

\* 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.  
\*

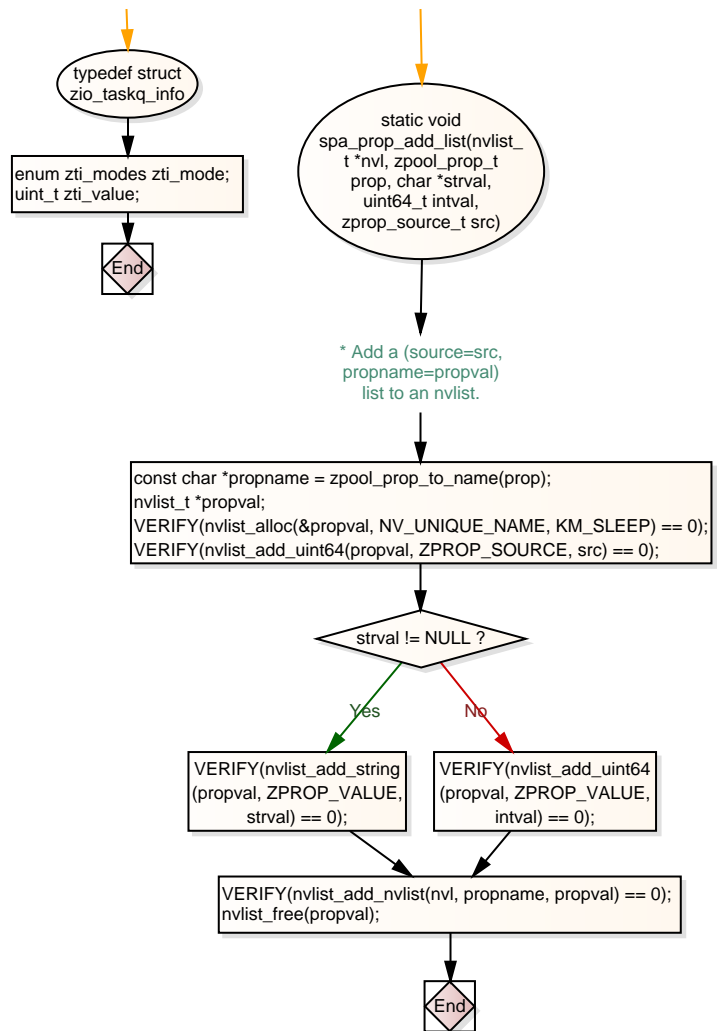
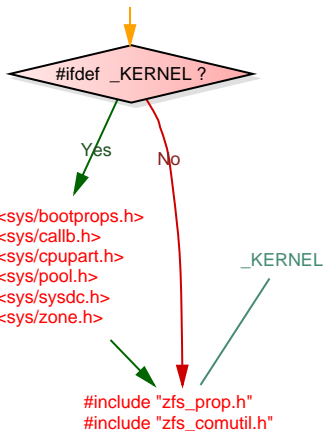
\* You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE  
\* or <http://www.opensolaris.org/os/licensing>.  
\* See the License for the specific language governing permissions  
\* and limitations under the License.  
\*

\* 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 2011 Nexenta Systems, Inc. All rights reserved.  
\* Copyright (c) 2012 by Delphix. All rights reserved.  
\* This file contains all the routines used when modifying on-disk SPA state.  
\* This includes opening, importing, destroying, exporting a pool, and syncing a  
\* pool.

```
#include <sys/zfs_context.h>
#include <sys/fm/fs/zfs.h>
#include <sys/spa_impl.h>
#include <sys/zio.h>
#include <sys/zio_checksum.h>
#include <sys/dmu.h>
#include <sys/dmu_tx.h>
#include <sys/zap.h>
#include <sys/zil.h>
#include <sys/ddt.h>
#include <sys/vdev_impl.h>
#include <sys/metastab.h>
#include <sys/metastab_impl.h>
#include <sys/uberblock_impl.h>
#include <sys/txg.h>
#include <sys/avl.h>
#include <sys/dmu_traverse.h>
#include <sys/dmu_objset.h>
#include <sys/unique.h>
#include <sys/dsl_pool.h>
#include <sys/dsl_dataset.h>
#include <sys/dsl_dir.h>
#include <sys/dsl_prop.h>
#include <sys/dsl_synctask.h>
#include <sys/fs/zfs.h>
#include <sys/arc.h>
#include <sys/callb.h>
#include <sys/systeminfo.h>
#include <sys/spa_boot.h>
#include <sys/zfs_ioctl.h>
#include <sys/dsl_scan.h>
#include <sys/zfeature.h>
```



typedef enum  
zti\_modes {...}

```
#define ZTI_FIX(n)      { zti_mode_fixed, (n) }  
#define ZTI_PCT(n)     { zti_mode_online_percent, (n) }  
#define ZTI_BATCH      { zti_mode_batch, 0 }  
#define ZTI_NULL       { zti_mode_null, 0 }  
#define ZTI_ONE        ZTI_FIX(1)
```

```
static const char *const  
zio_taskq_types[ZIO_  
TASKQ_TYPES] = {  
    "issue", "issue_high",  
    "intr", "intr_high";
```

\* Define the taskq threads for the following I/O types:  
\* NULL, READ, WRITE, FREE, CLAIM, and IOCTL  
ISSUE ISSUE\_HIGH INTR INTR\_HIGH

```
const zio_taskq_info_t zio_taskqs[ZIO_TYPES][ZIO_TASKQ_TYPES]  
= { { ZTI_ONE, ZTI_NULL, ZTI_ONE, ZTI_NULL },  
    ZTI_FIX(8), ZTI_NULL, ZTI_BATCH, ZTI_NULL  
    ,  
    ZTI_BATCH, ZTI_FIX(5), ZTI_FIX(8), ZTI_FIX(5)  
    ,  
    ZTI_FIX(100), ZTI_NULL, ZTI_ONE, ZTI_NULL  
    ,  
    ZTI_ONE, ZTI_NULL, ZTI_ONE, ZTI_NULL  
    ,  
    ZTI_ONE, ZTI_NULL, ZTI_ONE, ZTI_NULL  
    ,  
    static dsl_syncfunc_t spa_sync_version;  
    static dsl_syncfunc_t spa_sync_props;  
    static dsl_checkfunc_t spa_change_guid_check;  
    static dsl_syncfunc_t spa_change_guid_sync;  
    static boolean_t spa_has_active_shared_spare(spa_t *spa);  
    static int spa_load_impl(spa_t *spa, uint64_t, nvlist_t *config, spa_load_state_t  
    state, spa_import_type_t type, boolean_t mosconfig, char **ereport);  
    static void spa_vdev_resilver_done(spa_t *spa);
```

1 thread per cpu in pset

```
uint_t zio_taskq_batch_pct = 100;  
id_t zio_taskq_psrset_bind = PS_NONE;
```

use SDC scheduling class

```
boolean_t zio_taskq_sysdc  
= B_TRUE;
```

base duty cycle

```
uint_t zio_taskq_basedc  
= 80;
```

no process ==> no sysdc

```
boolean_t spa_create_process = B_TRUE;  
extern int zfs_sync_pass_deferred_free;
```

\* This (illegal) pool name is used when temporarily importing a spa\_t in order  
\* to get the vdev stats associated with the imported devices.

\* =====  
\* SPA properties routines  
\* =====

```
#define TRYIMPORT_NAME  
"$import"
```

```
#ifdef _KERNEL ?
```

\* Get the root pool information from the  
root disk, then import the root pool  
\* during the system boot up time.

Yes

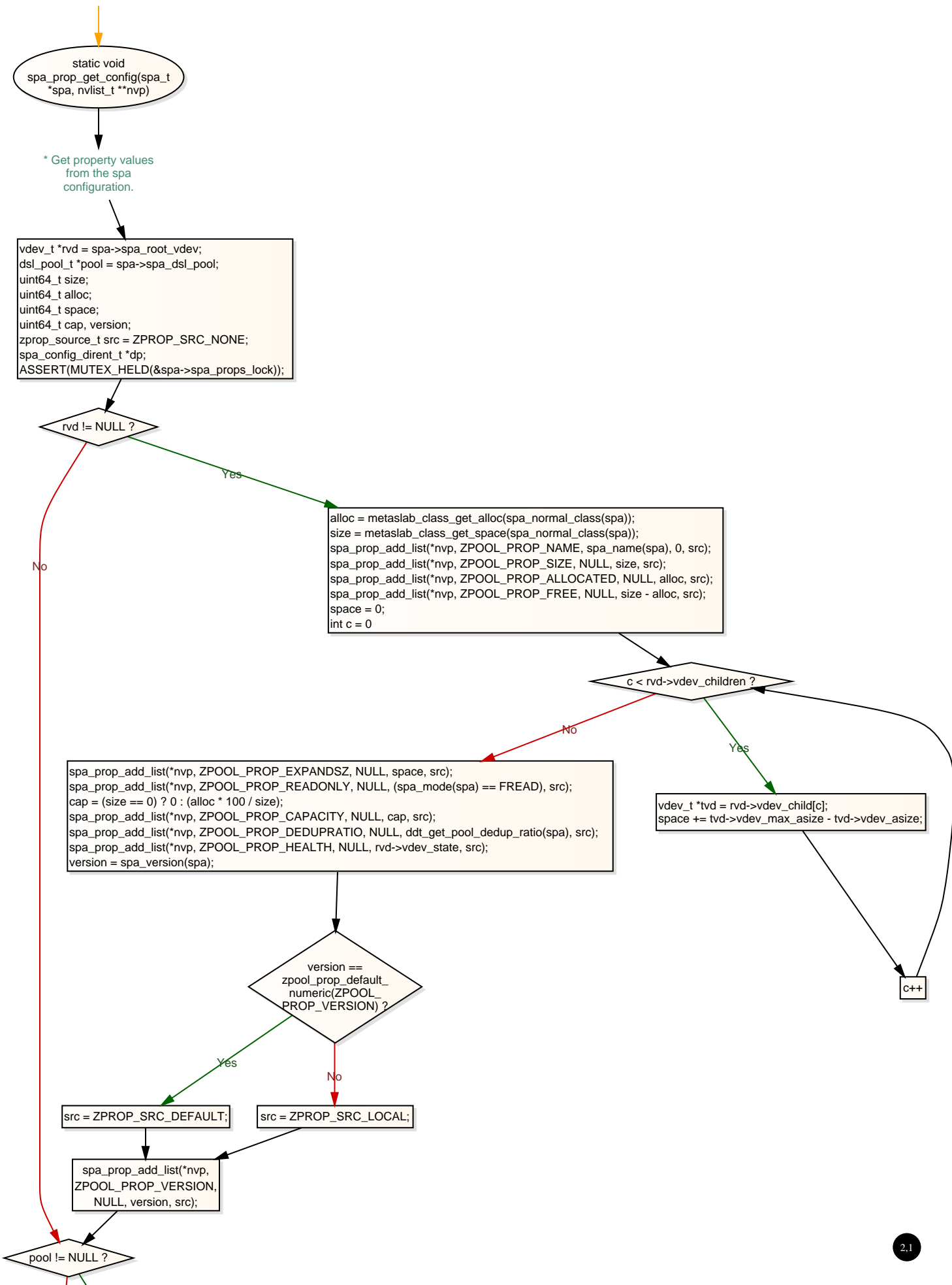
```
#ifdef _KERNEL ?
```

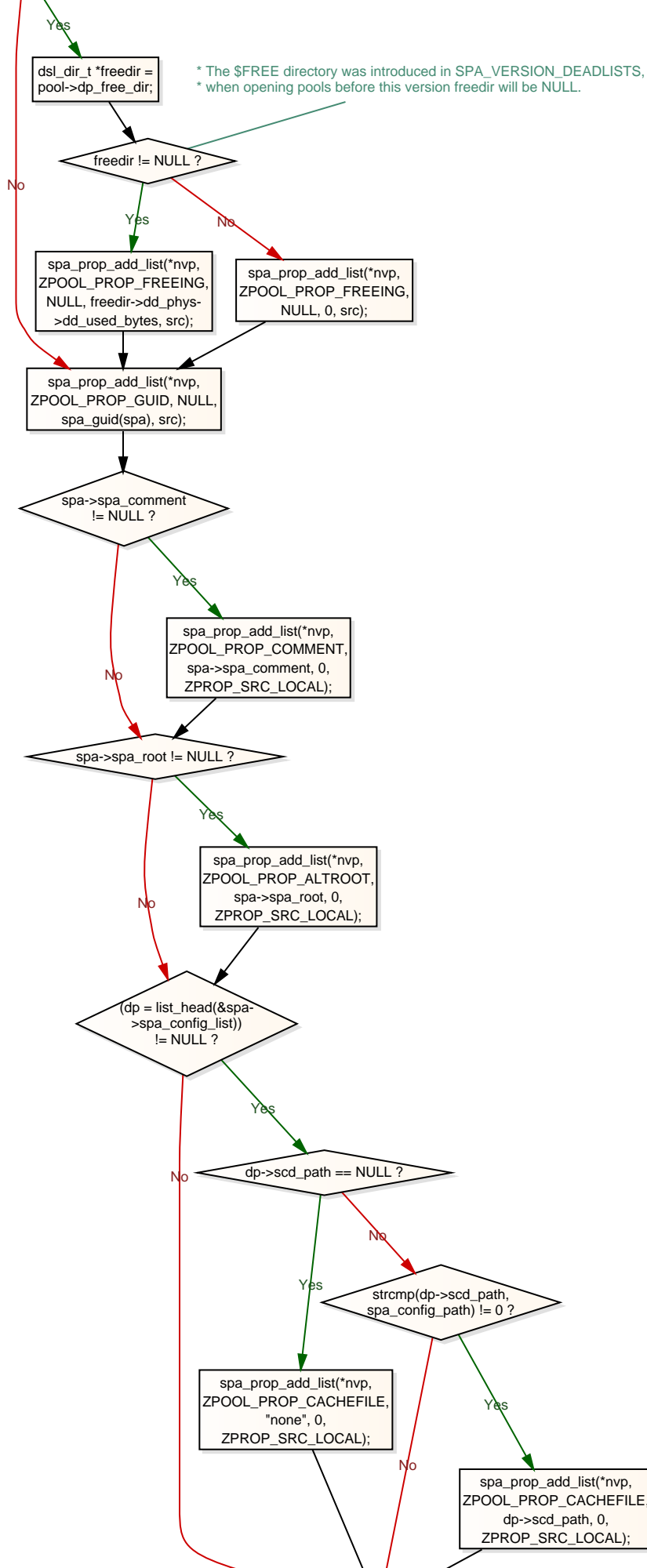
Yes

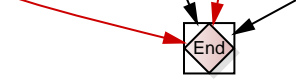
```
extern int  
vdev_disk_read_rootlabel  
(char *, char  
*, nvlist_t **);
```

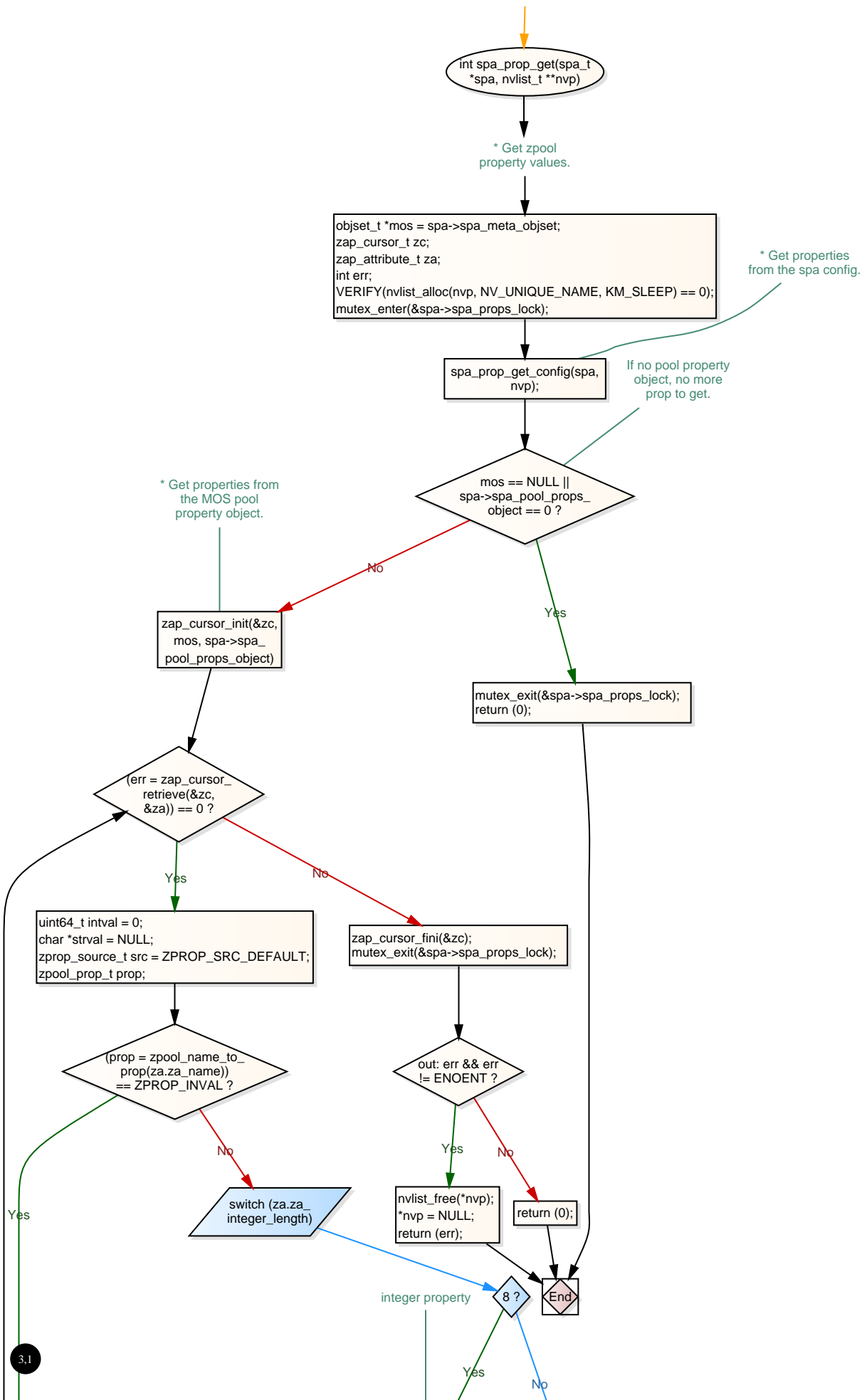
No

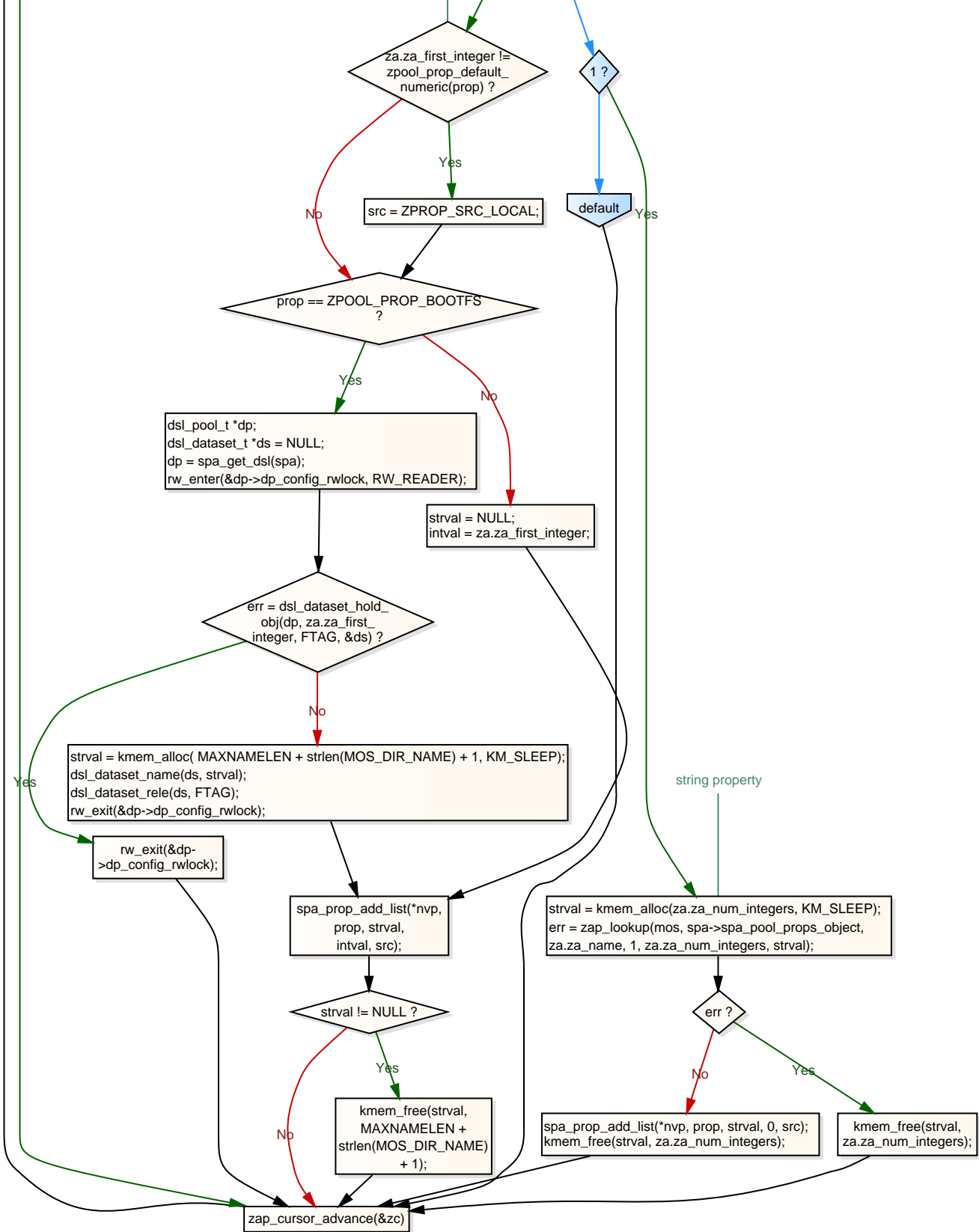






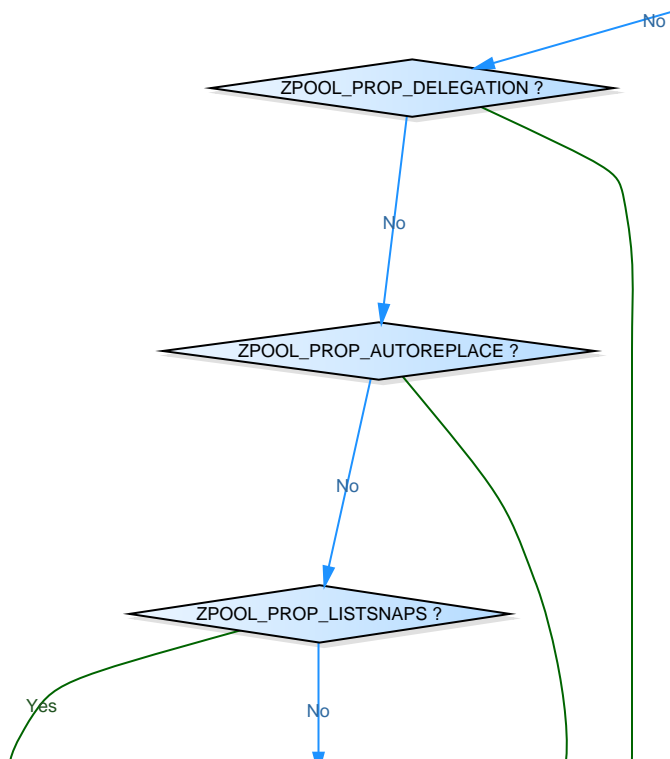


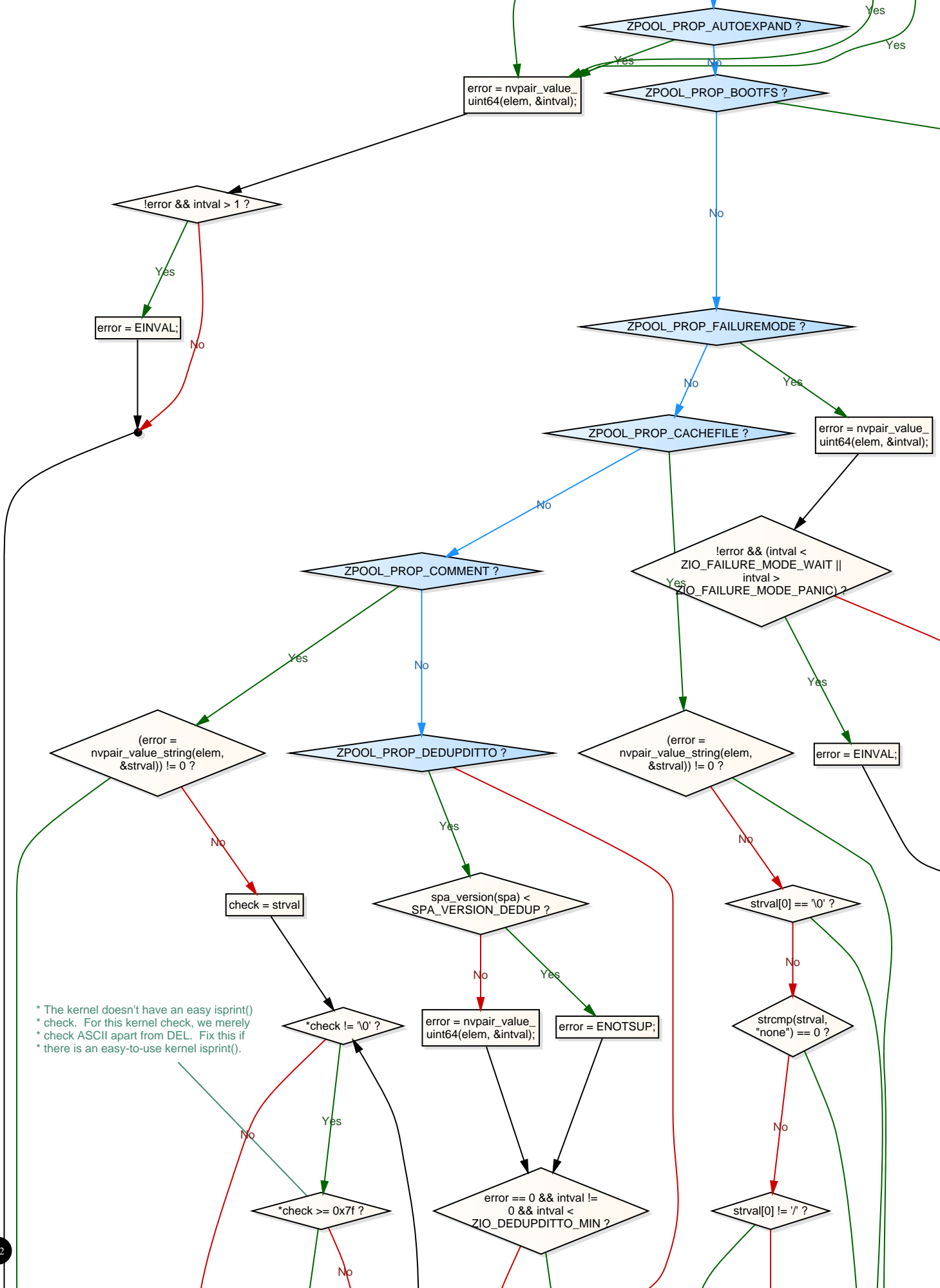


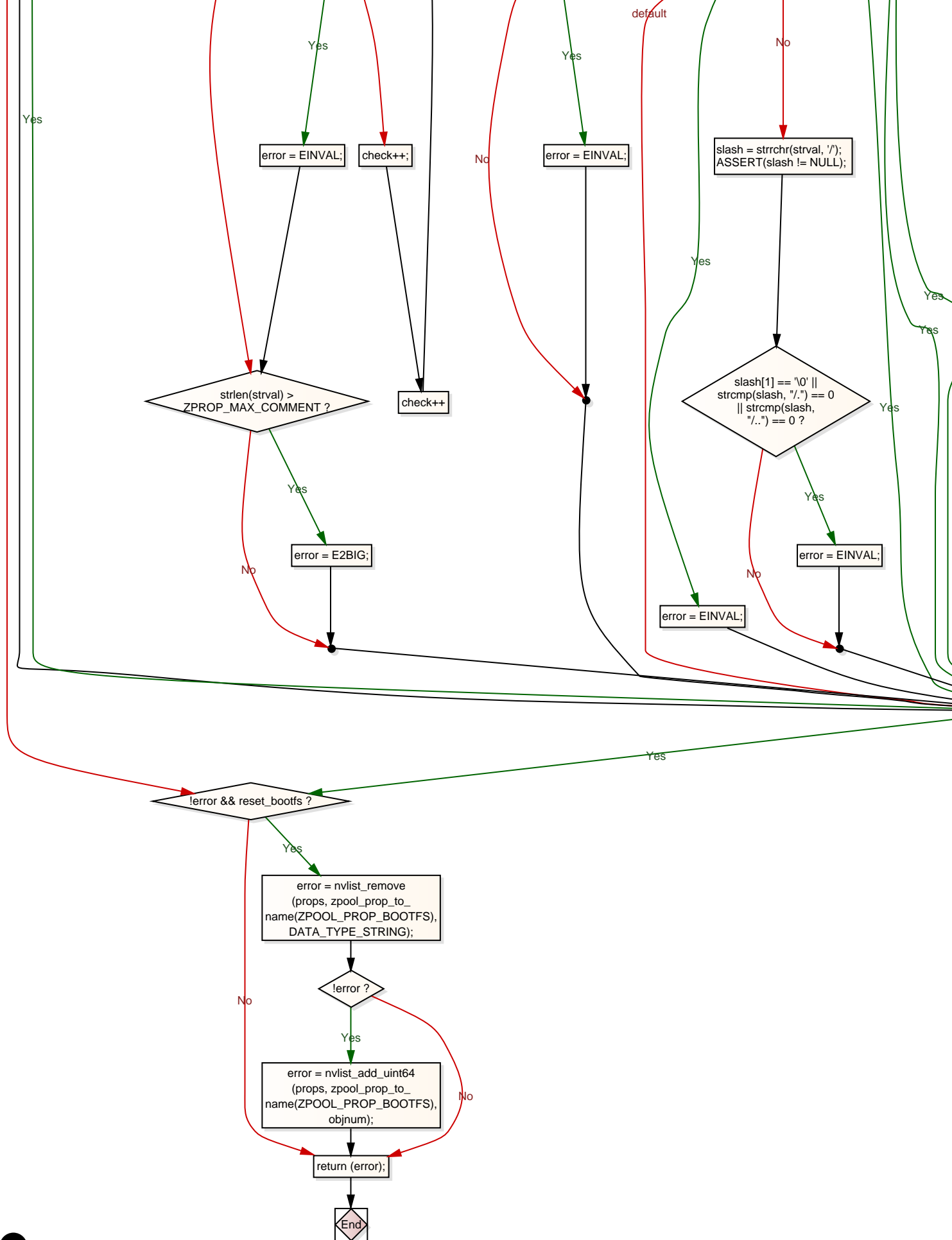


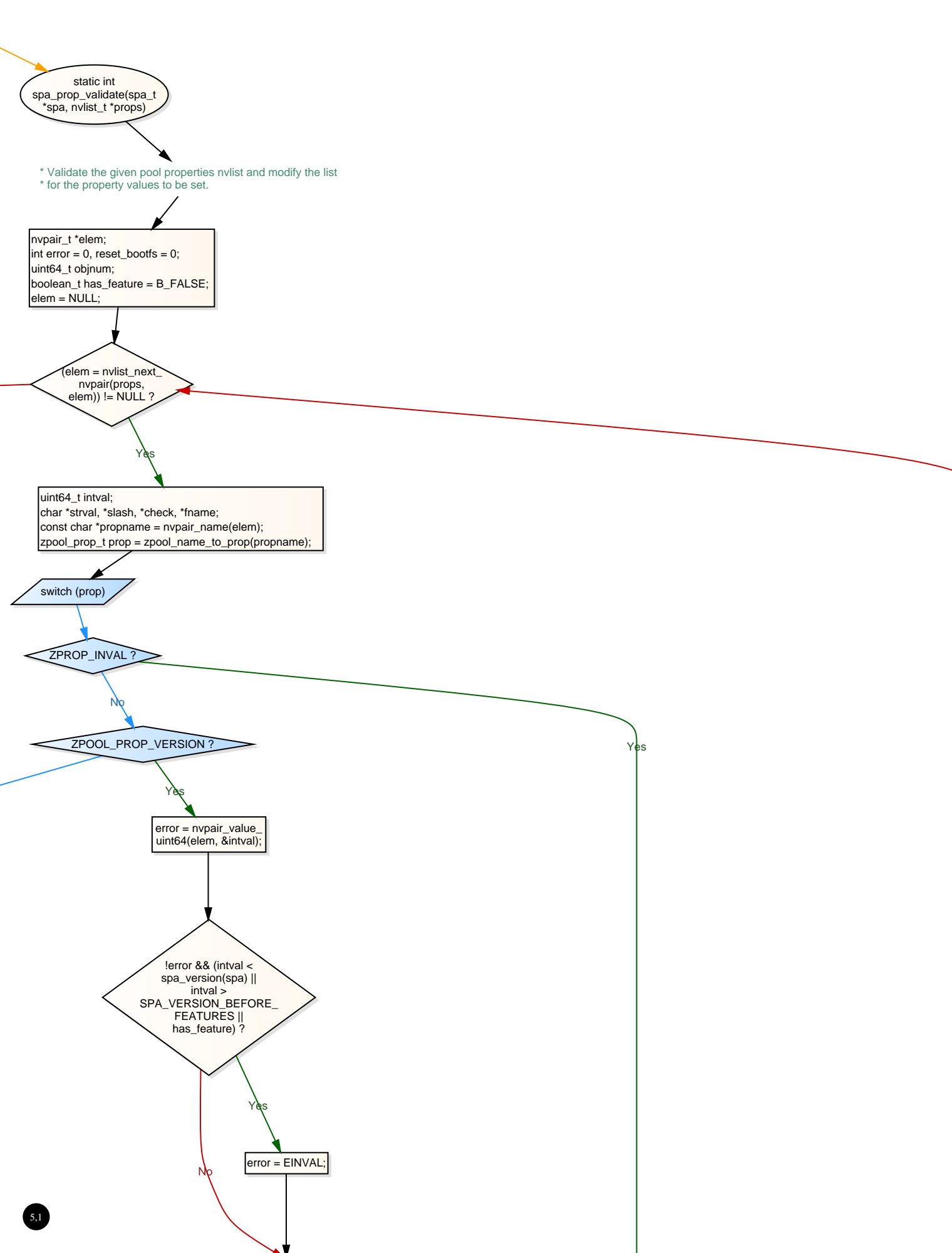


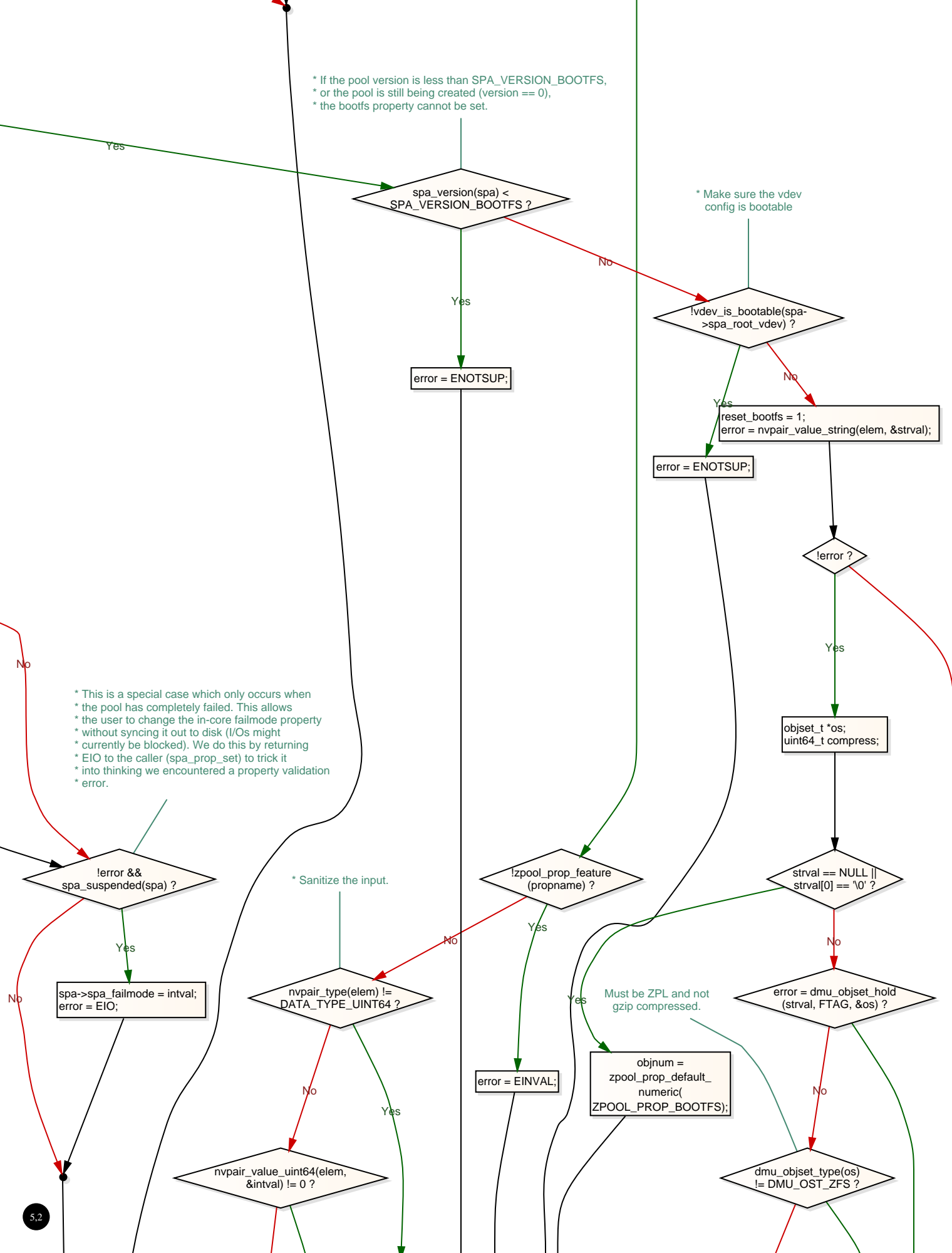
No

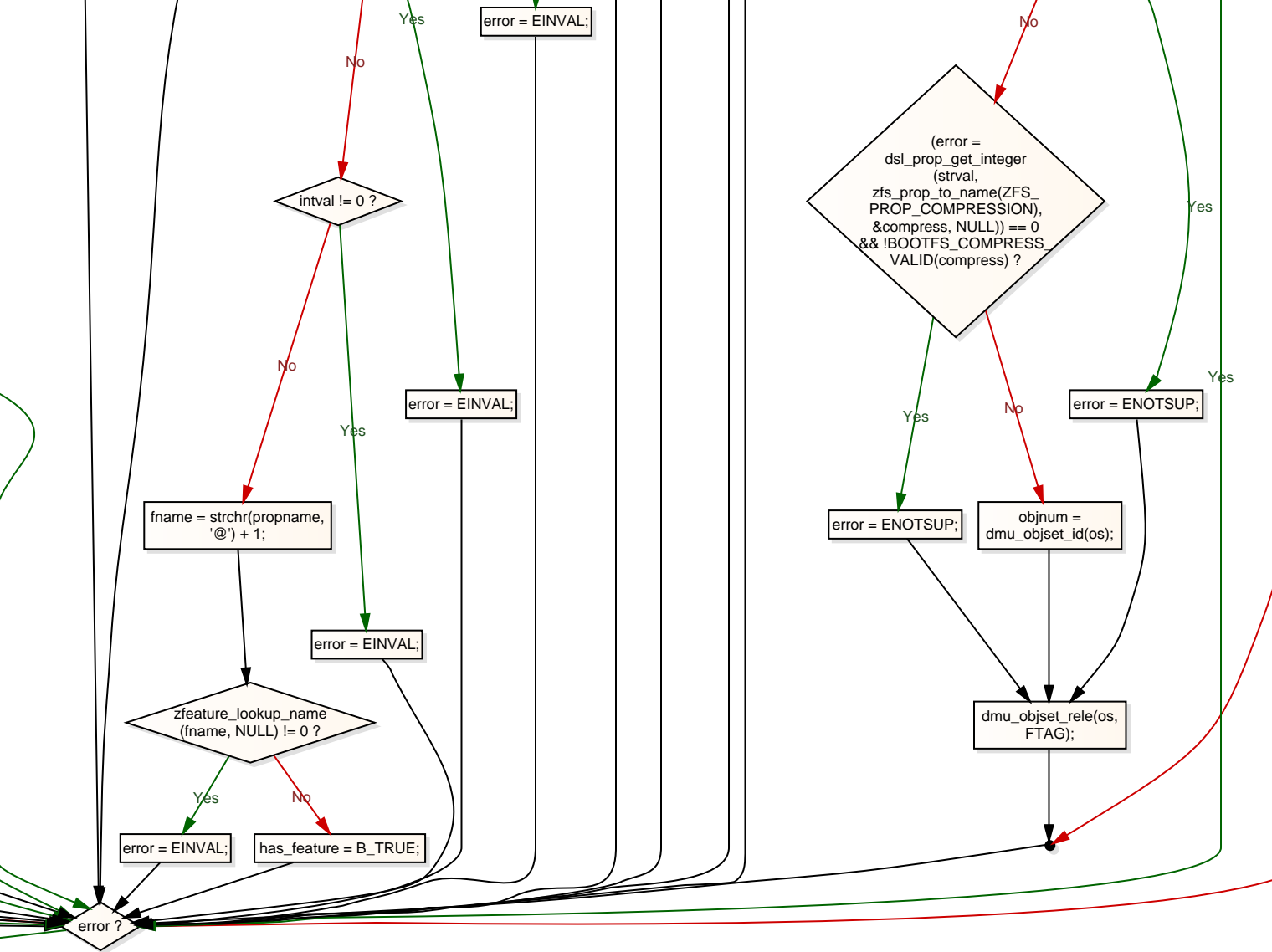


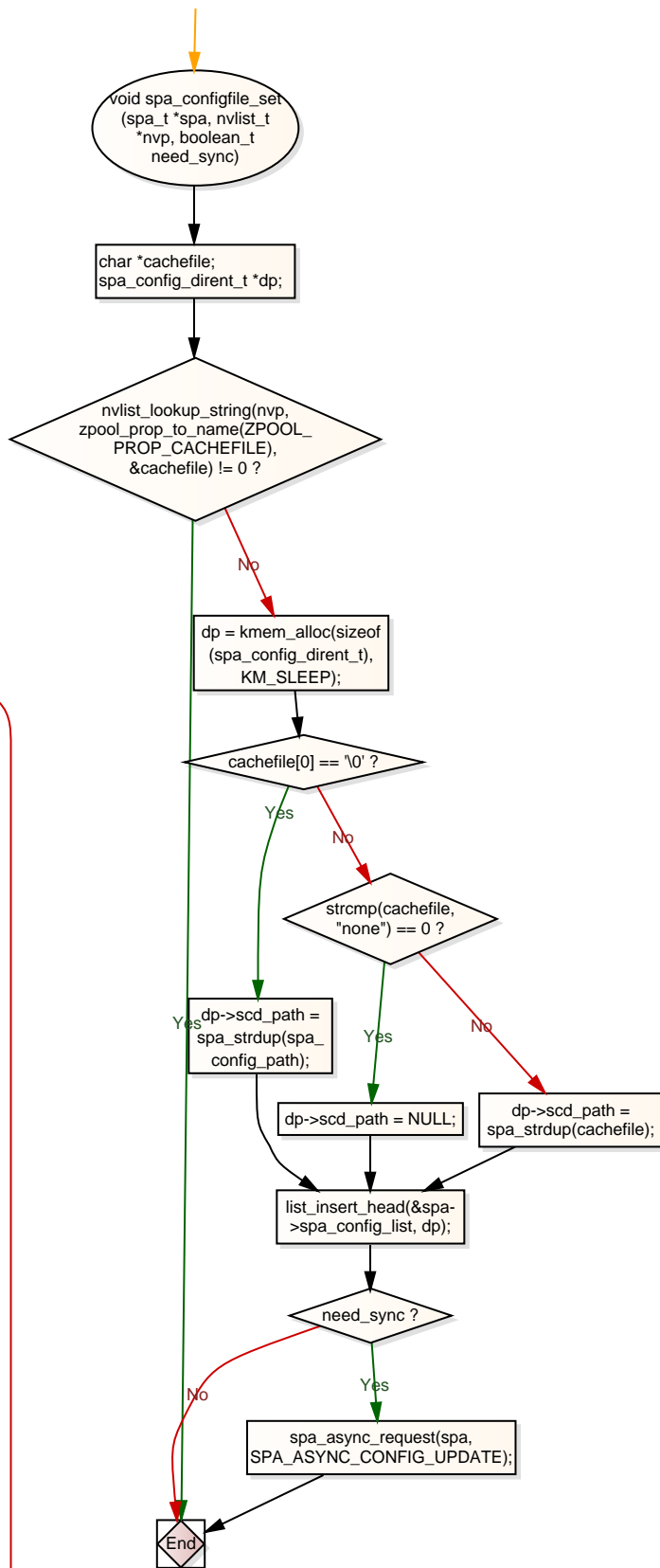










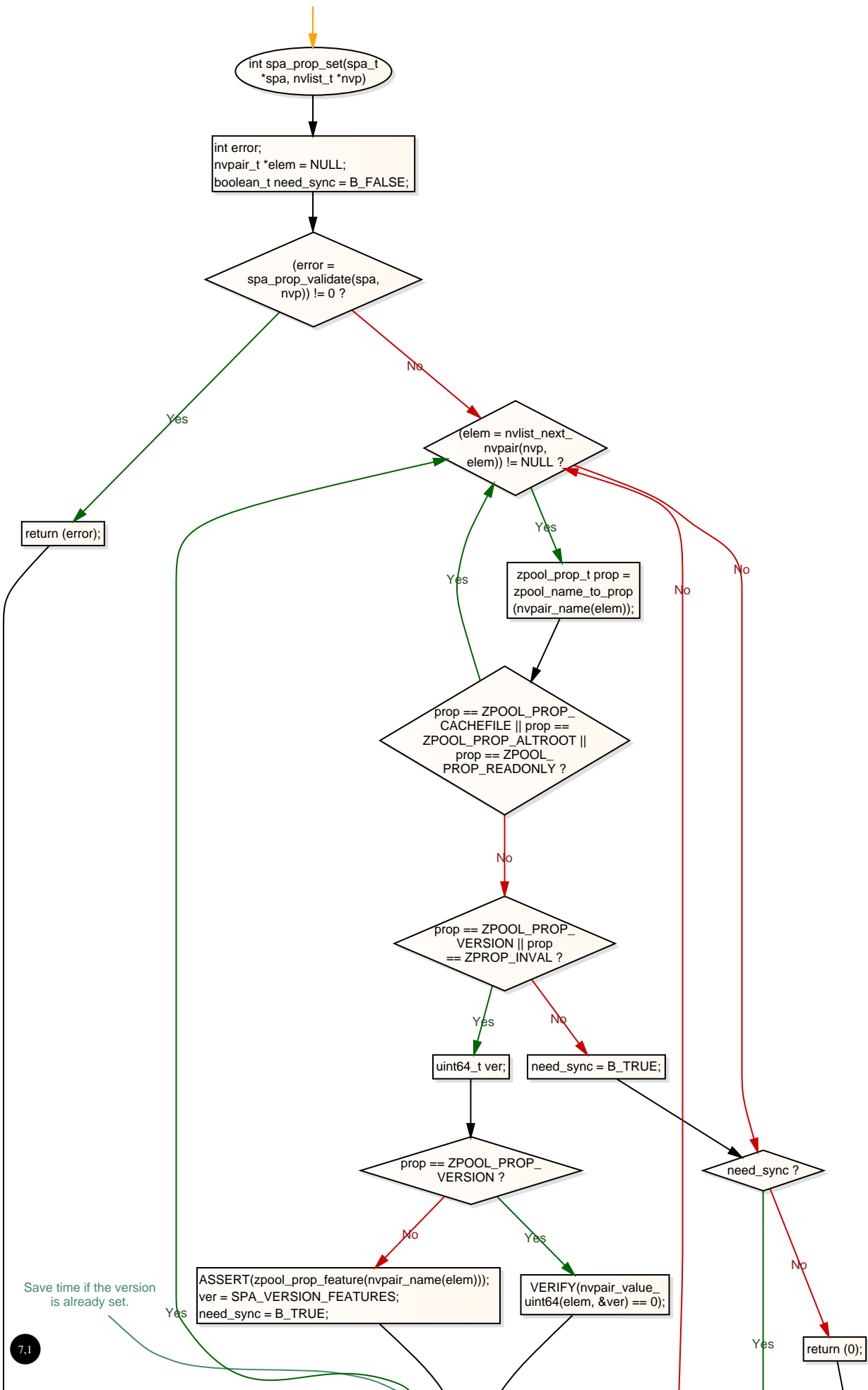




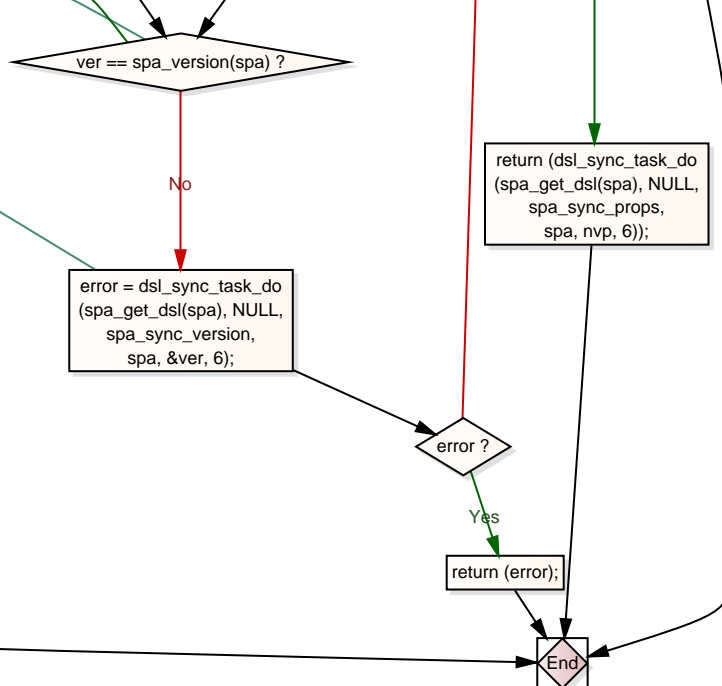


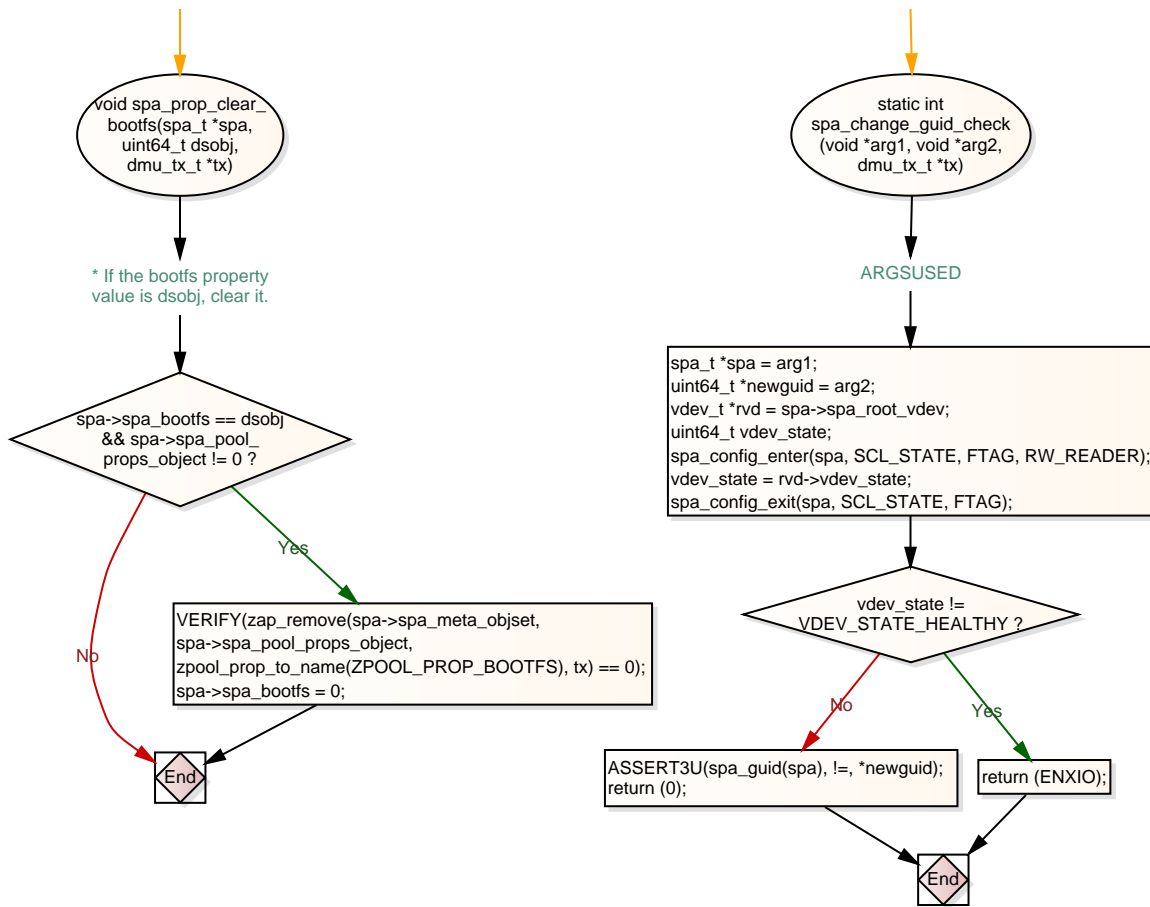


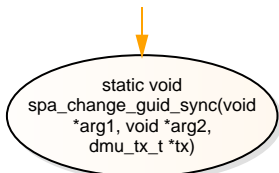
No



- \* In addition to the pool directory object, we might
- \* create the pool properties object, the features for
- \* read object, the features for write object, or the
- \* feature descriptions object.

