
#include <omp.h>
void main() {
    #pragma omp parallel
    {
        int id = omp_get_thread_num();
        printf(" hello(%d)", id);
        printf(" world(%d)\n", id);
    }
}
```

Spawn n threads

What is *my* thread ID (in 0..n-1)?

n threads execute this

Stop n threads

Output?

PARALLEL OUTPUT

```
hello(119) hello(69) hello(95) hello(15) hello(68) hello(32) hello(77) hello(61) hel.  
hello(138) hello(5) world(119)  
world(95)  
world(15)  
hello(67) world(138)  
hello(131) hello(27) hello(42) hello(80) world(115)  
hello(112) hello(78) hello(110) hello(113) world(131)  
hello(23) hello(109) hello(26) hello(12) world(110)  
hello(16) hello(35) hello(55) hello(13) hello(38) world(16)  
hello(45) hello(99) hello(105) world(45)  
hello(40) world(40)  
hello(4) world(4)  
hello(3) world(3)  
hello(82) world(82)  
hello(1) world(1)  
hello(11) world(11)  
hello(21) world(21)  
hello(39) world(39)  
hello(81) world(81)  
hello(22) world(22)  
hello(2) world(2)  
hello(9) world(9)  
hello(17) world(17)  
hello(83) world(83)  
hello(36) world(36)  
hello(10) world(10)
```

WHY IS THIS USEFUL?

Using OpenMP

```
#pragma omp parallel  
doSomething();
```

Equivalent code using *pthread*s

```
const int n = SomehowGetNumLogicalProcessors();  
pthread_t *threads = new pthread_t[n];  
for (int i=0;i<n;++i) {  
    if (pthread_create(&threads[i], NULL, doSomething)) {  
        std::cerr<<"ERROR: could not create thread"<<std::endl;  
        // ...  
    }  
}
```

Equivalent code using *std::thread*

```
1  const int n = SomehowGetNumLogicalProcessors();  
2  std::thread * threads = new std::thread[n];  
3  for (int i=0;i<n;++i) {  
4      threads[i] = new std::thread(doSomething)  
5  }  
6  for (int i=0;i<n;++i) {  
7      threads[i]->join();  
8  }  
9  delete[] threads;
```

```
{  
    thread"<<std::endl;  
}
```

2. PARALLEL FOR

```
8 matrix * multiply(matrix * o) {
9     matrix * ret = new matrix(h, o->w);
10    #pragma omp parallel for
11    for (int y=0; y < ret->h; ++y) {
12        for (int x=0; x < ret->w; ++x) {
13            for (int k=0; k < w; ++k) {
14                ret->data[y][x] += data[y][k] * o->data[k][x];
15            }
16        }
17    }
18    timer.split("multiply call finished");
19    return ret;
20 }
```

Implicit fork
(thread spawning)

Implicit join

OpenMP *automatically* decides
how many threads (*n*) to spawn

It splits work evenly between the
n threads (each does approx.
 $\frac{ret \rightarrow h}{n}$ iterations of this loop)

WHY IS THIS USEFUL?

Using OpenMP

