# MULTICORE PROGRAMMING

Epoch-based memory reclamation and experimental methodology
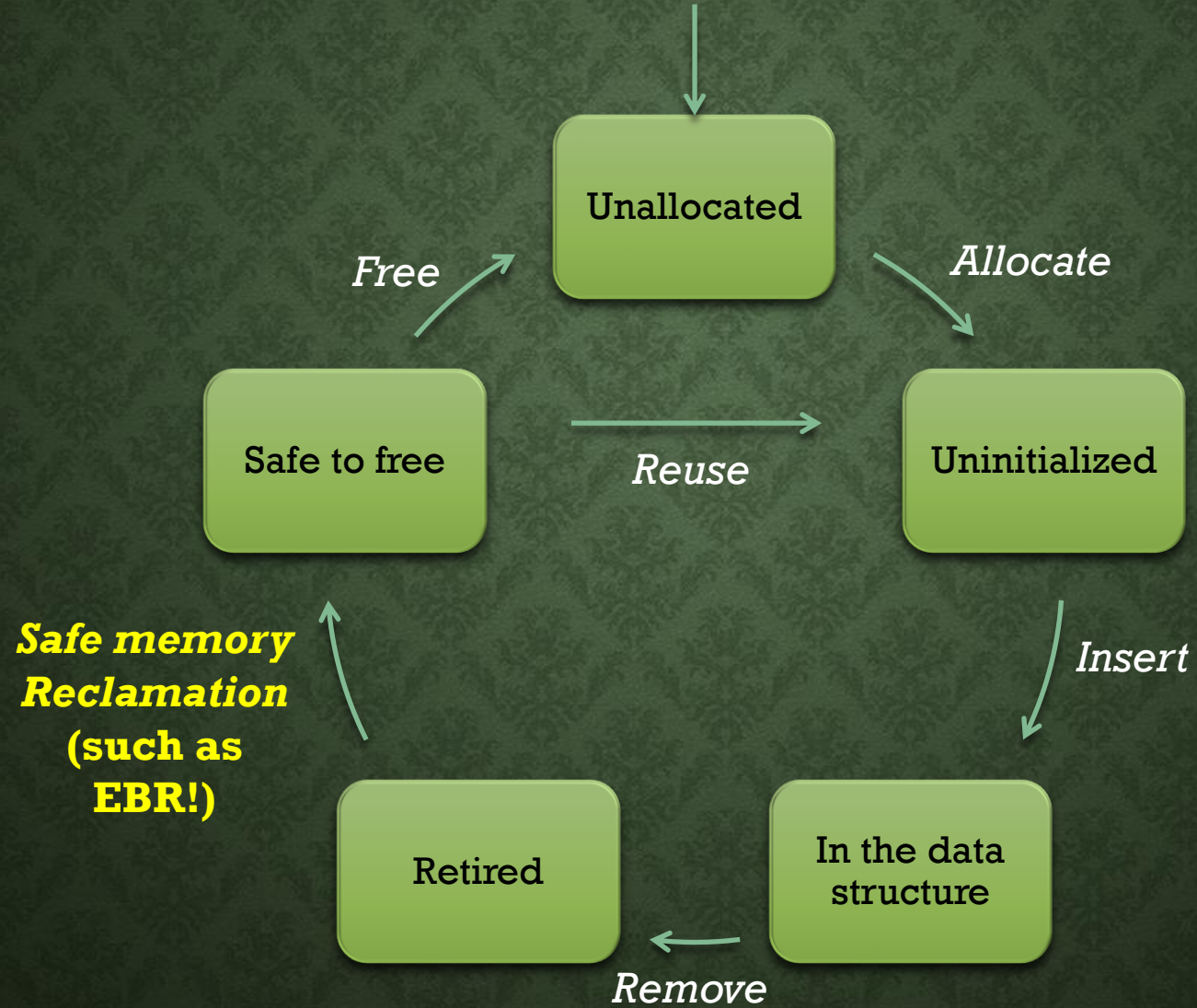
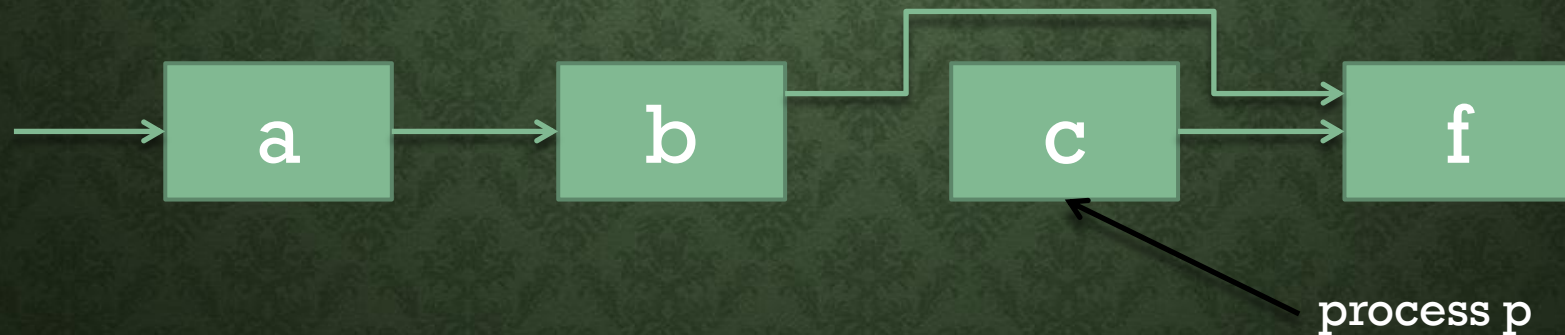## Lecture 17

