

* Note that buffers can be in one of 6 states:

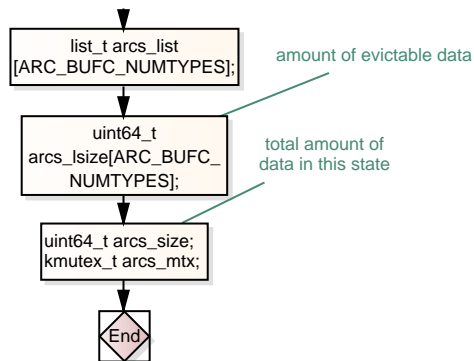
- * ARC_anon - anonymous (discussed below)
- * ARC_mru - recently used, currently cached
- * ARC_mru_ghost - recently used, no longer in cache
- * ARC_mfu - frequently used, currently cached
- * ARC_mfu_ghost - frequently used, no longer in cache
- * ARC_l2c_only - exists in L2ARC but not other states

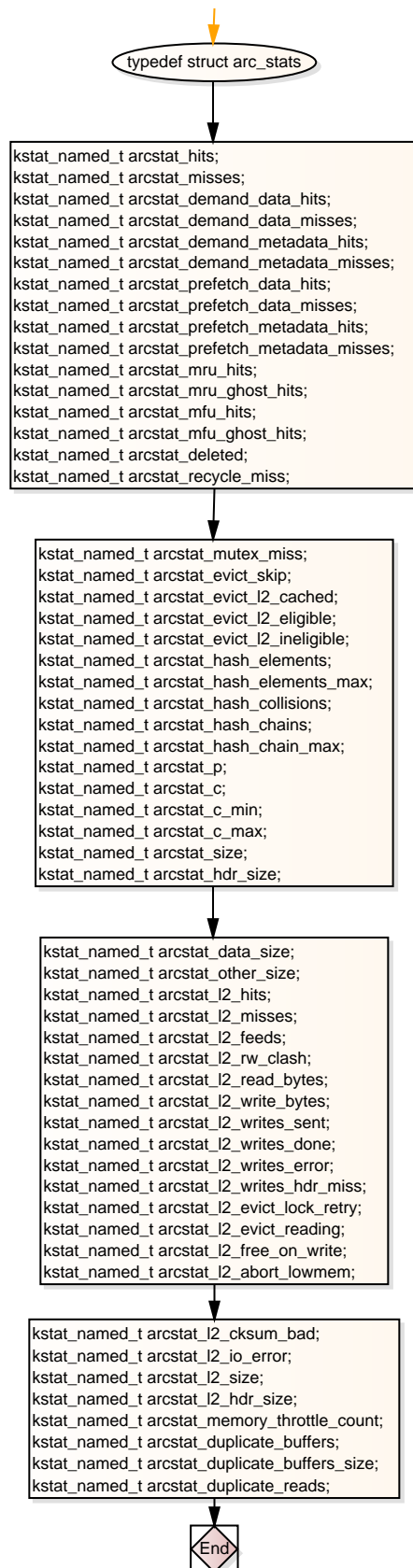
* When there are no active references to the buffer, they are
 * are linked onto a list in one of these arc states. These are
 * the only buffers that can be evicted or deleted. Within each
 * state there are multiple lists, one for meta-data and one for
 * non-meta-data. Meta-data (indirect blocks, blocks of dnodes,
 * etc.) is tracked separately so that it can be managed more
 * explicitly: favored over data, limited explicitly.

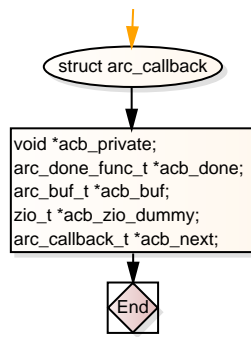
* Anonymous buffers are buffers that are not associated with
 * a DVA. These are buffers that hold dirty block copies
 * before they are written to stable storage. By definition,
 * they are "ref'd" and are considered part of arc_mru
 * that cannot be freed. Generally, they will acquire a DVA
 * as they are written and migrate onto the arc_mru list.

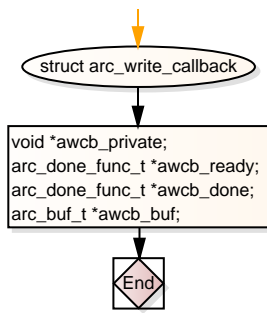
* The ARC_l2c_only state is for buffers that are in the second
 * level ARC but no longer in any of the ARC_m* lists. The second
 * level ARC itself may also contain buffers that are in any of
 * the ARC_m* states - meaning that a buffer can exist in two
 * places. The reason for the ARC_l2c_only state is to keep the
 * buffer header in the hash table, so that reads that hit the
 * second level ARC benefit from these fast lookups.