```
#pragma omp parallel for  
for (long i=0;i<n;++i)  
    loop_body(i);
```

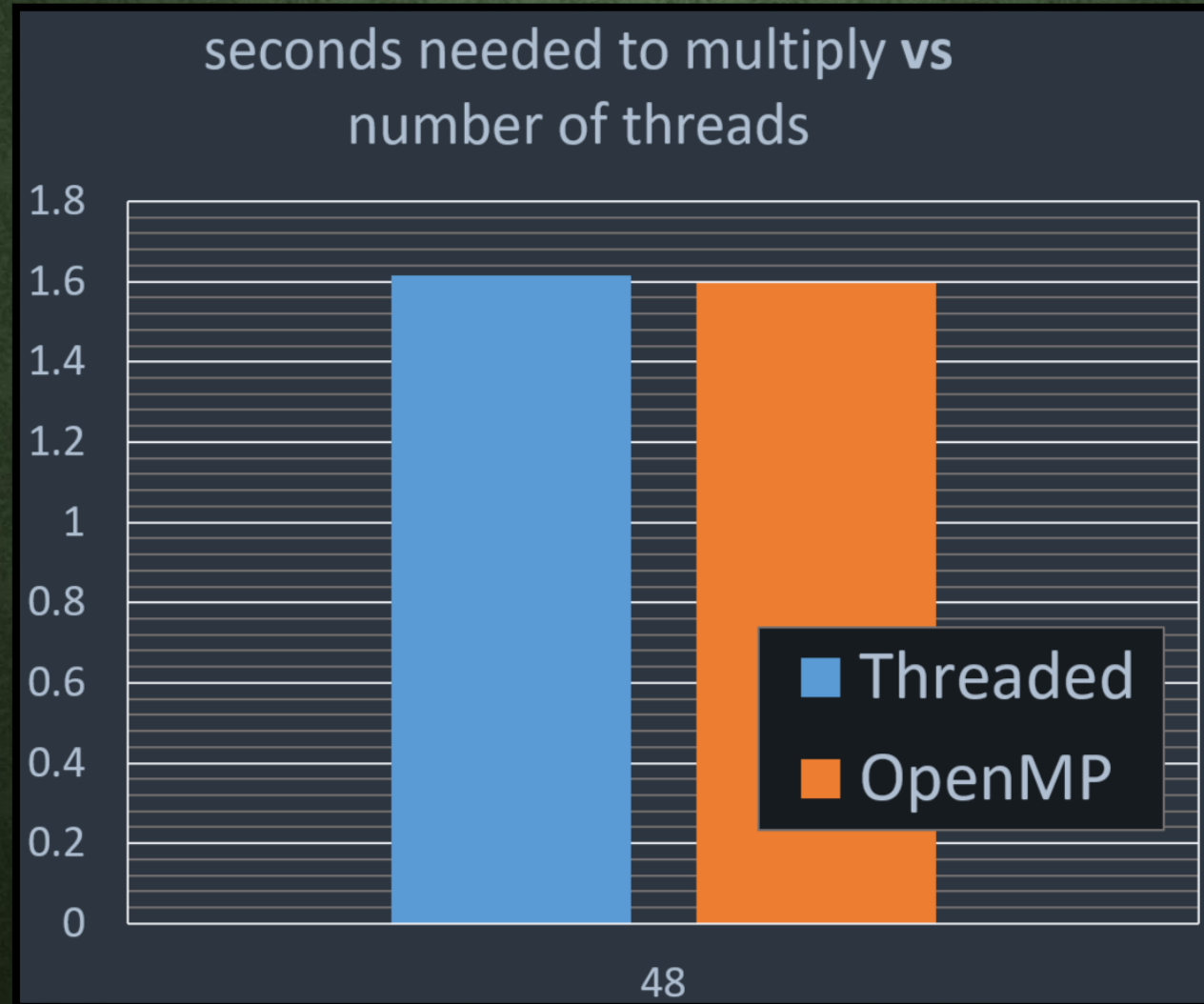
Without using OpenMP: starting threads to run loop_body

```
1  const int numThreads = GetNumberOfLogicalProcessors();  
2  std::thread * threads = new std::thread[numThreads];  
3  =for (int i=0;i<numThreads;++i) {  
4      threads[i] = new std::thread(loop_body);  
5  }  
6  =for (int i=0;i<numThreads;++i) {  
7      threads[i].join();  
8  }  
9  delete[] threads;
```

Without using OpenMP: loop_body

```
int id = getMyThreadID();  
int istart = id * n / numThreads;  
int iend = (id+1) * n / numThreads;  
if (id == numThreads-1) iend = n;  
for (int i=istart; i<iend; i++)  
    original_loop_body(i);
```

PERFORMANCE OF OPENMP VS MANUAL THREADING



3. REDUCTIONS

- A reduction takes a vector (array) and turns it into a scalar (single number)

Example: summing an array

```
long sum = 0;
for (long i=0;i<n;++i) {
    sum += array[i];
}
```

Naïve use of OpenMP

```
long sum = 0;
#pragma omp parallel for
for (long i=0;i<n;++i) {
    sum += array[i];
}
```

Problem: many threads do this, and this increment is not atomic!

How about using fetch&add?