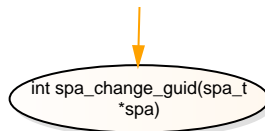






```

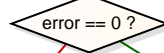
spa_t *spa = arg1;
uint64_t *newguid = arg2;
uint64_t oldguid;
vdev_t *rvd = spa->spa_root_vdev;
oldguid = spa_guid(spa);
spa_config_enter(spa, SCL_STATE, FTAG, RW_READER);
rvd->vdev_guid = *newguid;
rvd->vdev_guid_sum += (*newguid - oldguid);
vdev_config_dirty(rvd);
spa_config_exit(spa, SCL_STATE, FTAG);
spa_history_log_internal(spa, "guid change",
tx, "old=%lld new=%lld", oldguid, *newguid);
    
```



\* Change the GUID for the pool. This is done so that we can later  
 \* re-import a pool built from a clone of our own vdevs. We will modify  
 \* the root vdev's guid, our own pool guid, and then mark all of our  
 \* vdevs dirty. Note that we must make sure that all our vdevs are  
 \* online when we do this, or else any vdevs that weren't present  
 \* would be orphaned from our pool. We are also going to issue a  
 \* sysevent to update any watchers.

```

int error;
uint64_t guid;
mutex_enter(&spa_namespace_lock);
guid = spa_generate_guid(NULL);
error = dsl_sync_task_do(spa_get_dsl(spa),
spa_change_guid_check, spa_change_guid_sync, spa, &guid, 5);
    
```



Yes

No

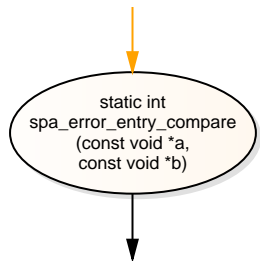
```

spa_config_sync(spa, B_FALSE, B_TRUE);
spa_event_notify(spa, NULL, ESC_ZFS_POOL_REGUID);
    
```

```

mutex_exit(&spa_namespace_lock);
return (error);
    
```

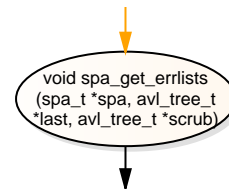
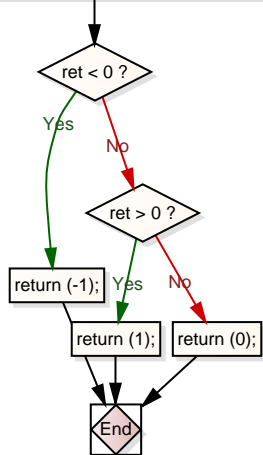




\* =====  
 \* SPA state manipulation (open/create/destroy/import/export)  
 \* =====

```

spa_error_entry_t *sa = (spa_error_entry_t *)a;
spa_error_entry_t *sb = (spa_error_entry_t *)b;
int ret;
ret = bcmp(&sa->se_bookmark, &sb->se_bookmark, sizeof (zbookmark_t));
  
```

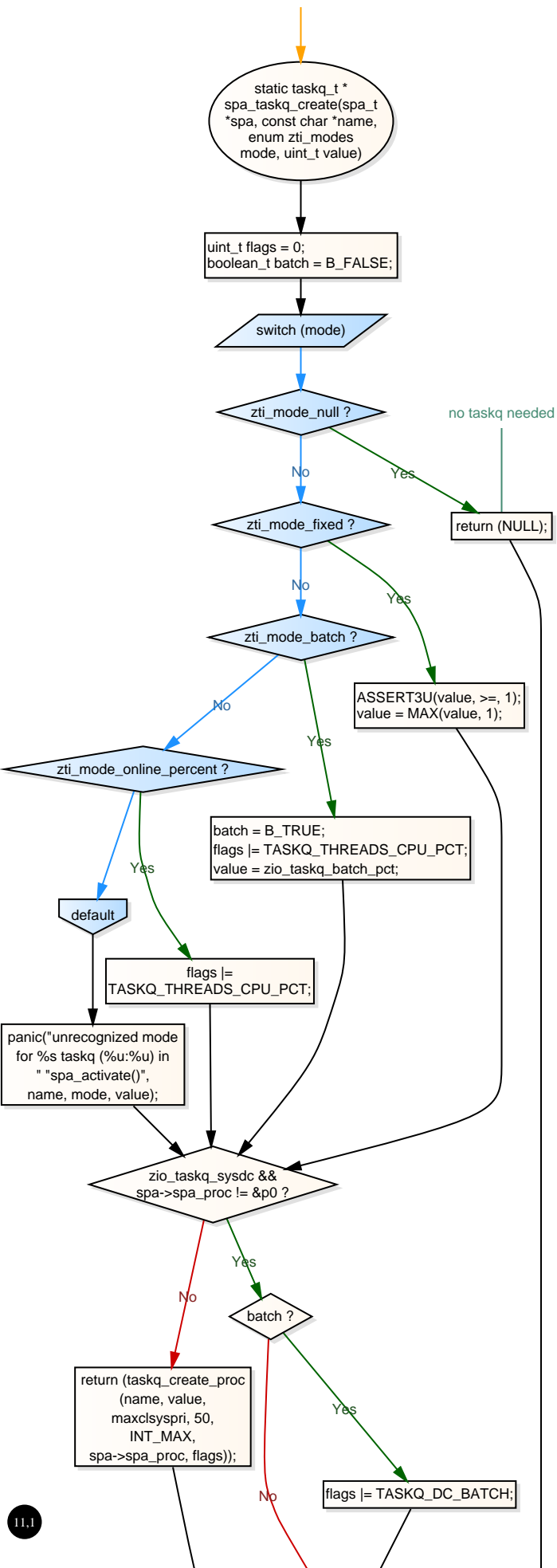


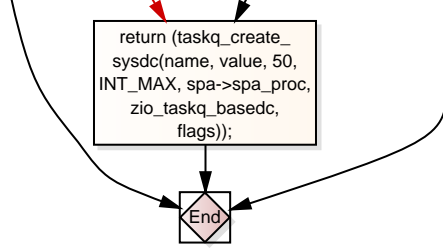
\* Utility function which retrieves copies of the current logs and  
 \* re-initializes them in the process.

```

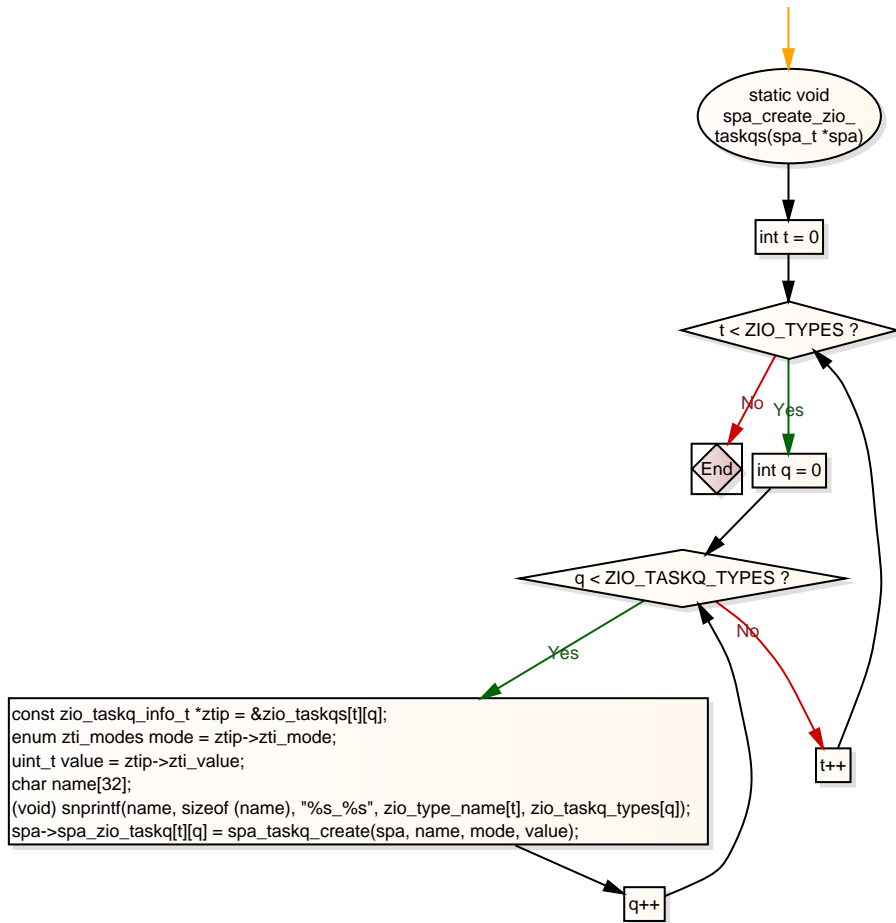
ASSERT(MUTEX_HELD(&spa->spa_errlist_lock));
bcopy(&spa->spa_errlist_last, last, sizeof (avl_tree_t));
bcopy(&spa->spa_errlist_scrub, scrub, sizeof (avl_tree_t));
avl_create(&spa->spa_errlist_scrub, spa_error_entry_compare,
sizeof (spa_error_entry_t), offsetof(spa_error_entry_t, se_avl));
avl_create(&spa->spa_errlist_last, spa_error_entry_compare,
sizeof (spa_error_entry_t), offsetof(spa_error_entry_t, se_avl));
  
```

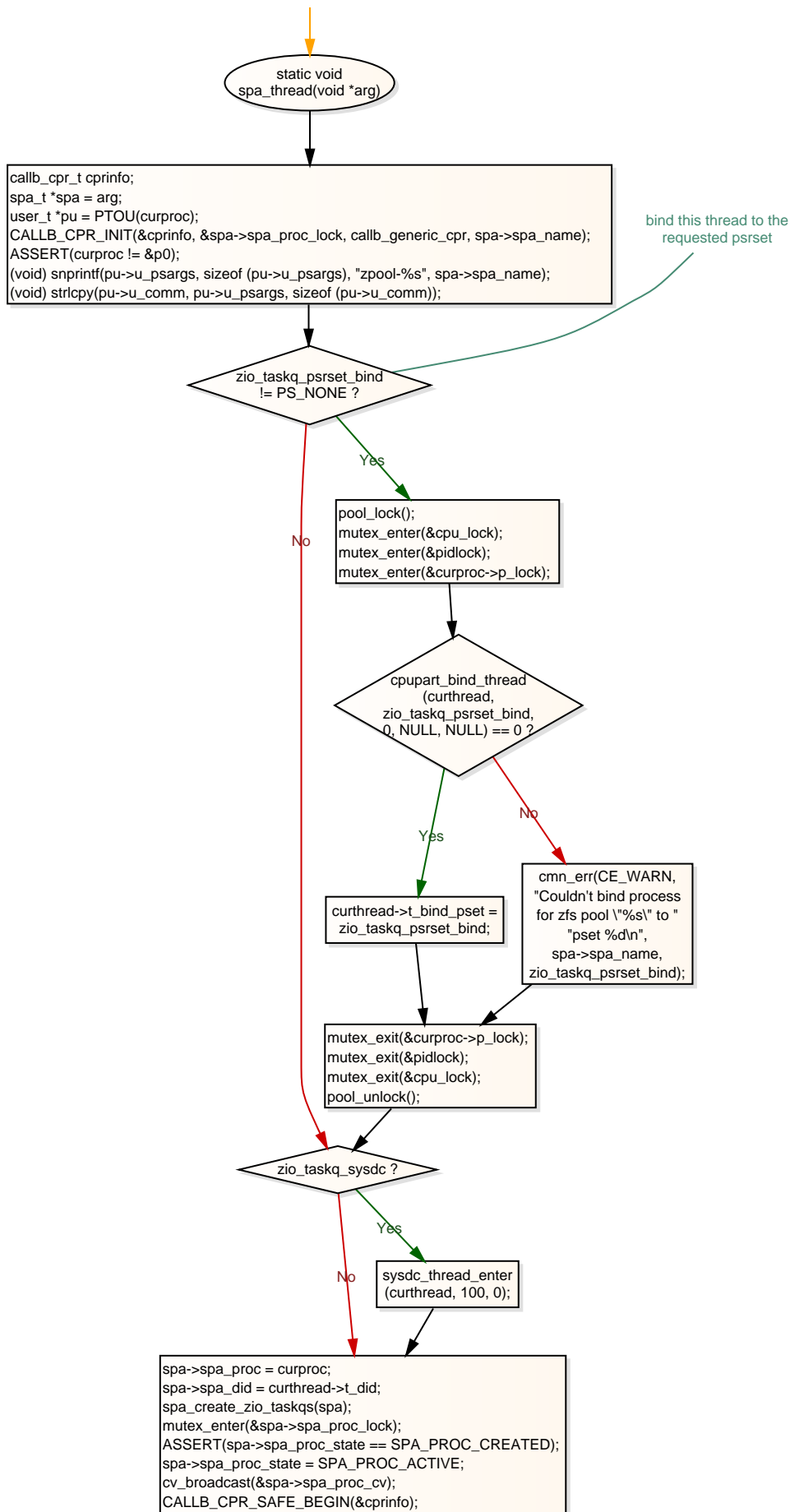


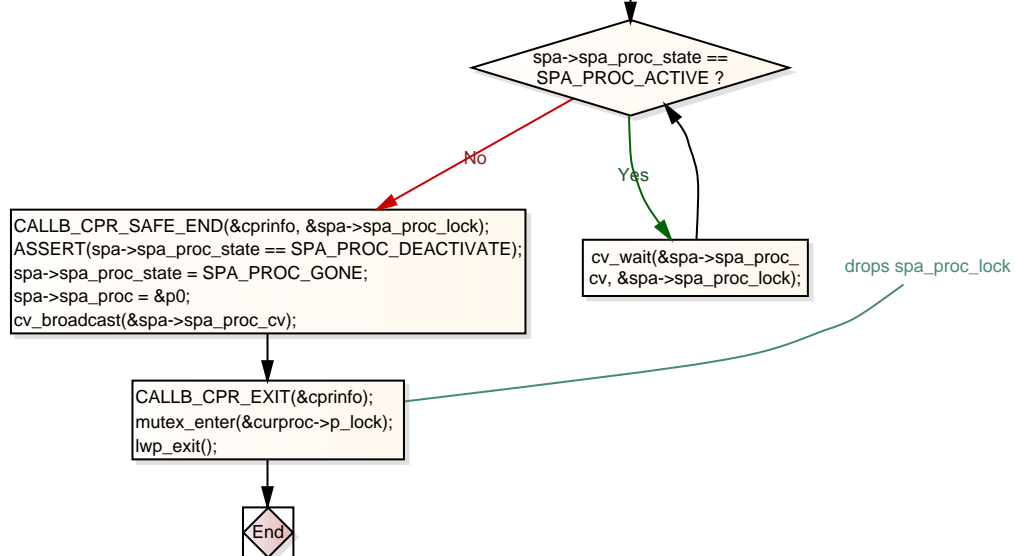


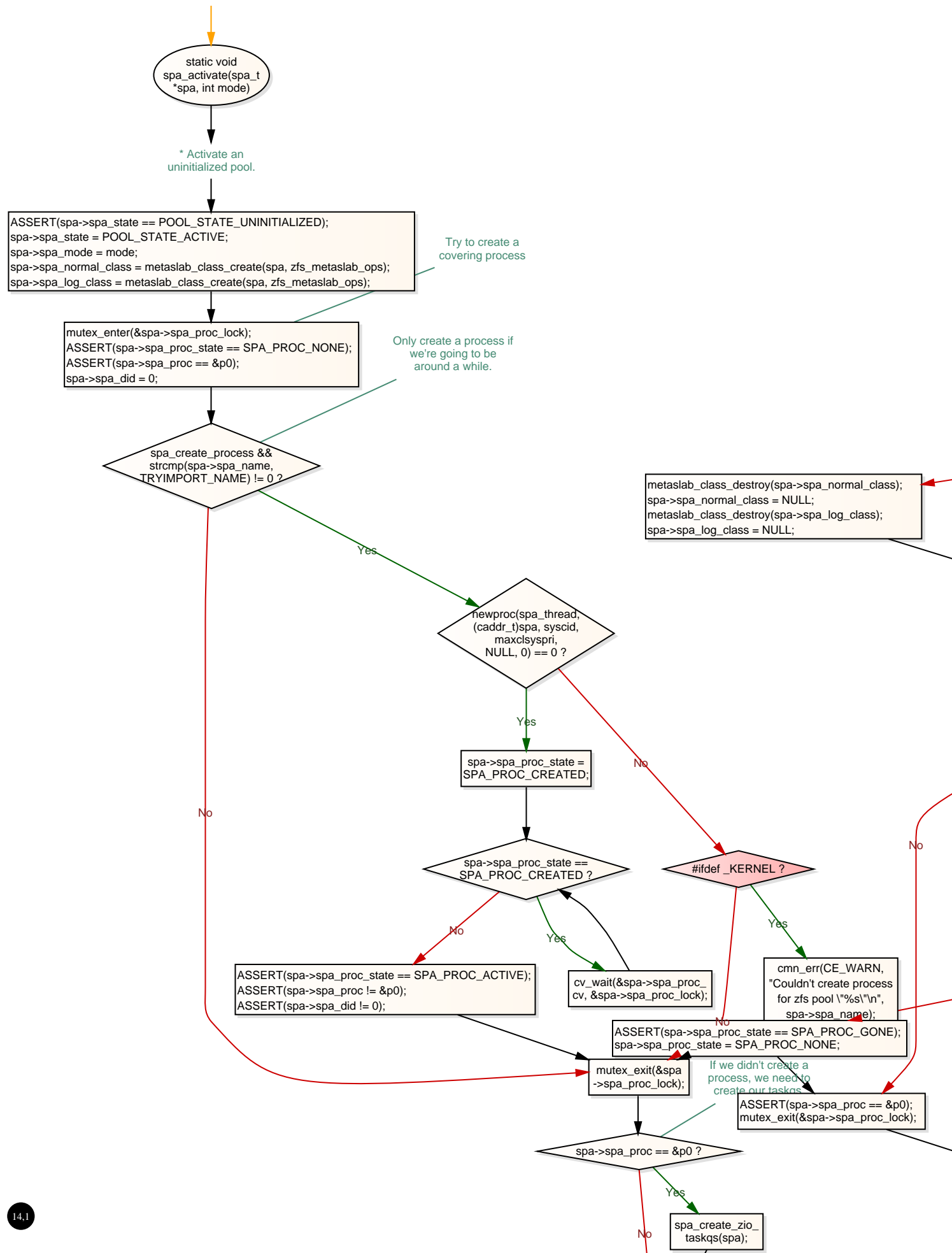








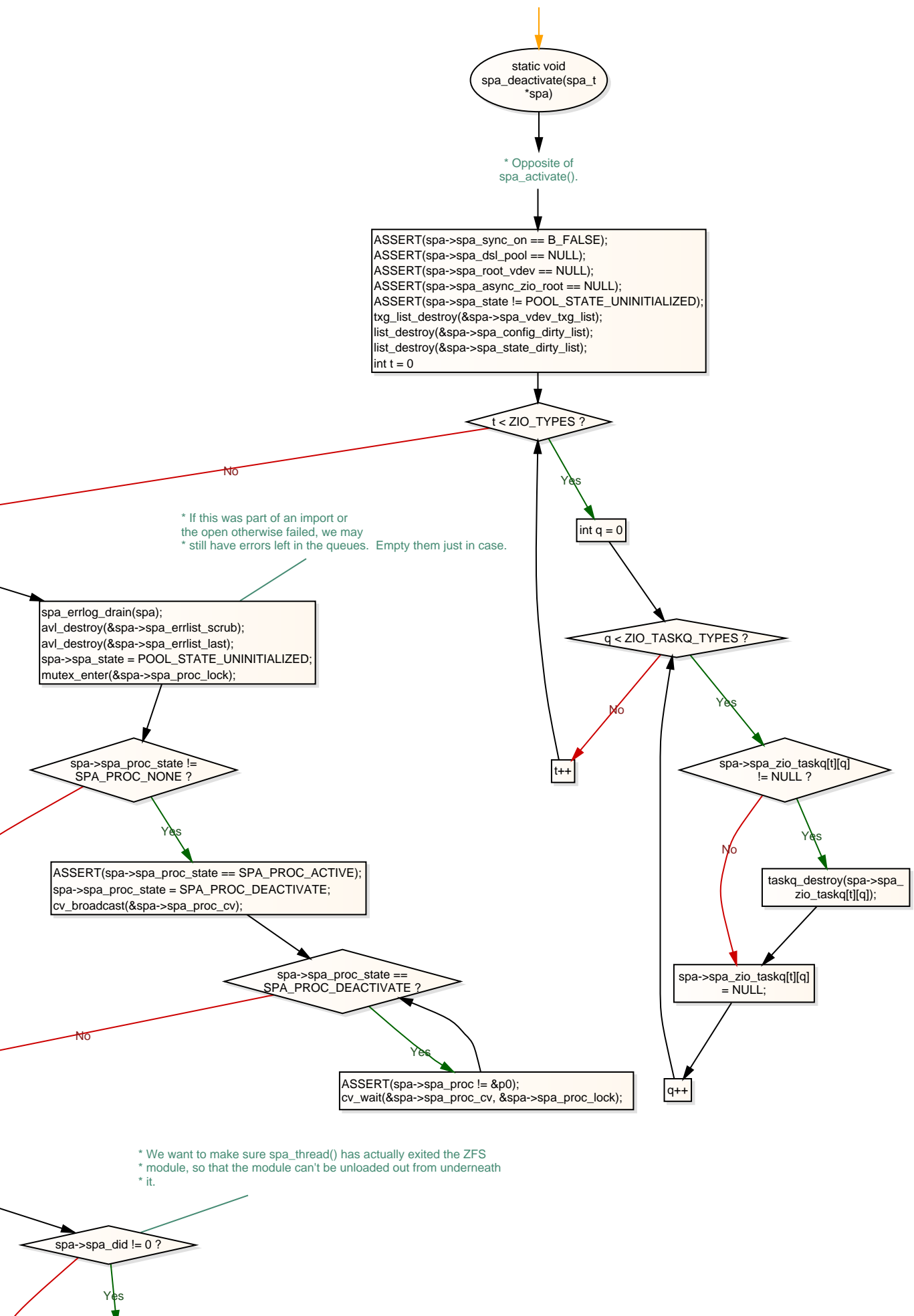


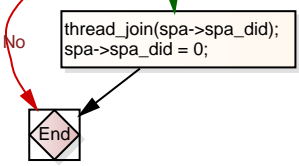


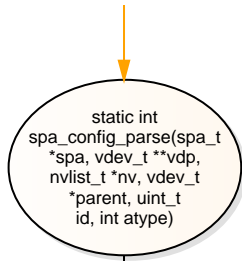
```
graph TD; Entry(( )) --> CodeBlock[Code Block]; CodeBlock --> End{End};
```

list\_create(&spa->spa\_config\_dirty\_list, sizeof (vdev\_t), offsetof(vdev\_t, vdev\_config\_dirty\_node));  
list\_create(&spa->spa\_state\_dirty\_list, sizeof (vdev\_t), offsetof(vdev\_t, vdev\_state\_dirty\_node));  
txg\_list\_create(&spa->spa\_vdev\_txg\_list, offsetof(struct vdev, vdev\_txg\_node));  
avl\_create(&spa->spa\_errlist\_scrub, spa\_error\_entry\_compare,  
sizeof (spa\_error\_entry\_t), offsetof(spa\_error\_entry\_t, se\_avl));  
avl\_create(&spa->spa\_errlist\_last, spa\_error\_entry\_compare,  
sizeof (spa\_error\_entry\_t), offsetof(spa\_error\_entry\_t, se\_avl));

End

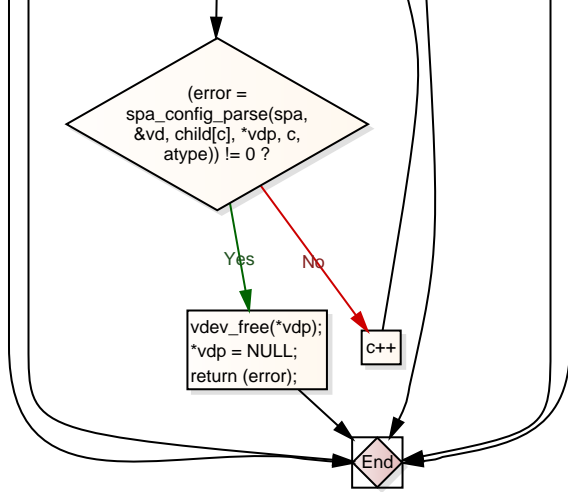


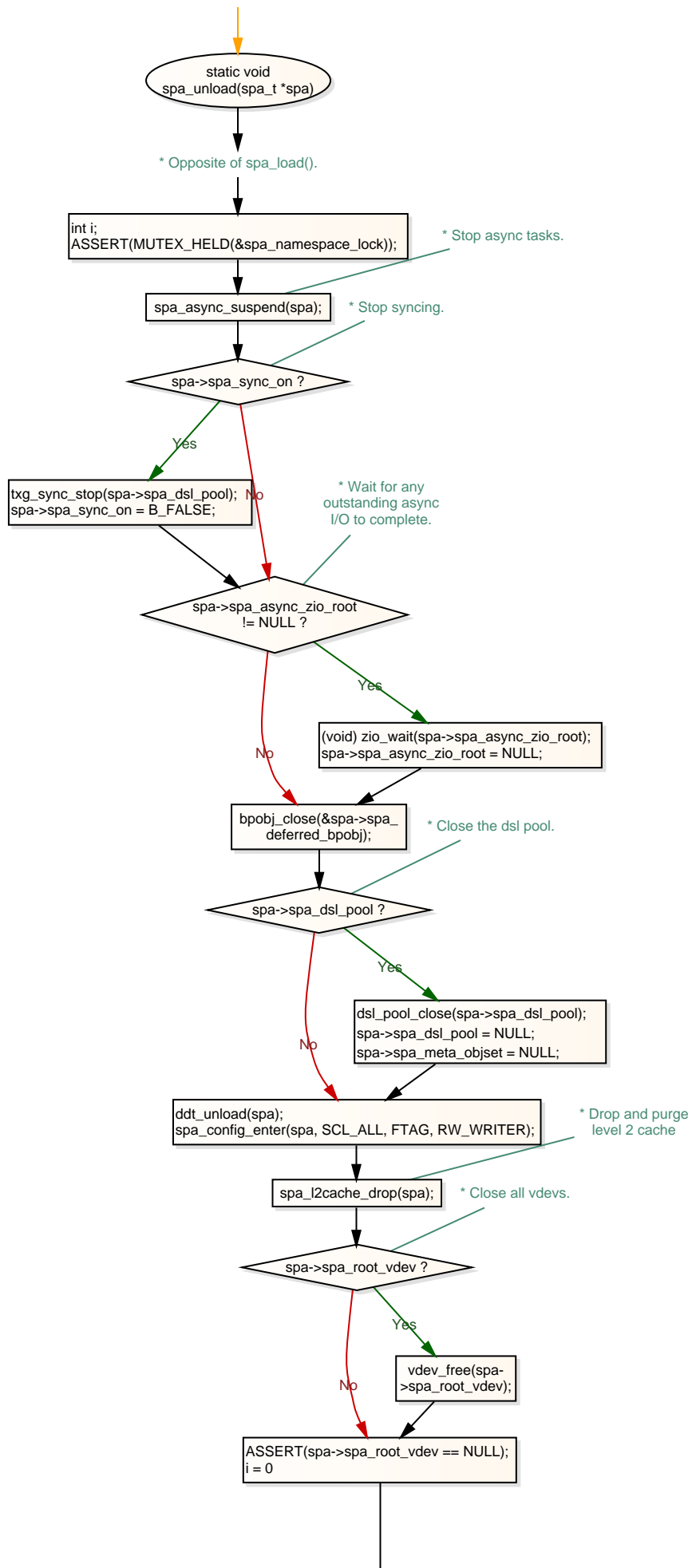


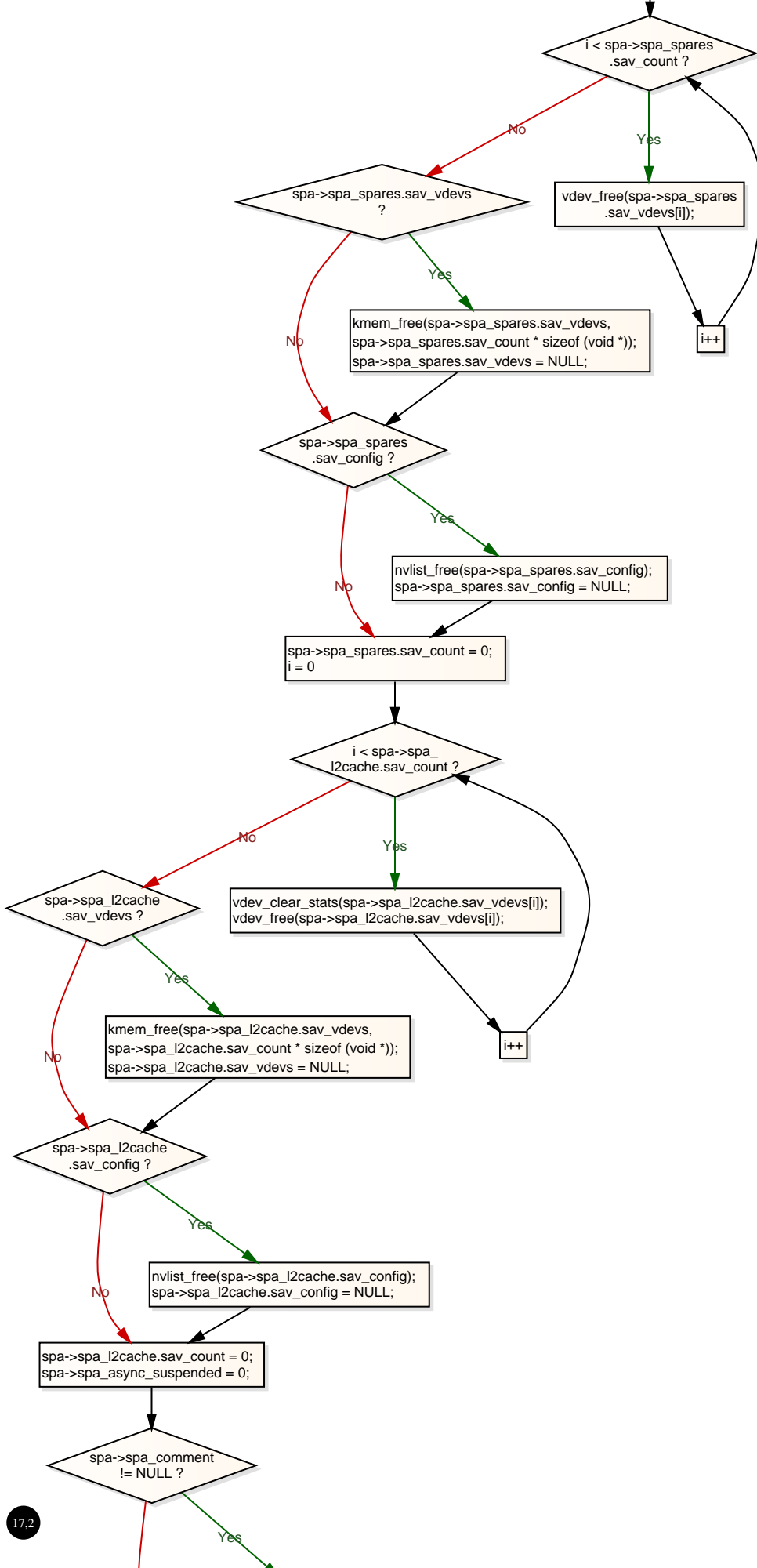


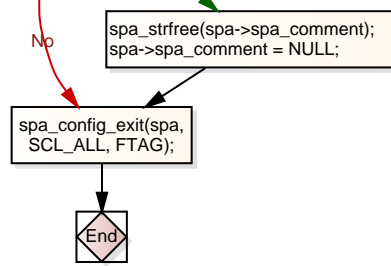
\* Verify a pool configuration, and construct the vdev tree appropriately. This  
 \* will create all the necessary vdevs in the appropriate layout, with each vdev  
 \* in the CLOSED state. This will prep the pool before open/creation/import.  
 \* All vdev validation is done by the vdev\_alloc() routine.

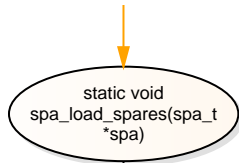












\* Load (or re-load) the current list of vdevs describing the active spares for this pool. When this is called, we have some form of basic information in \*spa\_spares.sav\_config'. We parse this into vdevs, try to open them, and \* then re-generate a more complete list including status information.

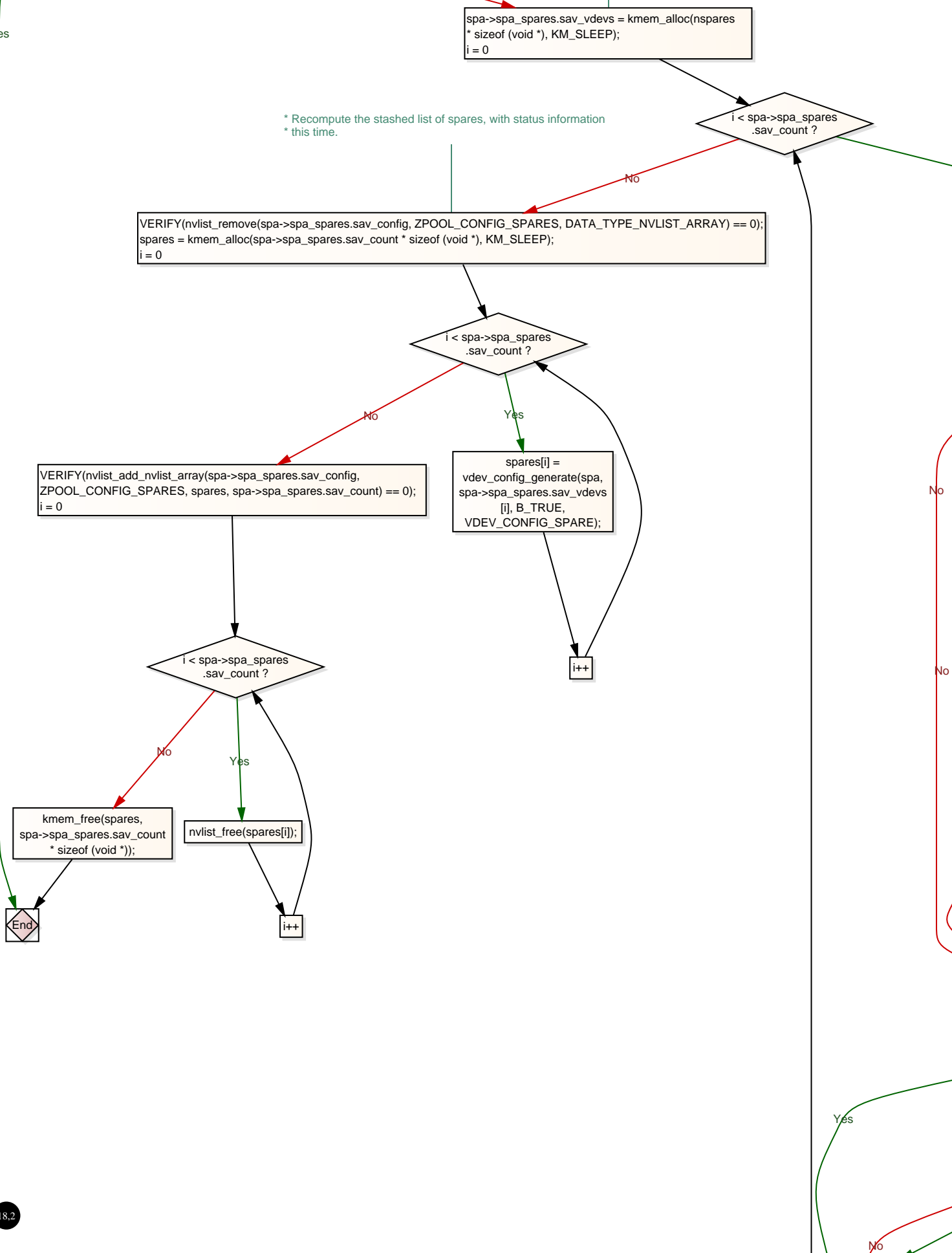
\* First, close and free any existing spare vdevs.

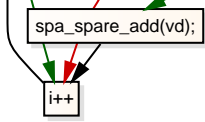
Undo the call to spa\_activate() below

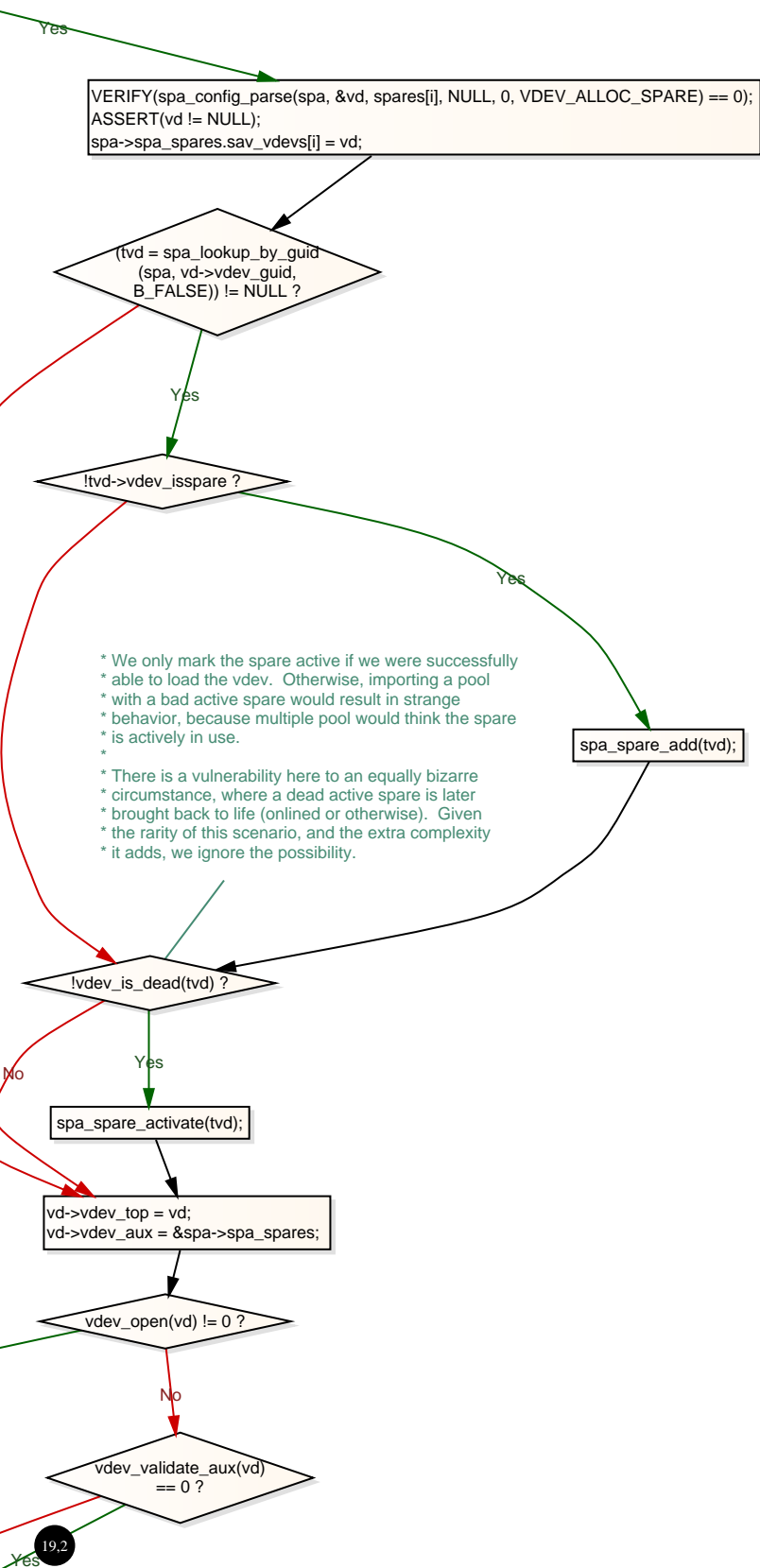
\* Construct the array of vdevs, opening them to get status in the process. For each spare, there is potentially two different vdev\_t structures associated with it: one in the list of spares (used only for basic validation purposes) and one in the active vdev configuration (if it's spared in). During this phase we open and validate each vdev on the spare list. If the vdev also exists in the active configuration, then we also mark this vdev as an active spare.

Yes

\* Recompute the stashed list of spares, with status information  
\* this time.











VERIFY(nvlist\_lookup\_nvlist\_array(sav->sav\_config,  
ZPOOL\_CONFIG\_L2CACHE, &l2cache, &nl2cache) == 0);  
newvdevs = kmem\_alloc(nl2cache \* sizeof(void \*), KM\_SLEEP);

oldvdevs = sav->sav\_vdevs;  
oldnvdevs = sav->sav\_count;  
sav->sav\_vdevs = NULL;  
sav->sav\_count = 0;

\* Purge vdevs  
that were dropped

i = 0

i < oldnvdevs ?

uint64\_t pool;  
vd = oldvdevs[i];

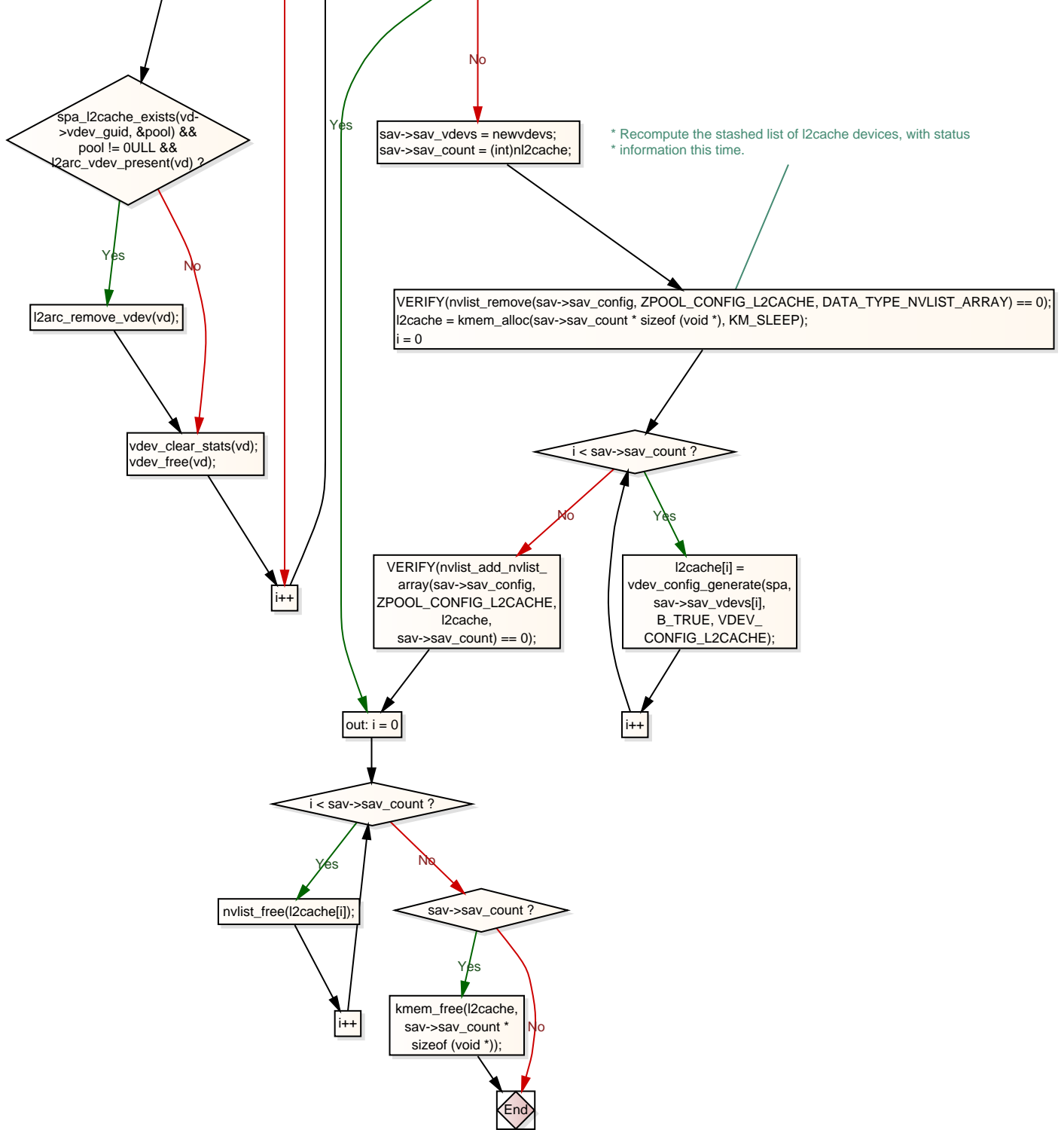
vd != NULL ?

ASSERT(vd->vdev\_  
isl2cache);

oldvdevs ?

kmem\_free(oldvdevs,  
oldnvdevs \*  
sizeof(void \*));

sav->sav\_config == NULL ?



```
static void
spa_load_l2cache(spa_t
*spa)
```

\* Load (or re-load) the current list of vdevs describing the active l2cache for this pool. When this is called, we have some form of basic information in \*spa\_l2cache.sav\_config\*. We parse this into vdevs, try to open them, and \* then re-generate a more complete list including status information.  
\* Devices which are already active have their details maintained, and are \* not re-opened.

```
nvlist_t **l2cache;
uint_t nl2cache;
int i, j, oldnvdevs;
uint64_t guid;
vdev_t *vd, **oldvdevs, **newvdevs;
spa_aux_vdev_t *sav = &spa->spa_l2cache;
ASSERT(spa_config_held(spa, SCL_ALL, RW_WRITER) == SCL_ALL);
```

sav->sav\_config != NULL ?

No

nl2cache = 0;

\* Process new nvlist of vdevs.

i = 0

i < nl2cache ?

Yes

```
VERIFY(nvlist_lookup_uint64(l2cache[i],
ZPOOL_CONFIG_GUID, &guid) == 0);
newvdevs[i] = NULL;
j = 0
```

j < oldnvdevs ?

Yes

vd = oldvdevs[j];

vd != NULL && guid == vd->vdev\_guid ?

\* Retain previous vdev for add/remove ops.

No

j++

```
newvdevs[i] = vd;
oldvdevs[j] = NULL;
```

```
static int
load_nvlist(spa_t *spa,
uint64_t obj,
nvlist_t **value)
```

```
dmu_buf_t *db;
char *packed = NULL;
size_t nvsize = 0;
int error;
*value = NULL;
VERIFY(0 == dmu_bonus_hold(spa->spa_meta_objset,
obj, FTAG, &db));
nvsize = *(uint64_t *)db->db_data;
dmu_buf_rele(db, FTAG);
packed = kmem_alloc(nvsize, KM_SLEEP);
error = dmu_read(spa->spa_meta_objset,
obj, 0, nvsize, packed, DMU_READ_PREFETCH);
```

error == 0 ?

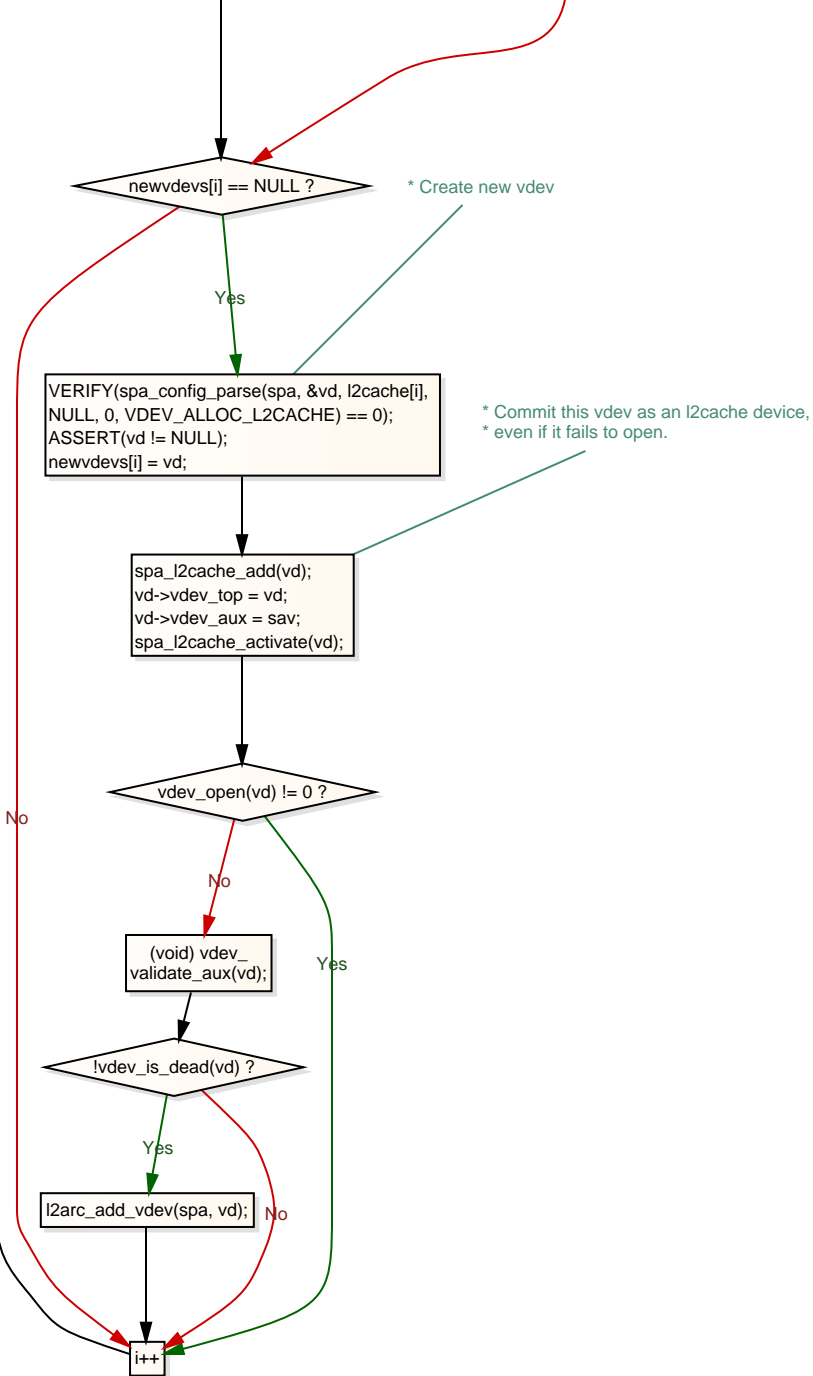
Yes

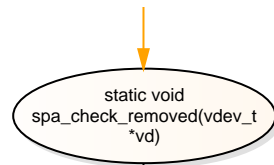
No

```
error = nvlist_unpack(
packed, nvsize,
value, 0);
```

```
kmem_free(packed, nvsize);
return (error);
```

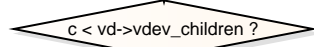
End





\* Checks to see if the given vdev could not be opened, in which case we post a sysevent to notify the autoreplace code that the device has been removed.

int c = 0

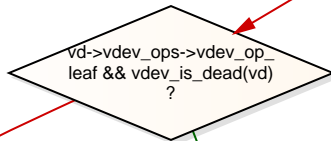


Yes

spa\_check\_removed(vd->vdev\_child[c]);

c++

No

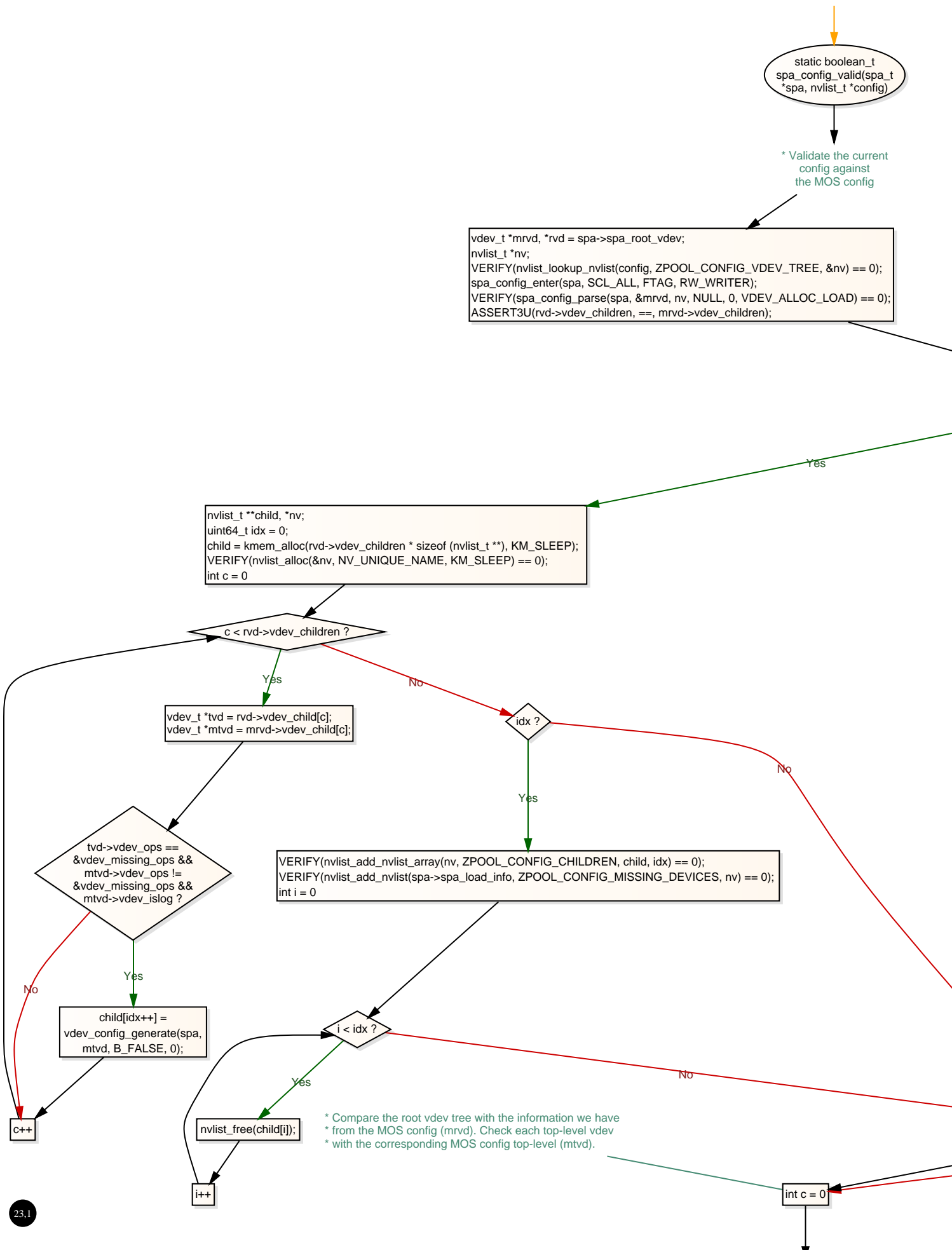


Yes

zfs\_post\_autoreplace(vd->vdev\_spa, vd);  
spa\_event\_notify(vd->vdev\_spa, vd, ESC\_ZFS\_VDEV\_CHECK);



No



\* Ensure we were able to validate the config.

```
vdev_free(mrvd);  
spa_config_exit(spa, SCL_ALL, FTAG);
```

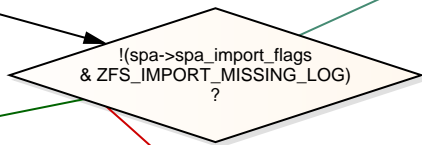
```
return (rvd->vdev_guid_  
sum == spa->spa_uberblock  
.ub_guid_sum);
```

End





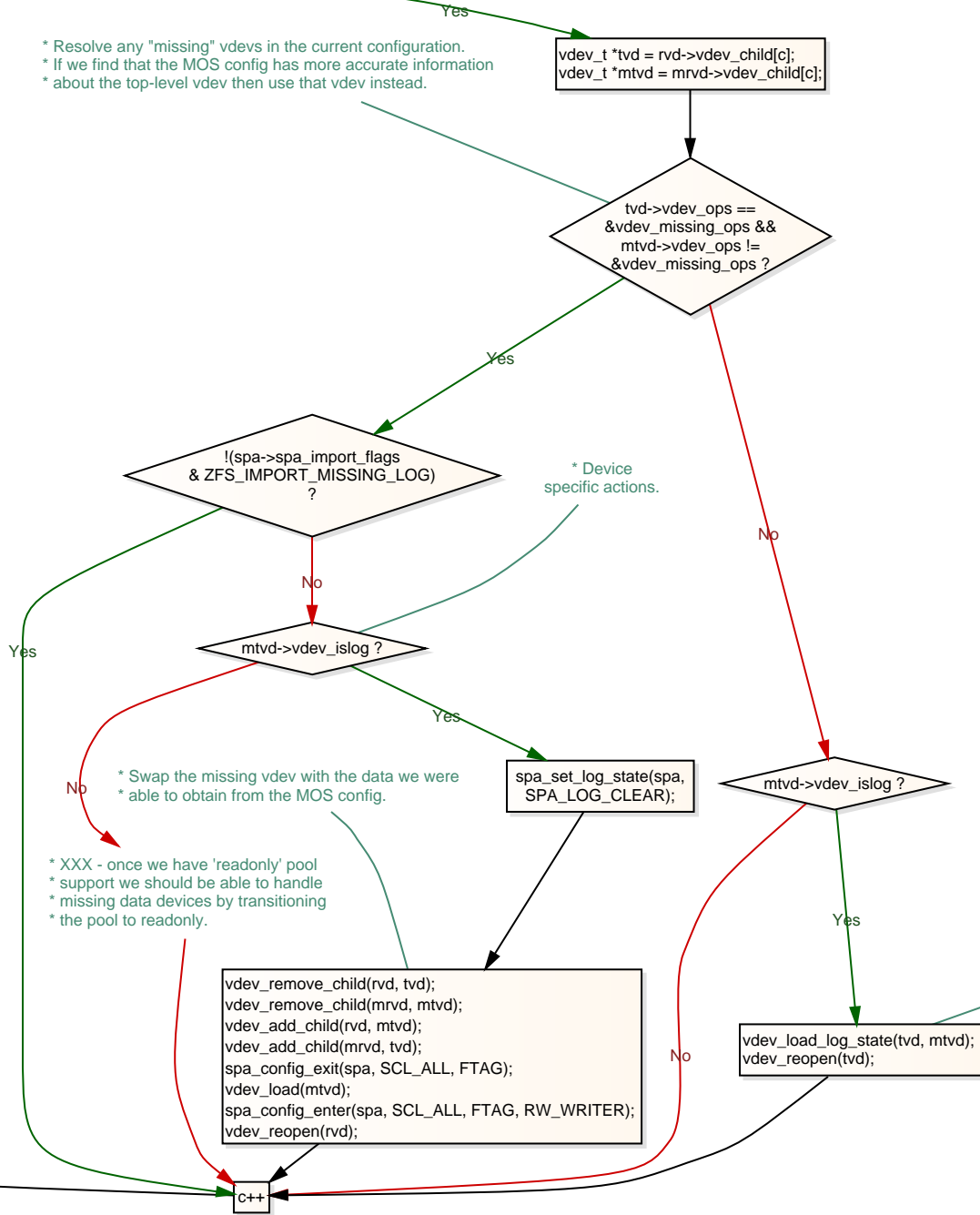
- \* If we're doing a normal import, then build up any additional
- \* diagnostic information about missing devices in this config.
- \* We'll pass this up to the user for further processing.

