

* 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/ddev_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;
```

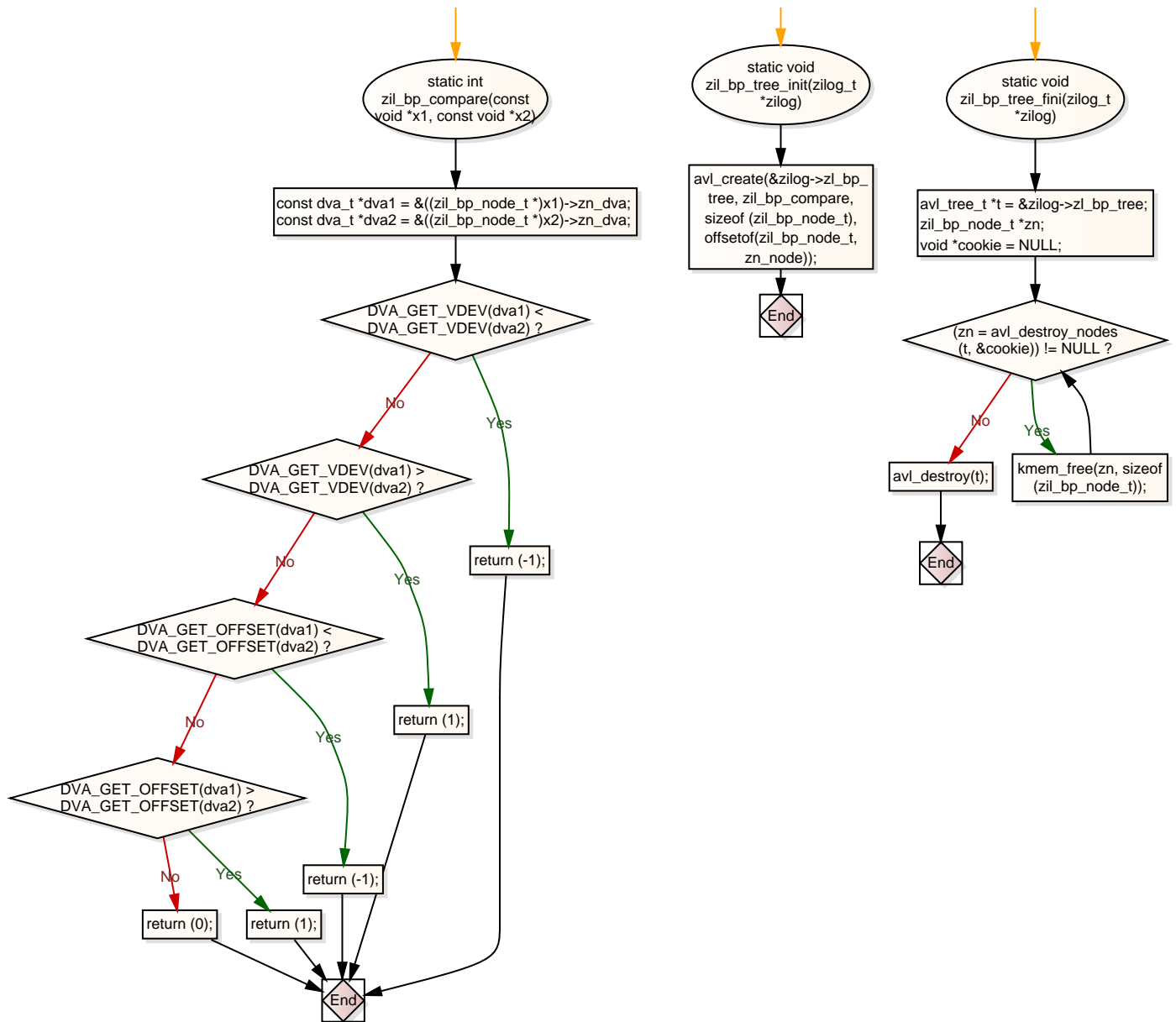
* Start a log block write and advance to the next log block.
* Calls are serialized.

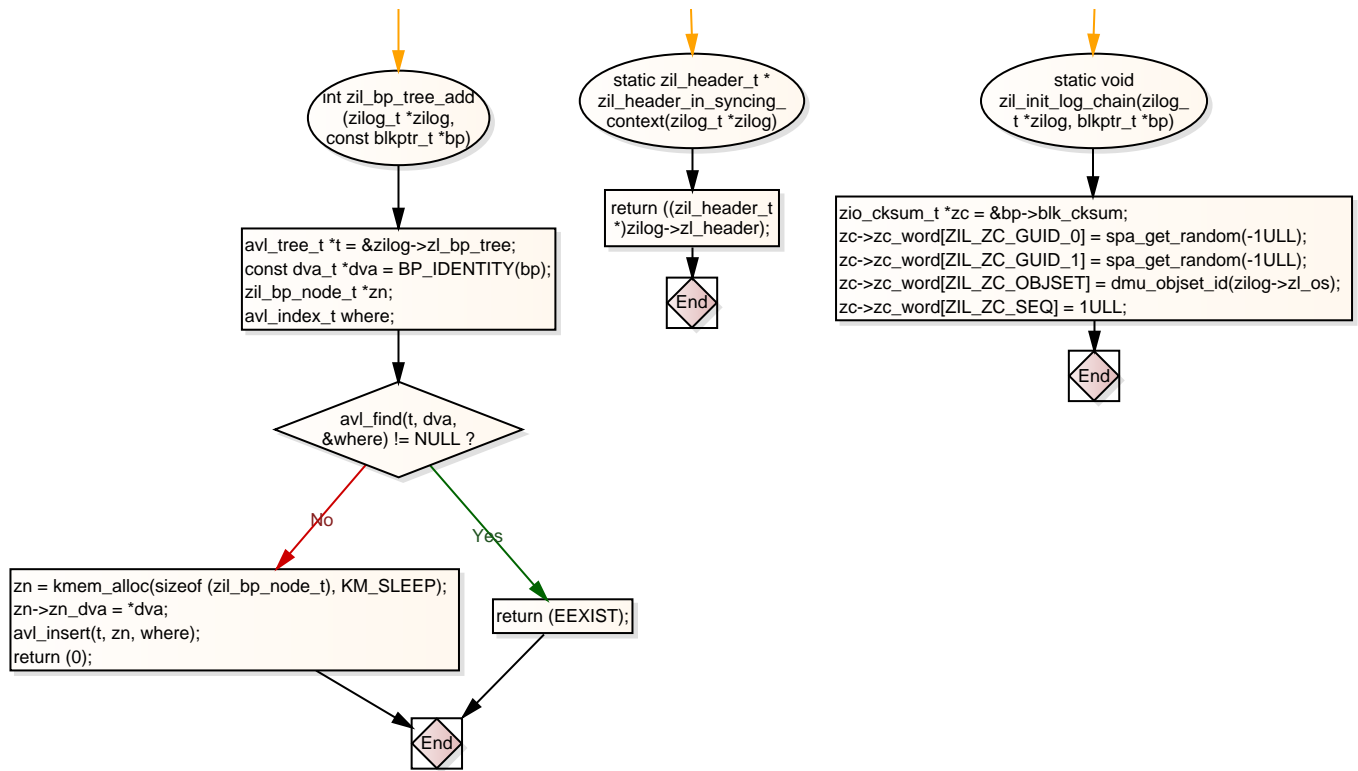
↓

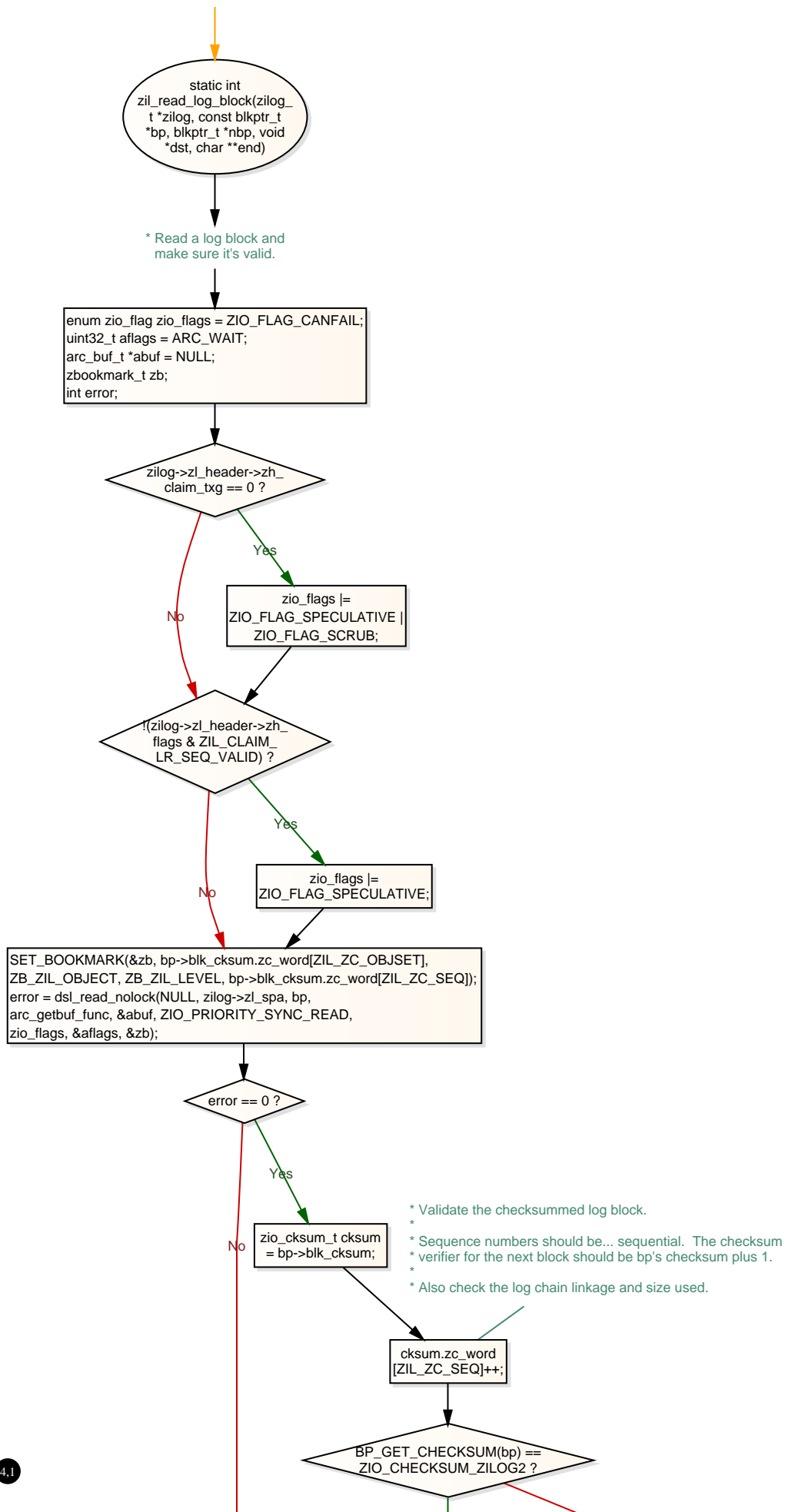
```
#define USE_SLOG(zilog)
(((zilog)->zl_logbias ==
ZFS_LOGBIAS_LATENCY) &&
(((zilog)->zl_cur_used <
zil_slog_limit) ||
((zilog)->zl_itx_list_sz
< (zil_slog_limit
<< 1))))
```

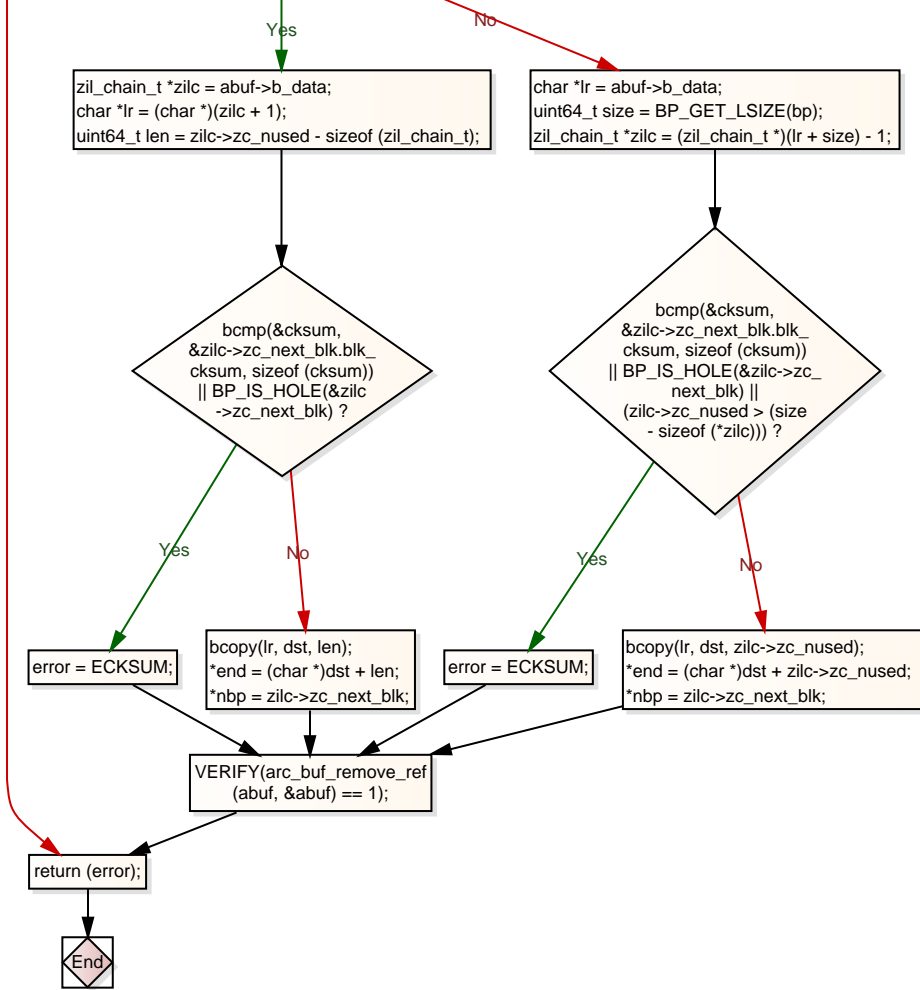
↙

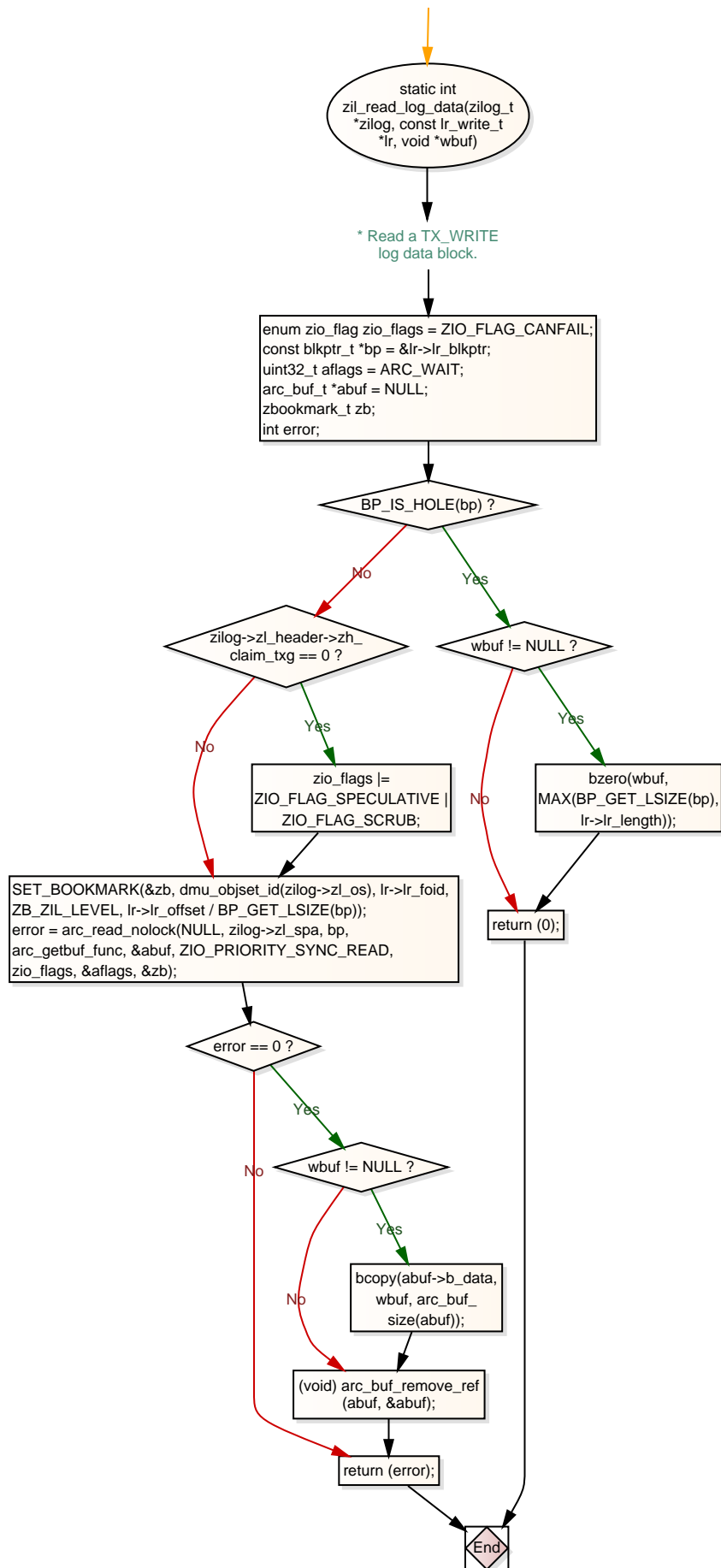


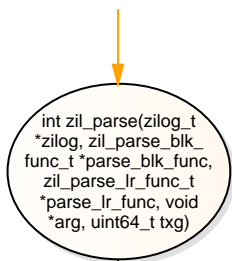








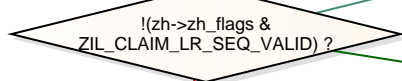




* Parse the intent log,
and call parse_func for
each valid record within.

