

\* CDDL HEADER START  
 \*  
 \* The contents of this file are subject to the terms of the  
 \* Common Development and Distribution License (the "License").  
 \* You may not use this file except in compliance with the License.  
 \*  
 \* You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE  
 \* or http://www.opensolaris.org/os/licensing.  
 \* See the License for the specific language governing permissions  
 \* and limitations under the License.  
 \*  
 \* When distributing Covered Code, include this CDDL HEADER in each  
 \* file and include the License file at usr/src/OPENSOLARIS.LICENSE.  
 \* If applicable, add the following below this CDDL HEADER, with the  
 \* fields enclosed by brackets "[]" replaced with your own identifying  
 \* information: Portions Copyright [yyyy] [name of copyright owner]  
 \*  
 \* CDDL HEADER END  
 \* Copyright (c) 2005, 2010, Oracle and/or its affiliates. All rights reserved.  
 \* Copyright (c) 2012 by Delphix. All rights reserved.  
 Portions Copyright 2010 Robert Milkowski

```
#include <sys/zfs_context.h>
#include <sys/spa.h>
#include <sys/dmu.h>
#include <sys/zap.h>
#include <sys/arc.h>
#include <sys/stat.h>
#include <sys/resource.h>
#include <sys/zil.h>
#include <sys/zil_impl.h>
#include <sys/dsl_dataset.h>
#include <sys/vdev_impl.h>
#include <sys/dmu_tx.h>
```

\* The zfs intent log (ZIL) saves transaction records of system calls  
 \* that change the file system in memory with enough information  
 \* to be able to replay them. These are stored in memory until  
 \* either the DMU transaction group (txg) commits them to the stable pool  
 \* and they can be discarded, or they are flushed to the stable log  
 \* (also in the pool) due to a fsync, O\_DSYNC or other synchronous  
 \* requirement. In the event of a panic or power fail then those log  
 \* records (transactions) are replayed.  
 \*

\* There is one ZIL per file system. Its on-disk (pool) format consists  
 \* of 3 parts:  
 \*  
 \* - ZIL header  
 \* - ZIL blocks  
 \* - ZIL records  
 \*

\* A log record holds a system call transaction. Log blocks can  
 \* hold many log records and the blocks are chained together.  
 \* Each ZIL block contains a block pointer (blkptr\_t) to the next  
 \* ZIL block in the chain. The ZIL header points to the first  
 \* block in the chain. Note there is not a fixed place in the pool  
 \* to hold blocks. They are dynamically allocated and freed as  
 \* needed from the blocks available. Figure X shows the ZIL structure:

```
#include <sys/dsl_pool.h>
```

\* This global ZIL switch affects all pools  
 disable intent logging replay

```
int zil_replay_disable  
= 0;
```

\* Tunable parameter for debugging or performance analysis. Setting  
 \* zfs\_nocacheflush will cause corruption on power loss if a volatile  
 \* out-of-order write cache is enabled.

```
boolean_t zfs_nocacheflush = B_FALSE;  
static kmem_cache_t *zil_lwb_cache;  
static void zil_async_to_sync(zilog_t *zilog, uint64_t foid);
```

\* ziltest is by and large an ugly hack, but very useful in  
 \* checking replay without tedious work.  
 \* When running ziltest we want to keep all itx's and so maintain  
 \* a single list in the zil\_itxg[] that uses a high txg: ZILTEST\_TXG  
 \* We subtract TXG\_CONCURRENT\_STATES to allow for common code.

```
#define LWB_EMPTY(lwb) ((BP_GET_LSIZE(&lwb->lwb_blk) -  
sizeof(zil_chain_t)) == (lwb->lwb_sz - lwb->lwb_nused))  
#define ZILTEST_TXG (UINT64_MAX - TXG_CONCURRENT_STATES)
```

\* Define a limited set of intent log block sizes.  
 \* These must be a multiple of 4KB. Note only the amount used (again  
 \* aligned to 4KB) actually gets written. However, we can't always just  
 \* allocate SPA\_MAXBLOCKSIZE as the slog space could be exhausted.  
 non TX\_WRITE  
 data base  
 NFS writes

```
uint64_t  
zil_block_buckets[] = {  
4096, 8192+4096, 32*1024  
+ 4096, UINT64_MAX};
```

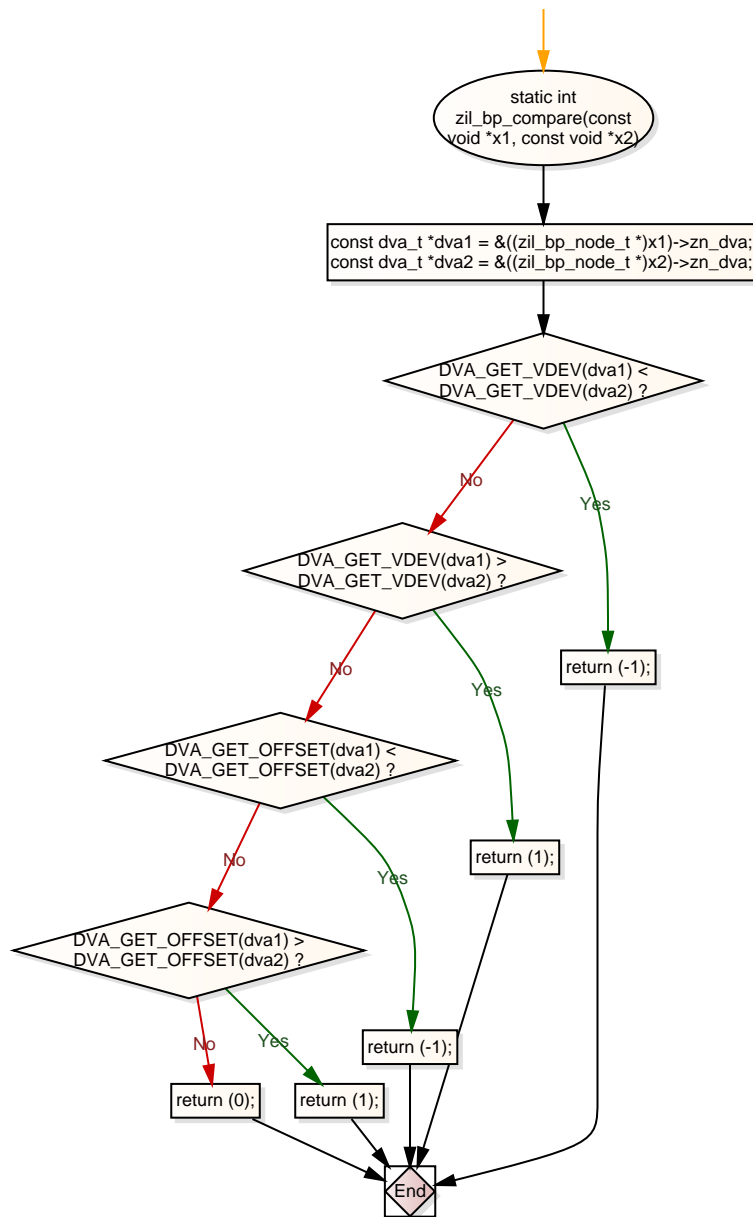
\* Use the slog as long as the logbias is 'latency' and the current commit size  
 \* is less than the limit or the total list size is less than 2X the limit.  
 \* Limit checking is disabled by setting zil\_slog\_limit to UINT64\_MAX.

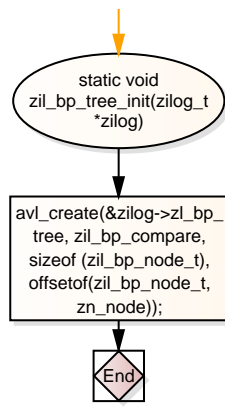
```
uint64_t zil_slog_limit  
= 1024 * 1024;
```

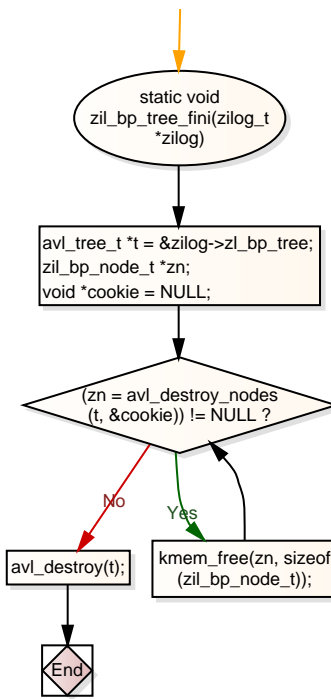
\* Start a log block write and advance to the next log block.  
 \* Calls are serialized.

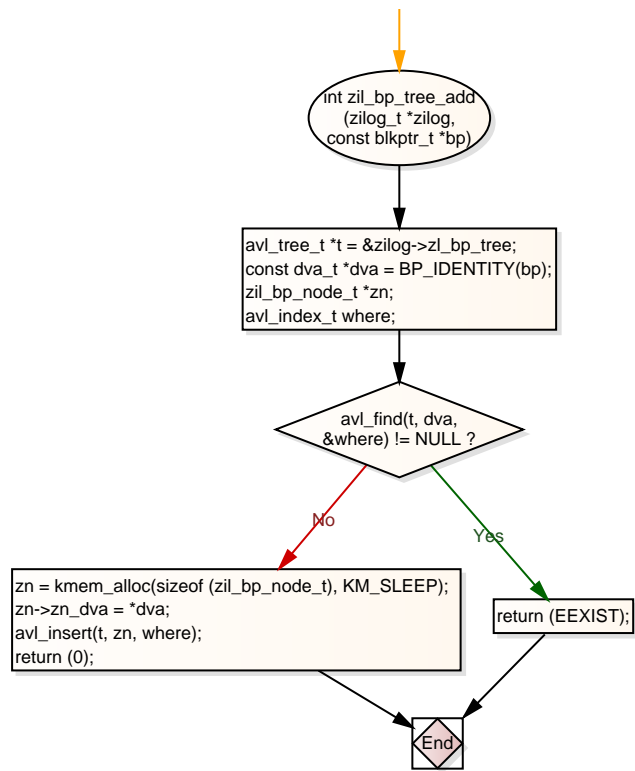
```
#define USE_SLOG(zilog)  
(((zilog->zil_logbias ==  
ZFS_LOGBIAS_LATENCY) &&  
((zilog->zil_cur_used <  
zil_slog_limit) ||  
(zilog->zil_itx_list_sz  
< (zil_slog_limit  
<< 1))))
```