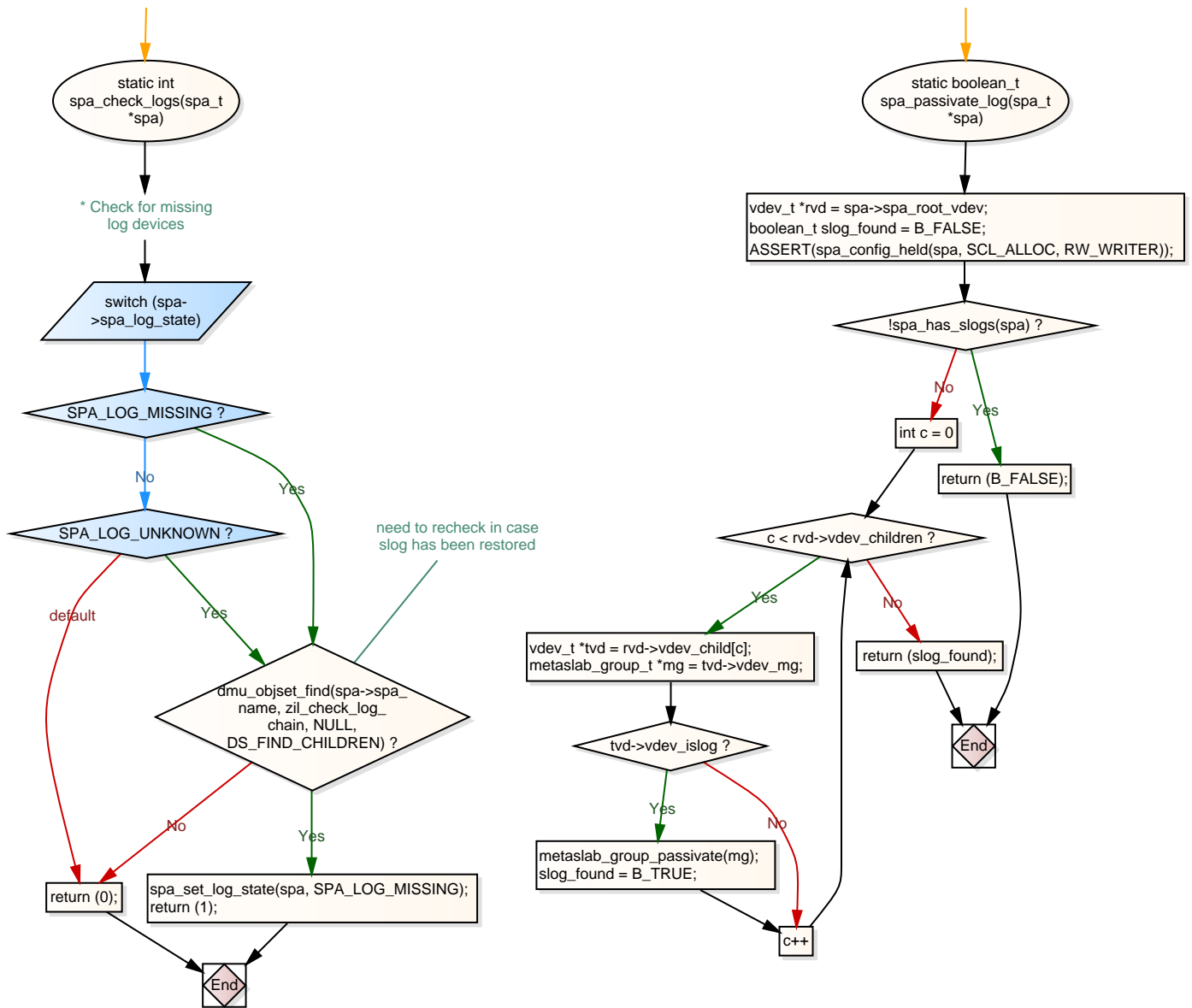



No

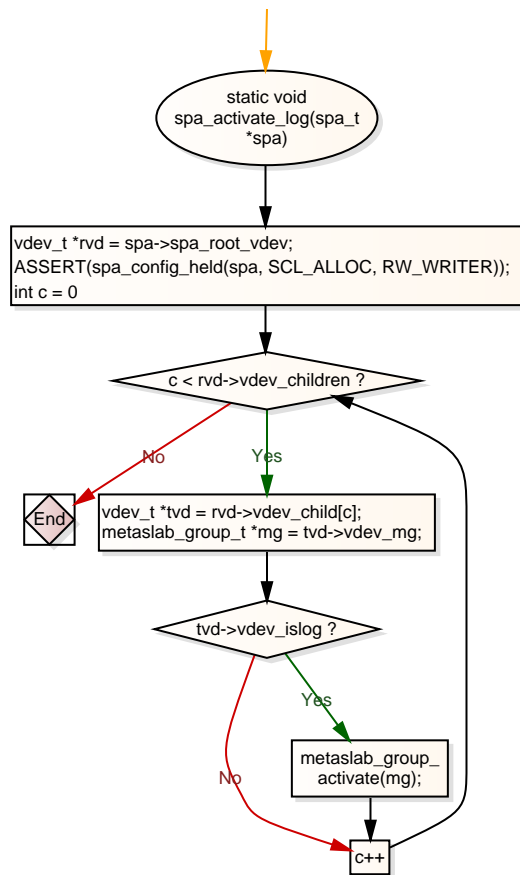
```
nvlist_free(nv);  
kmem_free(child, rvd->vdev_children * sizeof (char **));
```

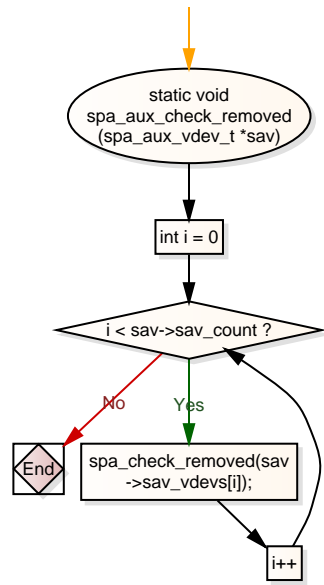
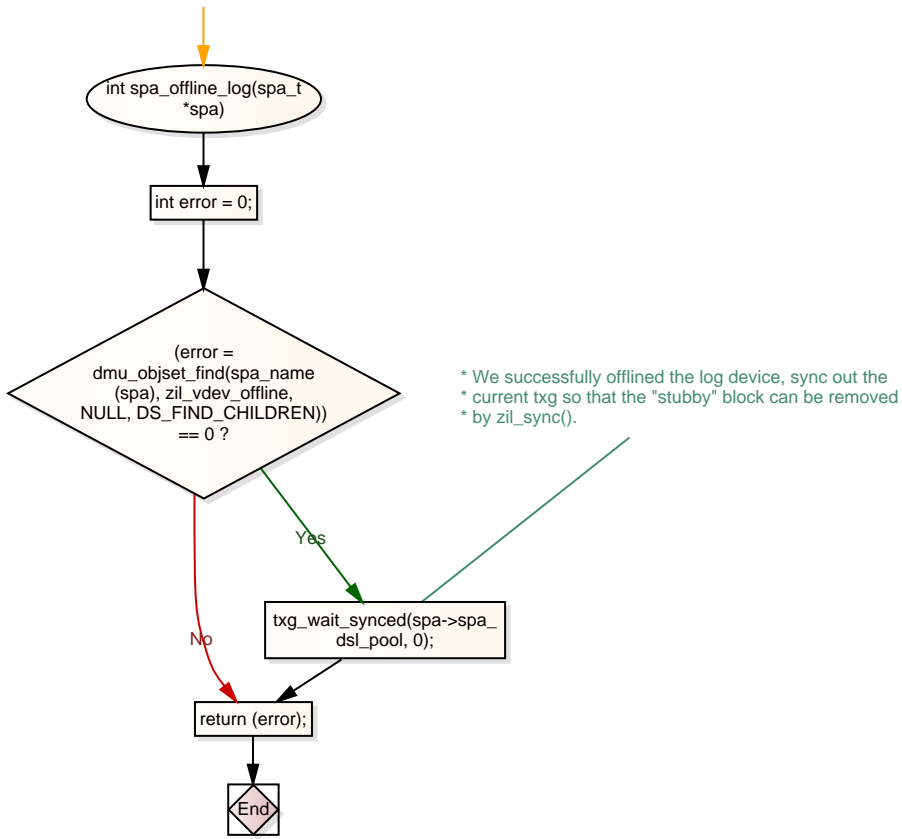
- \* Resolve any "missing" vdevs in the current configuration.
- \* If we find that the MOS config has more accurate information
- \* about the top-level vdev then use that vdev instead.

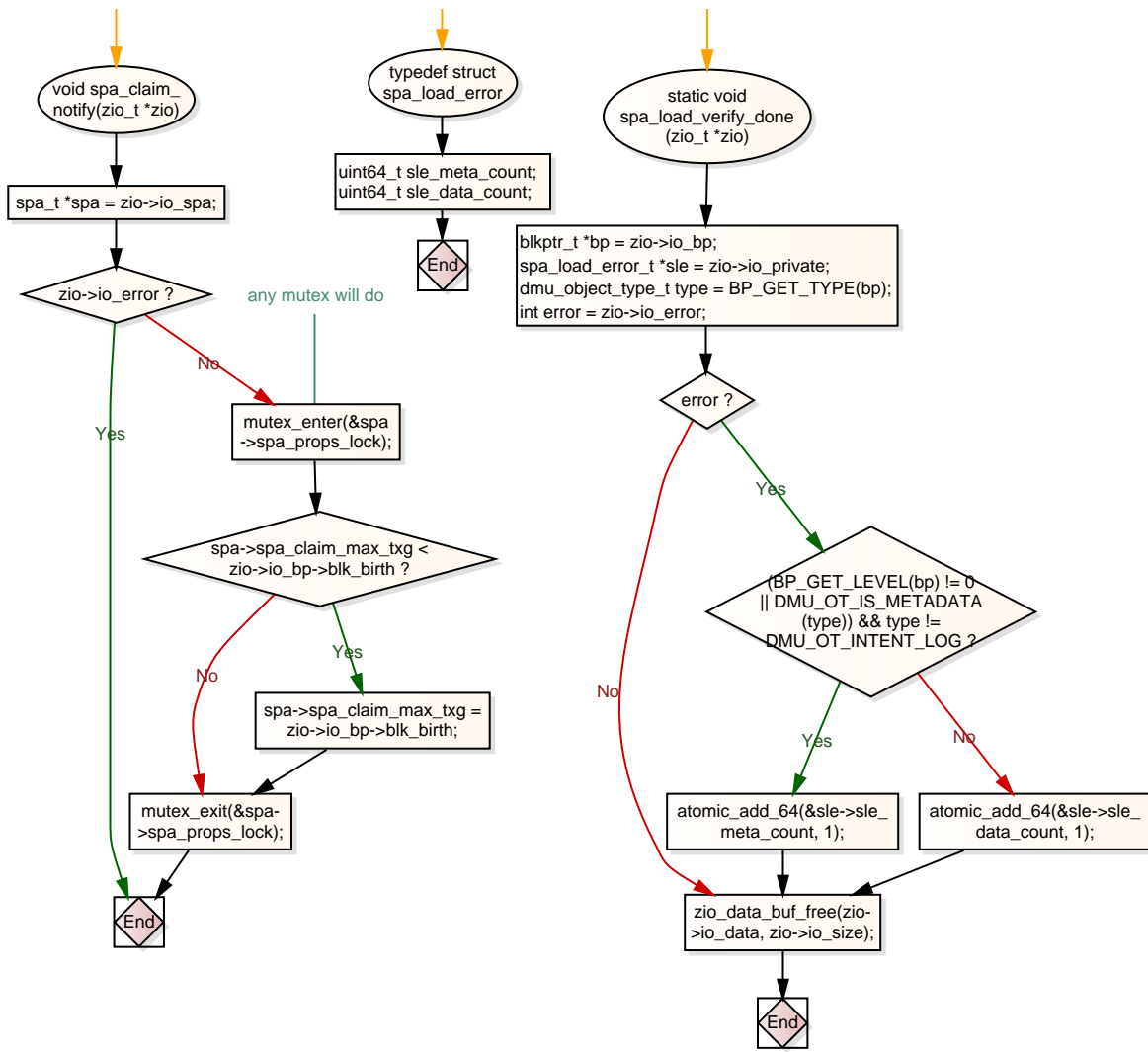


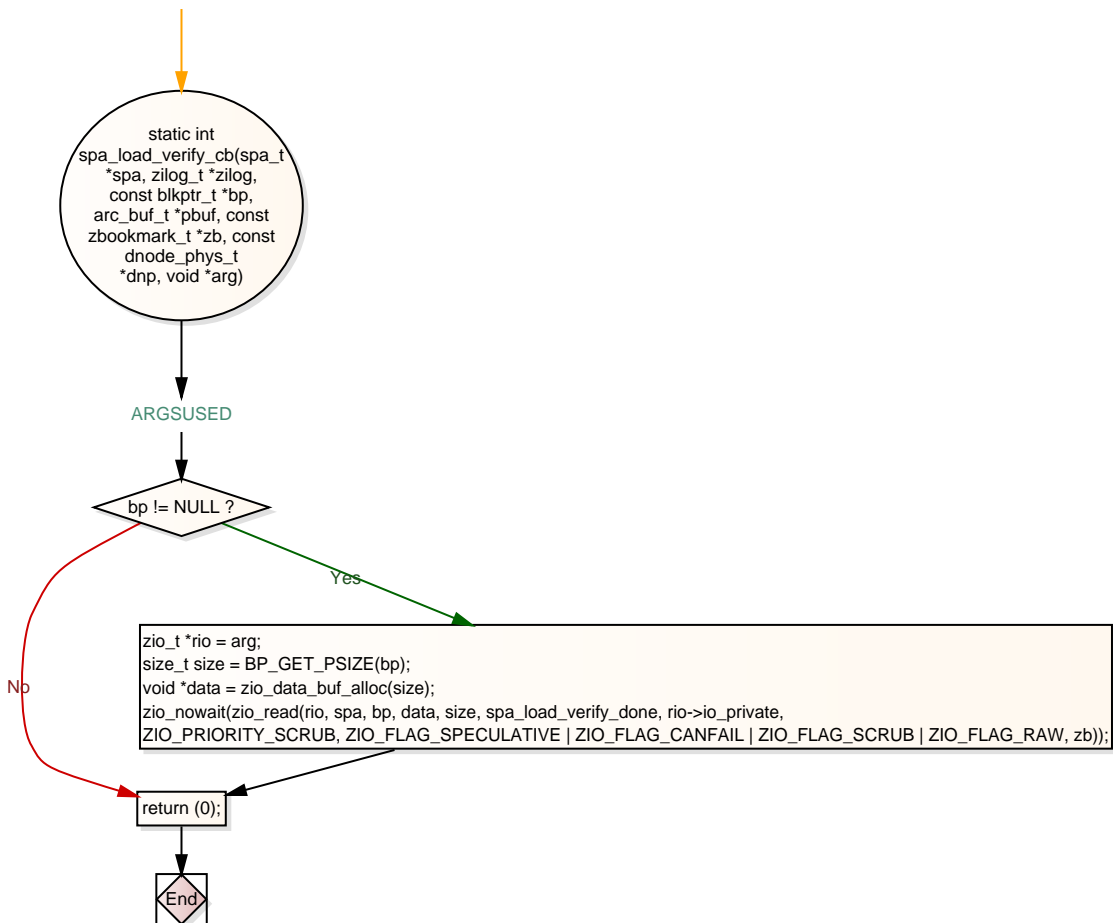


- \* Load the slog device's state from the MOS config
  - \* since it's possible that the label does not
  - \* contain the most up-to-date information.
- 

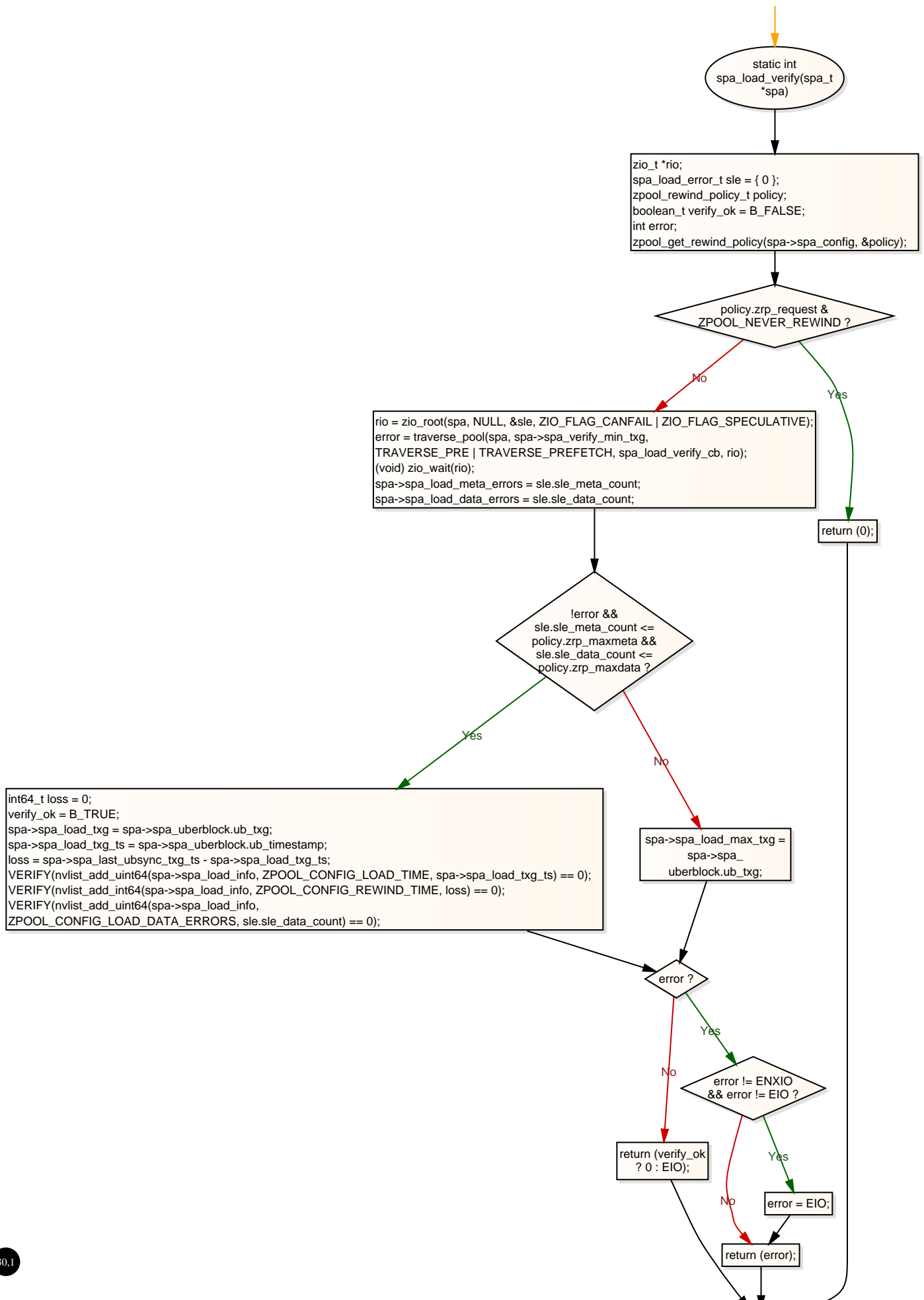


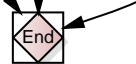


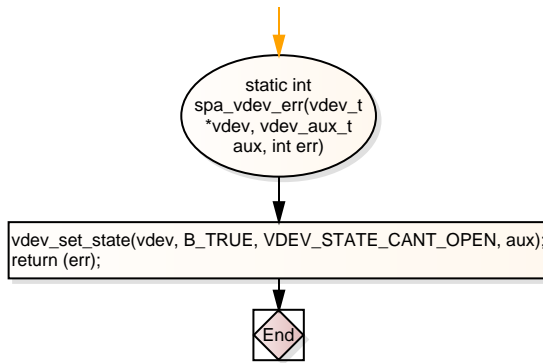
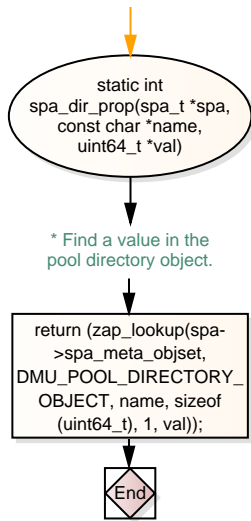
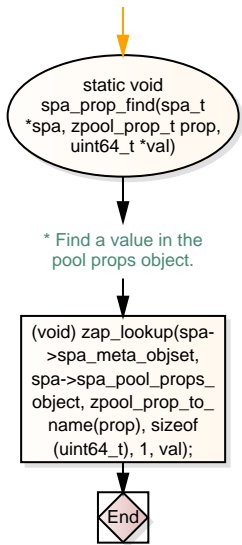


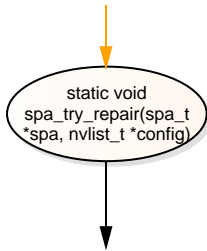




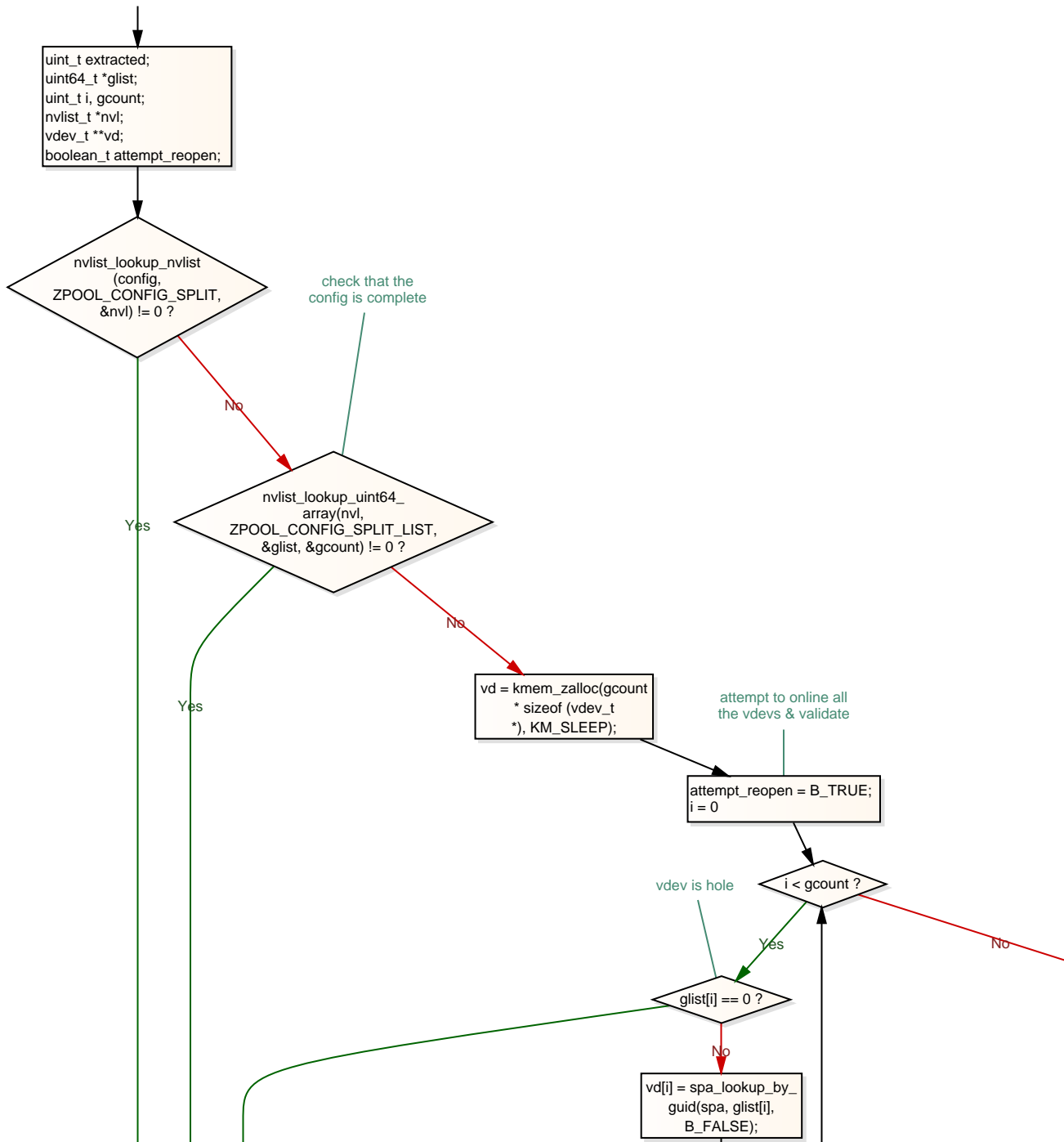


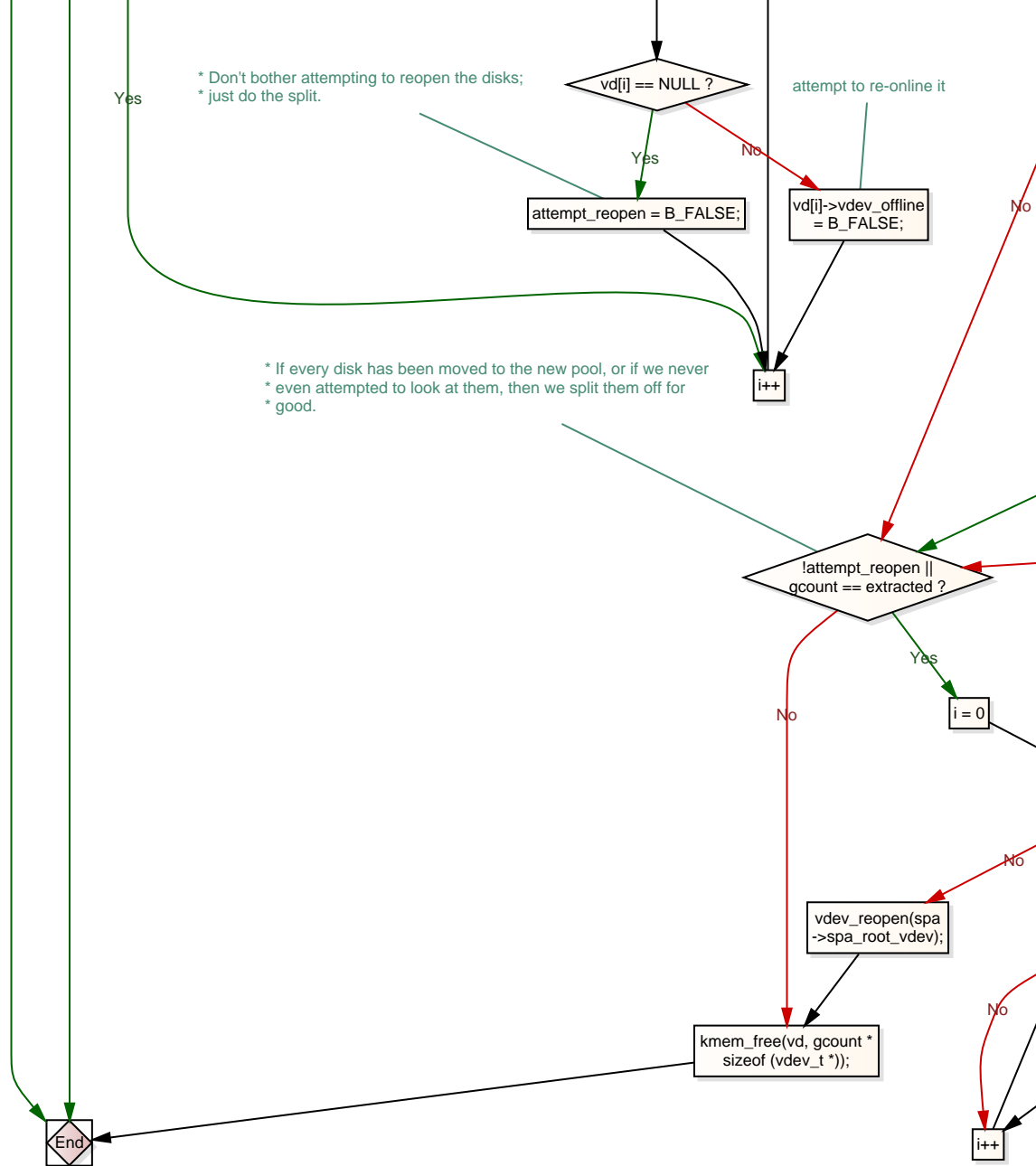


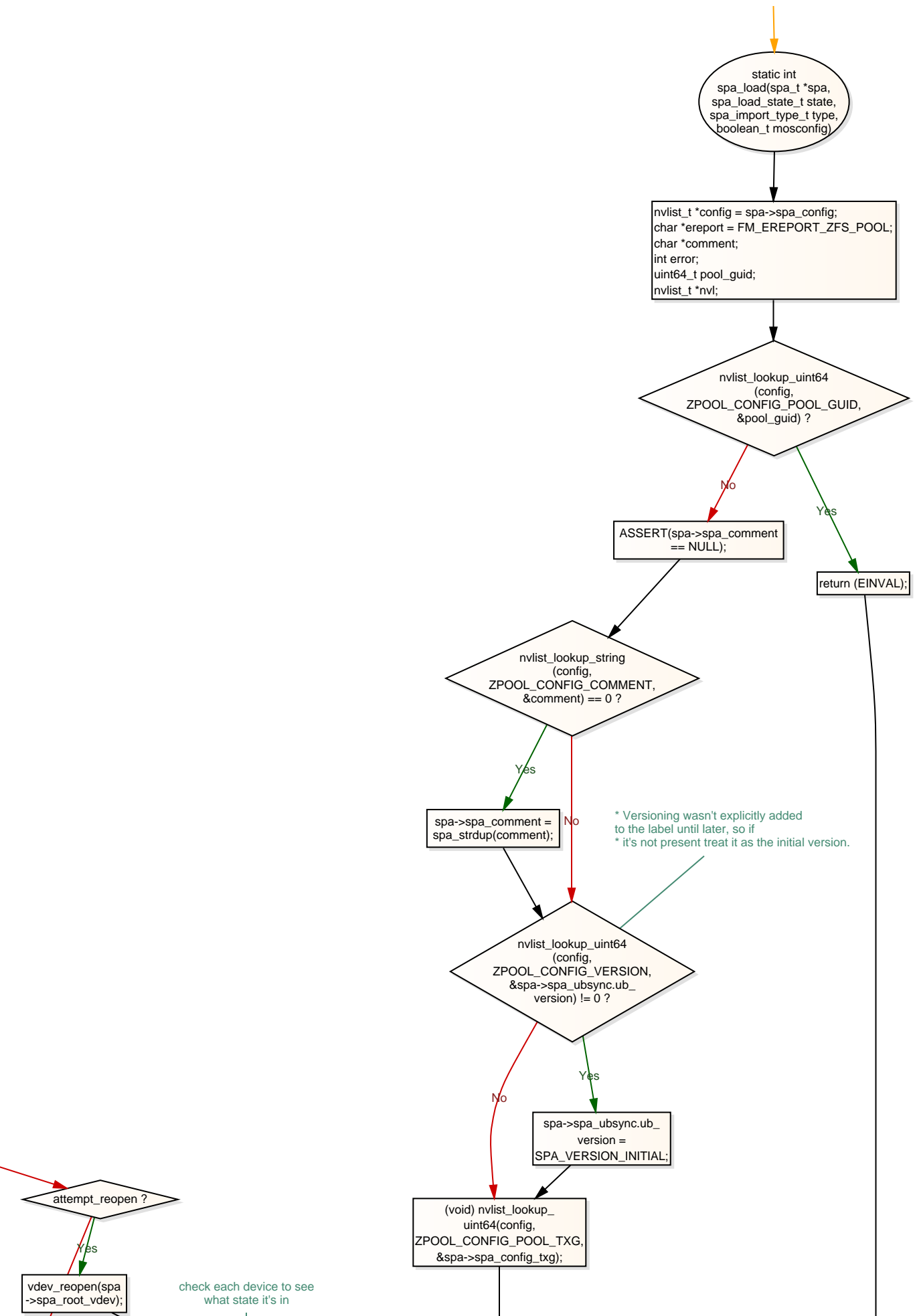


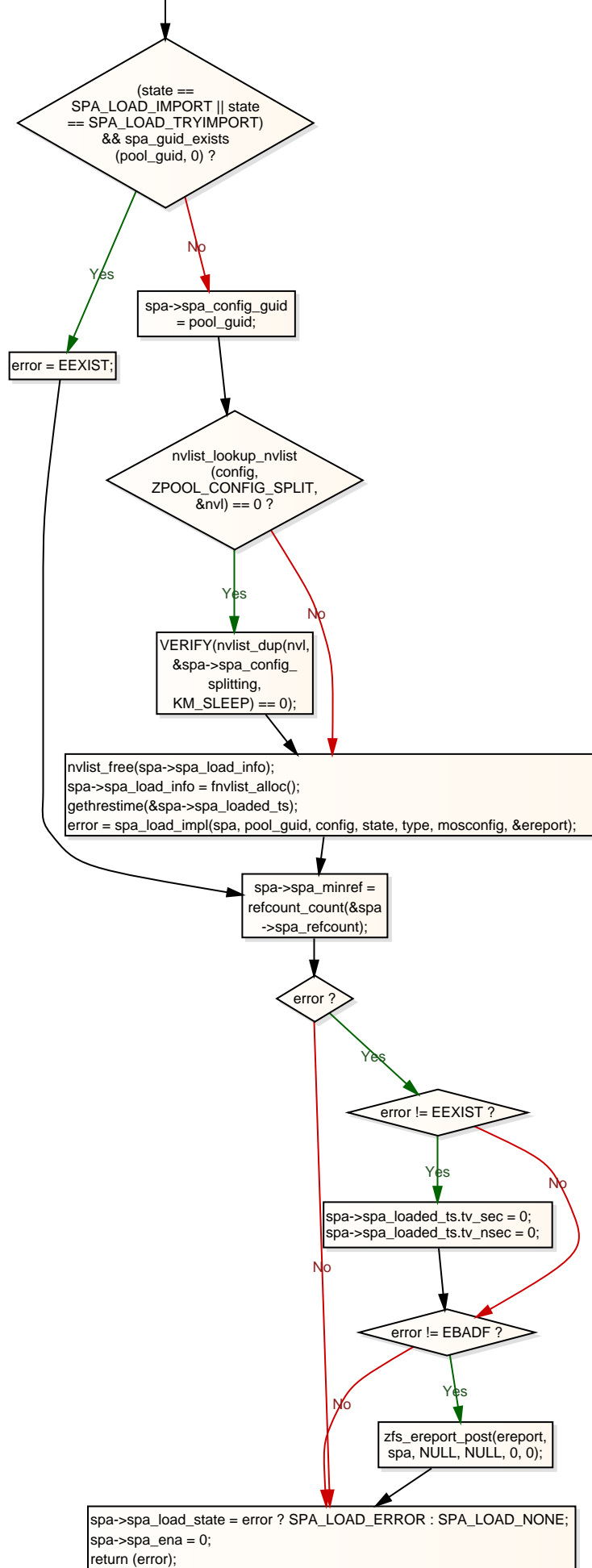
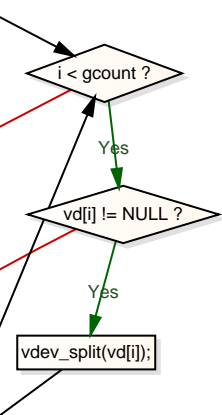
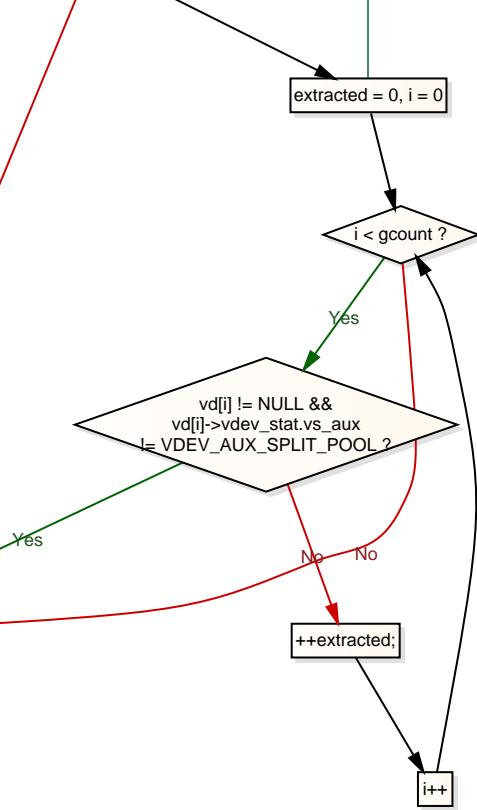


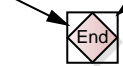
- \* Fix up config after a partly-completed split. This is done with the ZPOOL\_CONFIG\_SPLIT nvlist. Both the splitting pool and the split-off pool have that entry in their config, but only the splitting one contains a list of all the guides of the vdevs that are being split off.
- \* This function determines what to do with that list: either rejoin all the disks to the pool, or complete the splitting process. To attempt the rejoin, each disk that is offlined is marked online again, and we do a reopen() call. If the vdev label for every disk that was marked online indicates it was successfully split off (VDEV\_AUX\_SPLIT\_POOL) then we call vdev\_split() on each disk, and complete the split.
- \* Otherwise we leave the config alone, with all the vdevs in place in the original pool.



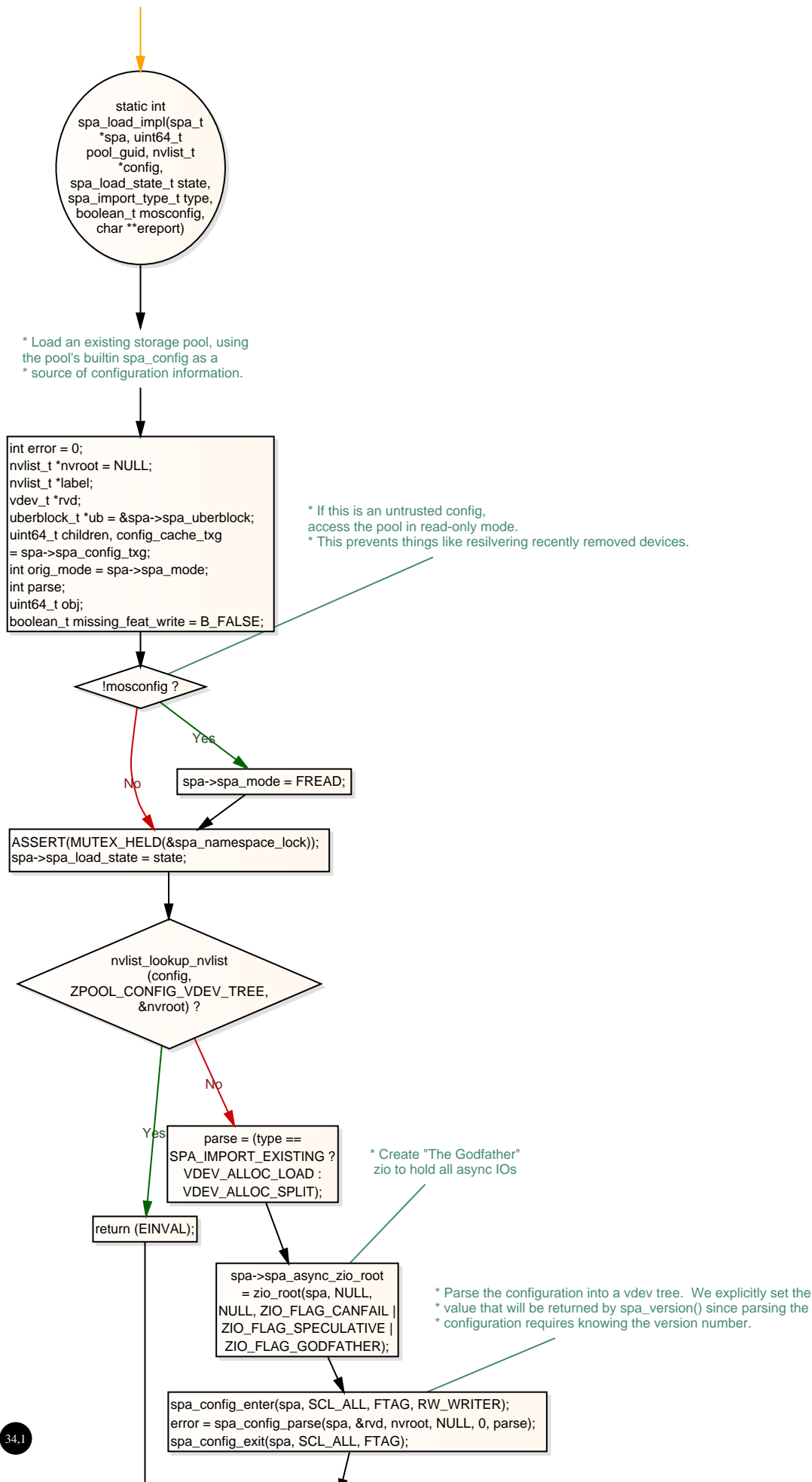


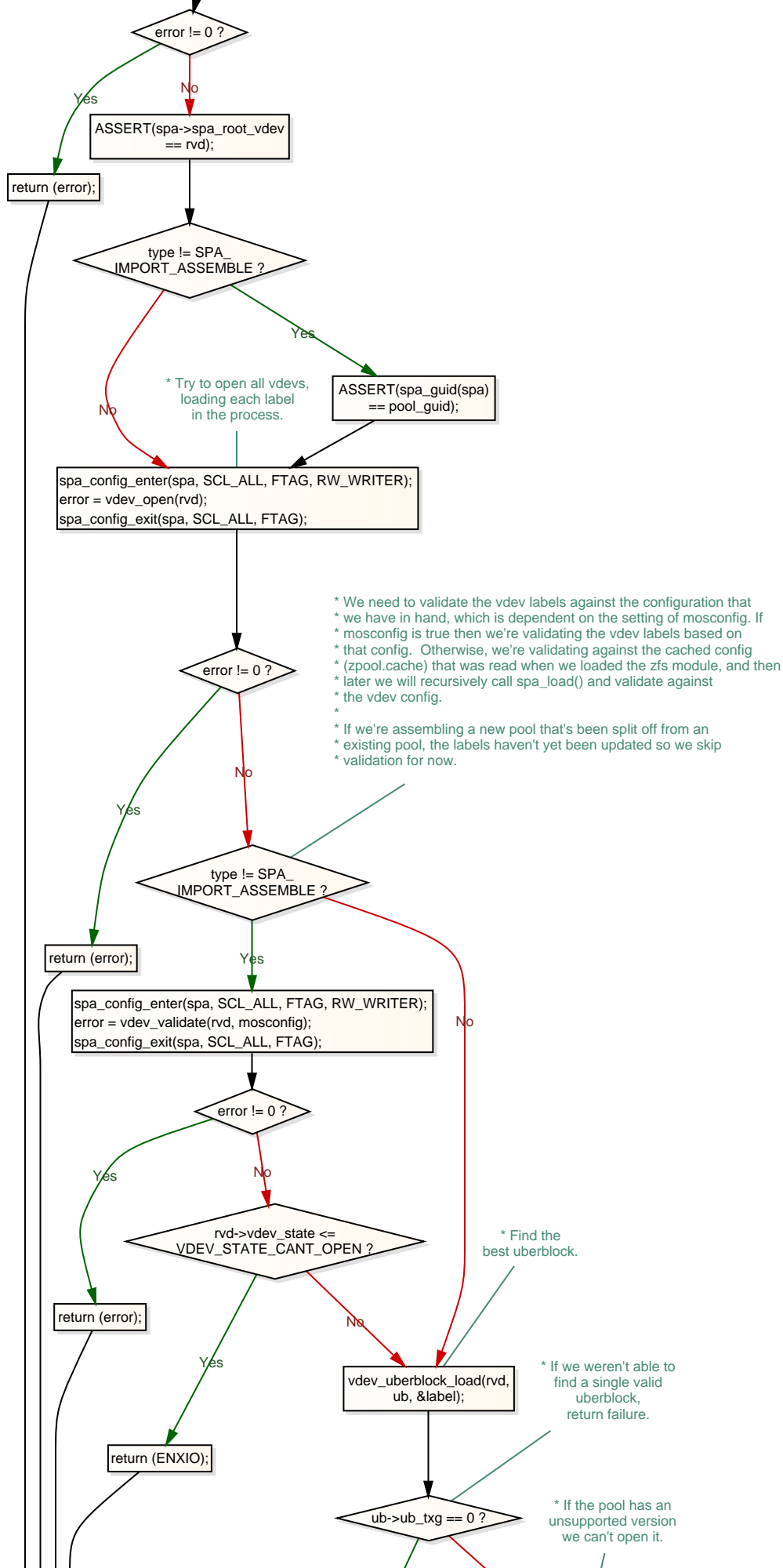


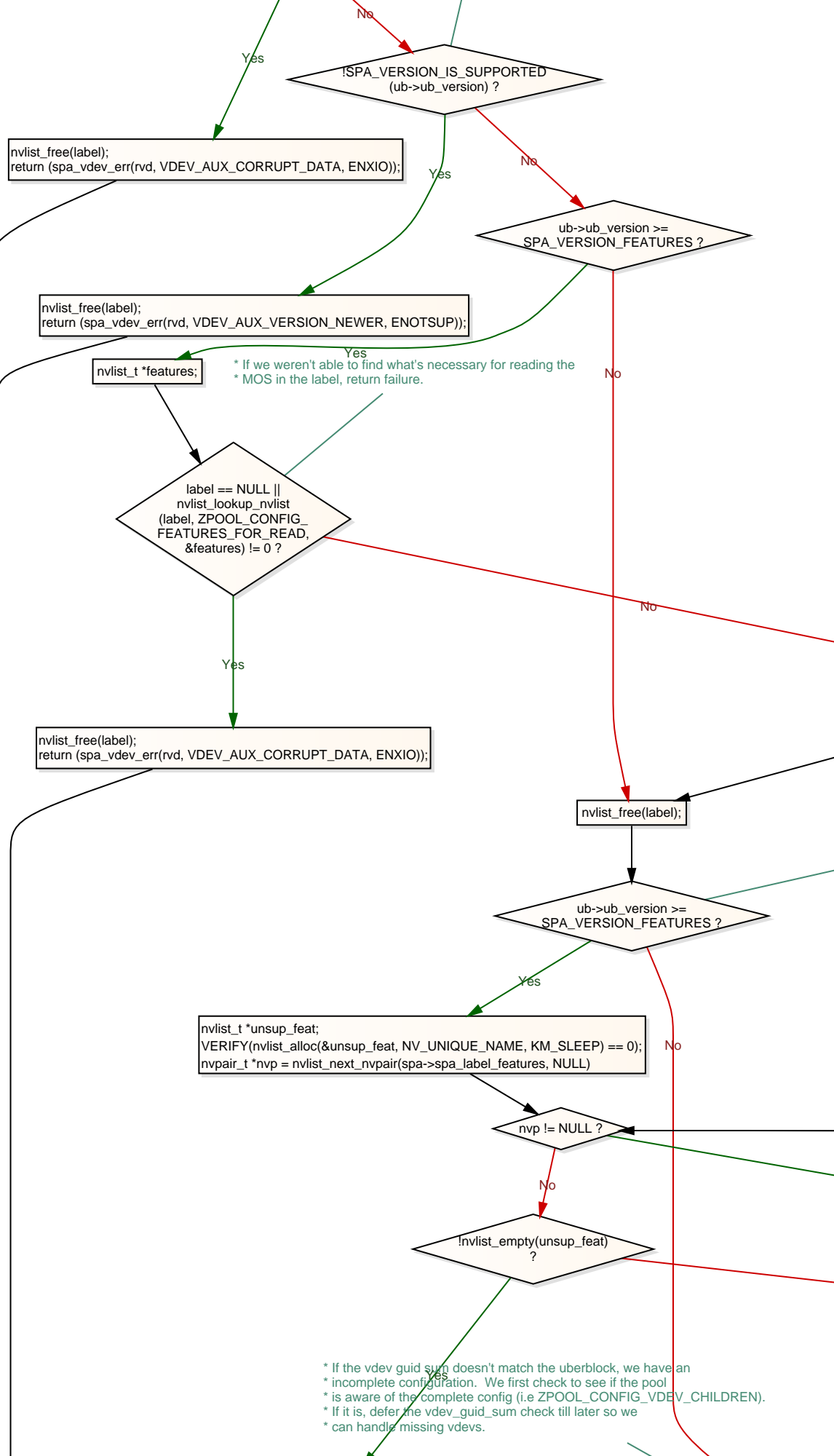








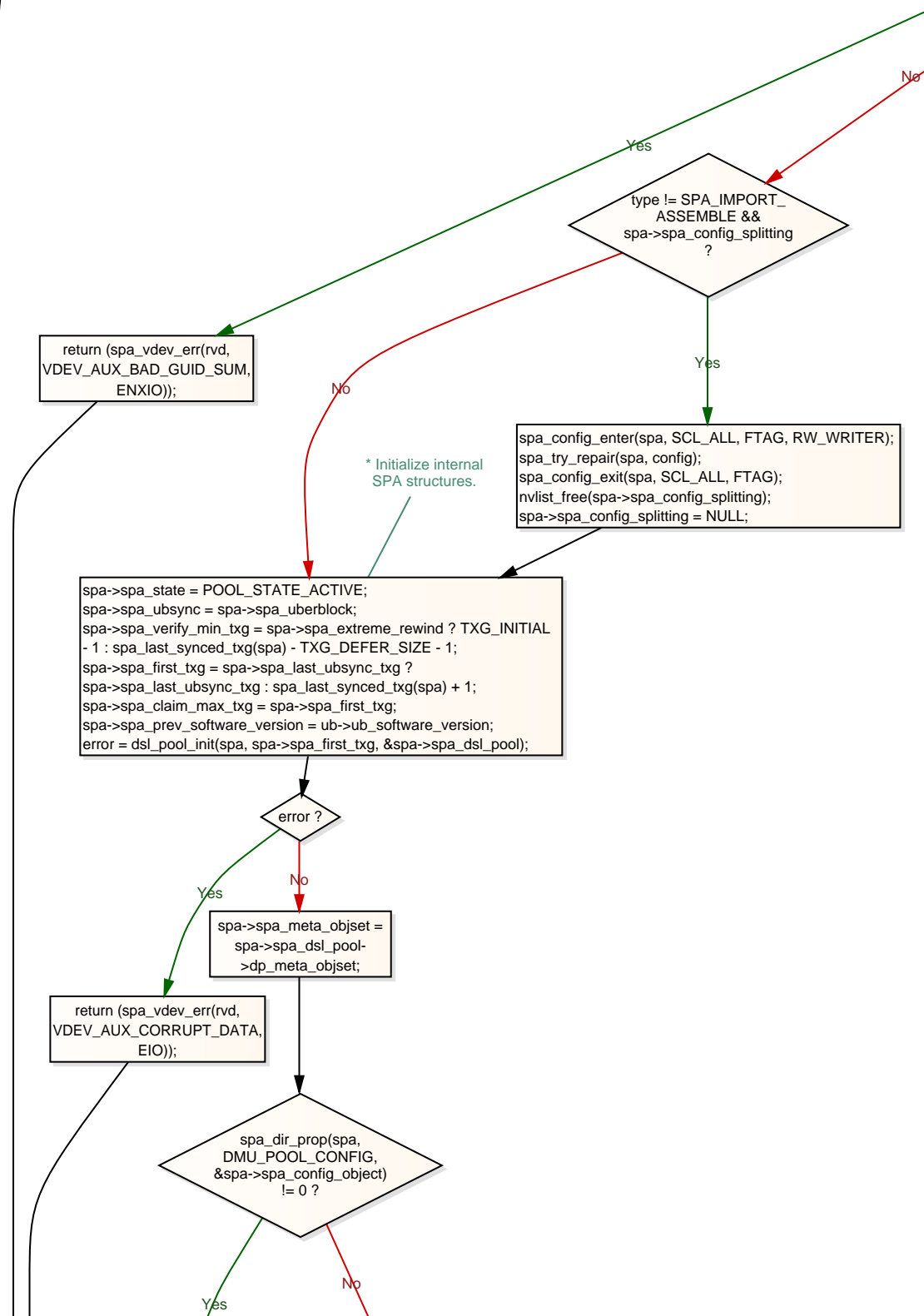


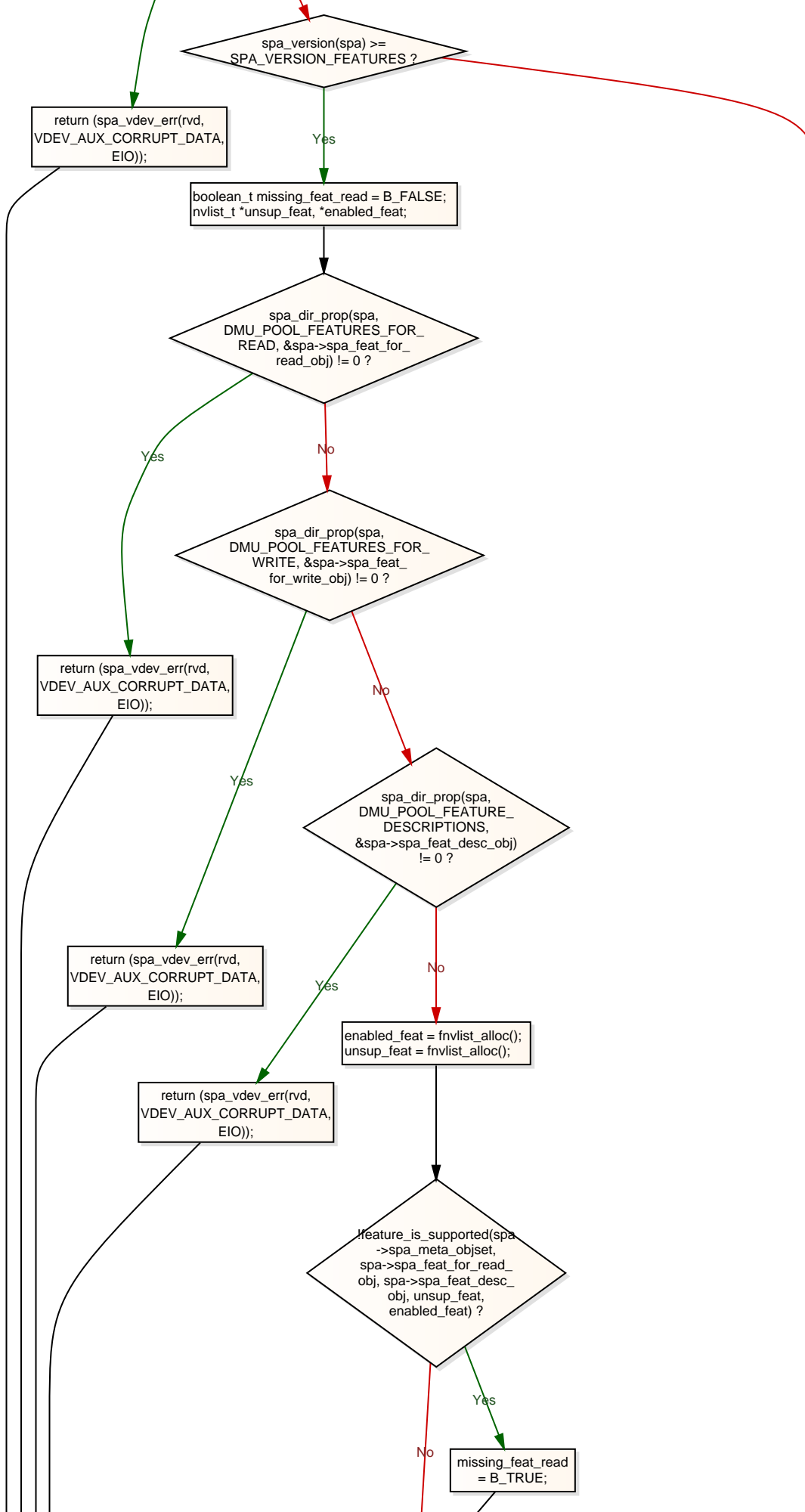


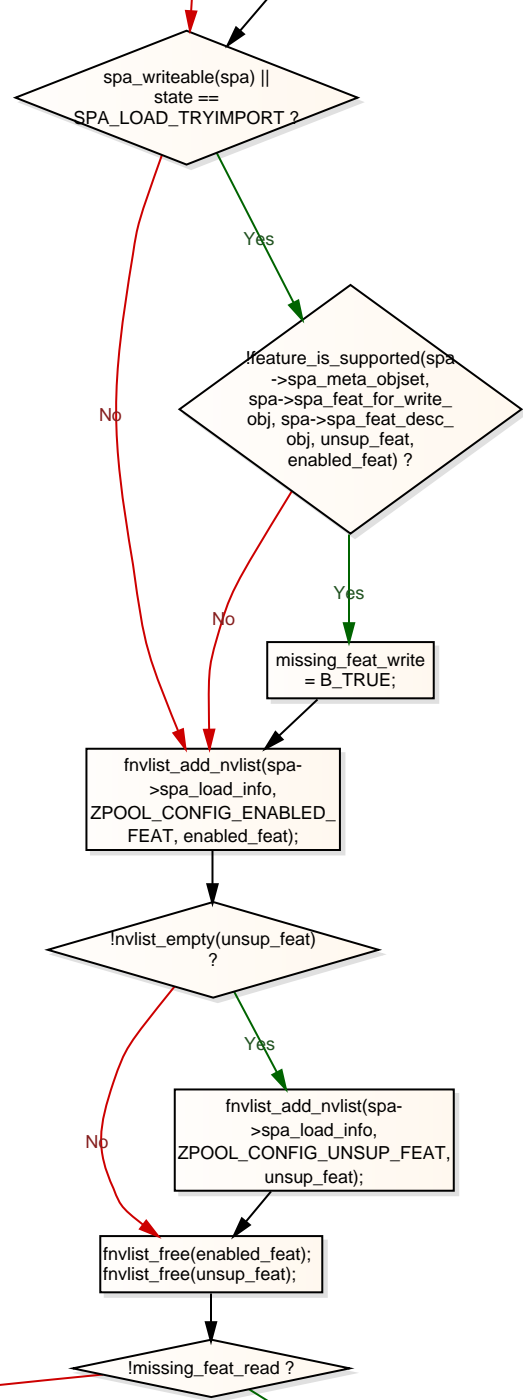
```

VERIFY(nvlist_add_nvlist(spa->spa_load_info, ZPOOL_CONFIG_UNSUP_FEAT, unsup_feat) == 0);
nvlist_free(unsup_feat);
return (spa_vdev_err(rvd, VDEV_AUX_UNSUP_FEAT, ENOTSUP));