```

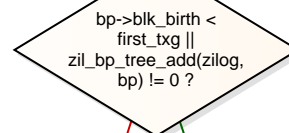
const zil_header_t *zh = zilog->zl_header;
boolean_t claimed = !zh->zh_claim_txg;
uint64_t claim_blk_seq = claimed ?
zh->zh_claim_blk_seq : UINT64_MAX;
uint64_t claim_lr_seq = claimed ? zh->zh_claim_lr_seq : UINT64_MAX;
uint64_t max_blk_seq = 0;
uint64_t max_lr_seq = 0;
uint64_t blk_count = 0;
uint64_t lr_count = 0;
blkptr_t blk, next_blk;
char *lrbuf, *lrp;
int error = 0;
  
```



* Old logs didn't record
the maximum
zh_claim_lr_seq.

No

* Claim log block if not already
committed and not already claimed.
* If tx == NULL, just verify that the block is claimable.



No

Yes

```

return (zio_wait(zio_
claim(NULL,
zilog->zl_spa, tx ==
NULL ? 0 : first_txg,
bp, spa_claim_notify,
NULL, ZIO_FLAG_CANFAIL |
ZIO_FLAG_SPECULATIVE |
ZIO_FLAG_SCRUB));
  
```

return (0);



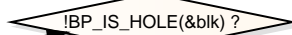
```

claim_lr_seq
= UINT64_MAX;
  
```

* Starting at the block pointed to by zh_log we read the log chain.
* For each block in the chain we strongly check that block to