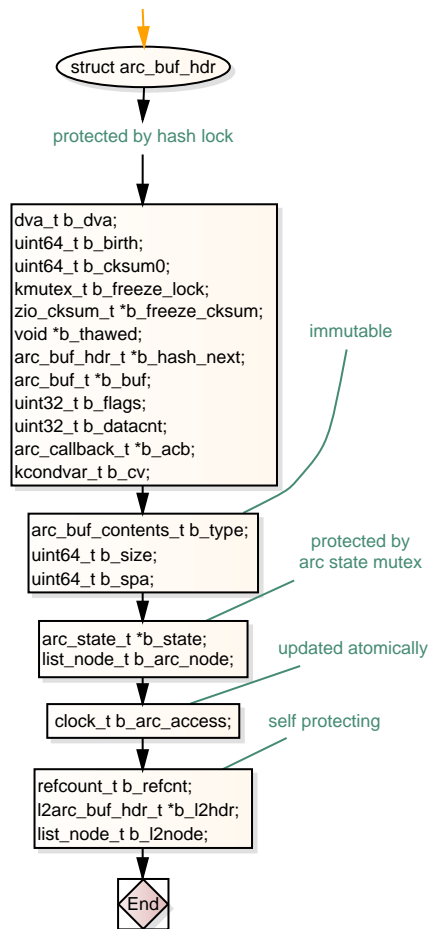
list of evictable buffers

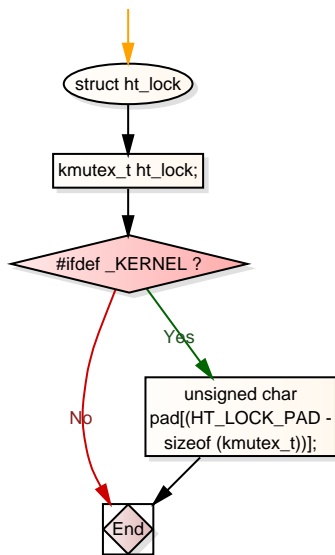


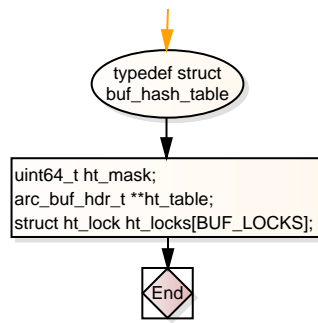


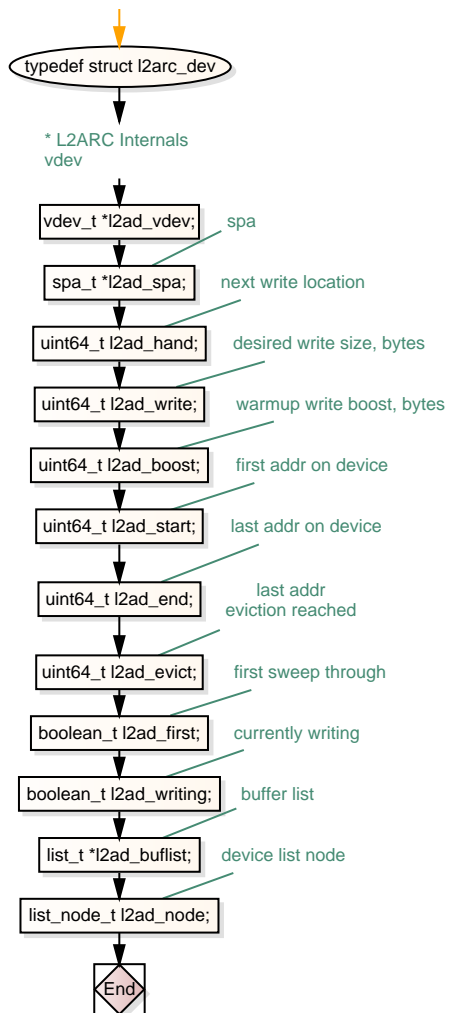


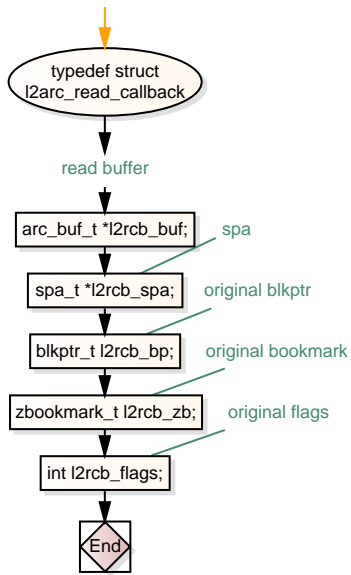


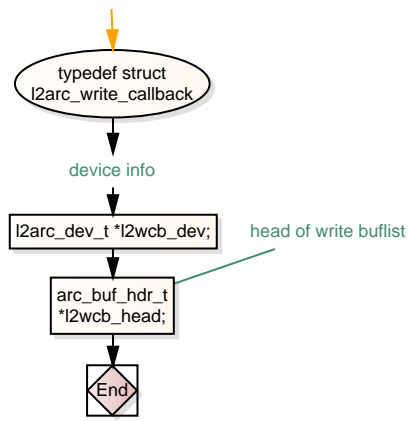


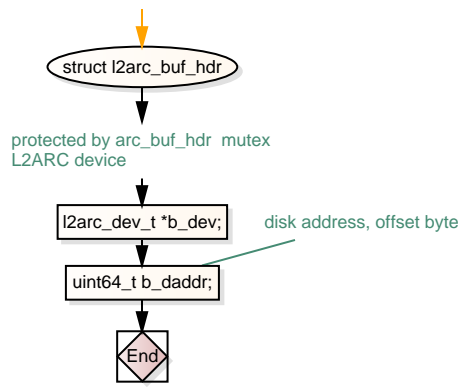


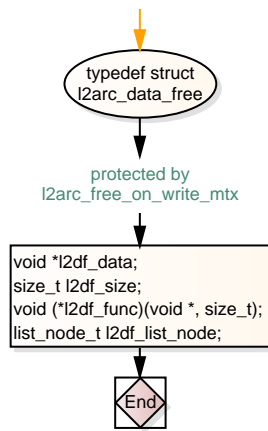


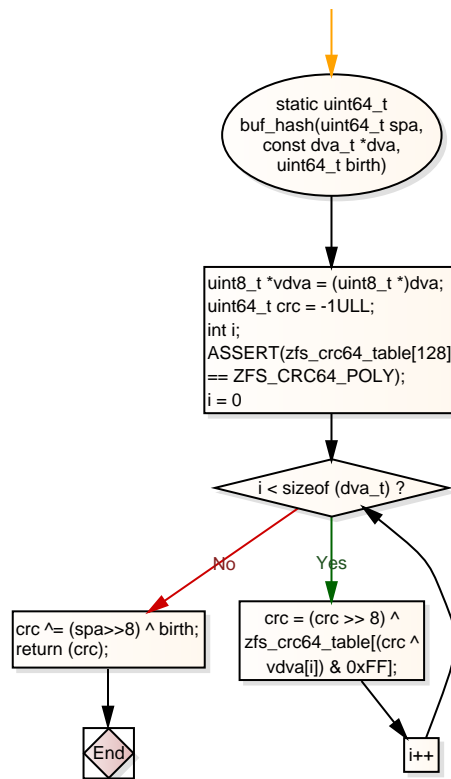


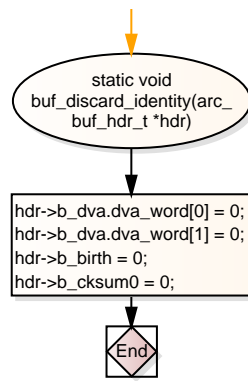


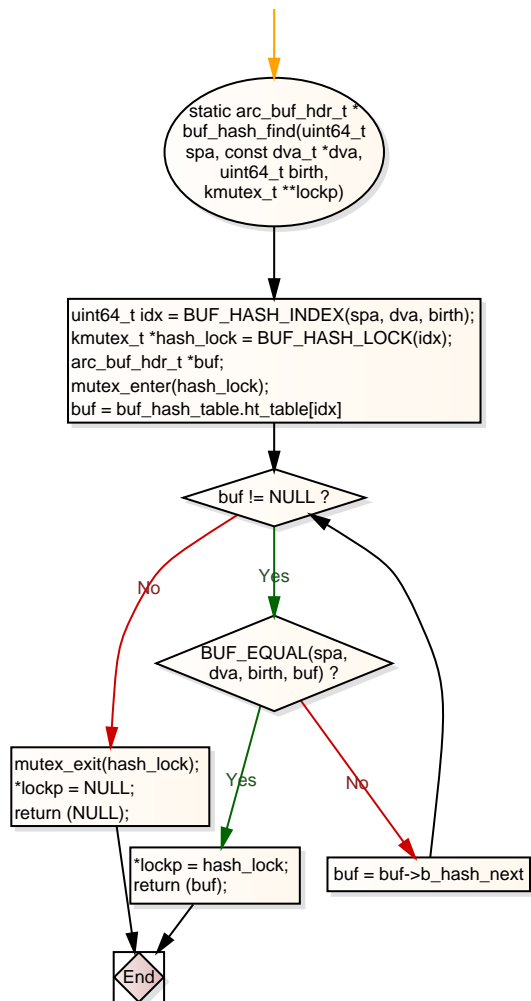


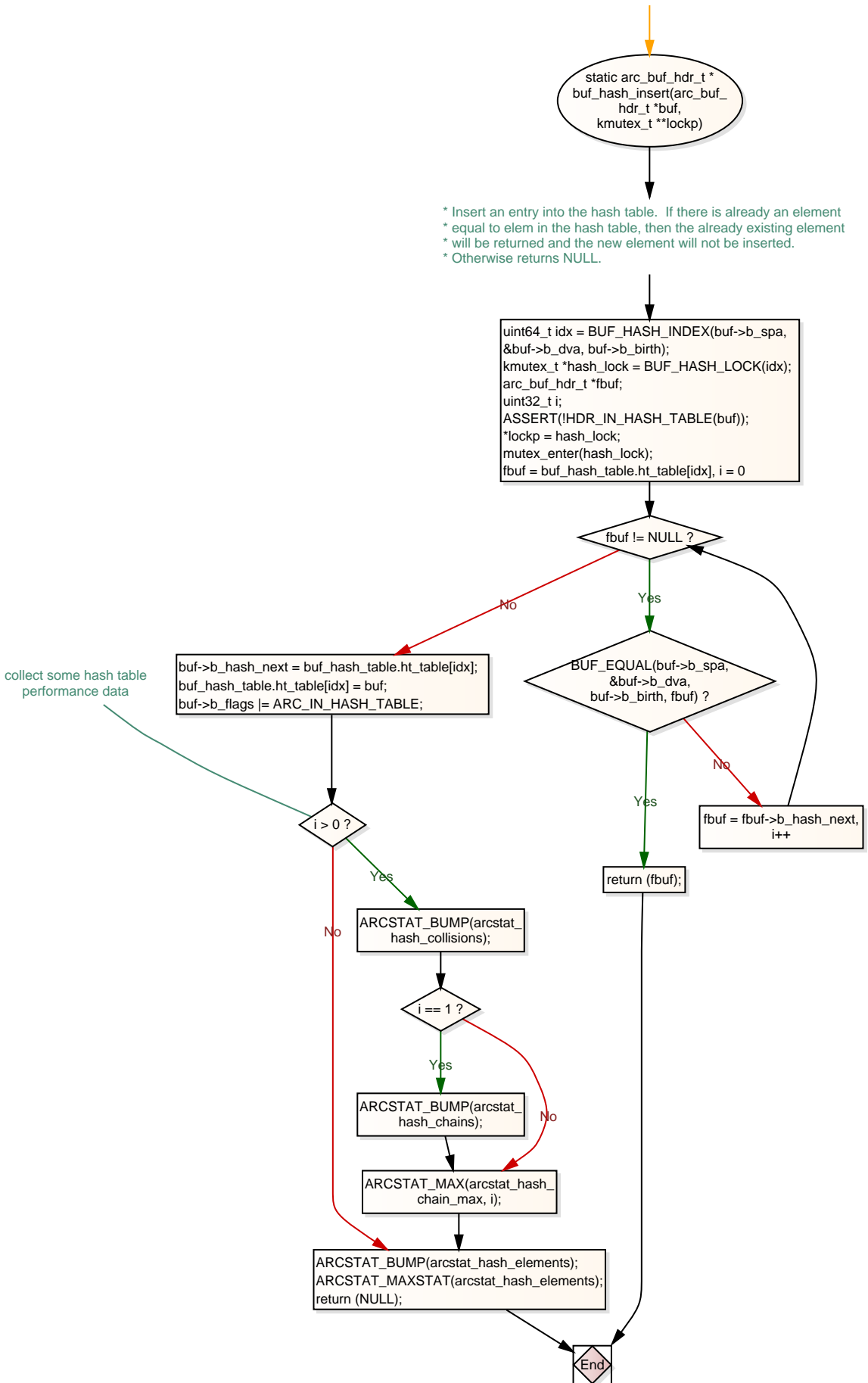


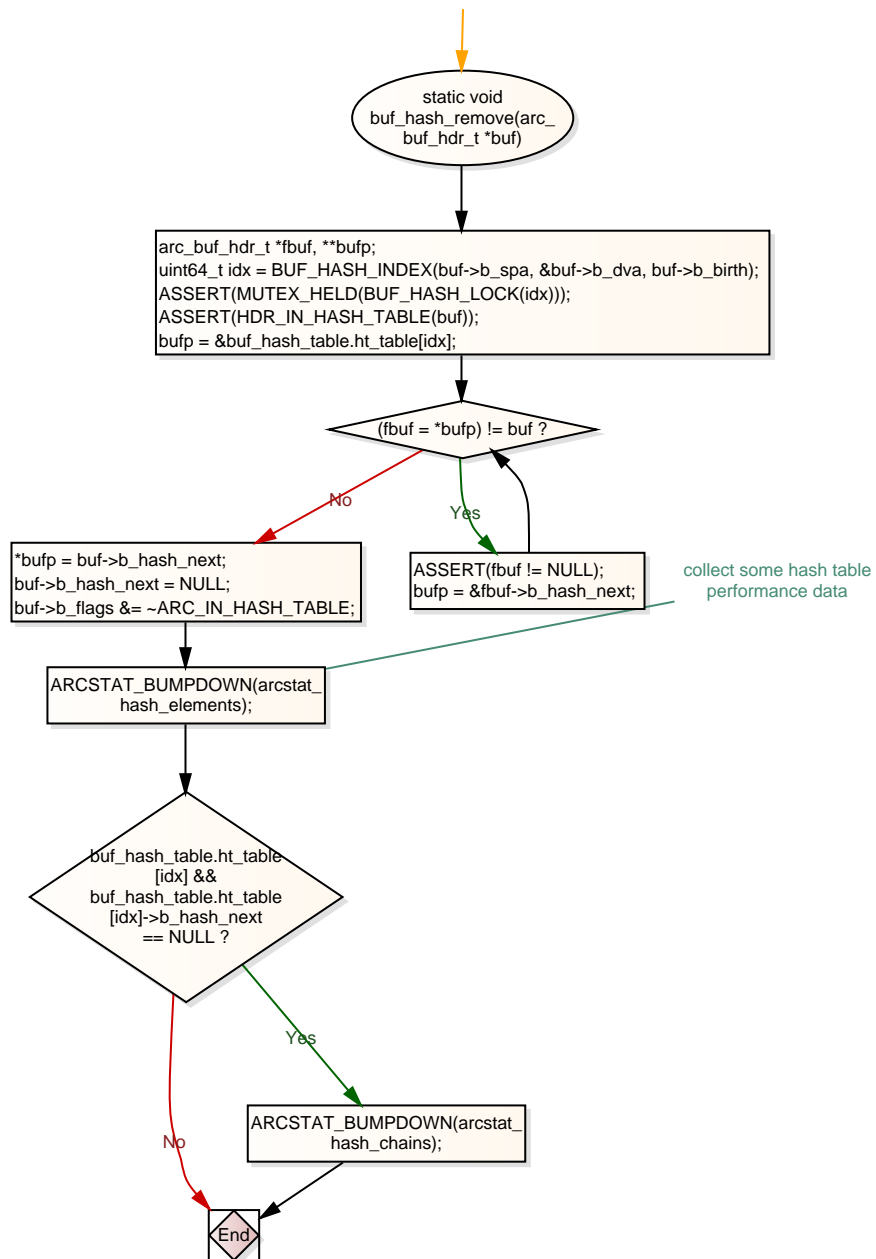


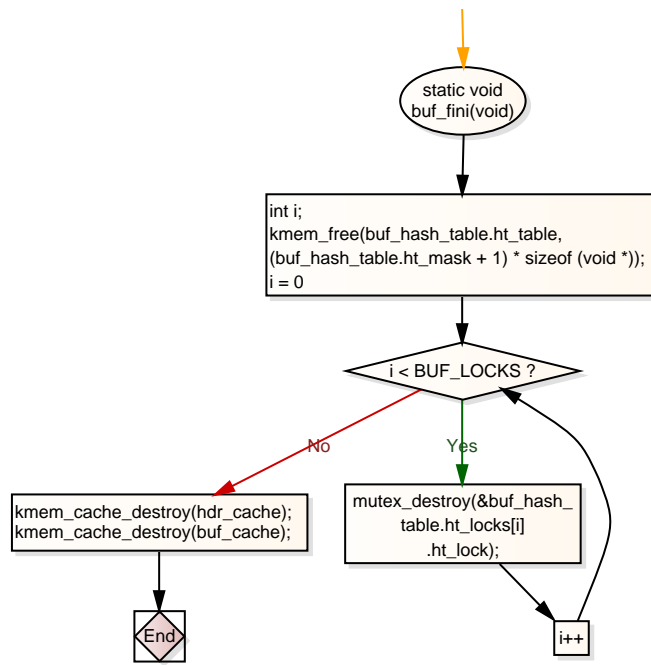


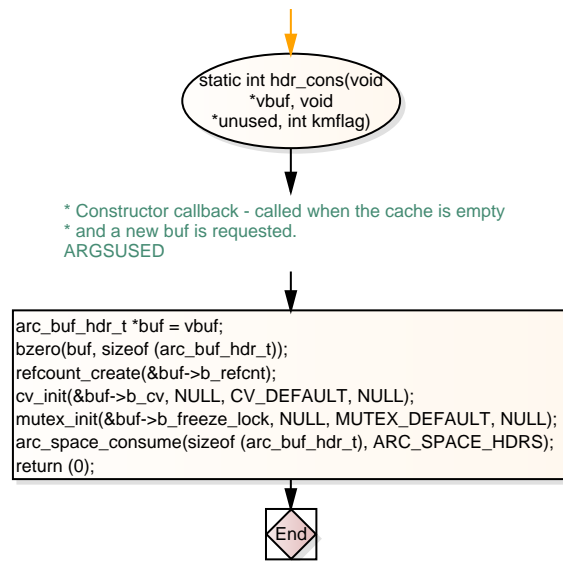


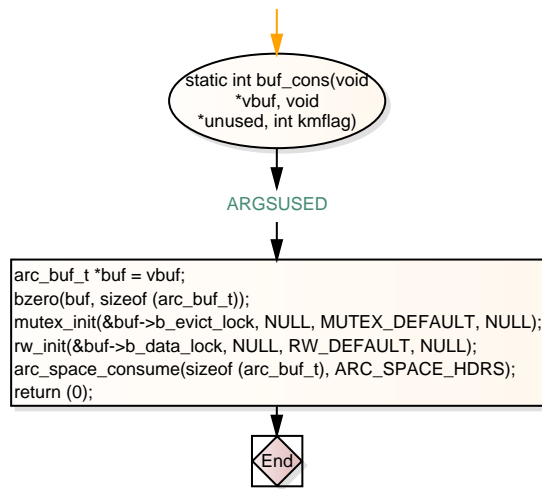


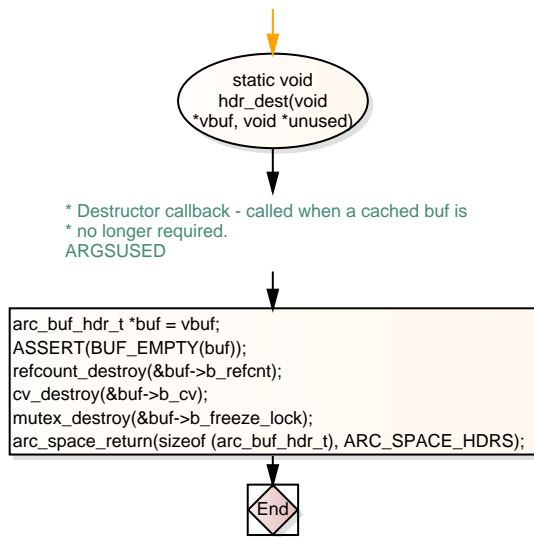


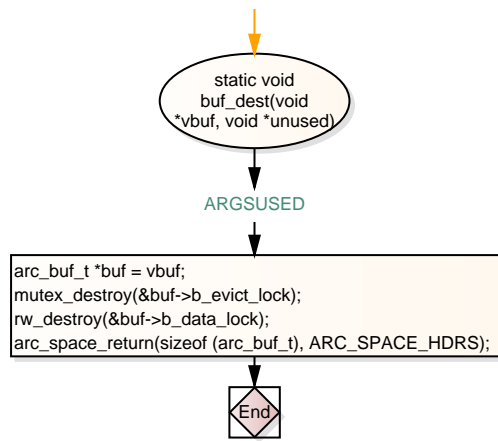




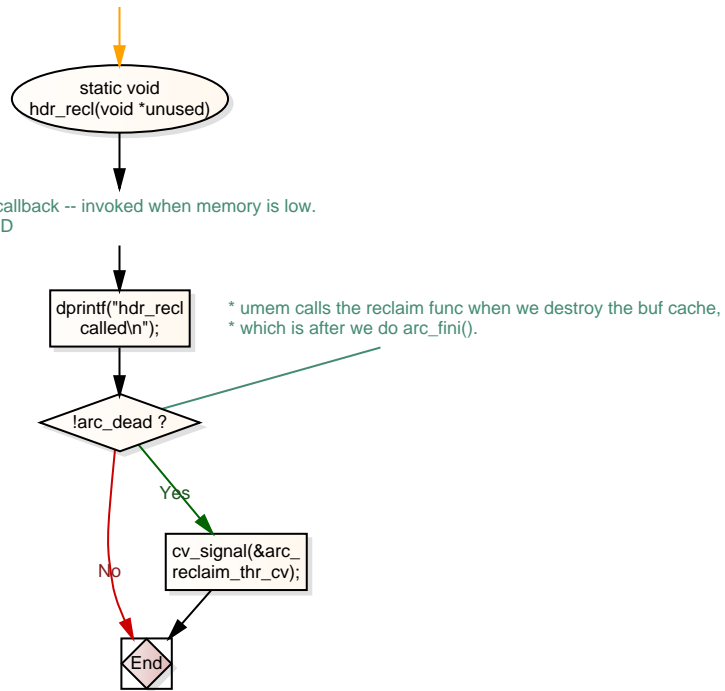


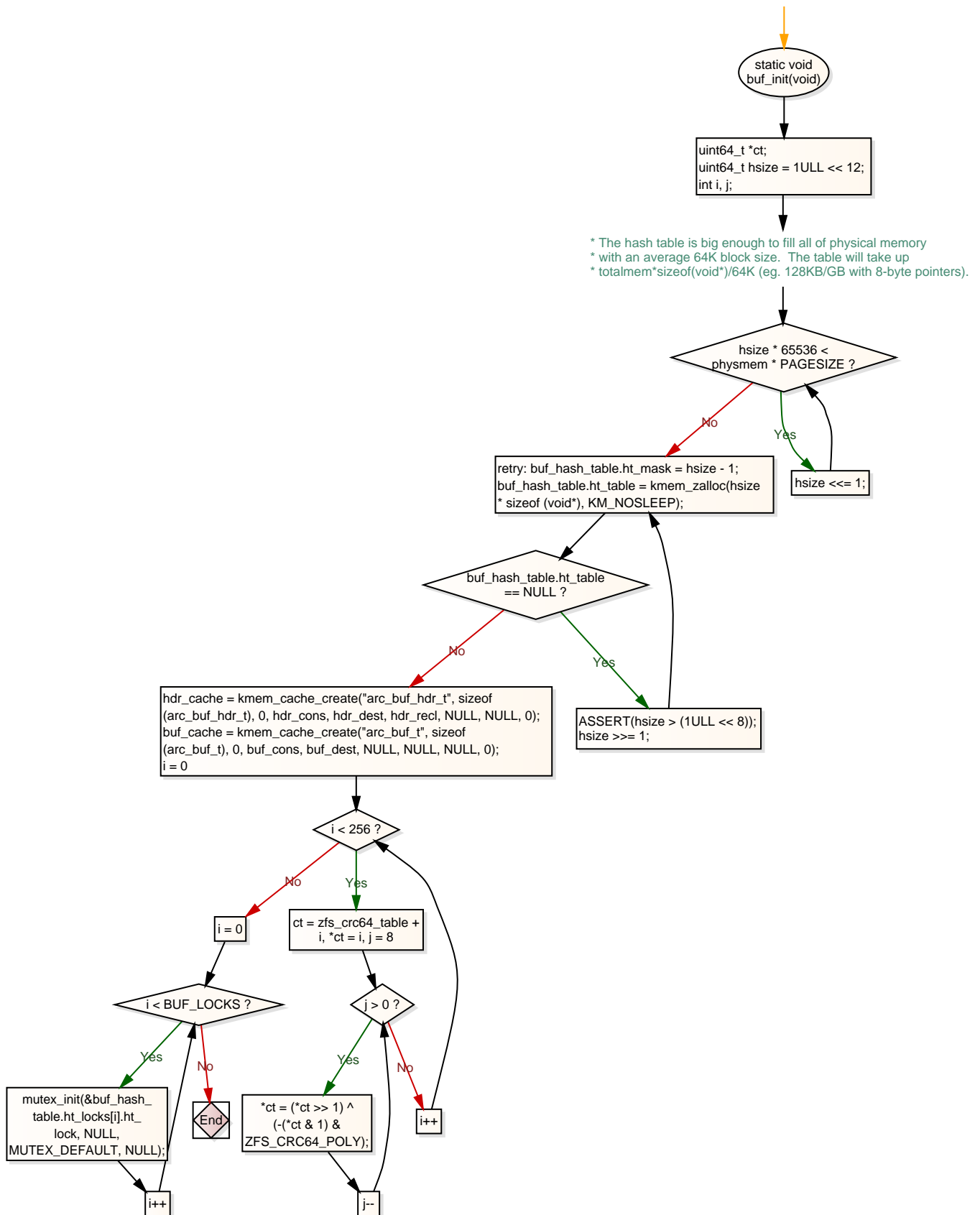


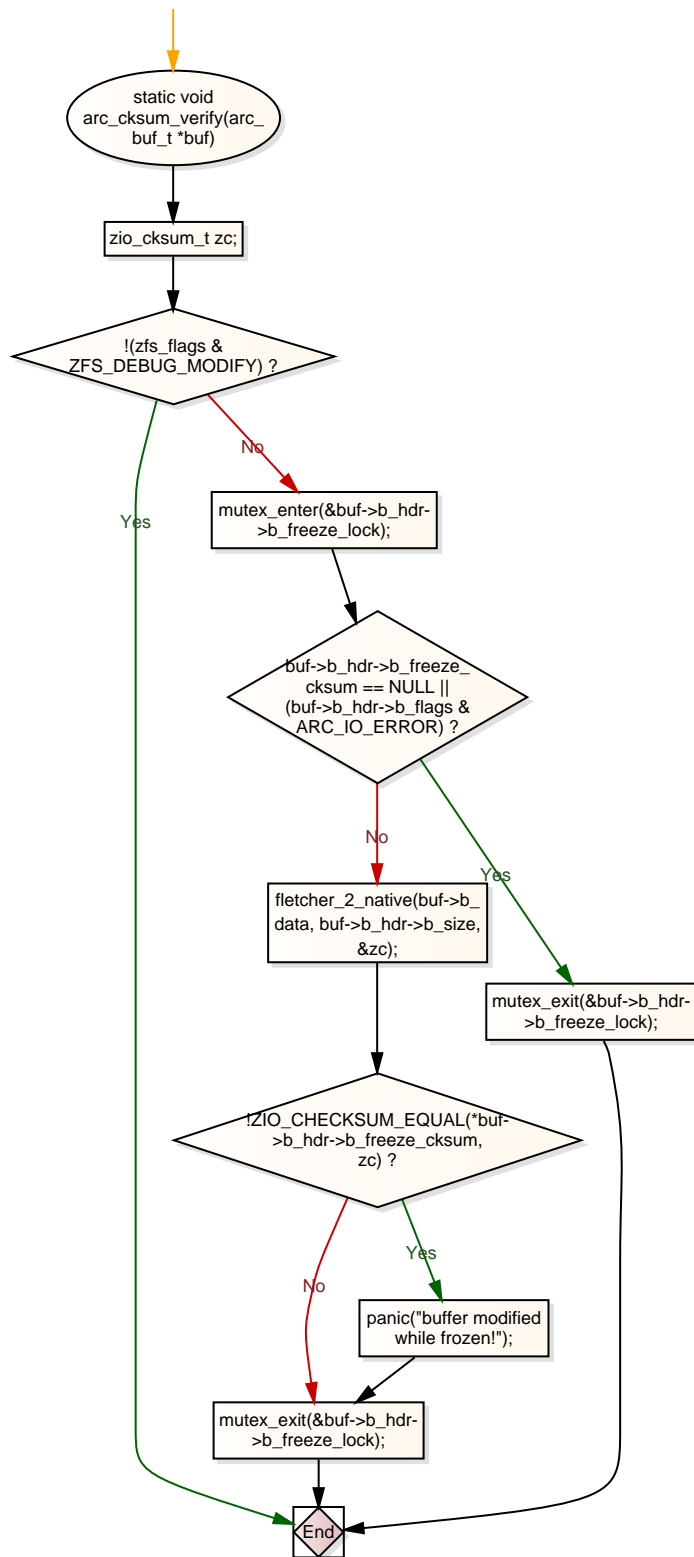


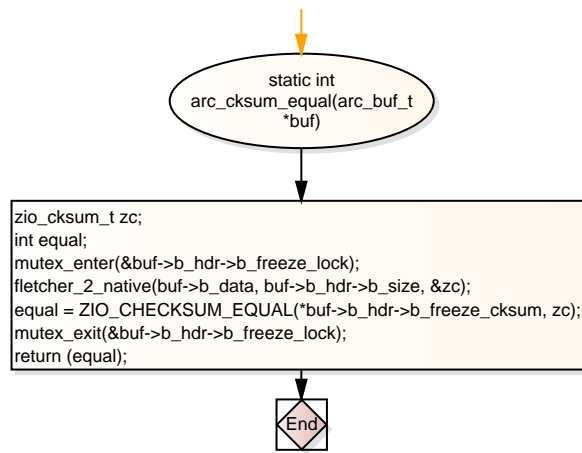


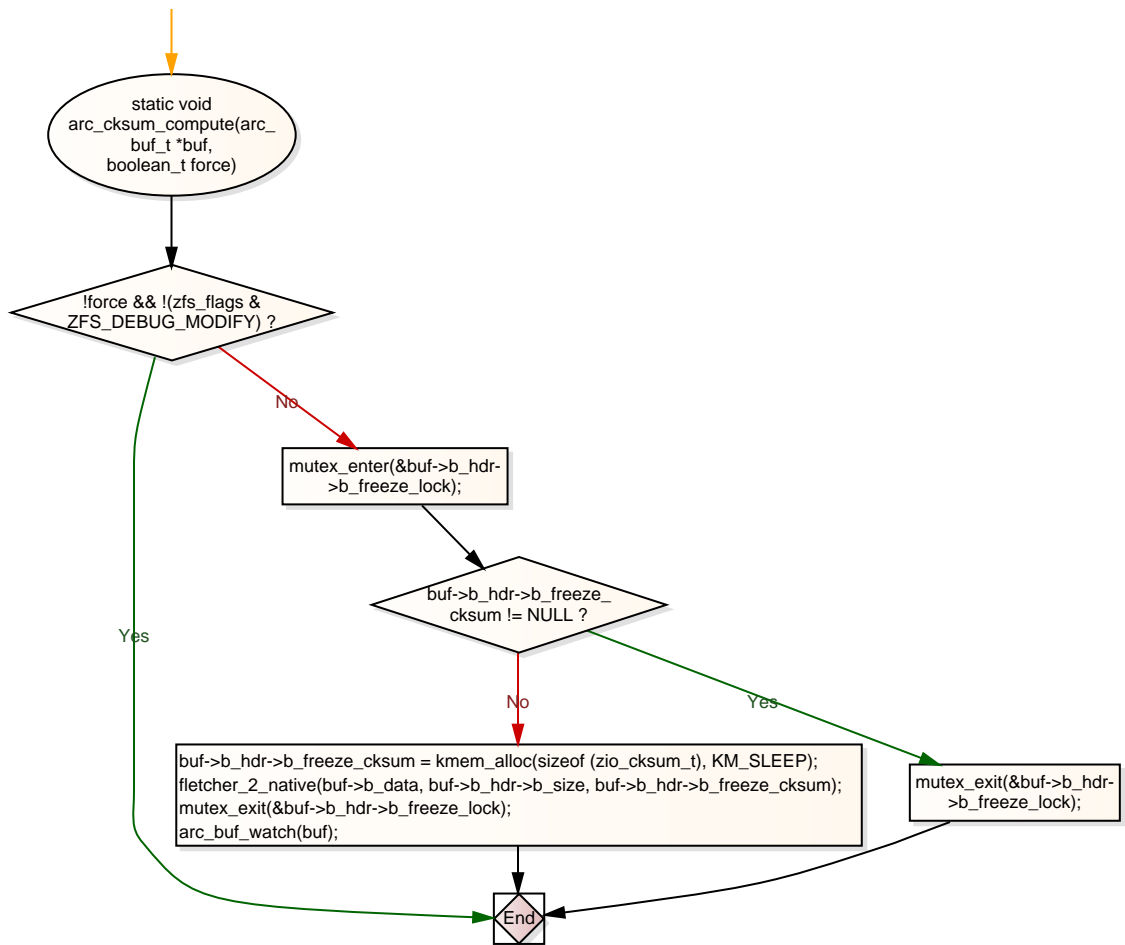
* Reclaim callback -- invoked when memory is low.
ARGSUSED