Trevor Brown

# THIS TIME

- **Epoch-based memory reclamation**

    - (**The <u>algorithm</u> itself**, as well as a bit more on usage)

- Time permitting:

    - A discussion on experimental methodology
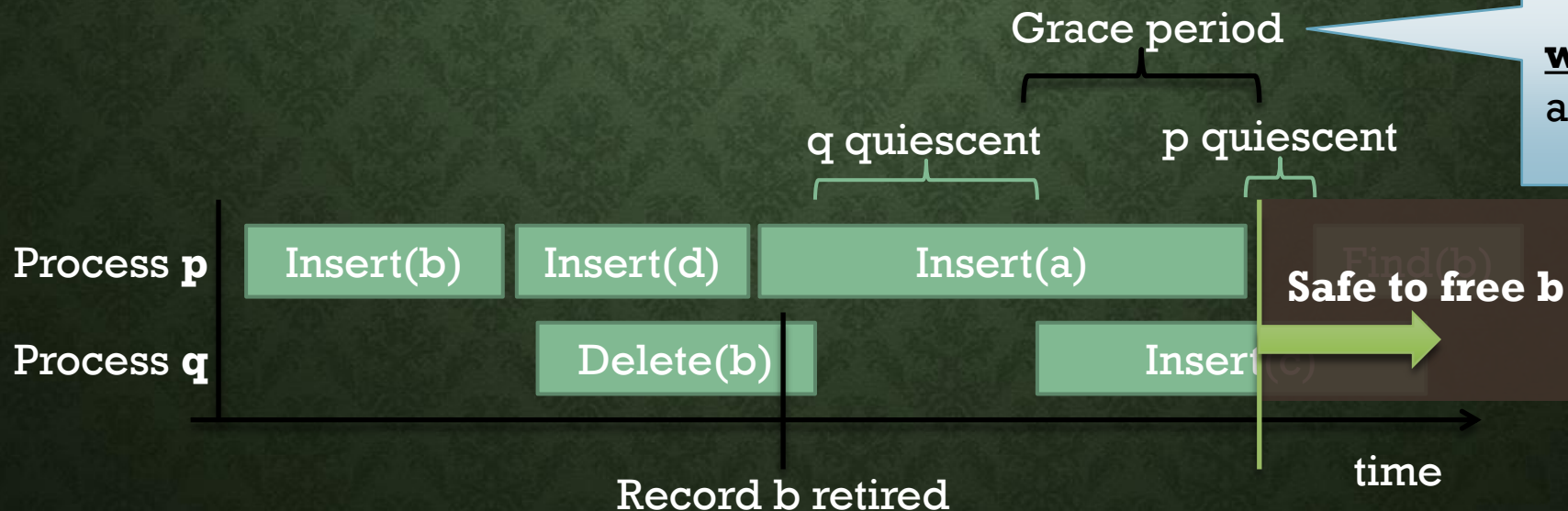
# LIFECYCLE OF A RECORD (OBJECT)