```

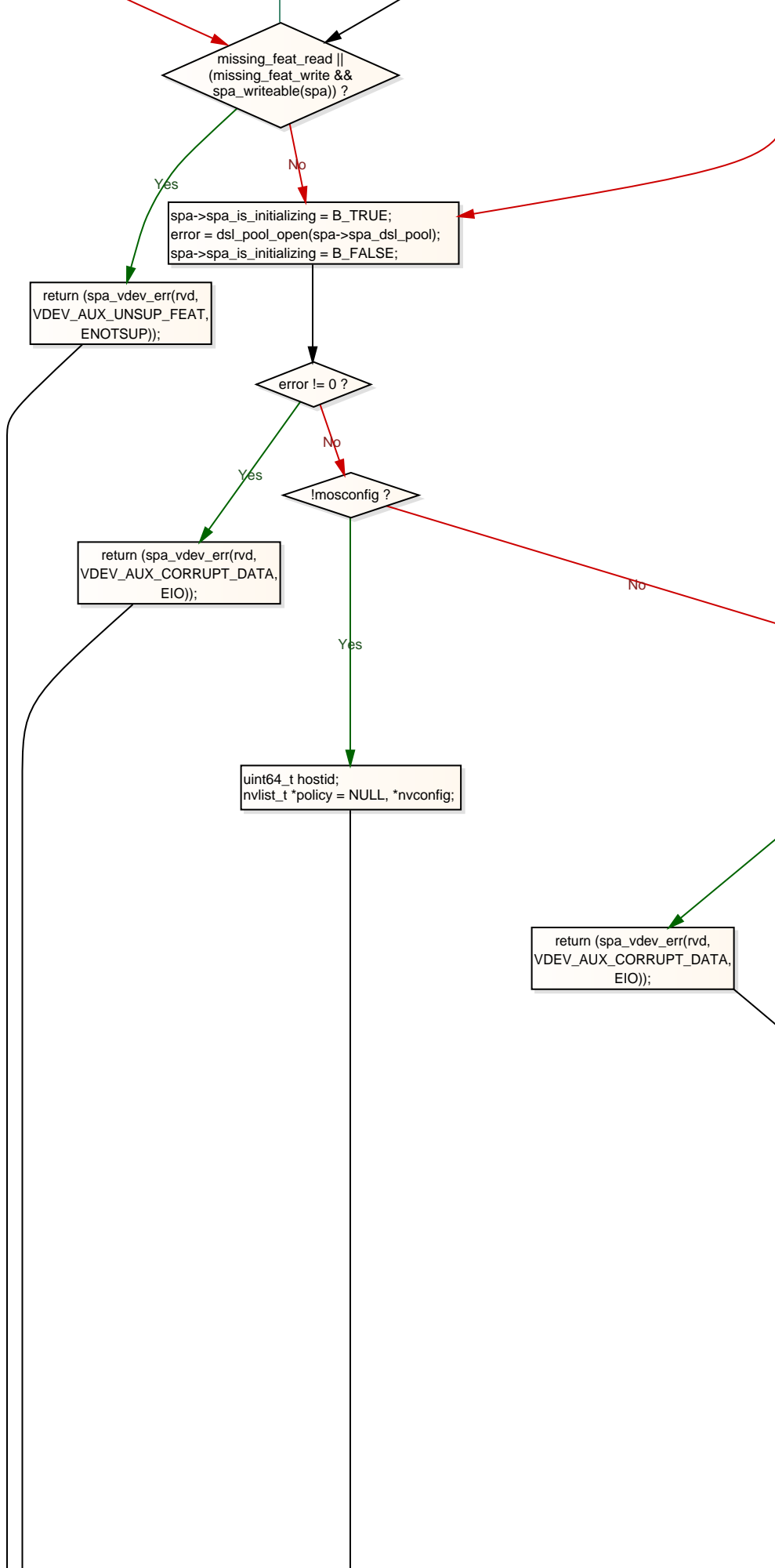


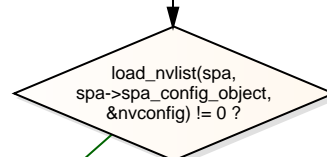




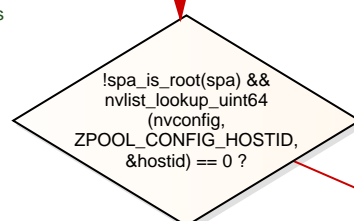
\* If the state is SPA\_LOAD\_TRYIMPORT, our objective is  
 \* twofold: to determine whether the pool is available for  
 \* import in read-write mode and (if it is not) whether the  
 \* pool is available for import in read-only mode. If the pool  
 \* is available for import in read-write mode, it is displayed  
 \* as available in userland; if it is not available for import  
 \* in read-only mode, it is displayed as unavailable in  
 \* userland. If the pool is available for import in read-only  
 \* mode but not read-write mode, it is displayed as unavailable  
 \* in userland with a special note that the pool is actually  
 \* available for open in read-only mode.  
 \*

\* As a result, if the state is SPA\_LOAD\_TRYIMPORT and we are  
 \* missing a feature for write, we must first determine whether  
 \* the pool can be opened read-only before returning whether  
 \* userland in order to know whether to display the  
 \* abovementioned note.



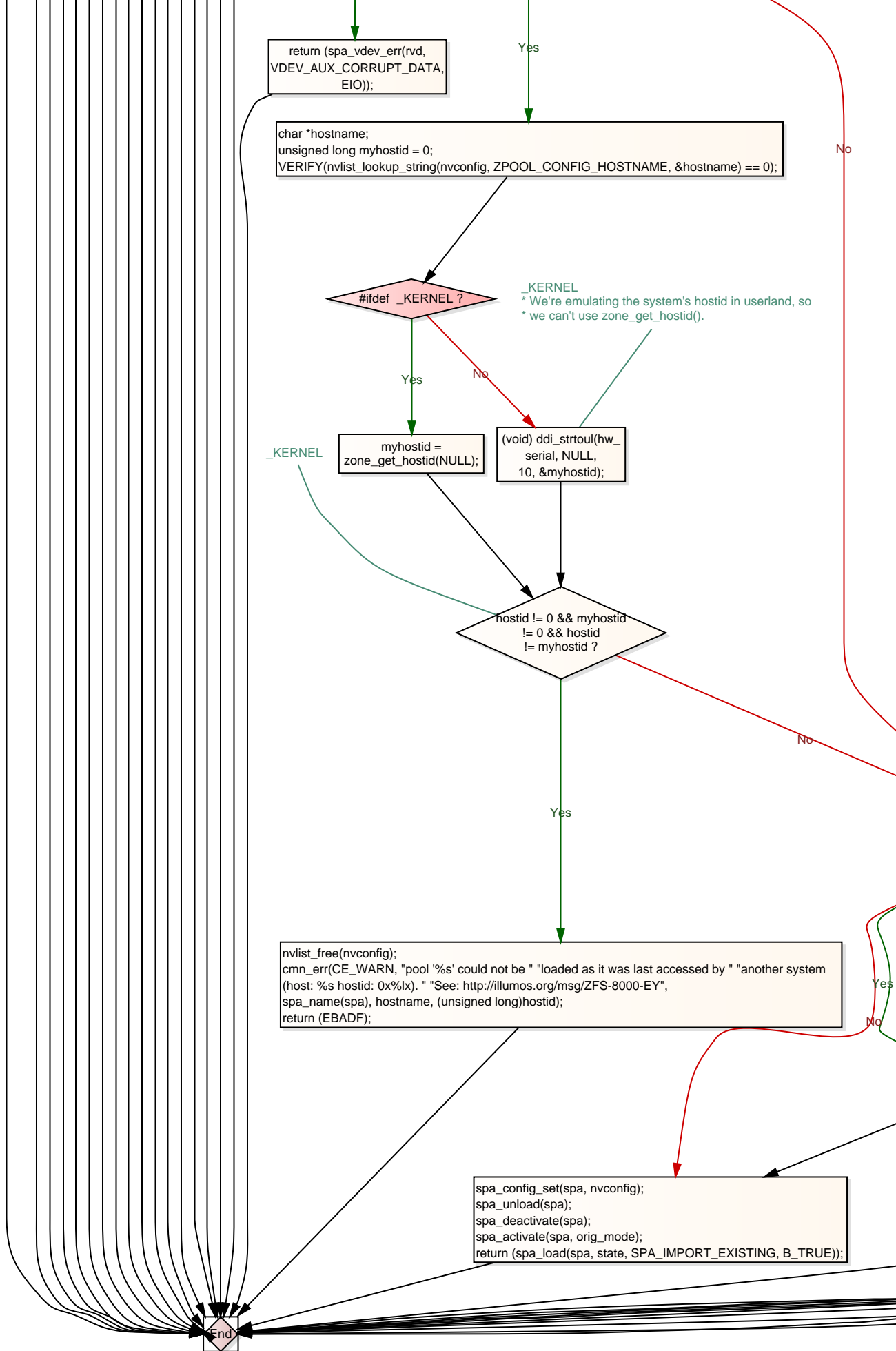


No



Yes





```
uint64_t autoreplace;
spa_prop_find(spa, ZPOOL_PROP_BOOTFS, &spa->spa_bootfs);
spa_prop_find(spa, ZPOOL_PROP_AUTOREPLACE, &autoreplace);
spa_prop_find(spa, ZPOOL_PROP_DELEGATION, &spa->spa_delegation);
spa_prop_find(spa, ZPOOL_PROP_FAILUREMODE, &spa->spa_failmode);
spa_prop_find(spa, ZPOOL_PROP_AUTOEXPAND, &spa->spa_autoexpand);
spa_prop_find(spa, ZPOOL_PROP_DEDUPDITTO, &spa->spa_dedup_ditto);
spa->spa_autoreplace = (autoreplace != 0);
```

```
graph TD; A[return (spa_vdev_err(rvd, VDEV_AUX_CORRUPT_DATA, EIO));] -- green arrow --> B[nvlist_free(nvconfig);]; B -- red arrow --> A;
```

return (spa\_vdev\_err(rvd,  
VDEV\_AUX\_CORRUPT\_DATA,  
EIO));

nvlist\_free(nvconfig);

```
*ereport = FM_EREPORT_ZFS_LOG_REPLAY;  
return (spa_vdev_err(rvd, VDEV_AUX_BAD_LOG, ENXIO));
```

\* At this point, we know that we can open the pool in  
\* read-only mode but not read-write mode. We now have enough  
\* information and can return to userland.

```
return (spa_vdev_err(rvd,  
VDEV_AUX_UNSUP_FEAT,  
ENOTSUP));
```

\* Update our in-core representation with the definitive values  
\* from the label.

```
nvlist_free(spa->spa_label_features);  
VERIFY(nvlist_dup(features, &spa->spa_label_features, 0) == 0);
```

\* Look through entries in the label nvlist's features\_for\_read. If  
\* there is a feature listed there which we don't understand then we  
\* cannot open a pool.

\* Claim log blocks that haven't been committed yet.  
\* This must all happen in a single txg.  
\* Note: spa\_claim\_max\_txg is updated by spa\_claim\_notify(),  
\* invoked from zil\_claim\_log\_block()'s i/o done callback.  
\* Price of rollback is that we abandon the log.

\* Wait for all claims to sync. We sync up to the highest  
\* claimed log block birth time so that claimed log blocks  
\* don't appear to be from the future. spa\_claim\_max\_txg  
\* will have been set for us by either zil\_check\_log\_chain()  
\* (invoked from spa\_check\_logs()) or zil\_claim() above.

\* If the config cache is stale, or we have uninitialized  
\* metaslabs (see spa\_vdev\_add()), then update the config.  
\*  
\* If this is a verbatim import, trust the current  
\* in-core spa\_config and update the disk labels.

Yes

!zfeature\_is\_supported  
(nvpair\_name(nvp)) ?

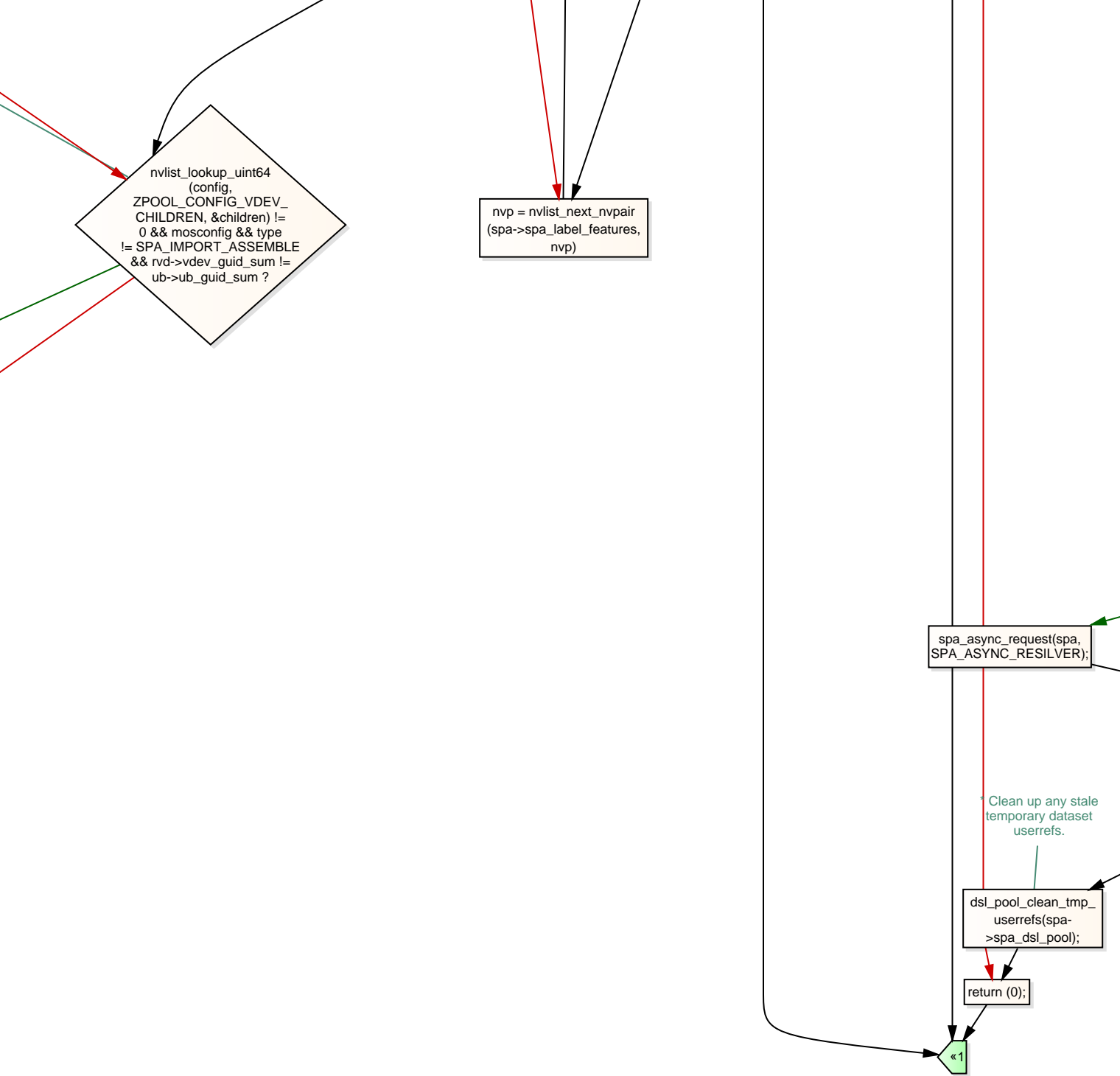
No

```
nvlist_free(unsup_feat);
```


Yes

```
VERIFY(nvlist_add_string  
(unsup_feat,  
nvpair_name(nvp),  
"") == 0);
```

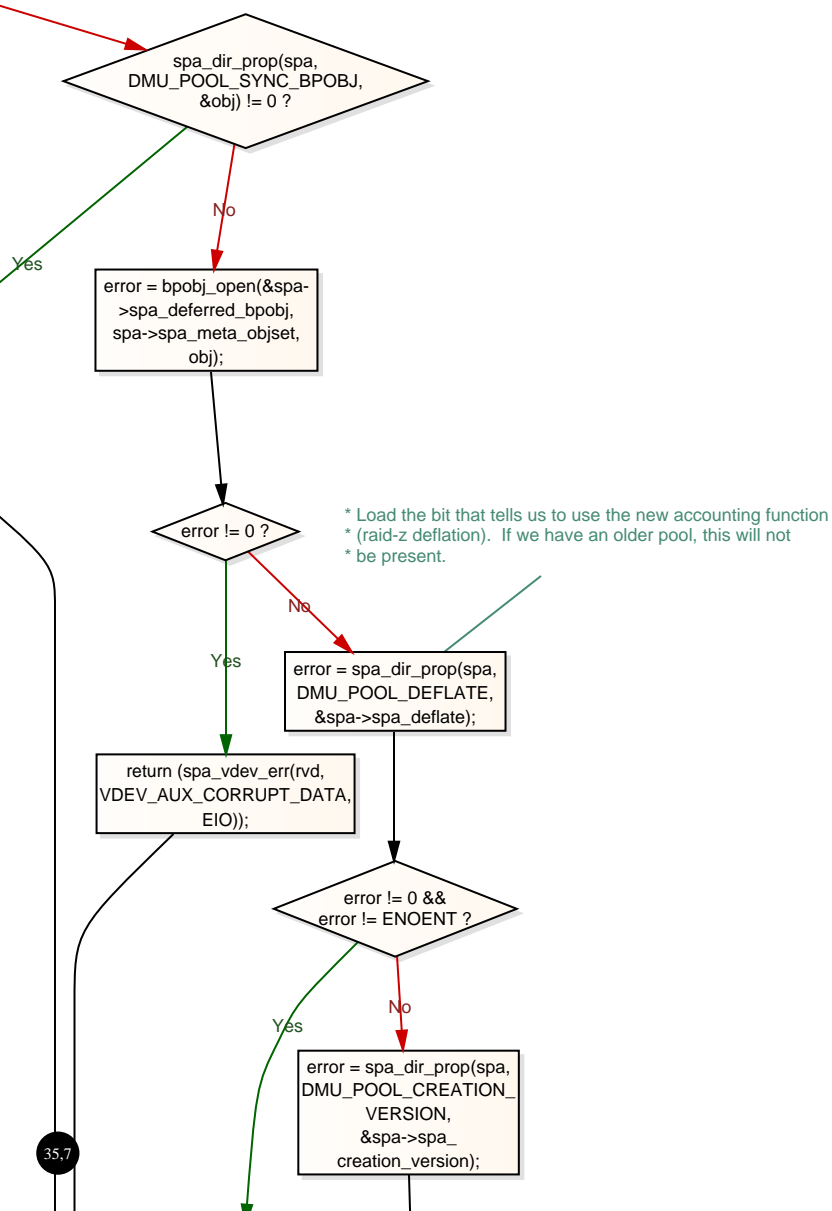
No



No



```
fnvlist_add_boolean(spa-  
    >spa_load_info,  
    ZPOOL_CONFIG_CAN_RDONLY);
```





```
return (spa_vdev_err(rvd,  
VDEV_AUX_CORRUPT_DATA,  
EIO));
```

error != 0 &&  
error != ENOENT ?

\* Load the persistent error log.  
If we have an older pool, this will  
\* not be present.

Yes

No

```
error = spa_dir_prop(spa,  
DMU_POOL_ERRLOG_LAST,  
&spa->spa_errlog_last);
```

```
return (spa_vdev_err(rvd,  
VDEV_AUX_CORRUPT_DATA,  
EIO));
```

error != 0 &&  
error != ENOENT ?

Yes

No

```
error = spa_dir_prop(spa,  
DMU_POOL_ERRLOG_SCRUB,  
&spa->spa_errlog_scrub);
```

```
return (spa_vdev_err(rvd,  
VDEV_AUX_CORRUPT_DATA,  
EIO));
```

error != 0 &&  
error != ENOENT ?

\* Load the history object. If we have an older pool, this  
\* will not be present.

Yes

No

```
error = spa_dir_prop(spa,  
DMU_POOL_HISTORY,  
&spa->spa_history);
```

```
return (spa_vdev_err(rvd,  
VDEV_AUX_CORRUPT_DATA,  
EIO));
```

error != 0 &&  
error != ENOENT ?

\* If we're assembling the pool from the split-off vdevs of  
\* an existing pool, we don't want to attach the spares & cache  
\* devices.  
\* Load any hot spares for this pool.

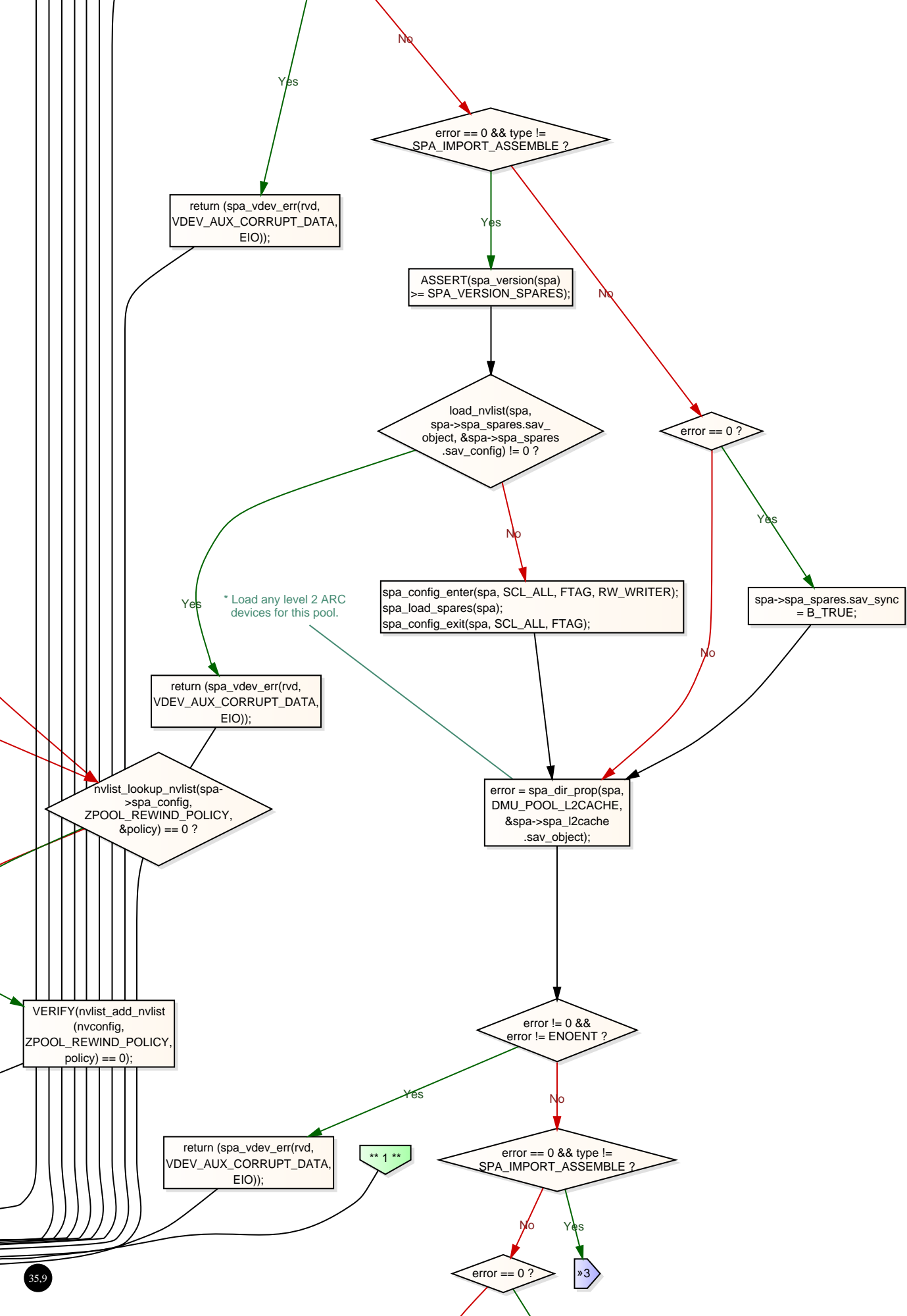
Yes

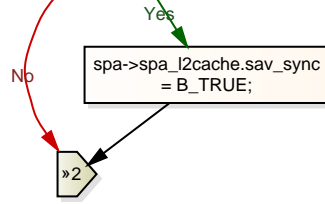
No

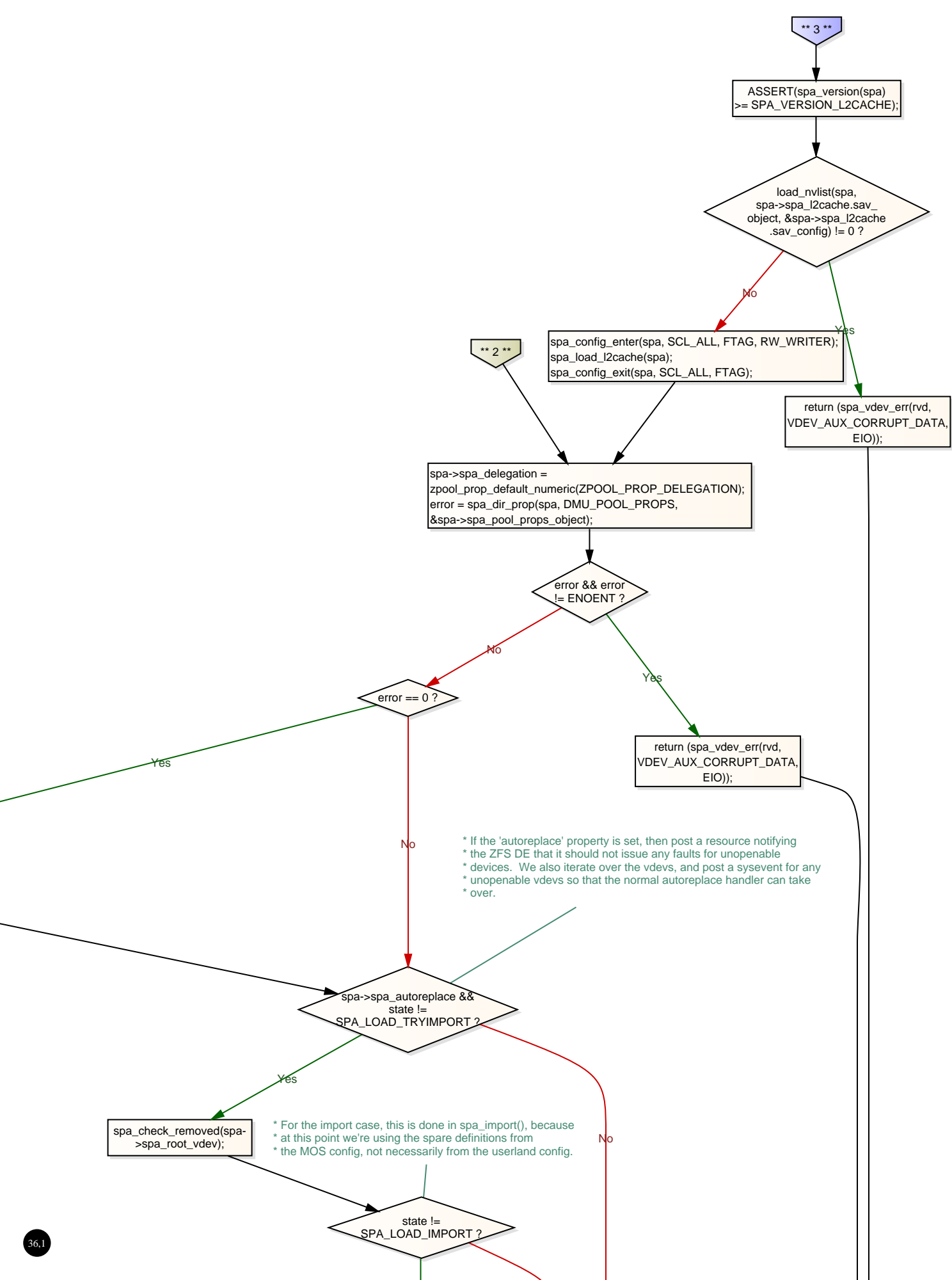
```
error = spa_dir_prop(spa,  
DMU_POOL_SPARES,  
&spa->spa_spares  
.sav_object);
```

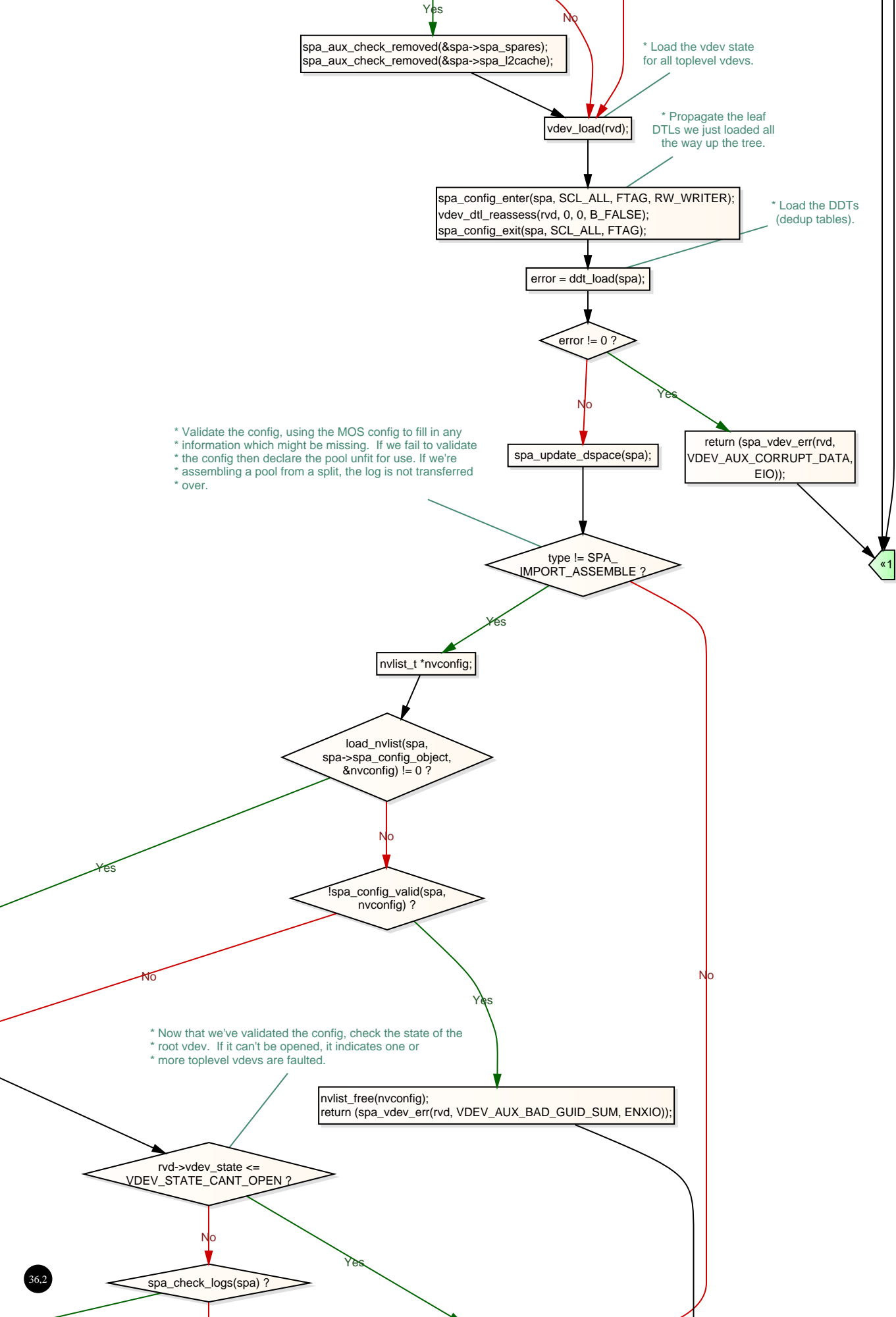
```
return (spa_vdev_err(rvd,  
VDEV_AUX_CORRUPT_DATA,  
EIO));
```

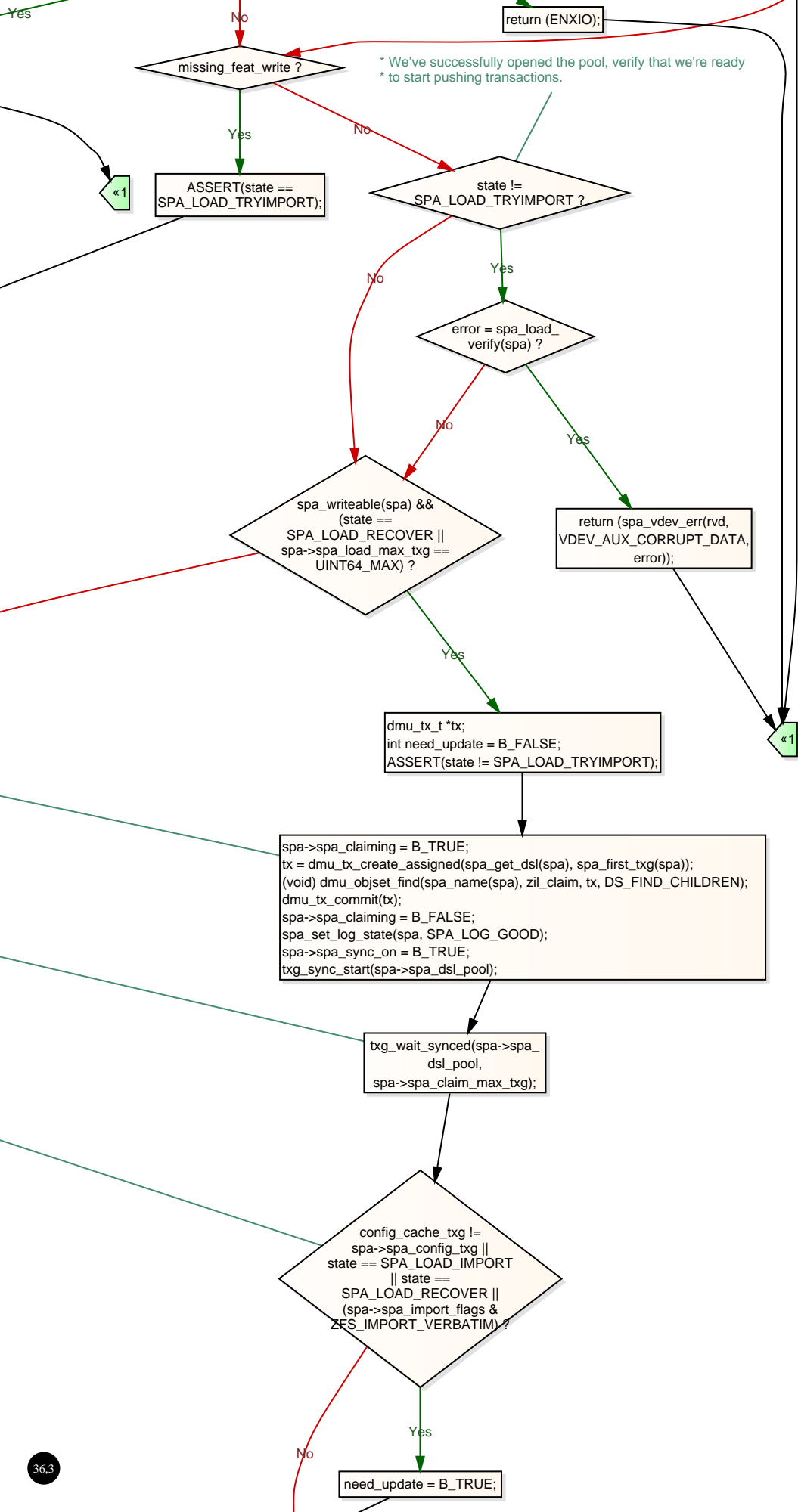
error != 0 &&  
error != ENOENT ?

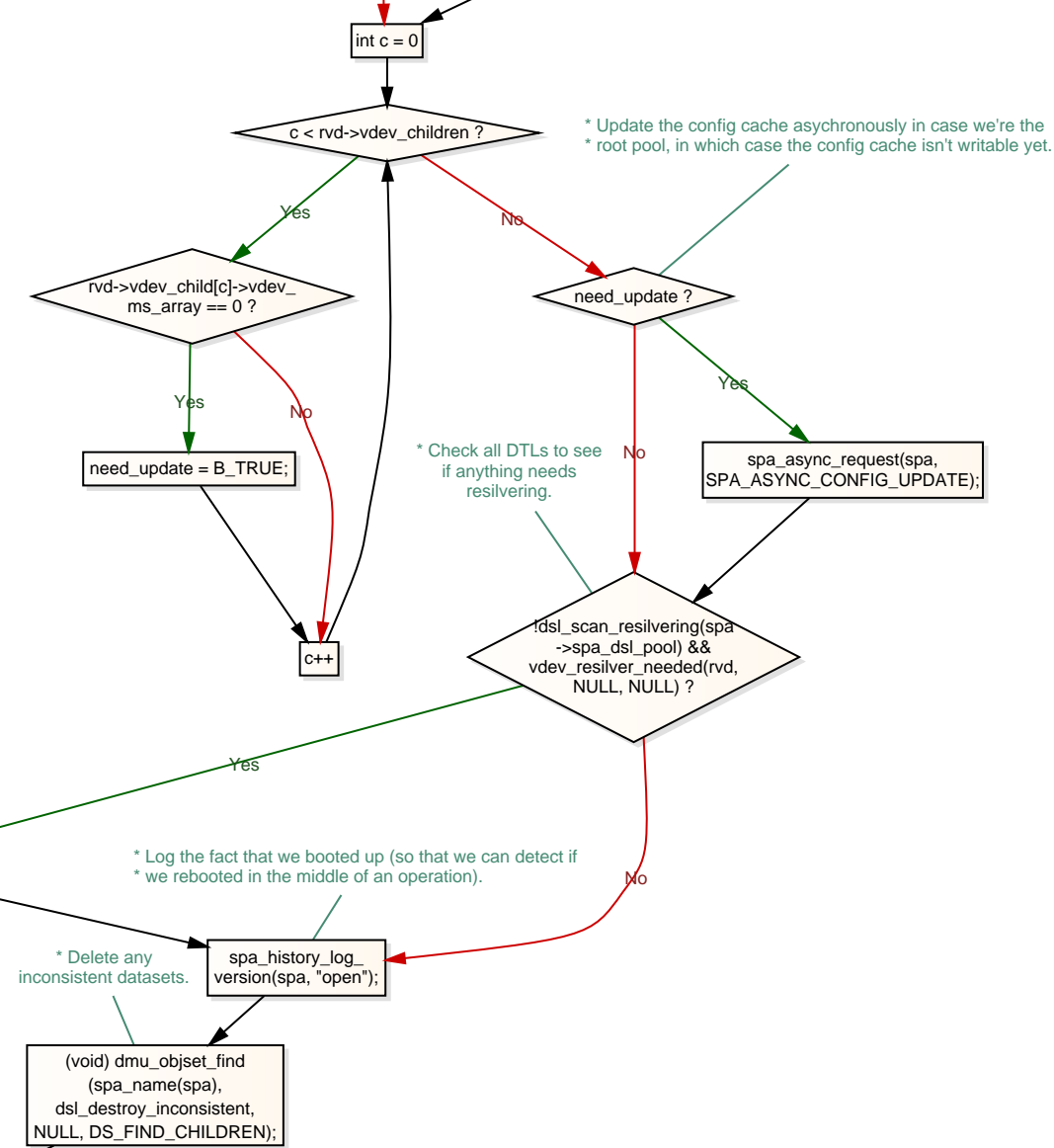


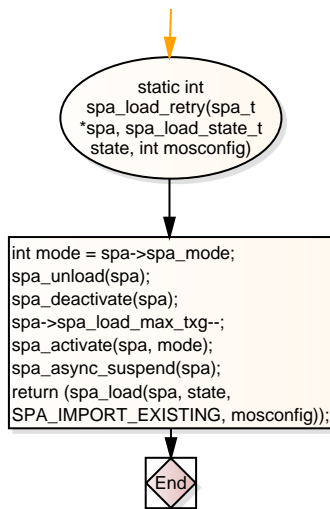




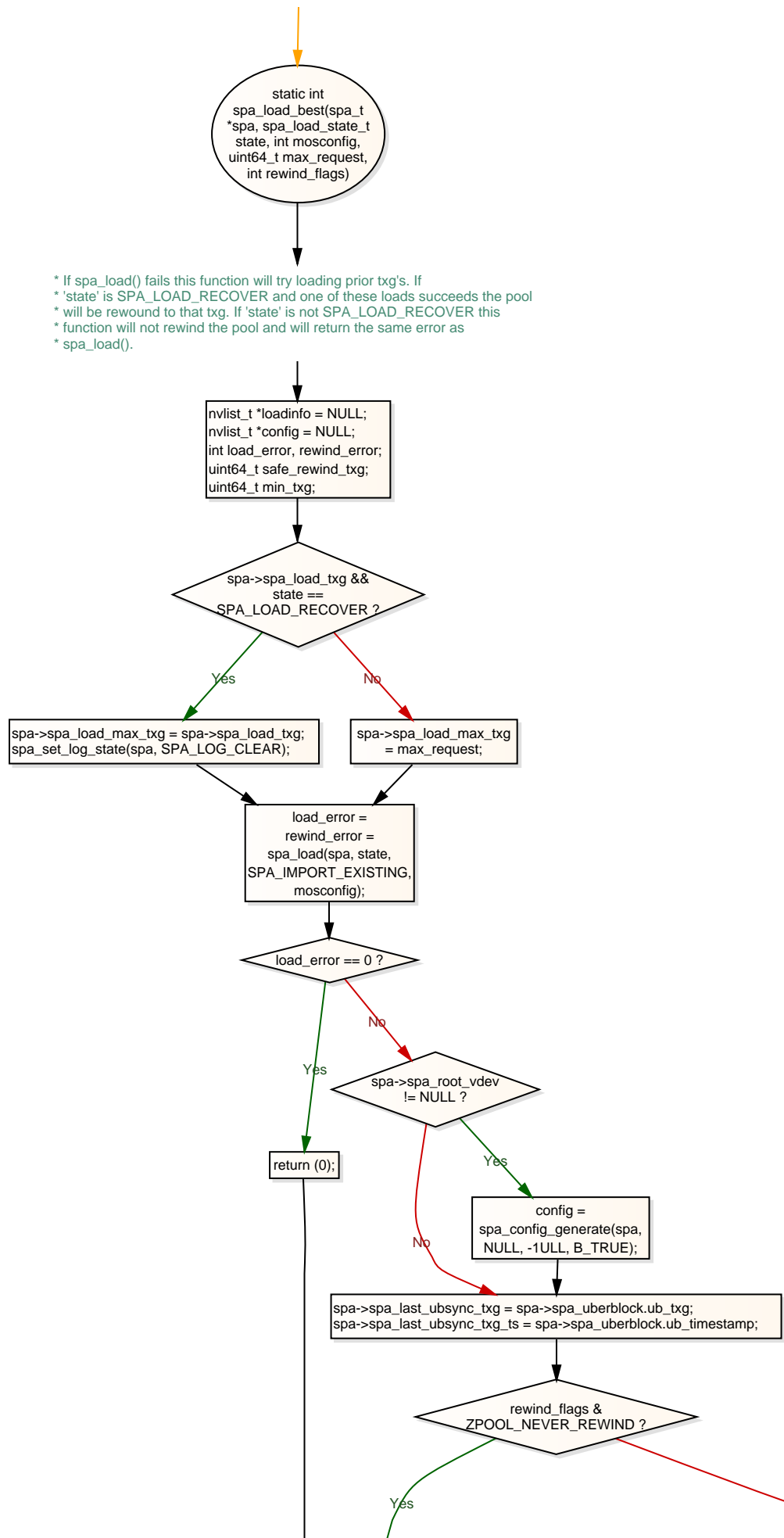








\* If spa\_load() fails this function will try loading prior txg's. If  
 \* 'state' is SPA\_LOAD\_RECOVER and one of these loads succeeds the pool  
 \* will be rewound to that txg. If 'state' is not SPA\_LOAD\_RECOVER this  
 \* function will not rewind the pool and will return the same error as  
 \* spa\_load().





\* If we aren't rolling back save the load info from our first  
\* import attempt so that we can restore it after attempting  
\* to rewind.

```
nvlist_free(config);  
return (load_error);
```

```
spa->spa_extreme_rewind = B_FALSE;  
spa->spa_load_max_txg = UINT64_MAX;
```

config && (rewind\_error  
|| state !=  
SPA\_LOAD\_RECOVER) ?

```
spa_config_set(spa,  
config);
```

state ==  
SPA\_LOAD\_RECOVER ?