```
long sum = 0;
#pragma omp parallel for
for (long i=0;i<n;++i) {
    __sync_fetch_and_add(&sum, array[i]);
}
```

Problem: correct, but not very scalable!

Ideally: want a thread to maintain a **local sum** while processing a **batch**, and fetch&add its local sum into the **global sum** at the end of the batch

OPENMP REDUCTIONS

- OpenMP natively supports reductions over numerous operators (+, *, &, |, ...)
- Must tell OpenMP which variable will be used to store the reduction

Example: summing an array

```
long sum = 0;
for (long i=0;i<n;++i) {
    sum += array[i];
}
```

1 thread, time to complete: **4330 ms**

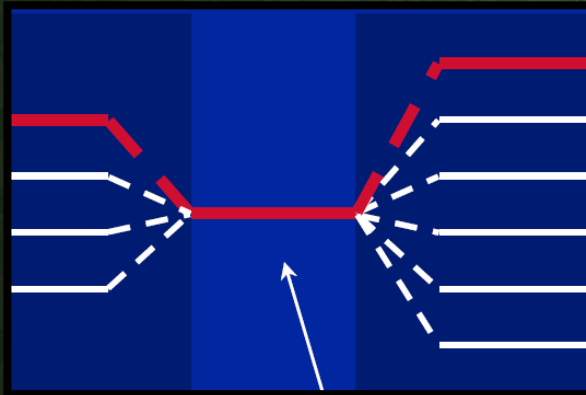
Correct OpenMP reduction

```
long sum = 0;
#pragma omp parallel for reduction (+:sum)
for (long i=0;i<n;++i) {
    sum += array[i];
}
```

48 threads, time to complete: **185ms**

4(A). OPENMP SINGLE

- Sometimes you want a **single threaded** computation in the middle of your parallel computation



```
1  #pragma omp parallel
2  = {
3      multiThreadedStuff1 ();
4
5      #pragma omp single
6      = {
7          ..... singleThreadedStuff ();
8      }
9
10     multiThreadedStuff2 ();
11 }
```

4(B). OPENMP TASK

- Define a **task** that should be completed by **any available thread** in a parallel section
- Common design pattern: one thread generates & launches tasks, tasks run in parallel

```
1  #pragma omp parallel
2  = {
3      #pragma omp single
4      = {
5      =     for (linked_list_node n = head; n != NULL; n = n->next) {
6          #pragma omp task
7          processTheNode (n);
8          }
9      }
10 }
```

Spin up **n** threads to run tasks

Single thread generates tasks

This function call, with this argument, becomes a **task**, and is run in the background (by one of the **n-1 other threads**)

This closing brace **waits** for all tasks to complete

Can also **manually** wait for all tasks, whenever you like, with `#pragma omp taskwait`

OMP task tutorial: <https://openmp.org/wp-content/uploads/sc13.tasking.ruud.pdf>

TOOLS FOR DEBUGGING AND PERFORMANCE

- Debugging and optimizing concurrent programs is **very** hard. Tools can help!

- Debugging

- GNU Debugger (GDB)
 - Segfaults, infinite loops
- Address Sanitizer (ASan)
 - Segfaults, memory leaks
 - 1~2x slowdown
- **Valgrind**
 - Segfaults, memory leaks, **memory access errors**
 - many-x slowdown
- Graphviz
 - **Visualizing** pointer based data structures

- Performance

- Linux Perftools (perf)
 - Studying cycles, cache misses, instructions, stalled cycles
 - At the whole-application level
- C/C++ Performance API (PAPI)
 - Precise information from perf, but recorded **within** your program
- VTune Amplifier
 - Powerful (and now free!) profiler

A lot of errors in concurrent programs manifest as memory access errors! For example, a thread may write a bad value into a pointer because of a concurrency bug, and another thread may then read it.

DEBUGGING TOOLS

USING VALGRIND TO FIND MEMORY ACCESS ERRORS