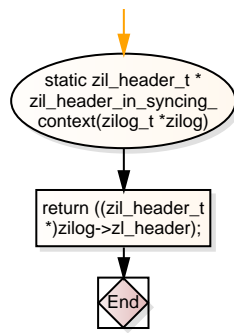


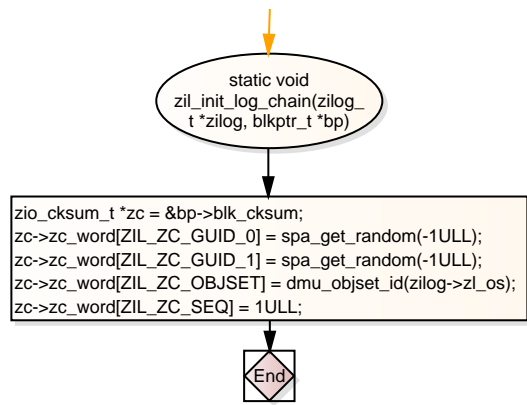






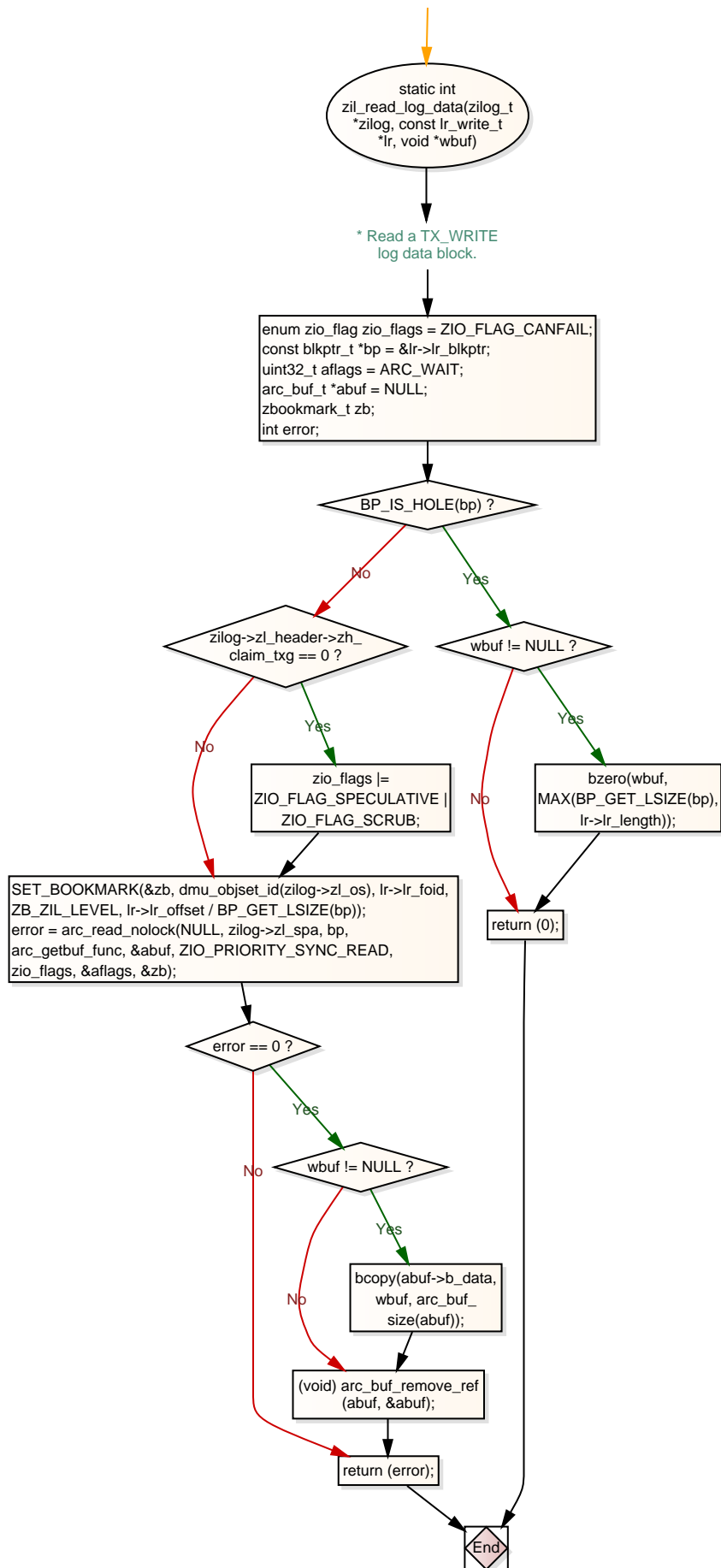


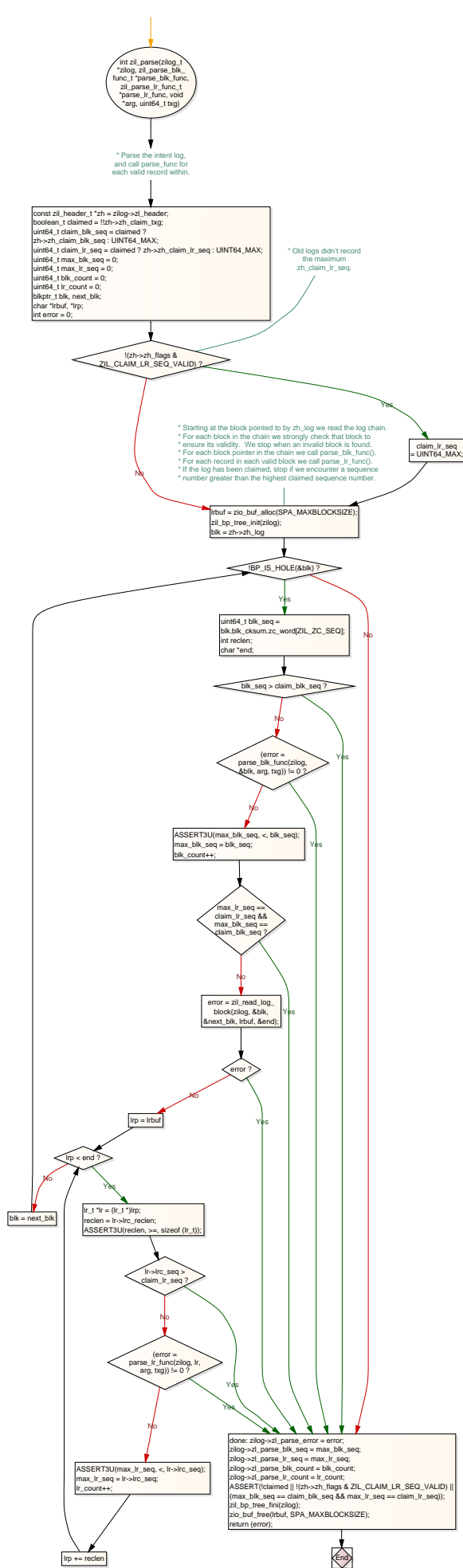


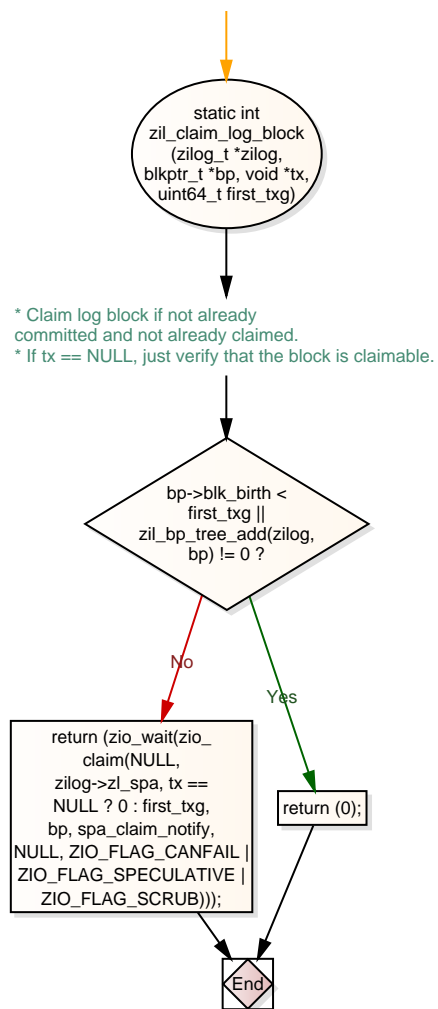


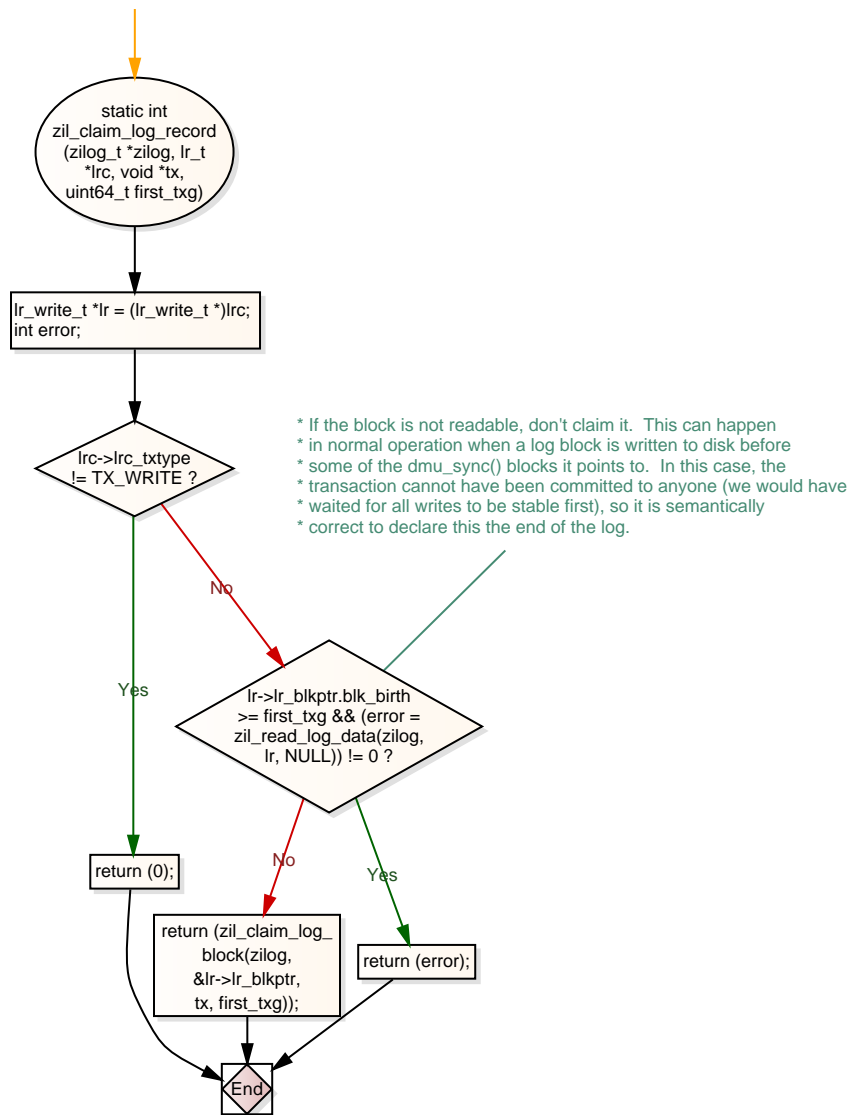


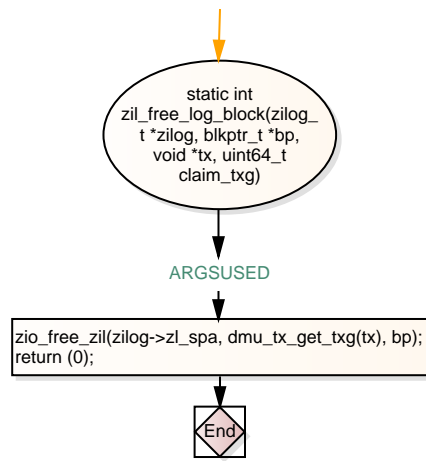


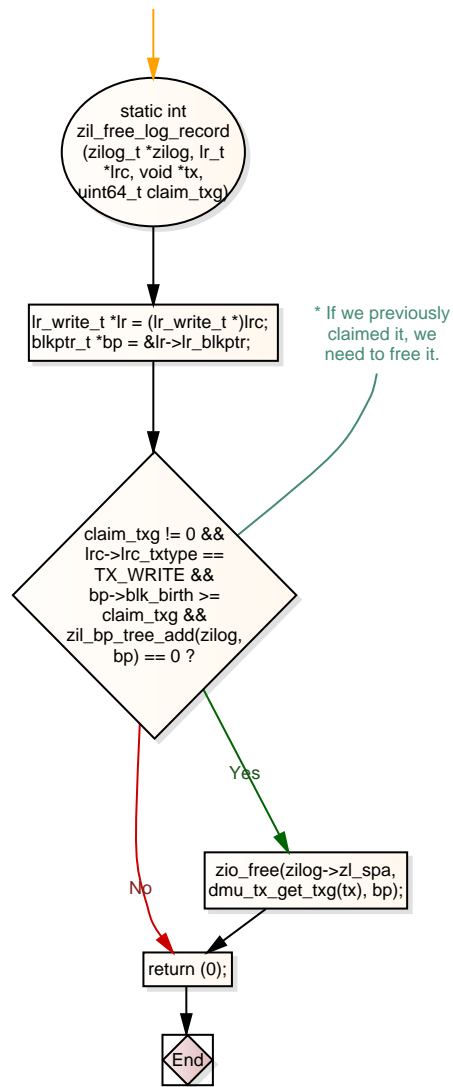


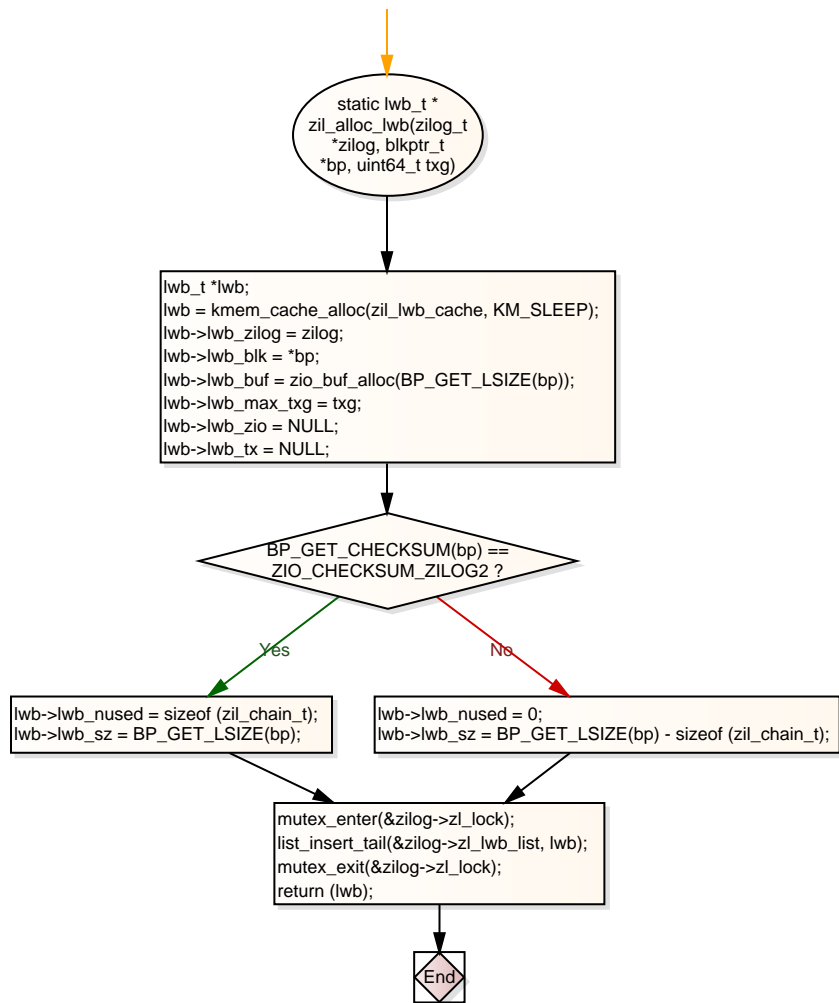