ensure its validity. We stop when an invalid block is found.
* For each block pointer in the chain we call parse_blk_func().
* For each record in each valid block we call parse_lr_func().
* If the log has been claimed, stop if we encounter a sequence
number greater than the highest claimed sequence number.

```

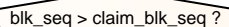
lrbuf = zio_buf_alloc(SPA_MAXBLOCKSIZE);
zil_bp_tree_init(zilog);
blk = zh->zh_log
  
```



Yes

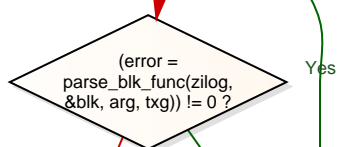
```

uint64_t blk_seq =
blk.blk_cksum.zc_word[ZIL_ZC_SEQ];
int reclen;
char *end;
  
```

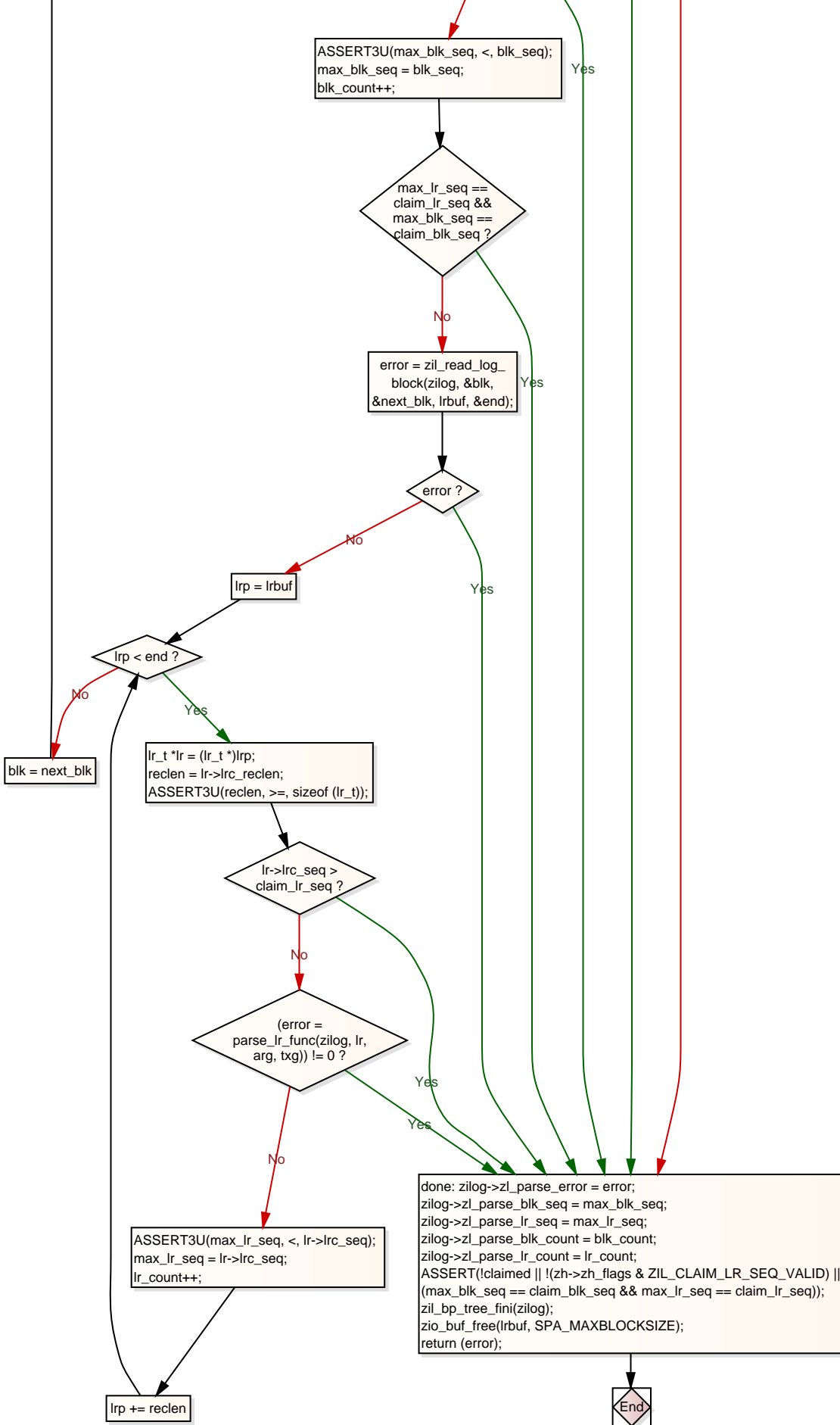


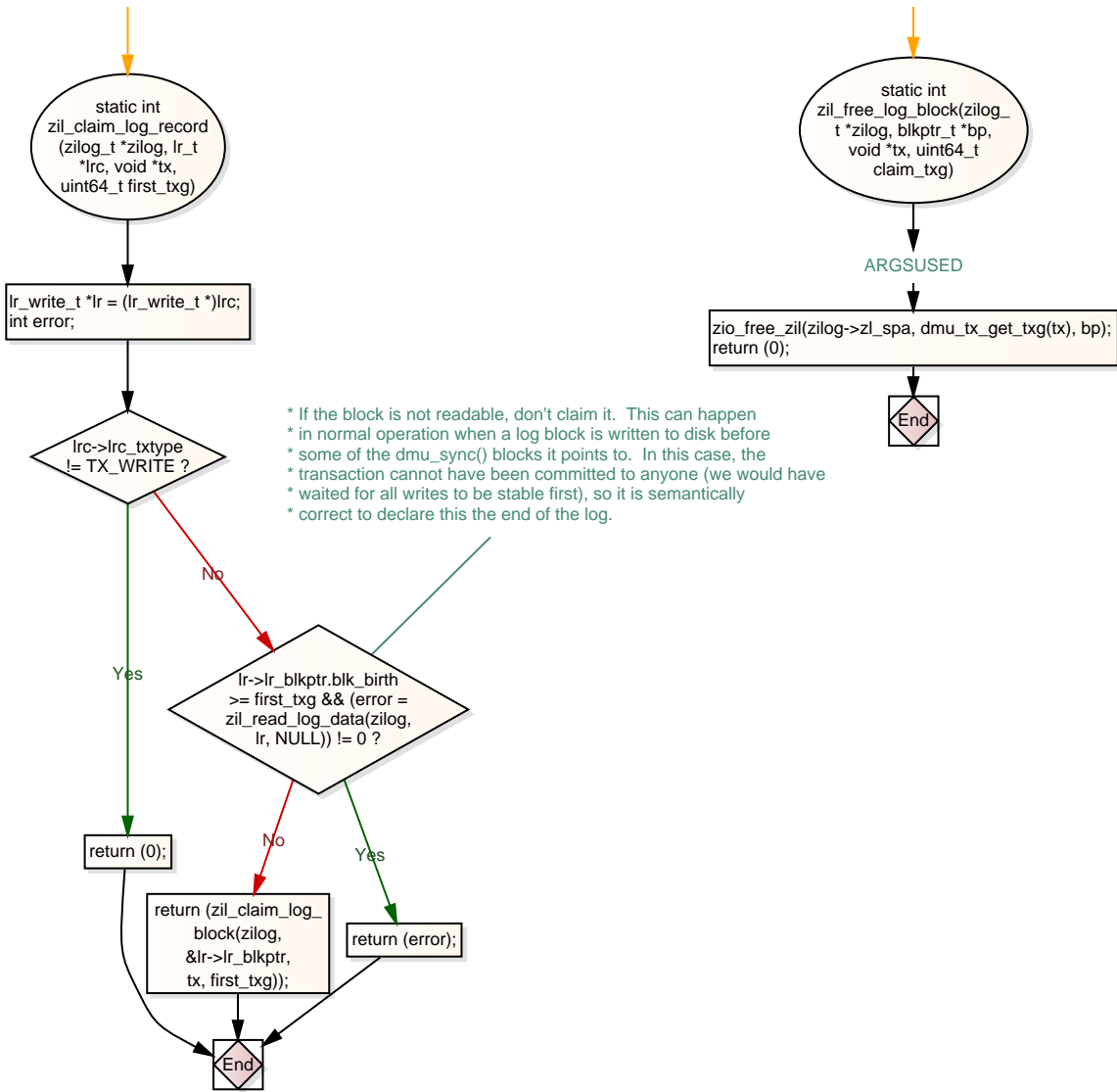
No

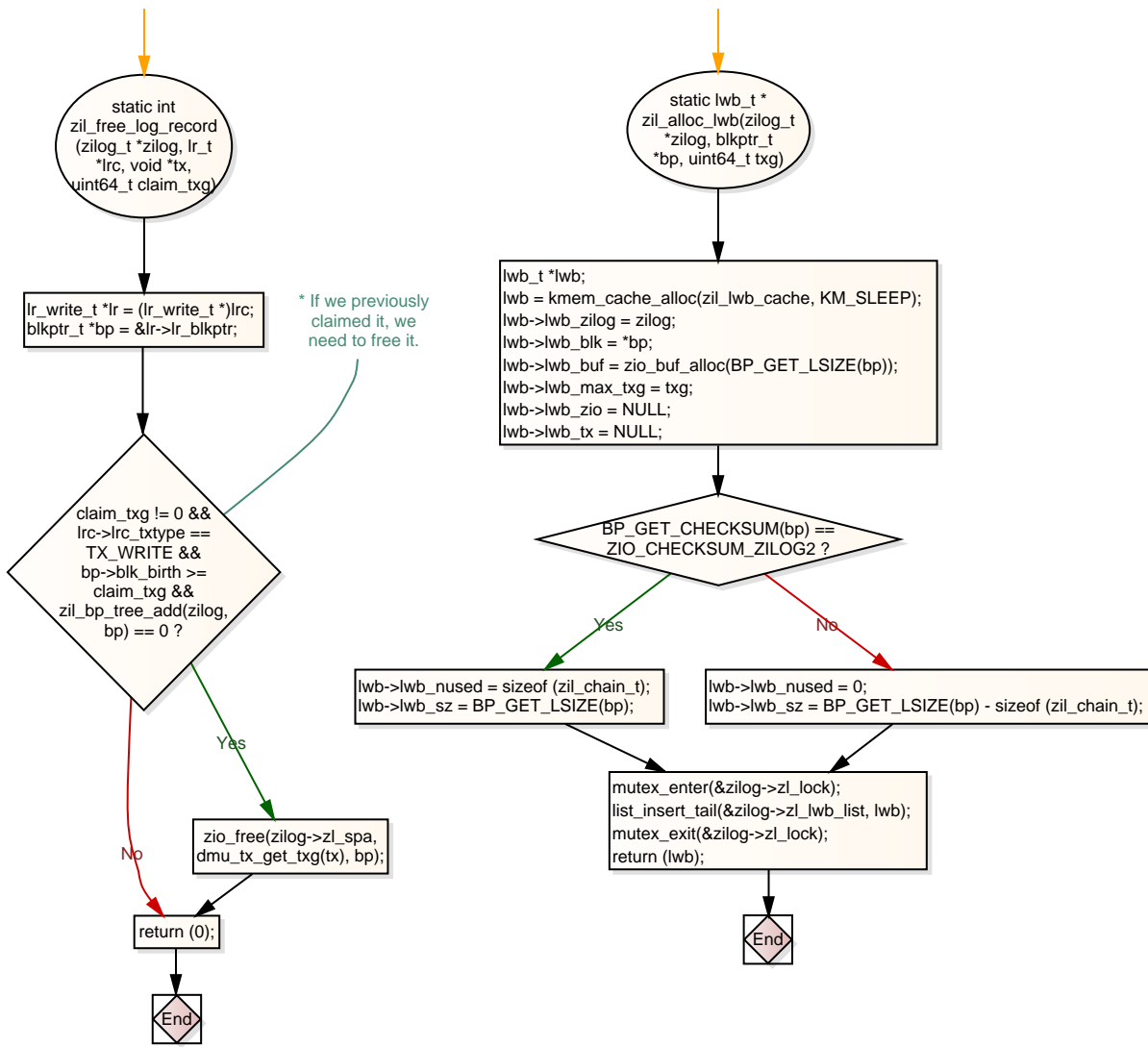
Yes

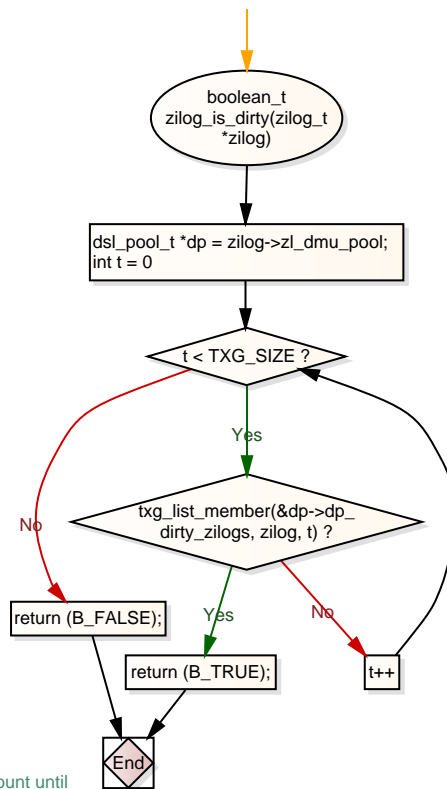
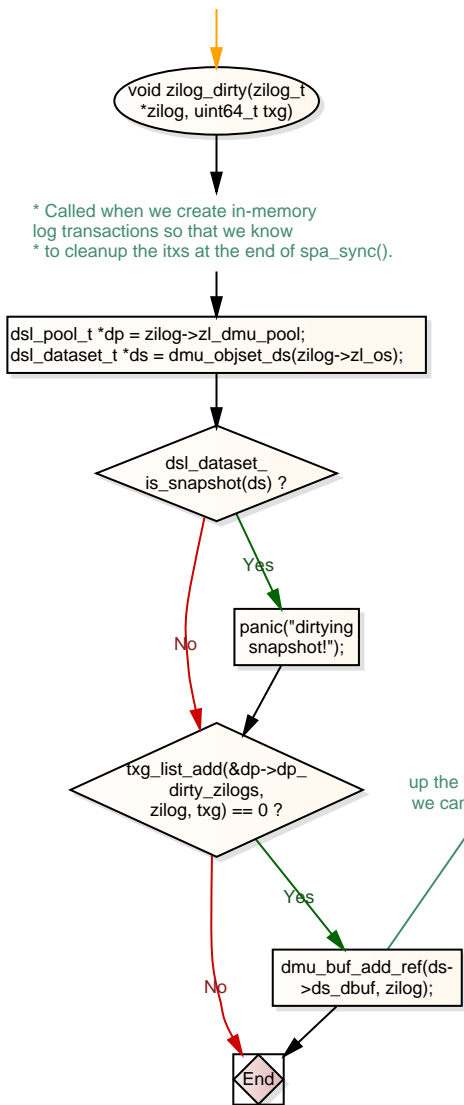


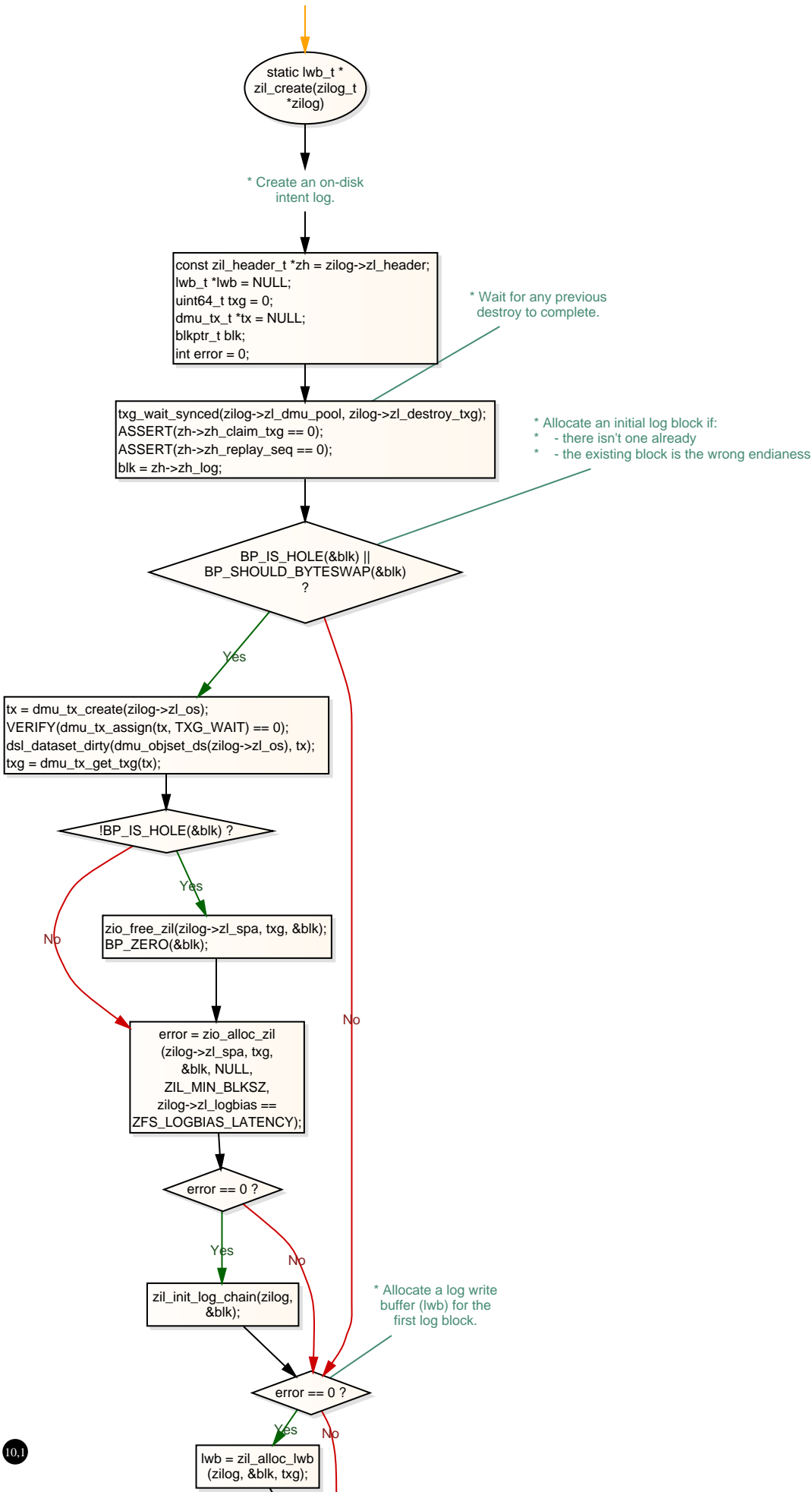
No

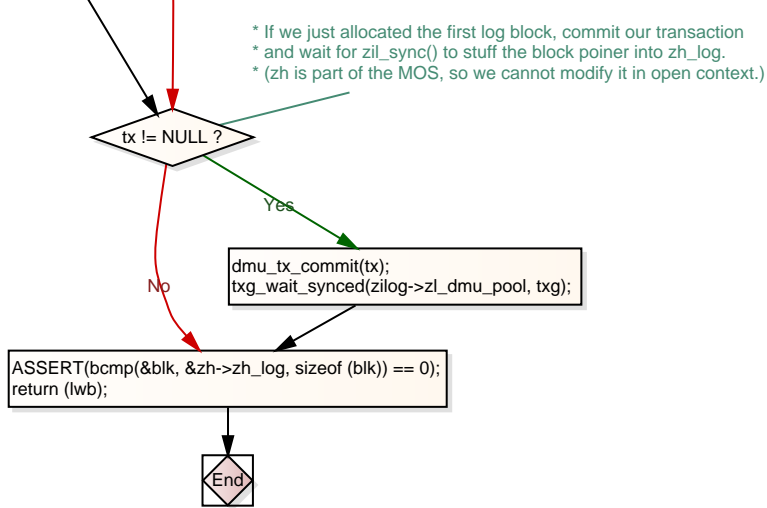


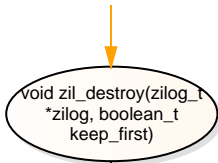








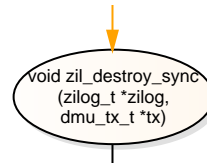




- * In one tx, free all log blocks and clear the log header.