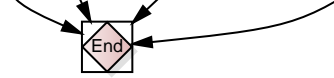
```
fnvlist_add_nvlist  
(loadinfo,  
ZPOOL_CONFIG_REWIND_INFO,  
spa->spa_load_info);
```

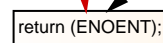
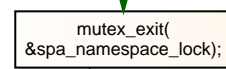
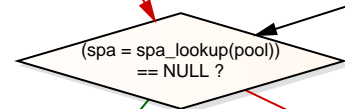
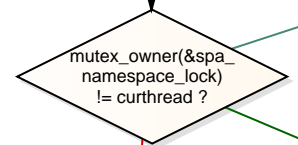
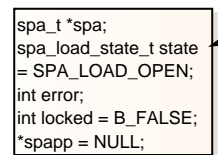
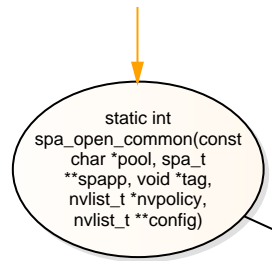
```
ASSERT3P(loadinfo, ==, NULL);  
return (rewind_error);
```

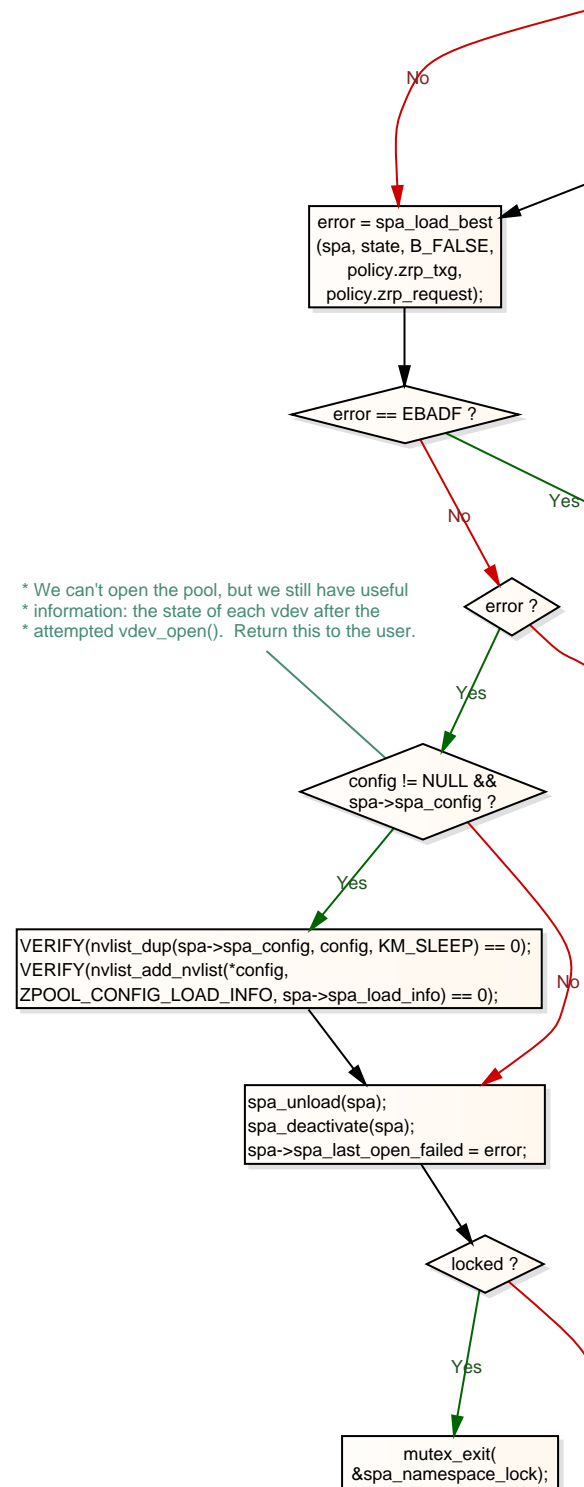
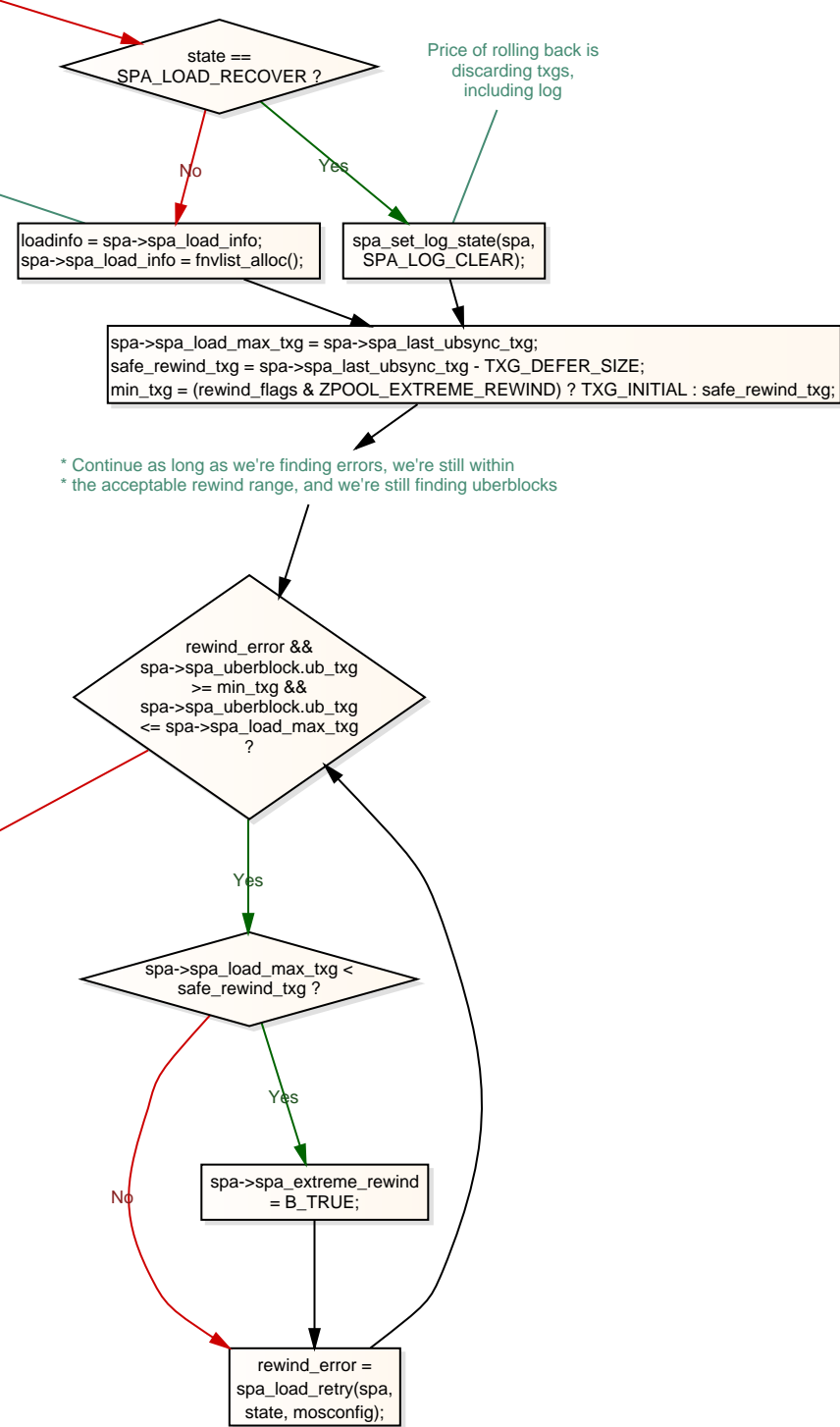
```
fnvlist_free(spa->spa_load_info);  
spa->spa_load_info = loadinfo;  
return (load_error);
```

Store the rewind info as  
part of the  
initial load info

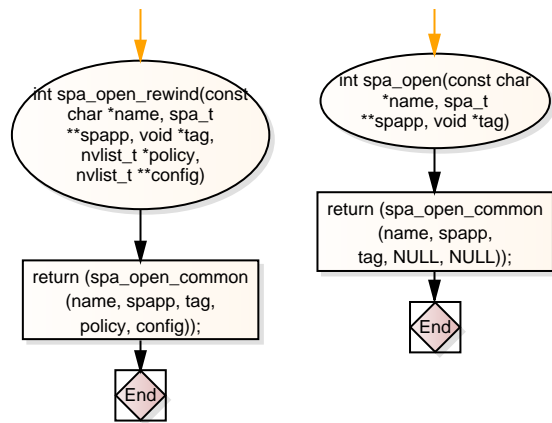
Restore the  
initial load info











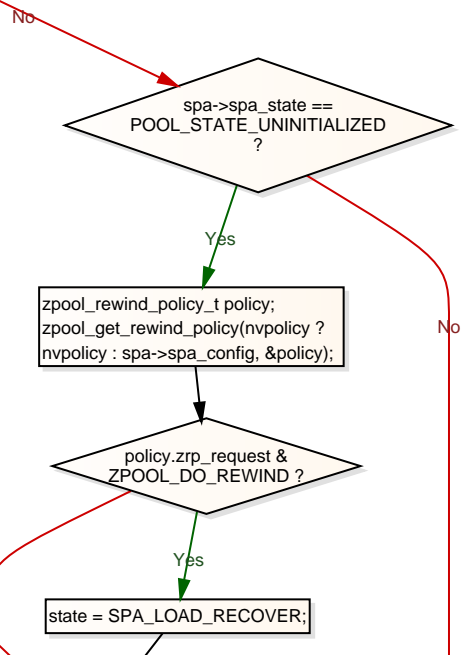
\* Pool Open/Import

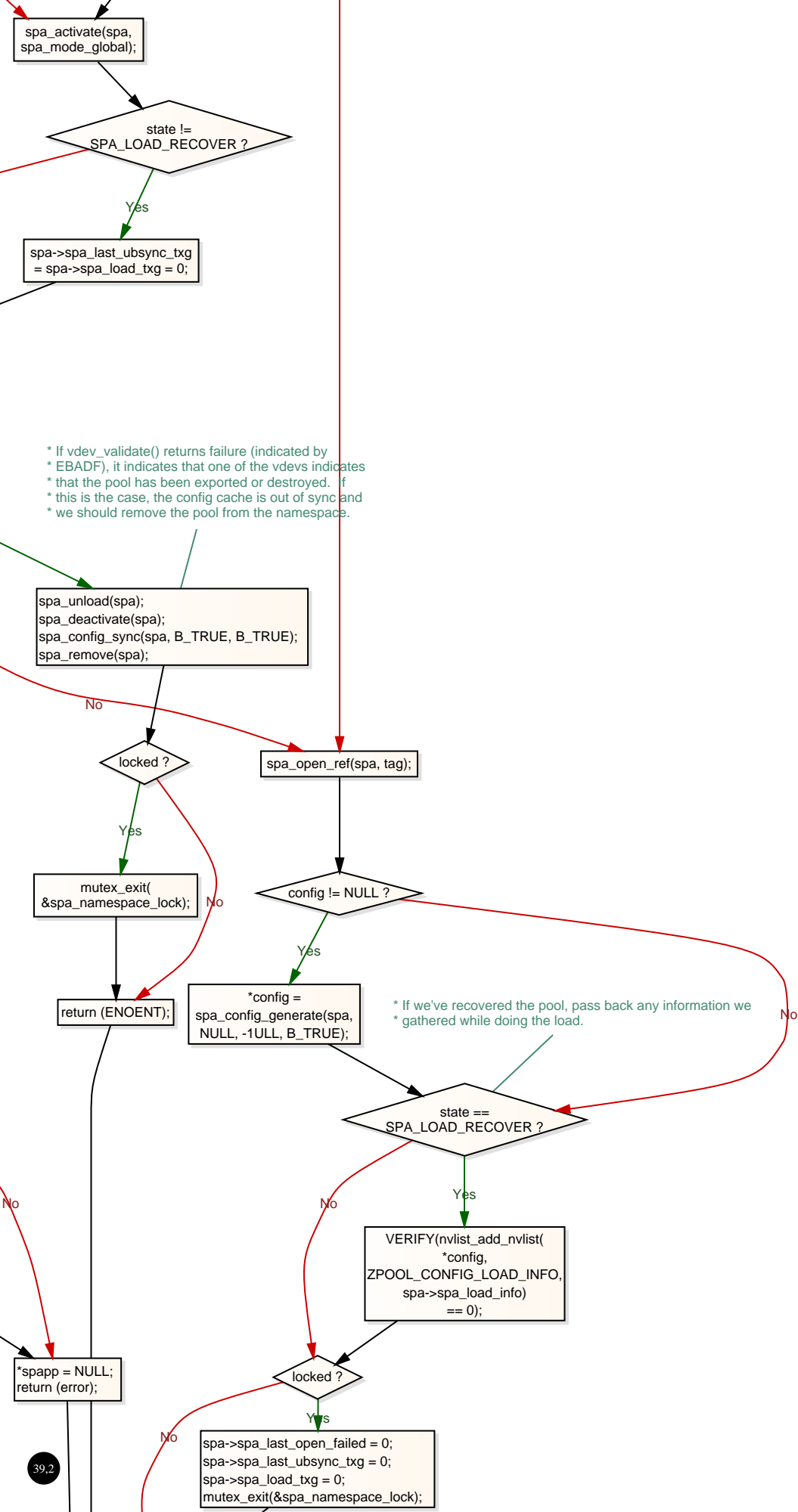
\* The import case is identical to an open except that the configuration is sent down from userland, instead of grabbed from the configuration cache. For the case of an open, the pool configuration will exist in the POOL\_STATE\_UNINITIALIZED state.

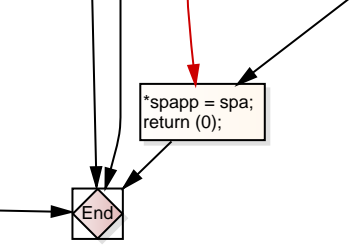
\* The stats information (gen/count/ustats) is used to gather vdev statistics at the same time open the pool, without having to keep around the spa\_t in some ambiguous state.

\* As disgusting as this is, we need to support recursive calls to this function because dsl\_dir\_open() is called during spa\_load(), and ends up calling spa\_open() again. The real fix is to figure out how to avoid dsl\_dir\_open() calling this in the first place.

```
mutex_enter(&spa_namespace_lock);
locked = B_TRUE;
```

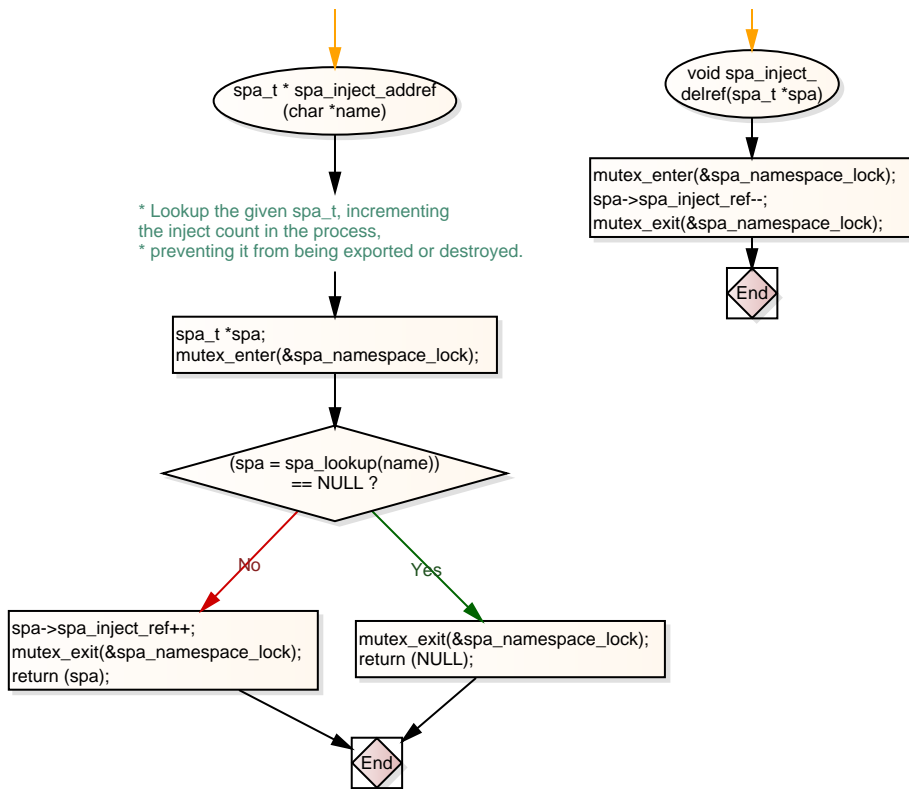








\* Lookup the given spa\_t, incrementing  
the inject count in the process,  
\* preventing it from being exported or destroyed.



static void  
spa\_add\_spares(spa\_t  
\*spa, nvlist\_t \*config)

\* Add spares device  
information  
to the nvlist.

```
nvlist_t **spares;  
uint_t i, nspares;  
nvlist_t *nvroot;  
uint64_t guid;  
vdev_stat_t *vs;  
uint_t vsc;  
uint64_t pool;  
ASSERT(spa_config_held(spa,  
SCL_CONFIG, RW_READER));
```

spa->spa\_spares.sav\_count  
== 0 ?

Yes

No

```
VERIFY(nvlist_lookup_nvlist(config,  
ZPOOL_CONFIG_VDEV_TREE, &nvroot) == 0);  
VERIFY(nvlist_lookup_nvlist_array(spa->spa_spares.sav_config,  
ZPOOL_CONFIG_SPARES, &spares, &nspares) == 0);
```

nspares != 0 ?

Yes

No

```
VERIFY(nvlist_add_nvlist_array(nvroot,  
ZPOOL_CONFIG_SPARES, spares, nspares) == 0);  
VERIFY(nvlist_lookup_nvlist_array(nvroot,  
ZPOOL_CONFIG_SPARES, &spares, &nspares) == 0);
```

\* Go through and find any spares which have since been  
\* repurposed as an active spare. If this is the case, update  
\* their status appropriately.

i = 0

i < nspares ?

Yes

```
VERIFY(nvlist_lookup_  
uint64(spares[i],  
ZPOOL_CONFIG_GUID,  
&guid) == 0);
```

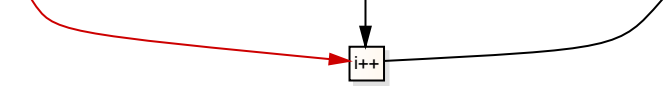
spa\_spare\_exists(guid,  
&pool, NULL) &&  
pool != 0ULL ?

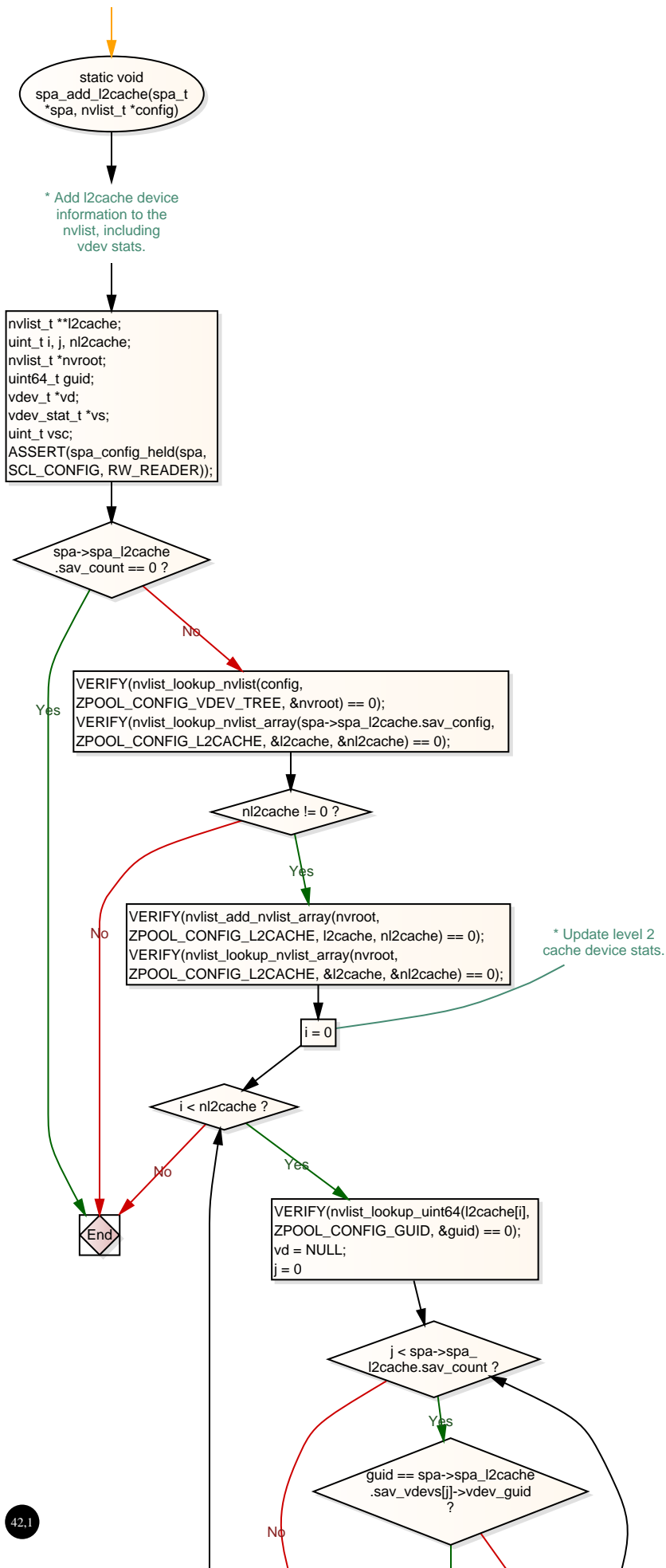
Yes

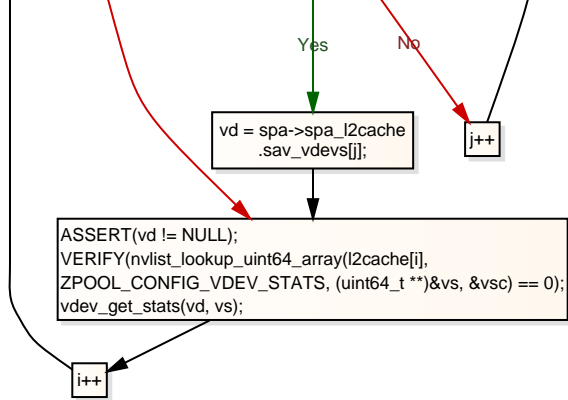
No

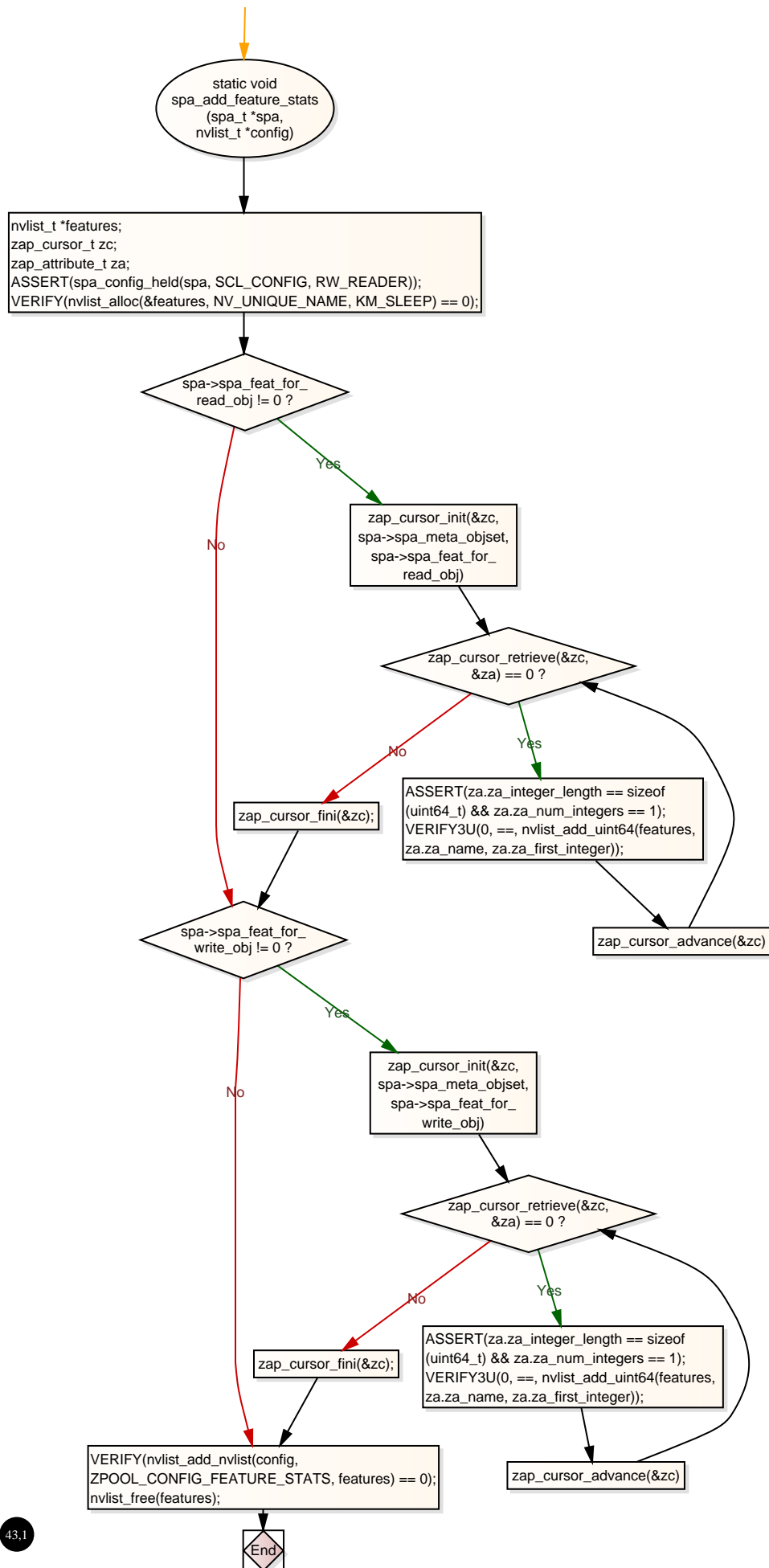
End

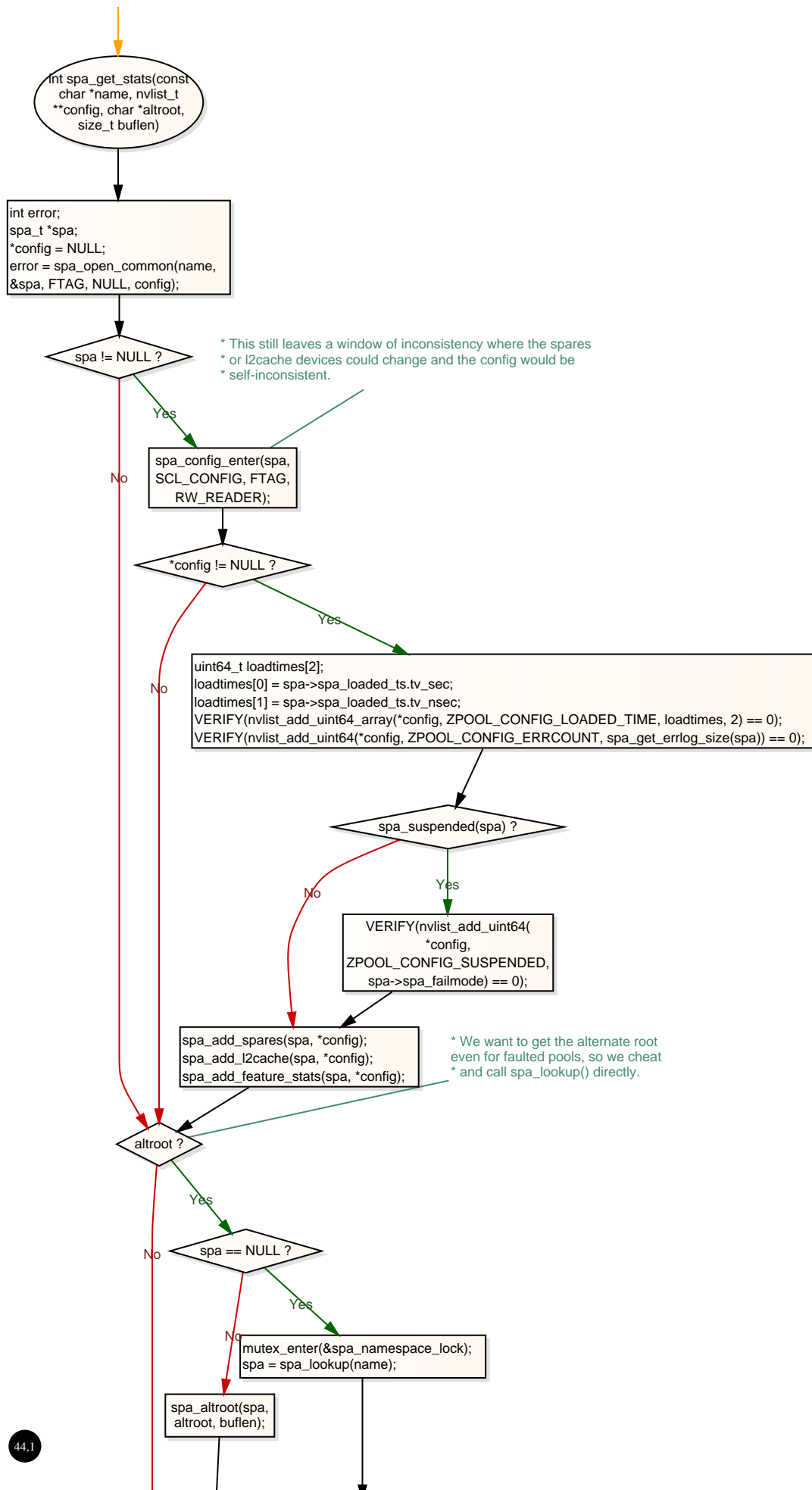
```
VERIFY(nvlist_lookup_uint64_array(spares[i],  
ZPOOL_CONFIG_VDEV_STATS, (uint64_t **)&vs, &vsc) == 0);  
vs->vs_state = VDEV_STATE_CANT_OPEN;  
vs->vs_aux = VDEV_AUX_SPARED;
```

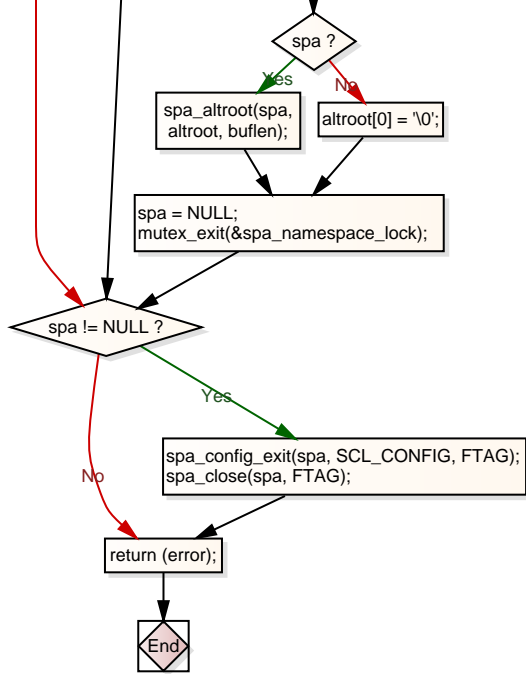




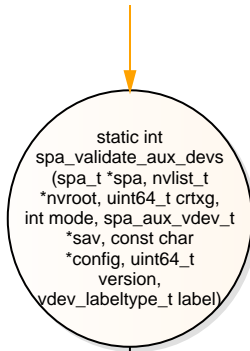




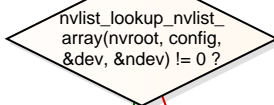
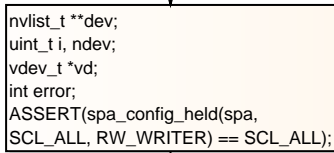




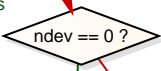




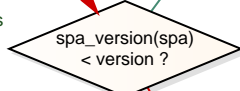
\* Validate that the auxiliary device array is well formed. We must have an  
 \* array of nvlists, each which describes a valid leaf vdev. If this is an  
 \* import (mode is VDEV\_ALLOC\_SPARE), then we allow corrupted spares to be  
 \* specified, as long as they are well-formed.



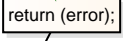
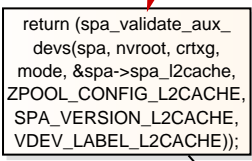
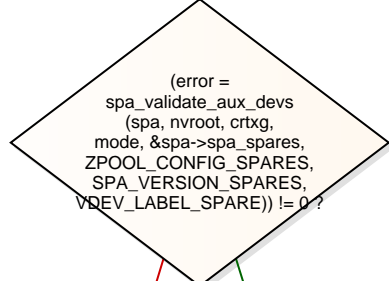
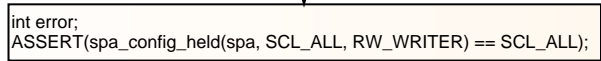
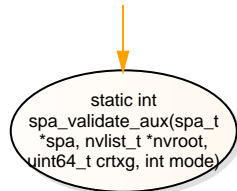
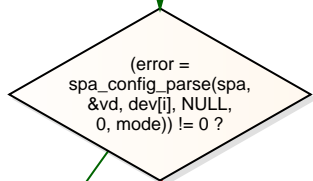
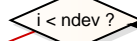
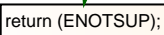
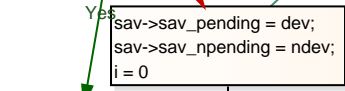
\* It's acceptable to have no devs specified.

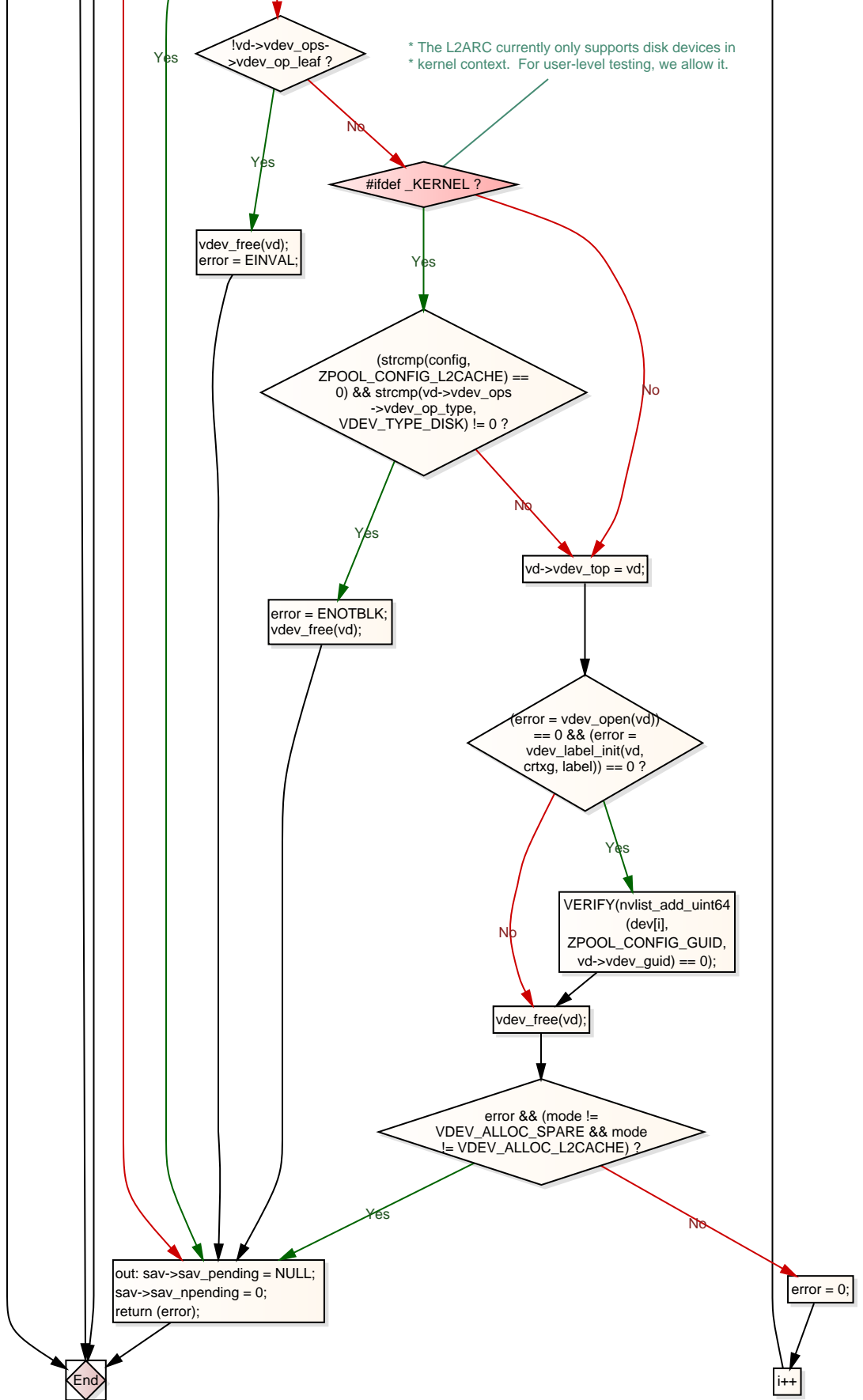


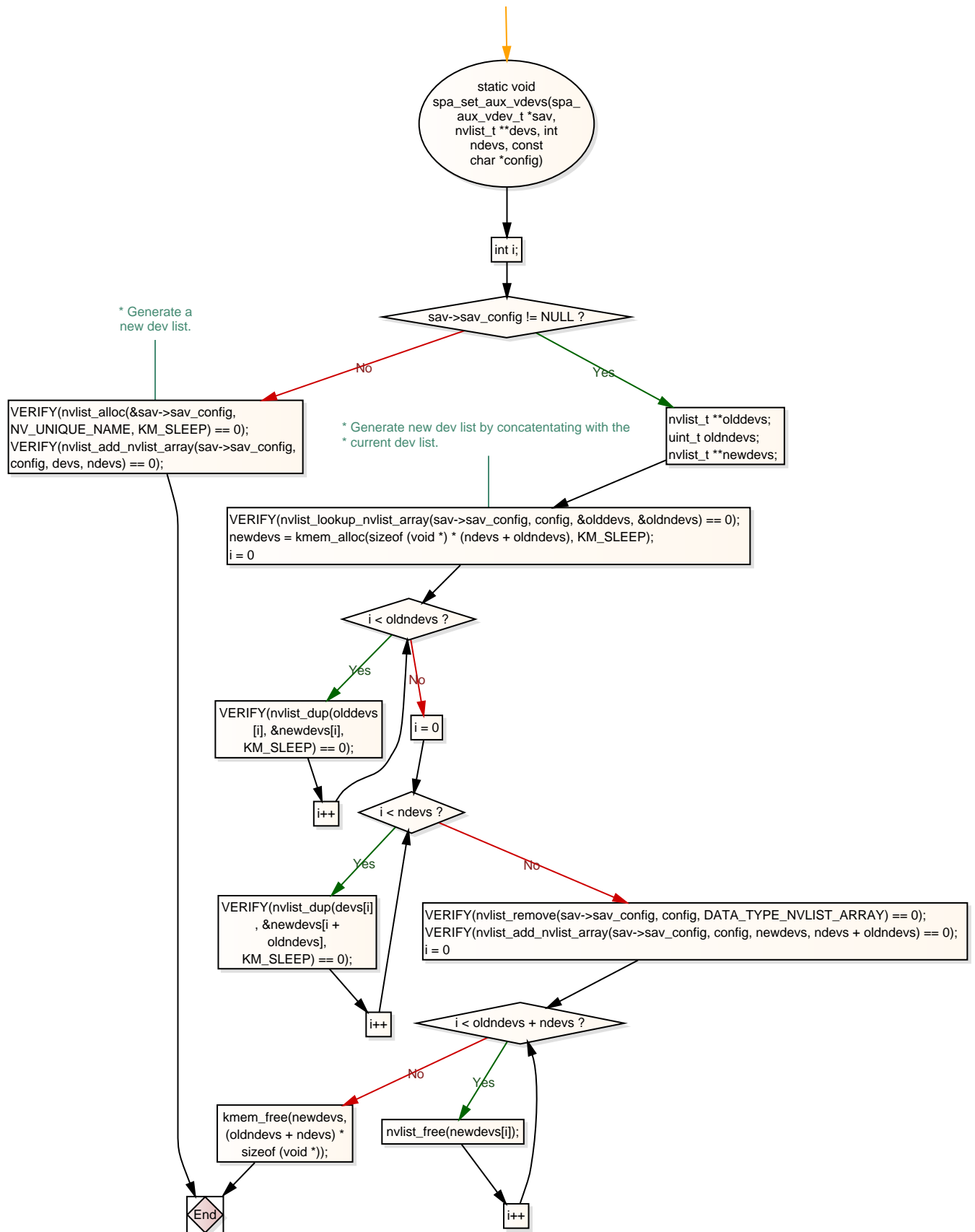
\* Make sure the pool is formatted with a version that supports this  
 \* device type.

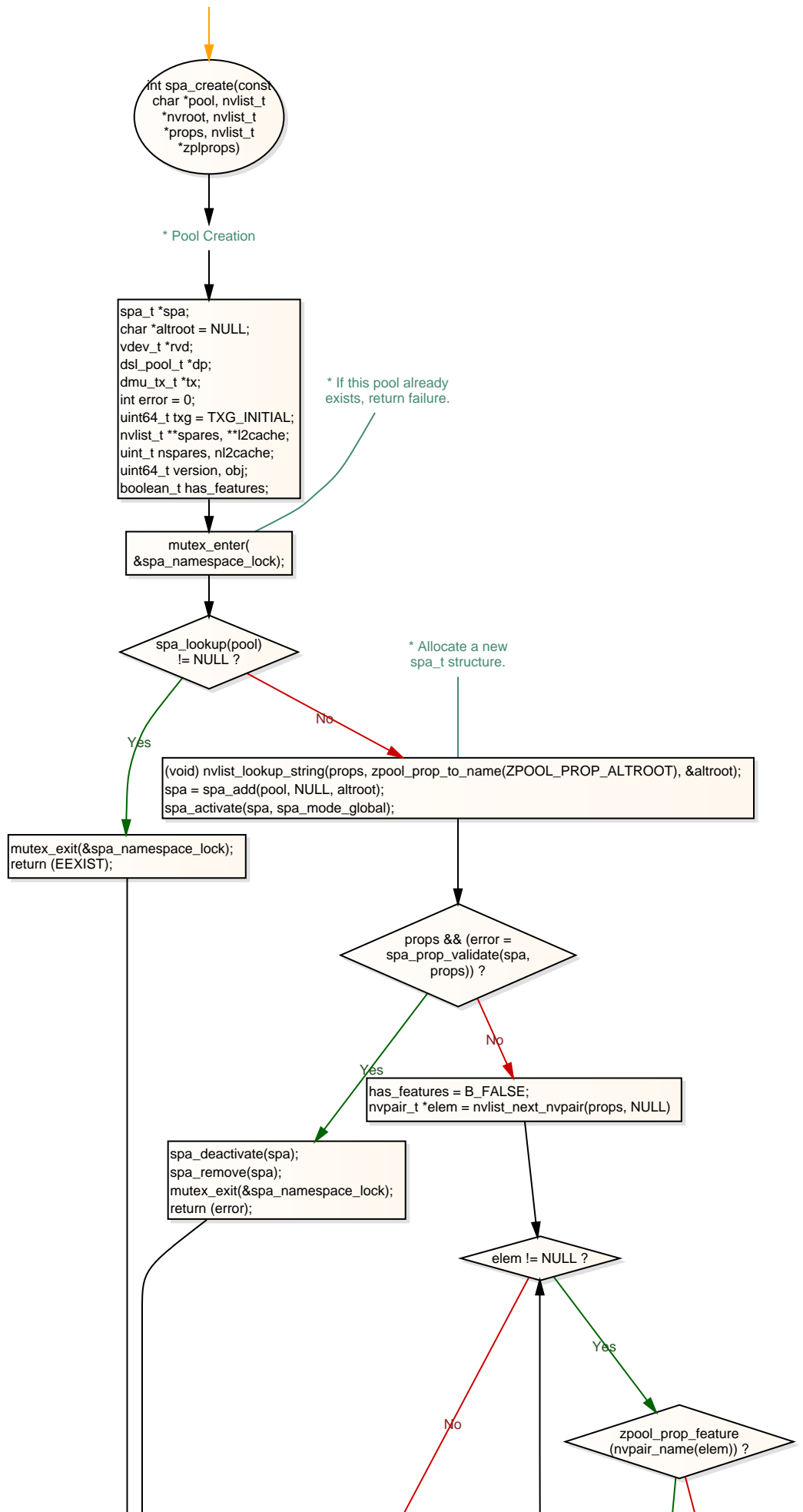
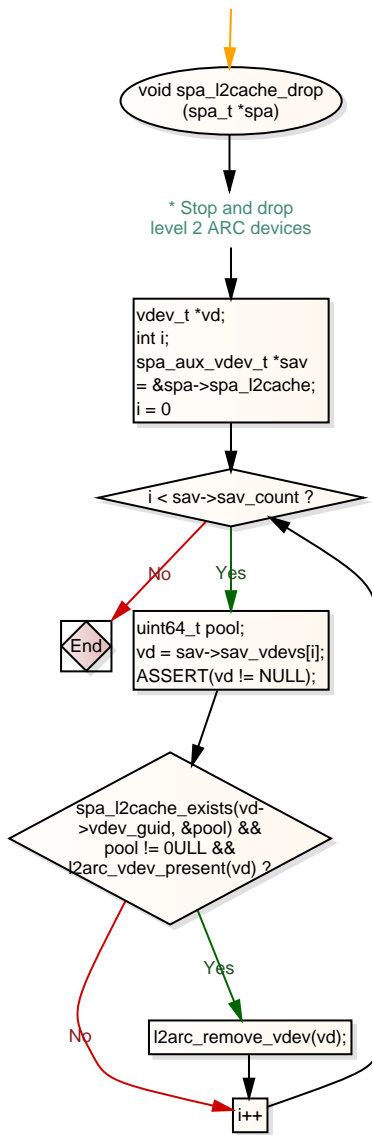


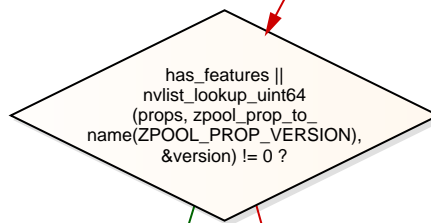
\* Set the pending device list so we correctly handle device in-use  
 \* checking.











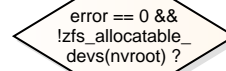
version = SPA\_VERSION;

ASSERT(SPA\_VERSION\_IS\_SUPPORTED(version));  
spa->spa\_first\_txg = txg;  
spa->spa\_uberblock.ub\_txg = txg - 1;  
spa->spa\_uberblock.ub\_version = version;  
spa->spa\_ubsync = spa->spa\_uberblock;

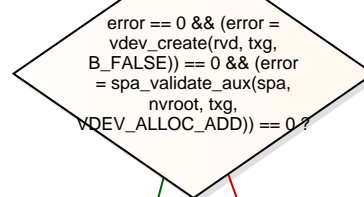
spa->spa\_async\_zio\_root  
= zio\_root(spa, NULL,  
NULL, ZIO\_FLAG\_CANFAIL |  
ZIO\_FLAG\_SPECULATIVE |  
ZIO\_FLAG\_GODFATHER);

\* Create the root vdev.

spa\_config\_enter(spa, SCL\_ALL, FTAG, RW\_WRITER);  
error = spa\_config\_parse(spa, &rvd, nvroot, NULL, 0, VDEV\_ALLOC\_ADD);  
ASSERT(error != 0 || rvd != NULL);  
ASSERT(error != 0 || spa->spa\_root\_vdev == rvd);



error = EINVAL;

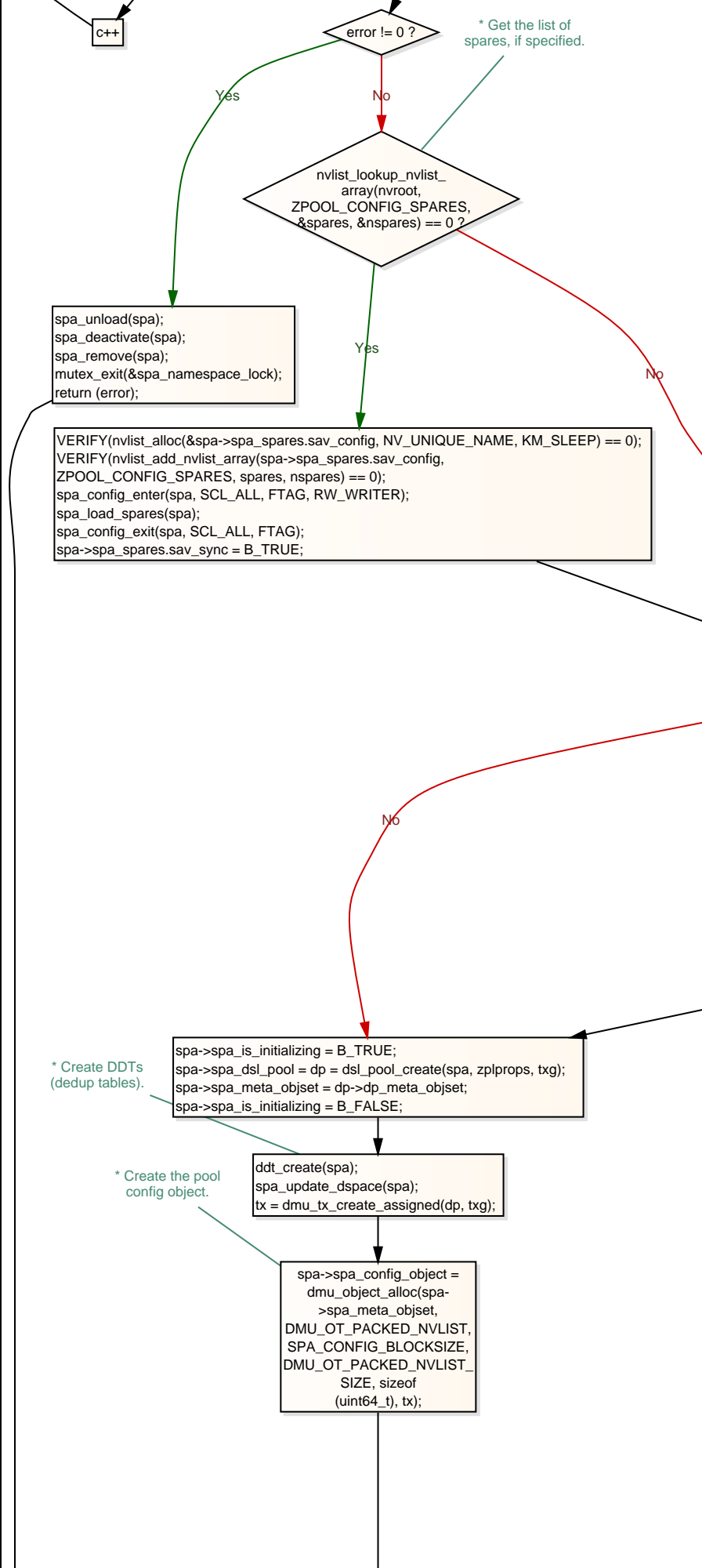


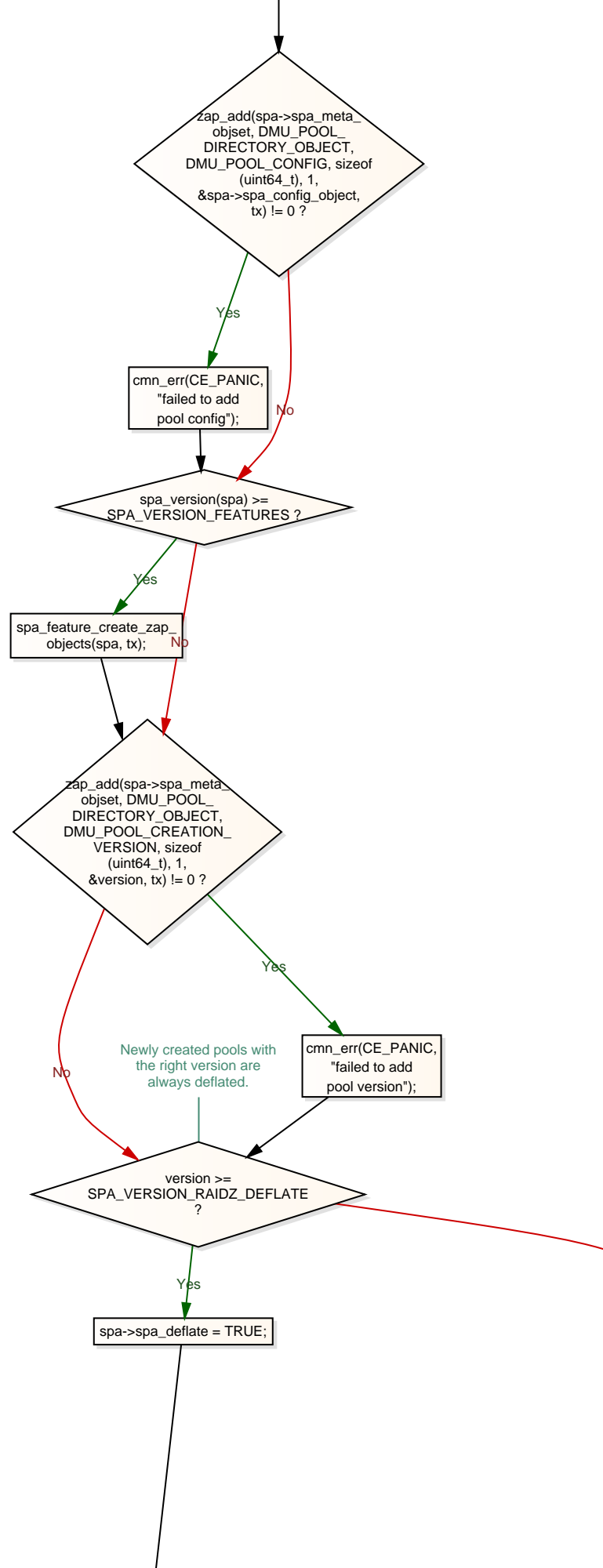
int c = 0

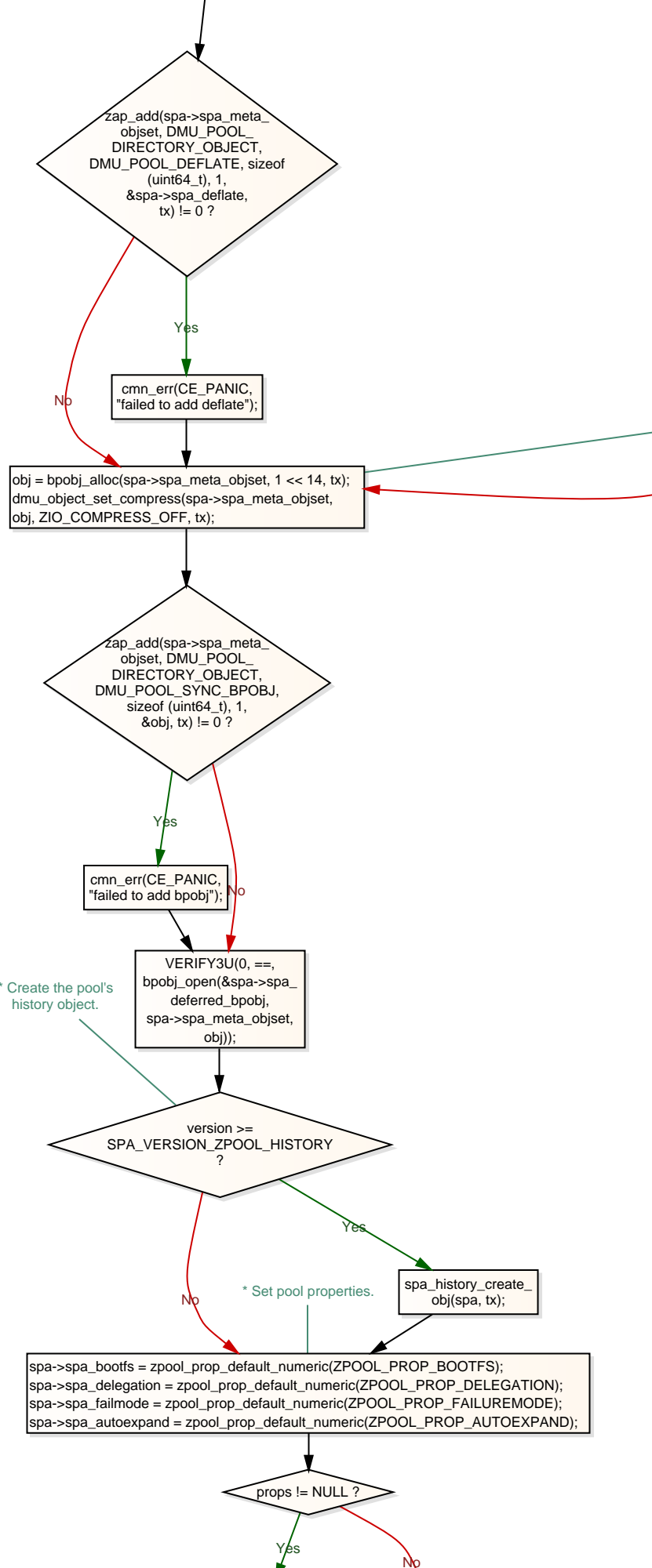


vdev metaslab\_set\_size(rvd->vdev\_child[c]);  
vdev\_expand(rvd->vdev\_child[c], txg);

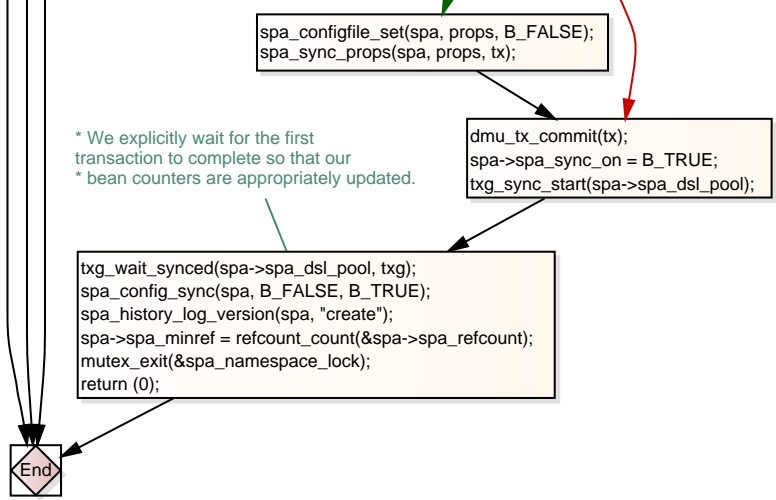
spa\_config\_exit(spa,  
SCL\_ALL, FTAG);

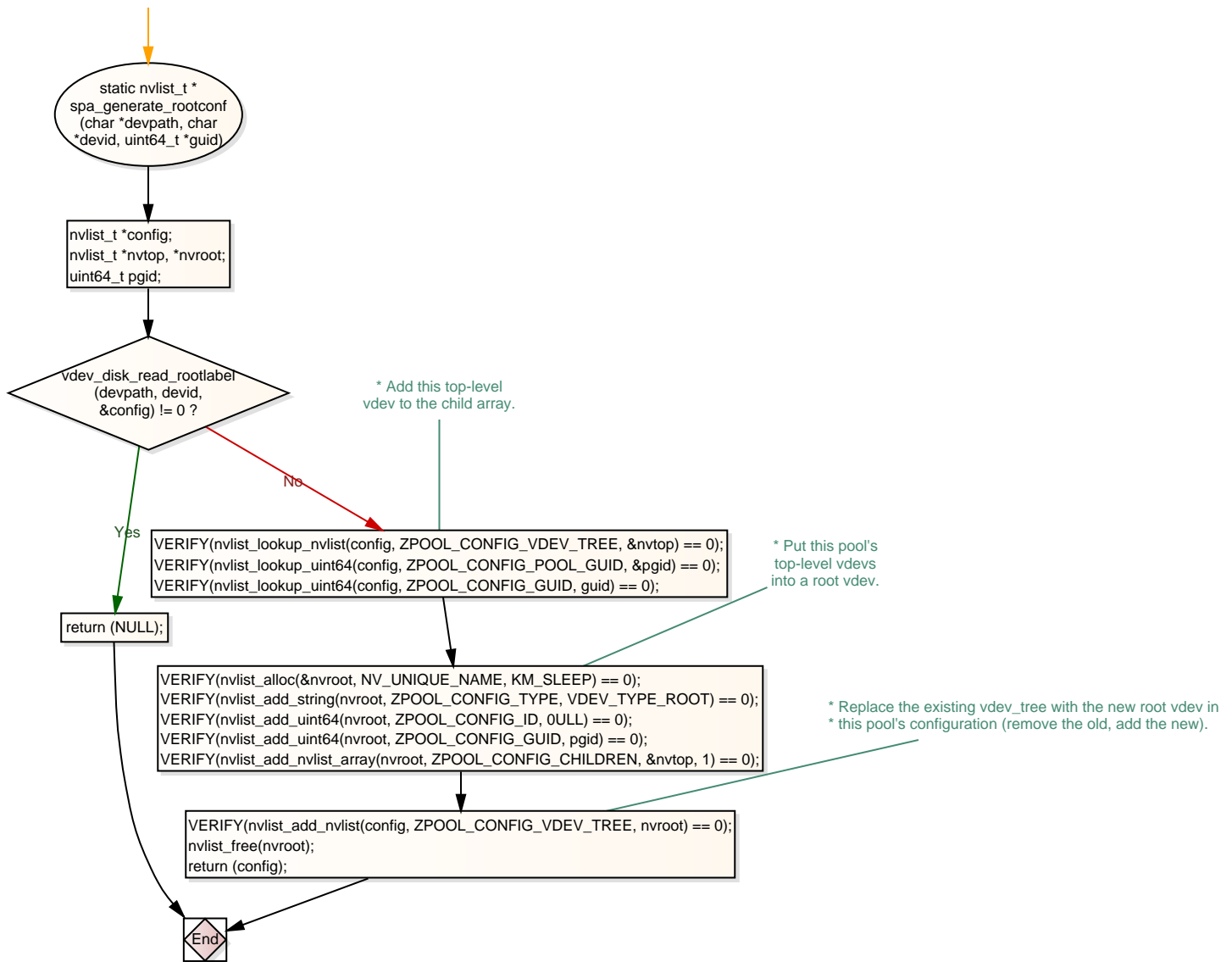












\* Create "The Godfather"  
zio to hold all async IOs

\* Get the list of level  
2 cache devices,  
if specified.