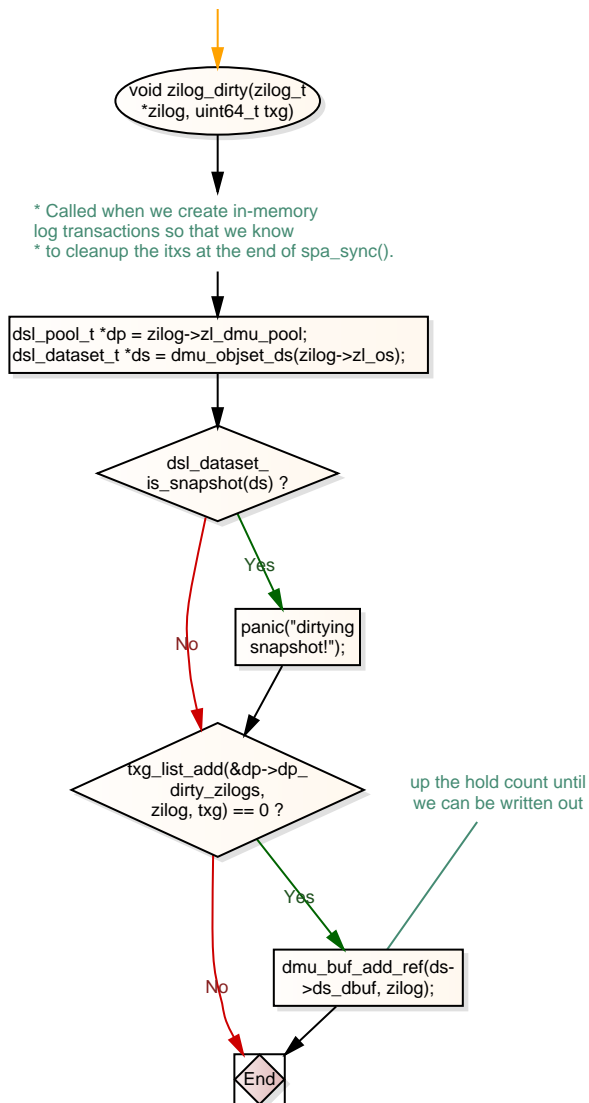




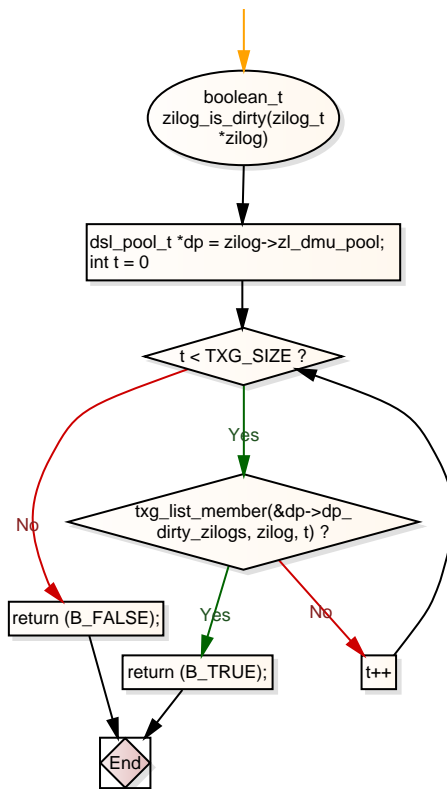


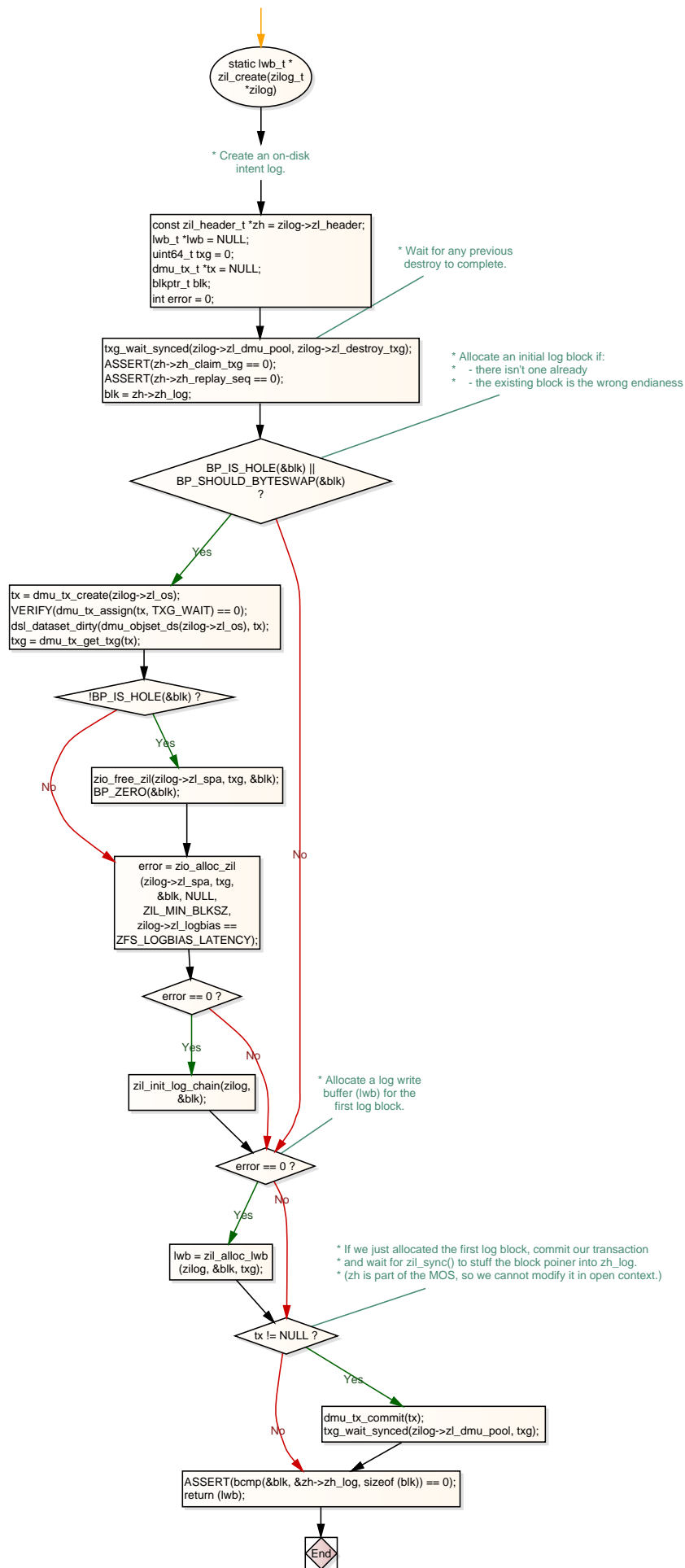




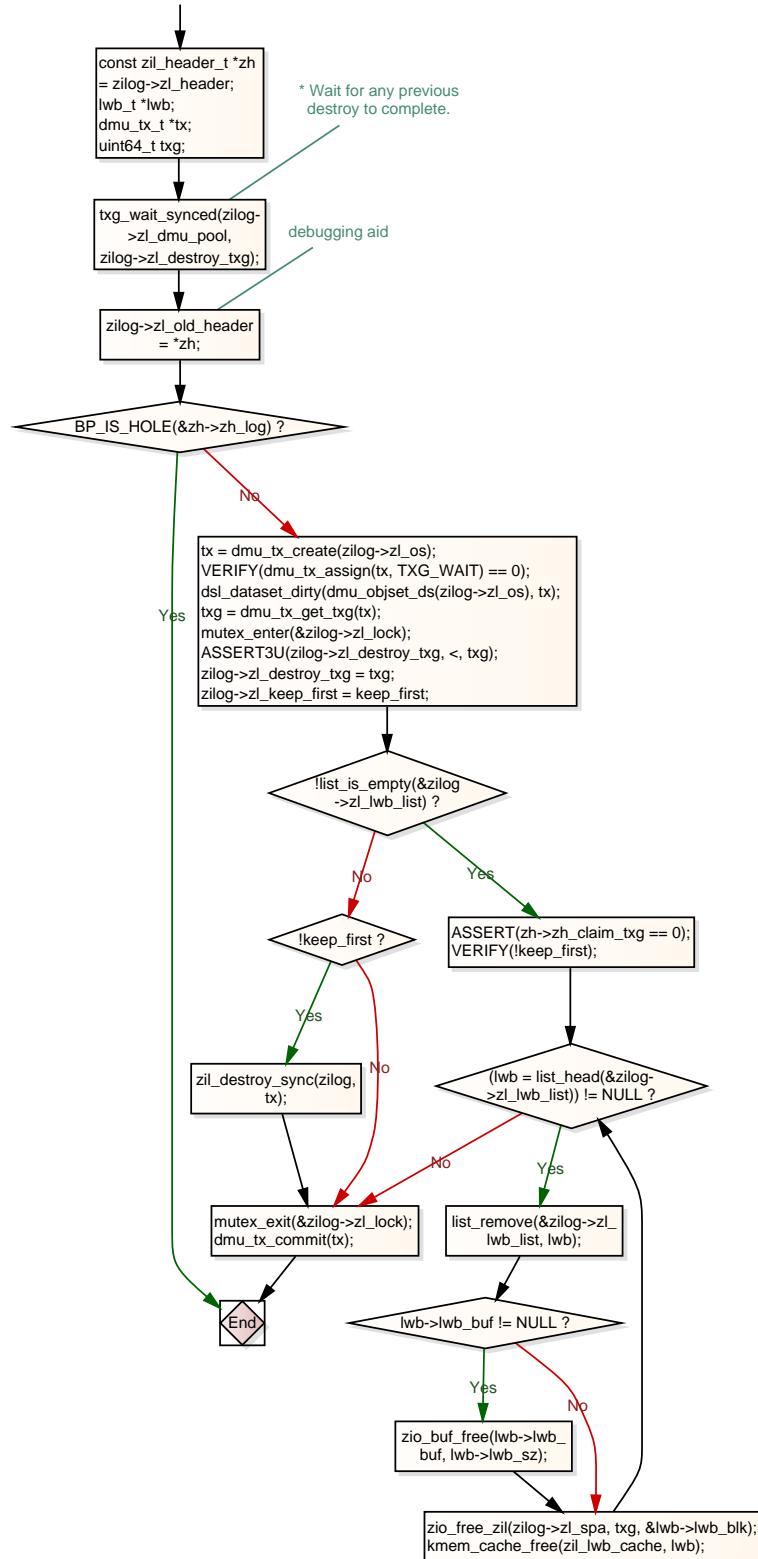


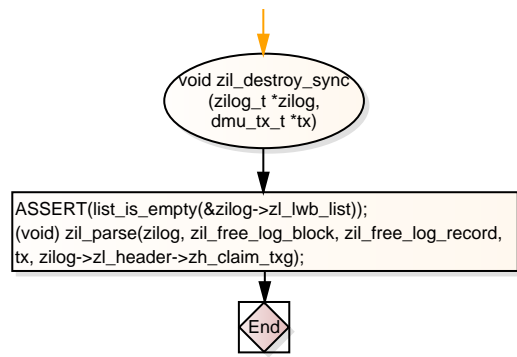


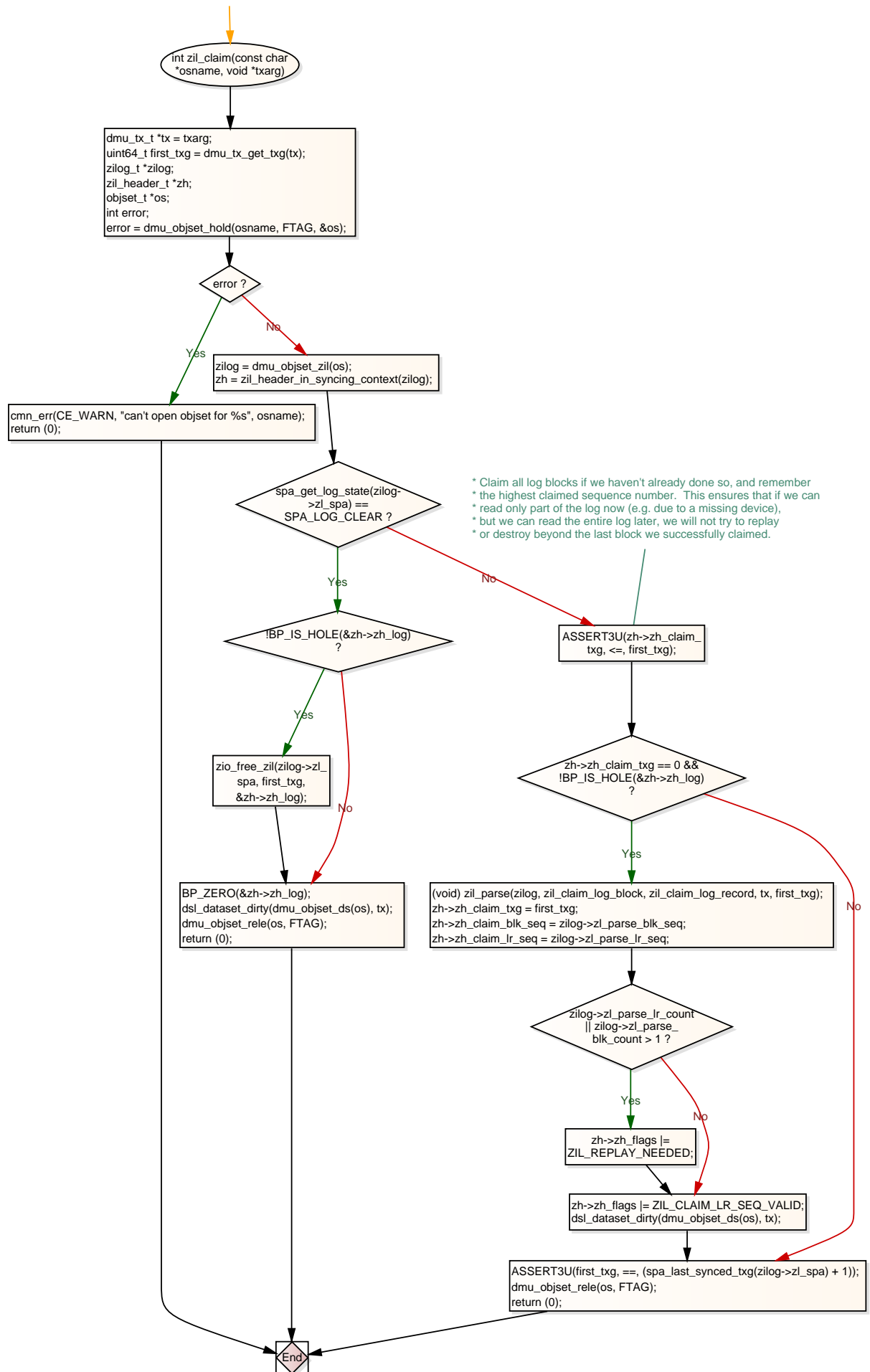


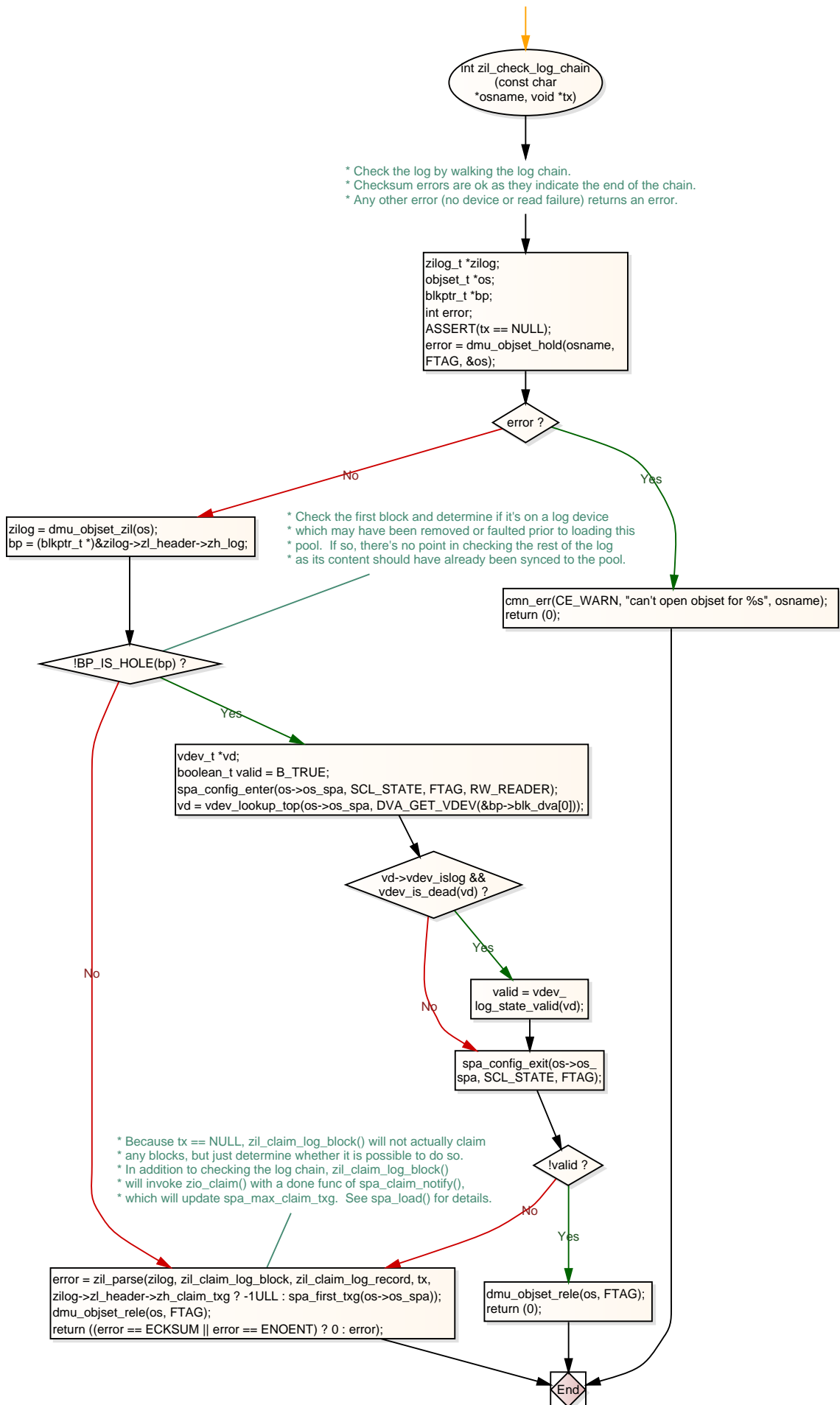


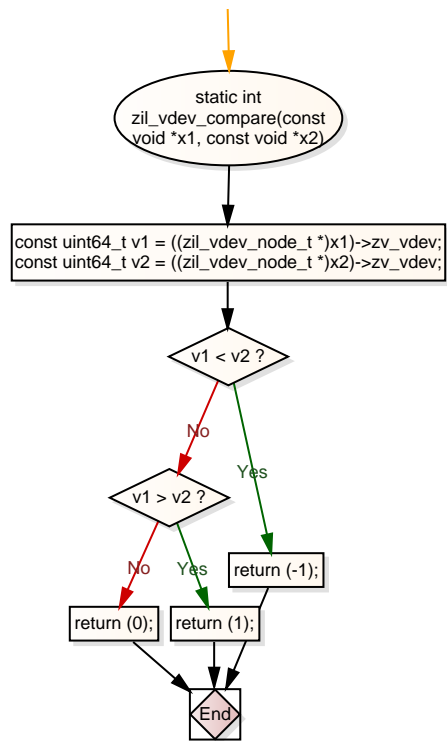
\* In one tx, free all log blocks and clear the log header.  
 \* If keep\_first is set, then we're replaying a log with no content.  
 \* We want to keep the first block, however, so that the first  
 \* synchronous transaction doesn't require a txg\_wait\_synced()  
 \* in zil\_create(). We don't need to txg\_wait\_synced() here either  
 \* when keep\_first is set, because both zil\_create() and zil\_destroy()  
 \* will wait for any in-progress destroys to complete.