- * If keep_first is set, then we're replaying a log with no content.
- * We want to keep the first block, however, so that the first synchronous transaction doesn't require a txg_wait_synced() in zil_create(). We don't need to txg_wait_synced()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

* Wait for any previous destroy to complete.

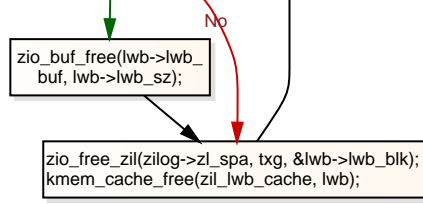
debugging aid

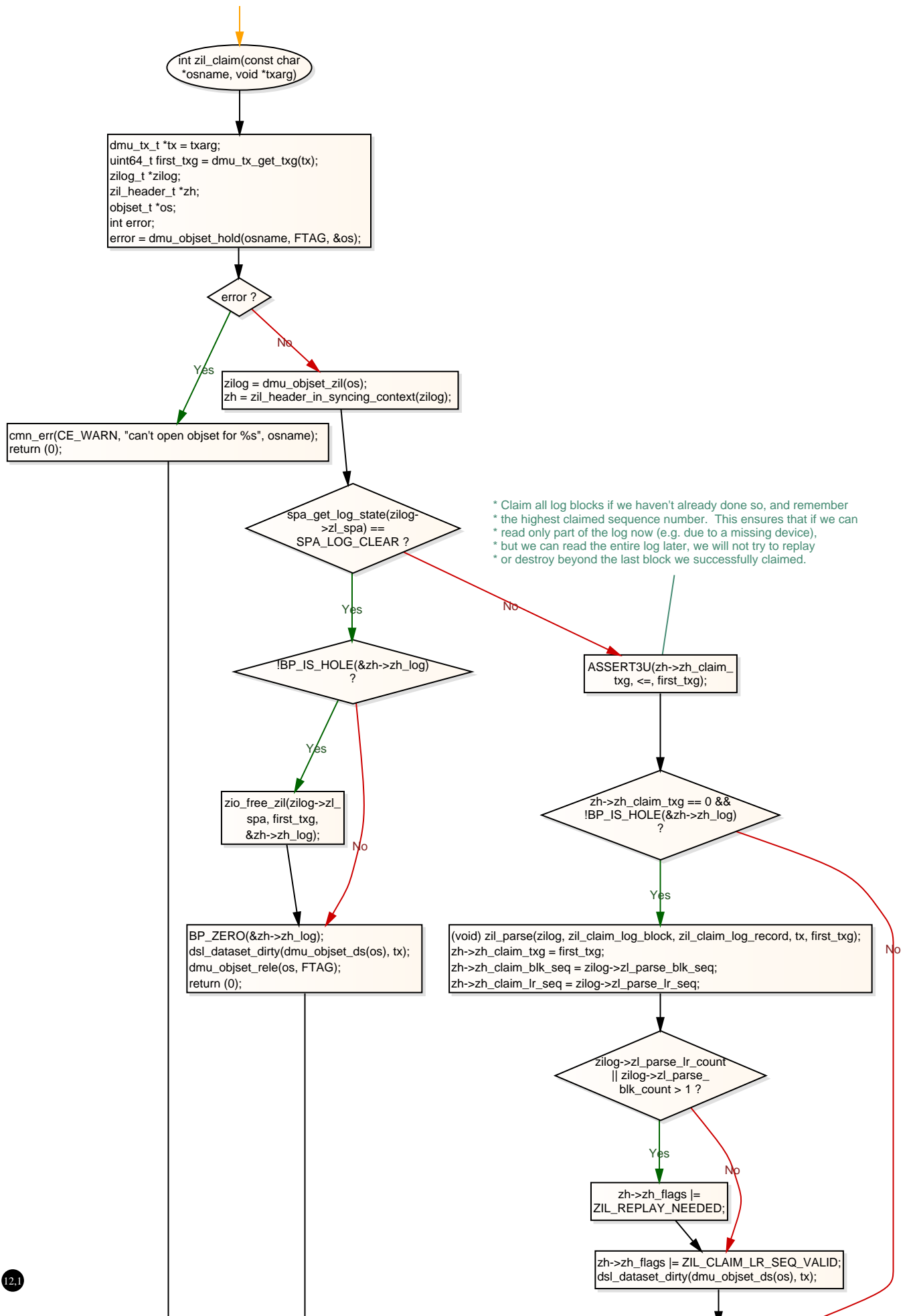


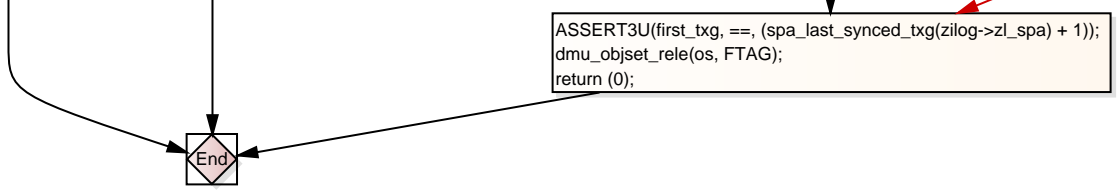
```

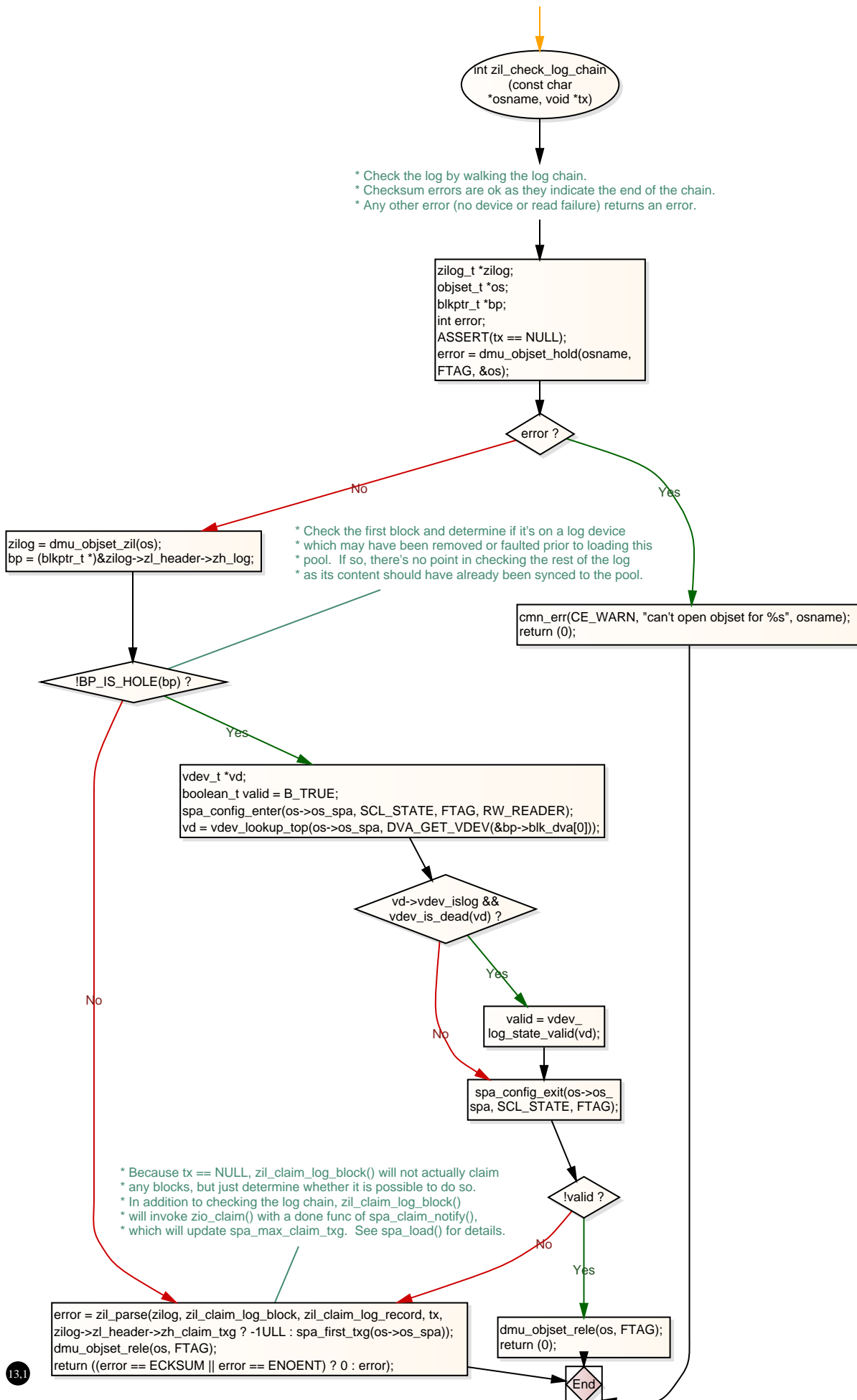
ASSERT(list_is_empty(&zilog->zl_lwb_list));
(void) zil_parse(zilog, zil_free_log_block, zil_free_log_record, tx, zilog->zl_header->zclaim_txg);
  
```

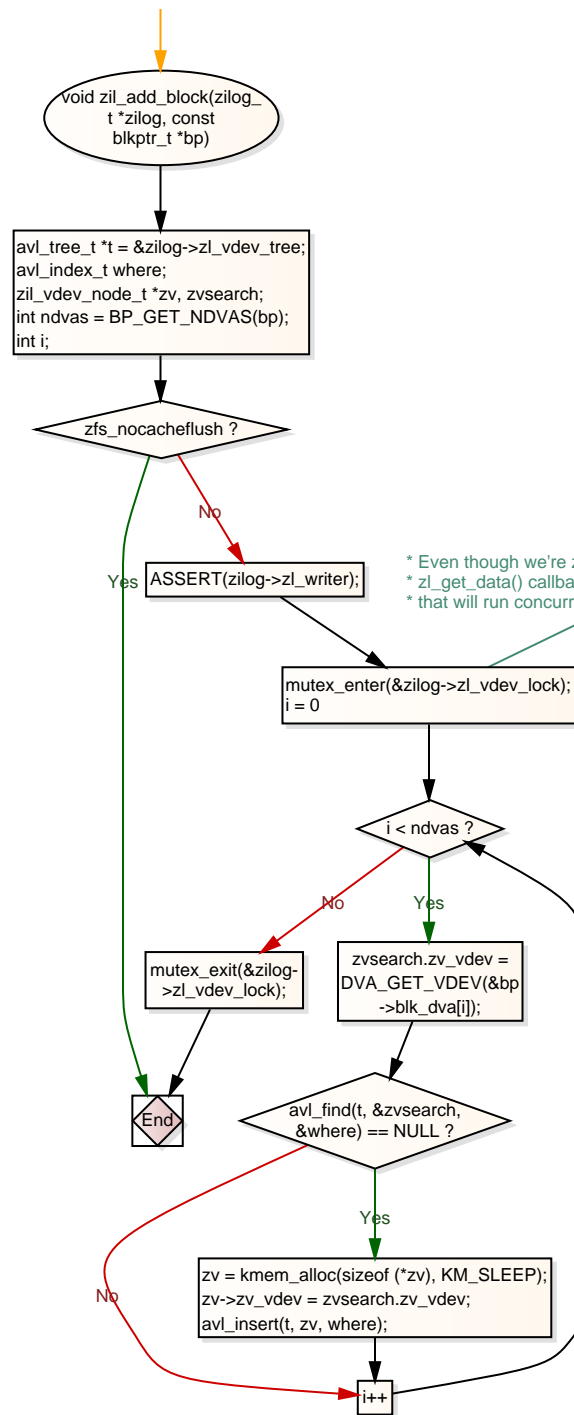
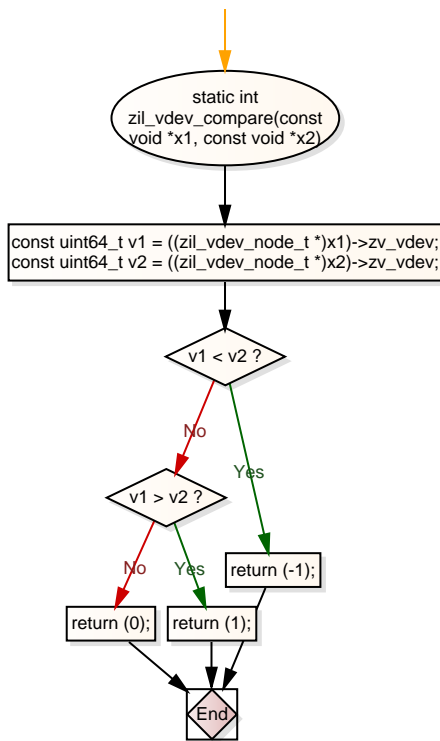


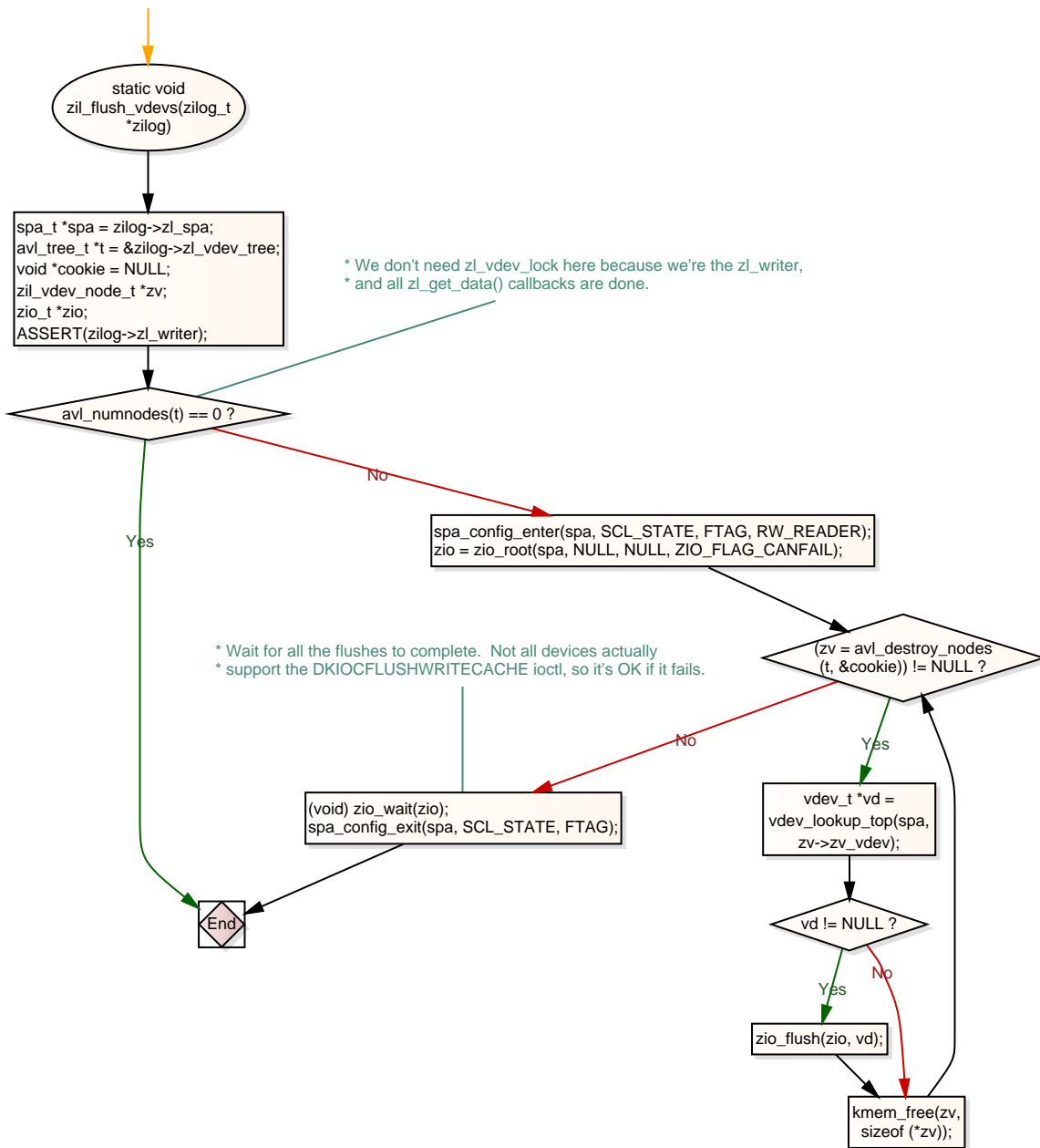


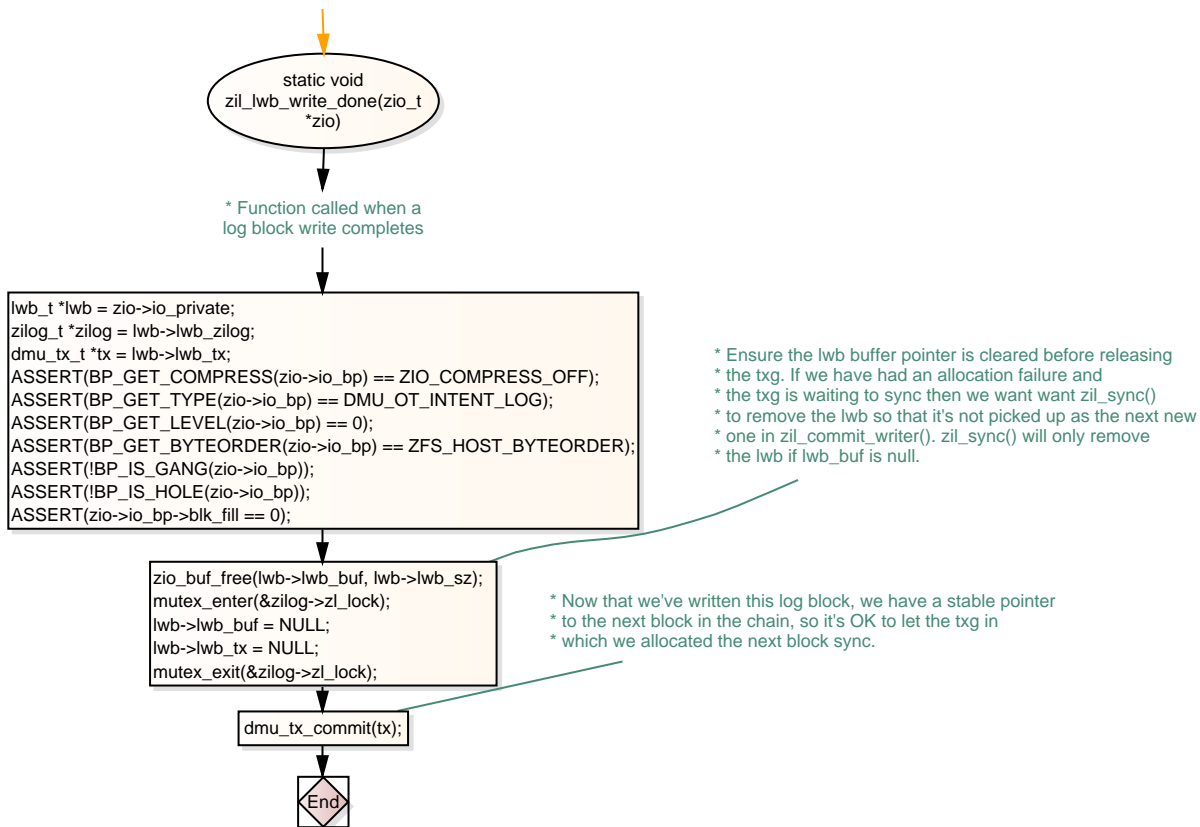


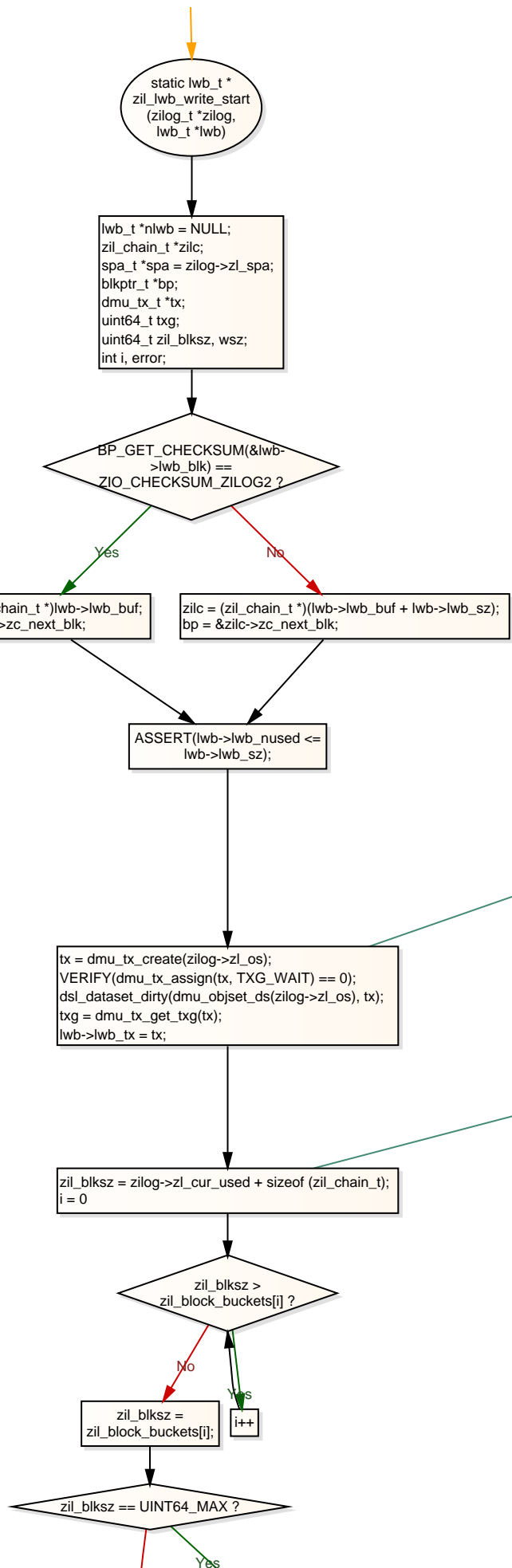
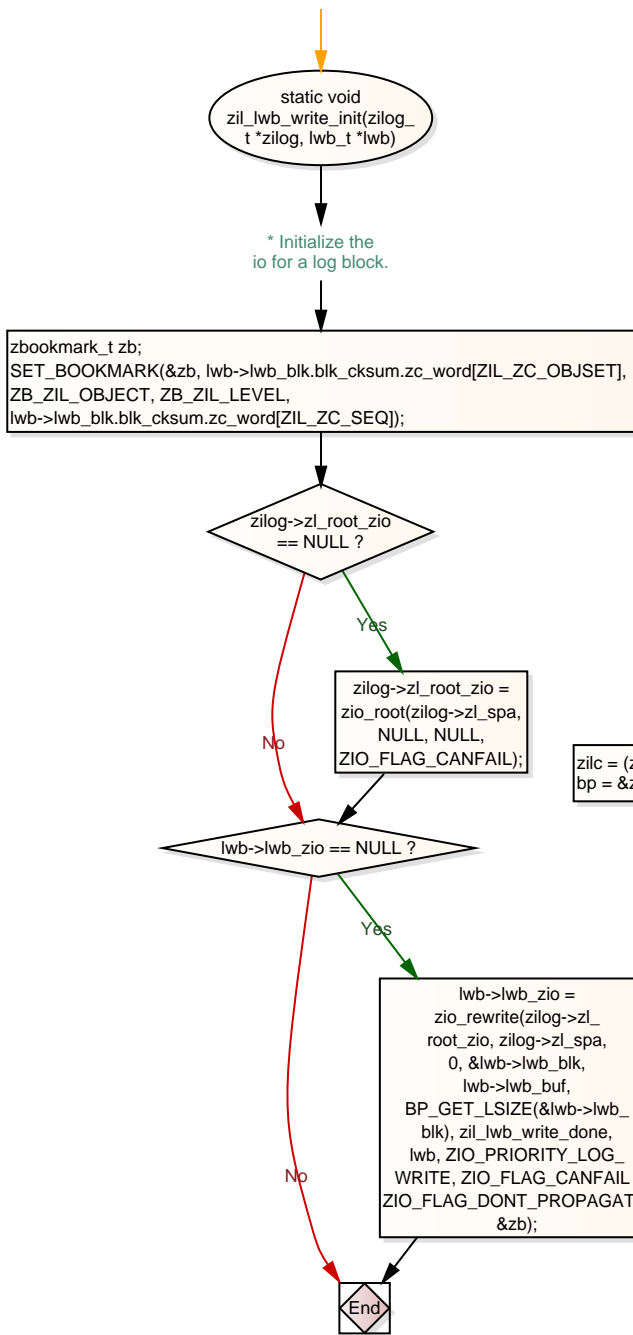


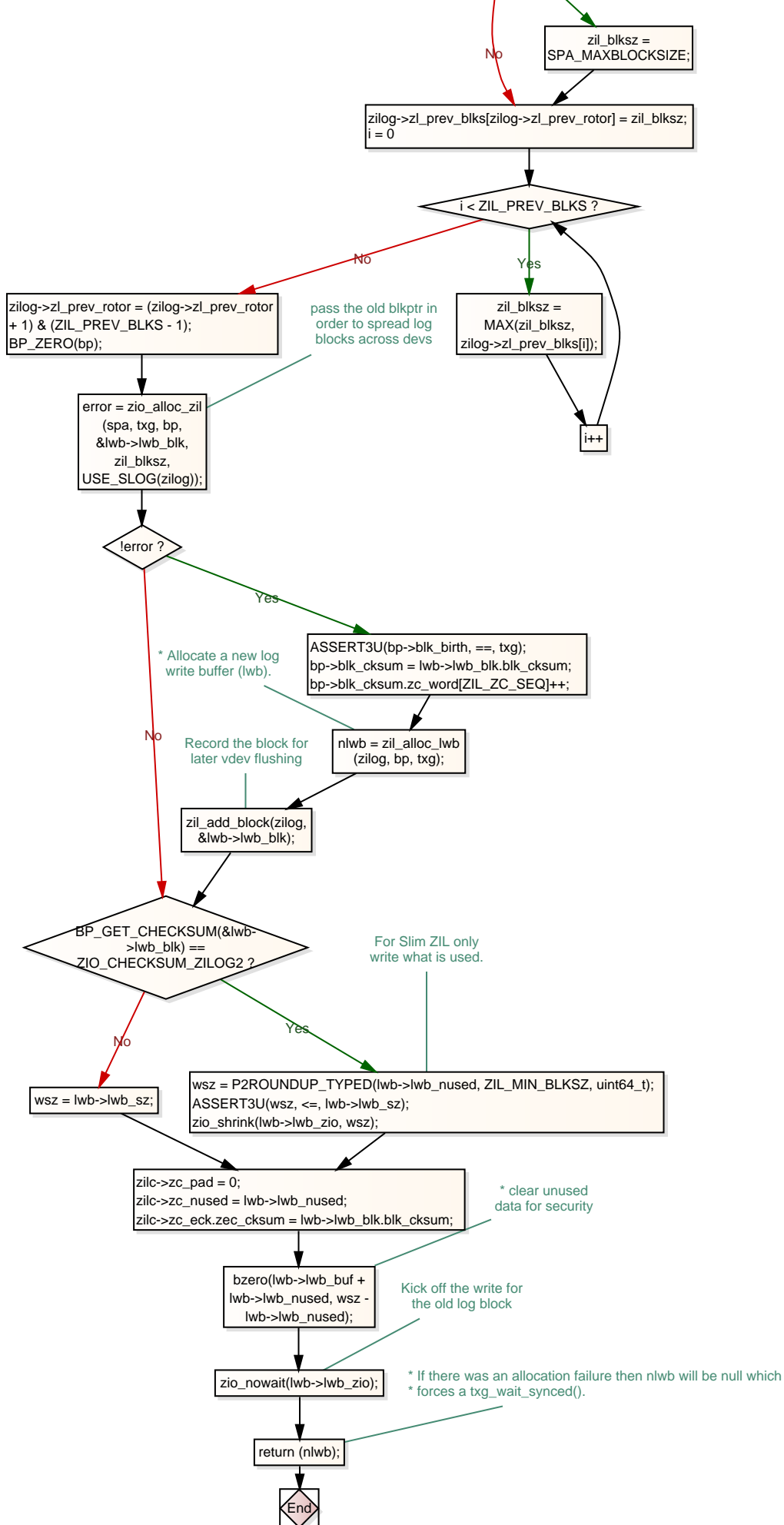










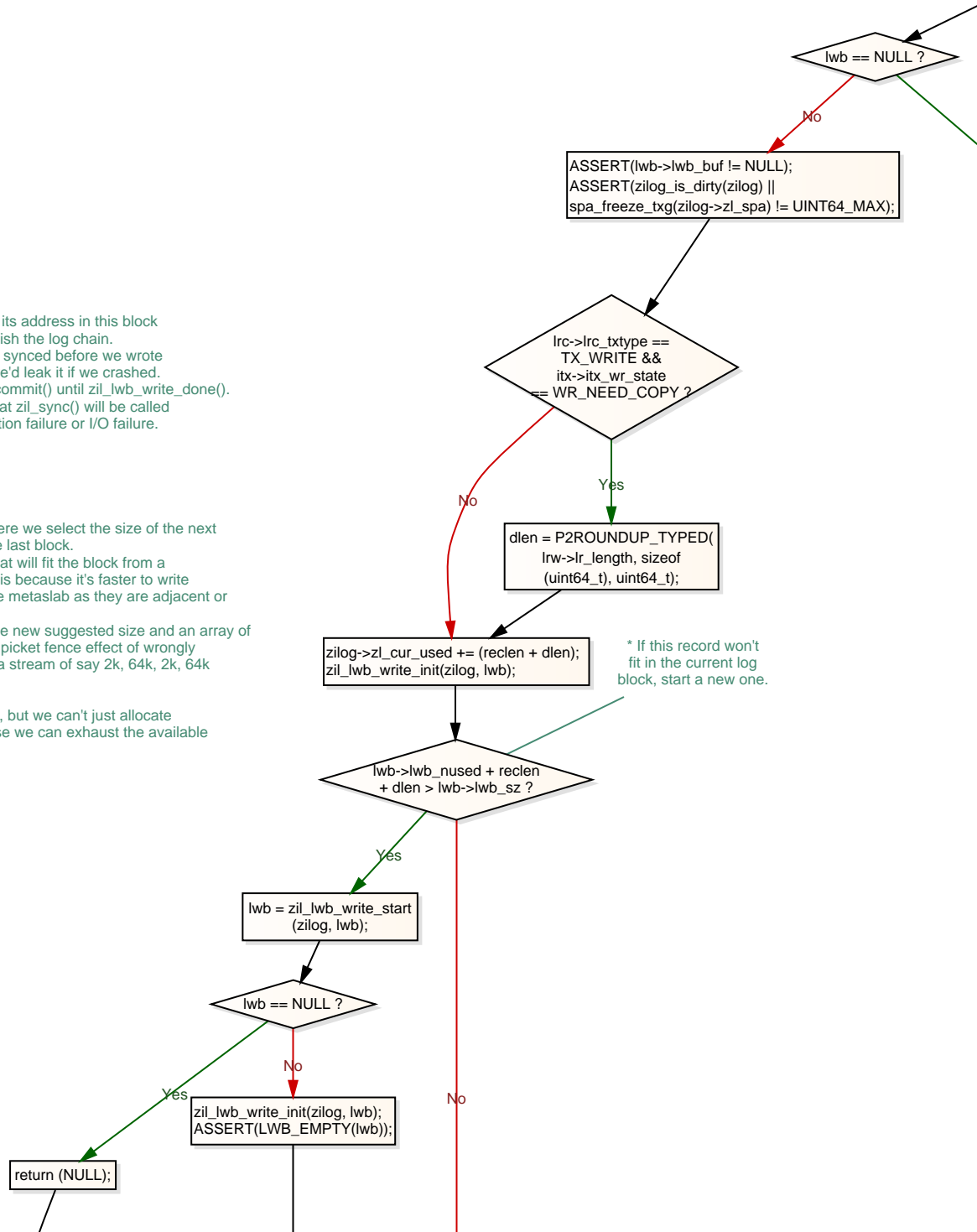