```
$ valgrind --fair-sched=yes ./alcode_segfault/workload_timed.out 4 1000 naive
==107893== Command: ./alcode_segfault/workload_timed.out 4 1000 naive
==107893== Use of uninitialised value of size 8
==107893==    at 0x510F0D4: std::thread::join() (in /.../x86_64-linux-gnu/libstdc++.so.6.0.22)
==107893==    by 0x1092DC: void runExperiment<CounterNaive>(...) (workload_timed.cpp:46)
==107893==    by 0x108E3B: main (workload_timed.cpp:70)
==107893==
==107893== Invalid read of size 8
==107893==    at 0x510F0D4: std::thread::join() (in /.../x86_64-linux-gnu/libstdc++.so.6.0.22)
==107893==    by 0x1092DC: void runExperiment<CounterNaive>(...) (workload_timed.cpp:46)
==107893==    by 0x108E3B: main (workload_timed.cpp:70)
==107893== Address 0x190 is not stack'd, malloc'd or (recently) free'd
...
```

- Typical first step in debugging any error that isn't obvious:
 - Ensure that valgrind runs without any such errors.
 - If there are such errors, **fix those first!**

Using Address Sanitizer to check for memory leaks

```
$ g++ -pthread -g -fsanitize=address -static-libasan -fopenmp -O3 ex2_mmult_threads.cpp
$ ./a.out 24
matrix created: 0.02s
randomize call finished: 0.10s
...
multiply call finished: 2.82s

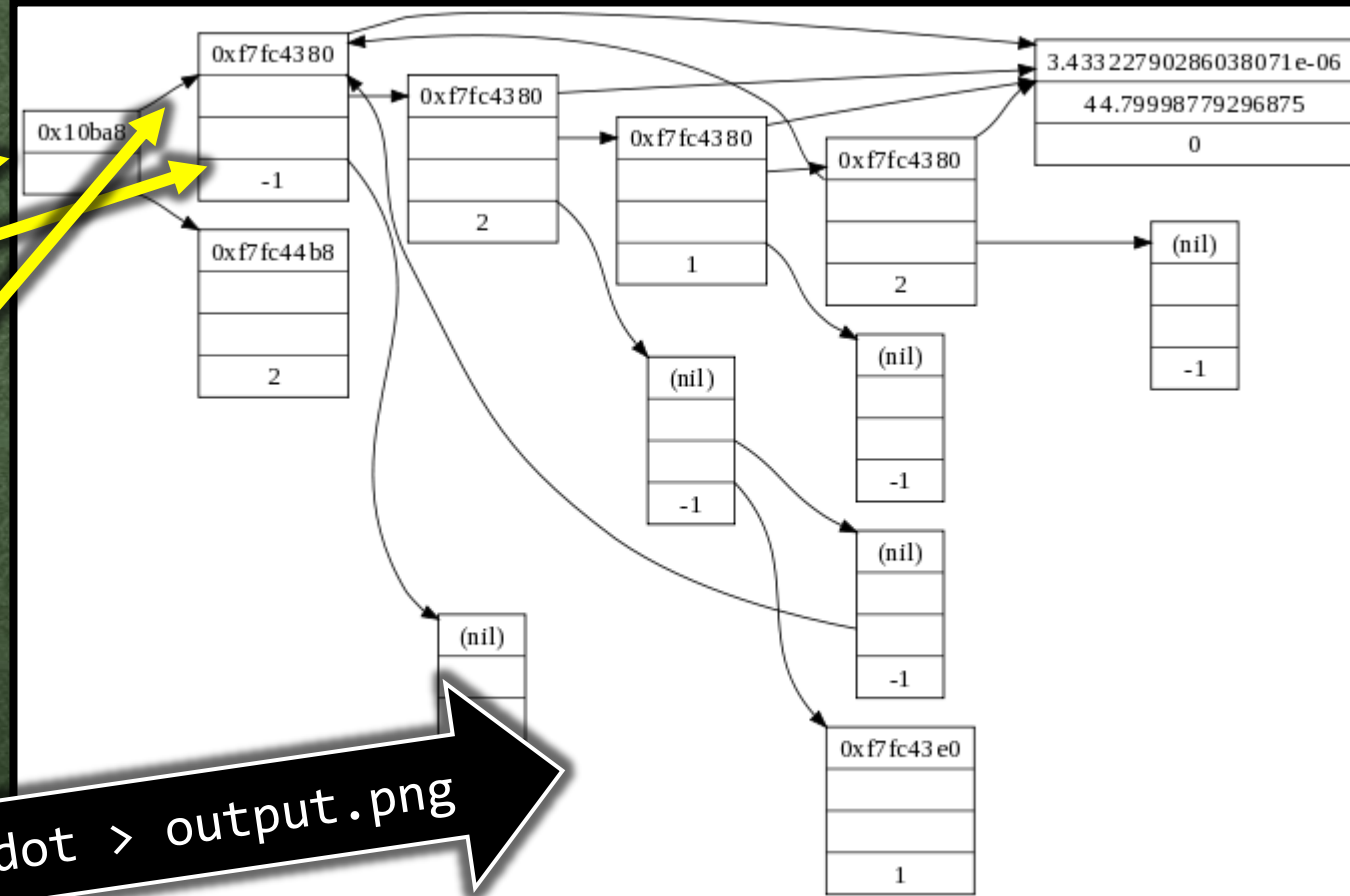
=====
==76549==ERROR: LeakSanitizer: detected memory leaks