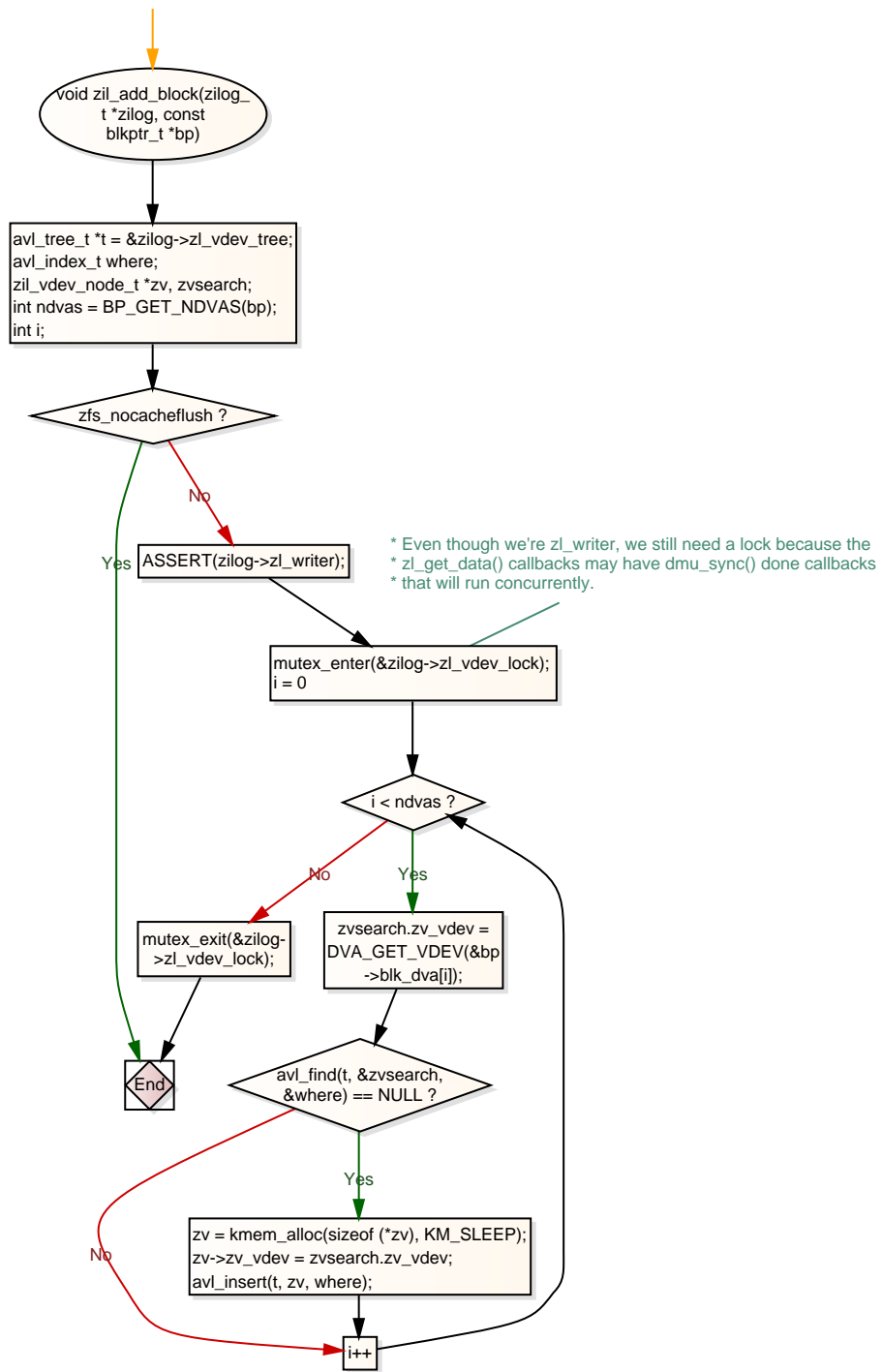




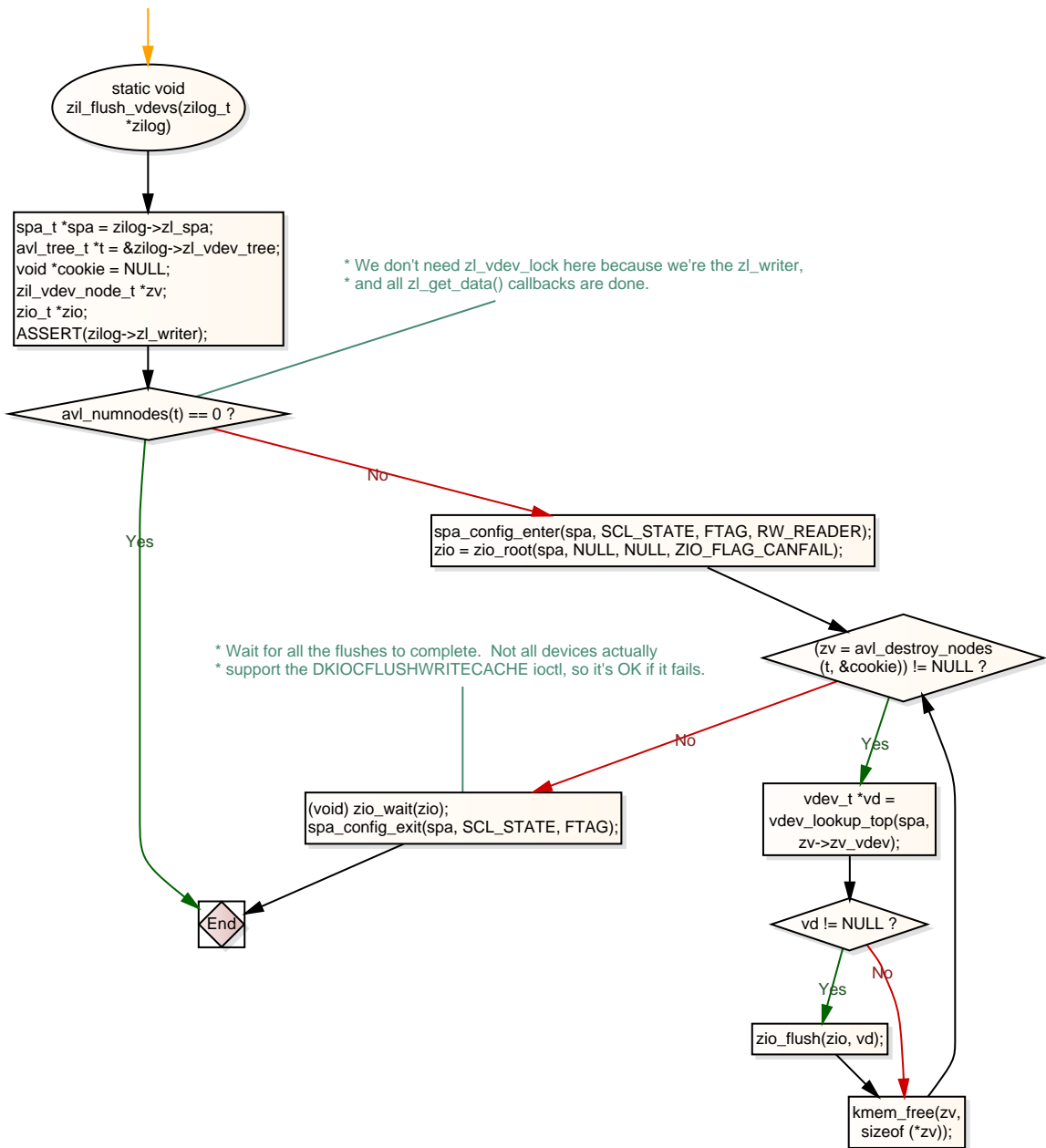


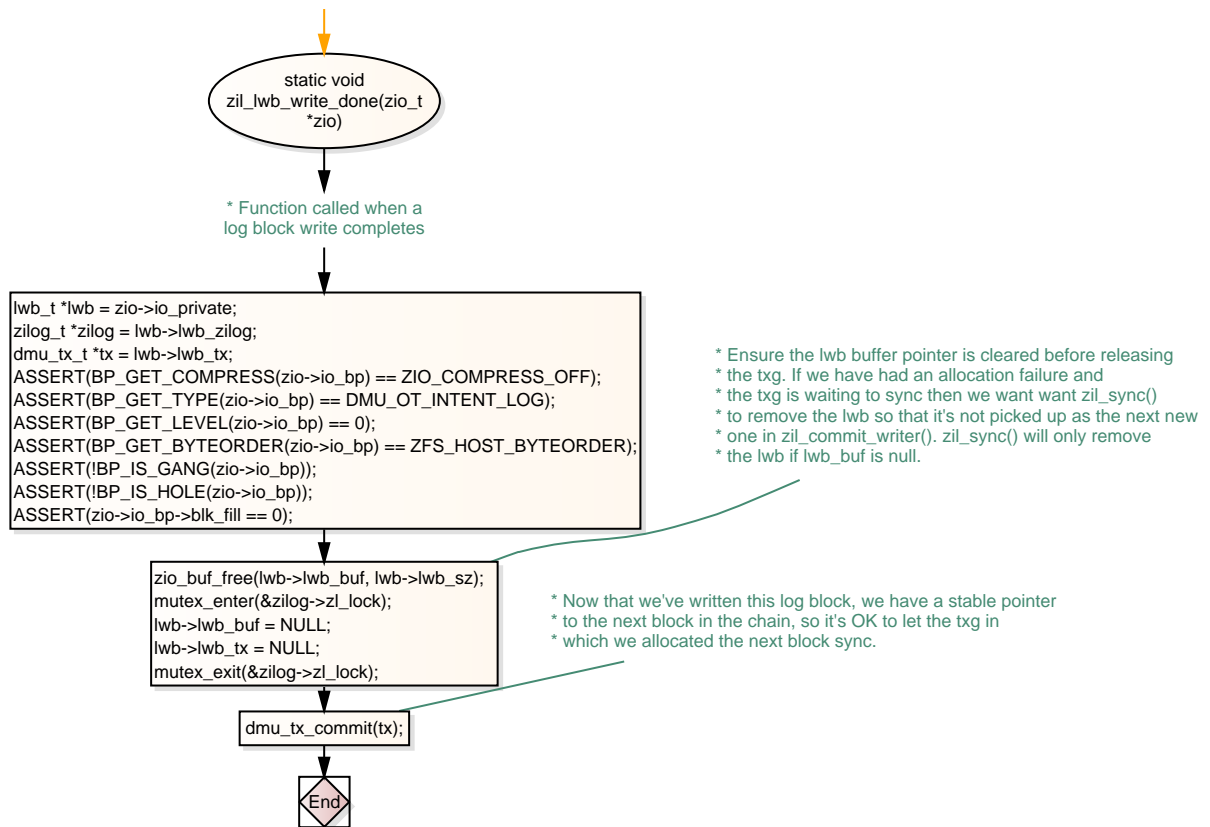


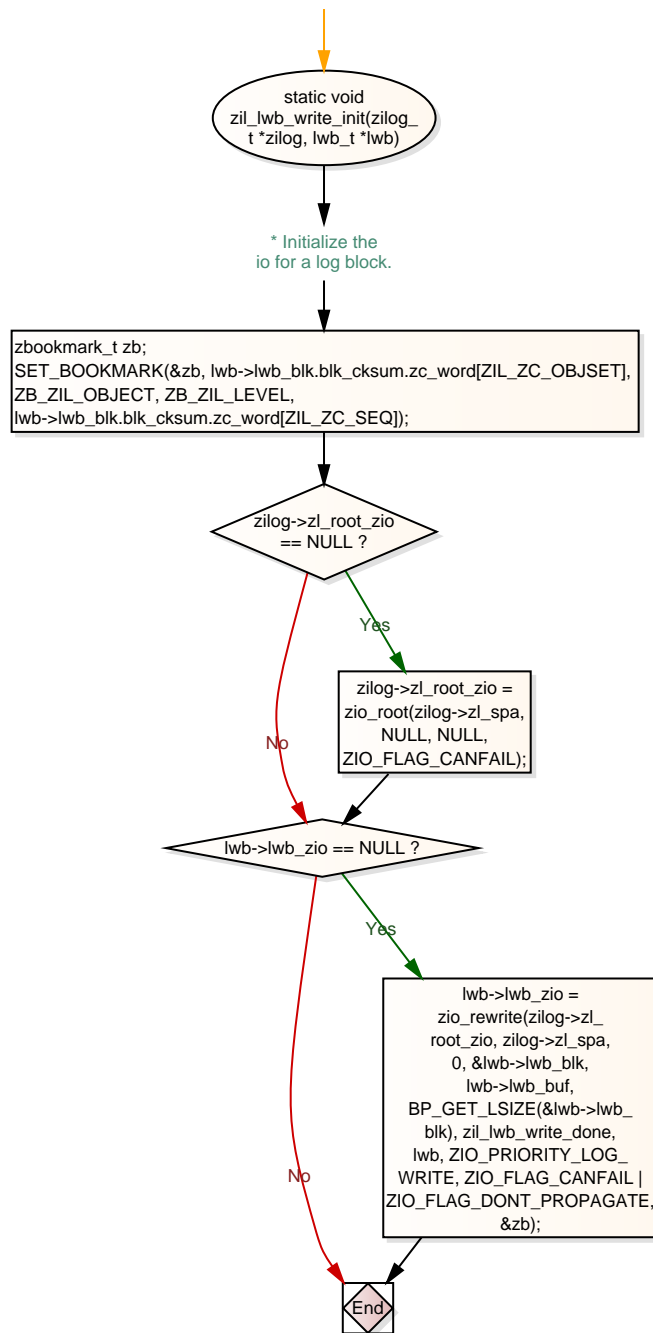


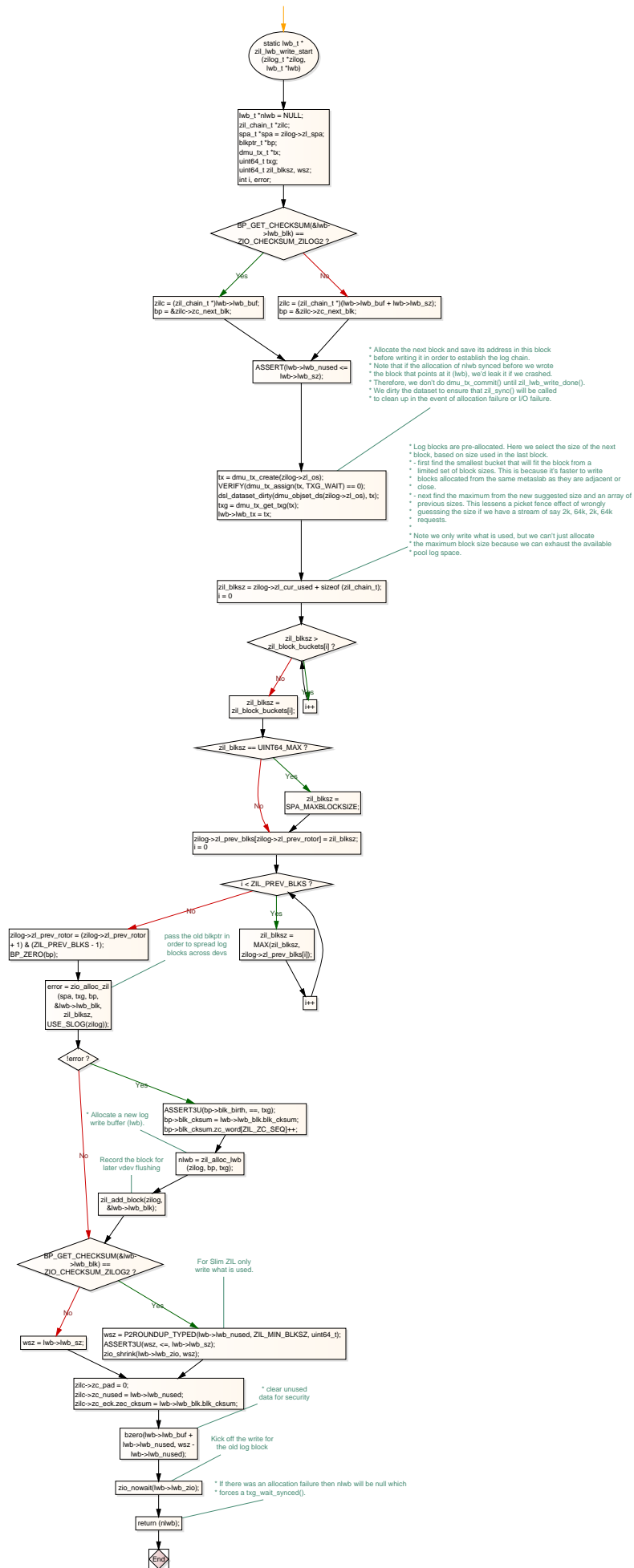




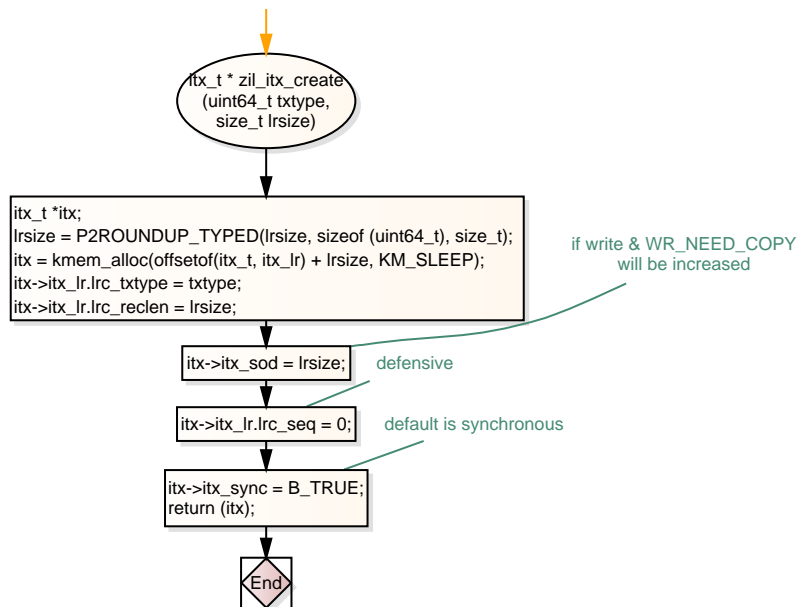


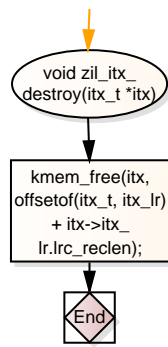


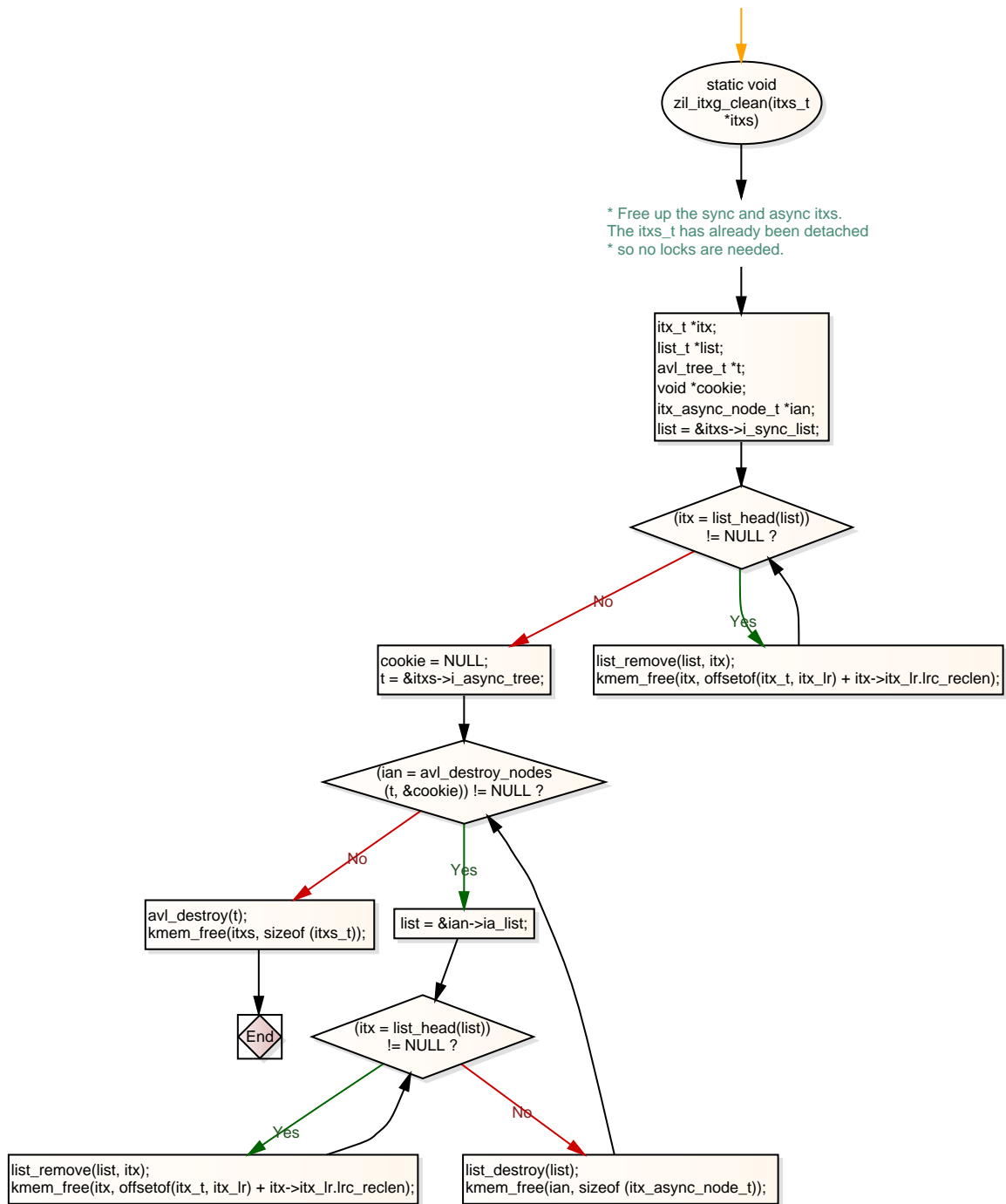




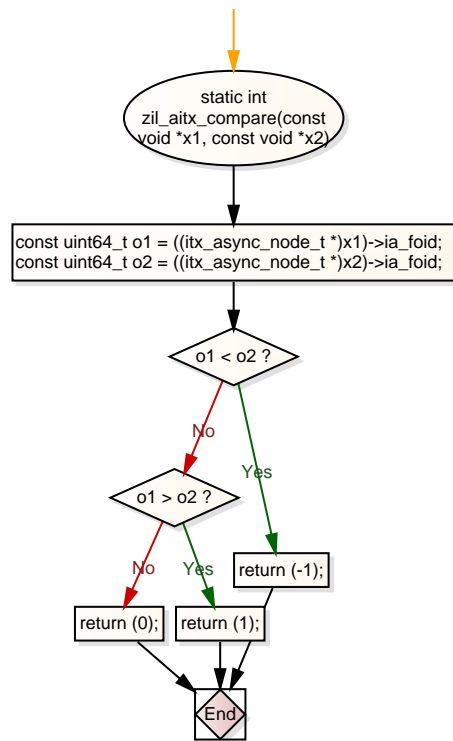


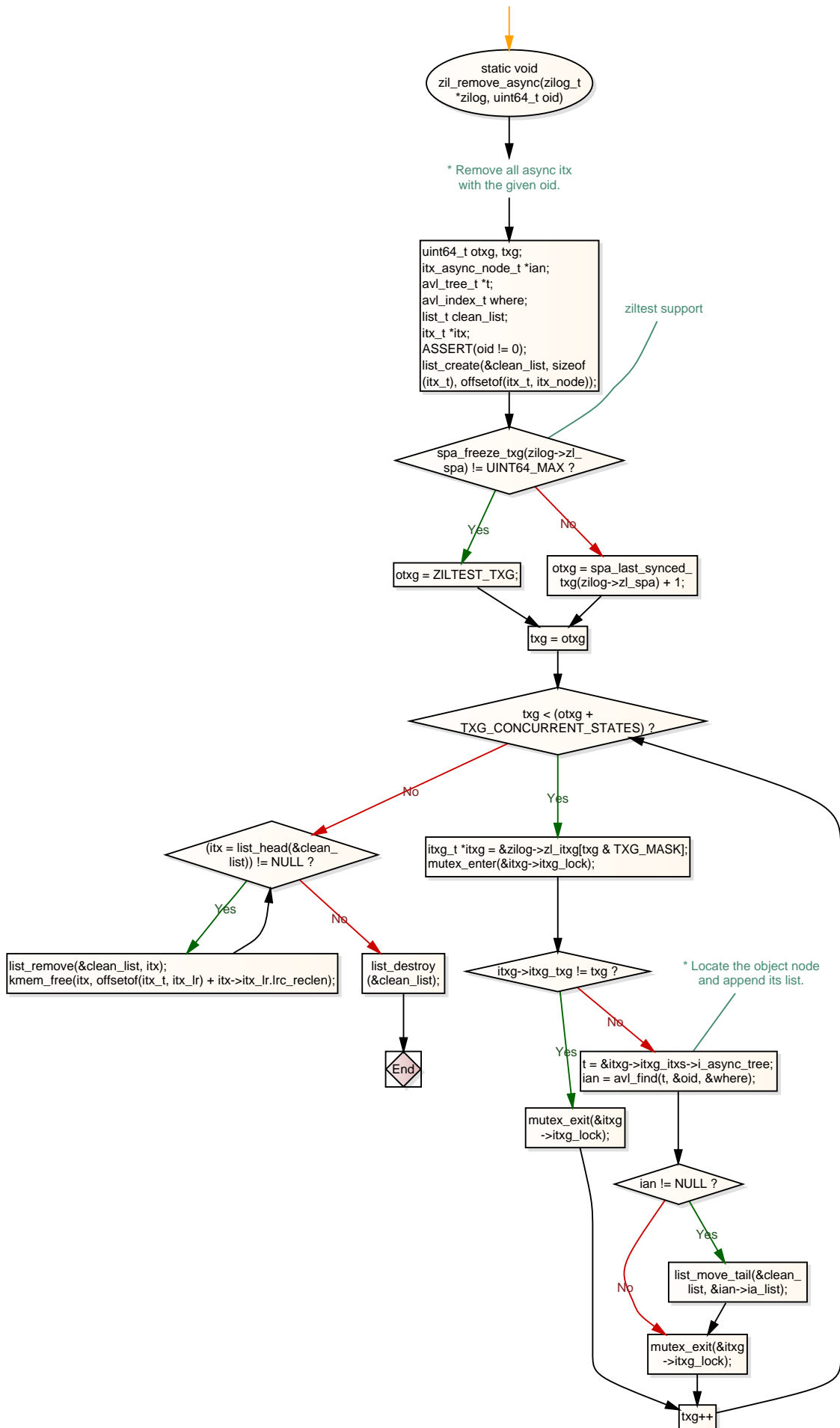