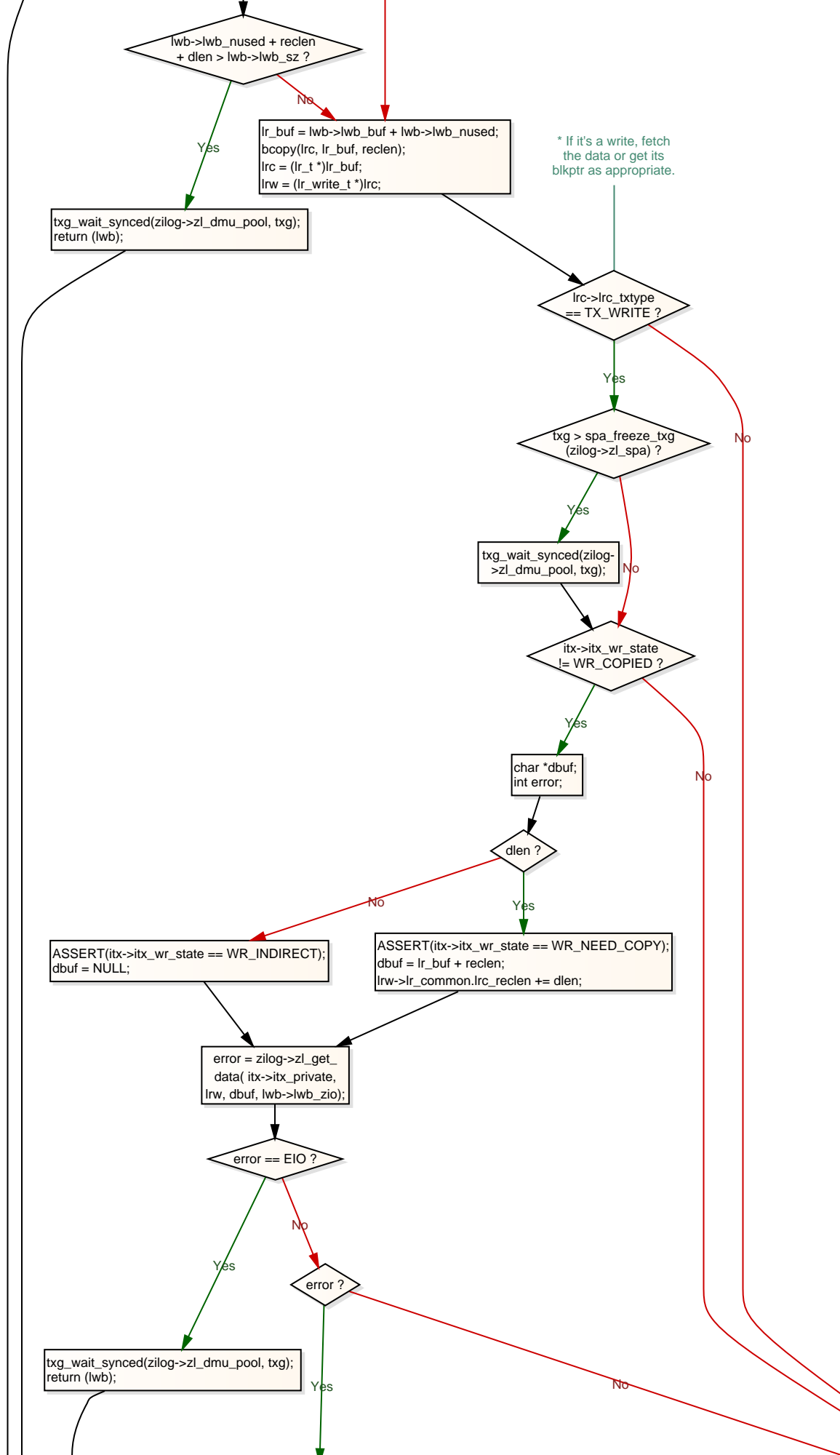
common log record

- * Allocate the next block and save its address in this block before writing it in order to establish the log chain.
- * Note that if the allocation of lwb synced before we wrote the block that points at it (lwb), we'd leak it if we crashed.
- * Therefore, we don't do dm_u_tx_commit() until zil_lwb_write_done().
- * We dirty the dataset to ensure that zil_sync() will be called to clean up in the event of allocation failure or I/O failure.

- * Log blocks are pre-allocated. Here we select the size of the next block, based on size used in the last block.
- * - first find the smallest bucket that will fit the block from a limited set of block sizes. This is because it's faster to write blocks allocated from the same metaslab as they are adjacent or close.
- * - next find the maximum from the new suggested size and an array of previous sizes. This lessens a picket fence effect of wrongly guessing the size if we have a stream of say 2k, 64k, 2k, 64k requests.