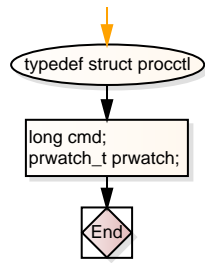


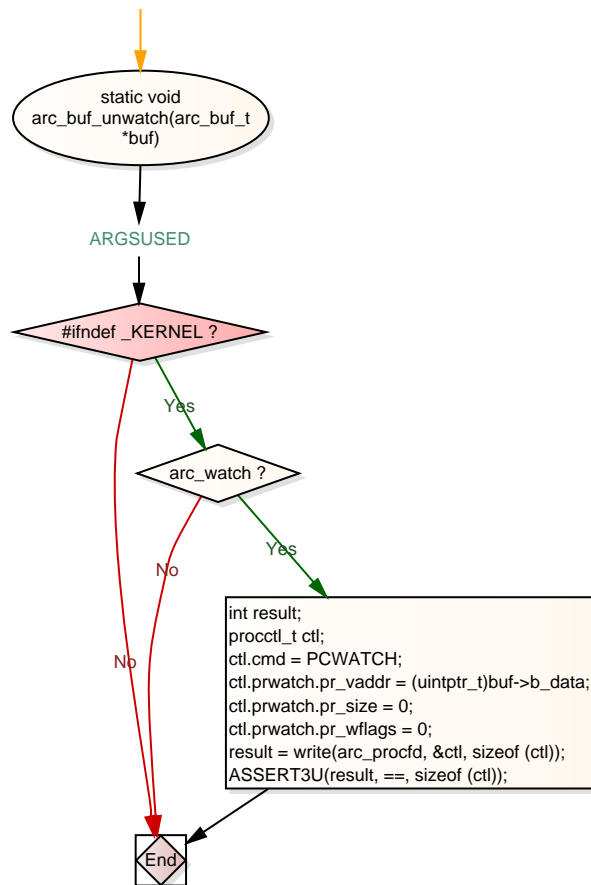


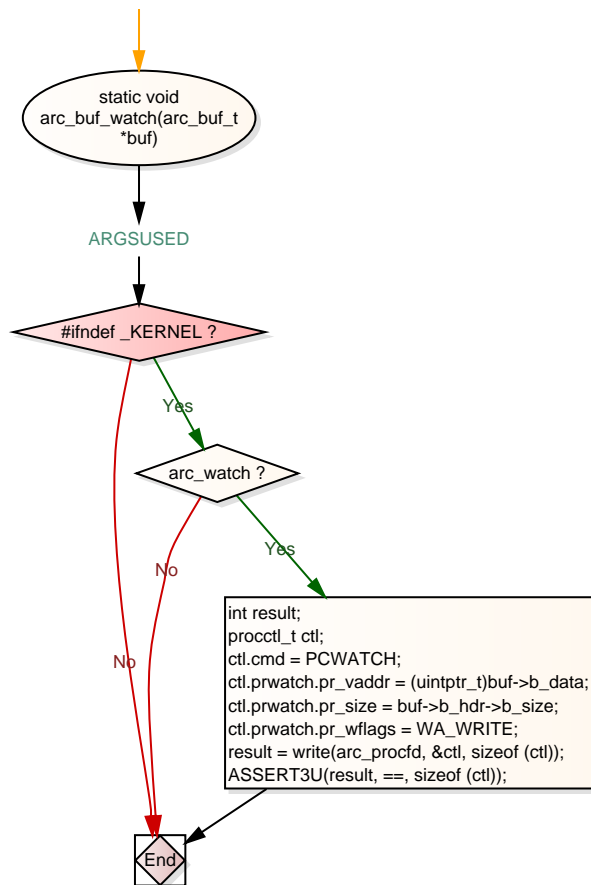


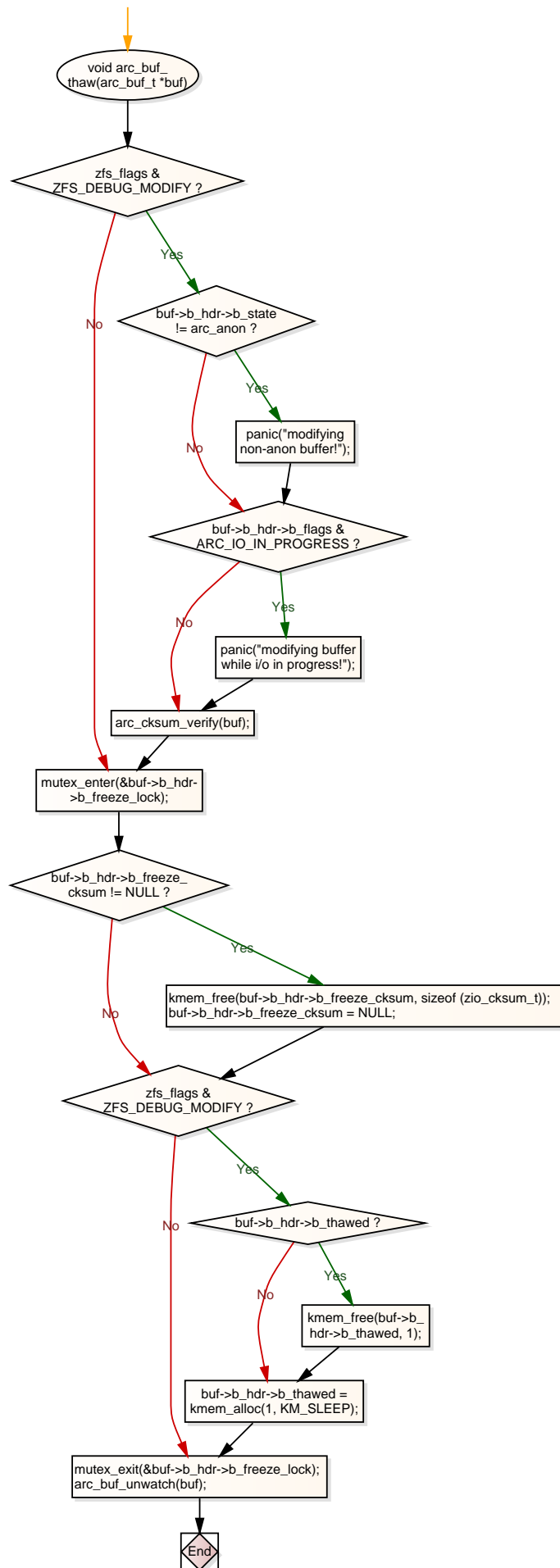


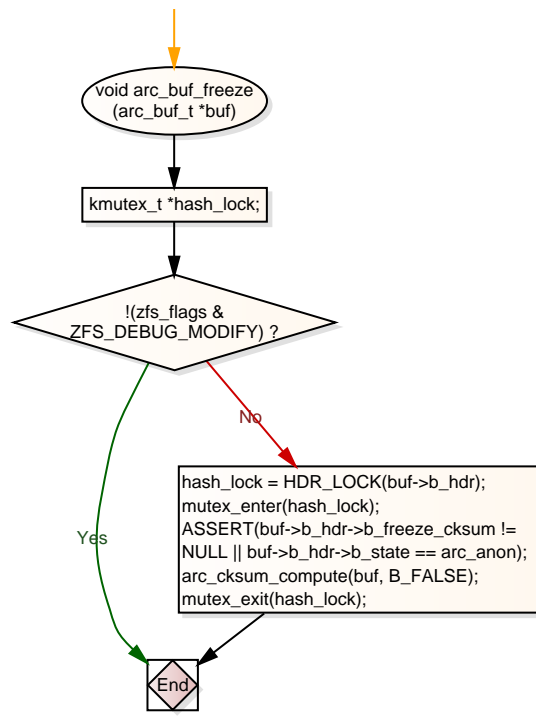


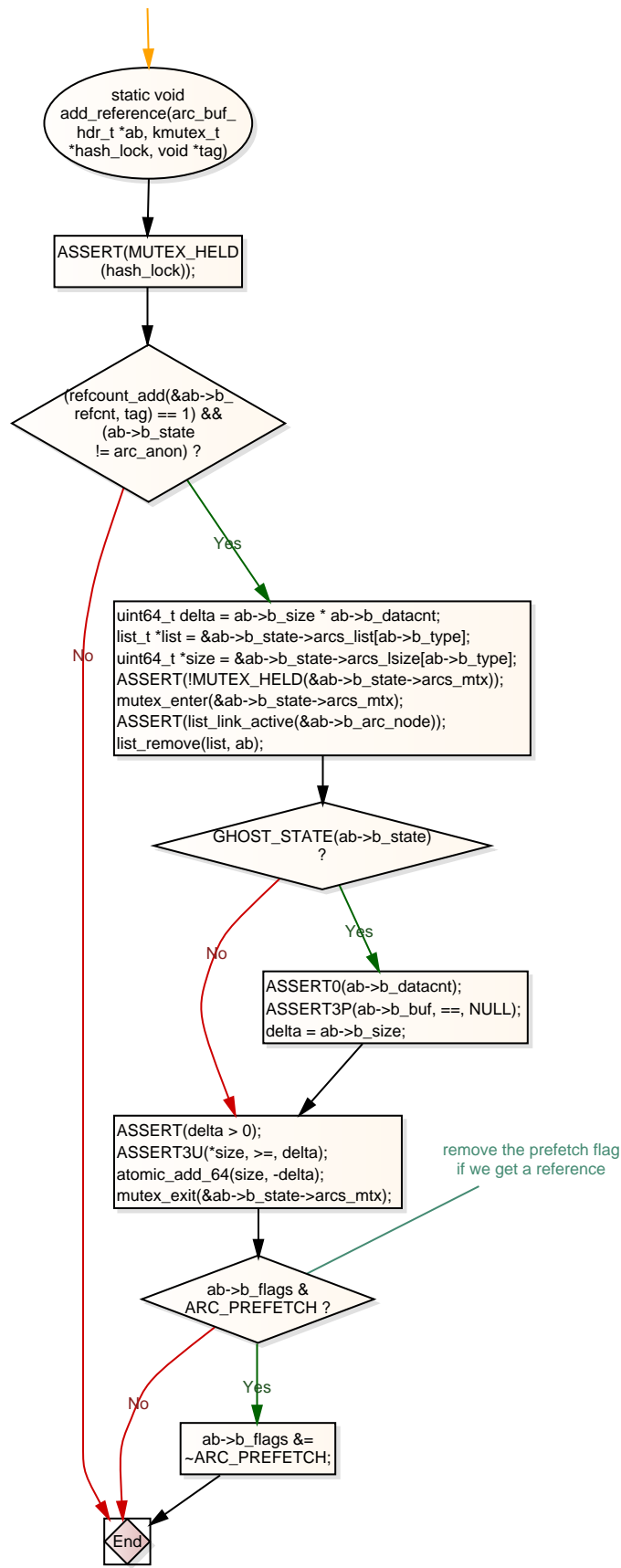


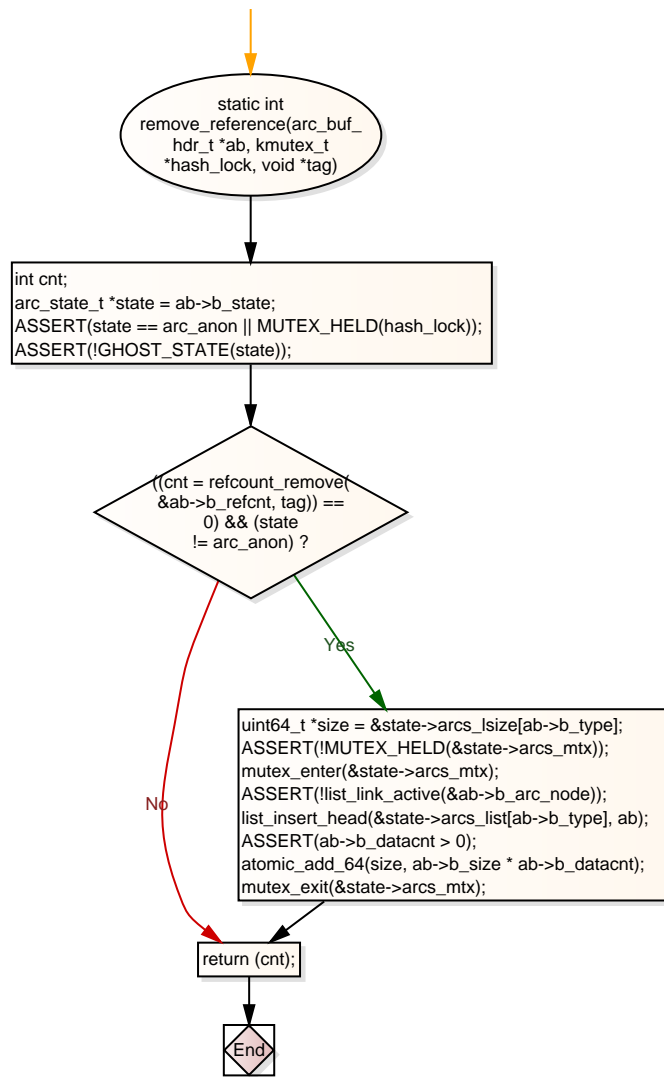


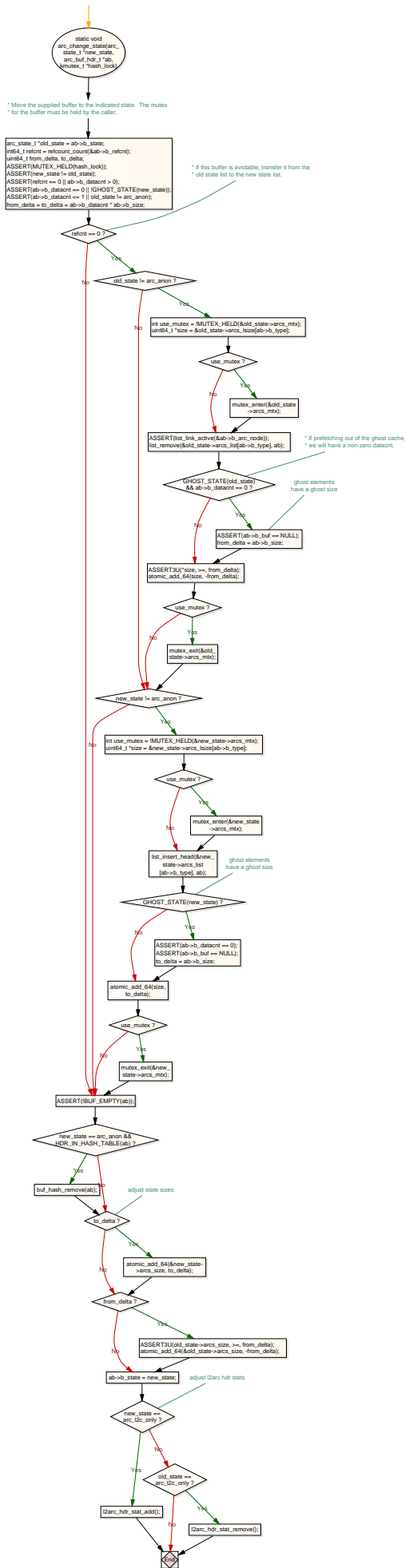


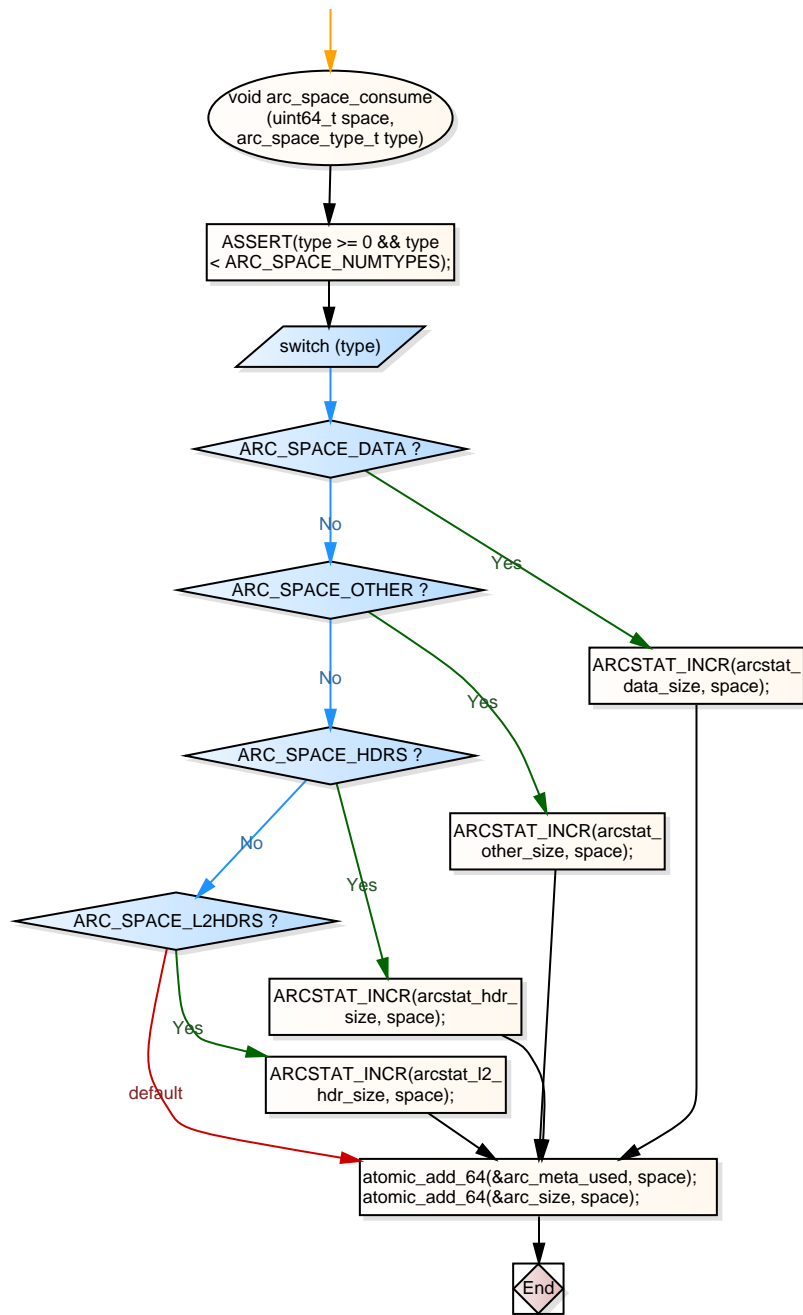


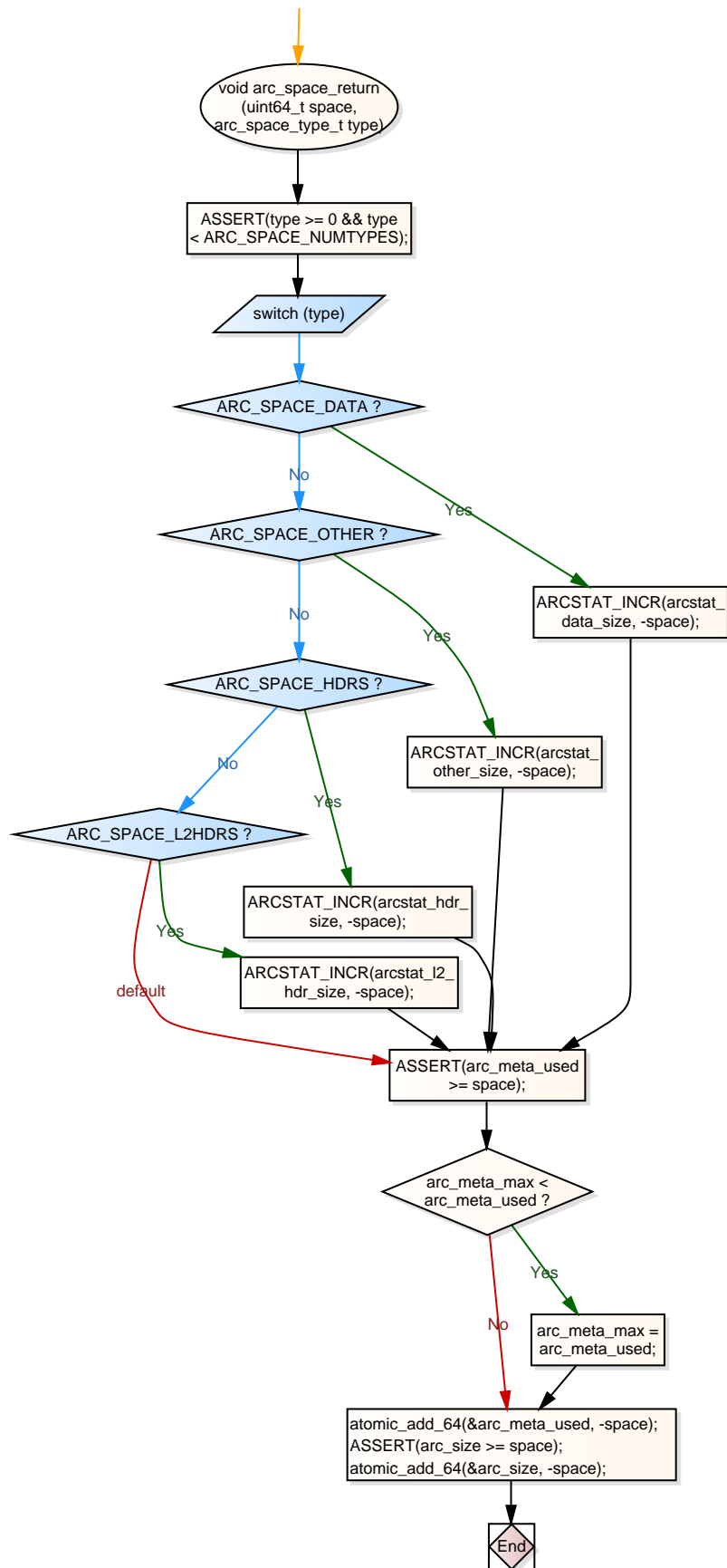


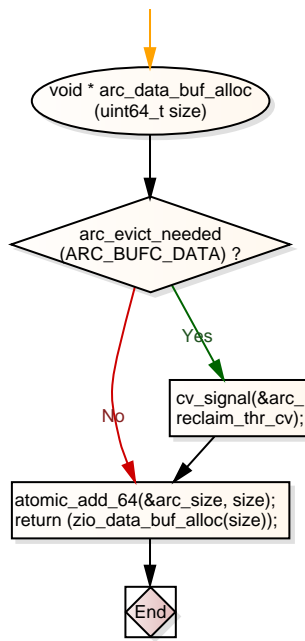


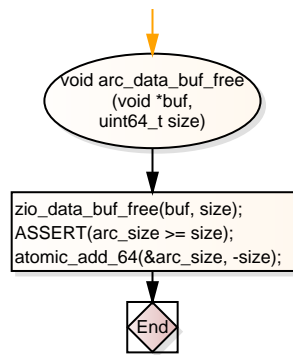


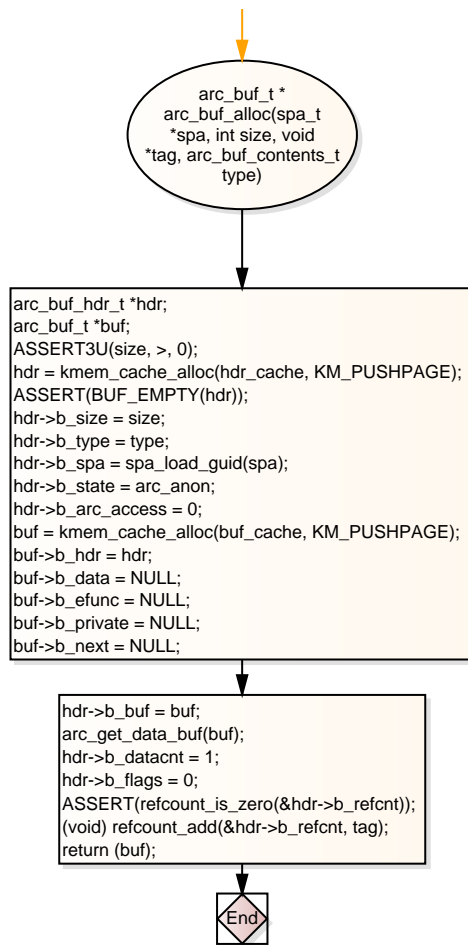










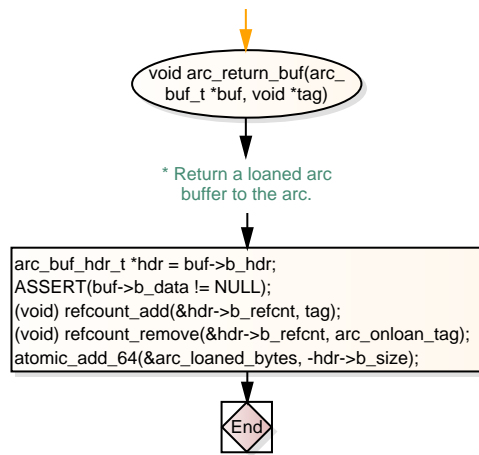


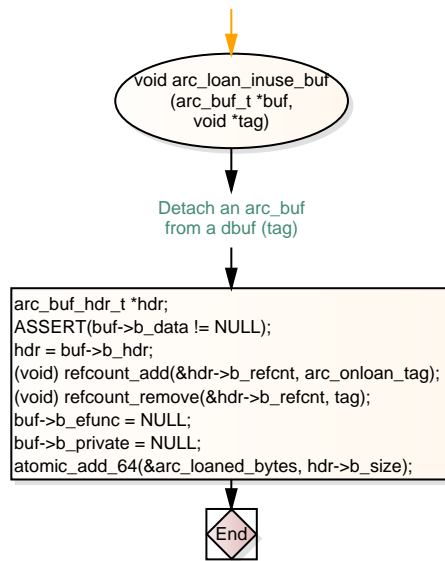
`arc_buf_t *`
`arc_loan_buf(spa_t`
`*spa, int size)`

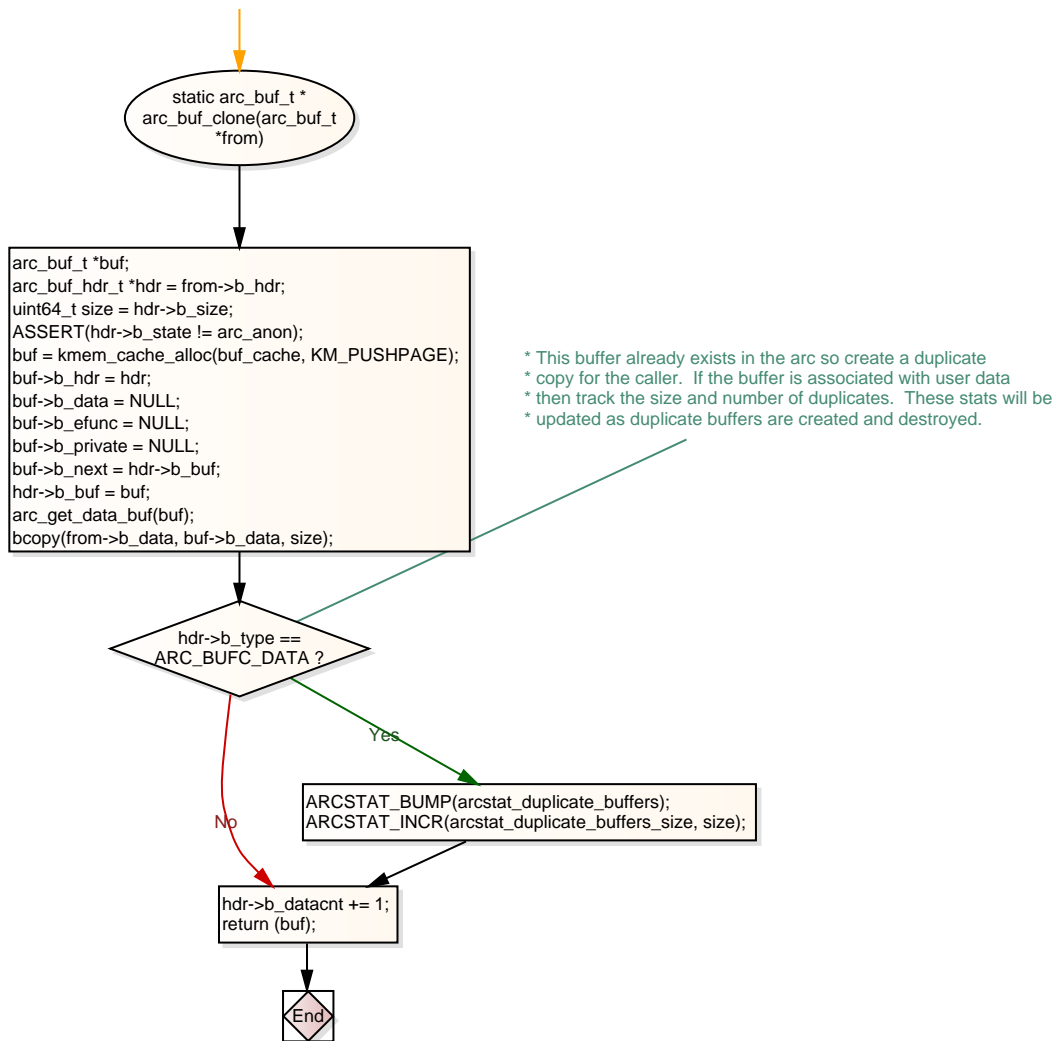
* Loan out an anonymous arc buffer. Loaned buffers are not counted as in
* flight data by `arc_tempreserve_space()` until they are "returned". Loaned
* buffers must be returned to the arc before they can be used by the DMU or
* freed.

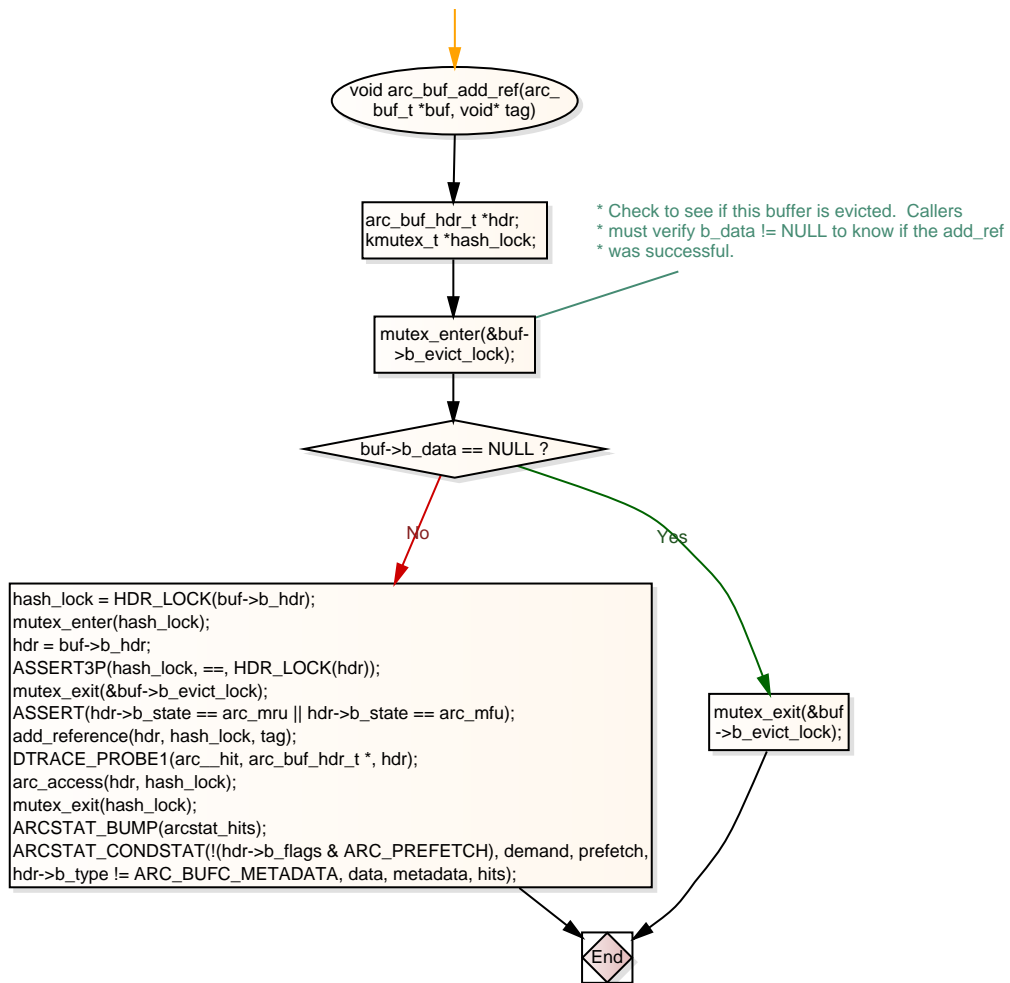
```
arc_buf_t *buf;  
buf = arc_buf_alloc(spa, size, arc_onloan_tag, ARC_BUFC_DATA);  
atomic_add_64(&arc_loaned_bytes, size);  
return (buf);
```