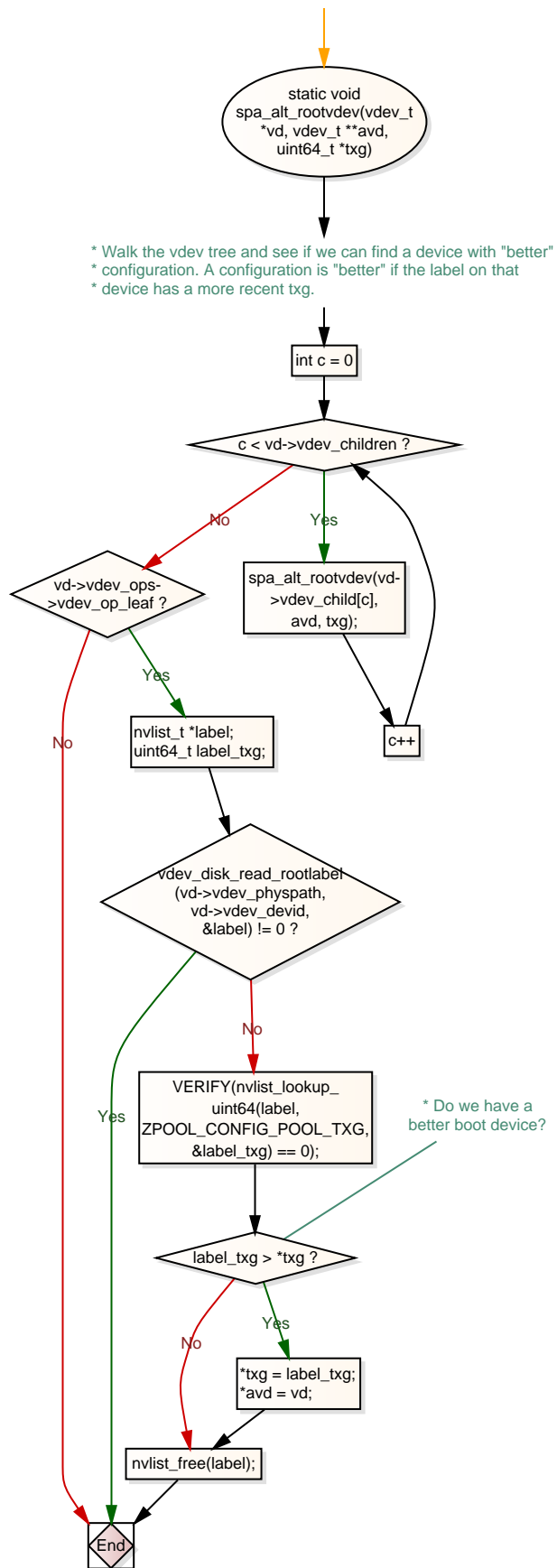
nvlist\_lookup\_nvlist\_  
array(nvroot,  
ZPOOL\_CONFIG\_L2CACHE,  
&l2cache,  
&nl2cache) == 0 ?

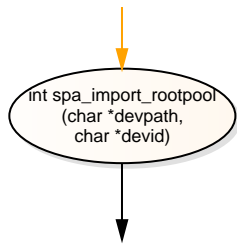
Yes

```
VERIFY(nvlist_alloc(&spa->spa_l2cache.sav_config, NV_UNIQUE_NAME, KM_SLEEP) == 0);  
VERIFY(nvlist_add_nvlist_array(spa->spa_l2cache.sav_config,  
ZPOOL_CONFIG_L2CACHE, l2cache, nl2cache) == 0);  
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);  
spa_load_l2cache(spa);  
spa_config_exit(spa, SCL_ALL, FTAG);  
spa->spa_l2cache.sav_sync = B_TRUE;
```



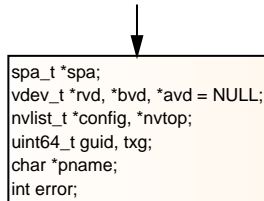
\* Create the deferred-free bpobj. Turn off compression  
\* because sync-to-convergence takes longer if the blocksize  
\* keeps changing.



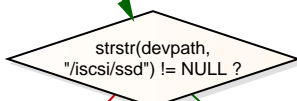
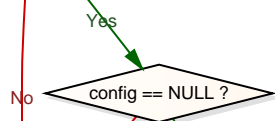
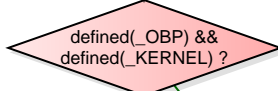
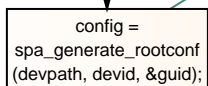


\* Import a root pool.

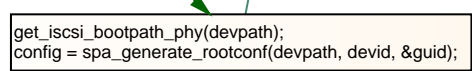
\*  
\* For x86. devpath\_list will consist of devid and/or physpath name of  
\* the vdev (e.g. "id1,sd@SSEAGATE..." or "/pci@1f,0/ide@d/disk@0,0:a").  
\* The GRUB "findroot" command will return the vdev we should boot.  
\*  
\* For Sparc, devpath\_list consists the physpath name of the booting device  
\* no matter the rootpool is a single device pool or a mirrored pool.  
\* e.g.  
\* "/pci@1f,0/ide@d/disk@0,0:a"



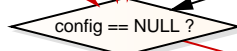
\* Read the label from  
the boot device and  
generate a configuration.



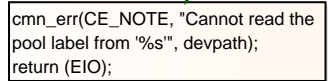
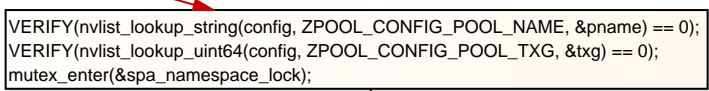
iscsi boot



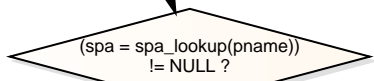
atim import - Take a pool and insert it into the namespace  
it had been loaded at boot.



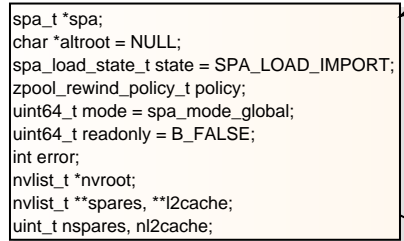
No



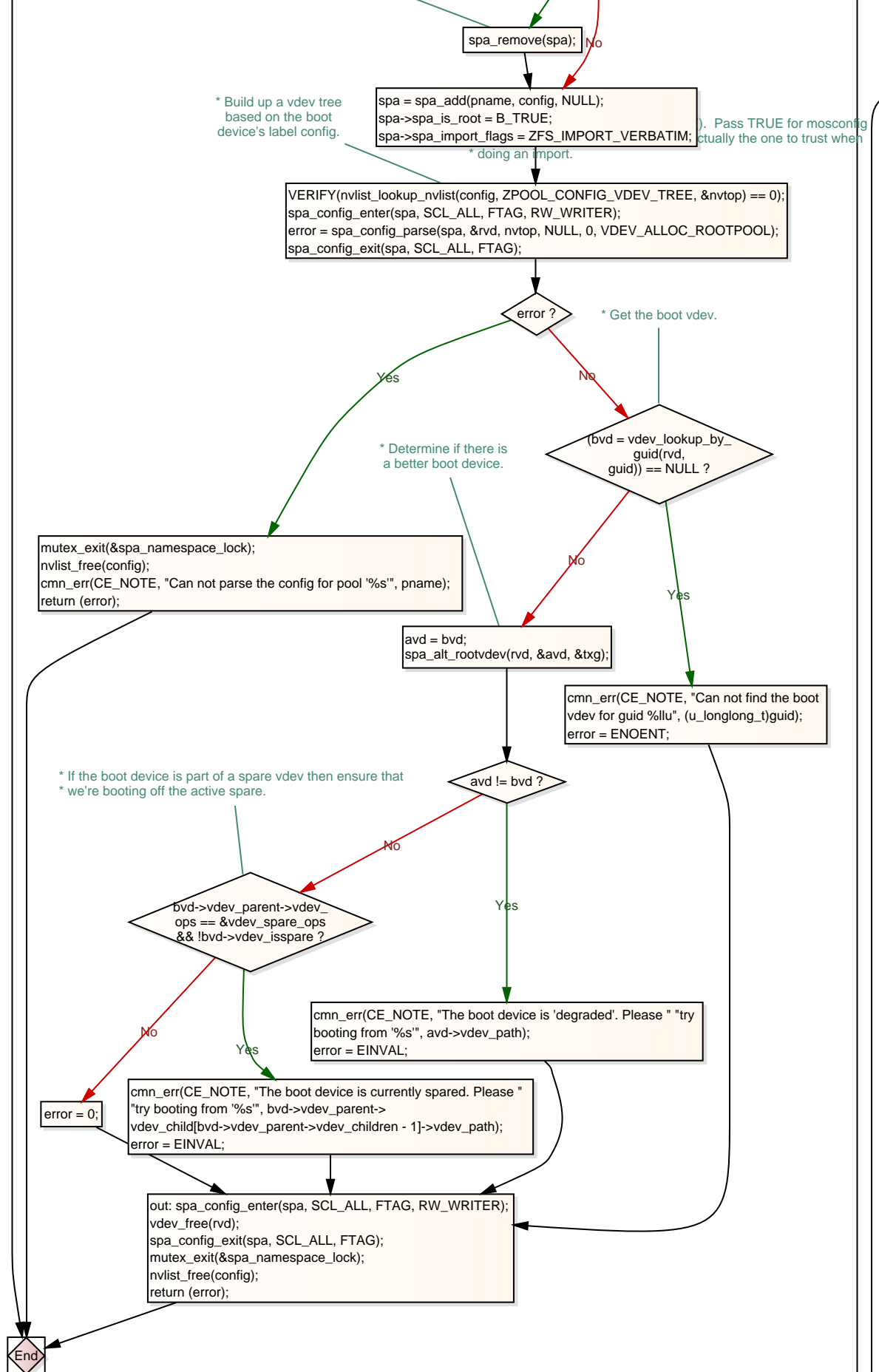
\* Remove the existing root pool from the namespace so that we  
\* can replace it with the correct config we just read in.



Yes









\* Check for any removed devices.

\* It's possible that the pool was expanded while it was exported.  
\* We kick off an async task to handle this for us.

```
int spa_import(const
char *pool, nvlist_t
*config, nvlist_t
*props, uint64_t flags)
```

\* Import a non-root pool  
into the system.

\* If a pool with this  
name exists,  
return failure.

```
mutex_enter(
&spa_namespace_lock);
```

```
spa_lookup(pool)
!= NULL ?
```

No

\* Create and initialize  
the spa structure.

```
(void) nvlist_lookup_string(props,
zpool_prop_to_name(ZPOOL_PROP_ALTRoot), &altroot);
(void) nvlist_lookup_uint64(props,
zpool_prop_to_name(ZPOOL_PROP_READONLY), &readonly);
```

```
mutex_exit(&spa_namespace_lock);
return (EEXIST);
```

```
readonly ?
```

Yes

No

```
mode = FREAD;
```

```
spa = spa_add(pool, config, altroot);
spa->spa_import_flags = flags;
```

```
spa->spa_import_flags &
ZFS_IMPORT_VERBATIM ?
```

Yes

No

```
spa_activate(spa, mode);
```

\* Don't start async  
tasks until we know  
everything is healthy.

```
props != NULL ?
```

Yes

No

```
spa_configfile_set(spa,
props, B_FALSE);
```

```
spa_async_suspend(spa);
zpool_get_rewind_policy(config, &policy);
```

```
nvlist_lookup_uint64
(tryconfig,
ZPOOL_CONFIG_POOL_STATE,
&state) ?
```

Yes

No

No

```
return (NULL);
```

Yes

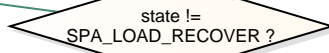
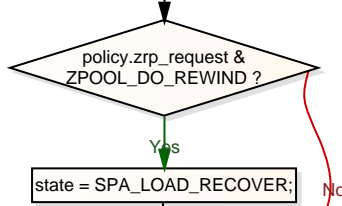
```
mutex_enter(&spa_namespace_lock);
spa = spa_add(TRYIMPORT_NAME, tryconfig, NULL);
spa_activate(spa, FREAD);
```

```
return (NULL);
```

```

spa_config_sync(spa, B_FALSE, B_TRUE);
mutex_exit(&spa_namespace_lock);
spa_history_log_version(spa, "import");
return (0);

```



```

spa->spa_last_ubsync_txg
= spa->spa_load_txg = 0;

```

```

error = spa_load_best
(spa, state, B_TRUE,
policy.zrp_txg,
policy.zrp_request);

```

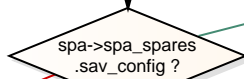
\* Propagate anything learned while loading the pool and pass it  
\* back to caller (i.e. rewind info, missing devices, etc).

```

VERIFY(nvlist_add_nvlist(config,
ZPOOL_CONFIG_LOAD_INFO, spa->spa_load_info) == 0);
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);

```

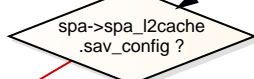
\* Toss any existing spareslist, as it doesn't have any validity  
\* anymore, and conflicts with spa\_has\_spare().



```

nvlist_free(spa->spa_spare.sav_config);
spa->spa_spare.sav_config = NULL;
spa_load_spare(spa);

```



```

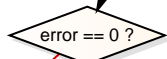
nvlist_free(spa->spa_l2cache.sav_config);
spa->spa_l2cache.sav_config = NULL;
spa_load_l2cache(spa);

```

```

VERIFY(nvlist_lookup_
nvlist(config,
ZPOOL_CONFIG_VDEV_TREE,
&nvroot) == 0);

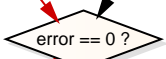
```

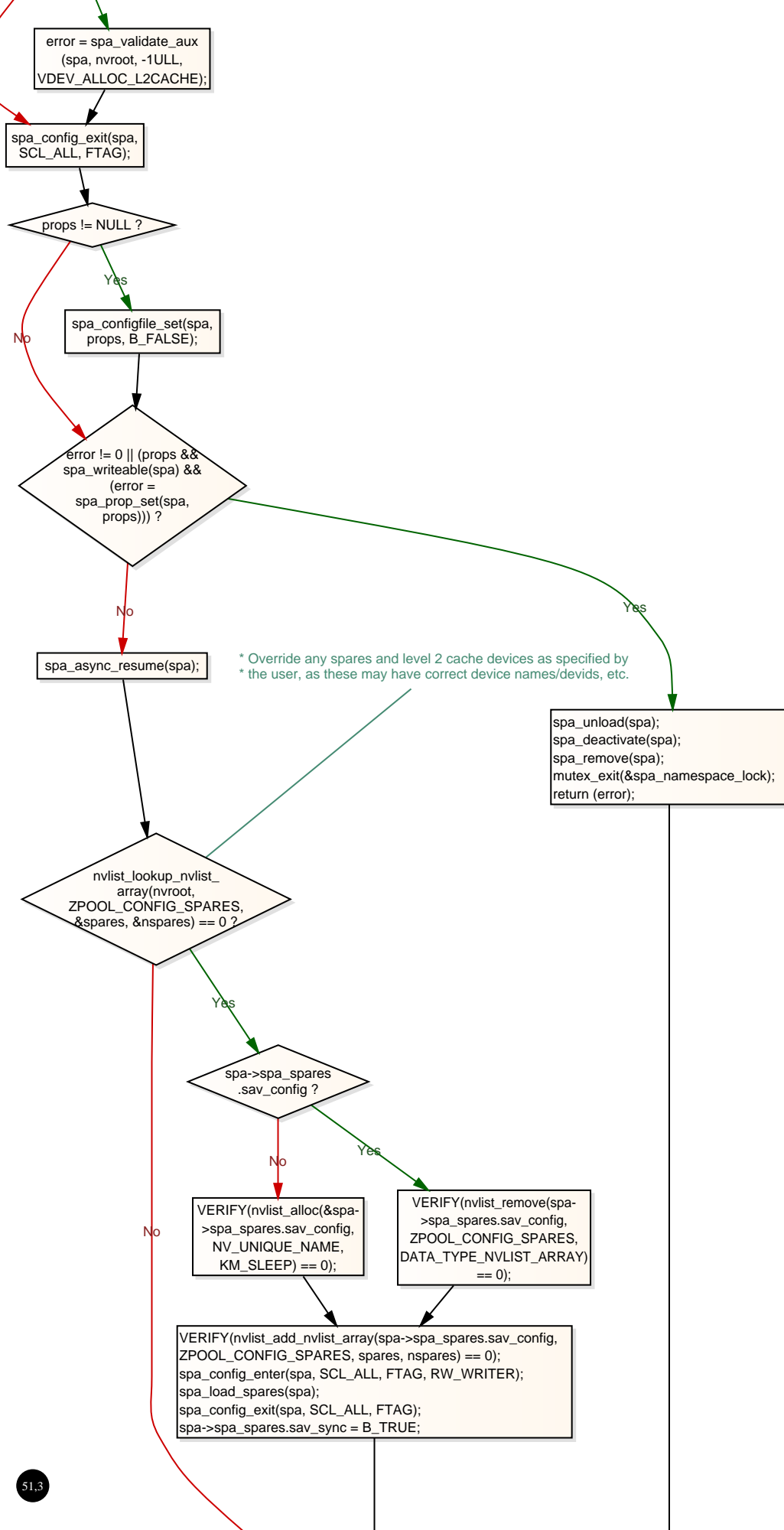


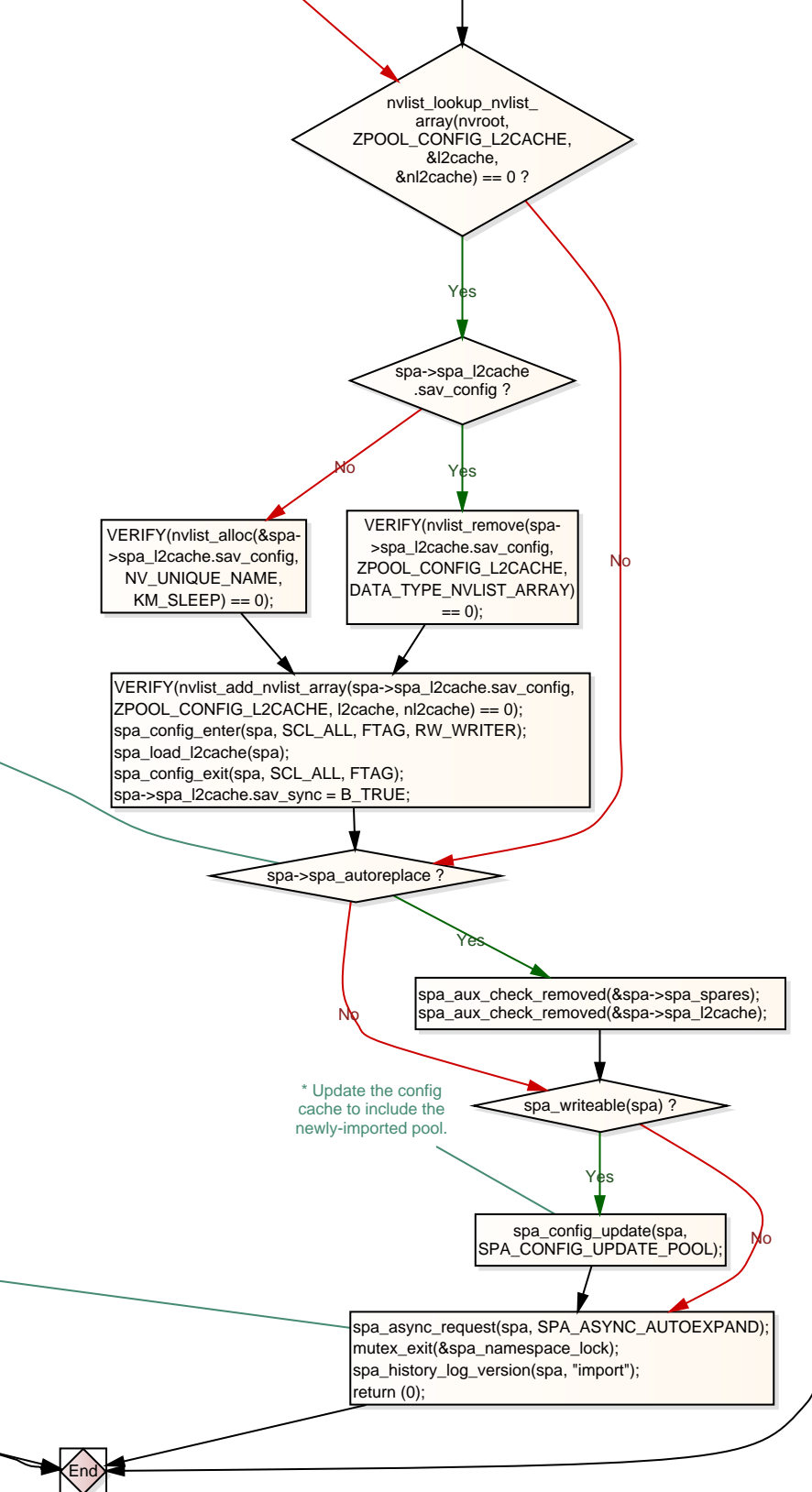
```

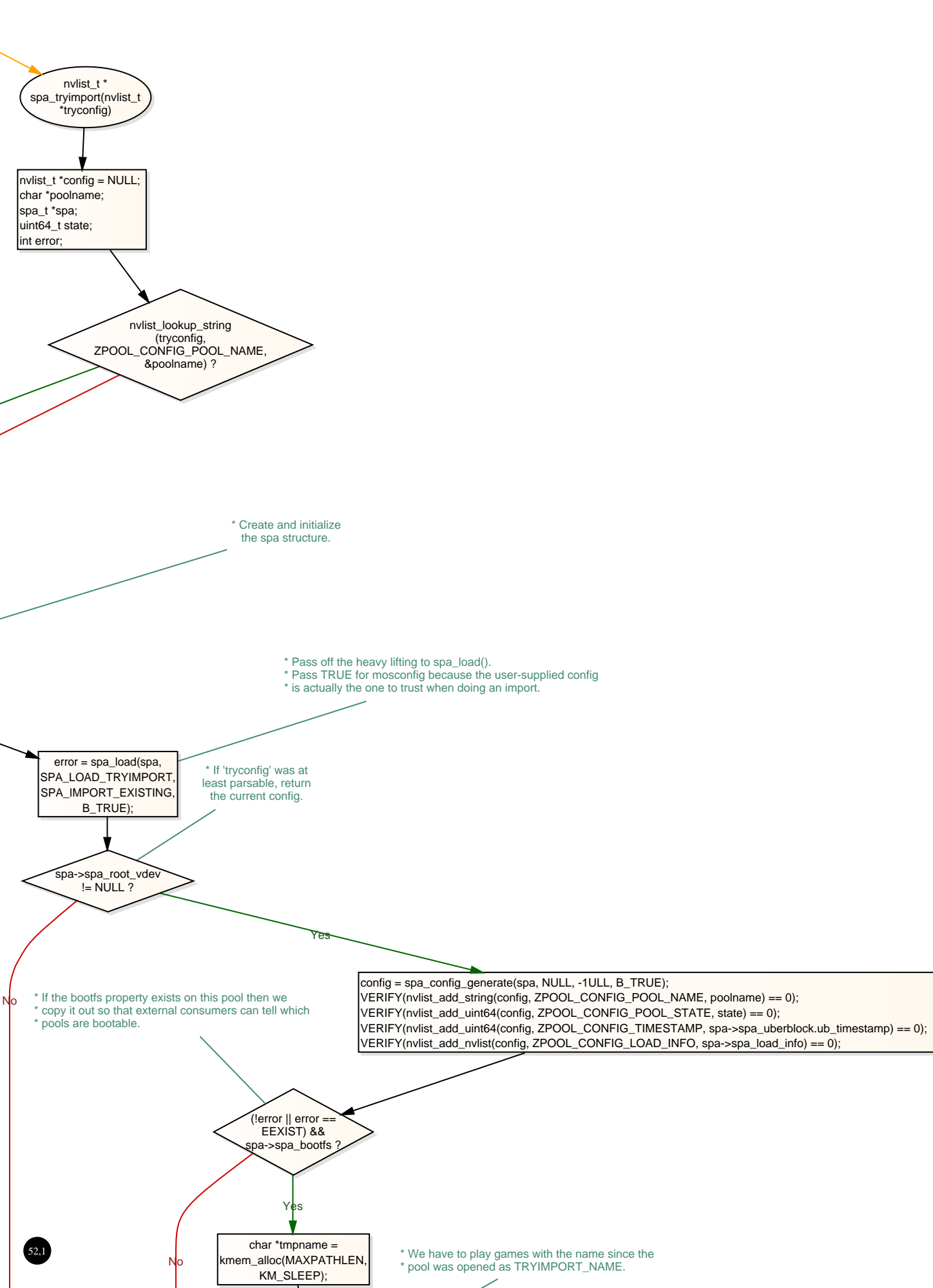
error = spa_validate_aux
(spa, nvroot, -1ULL,
VDEV_ALLOC_SPARE);

```

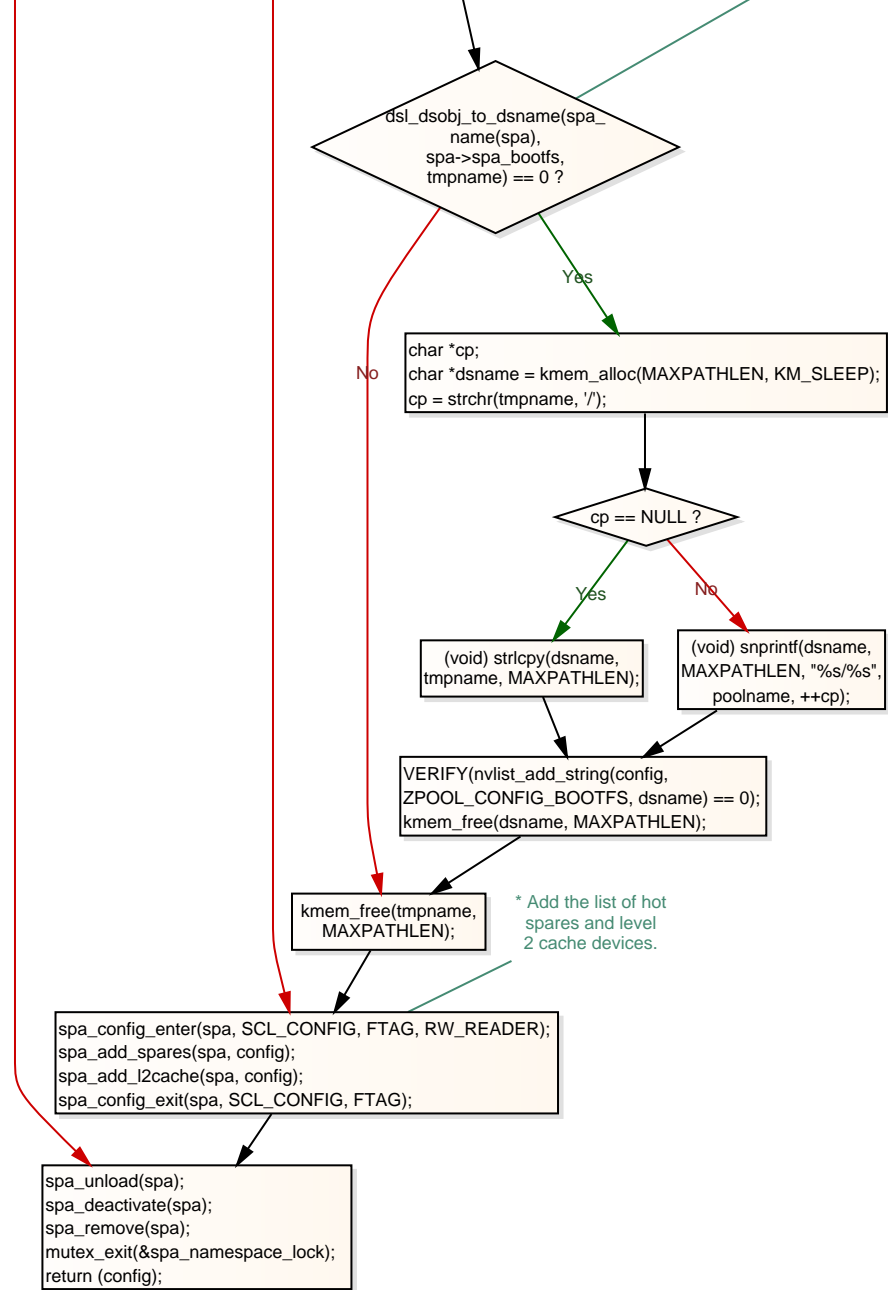


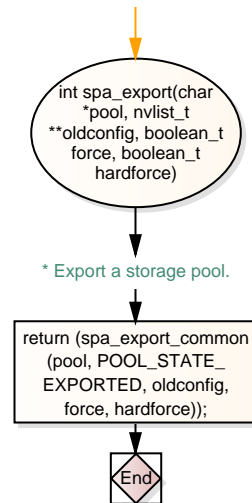
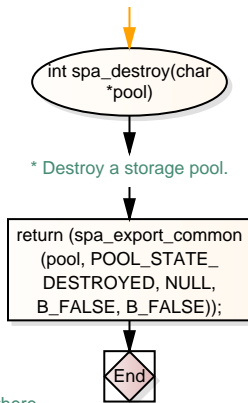
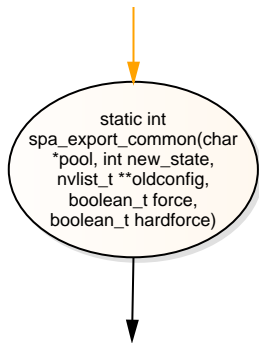






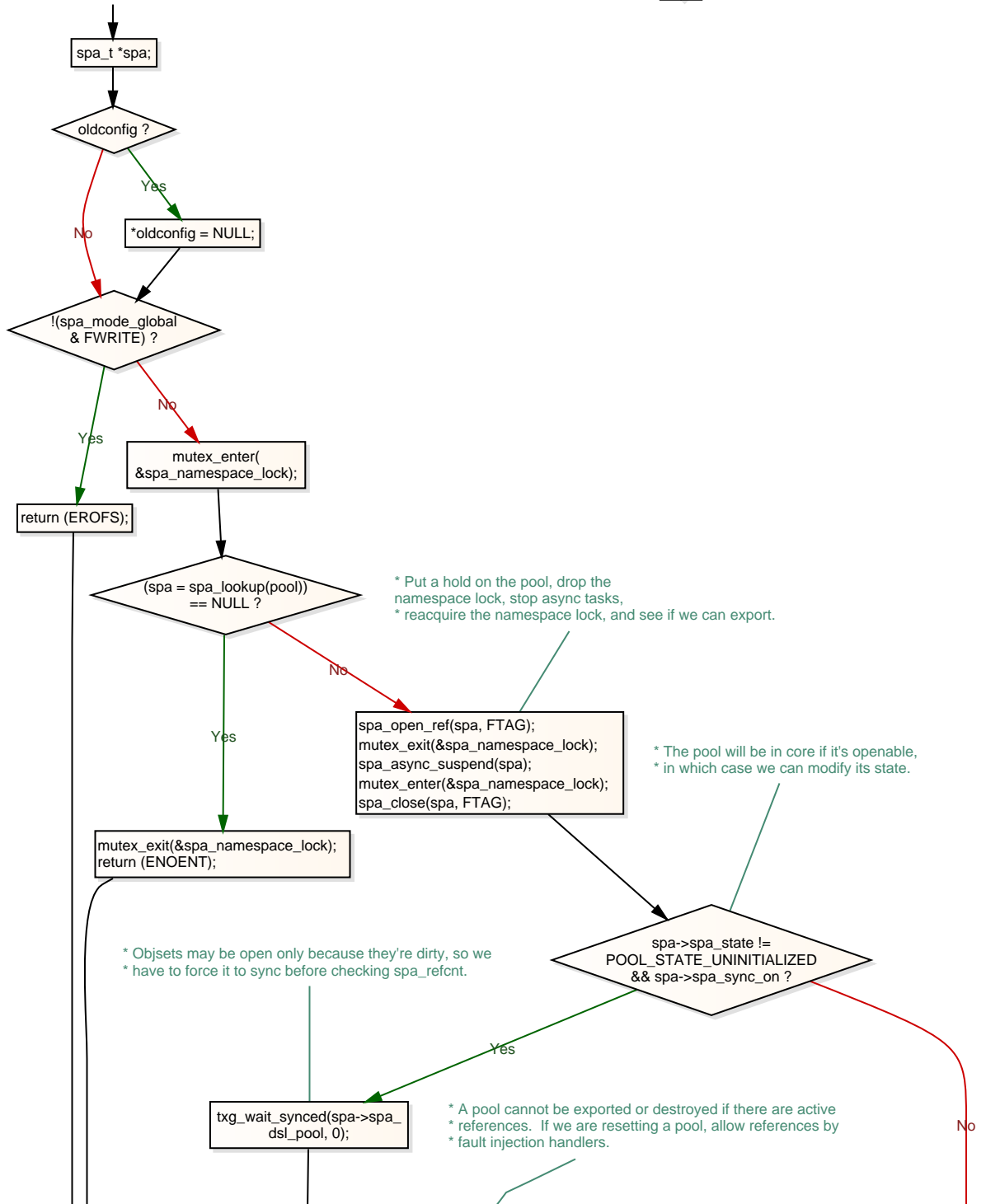


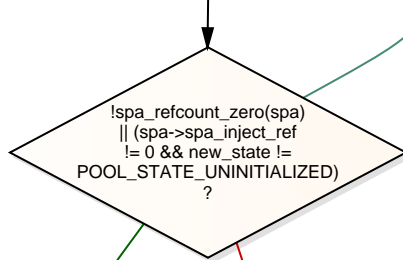




\* Pool export/destroy

\* The act of destroying or exporting a pool is very simple. We make sure there is no more pending I/O and any references to the pool are gone. Then, we update the pool state and sync all the labels to disk, removing the configuration from the cache afterwards. If the 'hardforce' flag is set, then we don't sync the labels or remove the configuration cache.



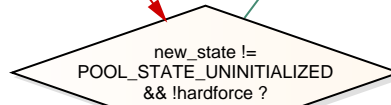


\* A pool cannot be exported if it has an active shared spare.  
\* This is to prevent other pools stealing the active spare  
\* from an exported pool. At user's own will, such pool can  
\* be forcibly exported.



\* We want this to be reflected on every label,  
\* so mark them all dirty. spa\_unload() will do the  
\* final sync that pushes these changes out.

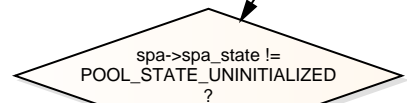
spa\_async\_resume(spa);  
mutex\_exit(&spa\_namespace\_lock);  
return (EBUSY);



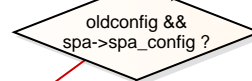
spa\_async\_resume(spa);  
mutex\_exit(&spa\_namespace\_lock);  
return (EXDEV);

spa\_config\_enter(spa, SCL\_ALL, FTAG, RW\_WRITER);  
spa->spa\_state = new\_state;  
spa->spa\_final\_txg = spa\_last\_synced\_txg(spa) + TXG\_DEFER\_SIZE + 1;  
vdev\_config\_dirty(spa->spa\_root\_vdev);  
spa\_config\_exit(spa, SCL\_ALL, FTAG);

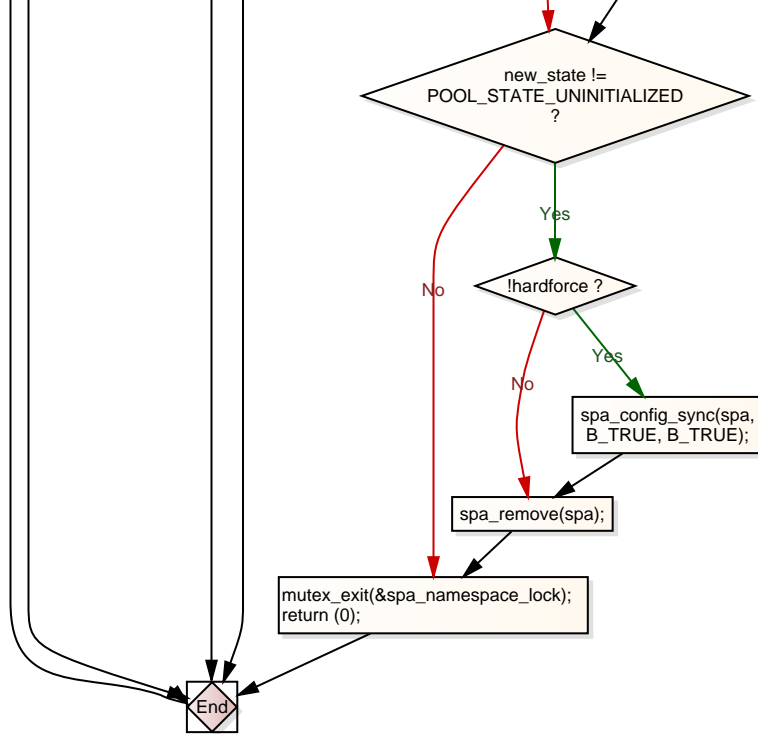
spa\_event\_notify(spa,  
NULL, ESC\_  
ZFS\_POOL\_DESTROY);

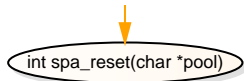


spa\_unload(spa);  
spa\_deactivate(spa);



VERIFY(nvlist\_dup(spa-  
>spa\_config,  
oldconfig, 0) == 0);

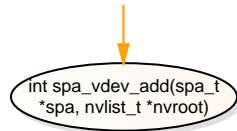




\* Similar to spa\_export(), this unloads the spa\_t without actually removing it from the namespace in any way.

```

return (spa_export_common
(pool, POOL_STATE_
UNINITIALIZED, NULL,
B_FALSE, B_FALSE));
  
```



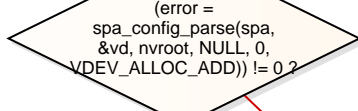
```

int spa_vdev_add(spa_t
*spa, nvlist_t *nvroot)
  
```

\* =====  
\* Device manipulation  
\* =====  
\* Add a device to a storage pool.

```

uint64_t txg, id;
int error;
vdev_t *vrd = spa->spa_root_vdev;
vdev_t *vd, *tvd;
nvlist_t **spares, **l2cache;
uint_t nspares, nl2cache;
ASSERT(spa_writeable(spa));
txg = spa_vdev_enter(spa);
  
```



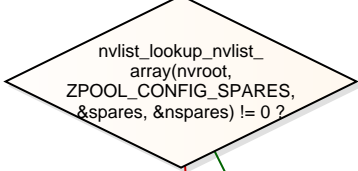
spa\_vdev\_exit()  
will clear this

```

spa->spa_pending_vdev
= vd;
  
```

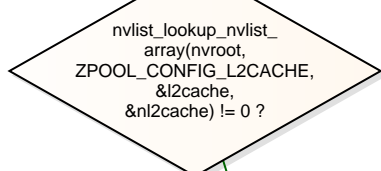
```

return (spa_vdev_exit
(spa, NULL, txg, error));
  
```



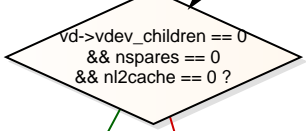
```

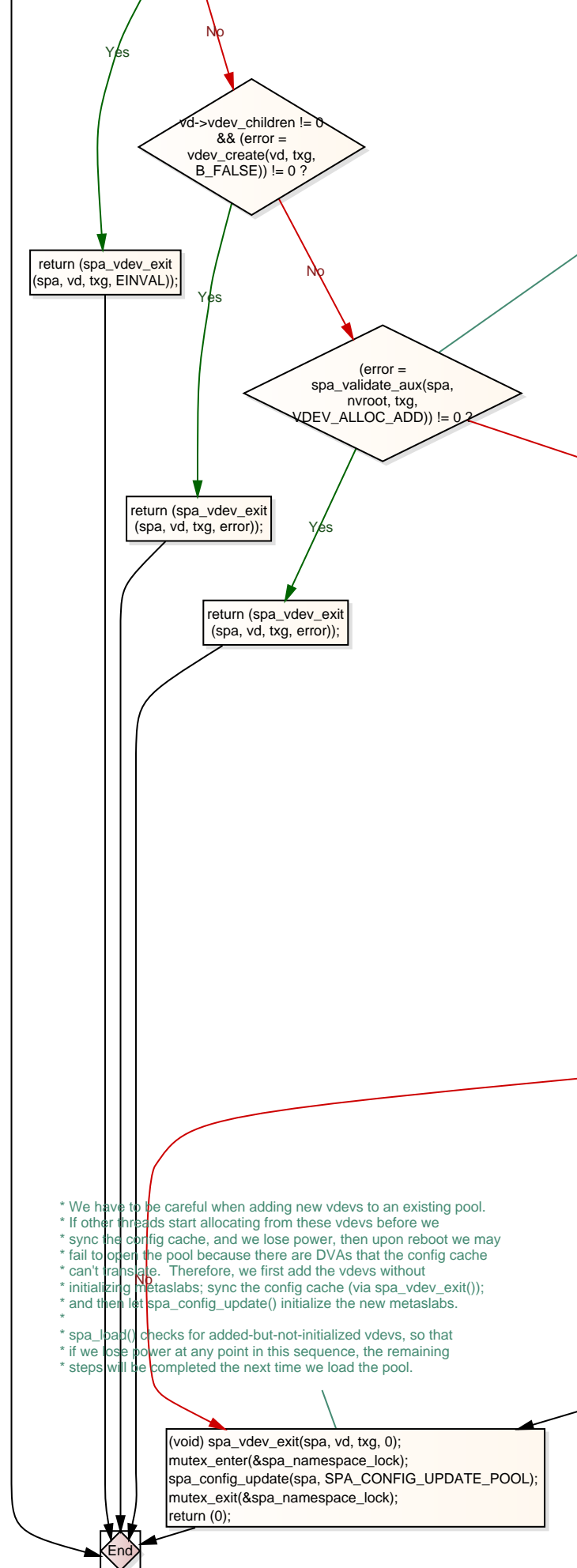
nspares = 0;
  
```



```

nl2cache = 0;
  
```

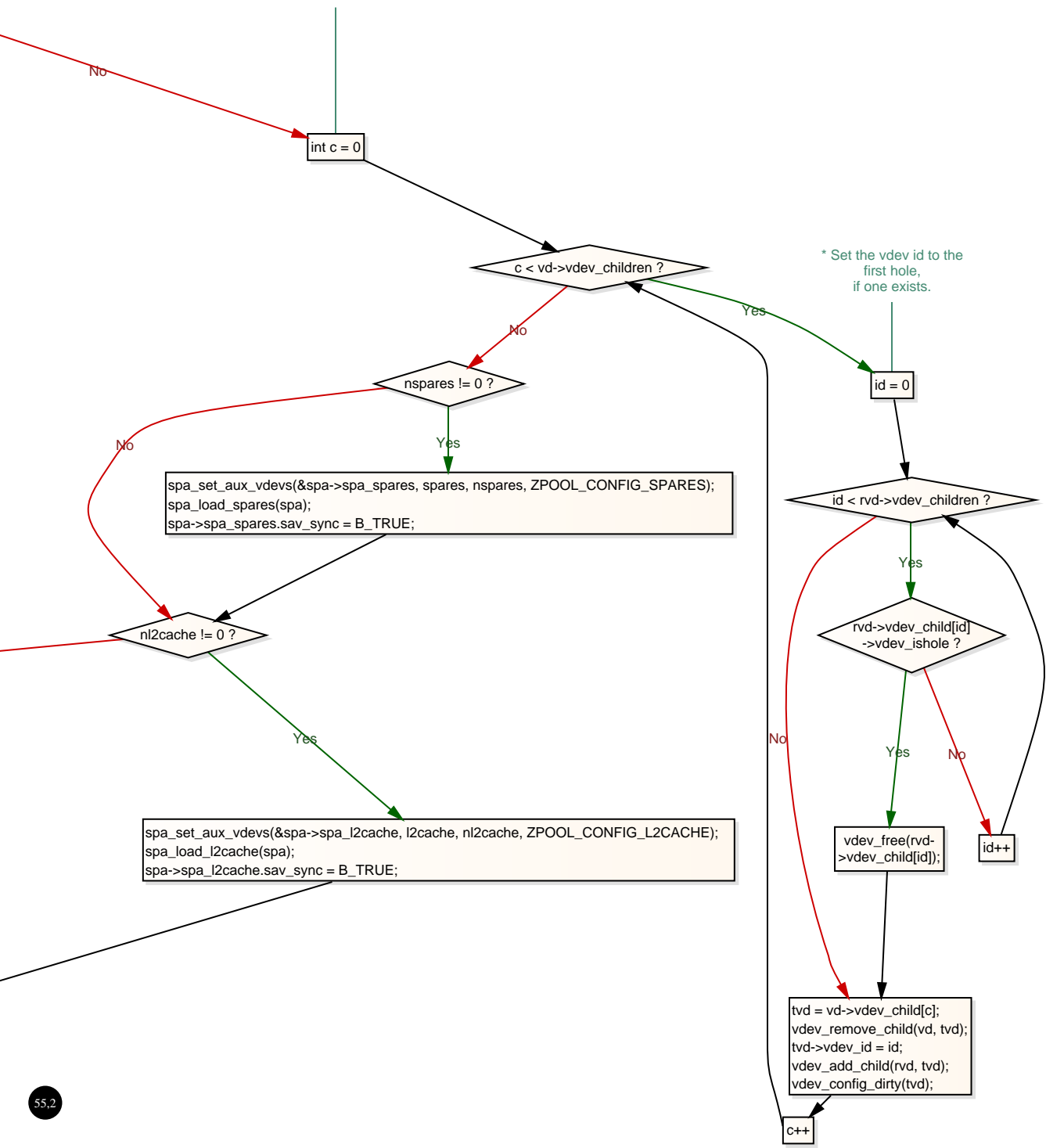


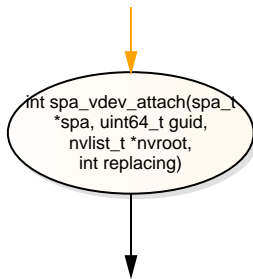


\* We must validate the spares and l2cache devices after checking the \* children. Otherwise, vdev\_inuse() will blindly overwrite the spare.

\* Transfer each new top-level vdev from vd to rvd.

\* Set the vdev id to the first hole, if one exists.



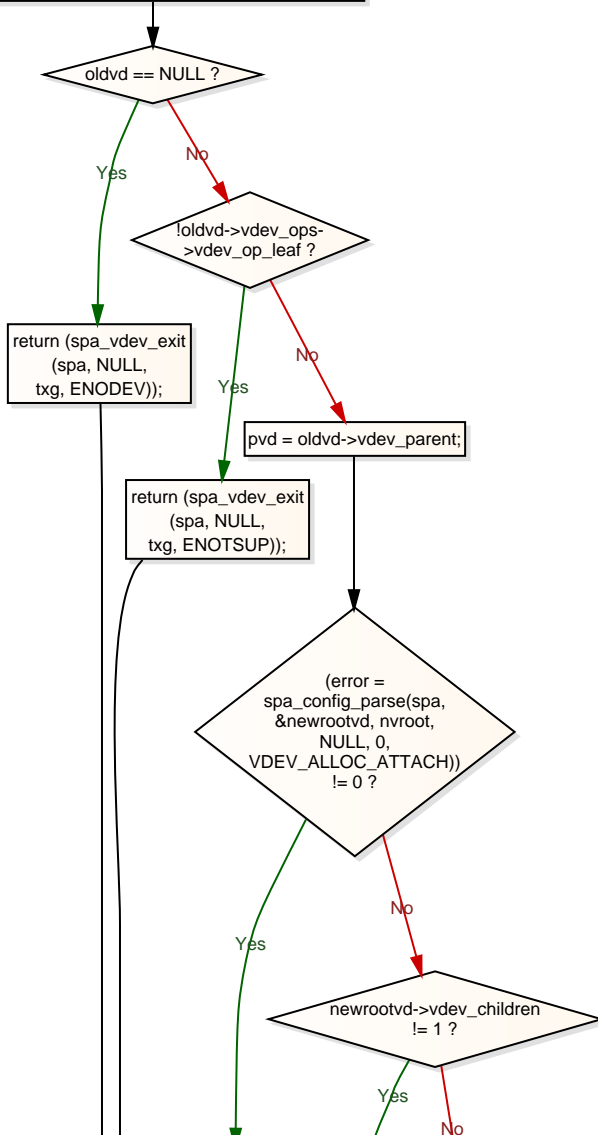


\* Attach a device to a mirror. The arguments are the path to any device in the mirror, and the nvroot for the new device. If the path specifies a device that is not mirrored, we automatically insert the mirror vdev.

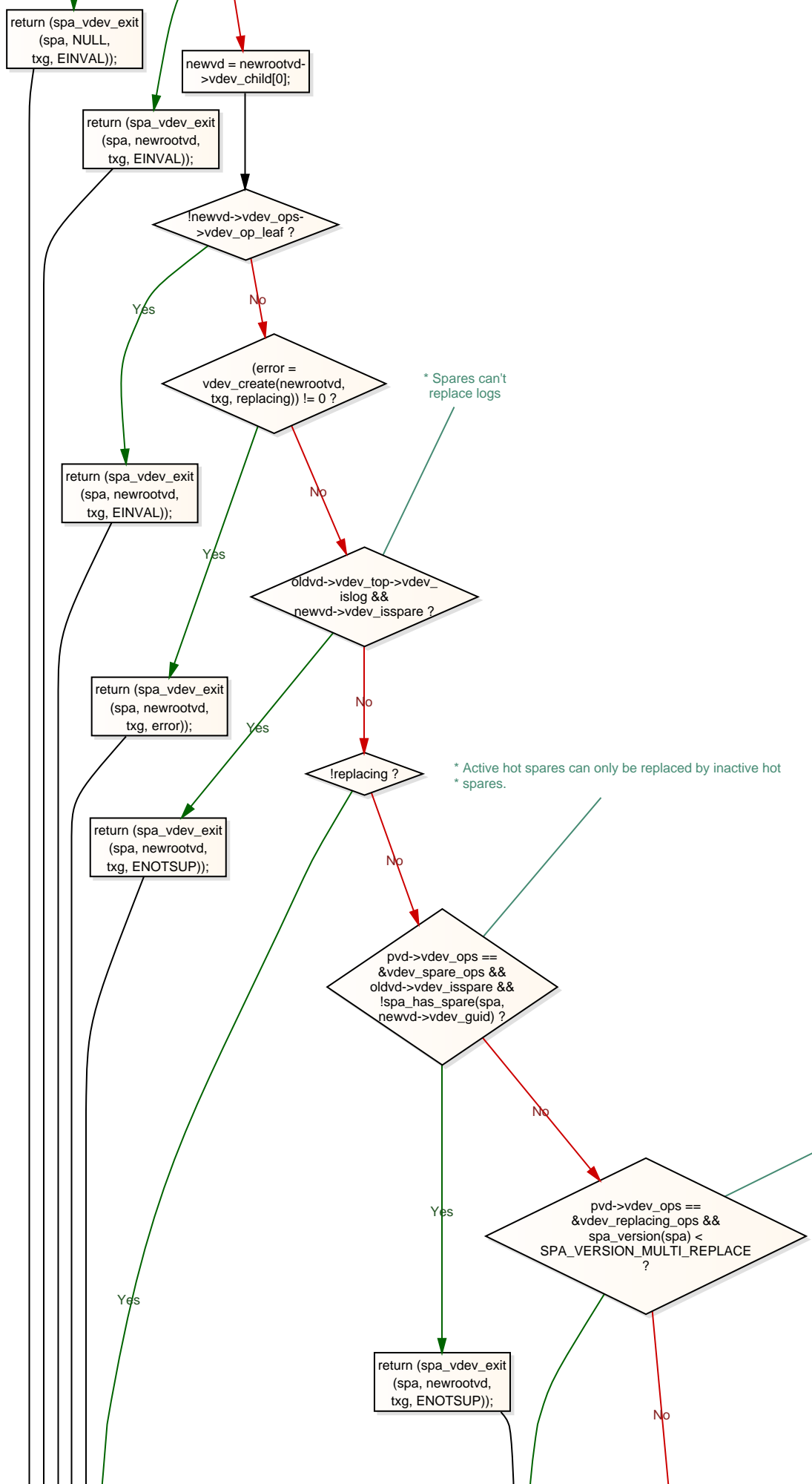
\* If 'replacing' is specified, the new device is intended to replace the existing device; in this case the two devices are made into their own mirror using the 'replacing' vdev, which is functionally identical to the mirror vdev (it actually reuses all the same ops) but has a few extra rules: you can't attach to it after it's been created, and upon completion of resilvering, the first disk (the one being replaced) is automatically detached.

```

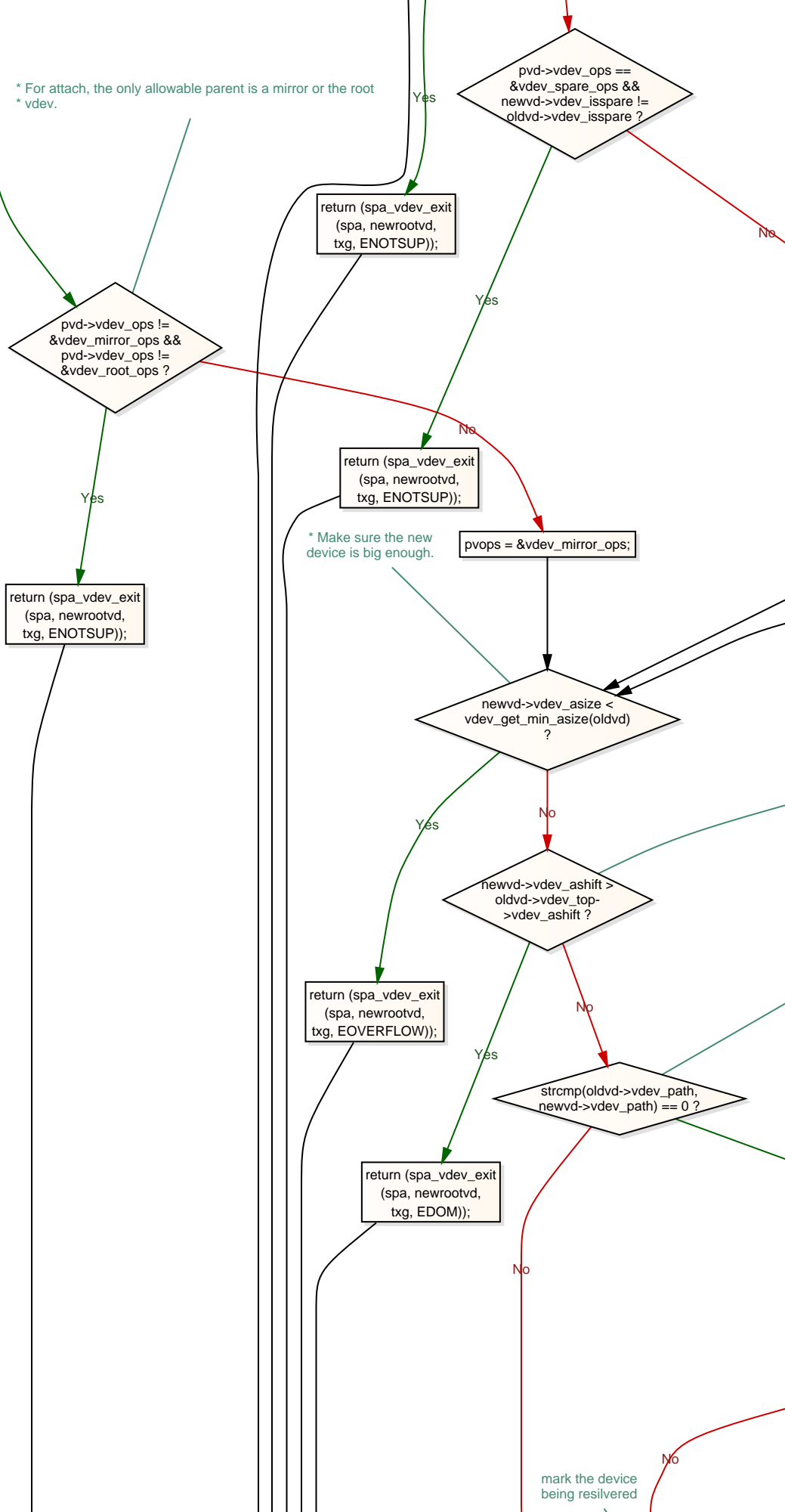
uint64_t txg, dtl_max_txg;
vdev_t *rvd = spa->spa_root_vdev;
vdev_t *oldvd, *newvd, *newrootvd, *pvd, *tvd;
vdev_ops_t *pprops;
char *oldvdpath, *newvdpath;
int newvd_isspare;
int error;
ASSERT(spa_writeable(spa));
txg = spa_vdev_enter(spa);
oldvd = spa_lookup_by_guid(spa, guid, B_FALSE);
  
```

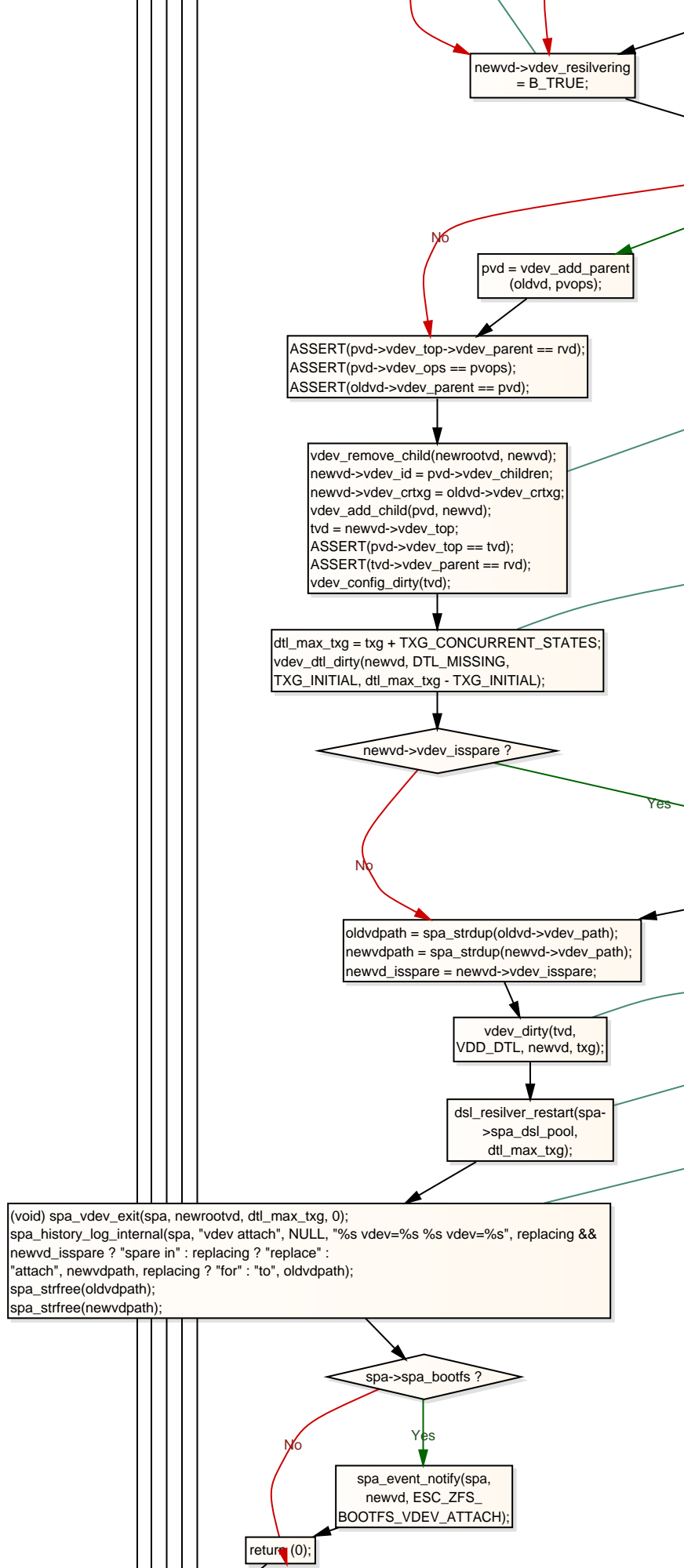







\* For attach, the only allowable parent is a mirror or the root  
 \* vdev.