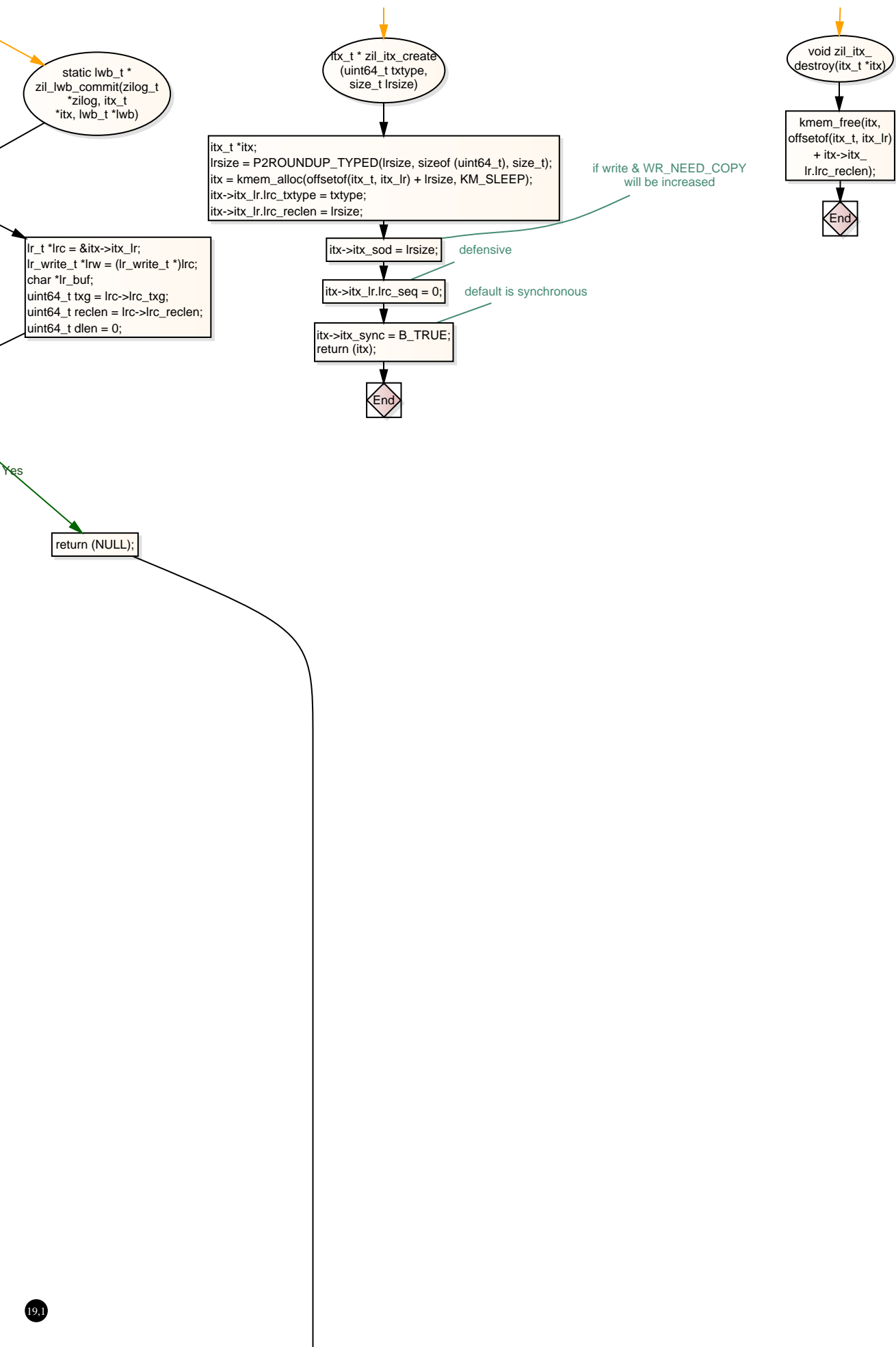
- * Note we only write what is used, but we can't just allocate the maximum block size because we can exhaust the available pool log space.



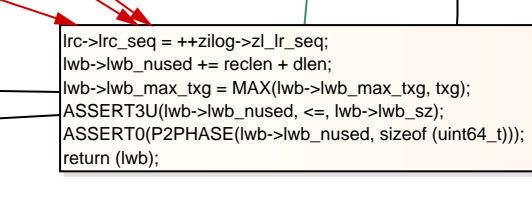


```
ASSERT(error == ENOENT || error == EEXIST || error == EALREADY);  
return (lwb);
```

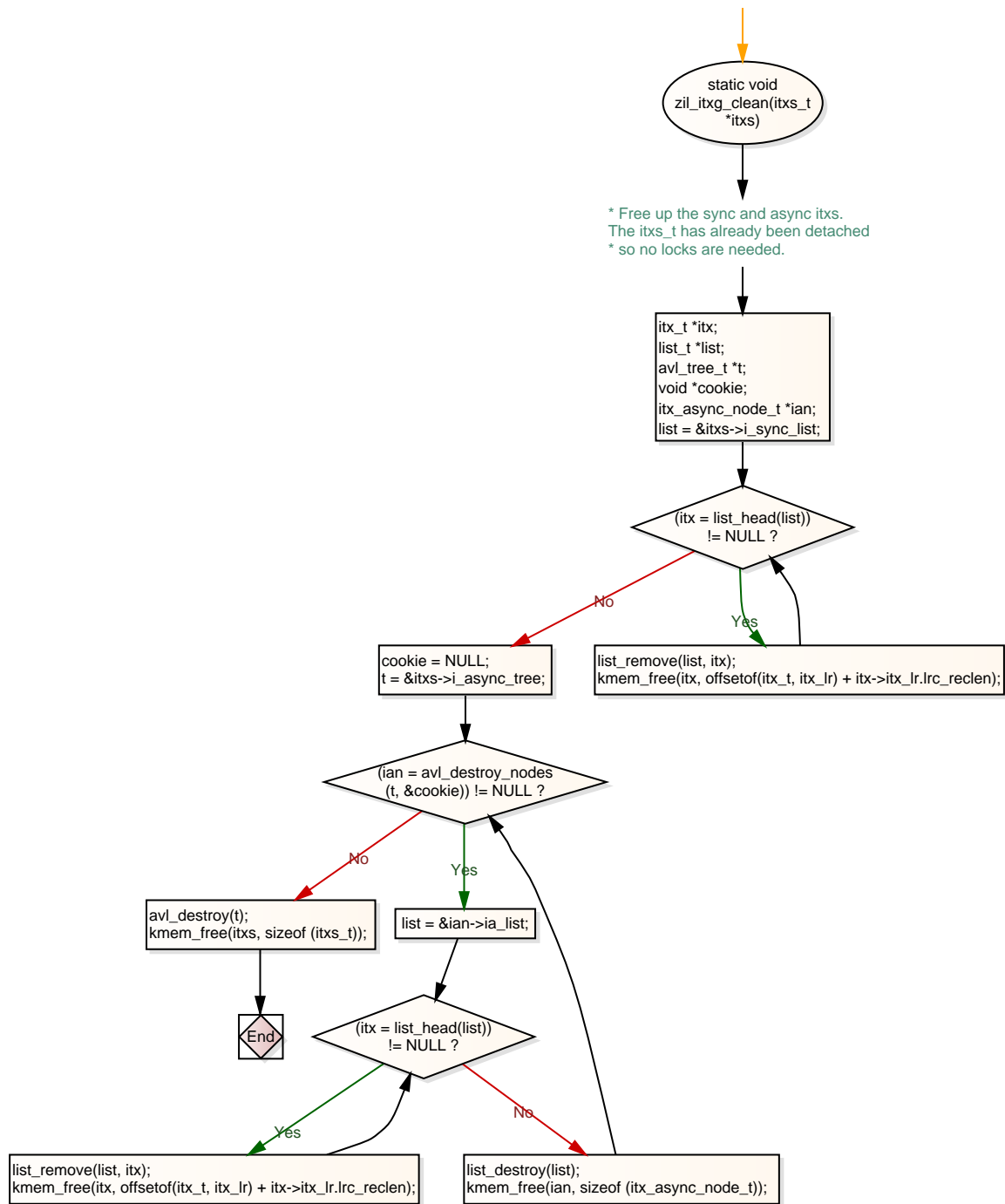


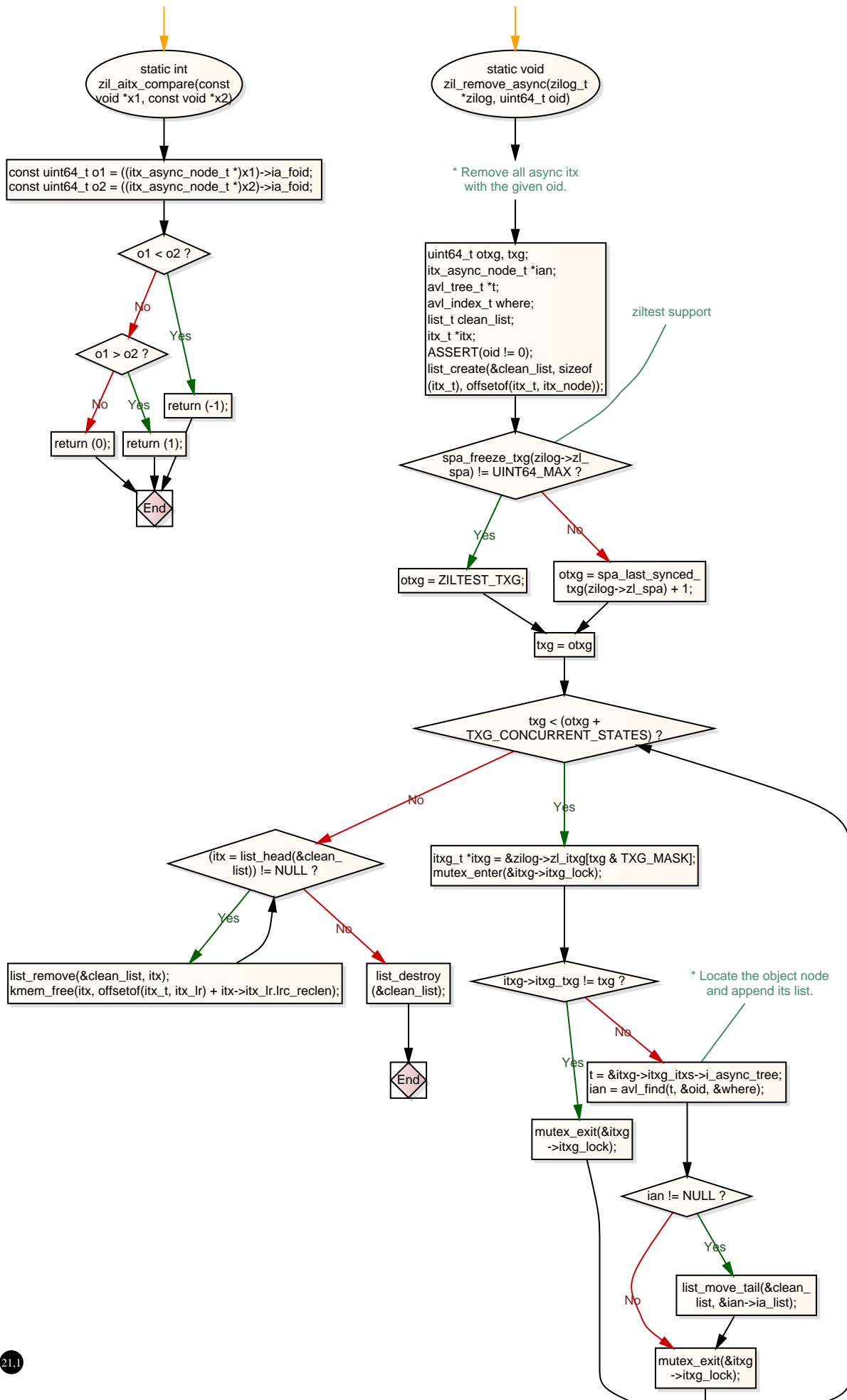
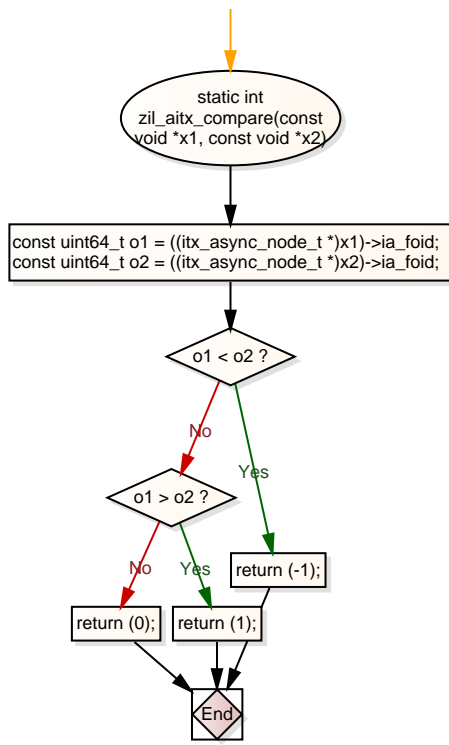


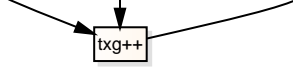
* We're actually making an entry, so update lrc_seq to be the
* log record sequence number. Note that this is generally not
* equal to the itx sequence number because not all transactions
* are synchronous, and sometimes spa_sync() gets there first.