# RECLAMATION WITH AND WITHOUT LOCKS

- Easy with locks: correct locking ensures that
  no process can access an unlinked record

- Hard without locks: processes must carefully coordinate to avoid
  accessing freed nodes

  - Challenge: any record you are about to free
    could be pointed to by another process



process p

# QUIESCENCE AND GRACE PERIODS

- **Definition:** a process is **quiescent** iff

  its private memory does **not** contain any pointers to records in the data structure

- **Definition:** a **grace period** is an interval **during which**

  each thread has **some time** when it is quiescent

  (different threads can be quiescent at different times---that is fine)

- **Fact:** a retired record can be freed after a **subsequent** grace period

Grace period

q quiescent          p quiescent

But how to **detect when** there has been a grace period **since** b was **retired**?

Process **p**    Insert(b)    Insert(d)    Insert(a)

**Safe to free b**

If we can **detect** this, we can know it is safe to **free** b

Process **q**    Delete(b)    Insert
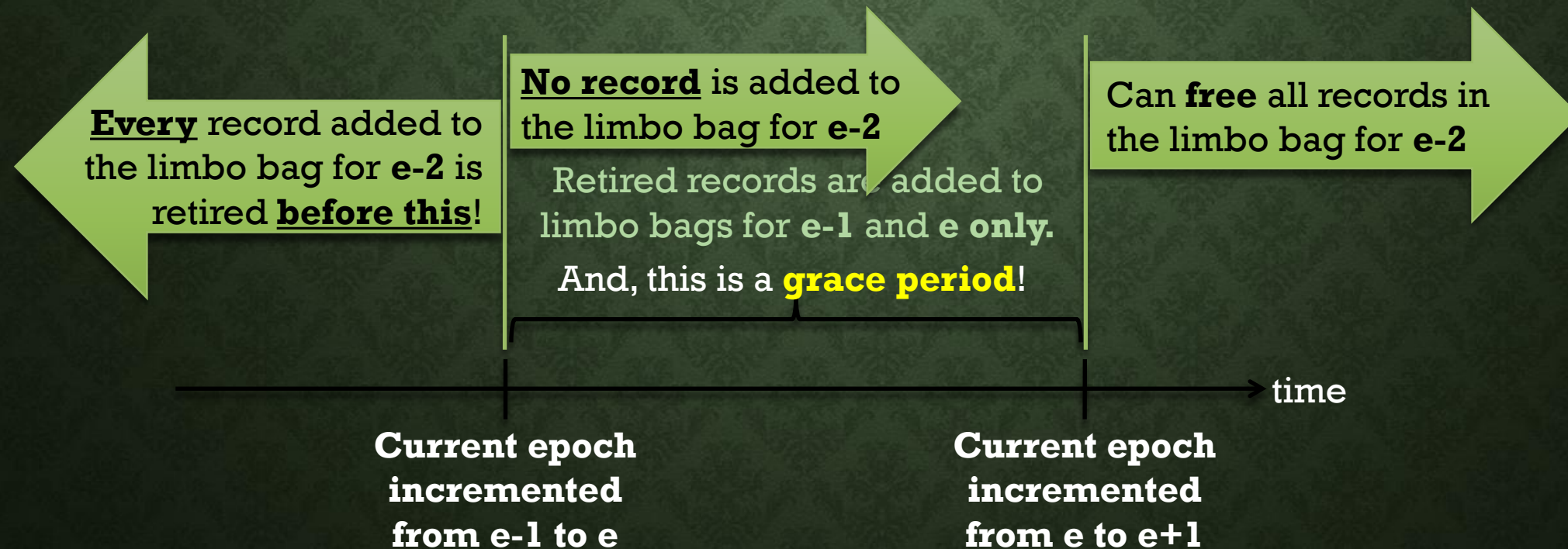
Record b retired

time

# USING EPOCHS TO DETECT GRACE PERIODS

- **Key assumption:** threads are **quiescent** when **not** executing data structure operations

- The execution is divided into *epochs*,
  and the current epoch number is stored (as a global) in shared memory

- At the start of **each** data structure operation:

  1. read current epoch and **announce** it (in a global array with one slot per thread)

  2. check whether all other threads have announced it

  3. if so, **increment** the current epoch

**Before** this time, **no** thread has announced epoch e-1

This is a grace period!