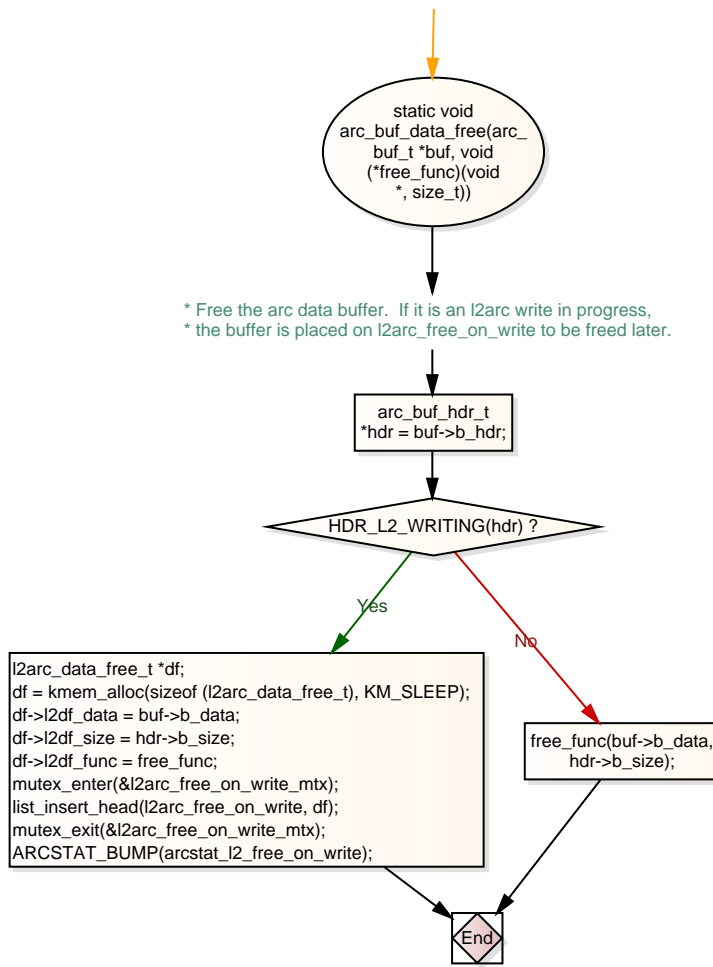
End

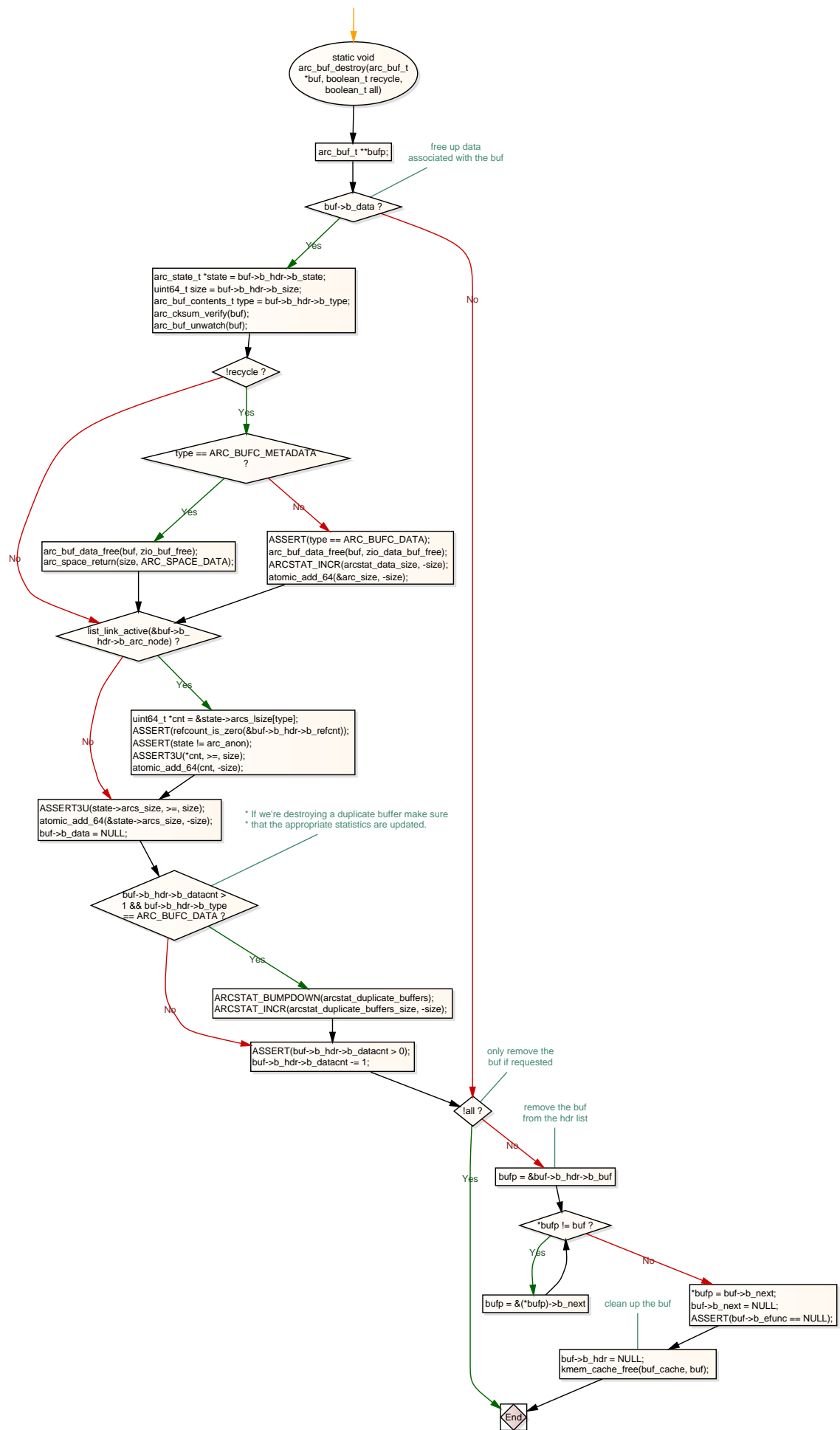


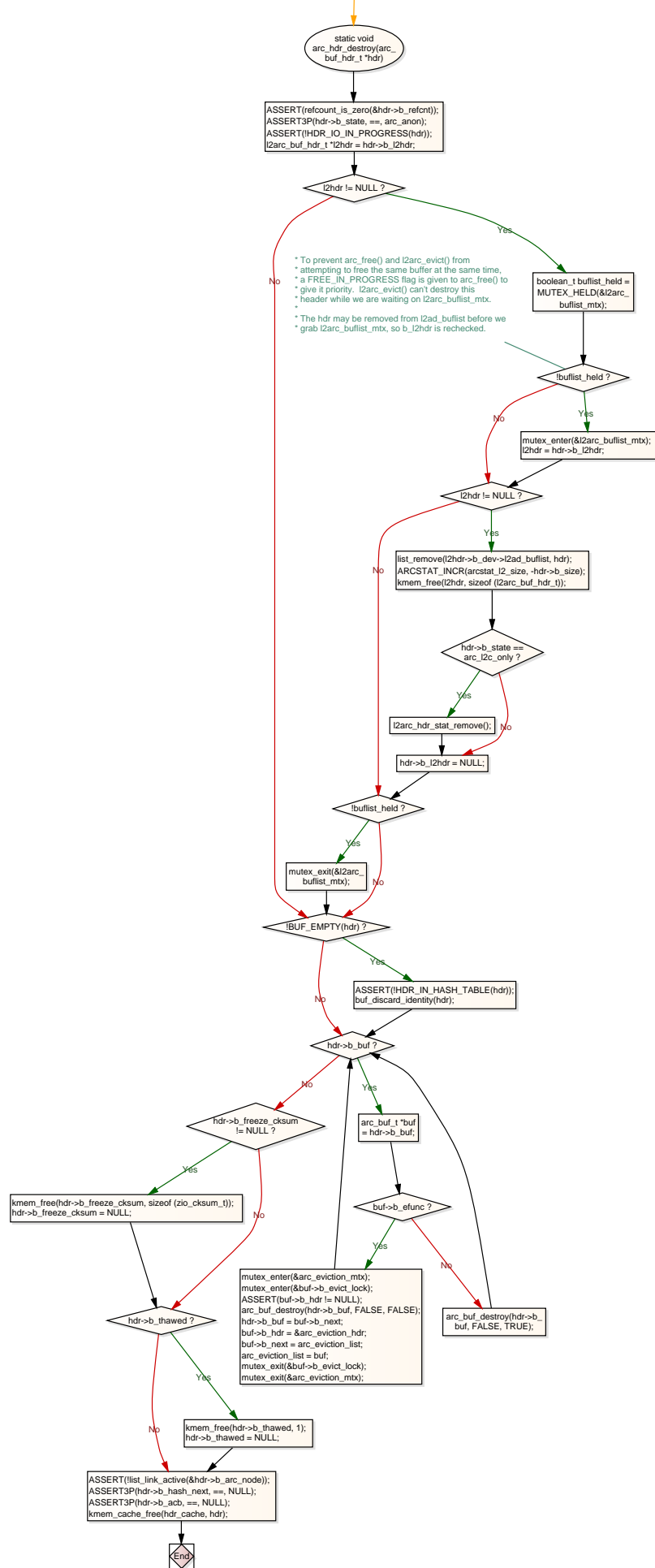


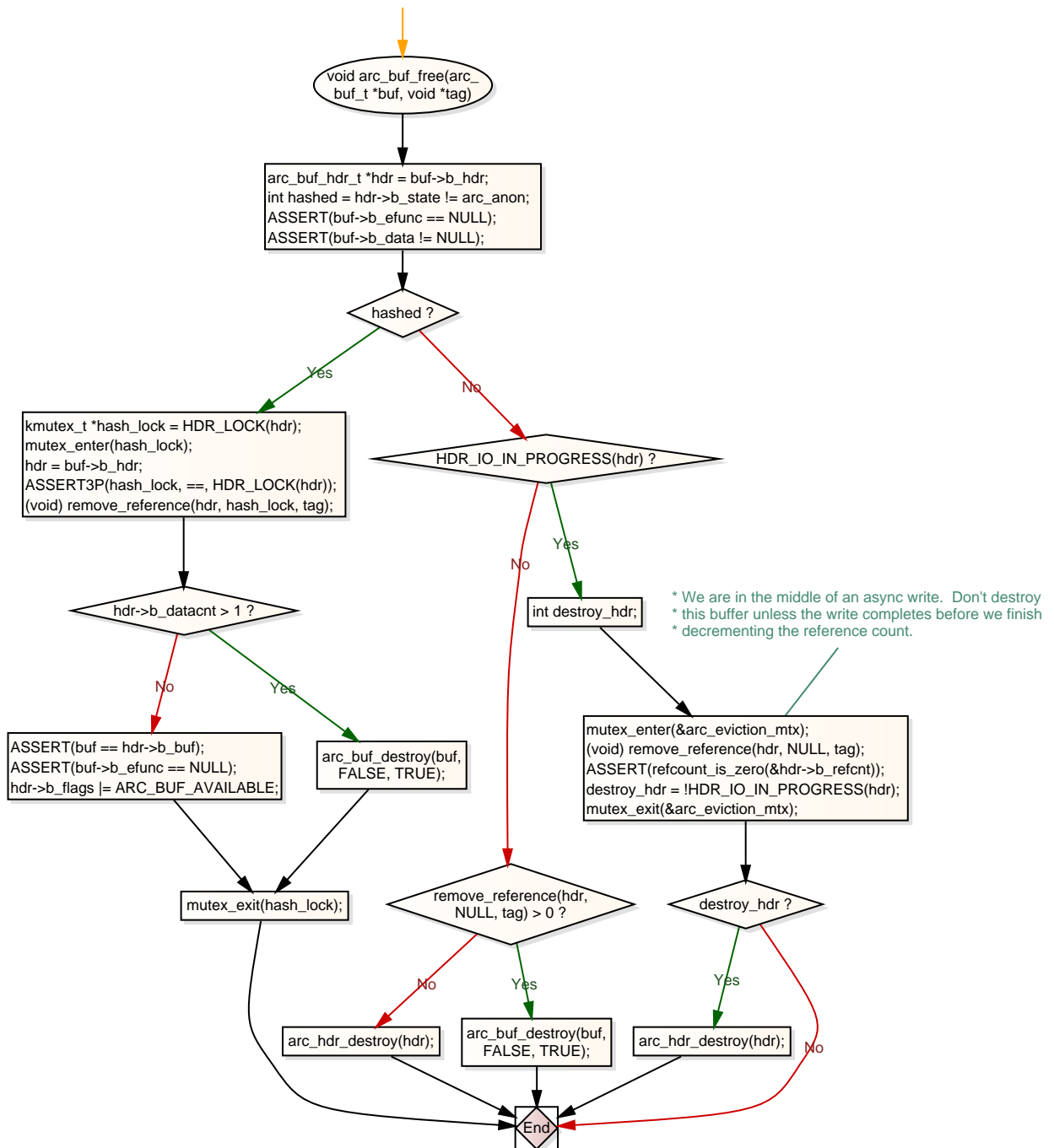


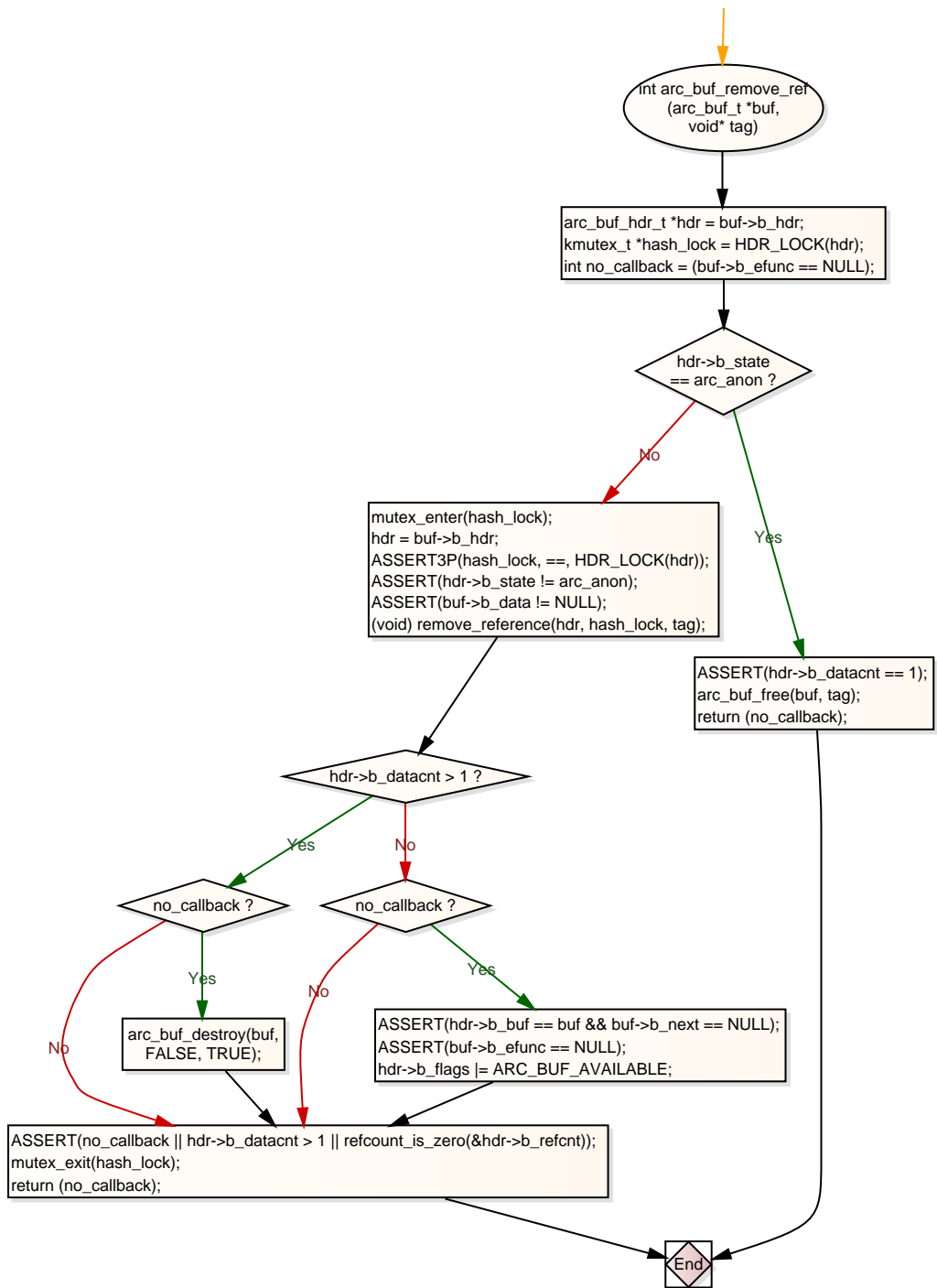


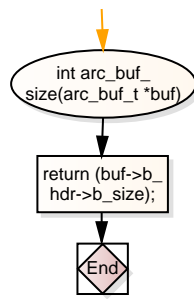


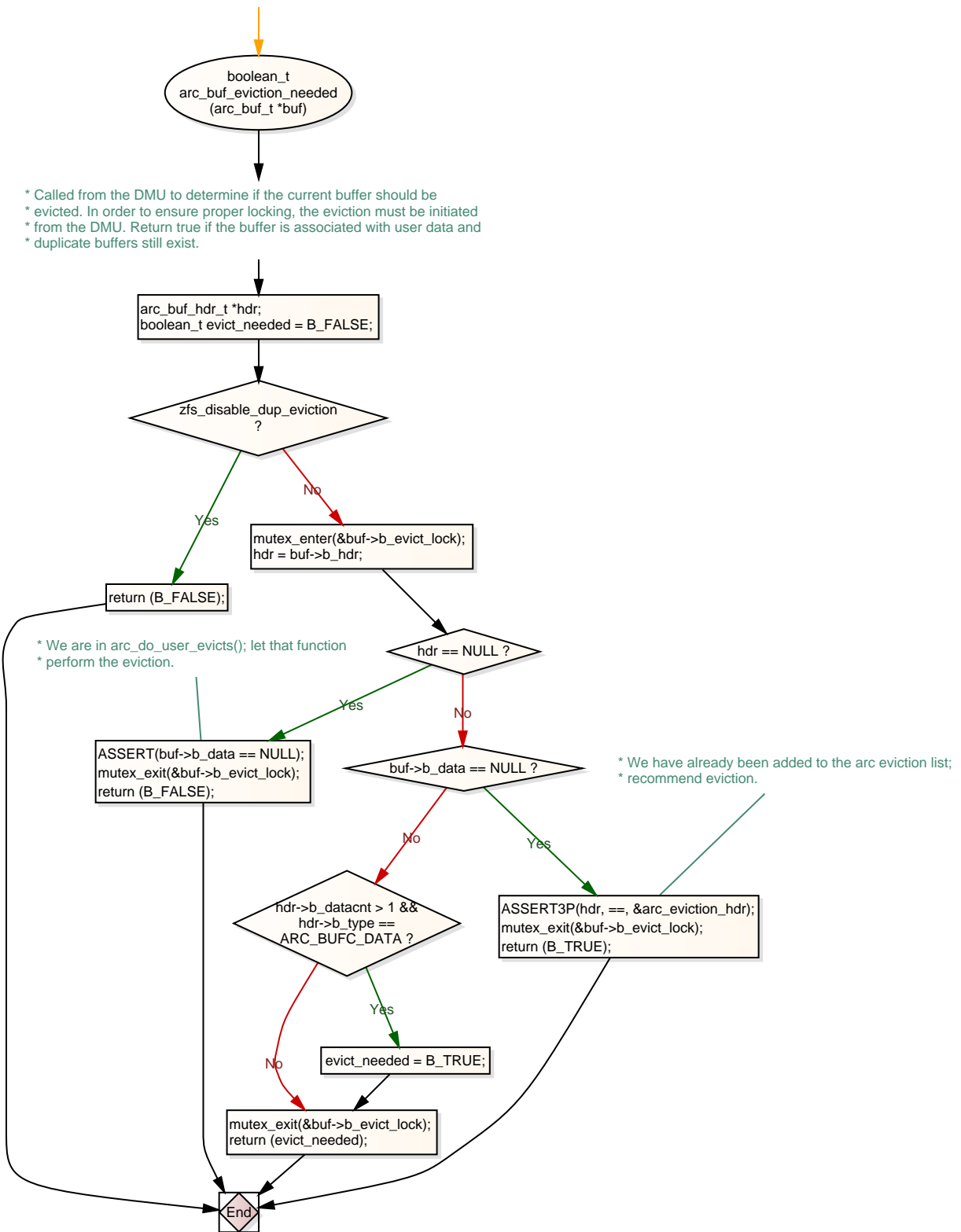


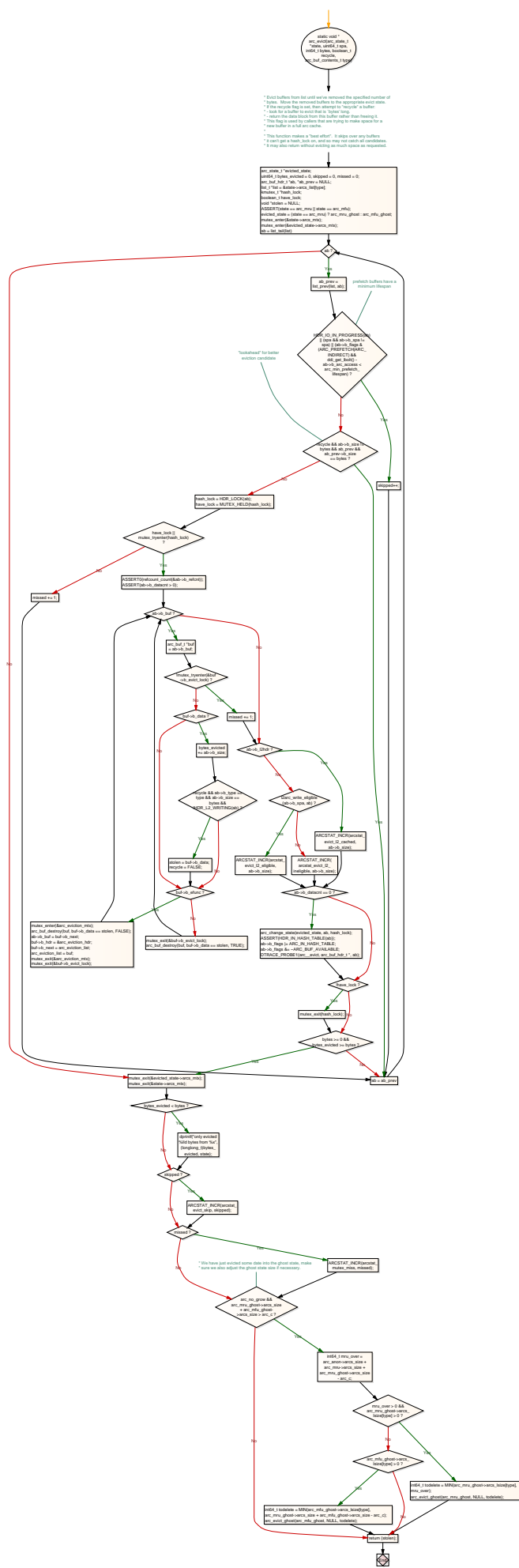


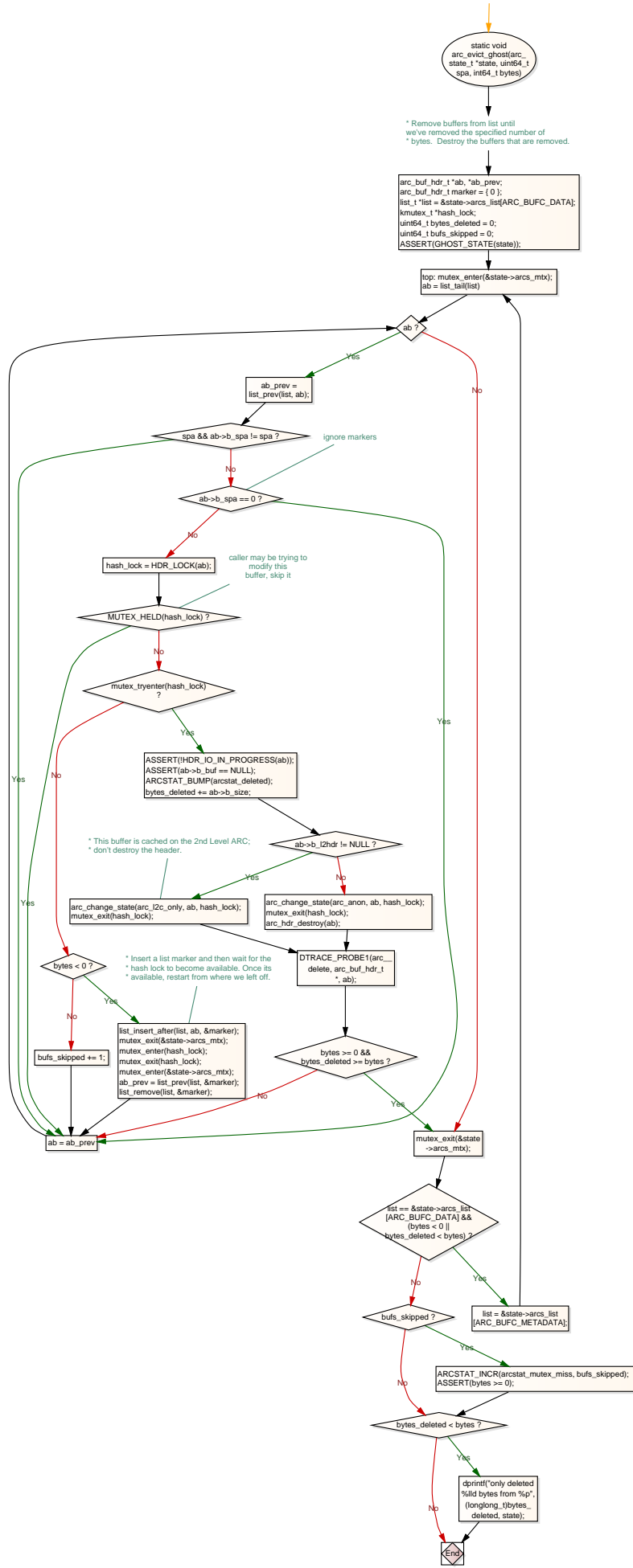


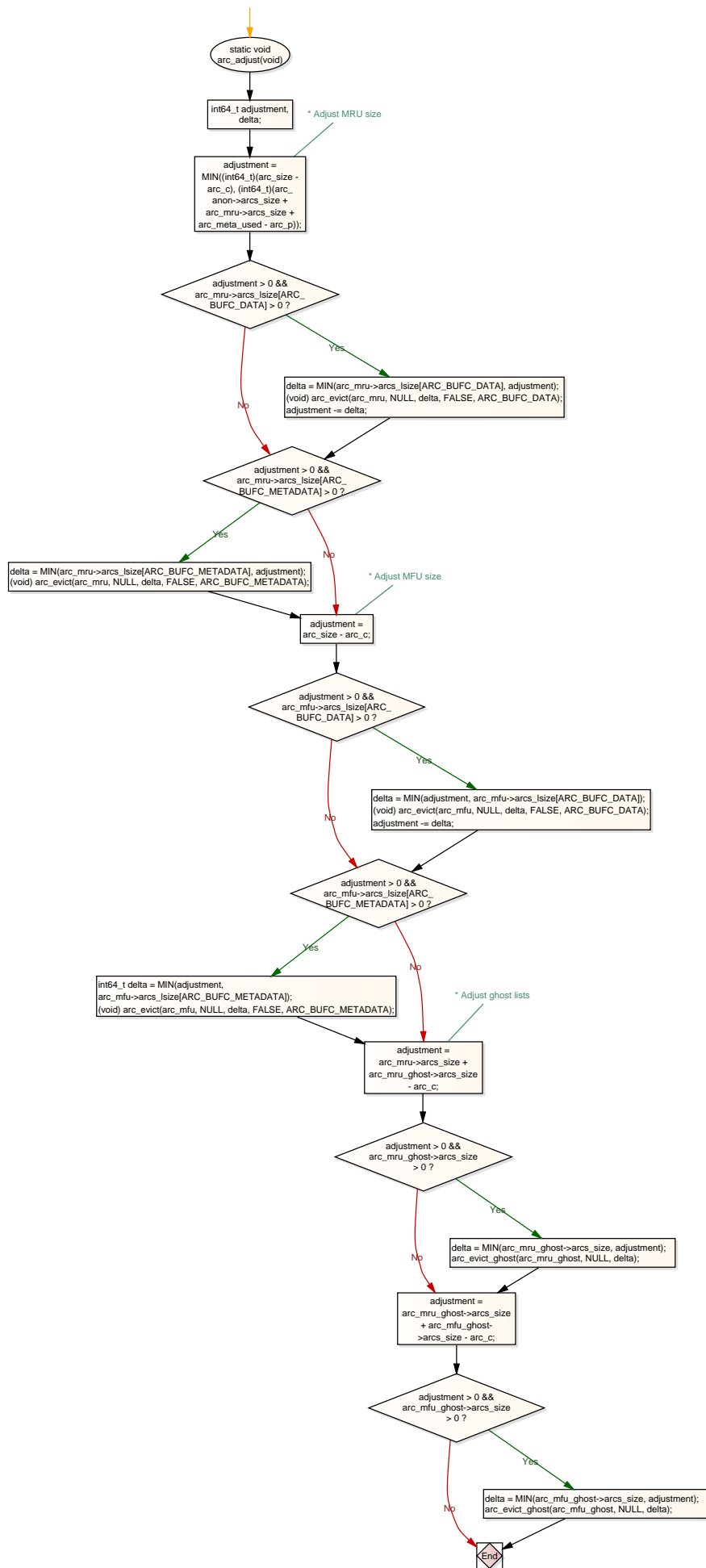


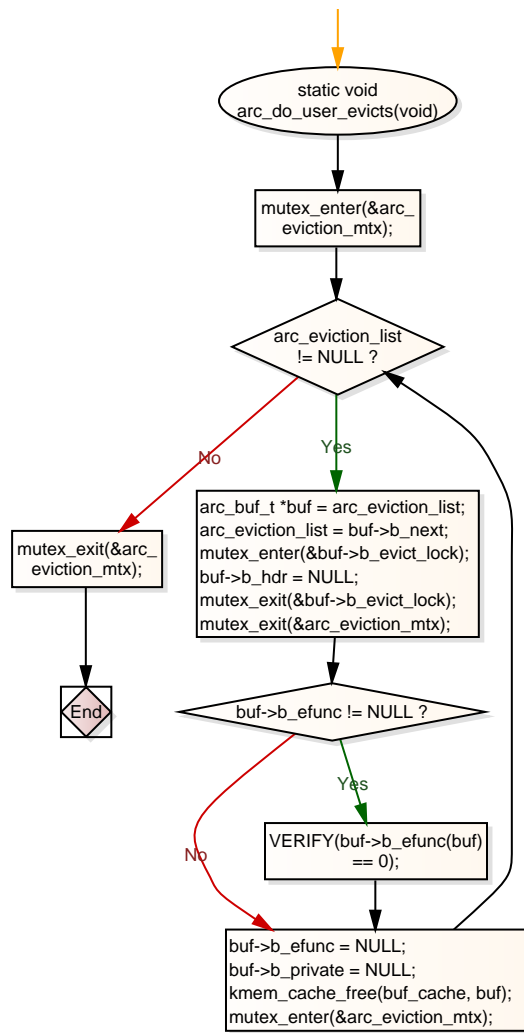


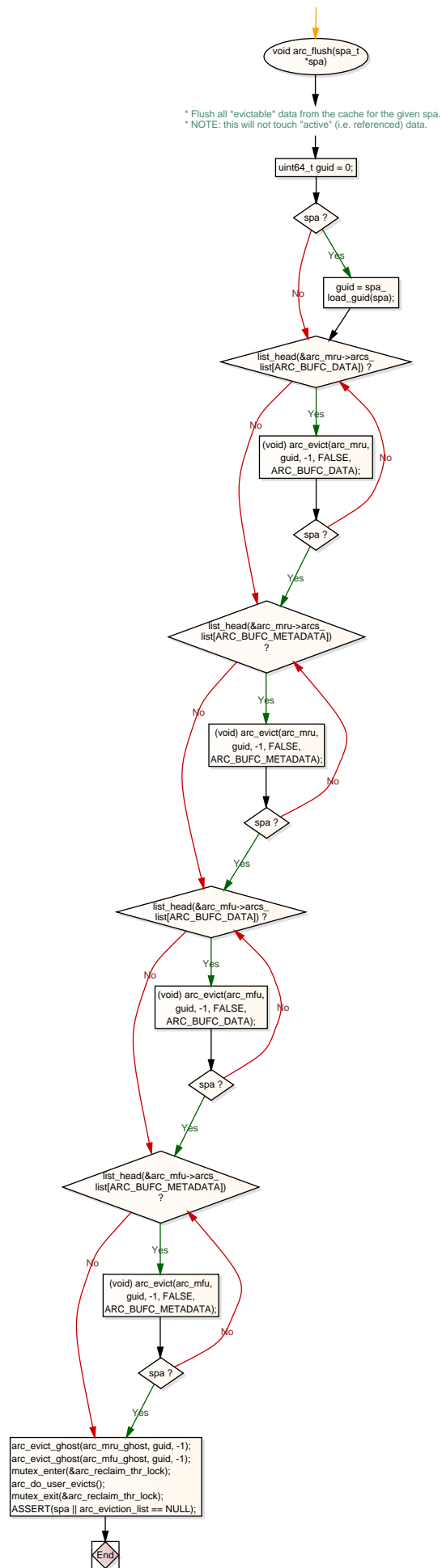


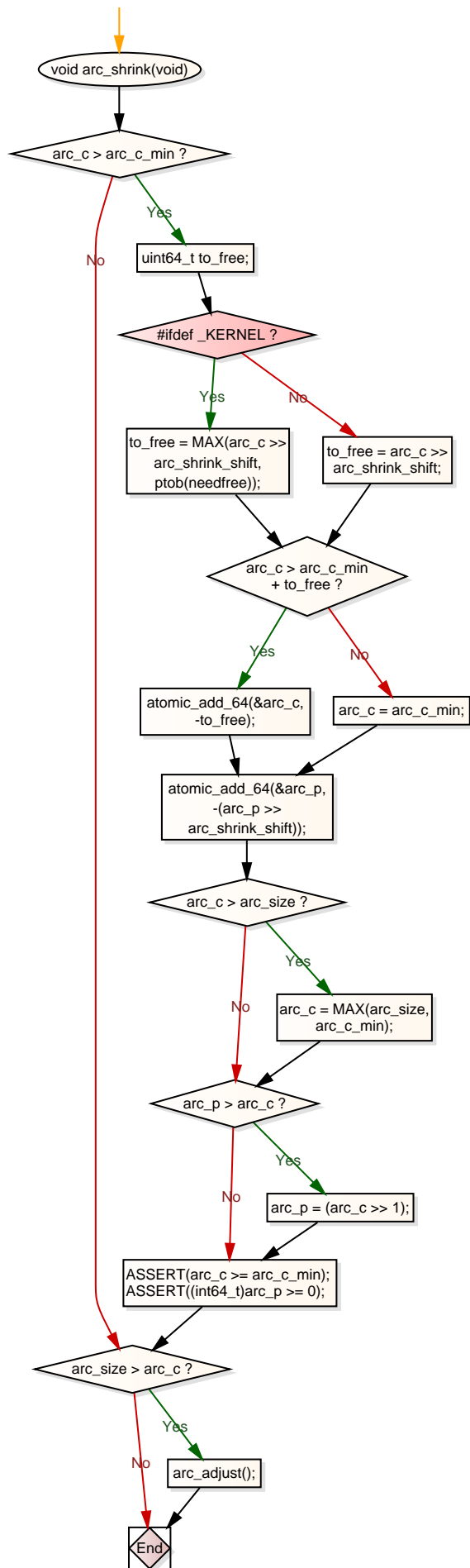












static int
arc_reclaim_needed(void)

* Determine if the system is under memory pressure and is asking
* to reclaim memory. A return value of 1 indicates that the system
* is under memory pressure and that the arc should adjust accordingly.

uint64_t extra;

#ifdef _KERNEL ?

needfree ?