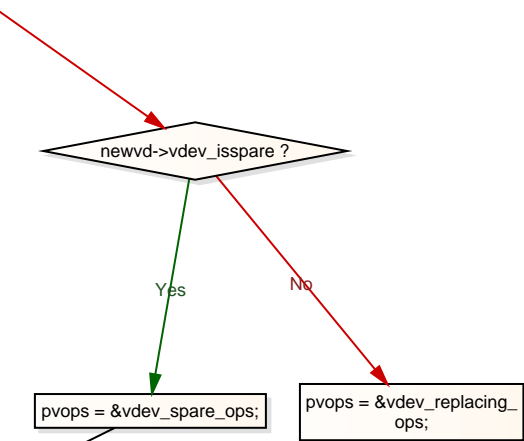






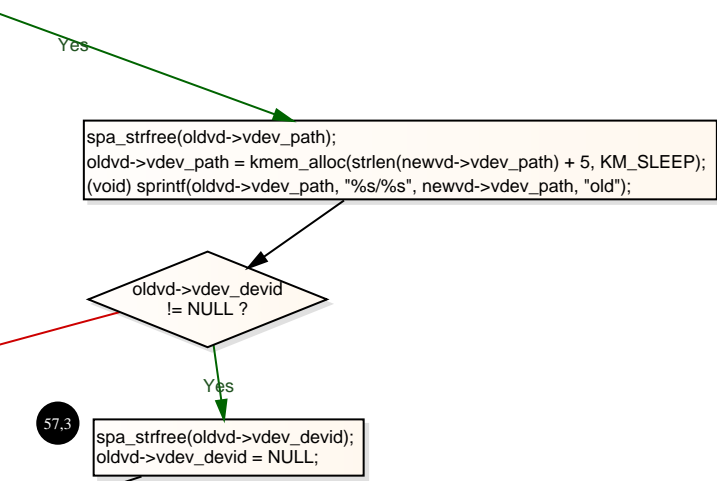
\* If the source is a hot spare, and the parent isn't already a  
\* spare, then we want to create a new hot spare. Otherwise, we  
\* want to create a replacing vdev. The user is not allowed to  
\* attach to a spared vdev child unless the 'isspare' state is  
\* the same (spare replaces spare, non-spare replaces  
\* non-spare).





\* The new device cannot have a higher alignment requirement  
\* than the top-level vdev.

\* If this is an in-place replacement,  
update oldvd's path and devid  
\* to make it distinguishable from  
newvd, and unopenable from now on.



\* If the parent is not a mirror, or  
if we're replacing, insert the new  
\* mirror/replacing/spare vdev above oldvd.

pvd->vdev\_ops != pvops ?

Yes

\* Extract the new device  
from its root and  
add it to pvd.

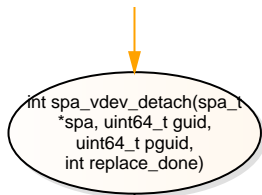
\* Set newvd's DTL to [TXG\_INITIAL, dtl\_max\_txg) so that we account  
\* for any dmu\_sync-ed blocks. It will propagate upward when  
\* spa\_vdev\_exit() calls vdev\_dtl\_reassess().

```
spa_spare_activate(newvd);  
spa_event_notify(spa, newvd, ESC_ZFS_VDEV_SPARE);
```

\* Mark newvd's DTL  
dirty in this txg.

\* Restart the resilver

\* Commit the config



\* Detach a device from a mirror or replacing vdev.  
 \* If 'replace\_done' is specified, only detach if the parent  
 \* is a replacing vdev.

```

uint64_t txg;
int error;
vdev_t *rvd = spa->spa_root_vdev;
vdev_t *vd, *pvd, *cvd, *tvd;
boolean_t unspare = B_FALSE;
uint64_t unspare_guid;
char *vdpath;
ASSERT(spa_writeable(spa));
txg = spa_vdev_enter(spa);
vd = spa_lookup_by_guid(spa, guid, B_FALSE);
    
```

vd == NULL ?

Yes

No

!vd->vdev\_ops->  
vdev\_op\_leaf ?

Yes

No

```

return (spa_vdev_exit
(spa, NULL,
txg, ENODEV));
    
```

pvd = vd->vdev\_parent;

```

return (spa_vdev_exit
(spa, NULL,
txg, ENOTSUP));
    
```

pvd->vdev\_guid != pguid  
&& pguid != 0 ?

Yes

No

\* Only 'replacing' or  
'spare' vdevs  
can be replaced.

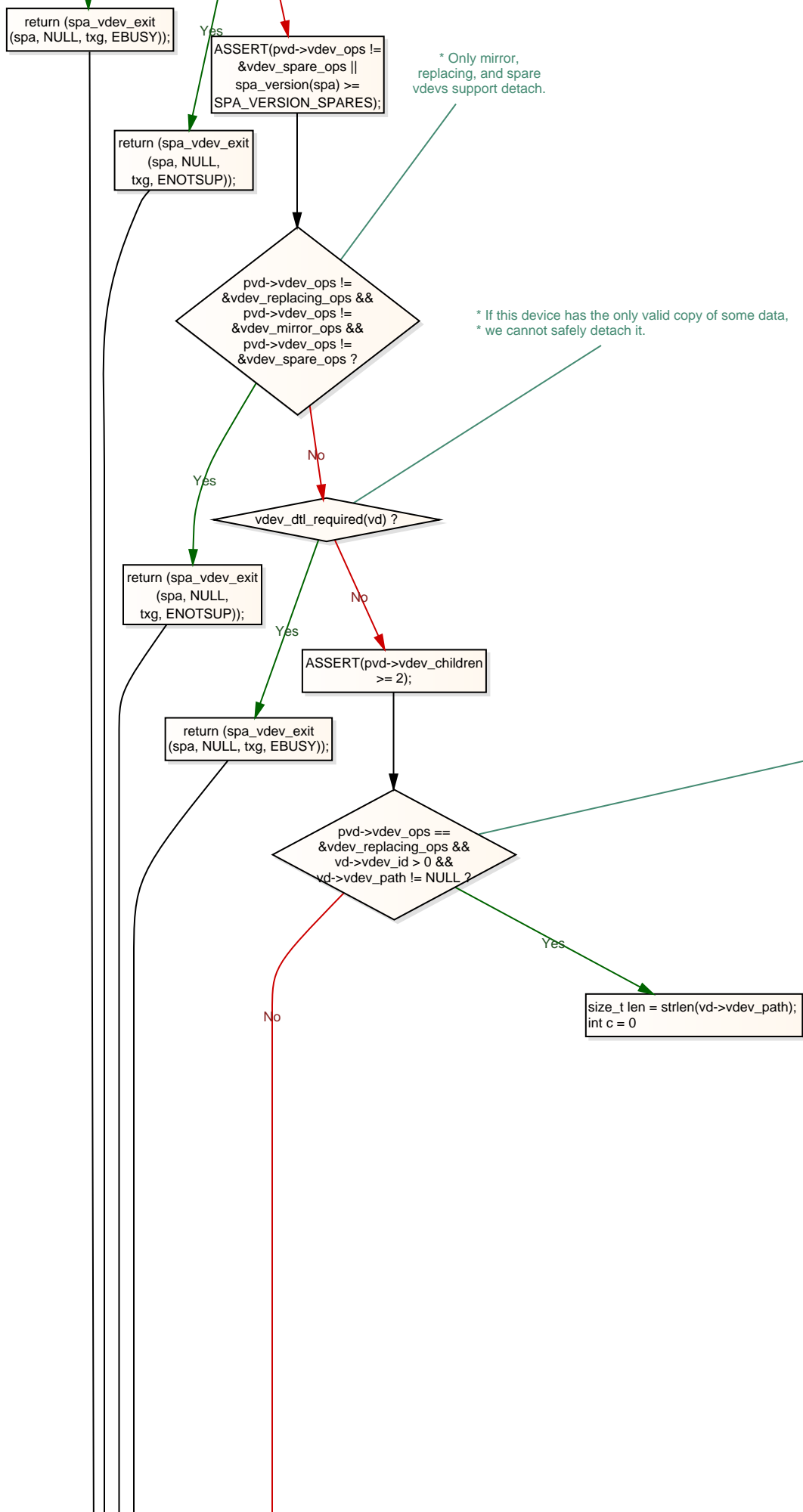
```

replace_done &&
pvd->vdev_ops !=
&vdev_replacing_ops &&
pvd->vdev_ops !=
&vdev_spare_ops ?
    
```

No

\* If the parent/child relationship is not as expected, don't do it.  
 \* Consider M(A,R(B,C)) -- that is, a mirror of A with a replacing  
 \* vdev that's replacing B with C. The user's intent in replacing  
 \* is to go from M(A,B) to M(A,C). If the user decides to cancel  
 \* the replace by detaching C, the expected behavior is to end up  
 \* M(A,B). But suppose that right after deciding to detach C,  
 \* the replacement of B completes. We would have M(A,C), and then  
 \* ask to detach C, which would leave us with just A -- not what  
 \* the user wanted. To prevent this, we make sure that the  
 \* parent/child relationship hasn't changed -- in this example,  
 \* that C's parent is still the replacing vdev R.





- \* If we are detaching the original disk from a spare, then it implies
- \* that the spare should become a real disk, and be removed from the
- \* active spare list for the pool.

```
vd->vdev_detached = B_TRUE;  
vdev_dirty(tvd, VDD_DTL, vd, txg);  
spa_event_notify(spa, vd, ESC_ZFS_VDEV_REMOVE);
```

hang on to the spa  
before we  
release the lock

```
spa_open_ref(spa, FTAG);  
error = spa_vdev_exit(spa, vd, txg, 0);  
spa_history_log_internal(spa, "detach", NULL, "vdev=%s", vdp_path);  
spa_strfree(vdp_path);
```

unspare ?

Yes

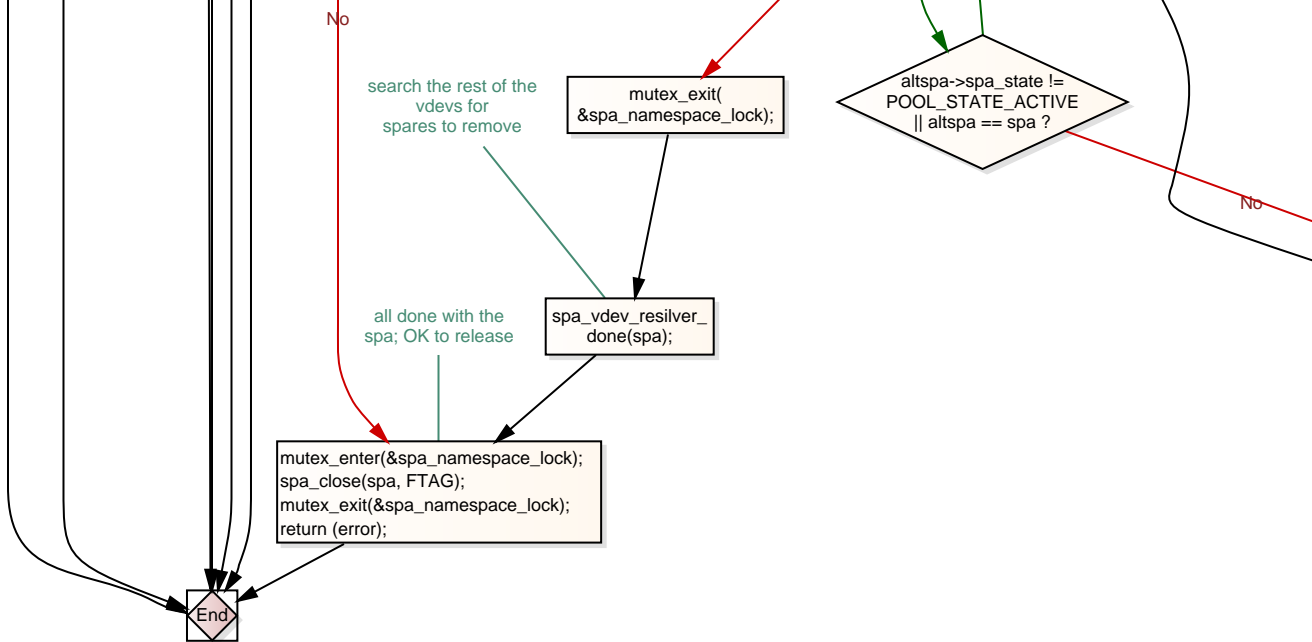
```
spa_t *altspa = NULL;  
mutex_enter(&spa_namespace_lock);
```

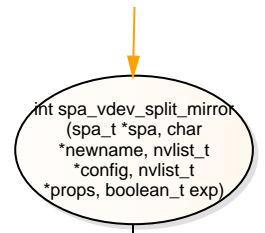
```
(altspa =  
spa_next(altspa))  
!= NULL ?
```

No

Yes

Yes





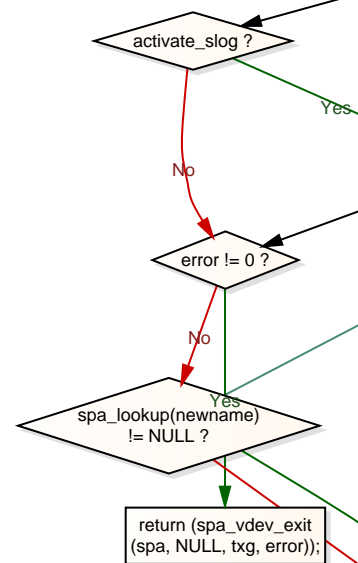
\* Split a set of devices from their mirrors, and create a new pool from them.

```

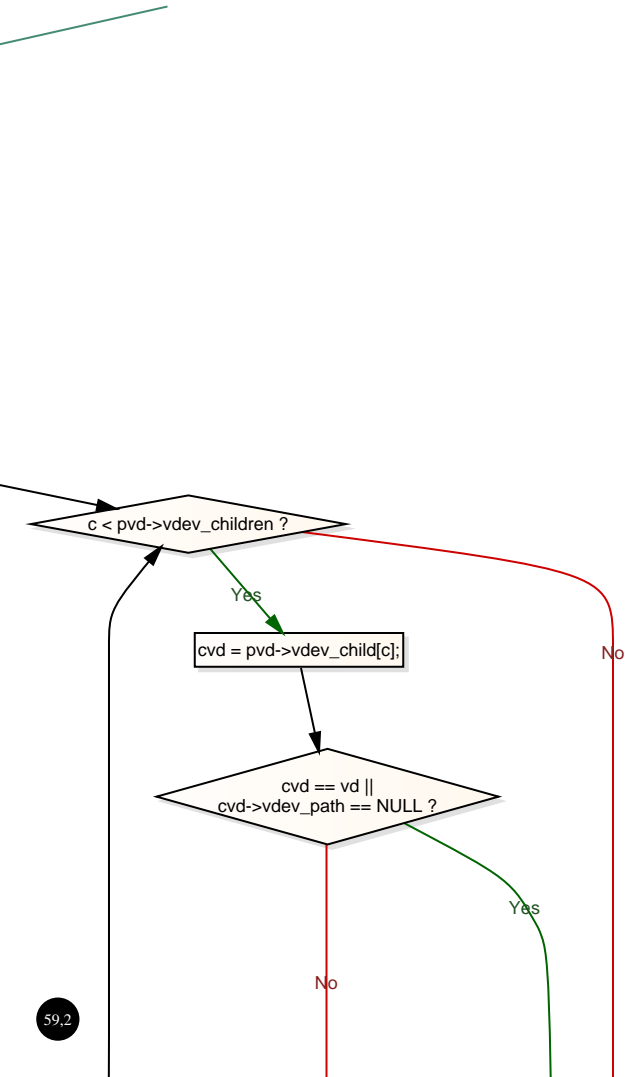
int error = 0;
uint64_t txg, *glist;
spa_t *newspa;
uint_t c, children, lastlog;
nvlist_t **child, *nvl, *tmp;
dmu_tx_t *tx;
char *altroot = NULL;
  
```

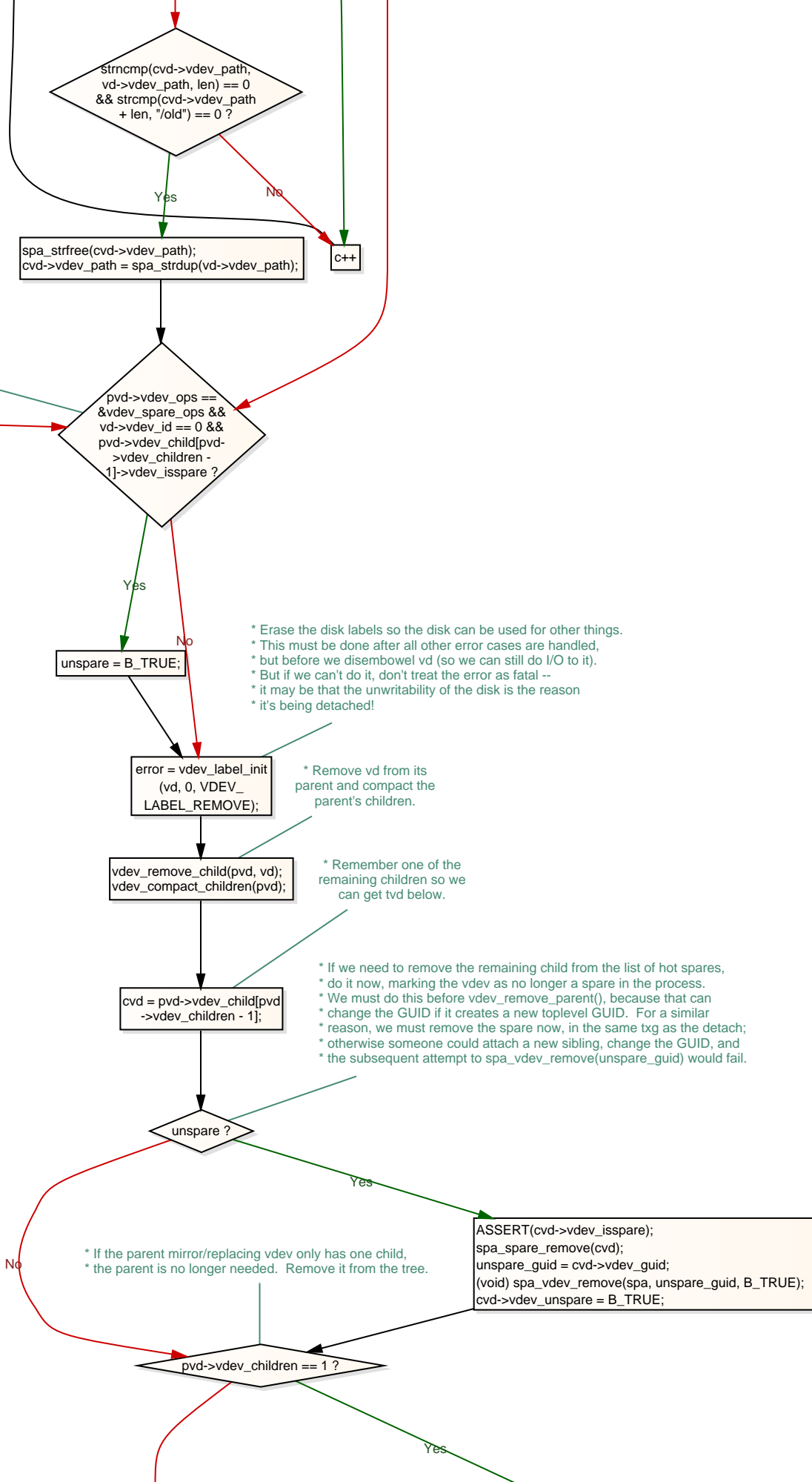
```

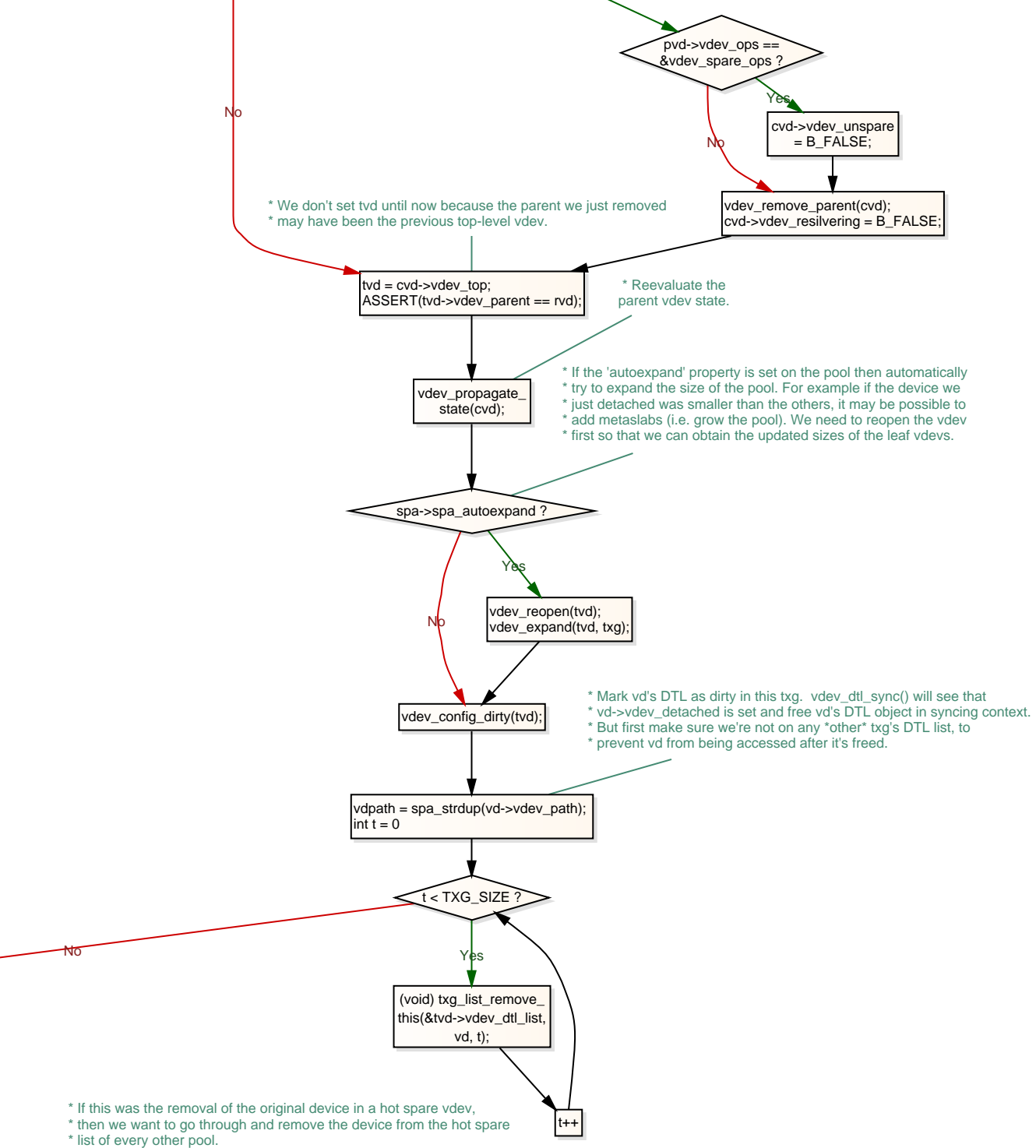
vdev_t *rvd, **vml = NULL;
boolean_t activate_slog;
ASSERT(spa_writeable(spa));
txg = spa_vdev_enter(spa);
  
```




\* If we are detaching the second disk from a replacing vdev, then  
\* check to see if we changed the original vdev's path to have "/old"  
\* at the end in spa\_vdev\_attach(). If so, undo that change now.





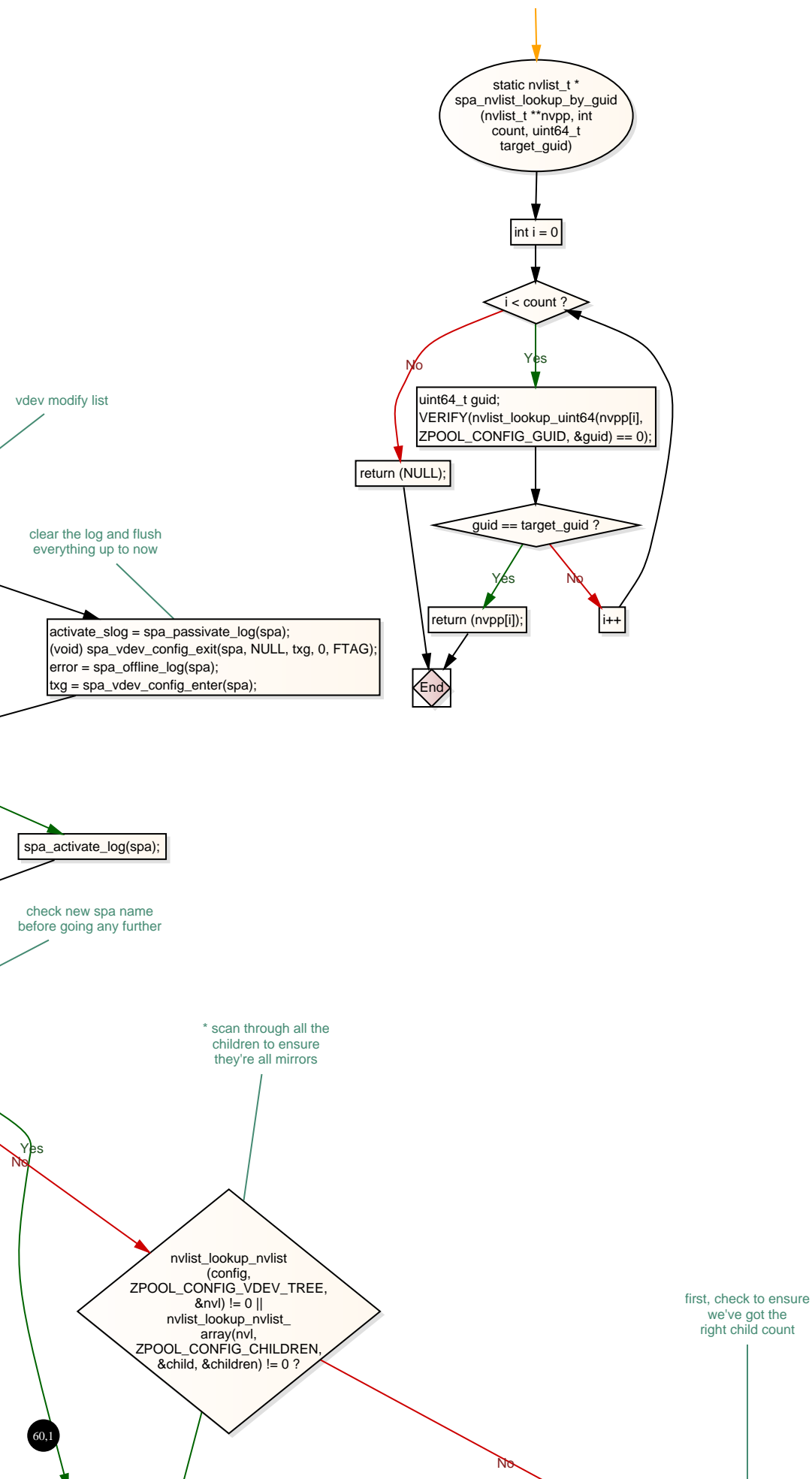


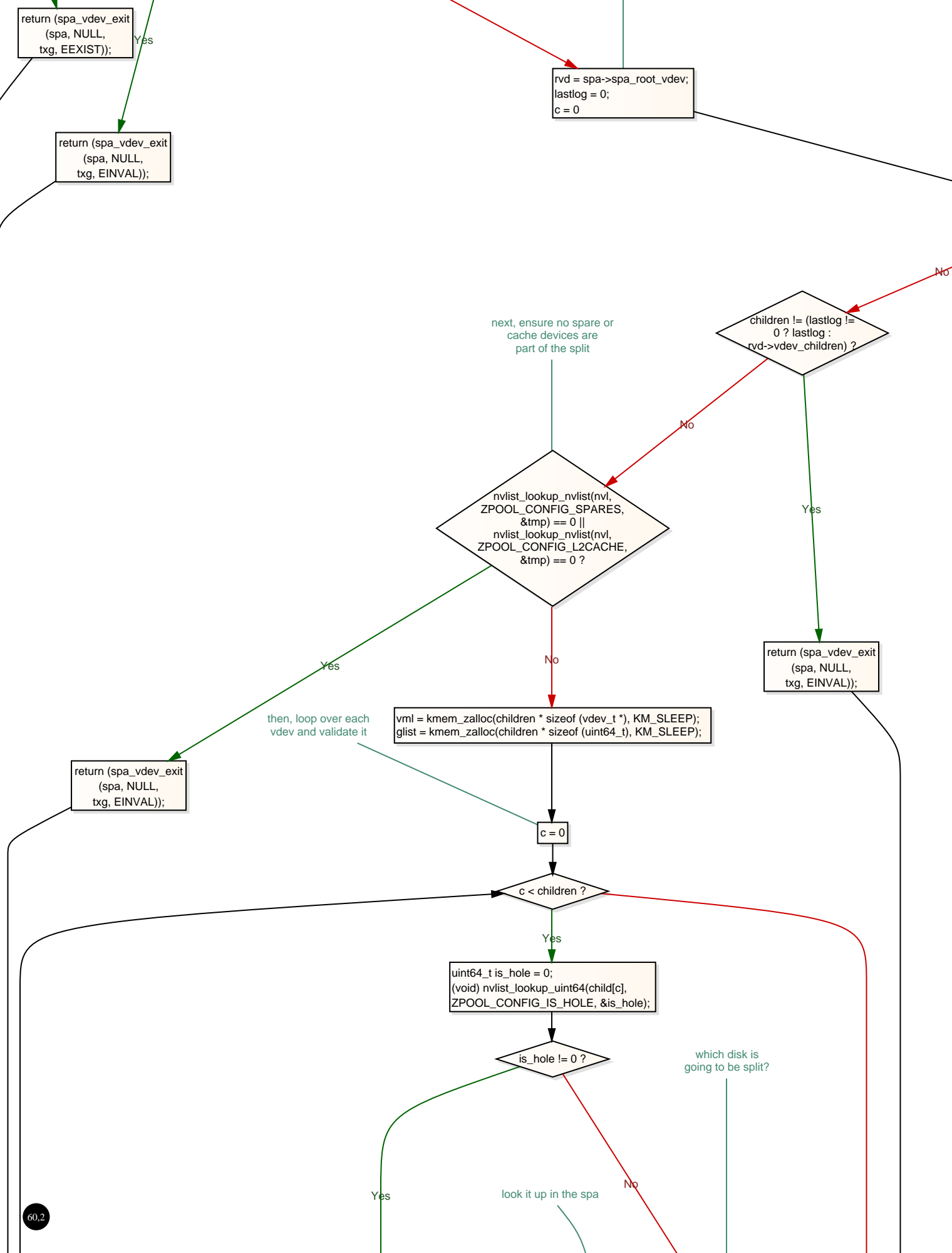


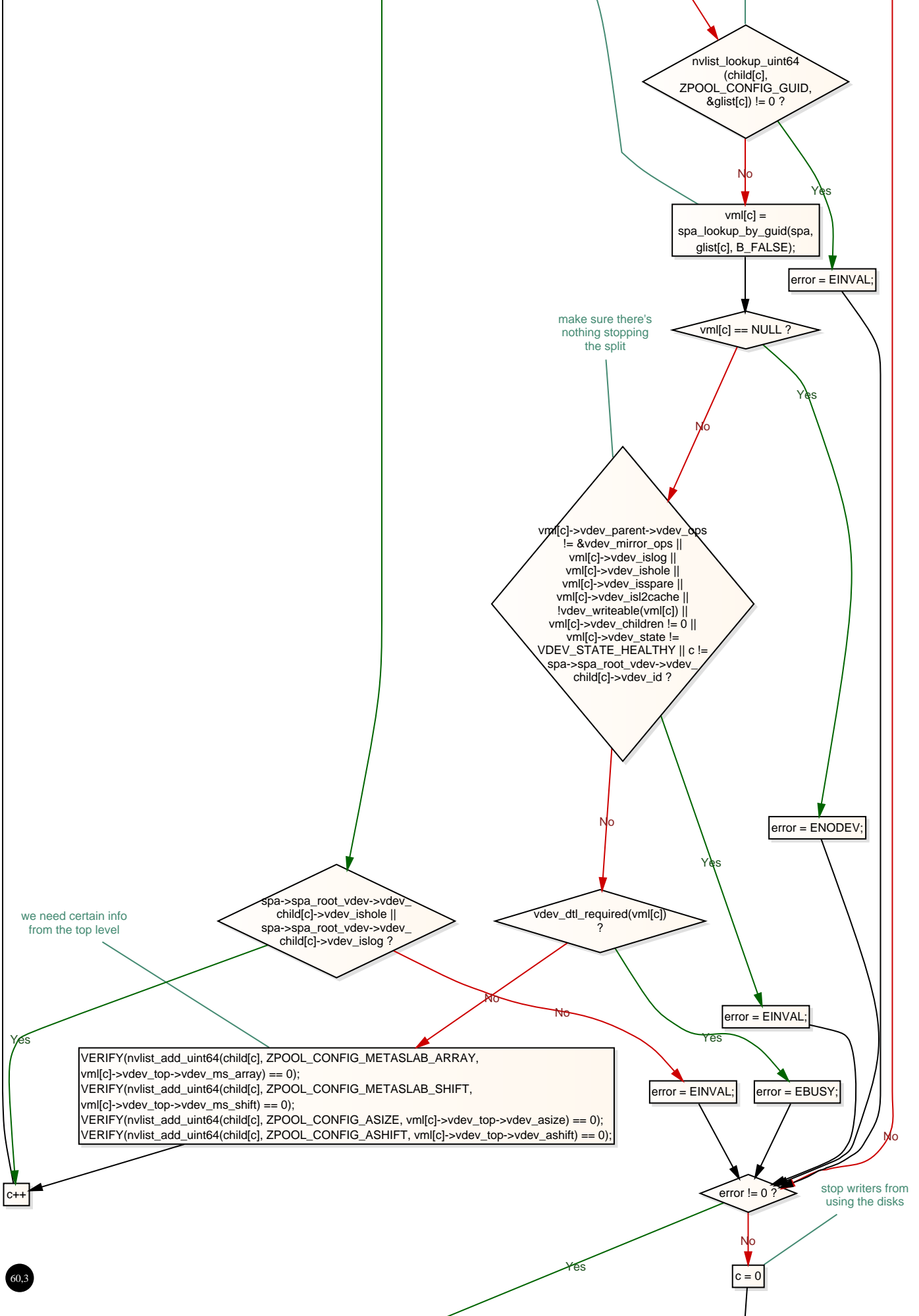


```
spa_open_ref(altspa, FTAG);  
mutex_exit(&spa_namespace_lock);  
(void) spa_vdev_remove(altspa, unspare_guid, B_TRUE);  
mutex_enter(&spa_namespace_lock);  
spa_close(altspa, FTAG);
```

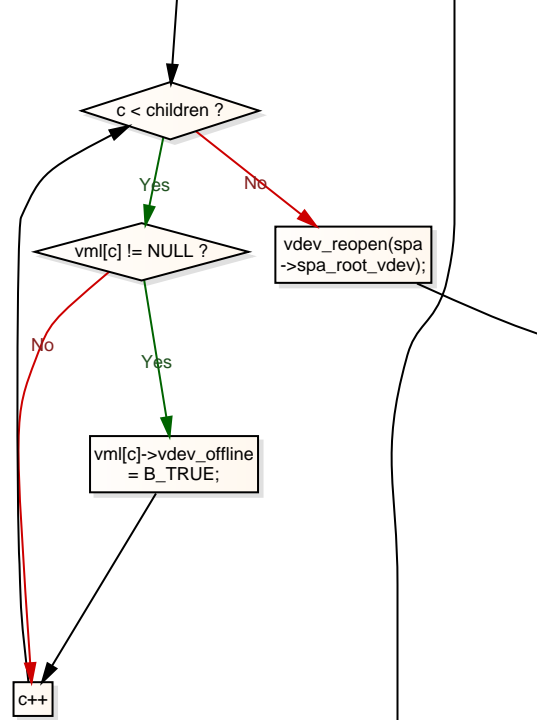




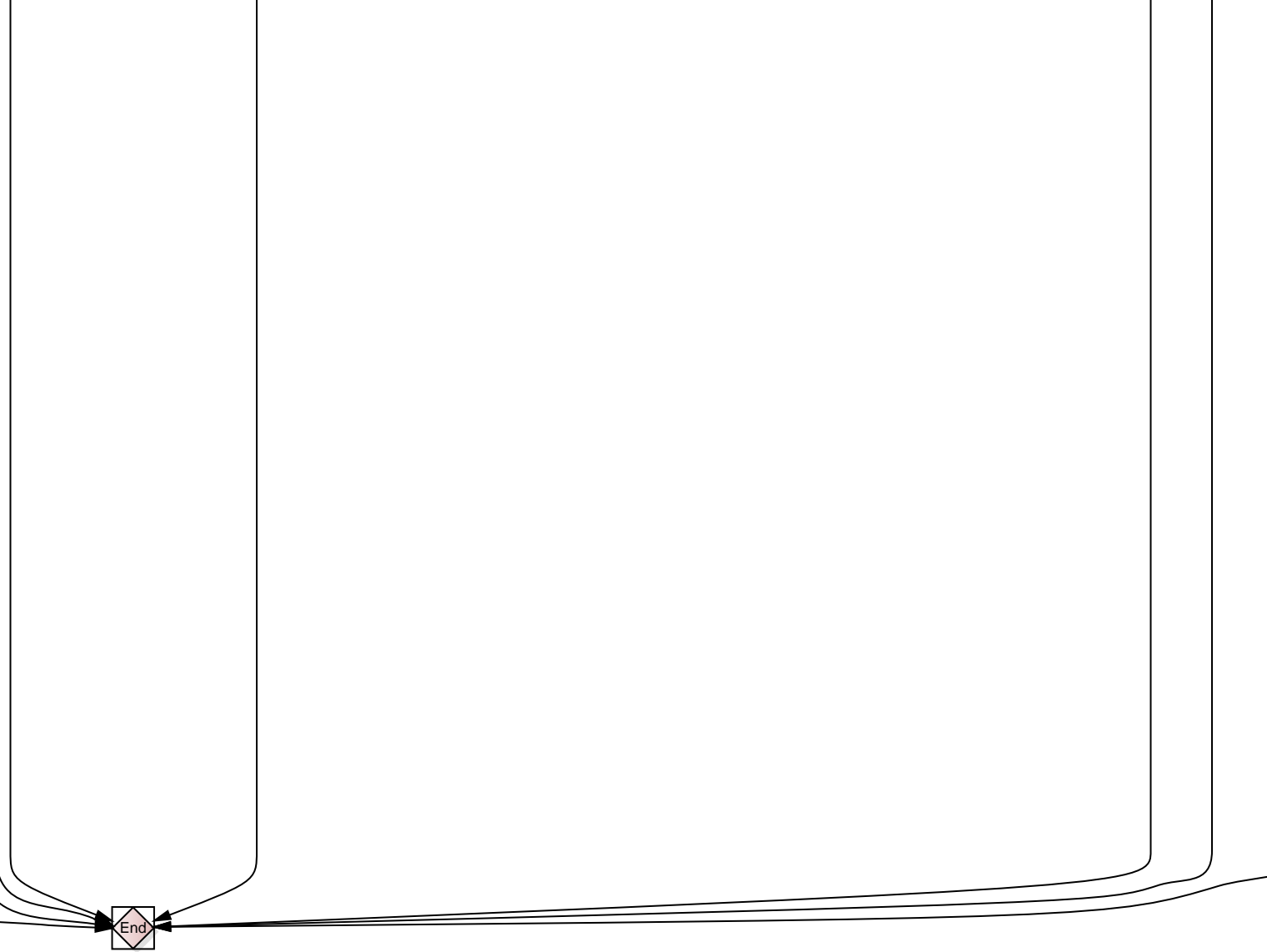




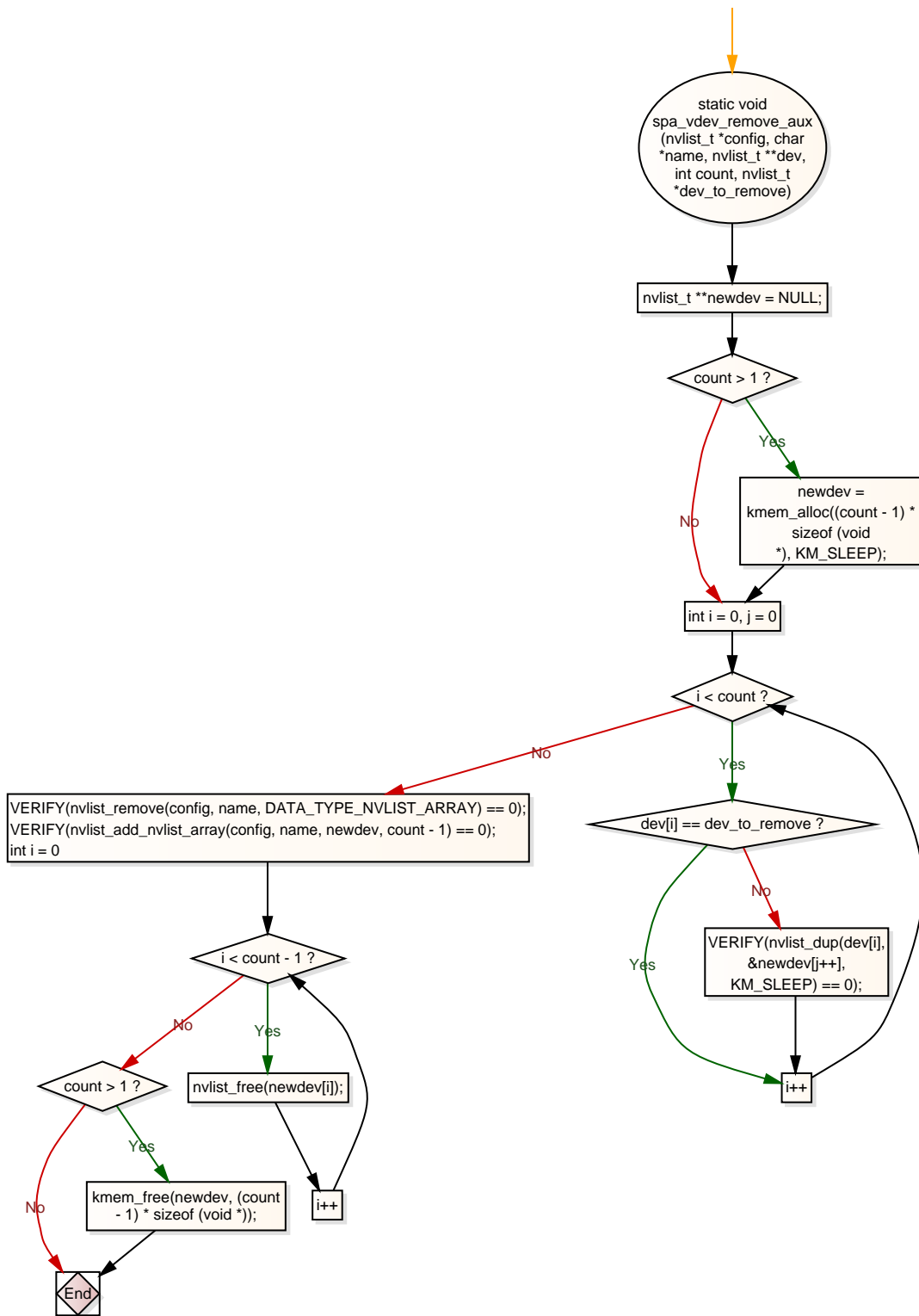
```
kmem_free(vml, children * sizeof (vdev_t *));  
kmem_free(glist, children * sizeof (uint64_t));  
return (spa_vdev_exit(spa, NULL, txg, error));
```

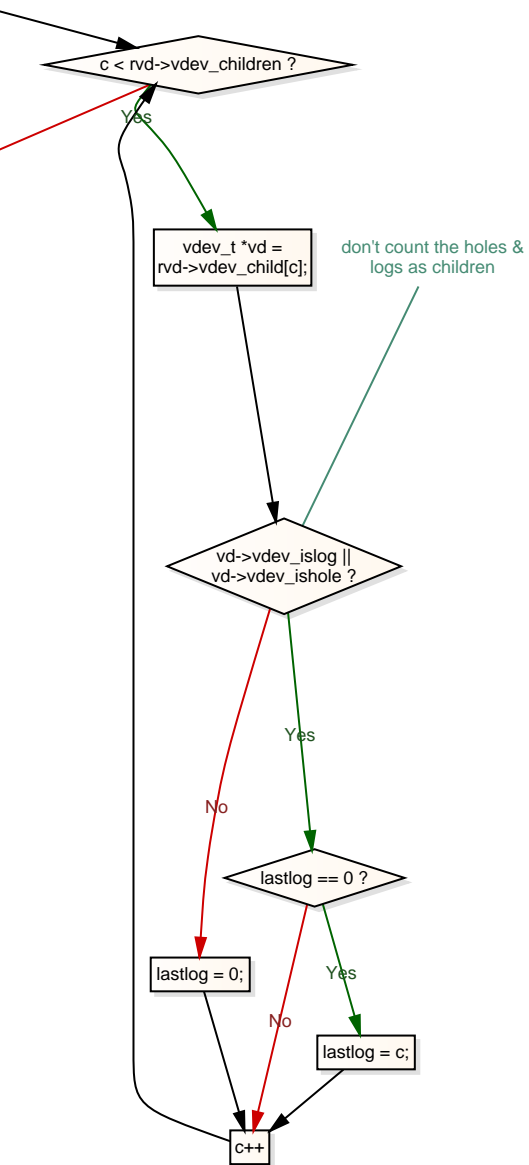












\* Temporarily record the splitting vdevs in the spa config. This  
\* will disappear once the config is regenerated.

```
VERIFY(nvlist_alloc(&nvl, NV_UNIQUE_NAME, KM_SLEEP) == 0);
VERIFY(nvlist_add_uint64_array(nvl, ZPOOL_CONFIG_SPLIT_LIST, glist, children) == 0);
kmem_free(glist, children * sizeof(uint64_t));
mutex_enter(&spa->spa_props_lock);
VERIFY(nvlist_add_nvlist(spa->spa_config, ZPOOL_CONFIG_SPLIT, nvl) == 0);
mutex_exit(&spa->spa_props_lock);
spa->spa_config_splitting = nvl;
vdev_config_dirty(spa->spa_root_vdev);
```

configure and  
create the new pool

```
VERIFY(nvlist_add_string(config, ZPOOL_CONFIG_POOL_NAME, newname) == 0);
VERIFY(nvlist_add_uint64(config, ZPOOL_CONFIG_POOL_STATE,
exp ? POOL_STATE_EXPORTED : POOL_STATE_ACTIVE) == 0);
VERIFY(nvlist_add_uint64(config, ZPOOL_CONFIG_VERSION, spa_version(spa)) == 0);
VERIFY(nvlist_add_uint64(config, ZPOOL_CONFIG_POOL_TXG, spa->spa_config_txg) == 0);
VERIFY(nvlist_add_uint64(config, ZPOOL_CONFIG_POOL_GUID, spa_generate_guid(NULL)) == 0);
(void) nvlist_lookup_string(props, zpool_prop_to_name(ZPOOL_PROP_ALTROOT), &altroot);
```

add the new pool  
to the namespace

```
newspa = spa_add(newname, config, altroot);
newspa->spa_config_txg = spa->spa_config_txg;
spa_set_log_state(newspa, SPA_LOG_CLEAR);
```

release the spa config  
lock, retaining  
the namespace lock

```
spa_vdev_config_exit(spa,
NULL, txg, 0, FTAG);
```

zio\_injection\_enabled ?

Yes

```
zio_handle_panic_injection(spa, FTAG, 1);
```

No

```
spa_activate(newspa, spa_mode_global);
spa_async_suspend(newspa);
```

create the new pool from  
the disks of the  
original pool

```
error = spa_load(newspa,
SPA_LOAD_IMPORT,
SPA_IMPORT_ASSEMBLE,
B_TRUE);
```

error ?