**Each** thread must have announced e-1 by this time

(And each thread is quiescent **just before** it starts this data structure operation…)

time

Current epoch incremented from e-2 to e-1

Current epoch incremented from e-1 to e

Therefore, **each** thread starts a new data structure operation and announces e-1 **at some point between these two times**!

How to figure out **which objects** were retired **before** this grace period?

Recall: anything **retired before** a grace period can be **freed after** it!

# DETERMINING **WHICH RECORDS** ARE SAFE TO FREE

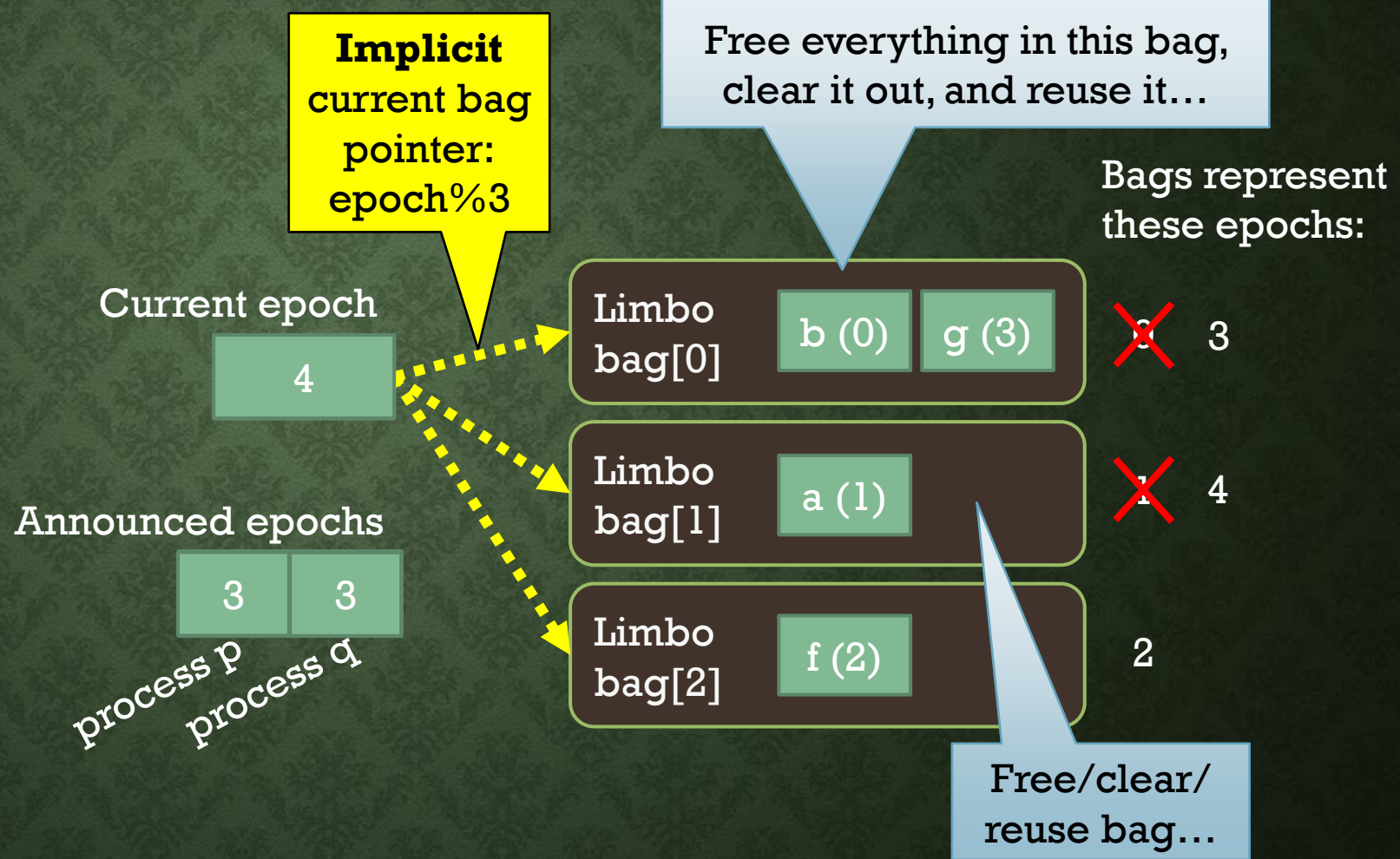- Maintain a shared **limbo bag** for <u>each</u> of the **last 3 epochs**

- Records retired by thread p are added to the bag for the last epoch <u>**announced**</u> by p

  - So, if p reads epoch 7, then retires a node u, it will place u in the limbo bag for epoch 7

  - Note: by the time p retires u, the current epoch might be 8 (but it **cannot** be 9 or larger... why?)