Yes

No

return (1);

Yes

return (1);

* check that we're out of range of the pageout scanner. It starts to
* schedule paging if freemem is less than lotsfree and needfree.
* lotsfree is the high-water mark for pageout, and needfree is the
* number of needed free pages. We add extra pages here to make sure
* the scanner doesn't start up while we're freeing memory.

* take 'desfree' extra
pages, so we reclaim
sooner, rather than later

extra = desfree;

freemem < lotsfree +
needfree + extra ?

Yes

return (1);

* If we're on an i386 platform, it's possible that we'll exhaust the
* kernel heap space before we ever run out of available physical
* memory. Most checks of the size of the heap_area compare against
* tune.t_minarmem, which is the minimum available real memory that we
* can have in the system. However, this is generally fixed at 25 pages
* which is so low that it's useless. In this comparison, we seek to
* calculate the total heap-size, and reclaim if more than 3/4ths of the
* heap is allocated. (Or, in the calculation, if less than 1/4th is
* free)

* check to make sure that swapfs has enough space so that anon
* reservations can still succeed. anon_resvmem() checks that the
* availmem is greater than swapfs_minfree, and the number of reserved
* swap pages. We also add a bit of extra here just to prevent
* circumstances from getting really dire.

availmem <
swapfs_minfree +
swapfs_reserve + extra ?

No

Yes

defined(_i386) ?

return (1);

* If zio data pages are being allocated out of a separate heap segment,
* then enforce that the size of available vmem for this arena remains
* above about 1/16th free.
*
* Note: The 1/16th arena free requirement was put in place
* to aggressively evict memory from the arc in order to avoid
* memory fragmentation issues.

vmem_size(heap_arena,
VMEM_FREE) <
(vmem_size(heap_arena,
VMEM_FREE |
VMEM_ALLOC) >> 2) ?

Yes

return (1);

zio_arena != NULL &&
vmem_size(zio_arena,
VMEM_FREE) <
(vmem_size(zio_arena,
VMEM_ALLOC) >> 4) ?

Yes

return (1);

No

return (0);

spa_get_random(100)
== 0 ?

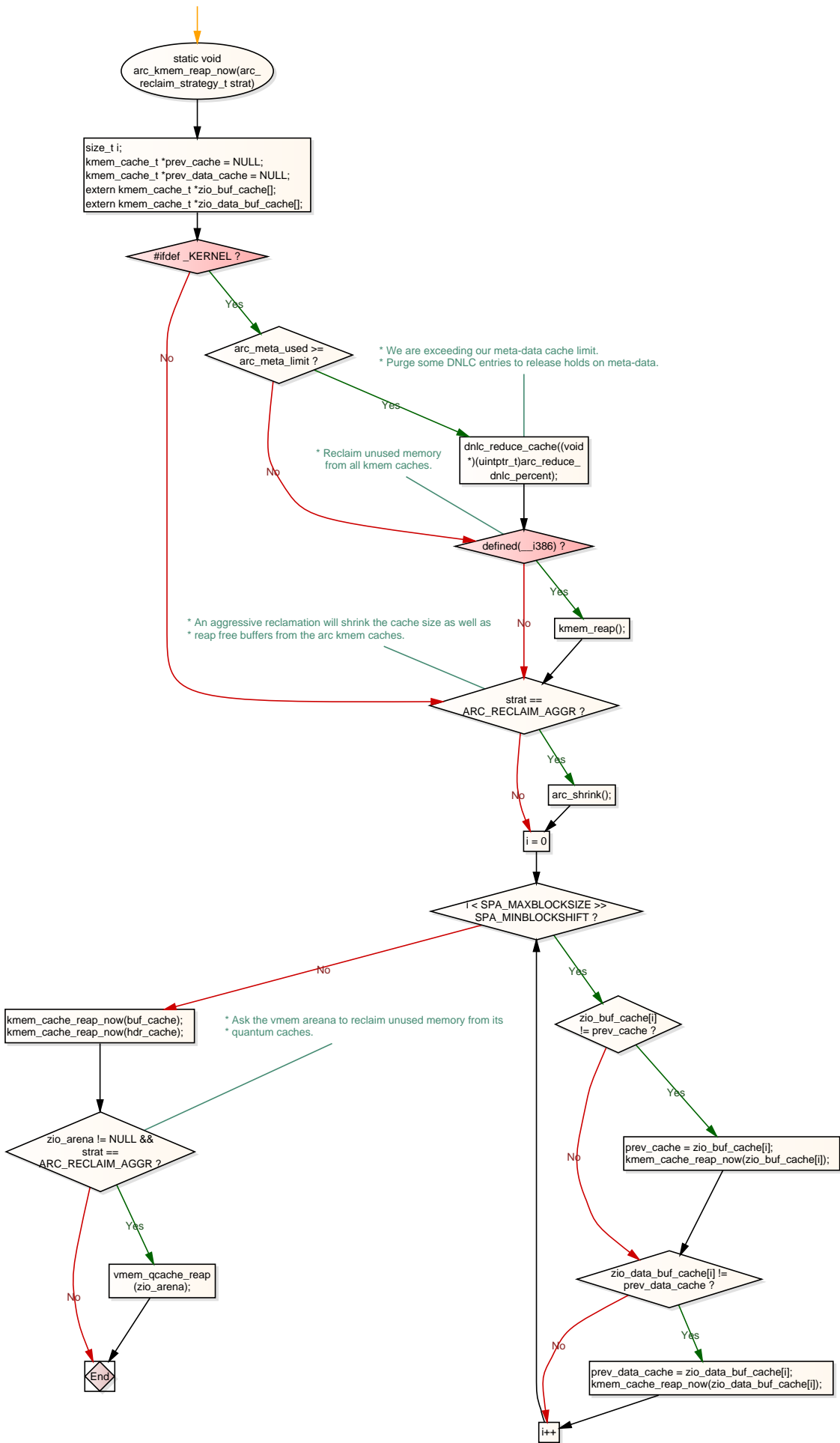
Yes

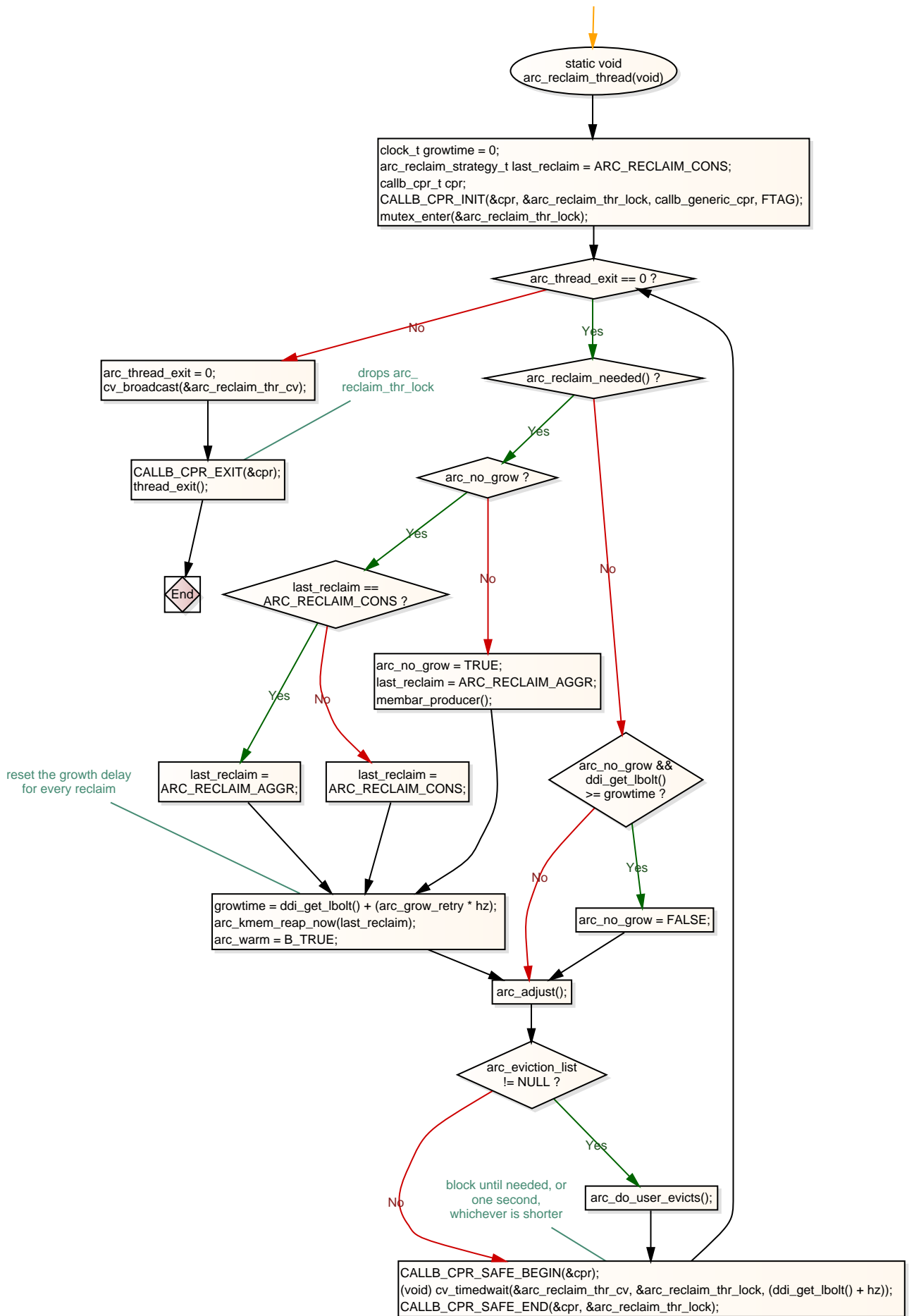
return (1);

No

return (1);

End





* Adapt arc info given the number of bytes we are trying to add and
 * the state that we are coming from. This function is only called
 * when we are adding new content to the cache.

```
int mult;
uint64_t arc_p_min = (arc_c >> arc_p_min_shift);
```

state == arc_l2c_only ?

Yes
ASSERT(bytes > 0);

* Adapt the target size of the MRU list:
 * - if we just hit in the MRU ghost list, then increase
 * the target size of the MRU list.
 * - if we just hit in the MFU ghost list, then increase
 * the target size of the MFU list by decreasing the
 * target size of the MRU list.

state == arc_mru_ghost ?

No

state == arc_mfu_ghost ?

Yes

```
mult = ((arc_mru_ghost->arcs_size >=
arc_mfu_ghost->arcs_size)
? 1 : (arc_mfu_ghost->arcs_size/arc_mru_ghost->arcs_size));
```

avoid wild arc_p
adjustment

```
uint64_t delta;
mult = ((arc_mfu_ghost->arcs_size >= arc_mru_ghost->arcs_size)
? 1 : (arc_mru_ghost->arcs_size/arc_mfu_ghost->arcs_size));
mult = MIN(mult, 10);
delta = MIN(bytes * mult, arc_p);
arc_p = MAX(arc_p_min, arc_p - delta);
```

```
mult = MIN(mult, 10);
arc_p = MIN(arc_c - arc_p_min, arc_p + bytes * mult);
```

ASSERT((int64_t)arc_p
 >= 0);

arc_reclaim_needed() ?

Yes

cv_signal(&arc_reclaim_thr_cv);

arc_no_grow ?

No

arc_c >= arc_c_max ?

Yes

* If we're within (2 * maxblocksize) bytes of the target
 * cache size, increment the target cache size

arc_size > arc_c - (2ULL
 << SPA_MAXBLOCKSHIFT) ?

No

```
atomic_add_64(&arc_c,
(int64_t)bytes);
```

arc_c > arc_c_max ?

No

state == arc_anon ?

Yes

```
atomic_add_64(&arc_p,
(int64_t)bytes);
```

```
arc_c = arc_c_max;
```

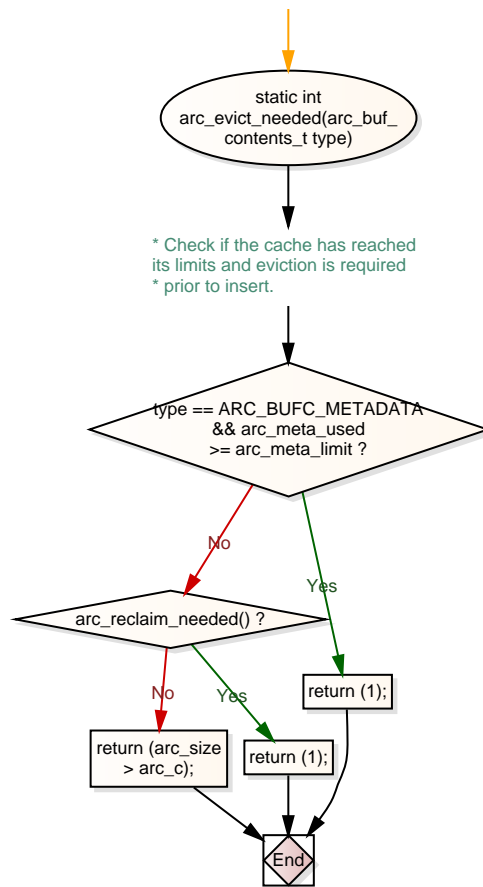
arc_p > arc_c ?

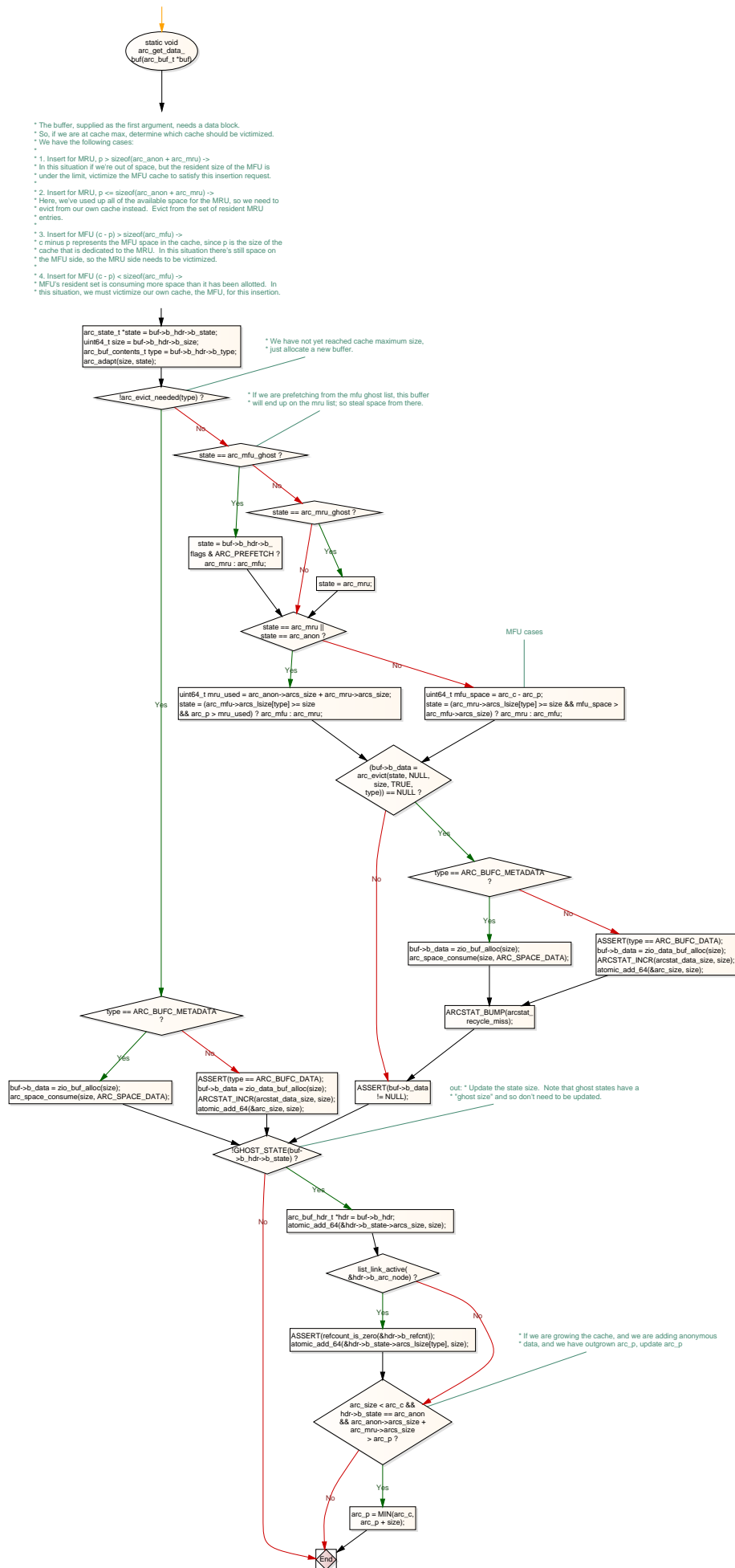
Yes

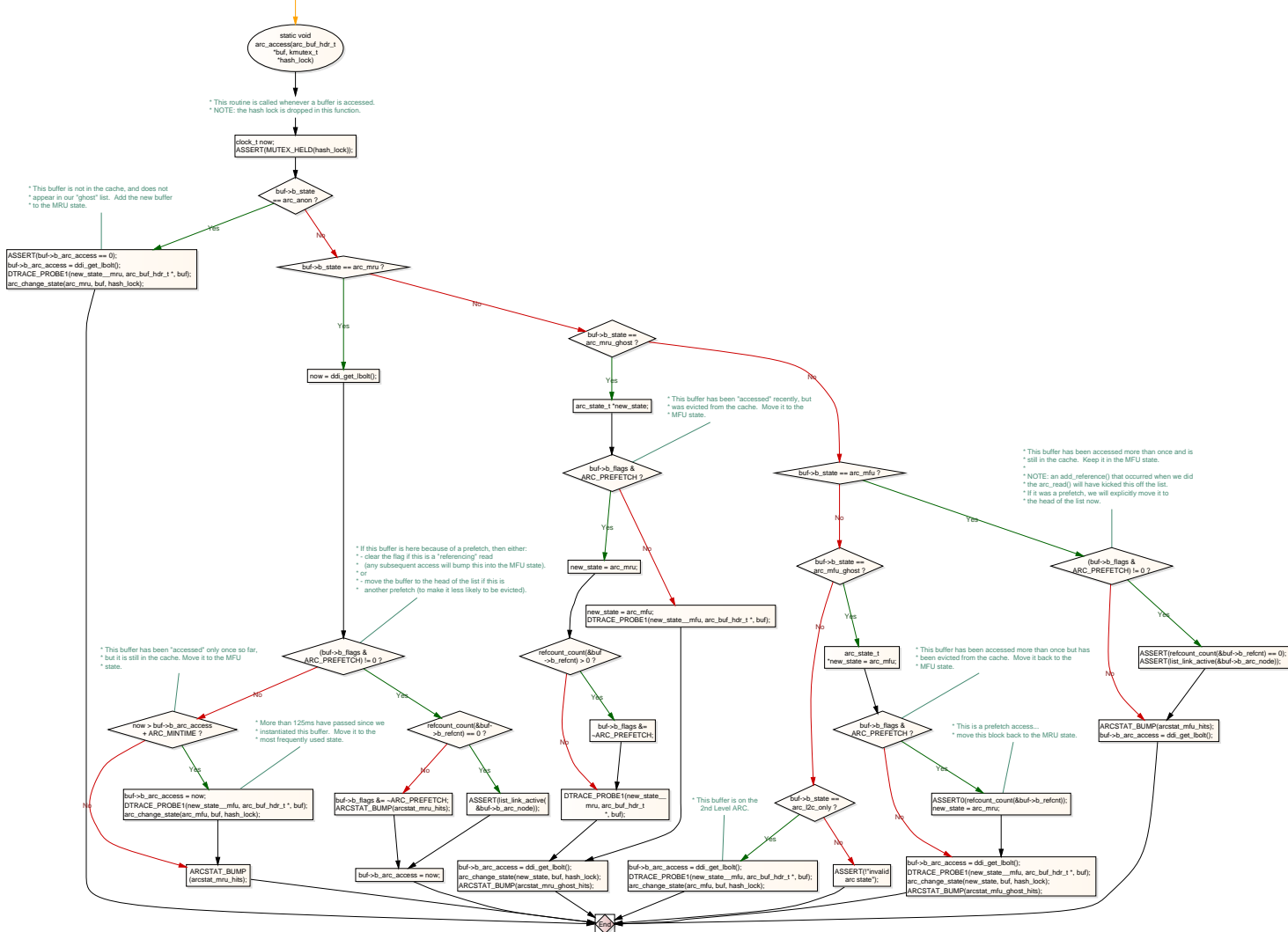
```
arc_p = arc_c;
```

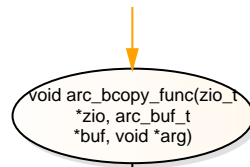
ASSERT((int64_t)arc_p
 >= 0);