Direct leak of 192 byte(s) in 24 object(s) allocated from:
    #0 0x555b294992d8 in operator new(unsigned long) (/home/.../a.out+0xc82d8)
    #1 0x555b294dab9e in matrix::multiply(matrix*, int) (/home/.../ex2_mmult_threads.cpp:12)

SUMMARY: AddressSanitizer: 192 byte(s) leaked in 24 allocation(s).
```

GRAPHVIZ: WHEN YOU JUST NEED TO SEE IT

```
digraph g {
  node [
    fontsize = "16"
    shape = "record"
  ];
  edge [];
  "node0" [
    label = "<f0> 0x10ba8|<f1>"
  ];
  "node1" [
    label = "<f0> 0xf7fc4380|<f1>|<f2>|-1"
  ];
  [...]
  "node0":f0 -> "node1":f0 [
    id = 0
  ];
  [...]
}
```



`$ dot -Tpng input.dot > output.png`

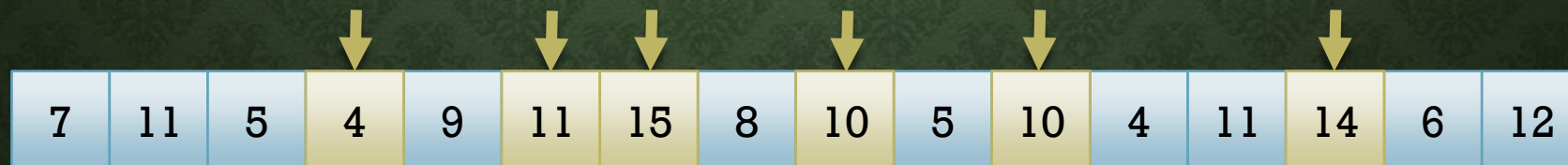
Tip: try to get your toughest bugs to happen in SMALL data structures, so you can graphviz them

SANITY CHECKING: EXPERIMENT CHECKSUMS

- Important to perform sanity checks wherever you can!
 - Helps to catch obvious (and non-obvious) mistakes
- One good sanity check: checksum based validation
 - Reduce the **data structure** to a number (a **data structure checksum**)
 - Reduce each threads' completed operations to a number (a **thread checksum**)
 - verify that thread checksums "**match**" the data structure checksum
 - (I.e., the work the threads **think** they've done is reflected **in the data structure!**)
- **Creativity needed to come up with good checksum functions**

EXAMPLE: SYNTHETIC KCAS BENCHMARK

- **n** threads repeatedly do the following for 3 seconds
 - Pick **K** uniform random slots in an array
 - Read integers stored in those slots
 - Do a KCAS to change each of the **K** slots from the value **exp** that we read, to a new value of **exp + 1**
- Report average throughput (KCAS operations/sec) over all trials



EXAMPLE: CHECKSUM VALIDATION FOR SUCH A BENCHMARK

- **Data structure checksum**

- Sum of all array entries
- Each successful KCAS increments k array slots by 1
 - Should add k to the **data structure checksum**

- **Thread checksum**

- kX where $X = \#$ of successful KCAS operations performed by the thread
 - Each successful KCAS should add k to the **thread checksum**

- **Validation**

- $\text{sum}(\text{thread checksums}) == \text{data structure checksum} ?$
- (If a KCAS operation is lost, or screws up the array, validation [hopefully] fails)