- When the current epoch changes from e to e+1,
  we free all records in the limbo bag for e-2 (lets see why...)

**<u>Every</u>** record added to the limbo bag for **e-2** is retired **<u>before this</u>**!

**<u>No record</u>** is added to the limbo bag for **e-2**

Retired records are added to limbo bags for **e-1** and **e only.**

And, this is a **grace period**!

Can **free** all records in the limbo bag for **e-2**

time

**Current epoch incremented from e-1 to e**

**Current epoch incremented from e to e+1**

# IMPLEMENTATION: DATA TYPES

```cpp
1   template <typename T>
2   class alignas(64) limbo_bag {
3       ... implementation details skipped ...
4   public:
5       void add(T * obj);
6       void freeAll();
7   };
8
9   struct alignas(64) padded_aint {
10      atomic<int>        v;
11  }
12
13  template <typename T, int NUM_THREADS>
14  class alignas(64) ebr_manager {
15  private:
16      padded_aint            announce[NUM_THREADS];
17      padded_aint            currEpoch;
18      limbo_bag<T>           bags[3];
19      char                   padding[64];
20  public:
21      ebr_manager();
22      void startOp();
23      void retire(T * obj);
24  };
25
26  thread_local int tid; // assuming this is set
```

**alignas** to serve as padding

**Sequential** implementation

Note: freeAll should cause object **destructors** to be called just before the objects are freed

Type T of objects to manage/reclaim is templated

Threads' epoch announcements represent the epoch they are **conceptually running in**

Limbo bags for the current epoch and the previous two epochs

API: let's see how these are implemented

Another option instead of passing tid to every function

# IMPLEMENTATION: EASY PARTS

```
28  template <typename T, int NUM_THREADS>
29  ebr_manager<T,NUM_THREADS>::ebr_manager() {
30      static_assert(alignof(announce[0]) == 64);
31      static_assert(alignof(announce[1]) == 64);
32      static_assert(alignof(currEpoch) == 64);
33      static_assert(alignof(bags[0]) == 64);
34      static_assert(alignof(bags[1]) == 64);
35  }
36
37  template <typename T, int NUM_THREADS>
38  void ebr_manager<T,NUM_THREADS>::retire(T * obj) {
39      assert(tid >= 0 && tid < NUM_THREADS);
40      int currBag = announce[tid] % 3;
41      bags[currBag].add(obj);
42  }
```

Paranoia: double check that alignment really implies data is **padded** the way we think it is

Paranoia: check that our tid is a valid one

announce[tid] is the epoch we are conceptually running in, and this modulo 3 is the limbo bag **for that epoch**. We put obj in that bag.

# THE HARD PART (NOT OPTIMIZED)

```cpp
44   template <typename T, int NUM_THREADS>
45   void ebr_manager::startOp() {
46       assert(tid >= 0 && tid < NUM_THREADS);
47
48       // read & announce the epoch that we will "run in"
49       int seen = -1; /* -1 means epoch is "locked" */
50       while (seen == -1) seen = currEpoch;
51       if (seen != announce[tid]) announce[tid] = seen;
52
53       // try to the advance currEpoch
54       for (int i=0;i<NUM_THREADS;++i) {
55           if (announce[i] != seen) return; // can't advance it
56       }
57       int exp = seen;
58       if (currEpoch.compare_exchange(exp, -1 /* locked */)) {
59           // we "win" the right to freeAll in the new bag
60           // (while we freeAll, the "lock" prevents add()s)
61           int newBag = (seen+1) % 3;
62           bags[newBag].freeAll();
63           currEpoch = seen+1; // "unlock" with the new epoch
64       }
65   }
```

Paranoia: check that our tid is a valid one

currEpoch can conceptually be "locked" by a thread, who temporarily sets it to -1. So, we repeatedly read until we see an "unlocked" value.