End

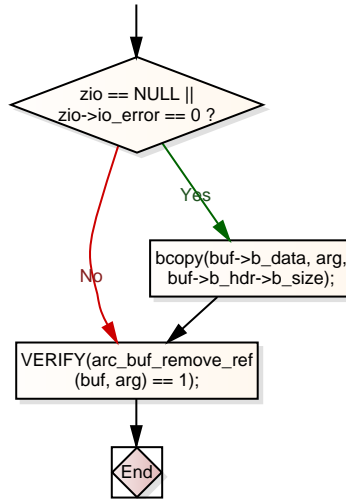


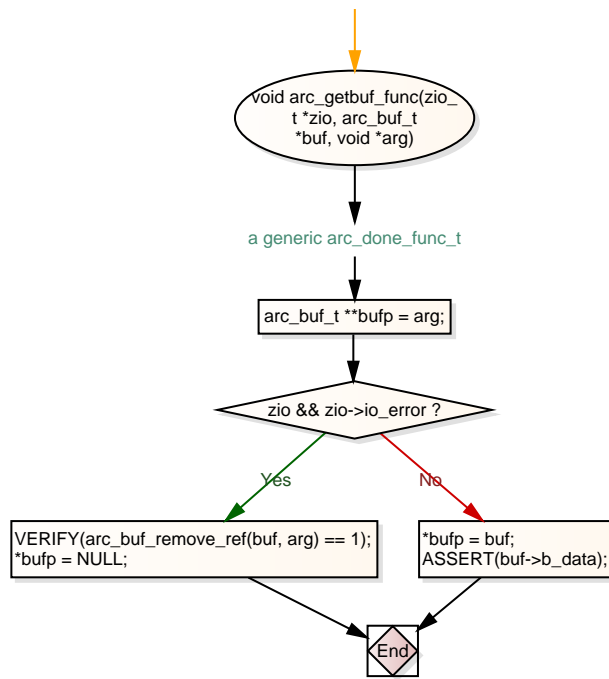


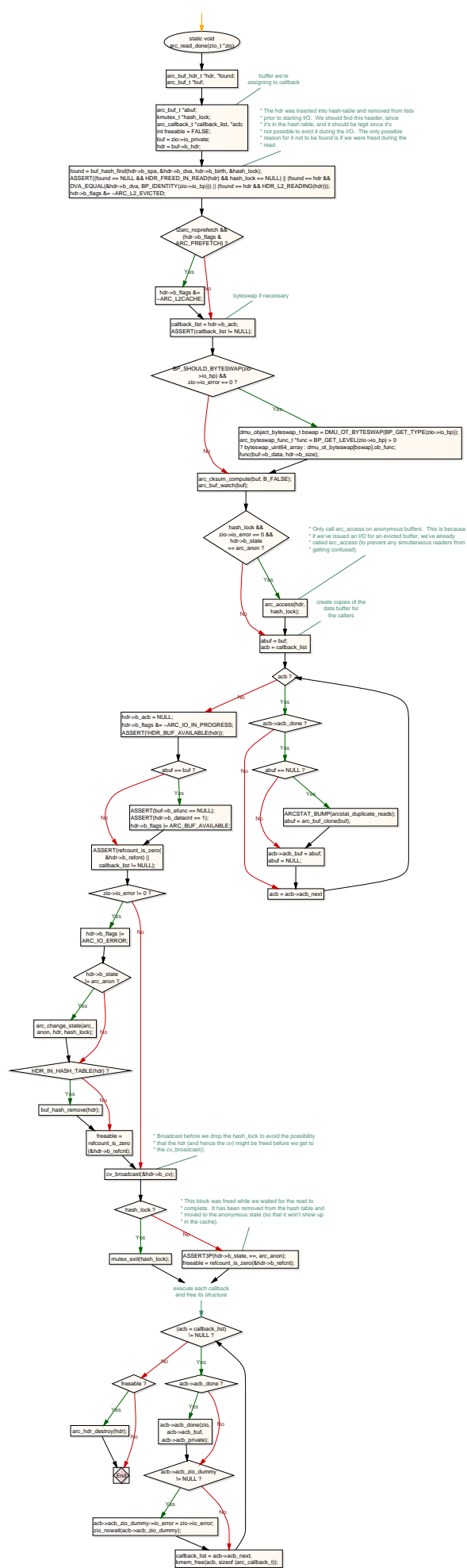


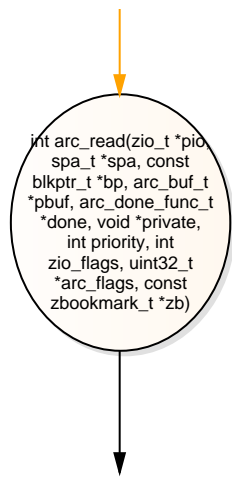


a generic `arc_done_func_t` which you can use
`ARGSUSED`

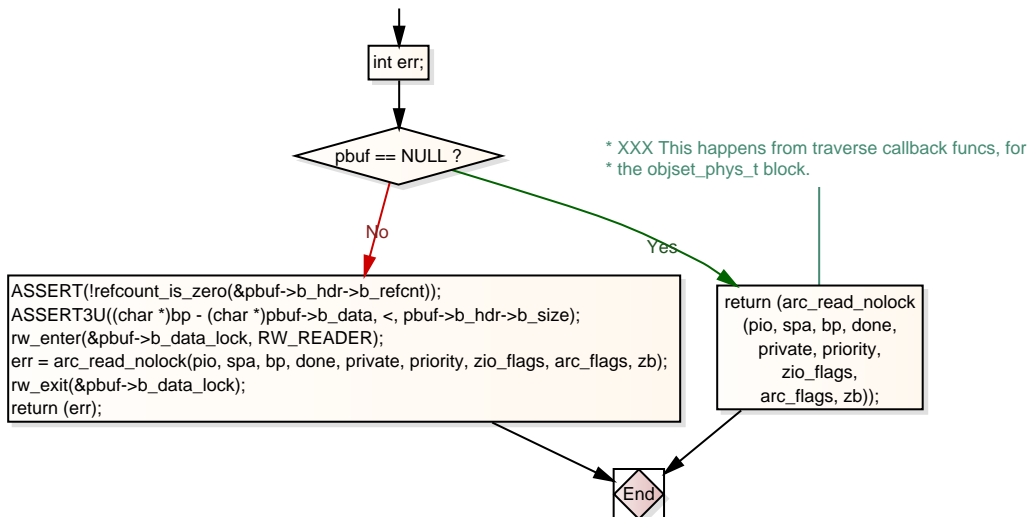


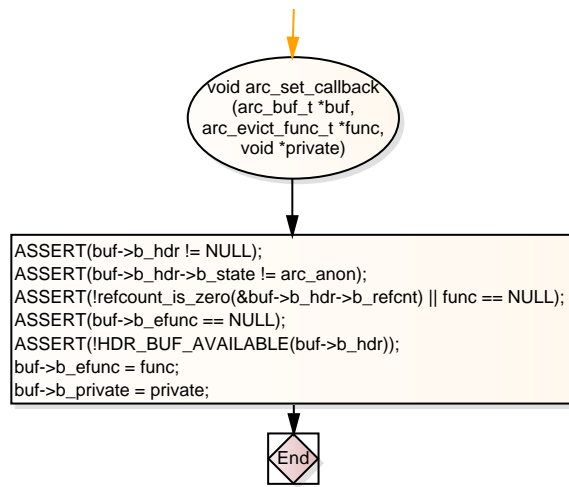


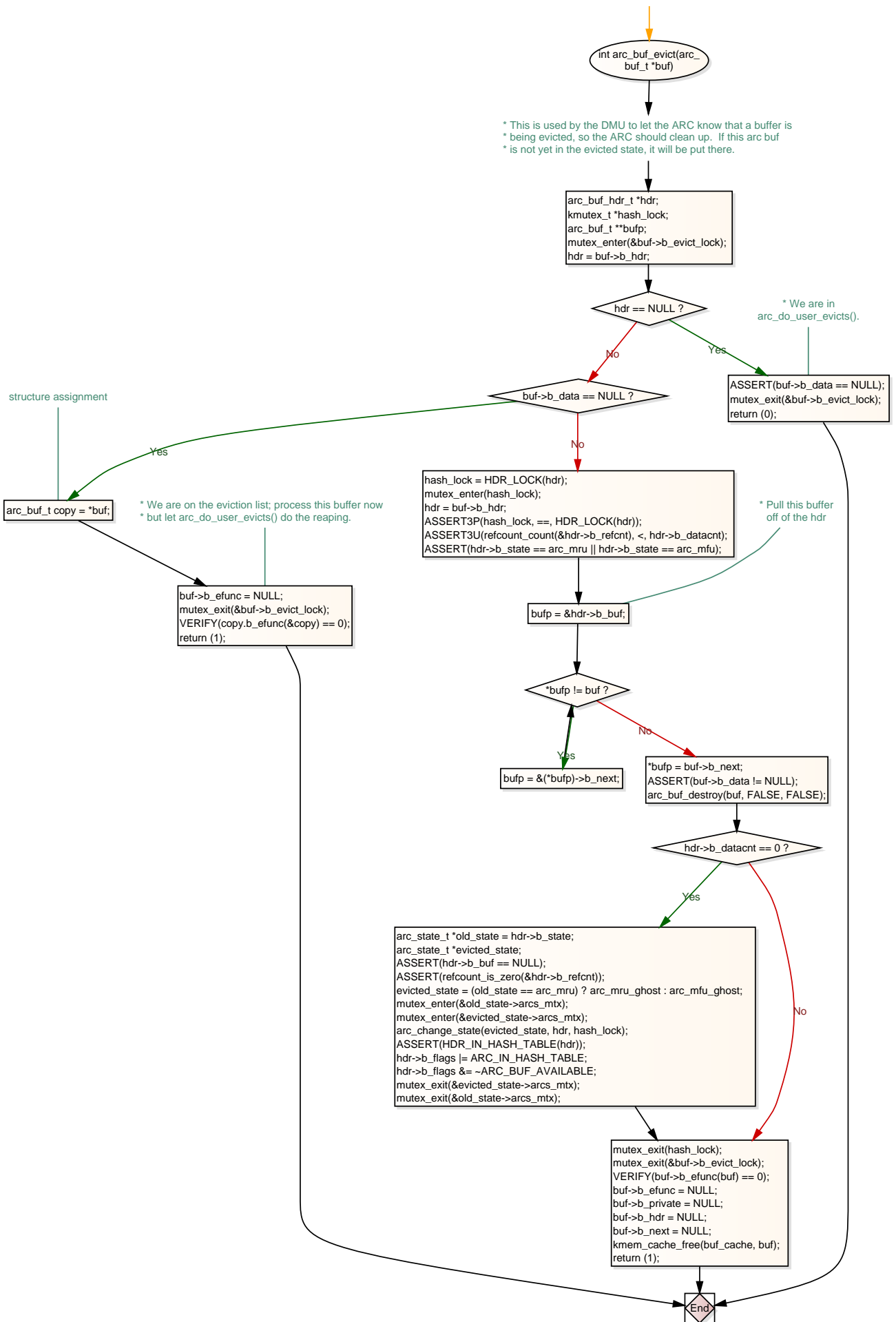


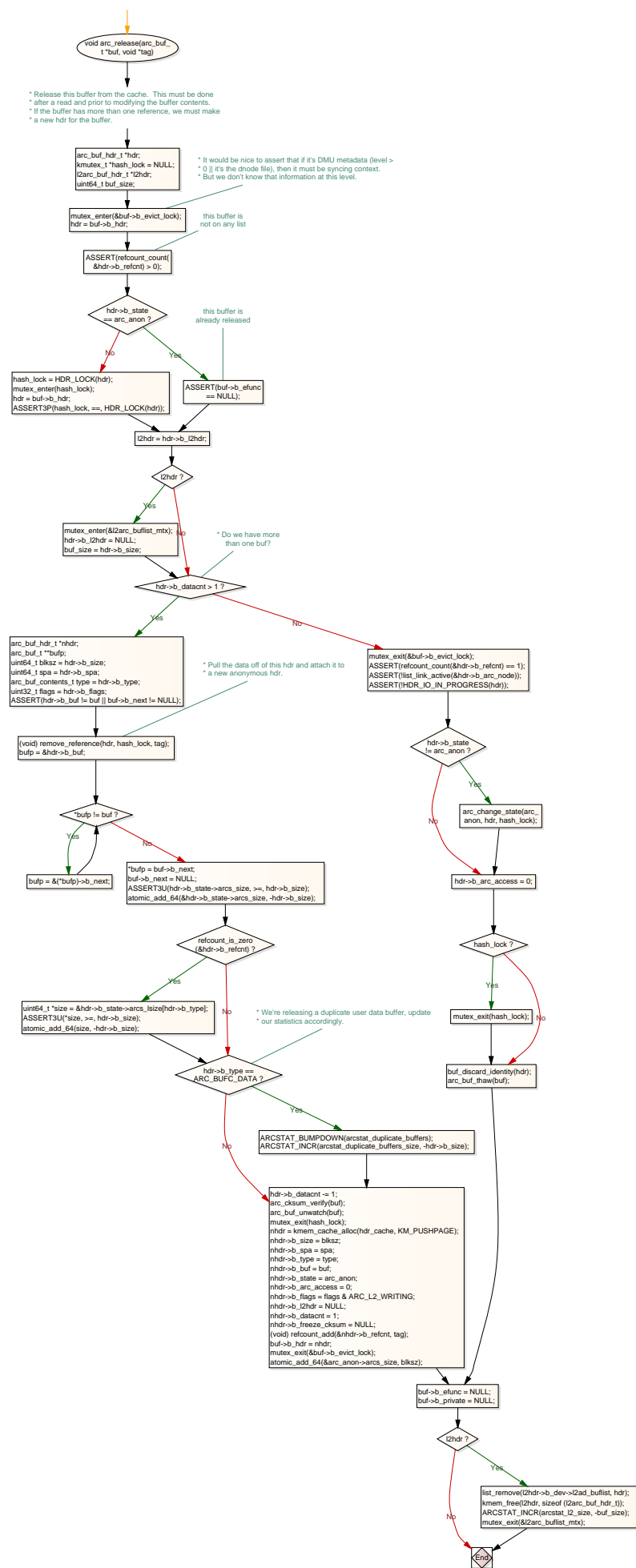


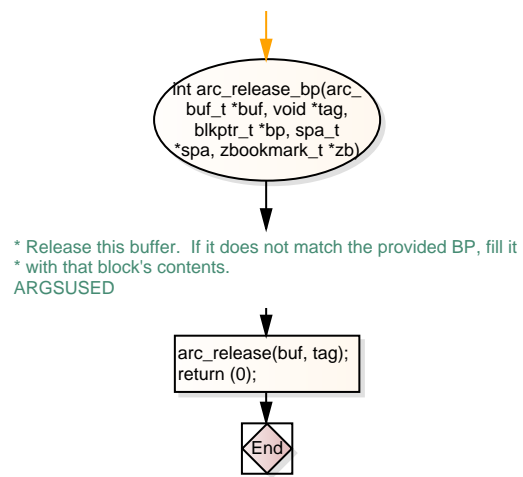
* "Read" the block at the specified DVA (in bp) via the
 * cache. If the block is found in the cache, invoke the provided
 * callback immediately and return. Note that the 'zio' parameter
 * in the callback will be NULL in this case, since no IO was
 * required. If the block is not in the cache pass the read request
 * on to the spa with a substitute callback function, so that the
 * requested block will be added to the cache.
 *
 * If a read request arrives for a block that has a read in-progress,
 * either wait for the in-progress read to complete (and return the
 * results); or, if this is a read with a "done" func, add a record
 * to the read to invoke the "done" func when the read completes,
 * and return; or just return.
 *
 * arc_read_done() will invoke all the requested "done" functions
 * for readers of this block.
 *
 * Normal callers should use arc_read and pass the arc buffer and offset
 * for the bp. But if you know you don't need locking, you can use
 * arc_read_bp.