we are single threaded

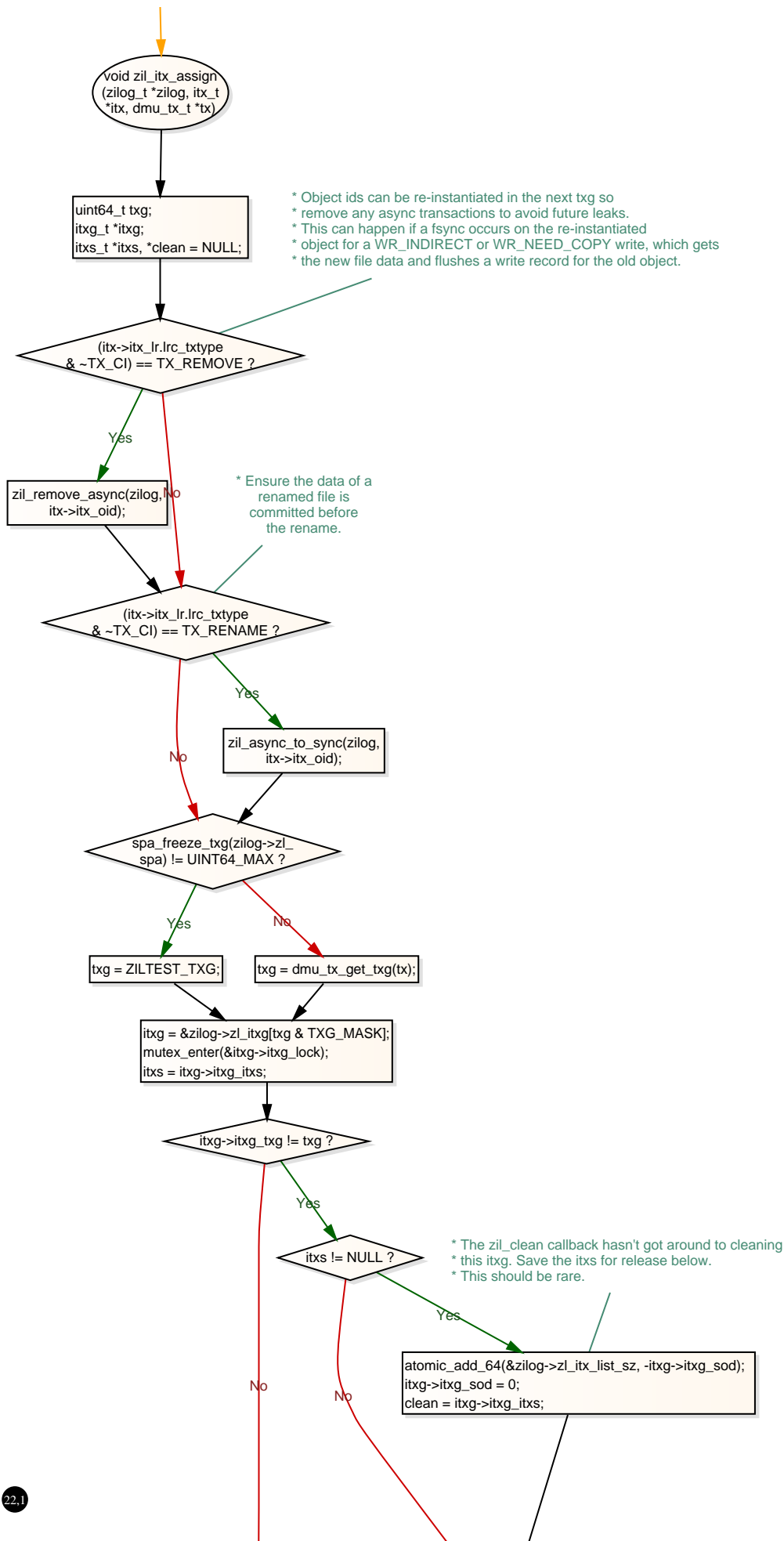


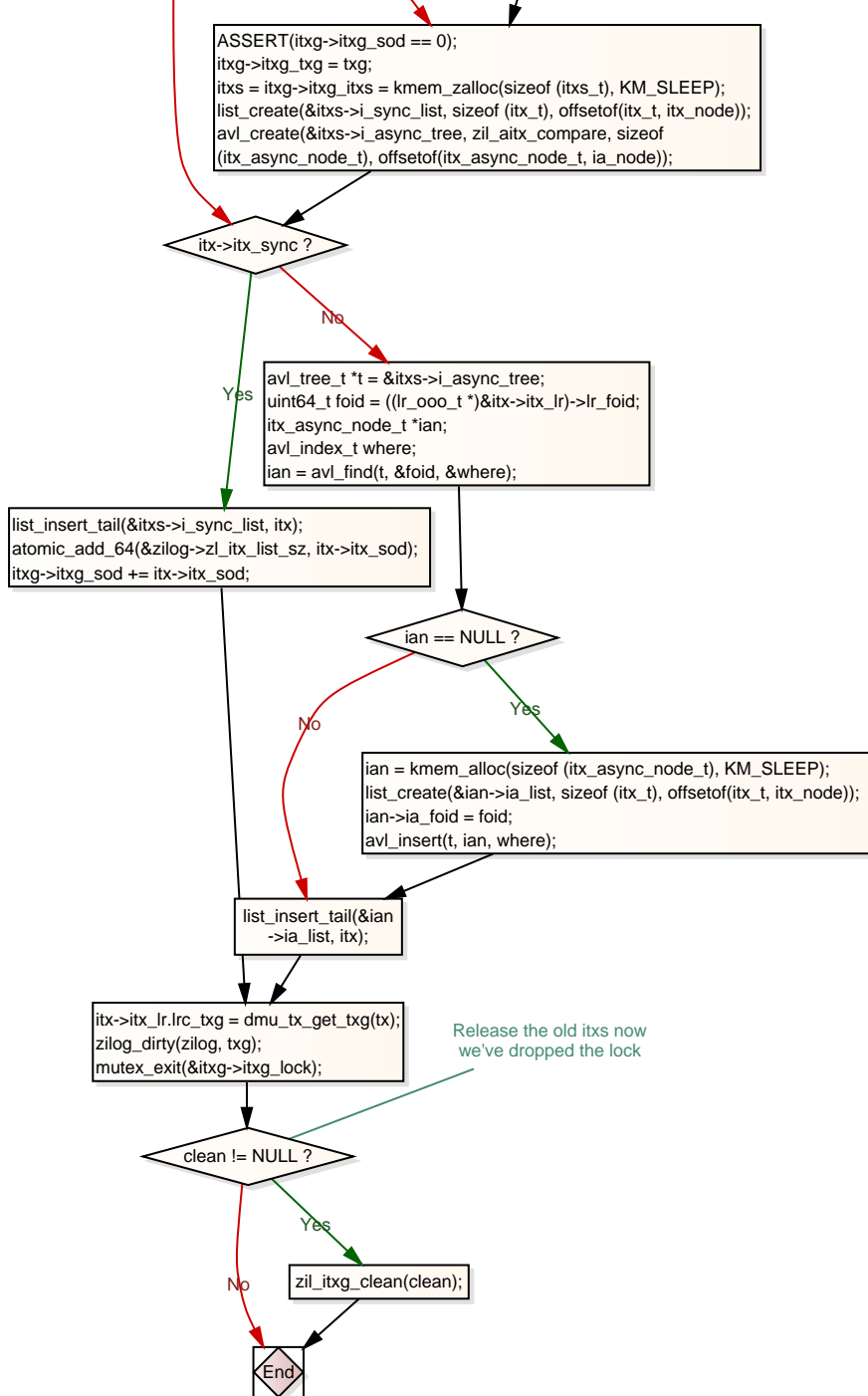
```
lrc->lrc_seq = ++zillog->zl_lr_seq;  
lwb->lwb_nused += reclen + dlen;  
lwb->lwb_max_txg = MAX(lwb->lwb_max_txg, txg);  
ASSERT3U(lwb->lwb_nused, <=, lwb->lwb_sz);  
ASSERT0(P2PHASE(lwb->lwb_nused, sizeof (uint64_t)));  
return (lwb);
```

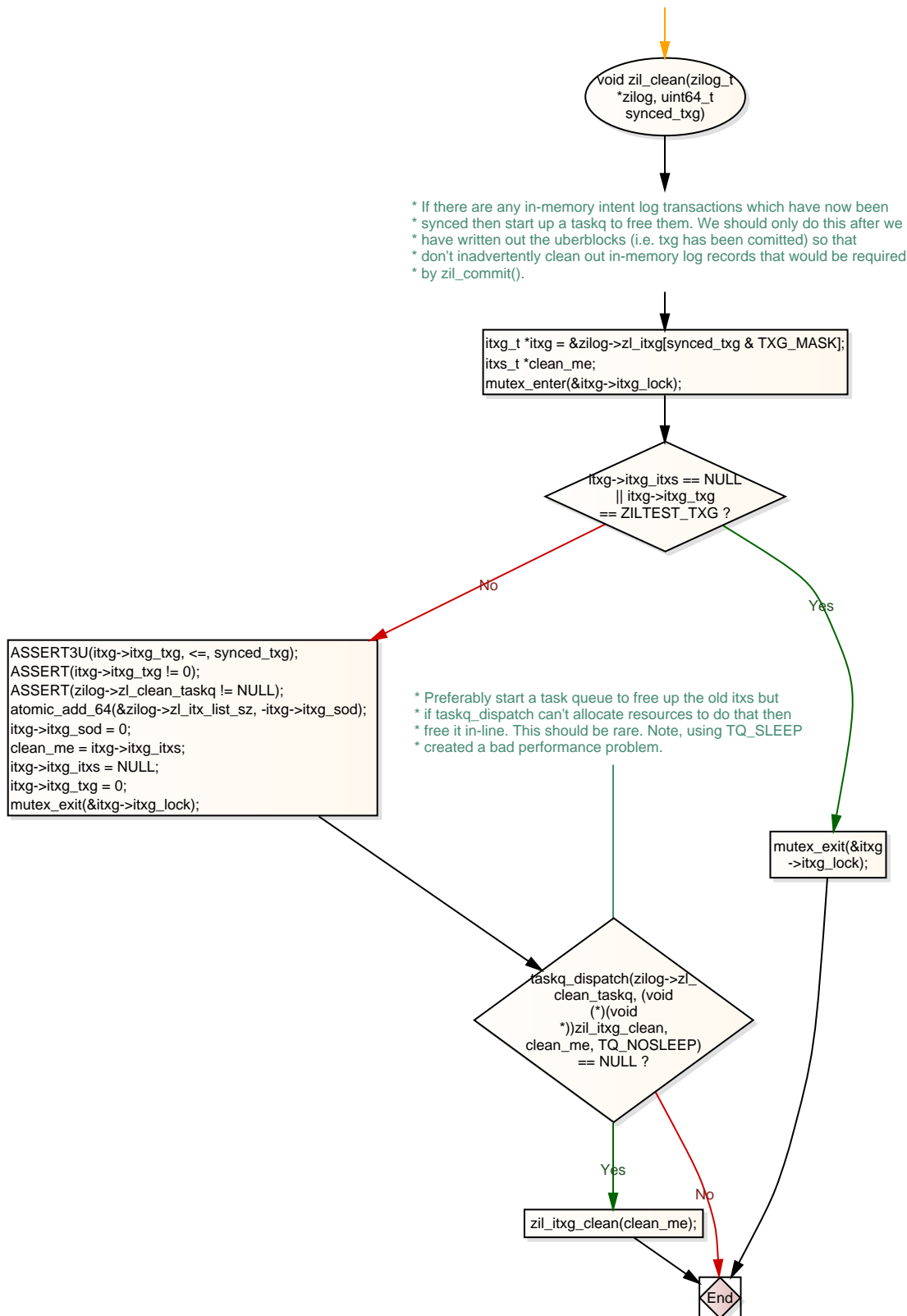


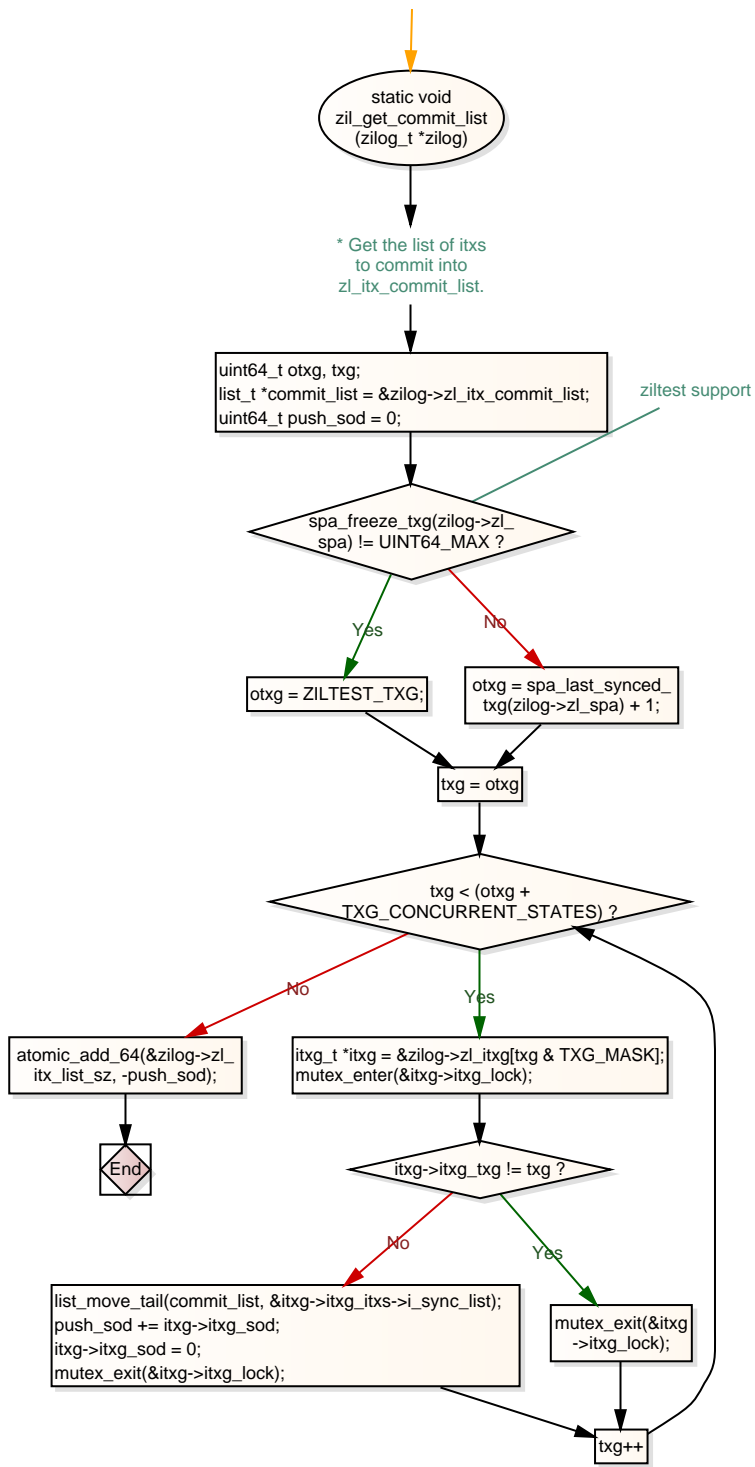


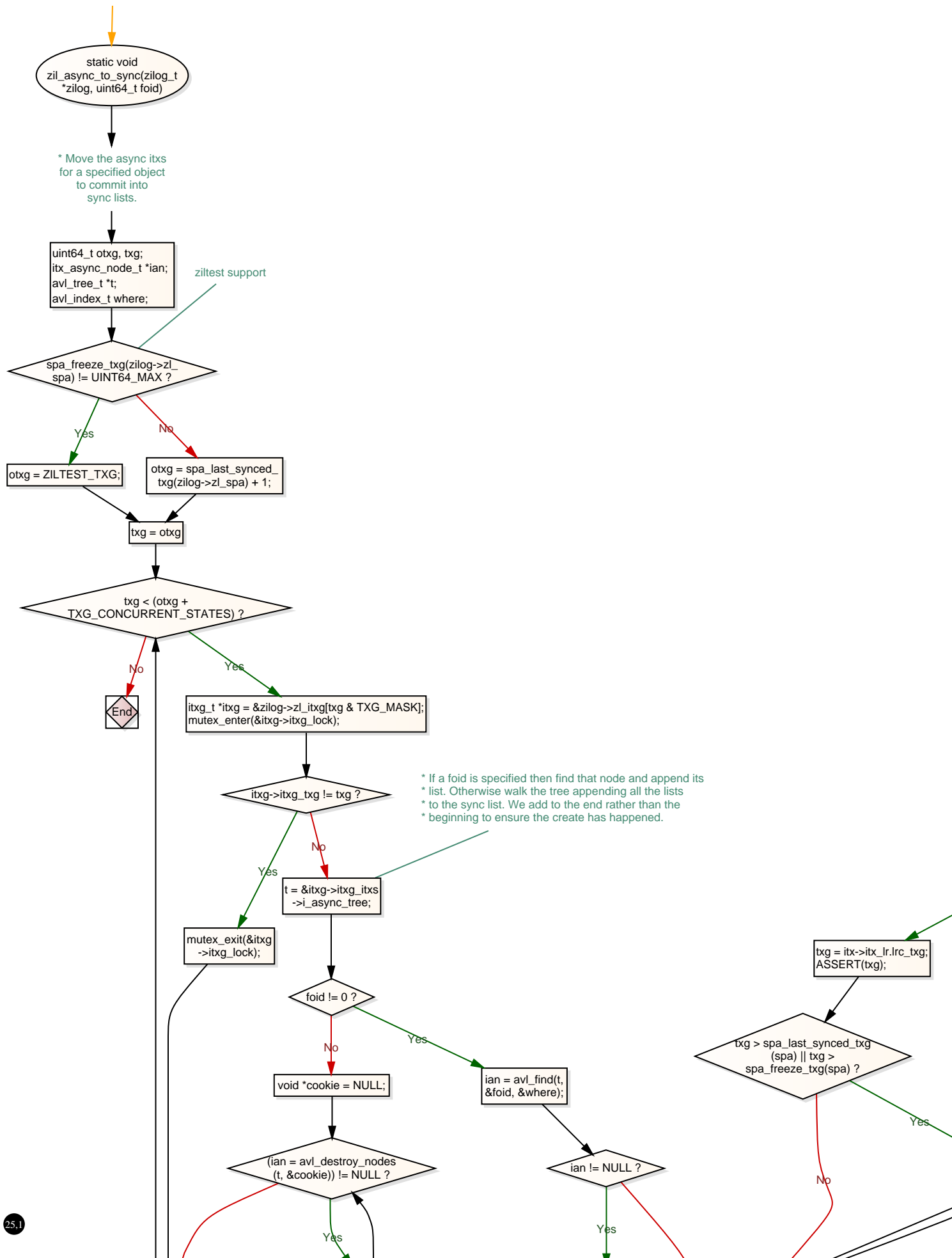


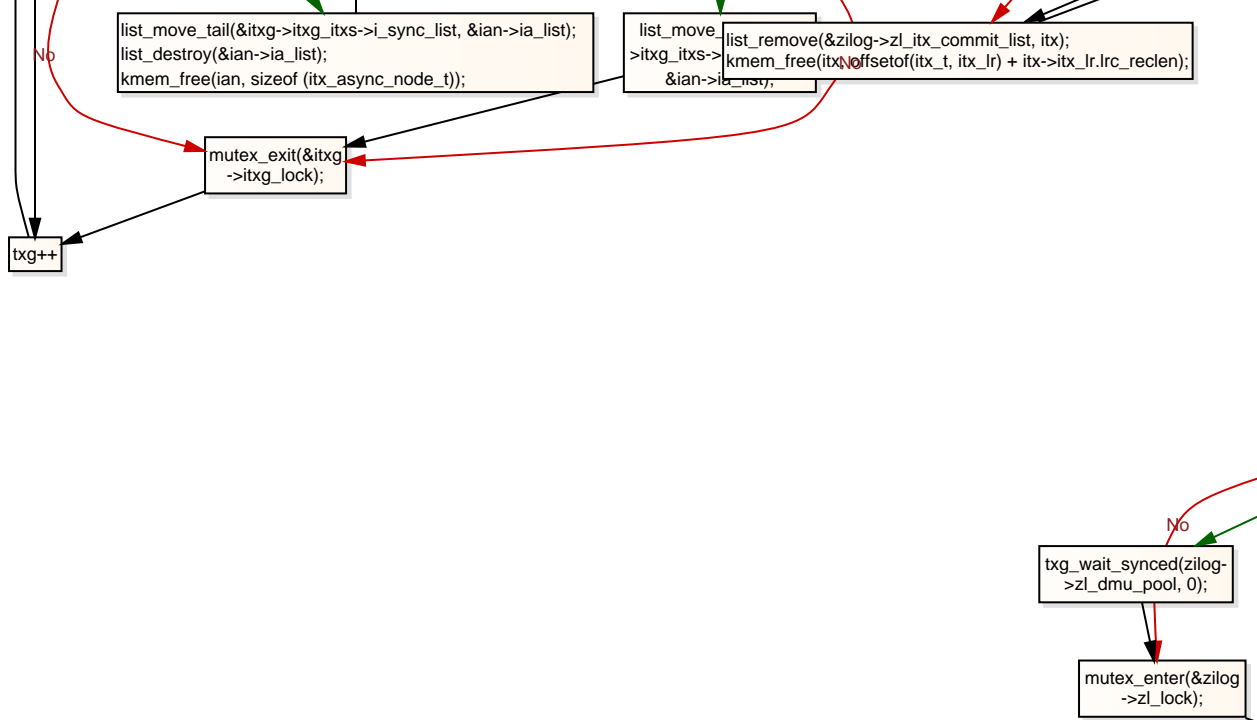


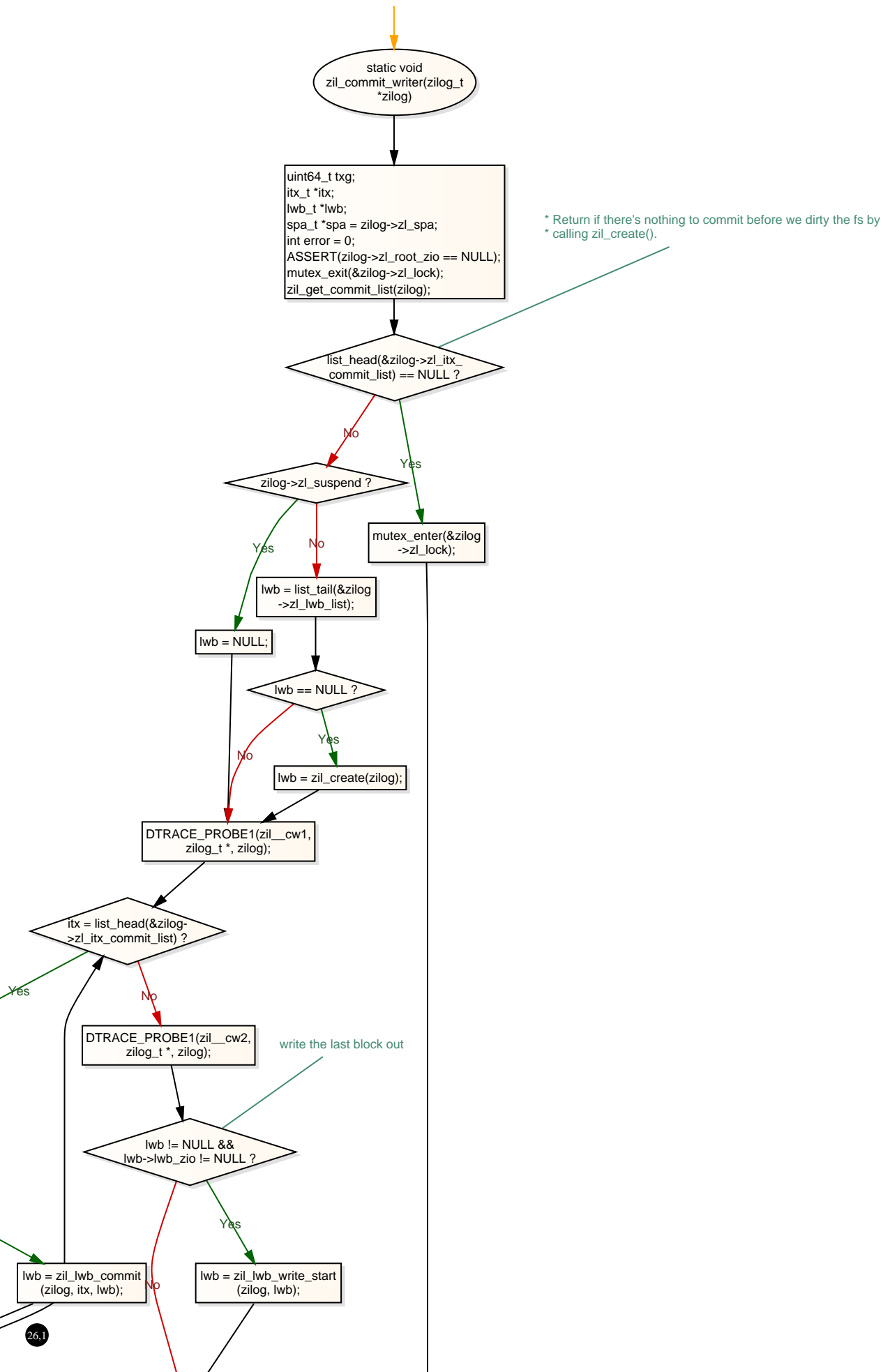












zilog->zl_cur_used = 0;

* Wait if necessary for
the log blocks to be on
stable storage.

zilog->zl_root_zio ?

Yes

error = zio_wait(zilog->zl_root_zio);
zilog->zl_root_zio = NULL;
zil_flush_vdevs(zilog);

error || lwb == NULL ?

No

Yes

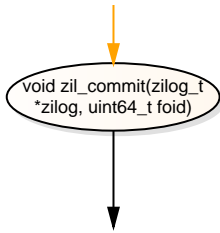
* Remember the highest committed log sequence number for ztest.
* We only update this value when all the log writes succeeded,