If we see a **new** epoch, we update our announcement

Check if all threads have announced **seen**

If so, we will try to increment currEpoch **and** reclaim memory

This involves "locking" currEpoch, freeing contents of the new current bag, and "unlocking" currEpoch to the new epoch value

```
44  template <typename T, int NUM_THREADS>
45  void ebr_manager::startOp() {
46      assert(tid >= 0 && tid < NUM_THREADS);
47
48      static thread_local numCalls = 0;
49      ++numCalls;
50
51      // read & announce the epoch that we will "run in"
52      int seen = -1; /* -1 means epoch is "locked" */
53      while (seen == -1) seen = currEpoch;
54      if (seen != announce[tid]) announce[tid] = seen;
55
56      // only try to advance the epoch after many operations
57      const int THRESHOLD = max(100, 10*NUM_THREADS);
58      if ((numCalls % THRESHOLD) == 0) {
59          for (int i=0;i<NUM_THREADS;++i) {
60              if (announce[i] != seen) return; // can't advance it
61          }
62          int exp = seen;
63          if (currEpoch.compare_exchange(exp, -1 /* locked */)) {
64              // we "win" the right to freeAll in the new bag
65              // (while we freeAll, the "lock" prevents add()s)
66              int newBag = (seen+1) % 3;
67              bags[newBag].freeAll();
68              currEpoch = seen+1; // "unlock" with the new epoch
69          }
70      }
71  }
```

Each thread keeps track of how many calls it has ever made to startOp()

**static:** accessible only within the function, but preserves its value across function calls ("= 0" happens **once**)

We try to advance the epoch only after performing THRESHOLD operations

**Tradeoff:** it takes longer to reclaim garbage, but we amortize the overhead of scanning announcements over many operations

Caveat: I have not run this code, and have not spent time tuning THRESHOLD

# ADDING AN RAII GETGUARD() OPERATION

```cpp
template <typename T, int NUM_THREADS>
class alignas(64) ebr_manager {
private:
    padded_aint            announce[NUM_THREADS];
    padded_aint            currEpoch;
    limbo_bag<T>           bags[3];
    char                   padding[64];
public:
    ebr_manager();
    void startOp();
    void retire(T * obj);

    class mem_guard {
        ebr_manager<T,NUM_THREADS> * mgr;
    public:
        mem_guard(ebr_manager<T,NUM_THREADS> * _mgr) : mgr(_mgr) {
            mgr->startOp();
        }
        ~mem_guard() {
            // if a "mgr->endOp()" were defined, we'd call it here
        }
    };
    mem_guard getGuard() {
        return mem_guard(this);
    }
};
```

Note: getGuard() has somewhat **limited usefulness** when there is no endOp() operation...

Just showing you how to do this RAII design pattern...

# USAGE EXAMPLE 1: RECALL THE TREIBER STACK

```
1  class stack {
2      class node {
3          int                 key;
4          node *              next;
5      };
6      static const int        EMPTY = -1;
7      char                    padding1[64];
8      atomic<node *>          top;
9      char                    padding2[64];
10 public:
11     stack() : top(NULL) {}
12     void push(int key);
13     int pop();
14 };
```

```
16 void stack::push(int key) {
17     node * n = new node();
18     n->key = key;
19     while (true) {
20         node * curr = top;
21         n->next = curr;
22         node * exp = curr;
23         if (top.compare_exchange(exp, n))
24             return;
25     }
26 }
27
28 void stack::pop() {
29     while (true) {
30         node * curr = top;
31         if (curr == NULL) return EMPTY;
32         node * next = curr->next;
33         node * exp = curr;
34         if (top.compare_exchange(exp, next))
35             return curr->key;
36     }
37 }
```