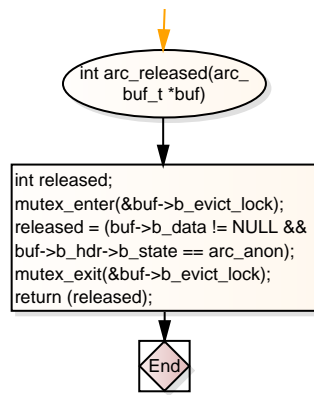


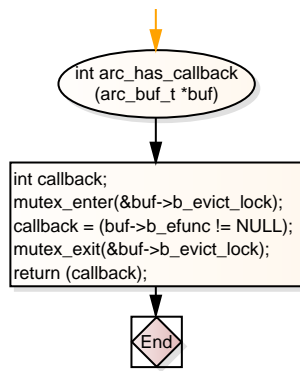


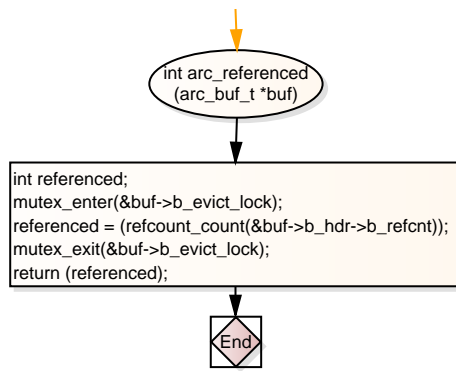


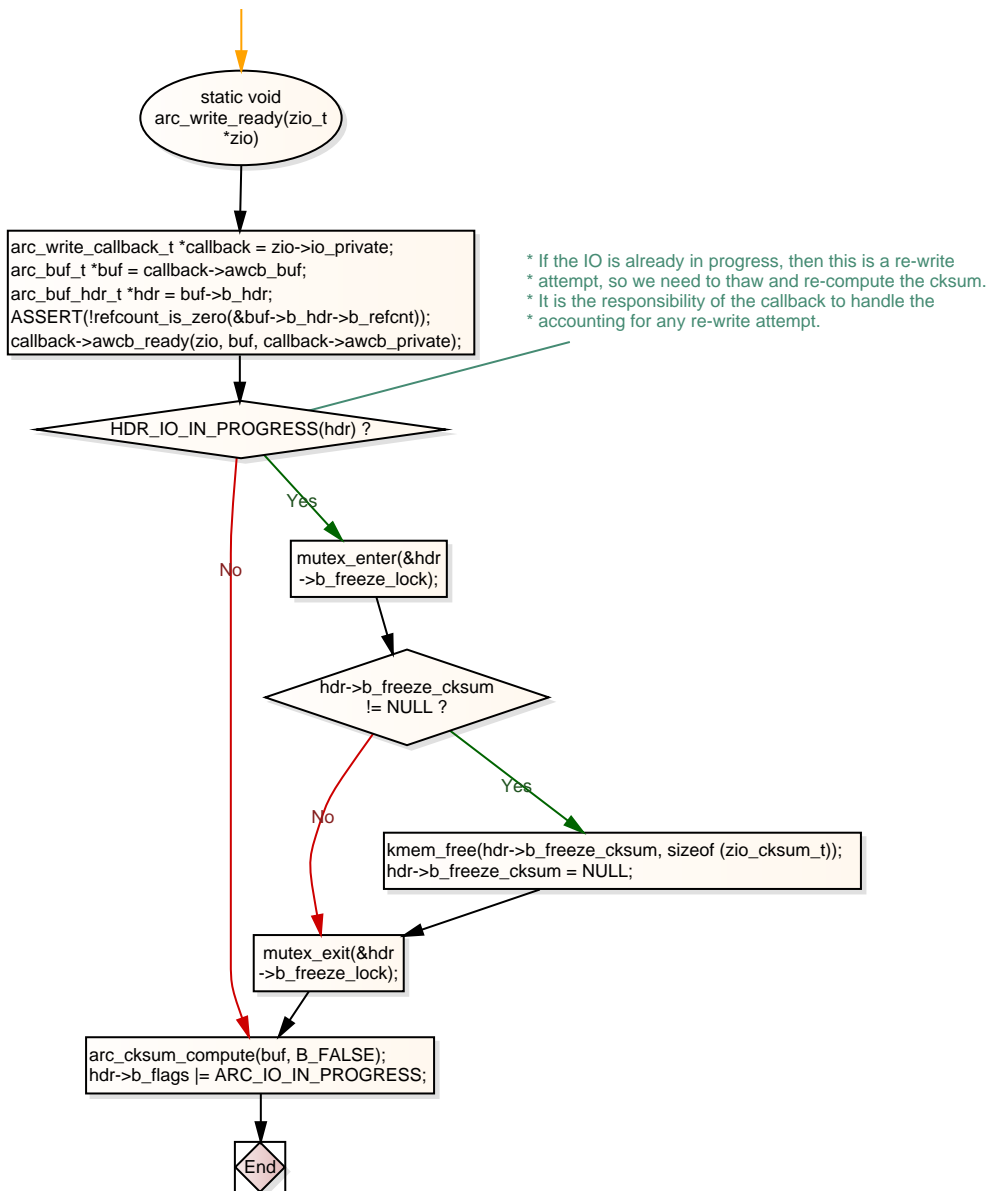


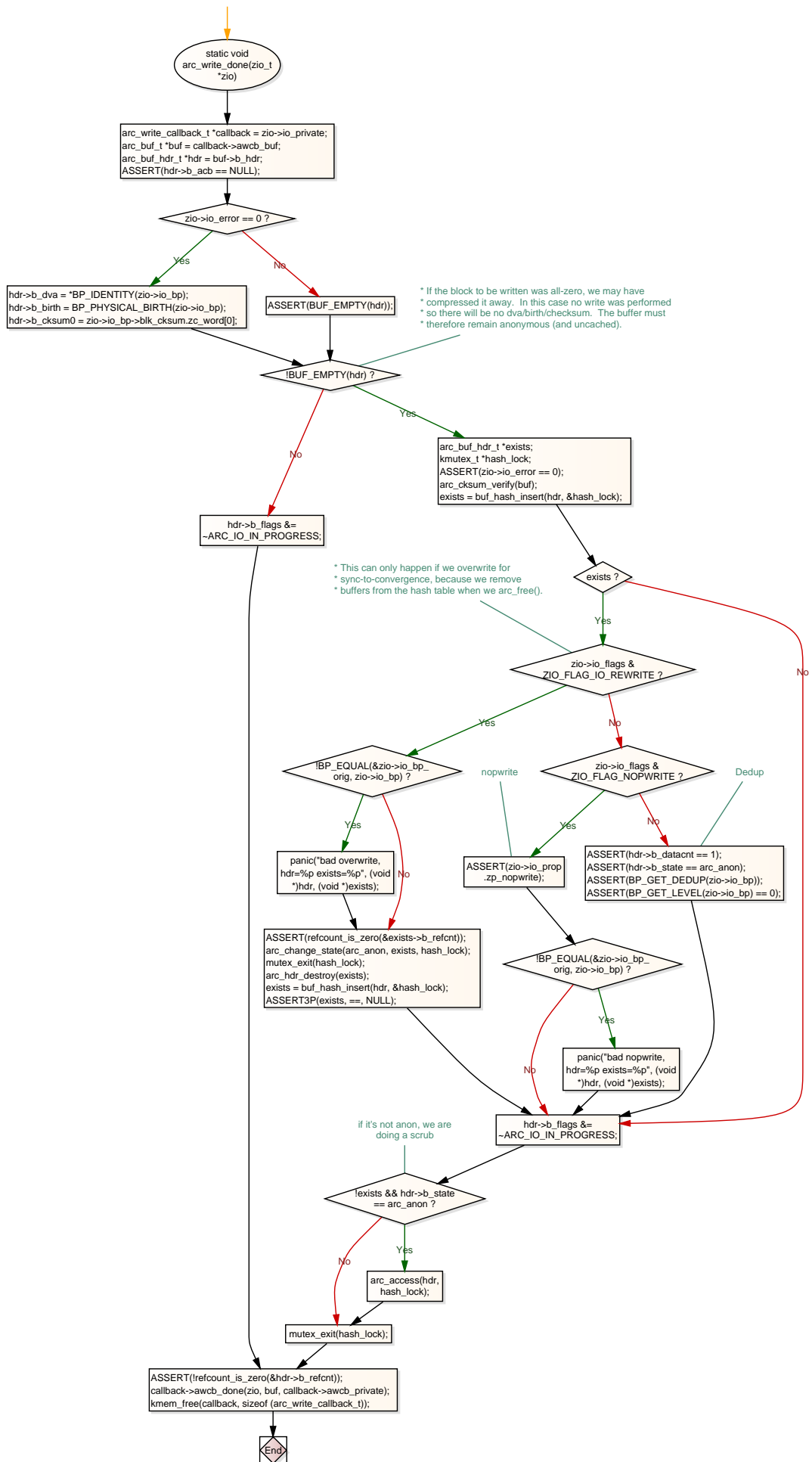


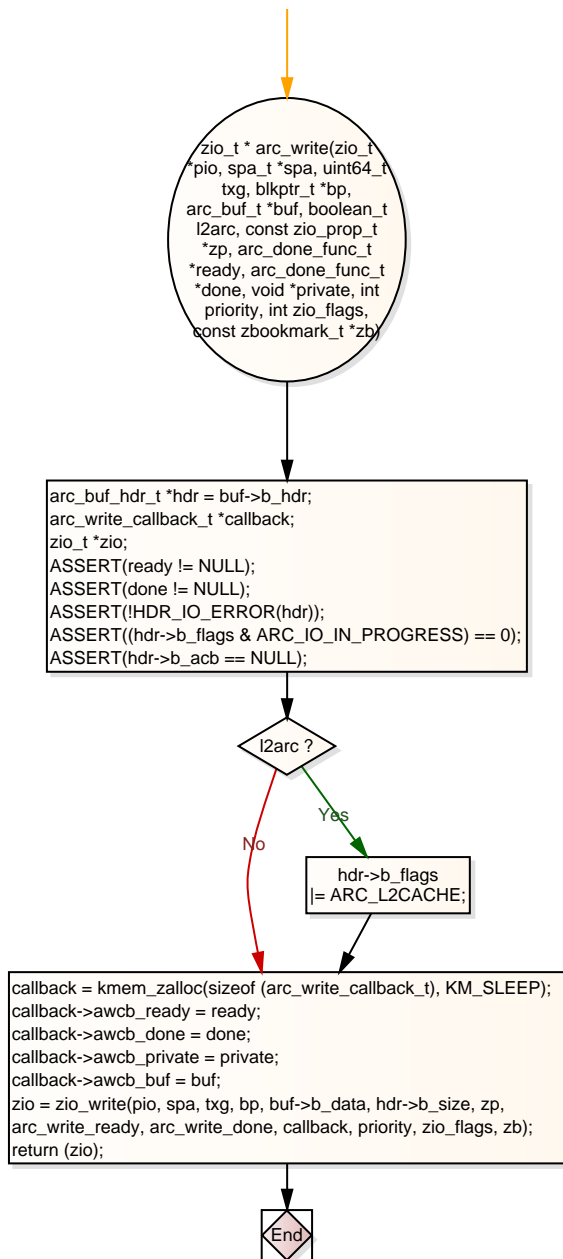


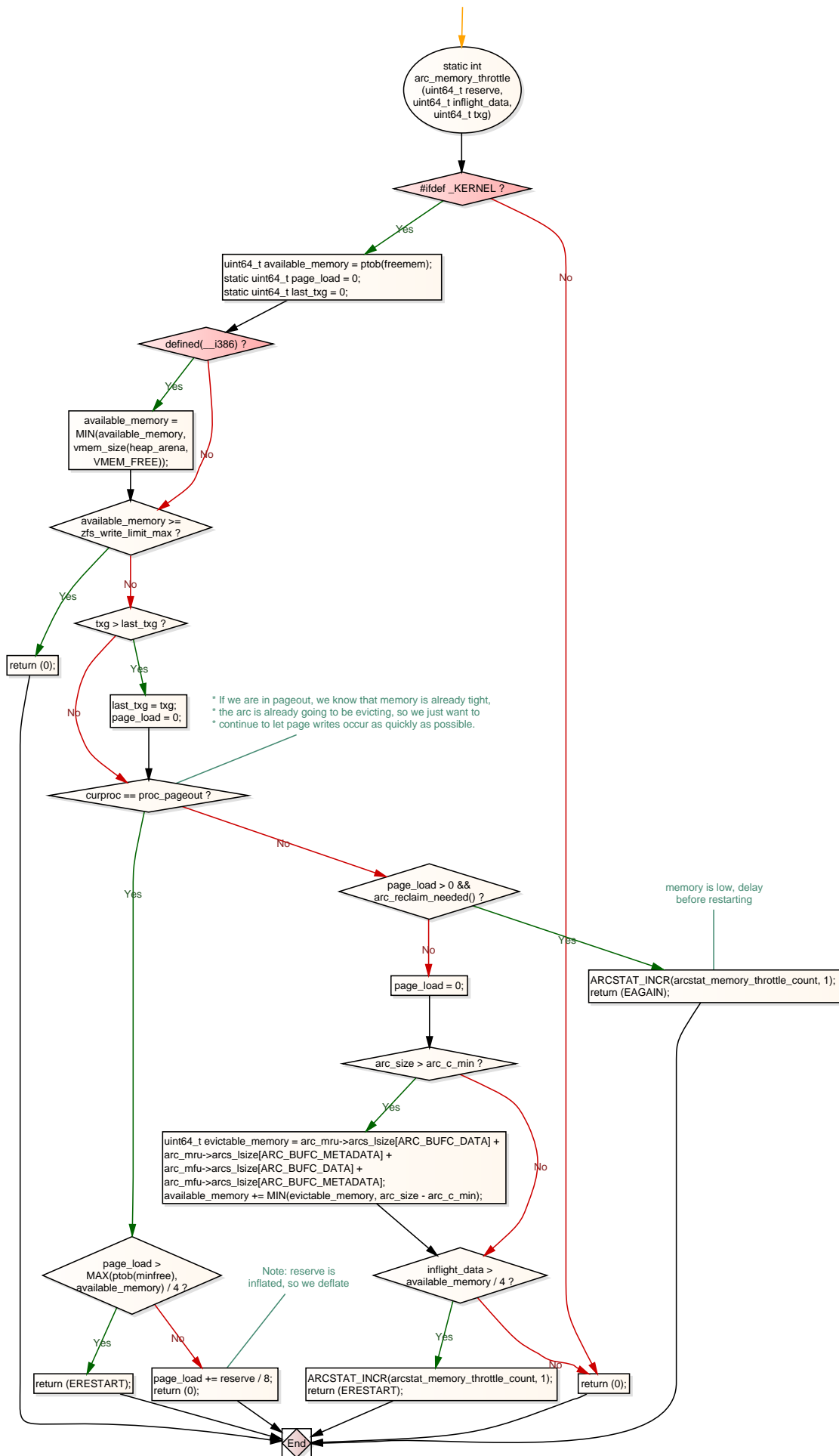


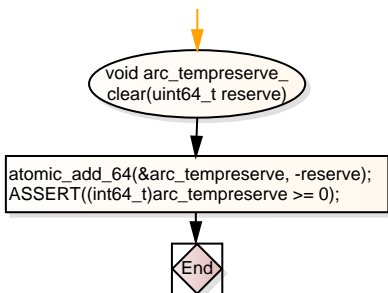


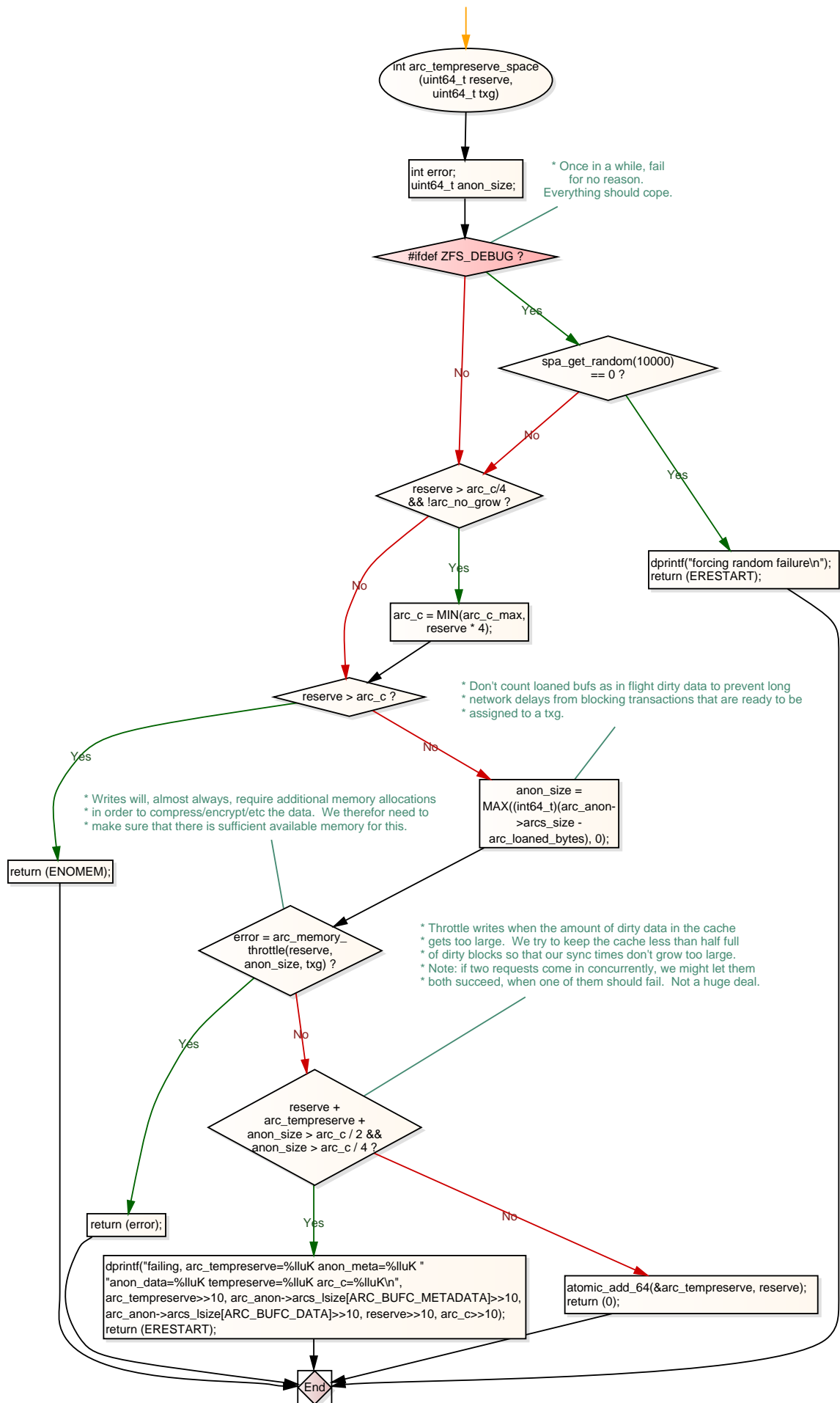


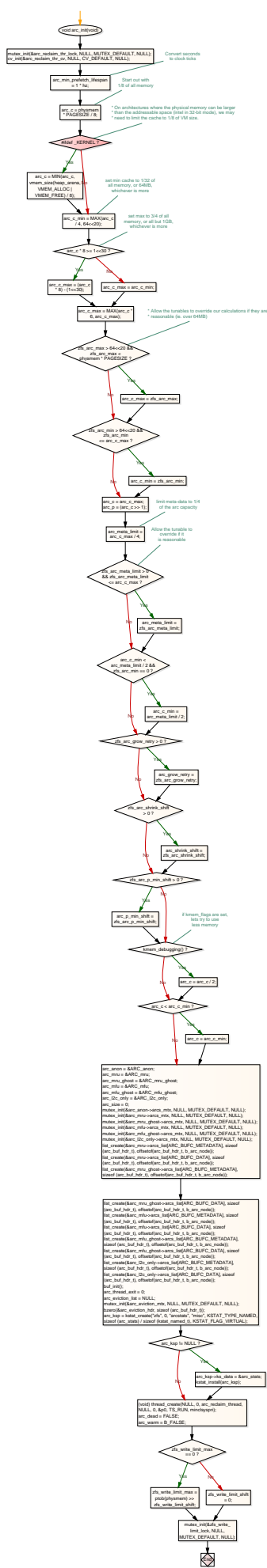


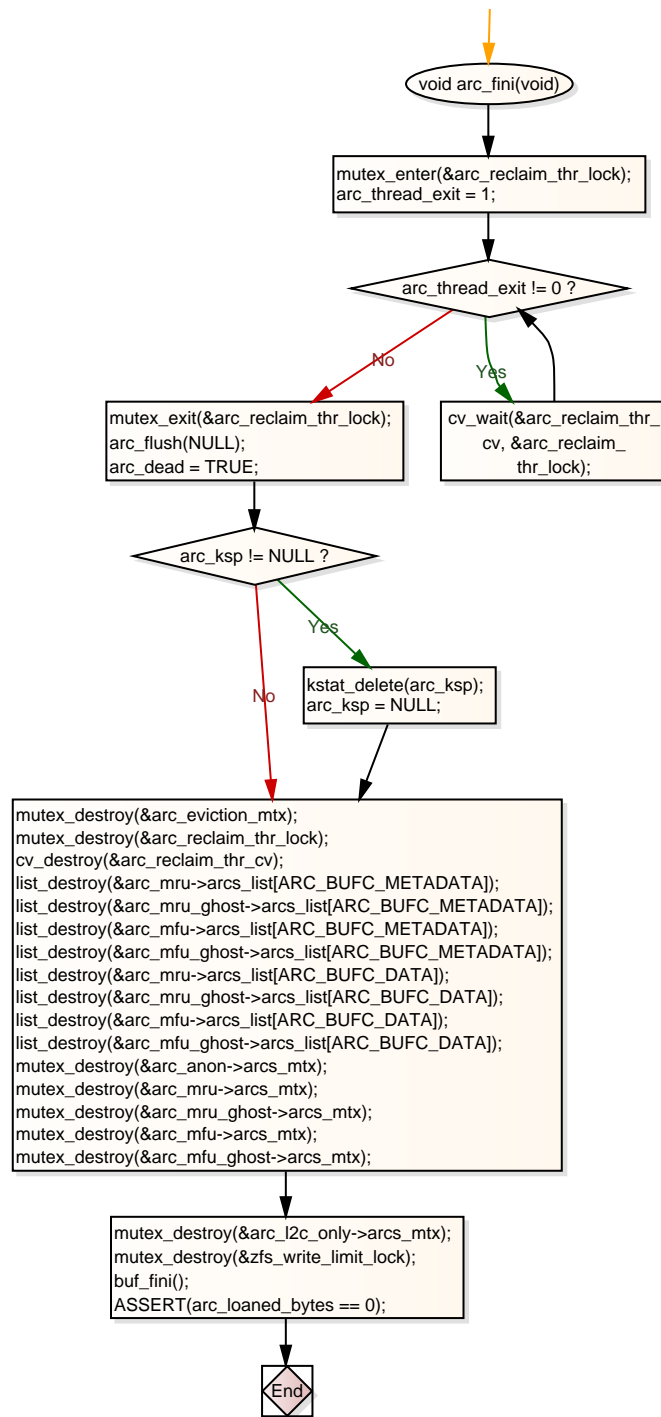












static boolean, i
l2arc_writ_eligible
(uint64_t spa_guid,
arc_buf_hdr_t *b0)

* Level 2 ARC

- * The level 2 ARC (L2ARC) is a cache layer in-between main memory and disk.
- * It uses dedicated storage devices to hold cached data, which are populated using large infrequent writes. The main role of this cache is to boost the performance of random read workloads. The intended L2ARC devices include short-stroked disks, solid state disks, and other media with substantially faster read latency than disk.



* Read requests are satisfied from the following sources, in order:

- 1) ARC
- 2) vdev cache of L2ARC devices
- 3) L2ARC devices
- 4) vdev cache of disks
- 5) disks

- * Some L2ARC device types exhibit extremely slow write performance.
- * To accommodate for this there are some significant differences between the L2ARC and traditional cache design:

1. There is no eviction path from the ARC to the L2ARC. Evictions from the ARC behave as usual, freeing buffers and placing headers on ghost lists. The ARC does not send buffers to the L2ARC during eviction as this would add inflated write latencies for all ARC memory pressure.

2. The L2ARC attempts to cache data from the ARC before it is evicted. It does this by periodically scanning buffers from the eviction-end of the MFU and MRU ARC lists, copying them to the L2ARC devices if they are not already there. It scans until a headroom of buffers is satisfied, which itself is a buffer for ARC eviction. The thread that does this is (l2arc_feed_thread), illustrated below, example sizes are included to provide a better sense of ratio than this diagram:



3. If an ARC buffer is copied to the L2ARC but then hit instead of evicted, then the L2ARC has cached a buffer much sooner than it probably needed to, potentially wasting L2ARC device bandwidth and storage. It is safe to say that this is an uncommon case, since buffers at the end of the ARC lists have moved there due to inactivity.

4. If the ARC evicts faster than the L2ARC can maintain a headroom, then the L2ARC simply misses copying some buffers. This serves as a pressure valve to prevent heavy read workloads from both stalling the ARC with waits and clogging the L2ARC with writes. This also helps prevent the potential for the L2ARC to churn if it attempts to cache content too quickly, such as during backups of the entire pool.

5. After system boot and before the ARC has filled main memory, there are no evictions from the ARC and so the tails of the ARC_mfu and ARC_mru lists can remain mostly static. Instead of searching from tail of these lists as pictured, the (l2arc_feed_thread) will search from the list heads for eligible buffers, greatly increasing its chance of finding them.

6. The L2ARC device write speed is also boosted during this time so that the L2ARC seems up faster. Since there have been no ARC evictions yet, there are no L2ARC reads, and no fear of degrading read performance through increased writes.

7. Writes to the L2ARC devices are grouped and sent in-sequence, so that the vdev queue can aggregate them into larger and fewer writes. Each device is written to in a rotor fashion, sweeping writes through available space then repeating.

8. The L2ARC does not store dirty content. It never needs to flush write buffers back to disk based storage.

9. If an ARC buffer is written (and dirtied) which also exists in the L2ARC, the now stale L2ARC buffer is immediately dropped.

- * The performance of the L2ARC can be tweaked by a number of tunables, which may be necessary for different workloads:

- | | |
|-------------------|---|
| l2arc_write_max | max write bytes per interval |
| l2arc_write_boost | extra write bytes during device warmup |
| l2arc_noprefetch | skip caching prefetched buffers |
| l2arc_headroom | number of max device writes to precache |
| l2arc_feed_secs | seconds between L2ARC writing |

- * Tunables may be removed or added as future performance improvements are integrated, and also may become pool properties.

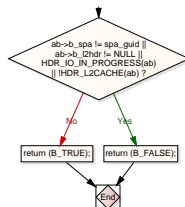
- * There are three key functions that control how the L2ARC warms up:

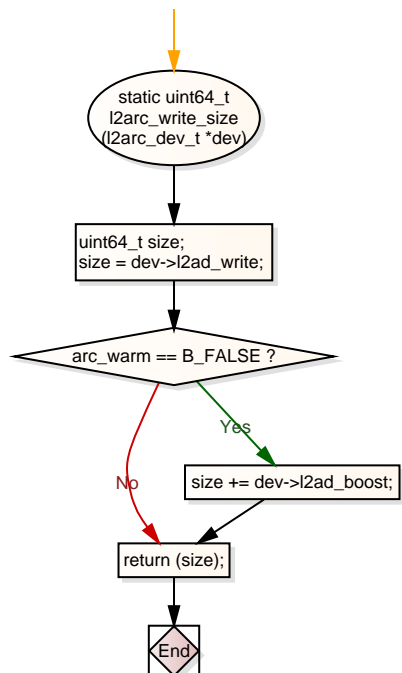
- | | |
|----------------------|--|
| l2arc_writ_eligible | check if a buffer is eligible to cache |
| l2arc_write_size | calculate how much to write |
| l2arc_write_interval | calculate sleep delay between writes |

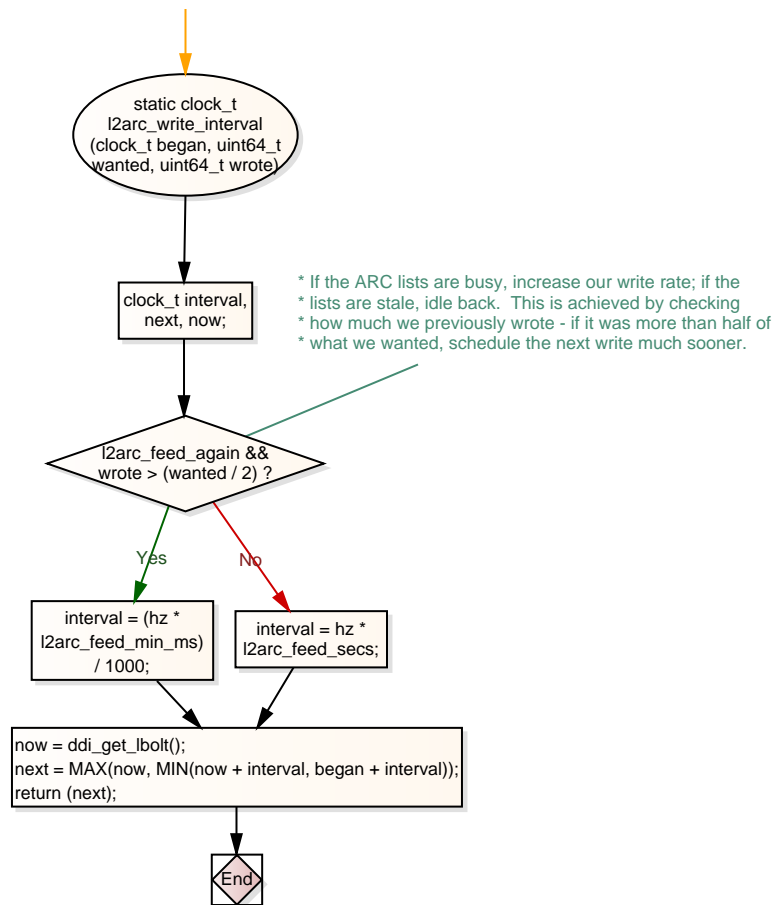
- * These three functions determine what to write, how much, and how quickly to send writes.

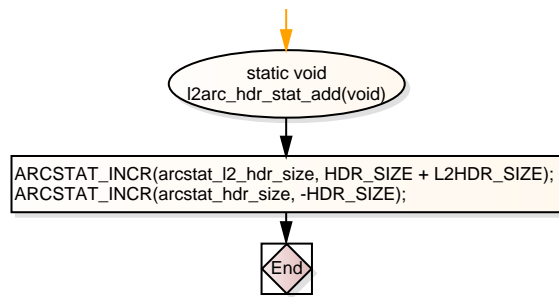
- * A buffer is "not" eligible for the L2ARC if it:

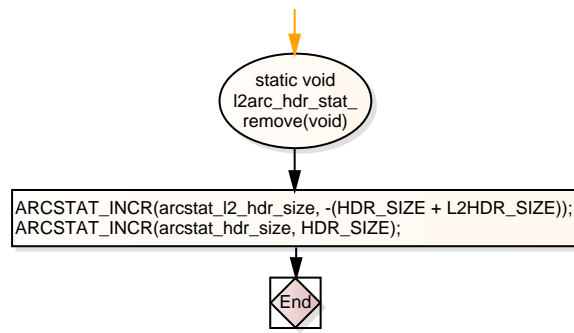
1. belongs to a different spa.
2. is already cached on the L2ARC.
3. has an I/O in progress (it may be an incomplete read).
4. is flagged not eligible (zfs property).