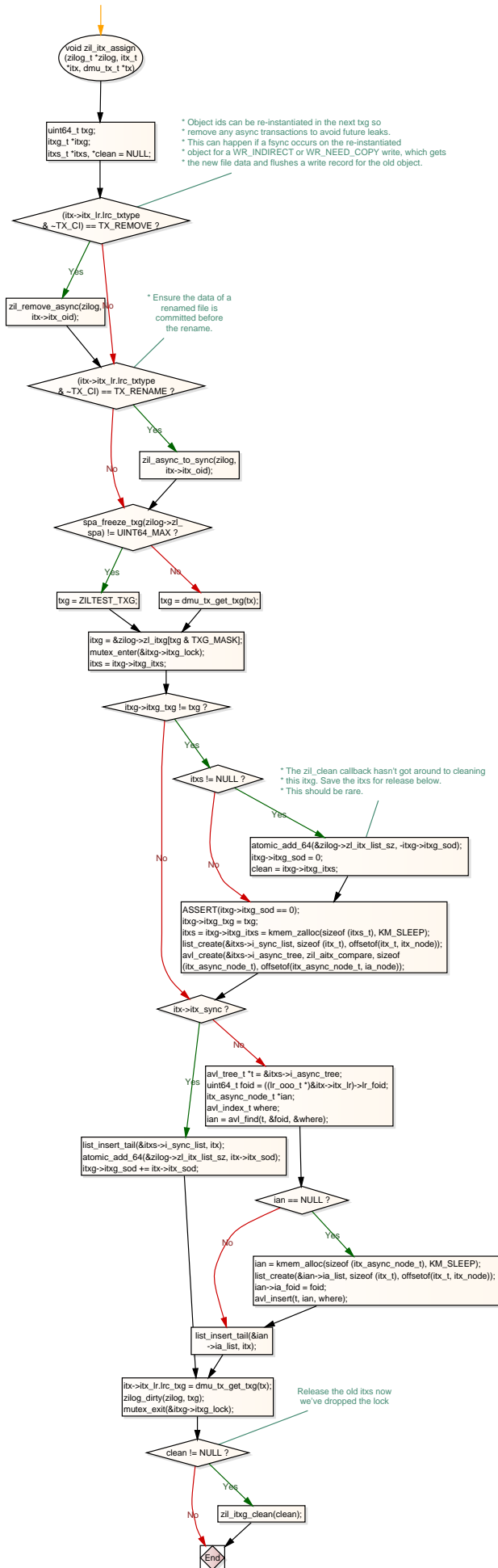


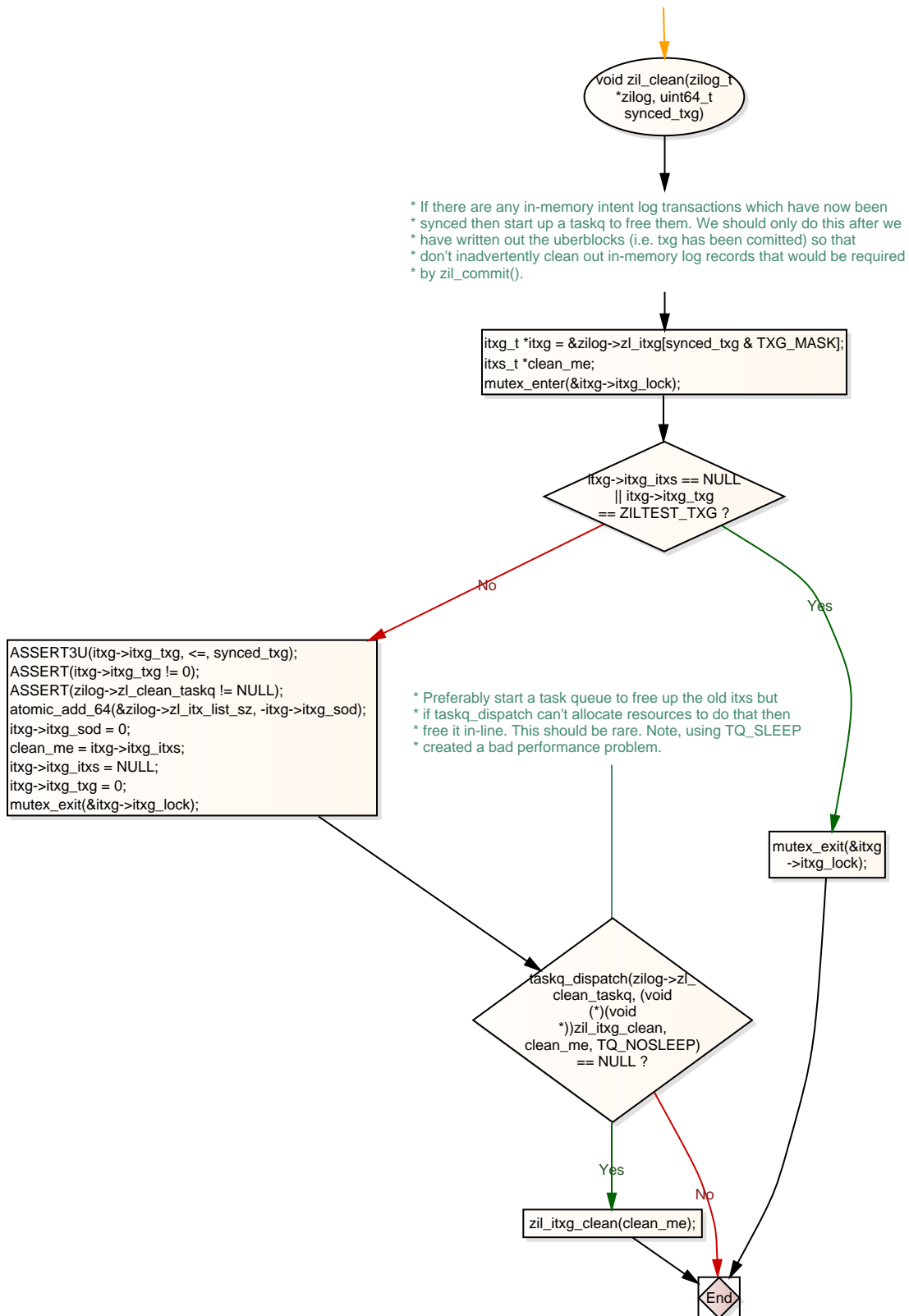


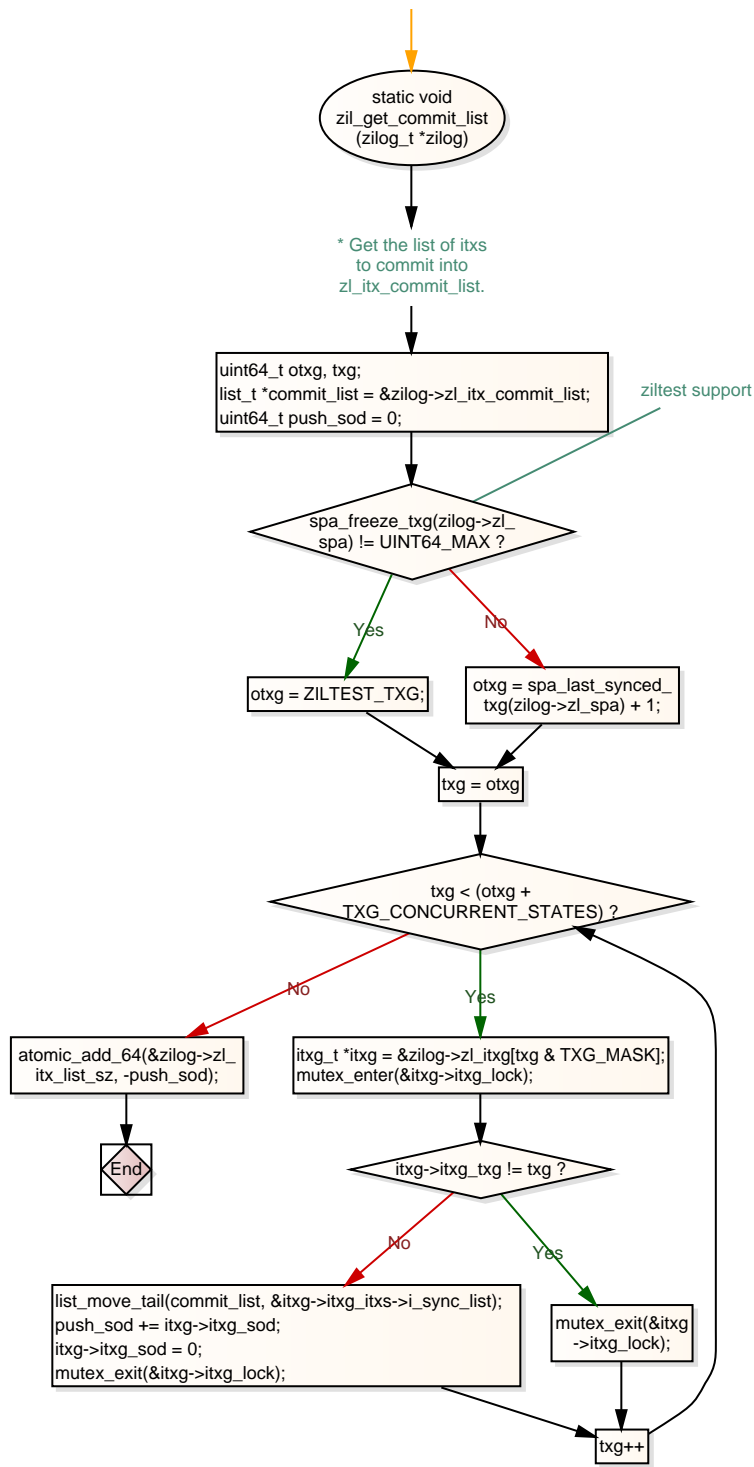


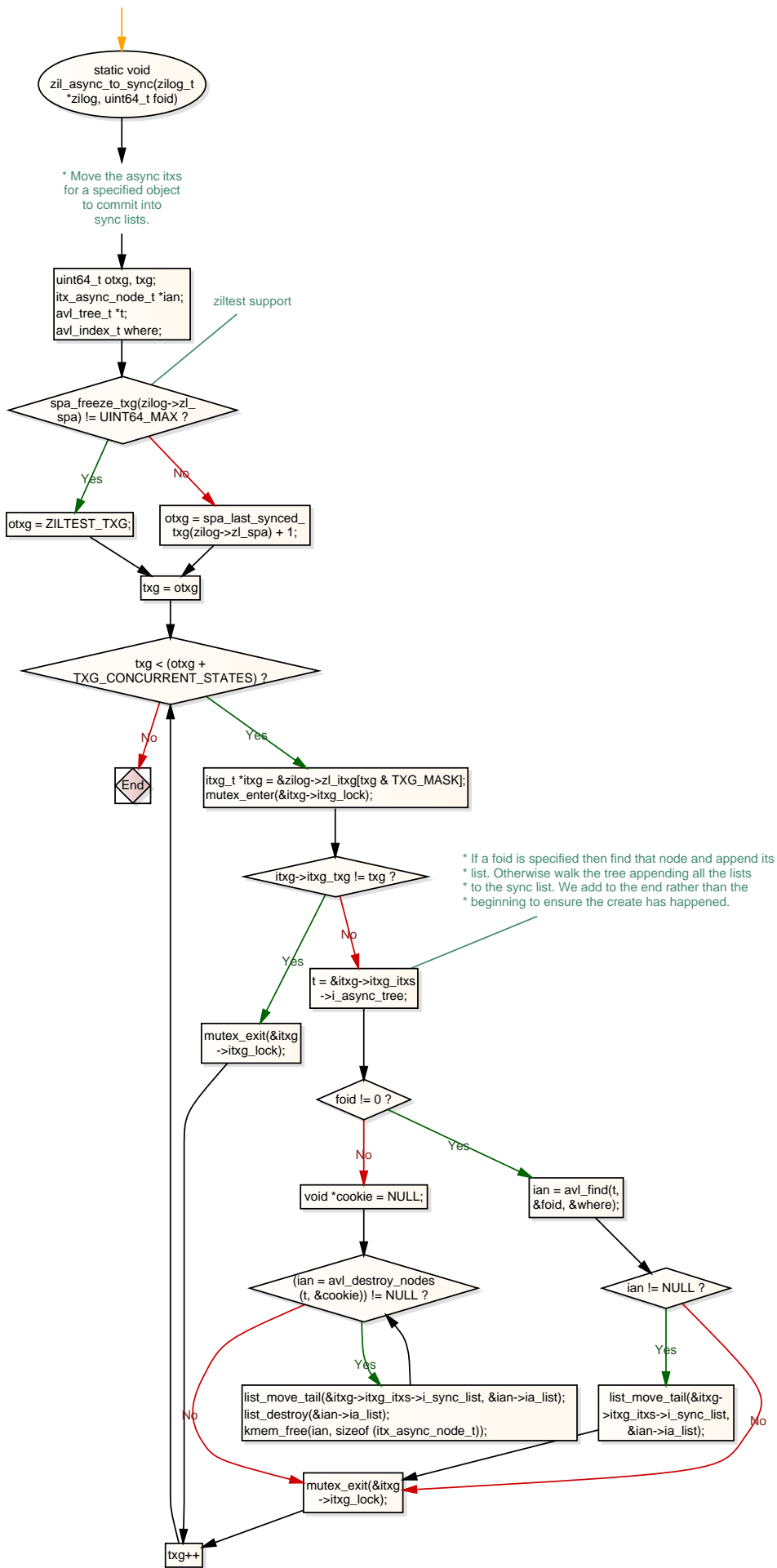


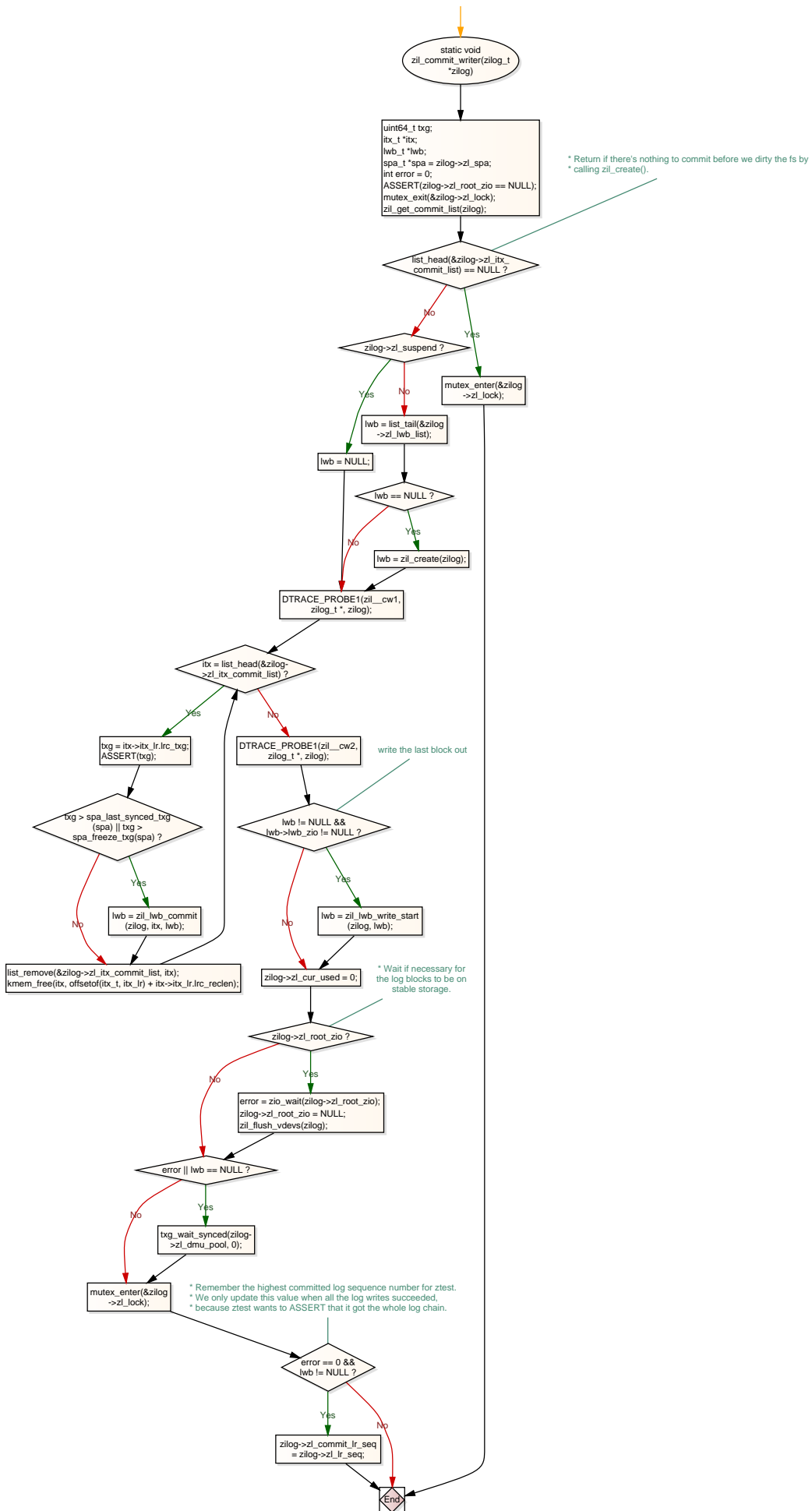


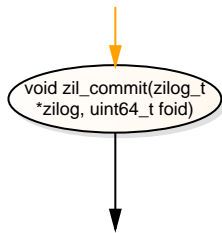












- \* Commit zfs transactions to stable storage.
- \* If foid is 0 push out all transactions, otherwise push only those
- \* for that object or might reference that object.
- \*
- \* itxs are committed in batches. In a heavily stressed zil there will be
- \* a commit writer thread who is writing out a bunch of itxs to the log
- \* for a set of committing threads (cthreads) in the same batch as the writer.
- \* Those cthreads are all waiting on the same cv for that batch.
- \*
- \* There will also be a different and growing batch of threads that are
- \* waiting to commit (qthreads). When the committing batch completes
- \* a transition occurs such that the cthreads exit and the qthreads become
- \* cthreads. One of the new cthreads becomes the writer thread for the
- \* batch. Any new threads arriving become new qthreads.
- \*
- \* Only 2 condition variables are needed and there's no transition
- \* between the two cvs needed. They just flip-flop between qthreads
- \* and cthreads.
- \*
- \* Using this scheme we can efficiently wakeup up only those threads
- \* that have been committed.