# TREIBER STACK WITH MEMORY RECLAMATION

```cpp
1  class stack {
2      class node {
3          int                 key;
4          node *              next;
5      };
6      static const int        EMPTY = -1;
7      char                    padding1[64];
8      atomic<node *>          top;
9      char                    padding2[64];
10     ebr_manager<node>       mgr;
11     char                    padding3[64];
12 public:
13     stack() : top(NULL) {}
14     void push(int key);
15     int pop();
16 };
```

```cpp
18  void stack::push(int key) {
19      node * n = new node();
20      n->key = key;
21      while (true) {
22          auto guard = mgr.getGuard();
23          node * curr = top;
24          n->next = curr;
25          node * exp = curr;
26          if (top.compare_exchange(exp, n))
27              return;
28      }
29  }
30
31  void stack::pop() {
32      while (true) {
33          auto guard = mgr.getGuard();
34          node * curr = top;
35          if (curr == NULL) return EMPTY;
36          node * next = curr->next;
37          node * exp = curr;
38          if (top.compare_exchange(exp, next))
39              mgr.retire(curr);
40              return curr->key;
41      }
42  }
```

```
20    void hashmap::startExpansion(t)
21  =    if (currentTable == t) {
22         t_new = createNewTableStruct(t);
23         if not CAS(&currentTable, t, t_new) delete t_new;
24       }
25       helpExpansion(currentTable);
```

```
struct hashmap
1    char padding1[64];
2    atomic<table *> currentTable;
3    char padding2[64];
```

Want to reclaim object types:
hashmap, table, counter, atomic<int> array

```
struct table
1     char padding1[64];
2     atomic<int> * data;
3     atomic<int> * old;
4     int capacity;
5     int oldCapacity;
6     counter * approxIns;
7     counter * approxDel;
8     char padding2[64];
9     atomic<int> chunksClaimed;
10    char padding3[64];
11    atomic<int> chunksDone;
12    char padding4[64];
```

Suppose we allocate
new counter objects
for each table object

```
20   void hashmap::startExpansion(t)
21 ▭   if (currentTable == t) {
22        t_new = createNewTableStruct(t);
23        if not CAS(&currentTable, t, t_new) delete t_new;
24      }
25      helpExpansion(currentTable);
```

If CAS **succeeds**, then **we unlinked t**, so we have the right to call mgr->retire(t)

```
struct hashmap
1    char padding1[64];
2    atomic<table *> currentTable;
3    char padding2[64];
4    ebr_manager<table> mgr;
5    char padding3[64];
```

```
struct table
1    char padding1[64];
2    atomic<int> * data;
3    atomic<int> * old;
4    int capacity;
5    int oldCapacity;
6    counter * approxIns;
7    counter * approxDel;
8    char padding2[64];
9    atomic<int> chunksClaimed;
10   char padding3[64];
11   atomic<int> chunksDone;
12   char padding4[64];
```

Observation: if a table object **t** is safe to free, then no thread will ever access **t->old** or **t->approxIns** or **t->approxDel**

So those objects can be directly passed to free() whenever a table object is freed (carefully do this in ~table(), which is invoked by EBR)

# OPTIMIZING INTO A "REAL" ALGORITHM (ONE WORTH IMPLEMENTING)

# SIGNIFICANT CHANGE FROM EBR

- Per-thread **quiescent bit** to allow reclamation to **continue** while a process is quiescent

  - Useful if some threads finish their work and stop, or work on something else

  - **Partial** fault tolerance

    - Crashing **while quiescent** does not block reclamation

# EASY CHANGE: SCANNING EPOCH ANNOUNCEMENTS

- Amortize cost over several operations

- Each operation checks **one** epoch announcement
  (or you could check **one** announcement per **K** operations)

- After checking **n** announcements, where n is the number of threads,
  *and seeing the announcements are up to date*, the epoch can be advanced

# EFFICIENT BAGS

- Per-process limbo bags
  - Each process rotates its limbo bags whenever its announcement changes

- Per-process free bags and one shared free bag
  - When rotating its limbo bags,
    a process appends its oldest limbo bag to its own free bag
  - Entire blocks moved to/from shared free bag

- More details (and optimizations) in the paper
  Brown, T. Reclaiming memory for lock-free data structures. PODC 2015.
  - Conference paper
  - Extended paper
  - Slides for that talk

# COMPLEXITY OF DEBRA

- leaveQstate: O(1) steps          ← called at start of operation

- enterQstate: O(1) steps          ← called at end of operation

- retire: O(1) steps          ← called after **unlinking** a record

- Reclamation operations are **wait-free!**

 

- However, this does **not** mean reclamation is fault tolerant!

    - A thread that crashes while non-quiescent can still block reclamation!

    - This is addressed in the **DEBRA+** algorithm (same paper as DEBRA)

    - … and more recently with Ajay Singh and Ali Mashtizadeh in **Neutralization Based Reclamation, PPoPP 2021. [paper] [talk]**