* because ztest wants to ASSERT that it got the whole log chain.

error == 0 &&
lwb != NULL ?

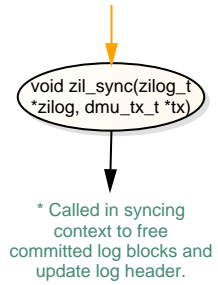
Yes

zilog->zl_commit_lr_seq
= zilog->zl_lr_seq;

End

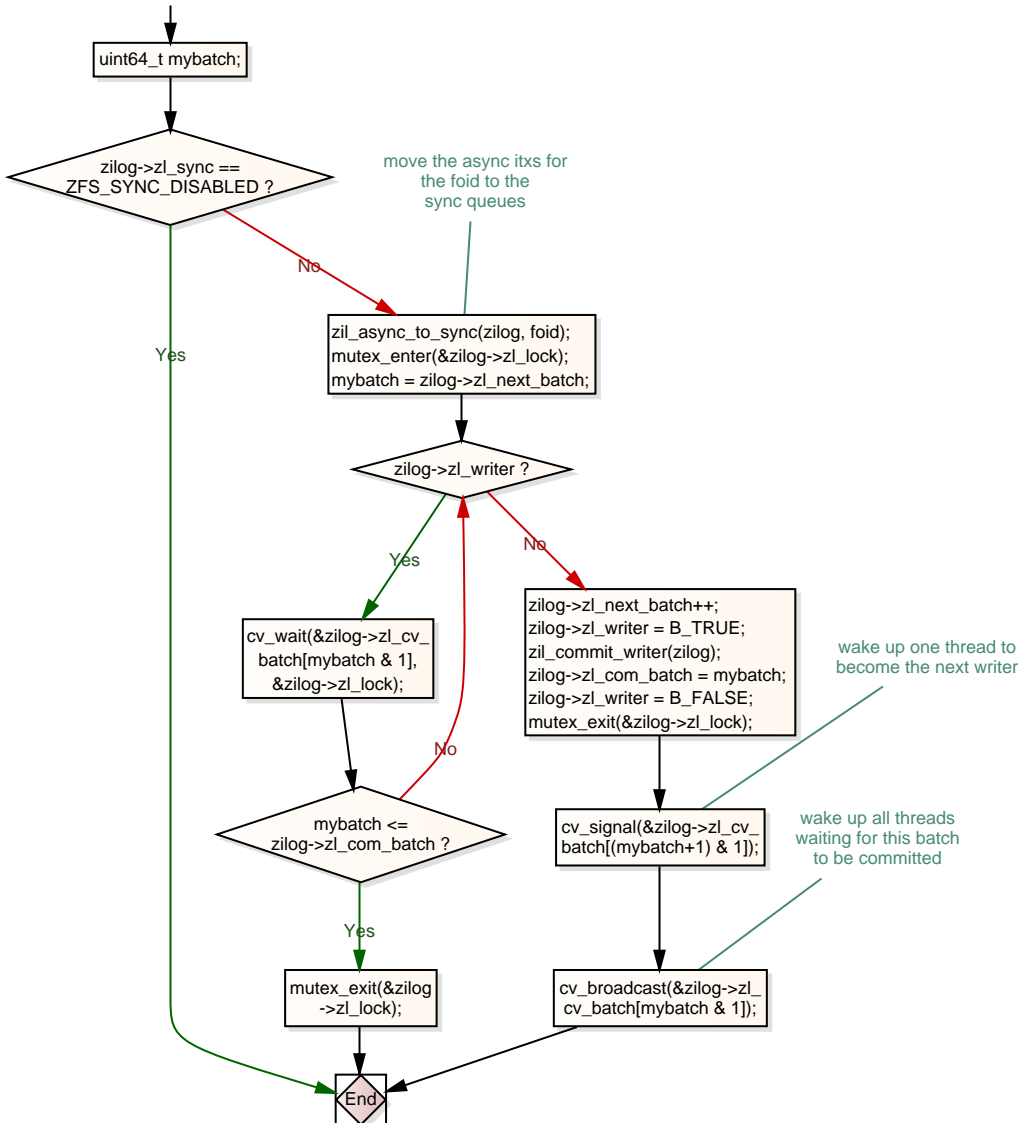
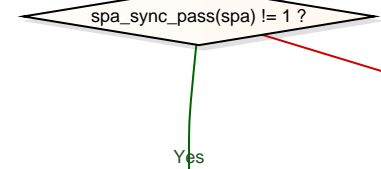


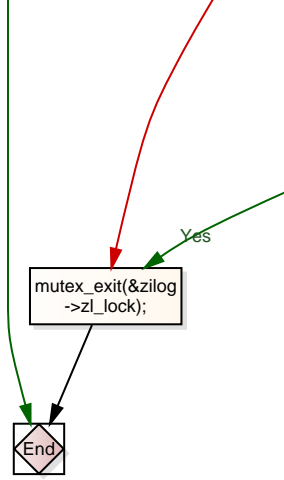
- * Commit zfs transactions to stable storage.
- * If foid is 0 push out all transactions, otherwise push only those for that object or might reference that object.
- *
- * Itxs are committed in batches. In a heavily stressed zil there will be a commit writer thread who is writing out a bunch of itxs 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.

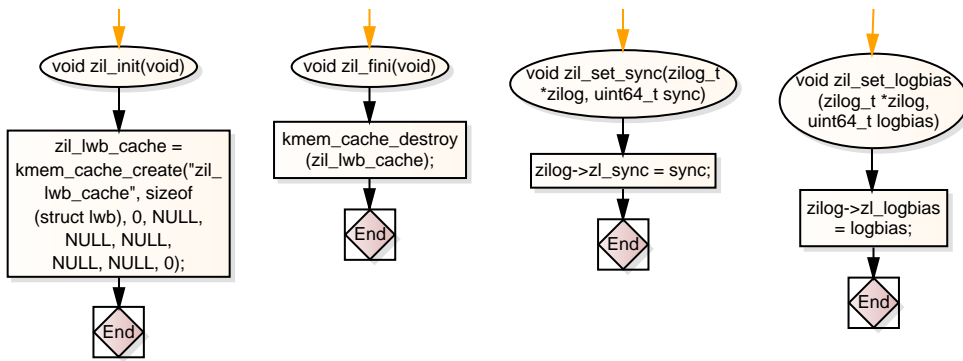


```

zil_header_t *zh = zil_header_in_syncing_context(zilog);
uint64_t txg = dmu_tx_get_txg(tx);
spa_t *spa = zilog->zl_spa;
uint64_t *replayed_seq = &zilog->zl_replayed_seq[txg & TXG_MASK];
lwb_t *lwb;
  
```







* We don't zero out zl_destroy_txg, so make sure we don't try to destroy it twice.