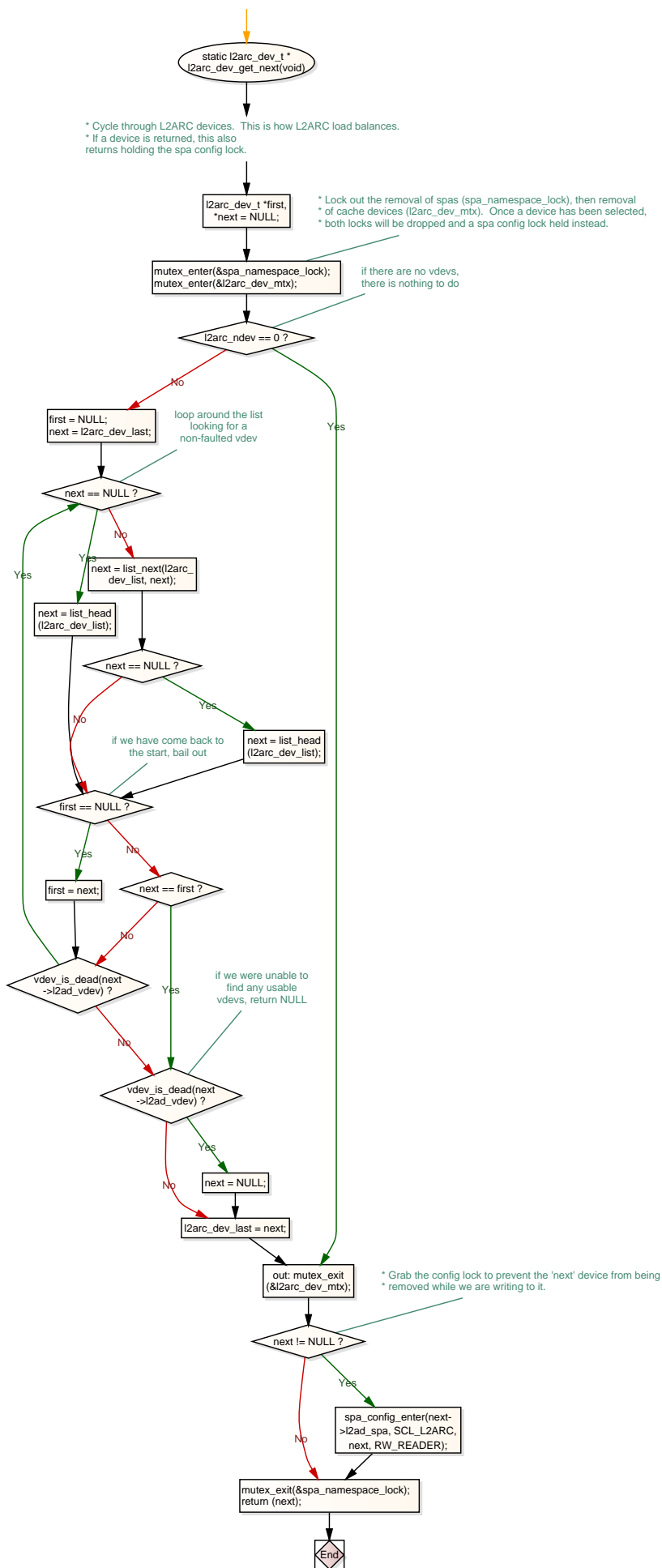


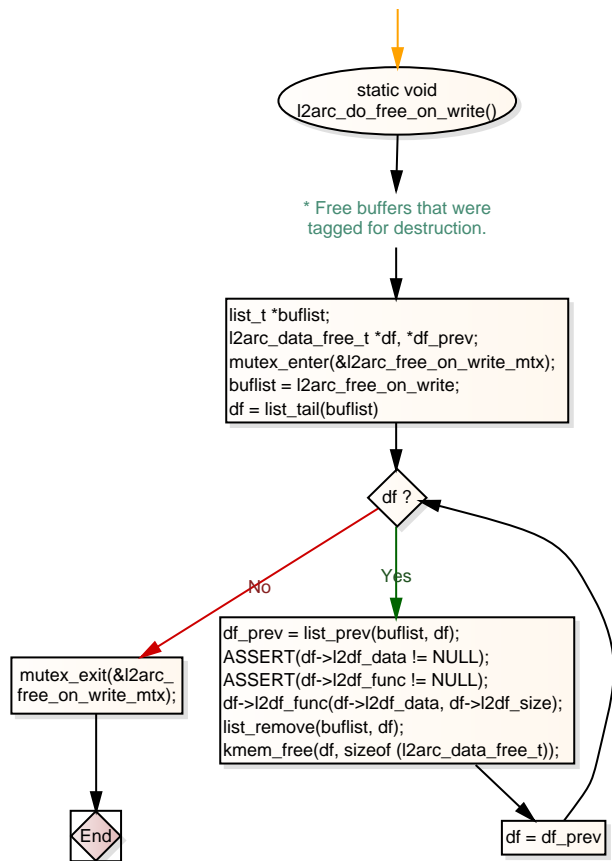


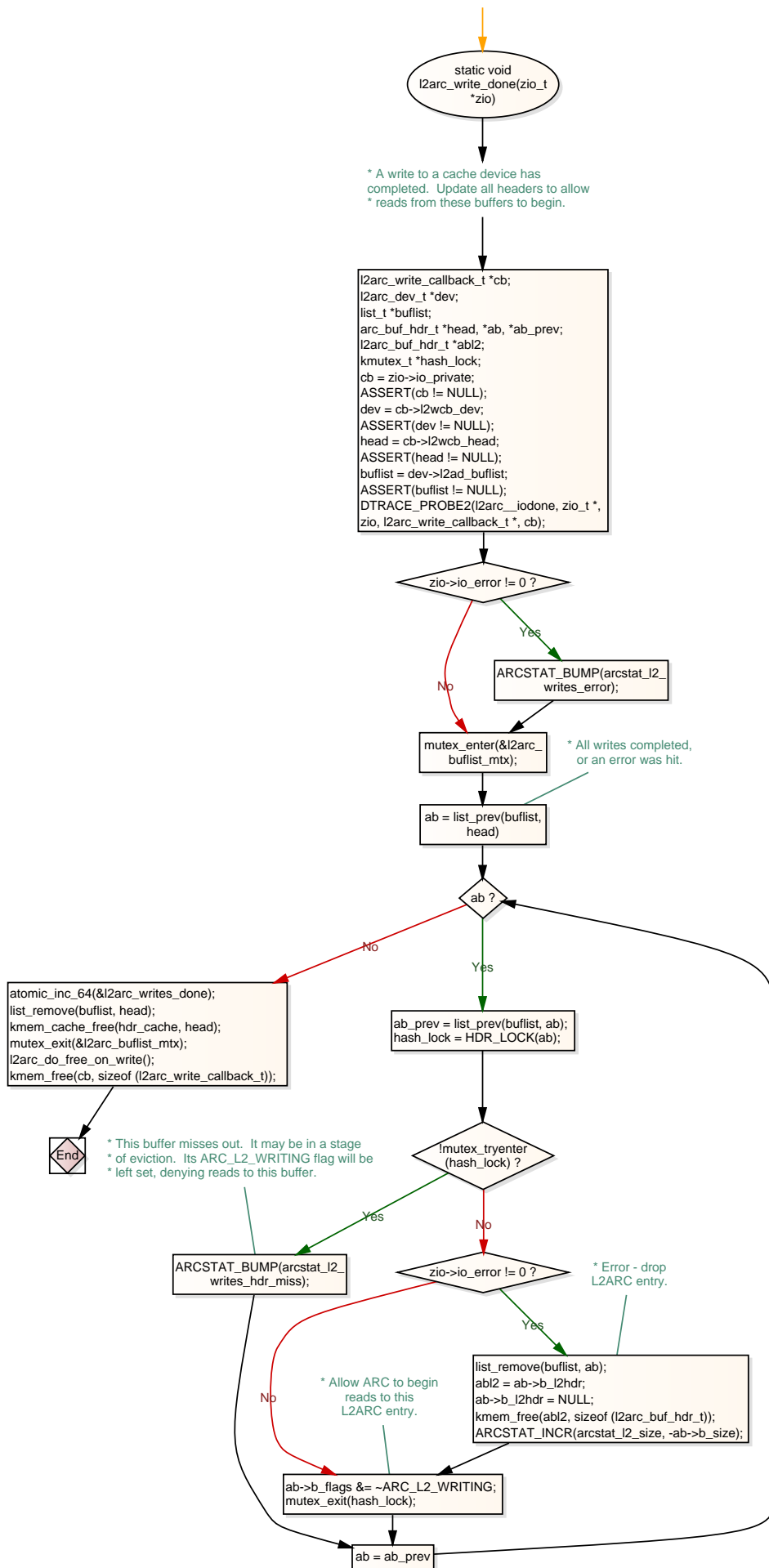


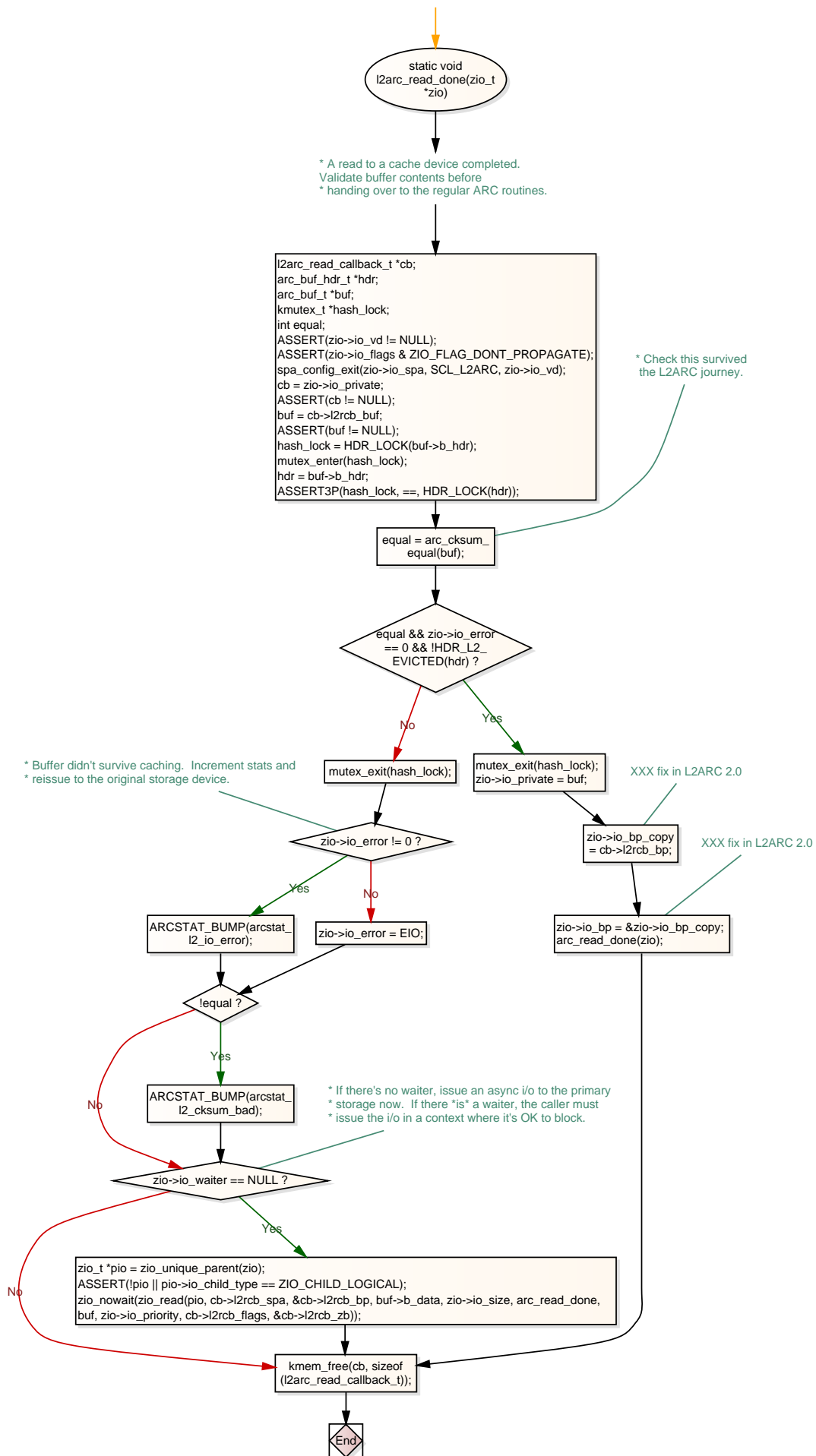


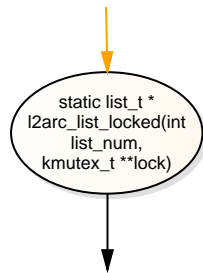




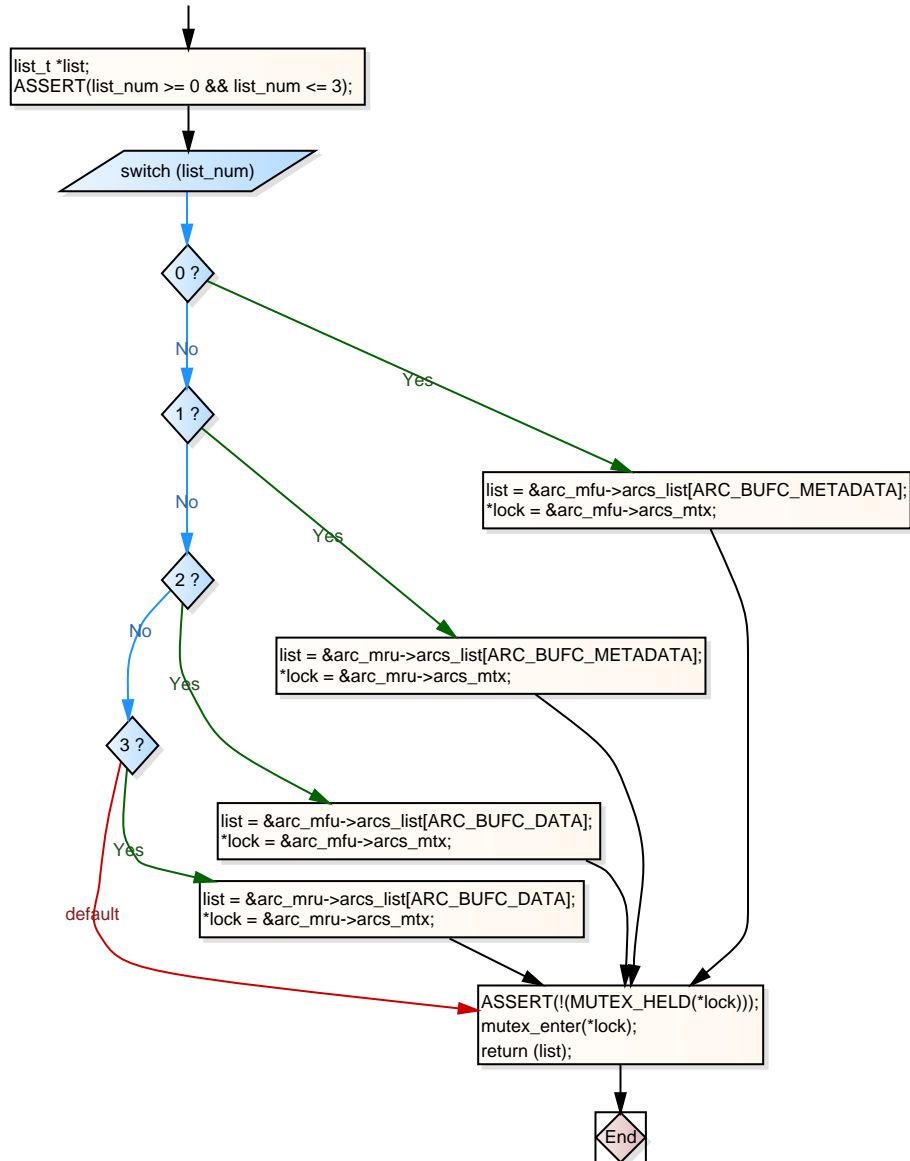


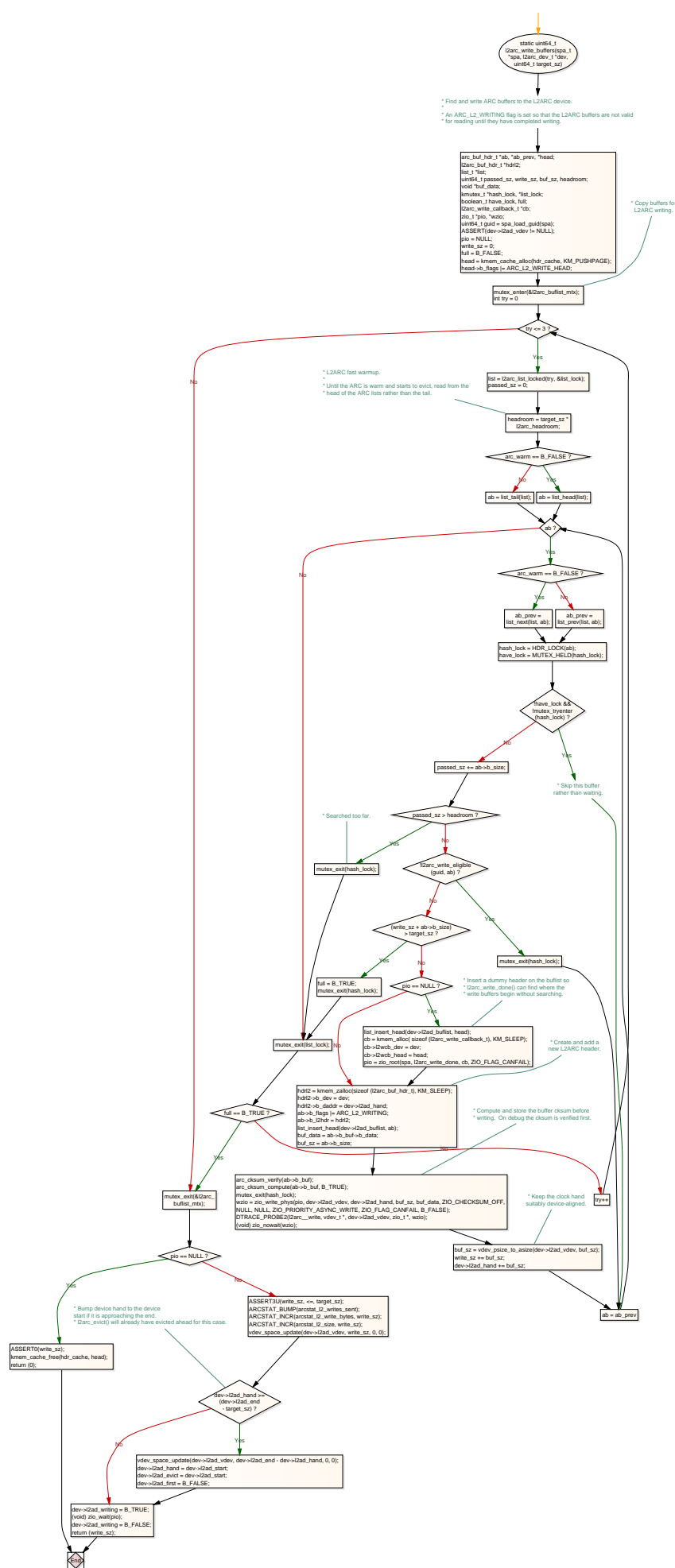


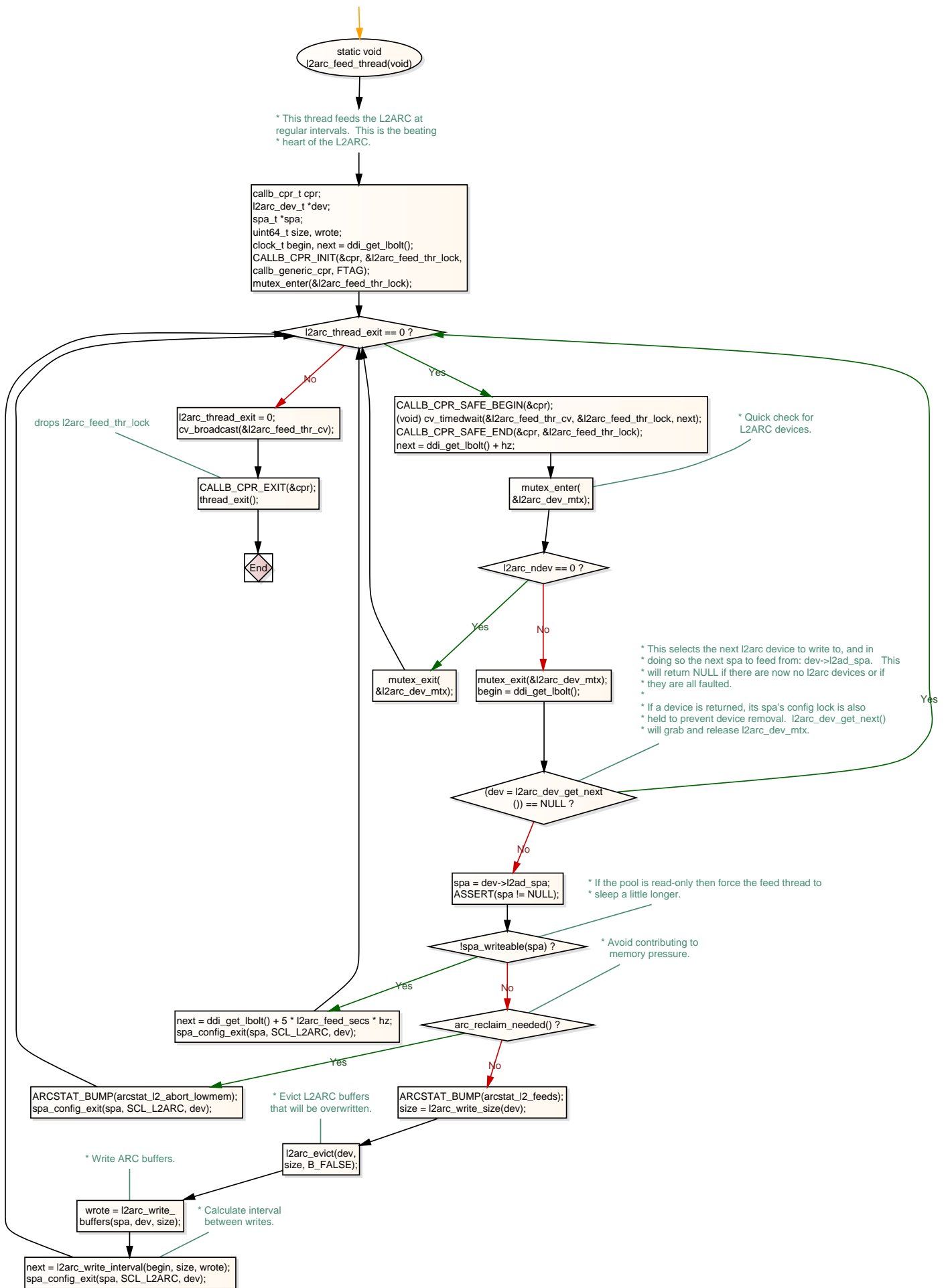


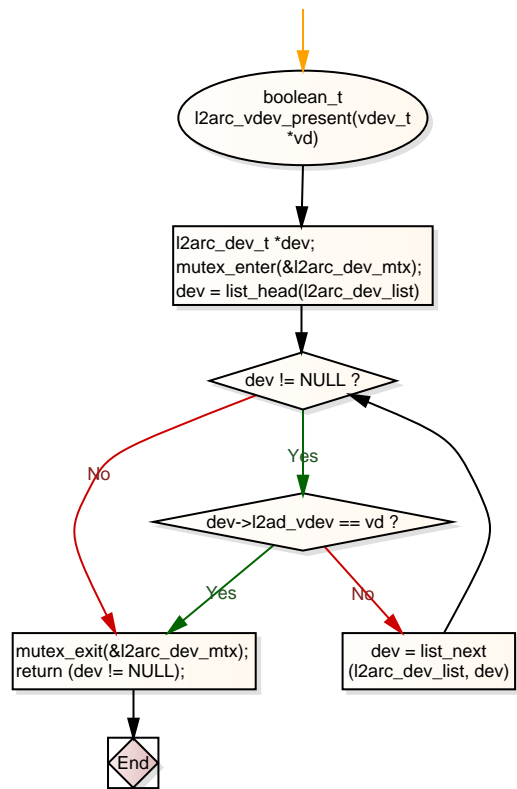


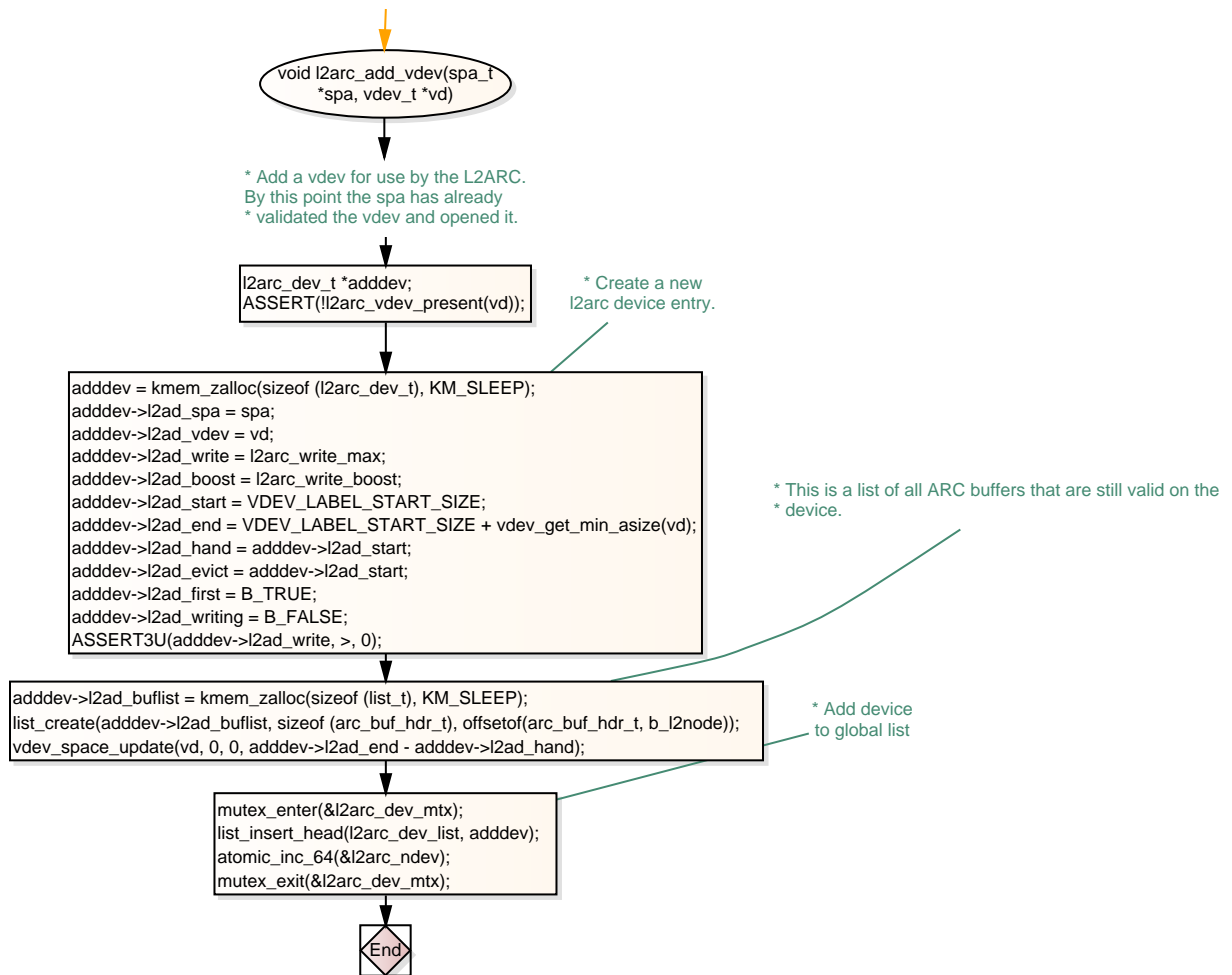
* This is the list priority from which the L2ARC will search for pages to
 * cache. This is used within loops (0..3) to cycle through lists in the
 * desired order. This order can have a significant effect on cache
 * performance.
 *
 * Currently the metadata lists are hit first, MFU then MRU, followed by
 * the data lists. This function returns a locked list, and also returns
 * the lock pointer.

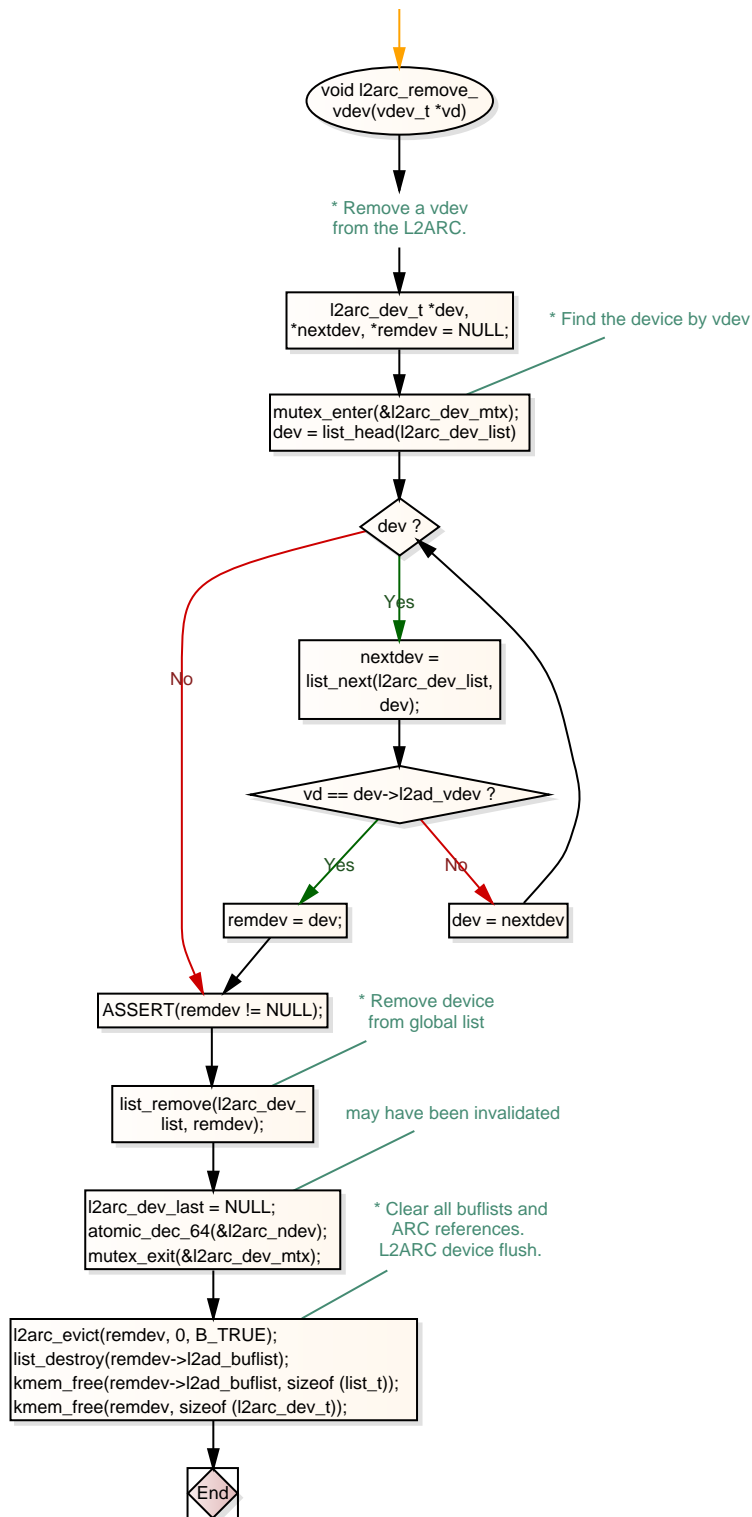


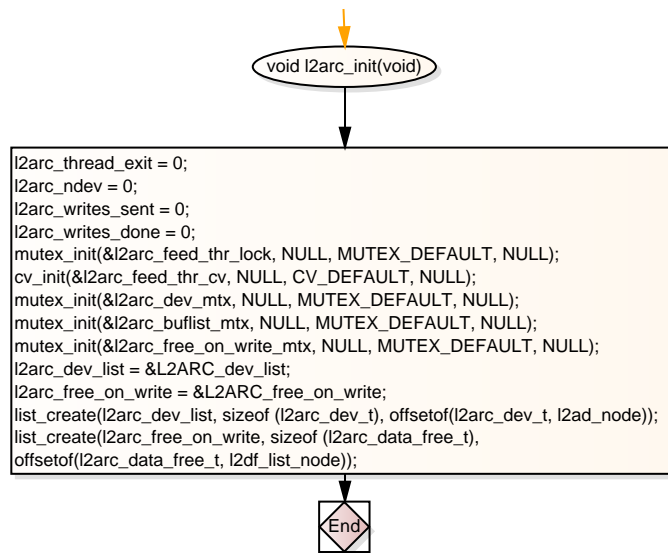














* This is called from `dmu_fini()`, which is called from `spa_fini()`;
* Because of this, we can assume that all `l2arc` devices have
* already been removed when the pools themselves were removed.

```
l2arc_do_free_on_write();
mutex_destroy(&l2arc_feed_thr_lock);
cv_destroy(&l2arc_feed_thr_cv);
mutex_destroy(&l2arc_dev_mtx);
mutex_destroy(&l2arc_buflist_mtx);
mutex_destroy(&l2arc_free_on_write_mtx);
list_destroy(l2arc_dev_list);
list_destroy(l2arc_free_on_write);
```