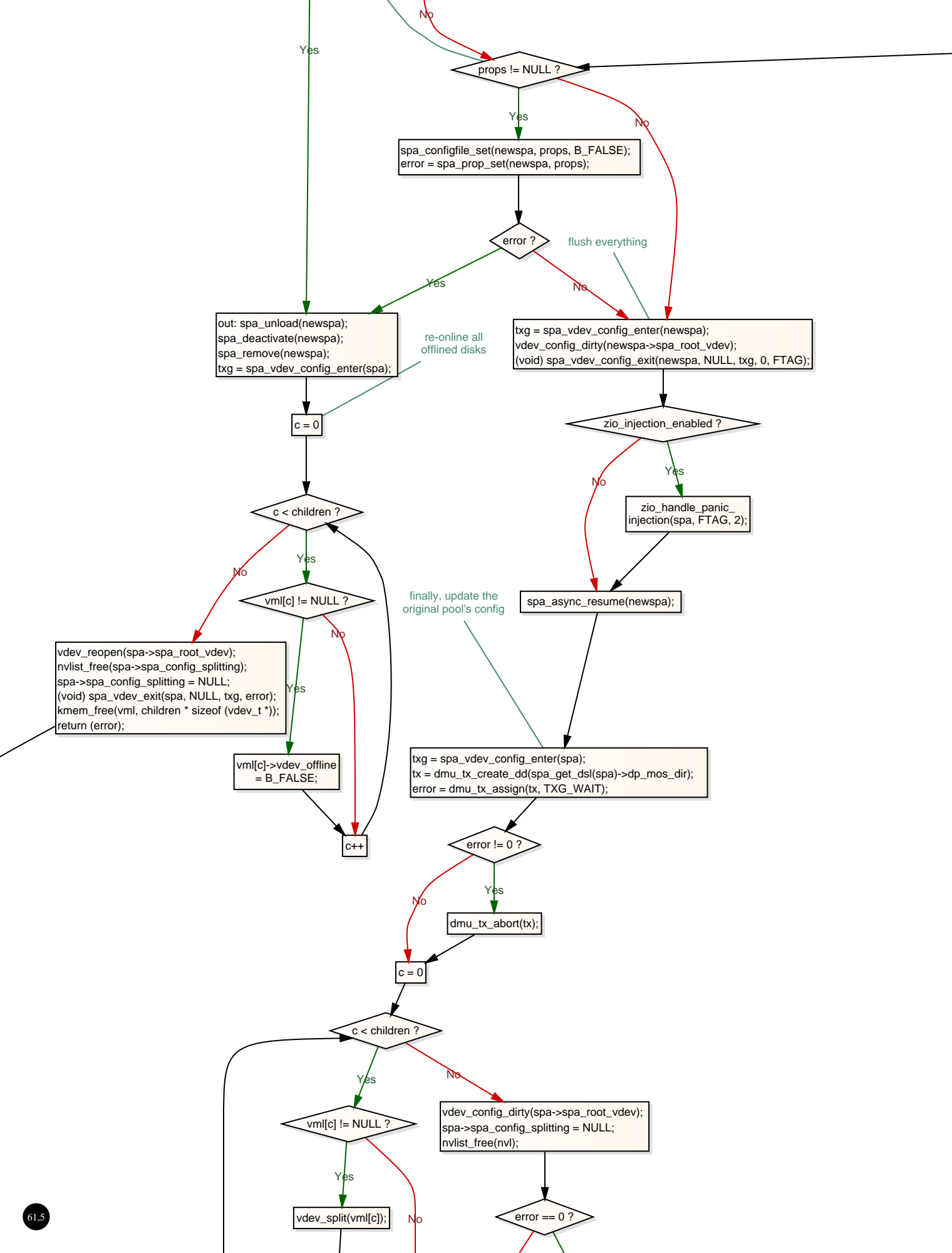
if that worked, generate  
a real config  
for the new pool

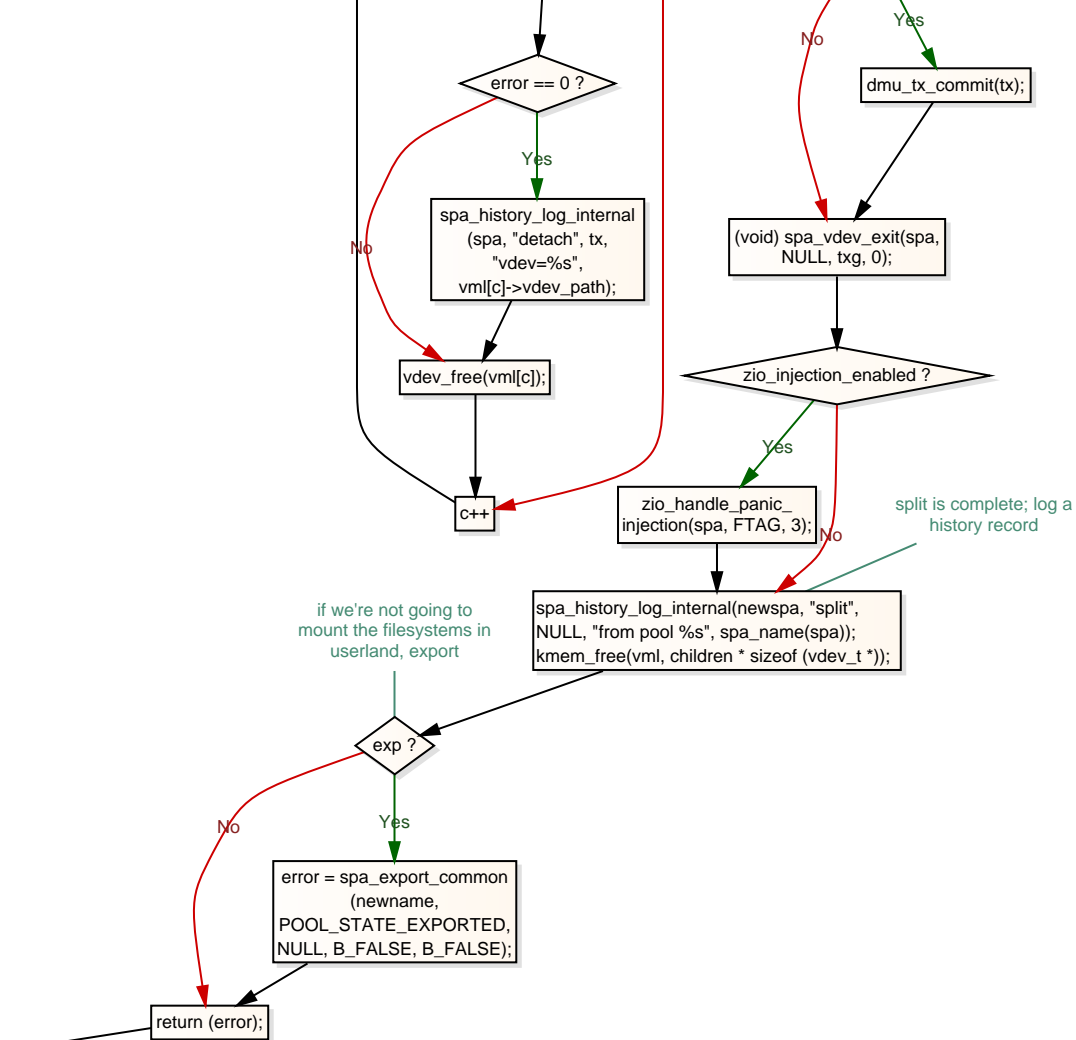
No

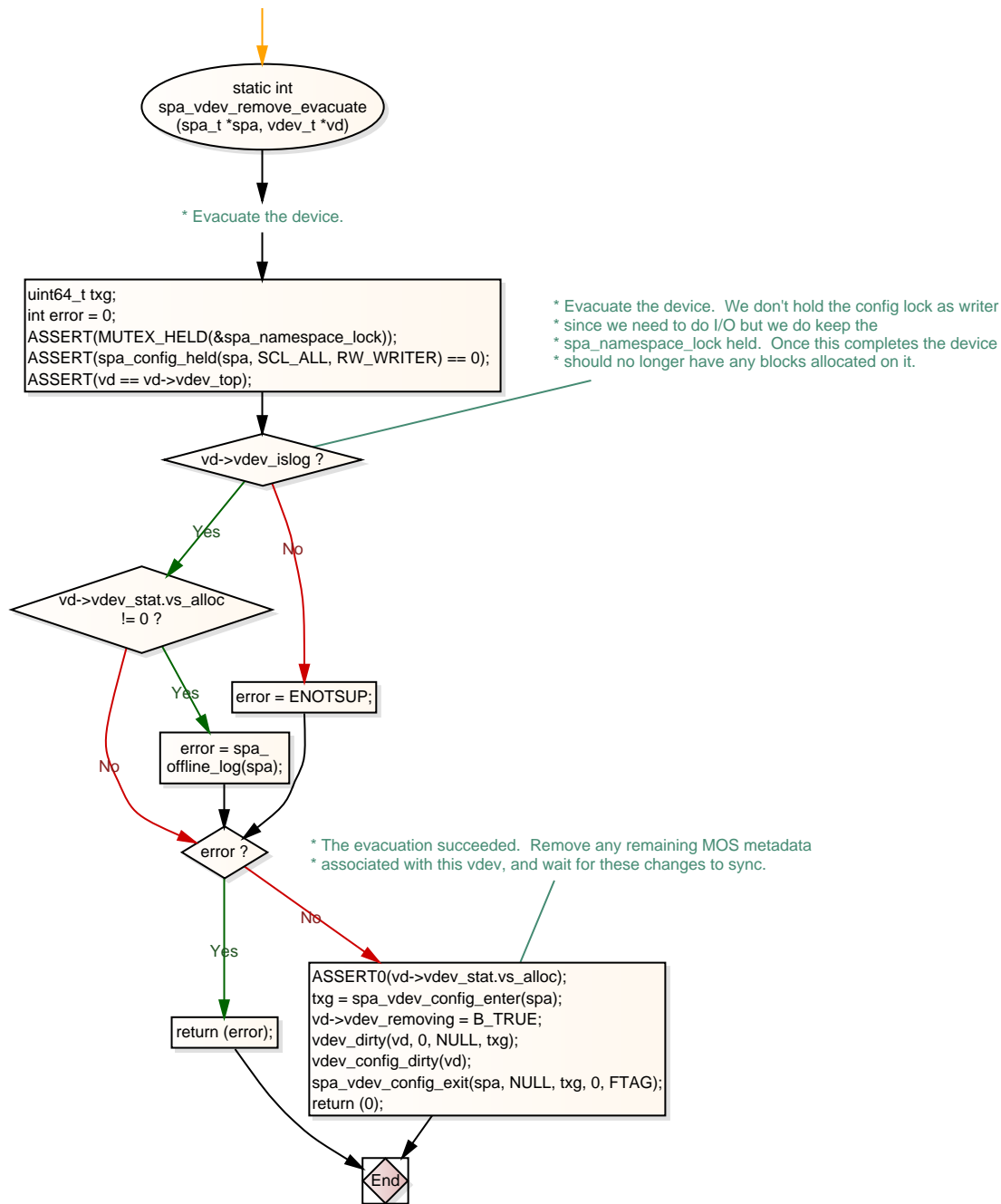
newspa->spa\_root\_vdev  
!= NULL ?

Yes


set the props





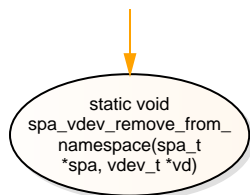






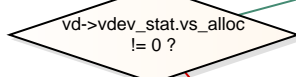
```
VERIFY(nvlist_alloc(&newspa->spa_config_splitting, NV_UNIQUE_NAME, KM_SLEEP) == 0);  
VERIFY(nvlist_add_uint64(newspa->spa_config_splitting,  
ZPOOL_CONFIG_SPLIT_GUID, spa_guid(spa)) == 0);  
spa_config_set(newspa, spa_config_generate(newspa, NULL, -1ULL, B_TRUE));
```





```

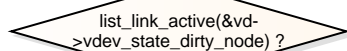
vdev_t *rvd = spa->spa_root_vdev;
uint64_t id = vd->vdev_id;
boolean_t last_vdev = (id == (rvd->vdev_children - 1));
ASSERT(MUTEX_HELD(&spa_namespace_lock));
ASSERT(spa_config_held(spa, SCL_ALL, RW_WRITER) == SCL_ALL);
ASSERT(vd == vd->vdev_top);
  
```



Yes

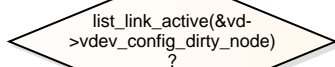
```

(void) vdev_label_init
(vd, 0, VDEV_
LABEL_REMOVE);
  
```



```

vdev_state_clean(vd);
  
```

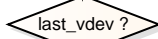


```

vdev_config_clean(vd);
  
```

```

vdev_free(vd);
  
```



```

vdev_compact_
children(rvd);
  
```

```

vd = vdev_alloc_common(spa, id, 0, &vdev_hole_ops);
vdev_add_child(rvd, vd);
  
```

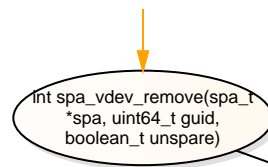
```

vdev_config_dirty(rvd);
  
```

\* Reassess the health of  
our root vdev.

```

vdev_reopen(rvd);
  
```



```

vdev_t *vd;
metaslab_group_t *mg;
nvlist_t **spares, **l2cache, *nv;
uint64_t txg = 0;
uint_t nspares, nl2cache;
int error = 0;
boolean_t locked =
MUTEX_HELD(&spa_namespace_lock);
ASSERT(spa_writeable(spa));
  
```

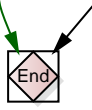


```

txg = spa_vdev_
enter(spa);
  
```

```

vd = spa_lookup_by_guid
(spa, guid, B_FALSE);
  
```



\* Cache devices can always be removed.

```
spa_vdev_remove_aux(spa->spa_l2cache.sav_config, ZPOOL_CONFIG_L2CACHE, l2cache, nl2cache, nv);  
spa_load_l2cache(spa);  
spa->spa_l2cache.sav_sync = B_TRUE;
```

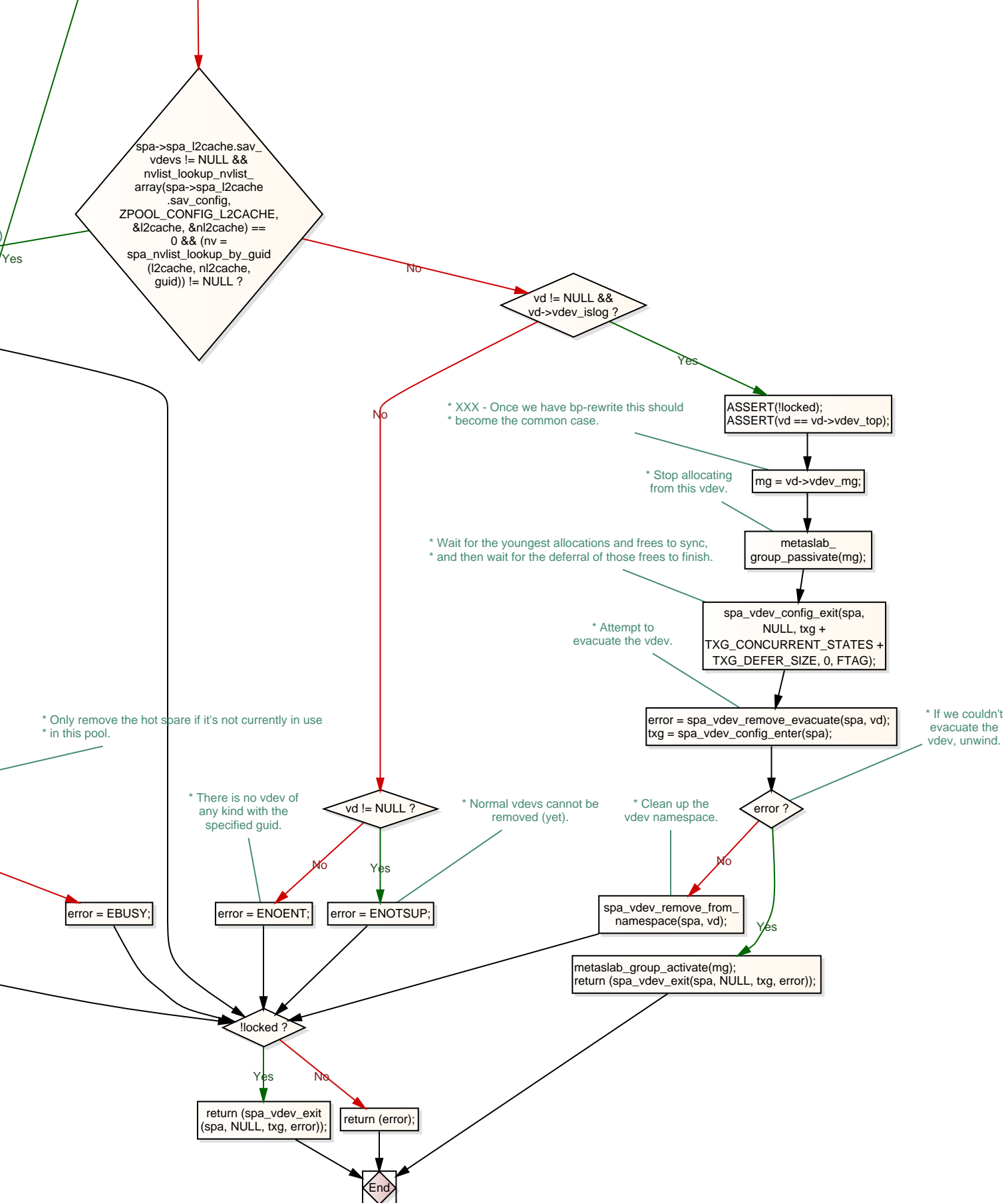


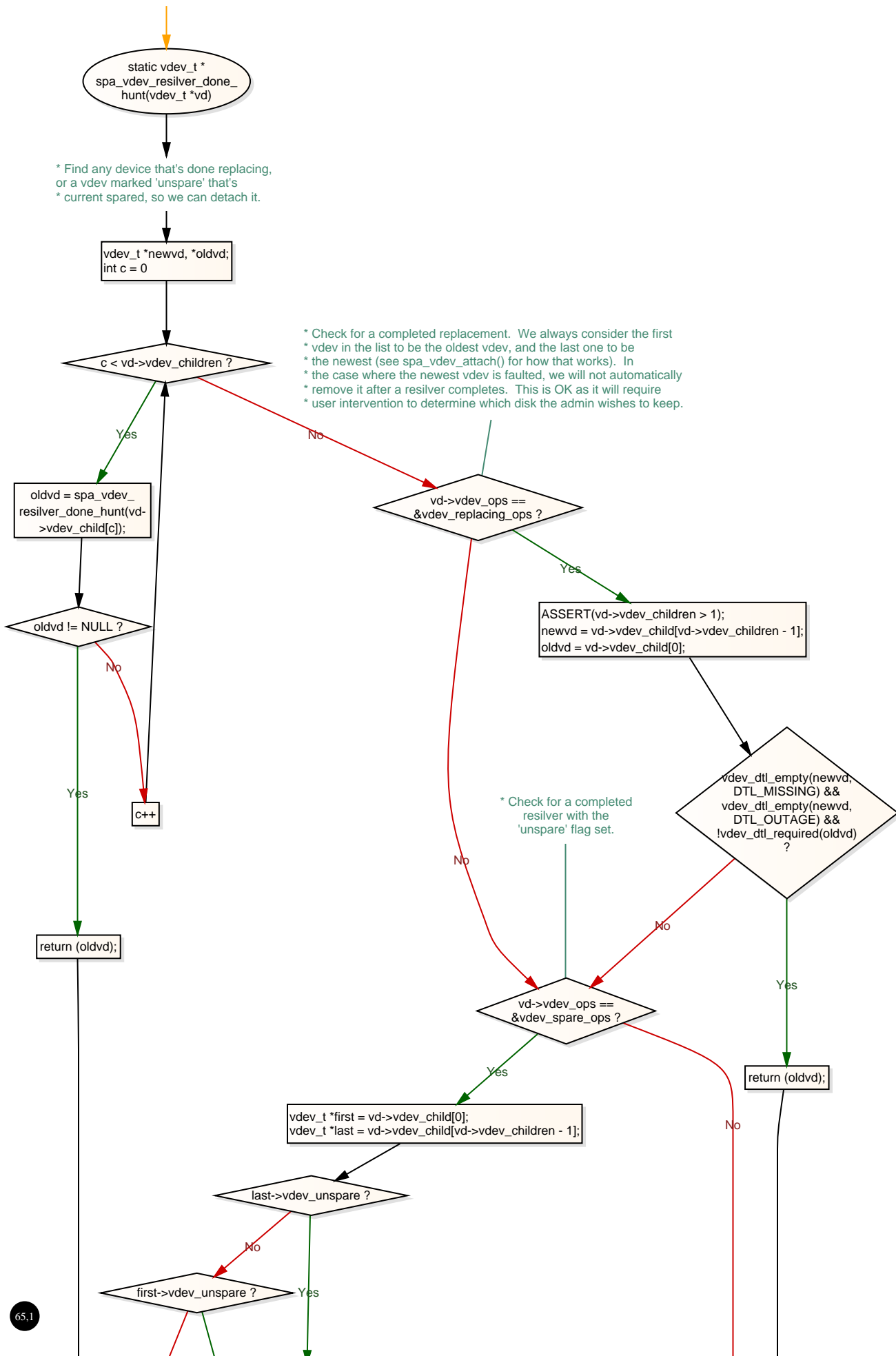
```
spa_vdev_remove_aux(spa->spa_spares.sav_config, ZPOOL_CONFIG_SPARES, spares, nspares, nv);  
spa_load_spares(spa);  
spa->spa_spares.sav_sync = B_TRUE;
```

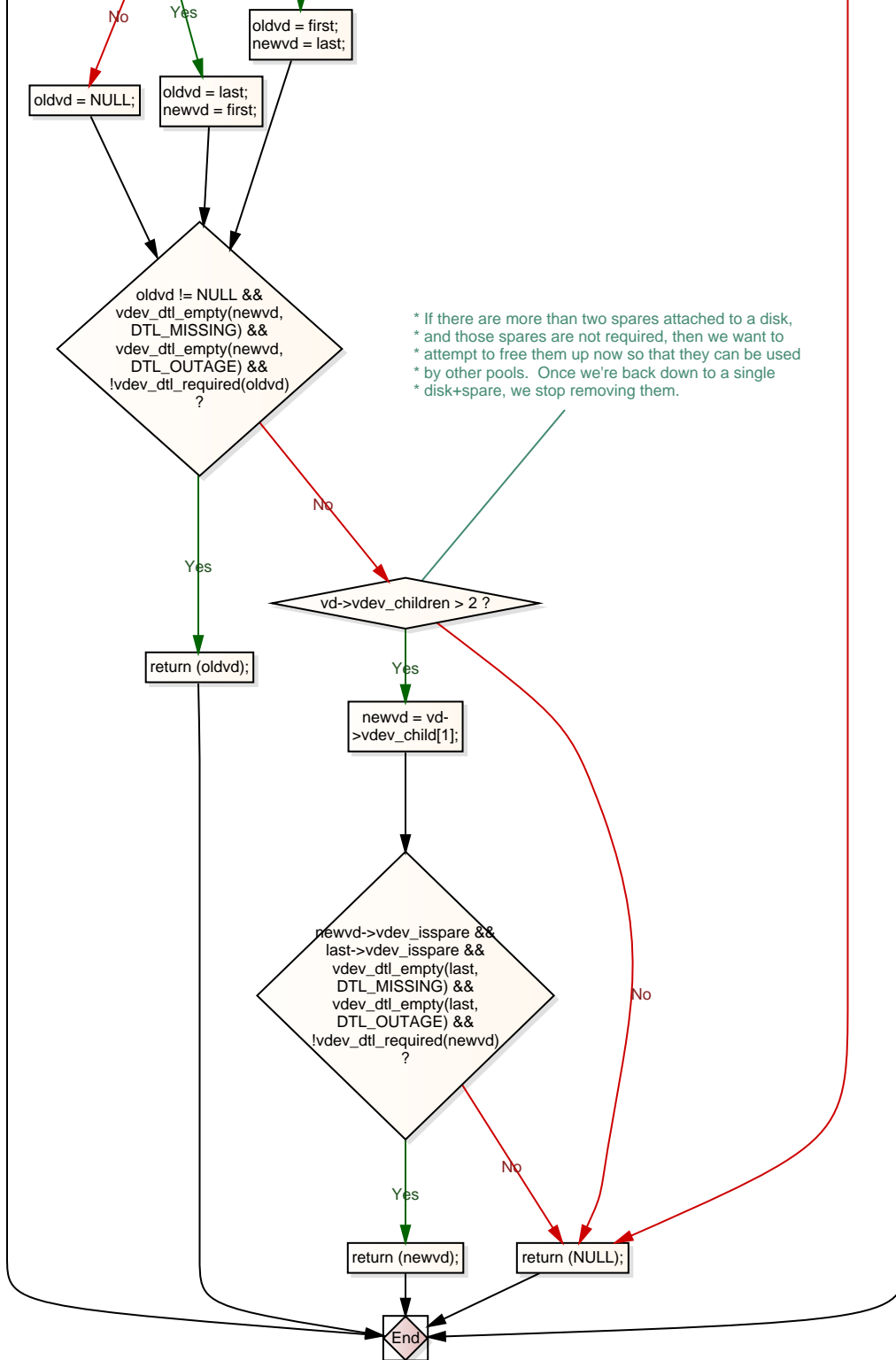
- \* Remove a device from the pool -
- \*
- \* Removing a device from the vdev namespace requires several steps and can take a significant amount of time. As a result we use
- \* the `spa_vdev_config_[enter/exit]` functions which allow us to
- \* grab and release the `spa_config_lock` while still holding the namespace
- \* lock. During each step the configuration is synced out.
- \* Remove a device from the pool. Currently, this supports removing only hot
- \* spares, slogs, and level 2 ARC devices.

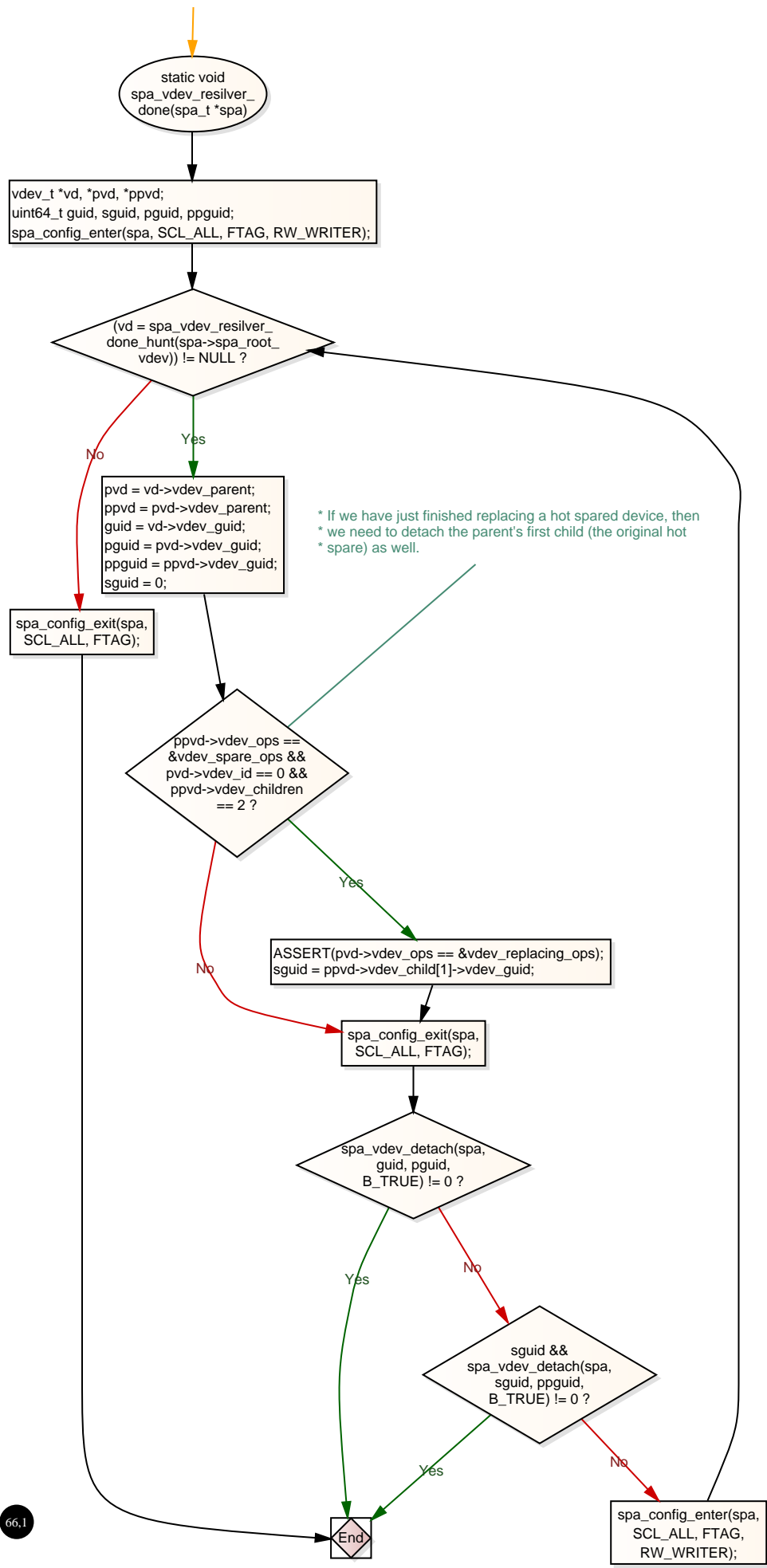
`spa->spa_spares.sav_vdevs`  
`!= NULL &&`  
`nvlist_lookup_nvlist_`  
`array(spa->spa_spares`  
`.sav_config,`  
`ZPOOL_CONFIG_SPARES,`  
`&spares, &nspares) == 0`  
`&& (nv =`  
`spa_nvlist_lookup_by_guid`  
`(spares, nspares,`  
`guid)) != NULL ?`

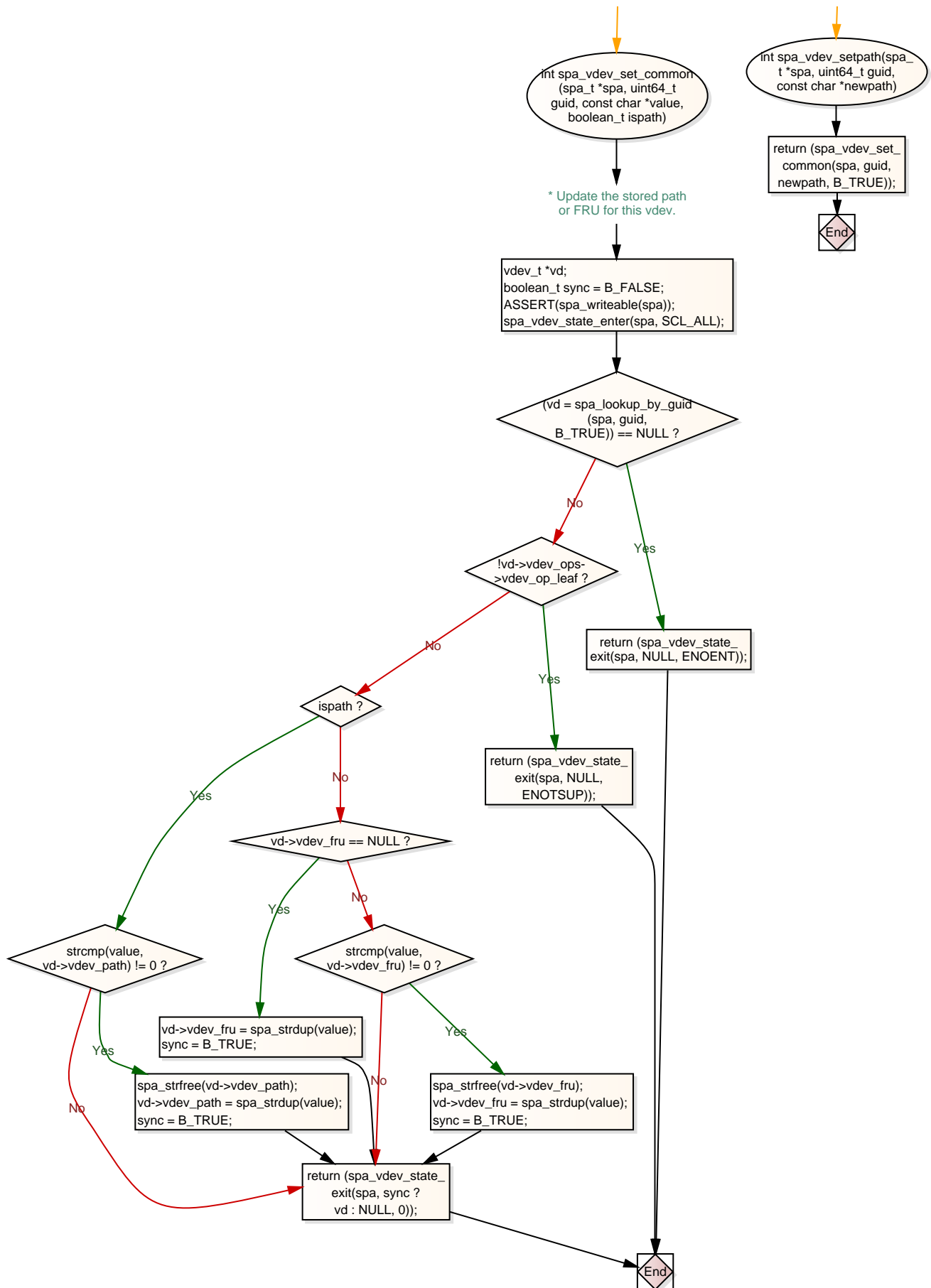
No



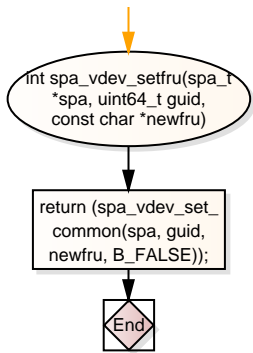








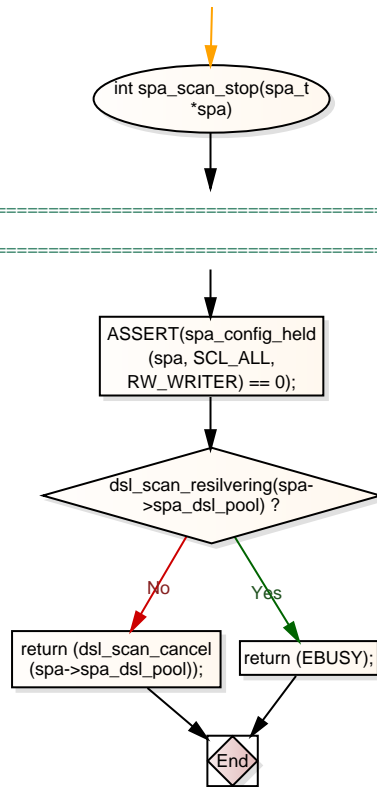


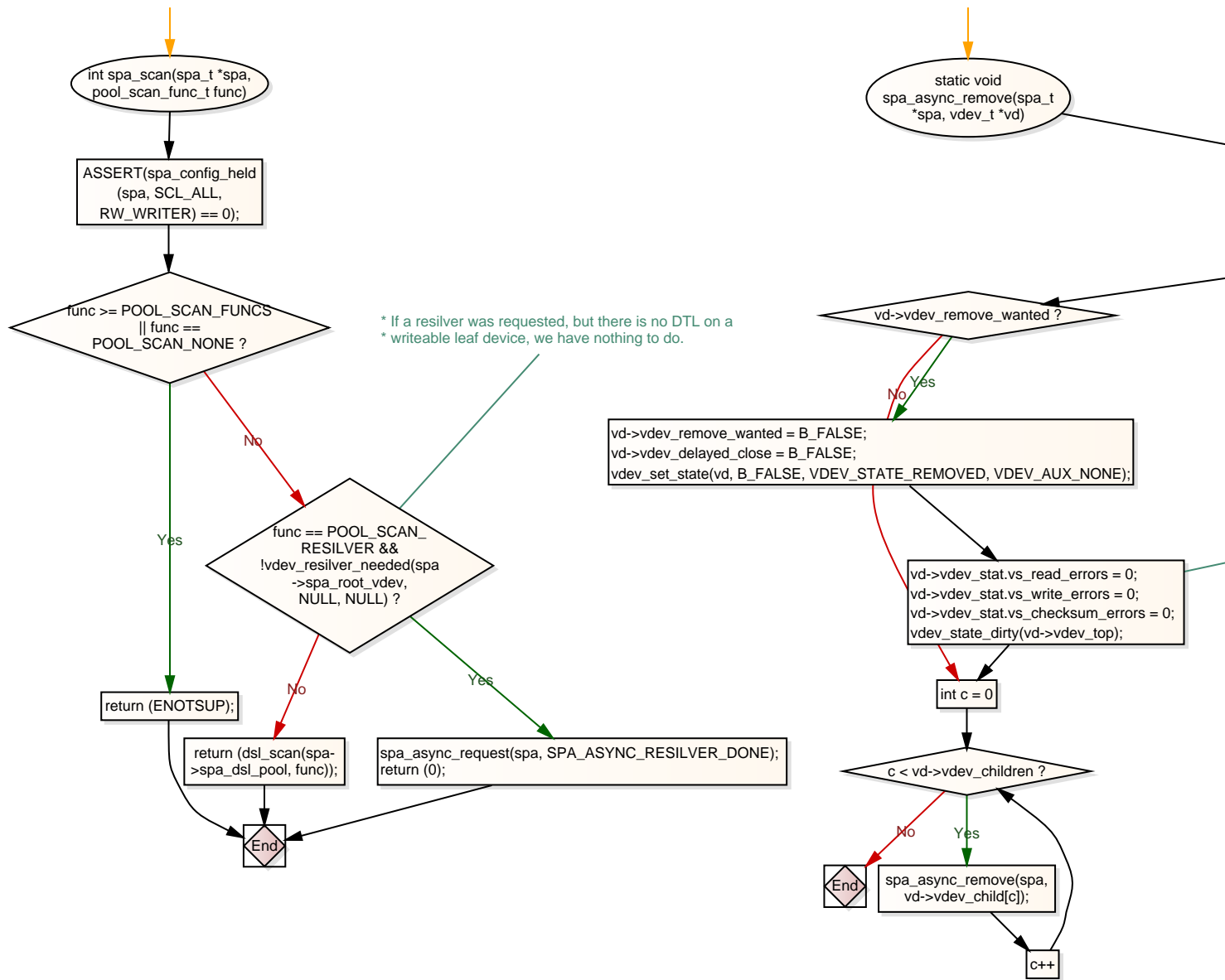


=====

\* SPA Scanning

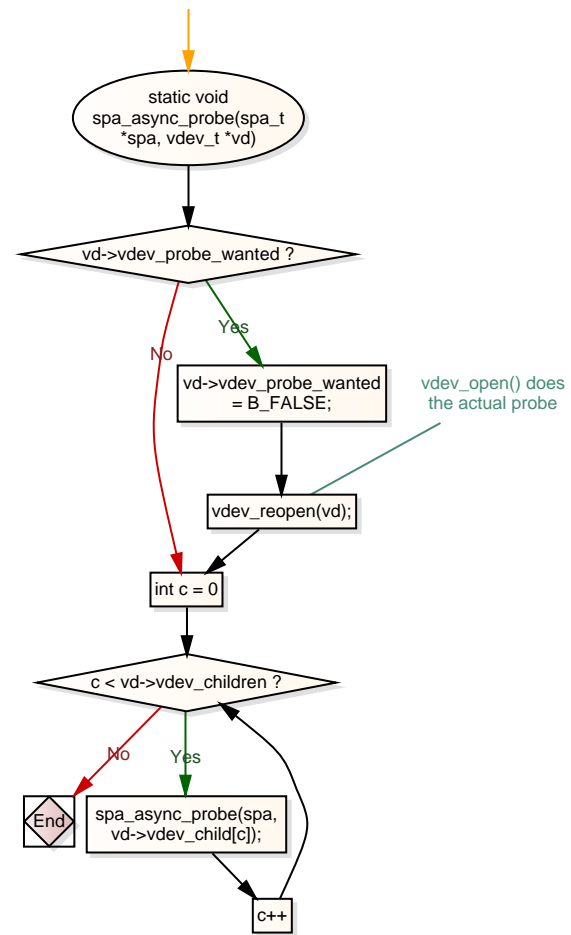
=====

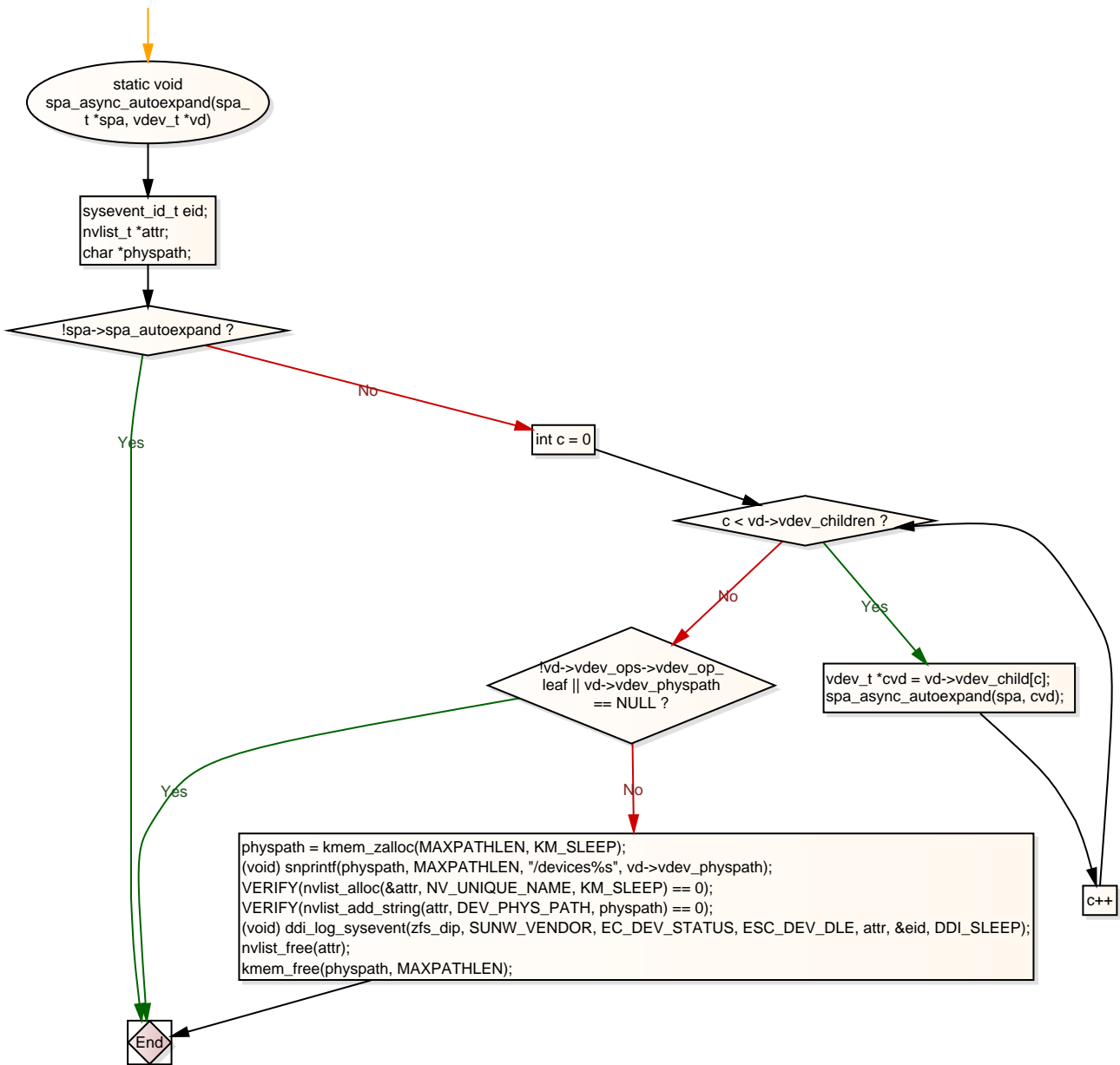


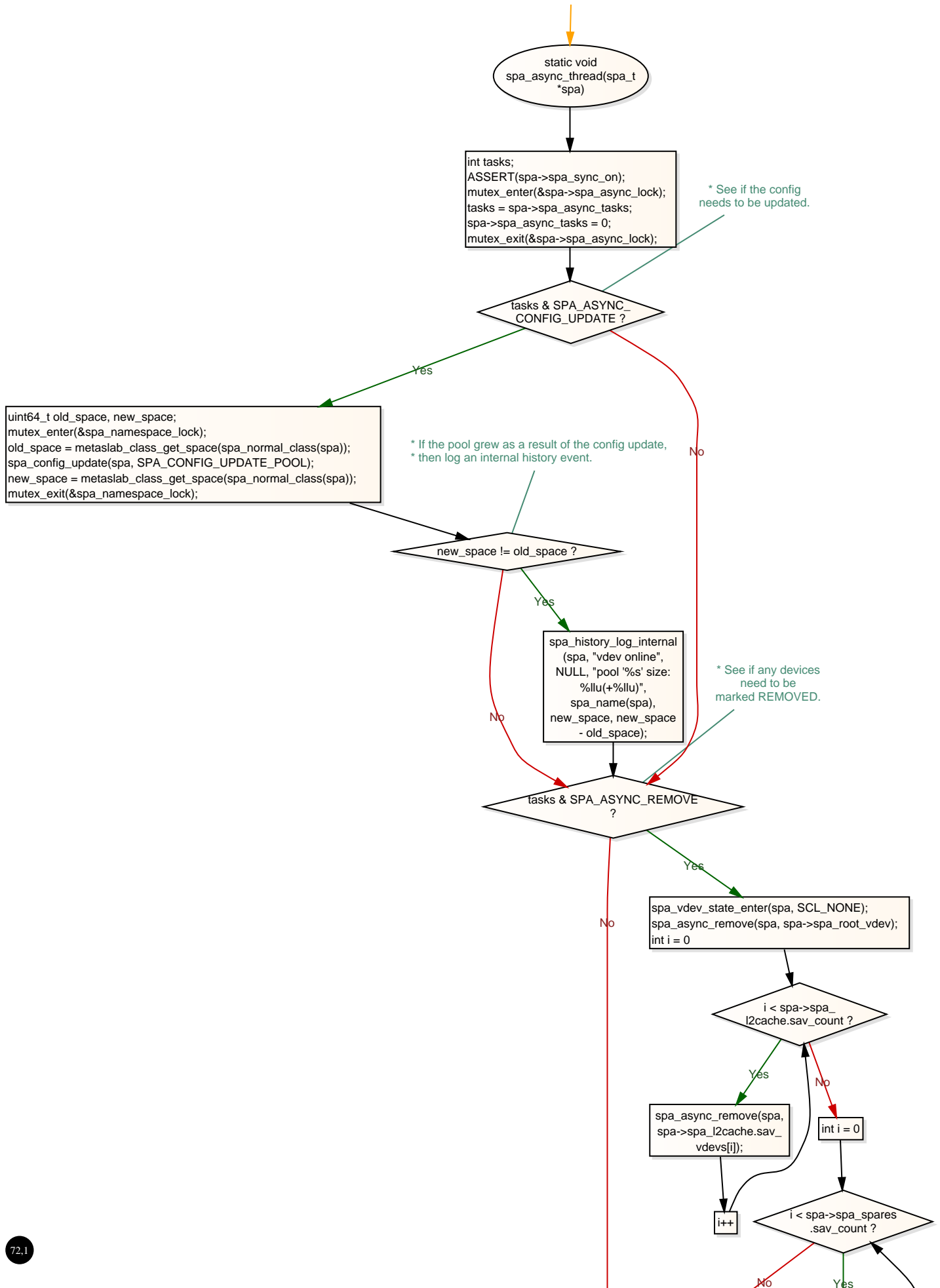


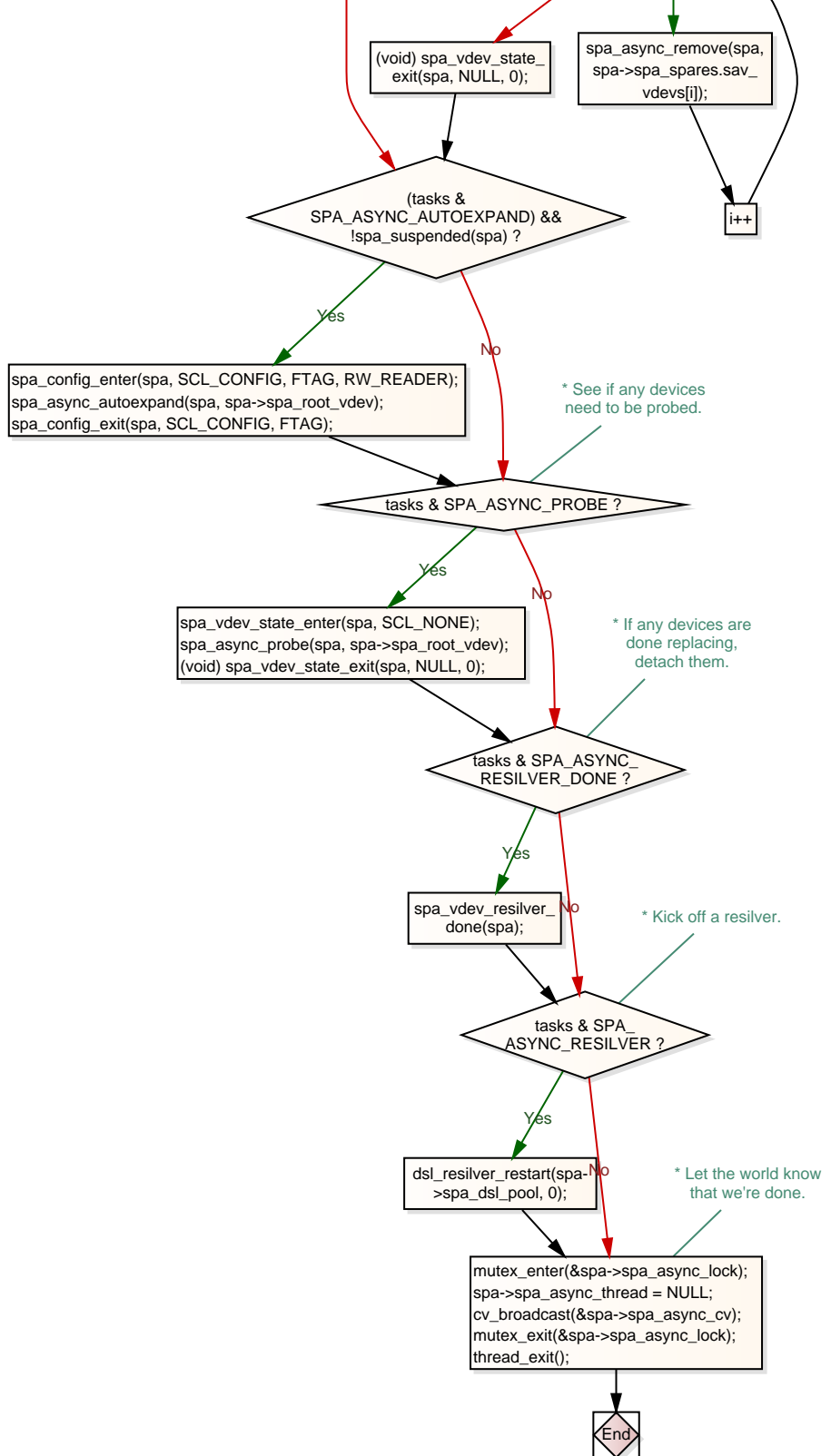
\* SPA async task processing

\* We want to clear the stats, but we don't want to do a full  
\* vdev\_clear() as that will cause us to throw away  
\* degraded/faulted state as well as attempt to reopen the  
\* device, all of which is a waste.









void spa\_async\_  
suspend(spa\_t \*spa)

mutex\_enter(&spa->spa\_async\_lock);  
spa->spa\_async\_suspended++;

spa->spa\_async\_thread  
!= NULL ?

No

Yes

mutex\_exit(&spa->spa\_async\_lock);

End

cv\_wait(&spa->spa\_async\_  
cv, &spa->spa\_  
async\_lock);

void spa\_async\_  
resume(spa\_t \*spa)

mutex\_enter(&spa->spa\_async\_lock);  
ASSERT(spa->spa\_async\_suspended != 0);  
spa->spa\_async\_suspended--;  
mutex\_exit(&spa->spa\_async\_lock);

End

static void  
spa\_async\_dispatch(spa\_t  
\*spa)

mutex\_enter(&spa  
->spa\_async\_lock);

spa->spa\_async\_tasks &&  
!spa->spa\_async\_suspended  
&& spa->spa\_async\_thread  
== NULL && rootdir !=  
NULL && !vn\_is\_  
readonly(rootdir) ?

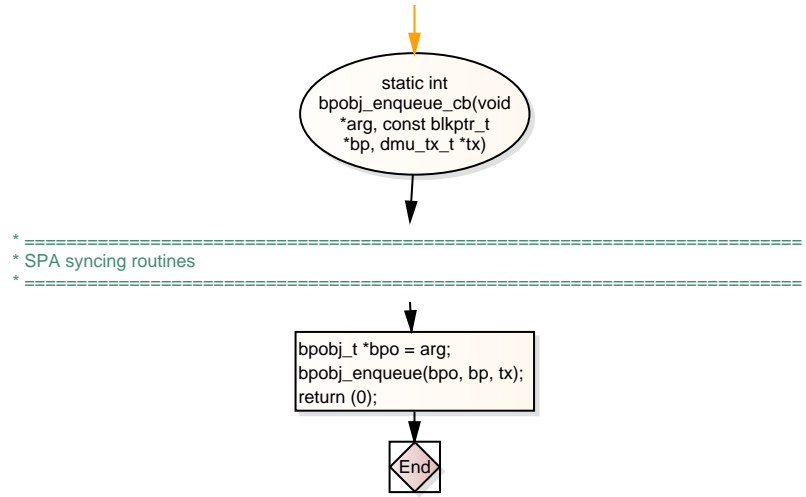
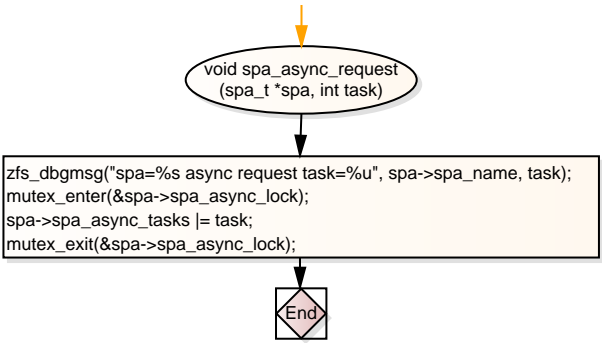
Yes

No

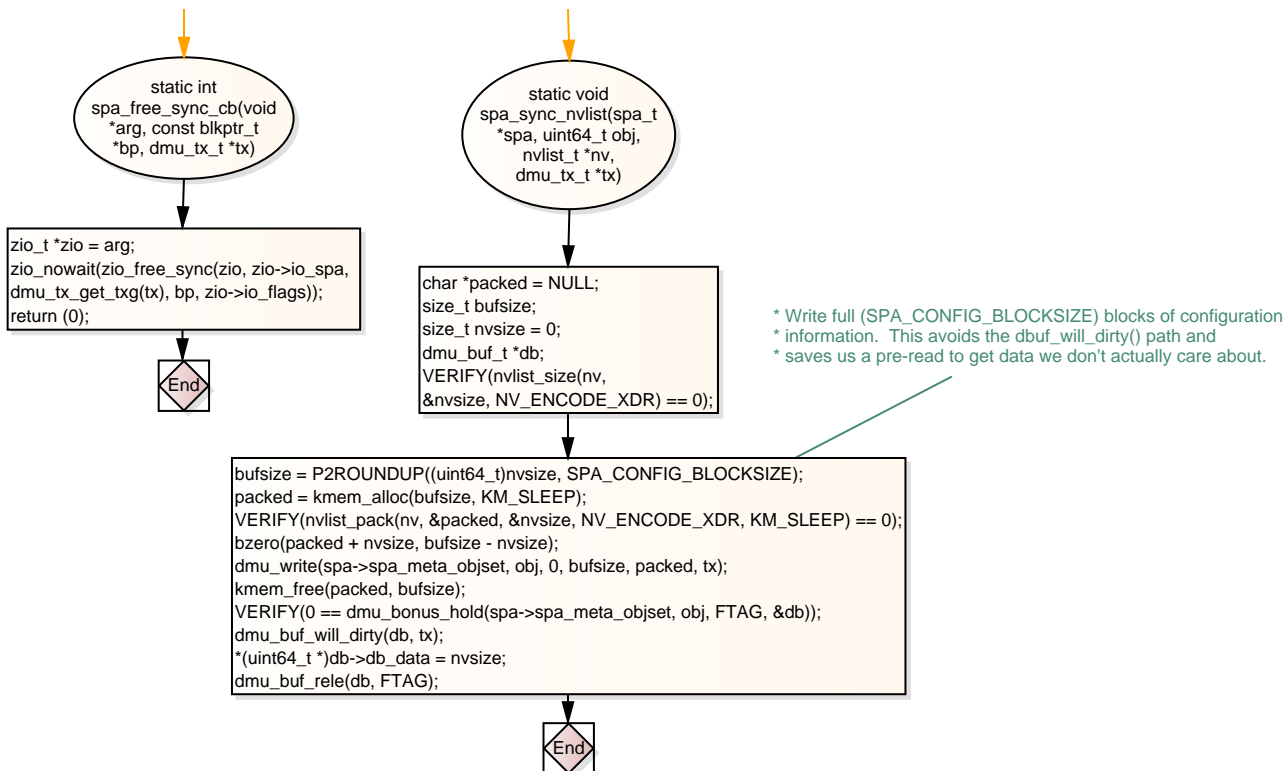
spa->spa\_async\_thread =  
thread\_create(NULL, 0,  
spa\_async\_thread, spa,  
0, &p0, TS\_RUN,  
maxclsypri);

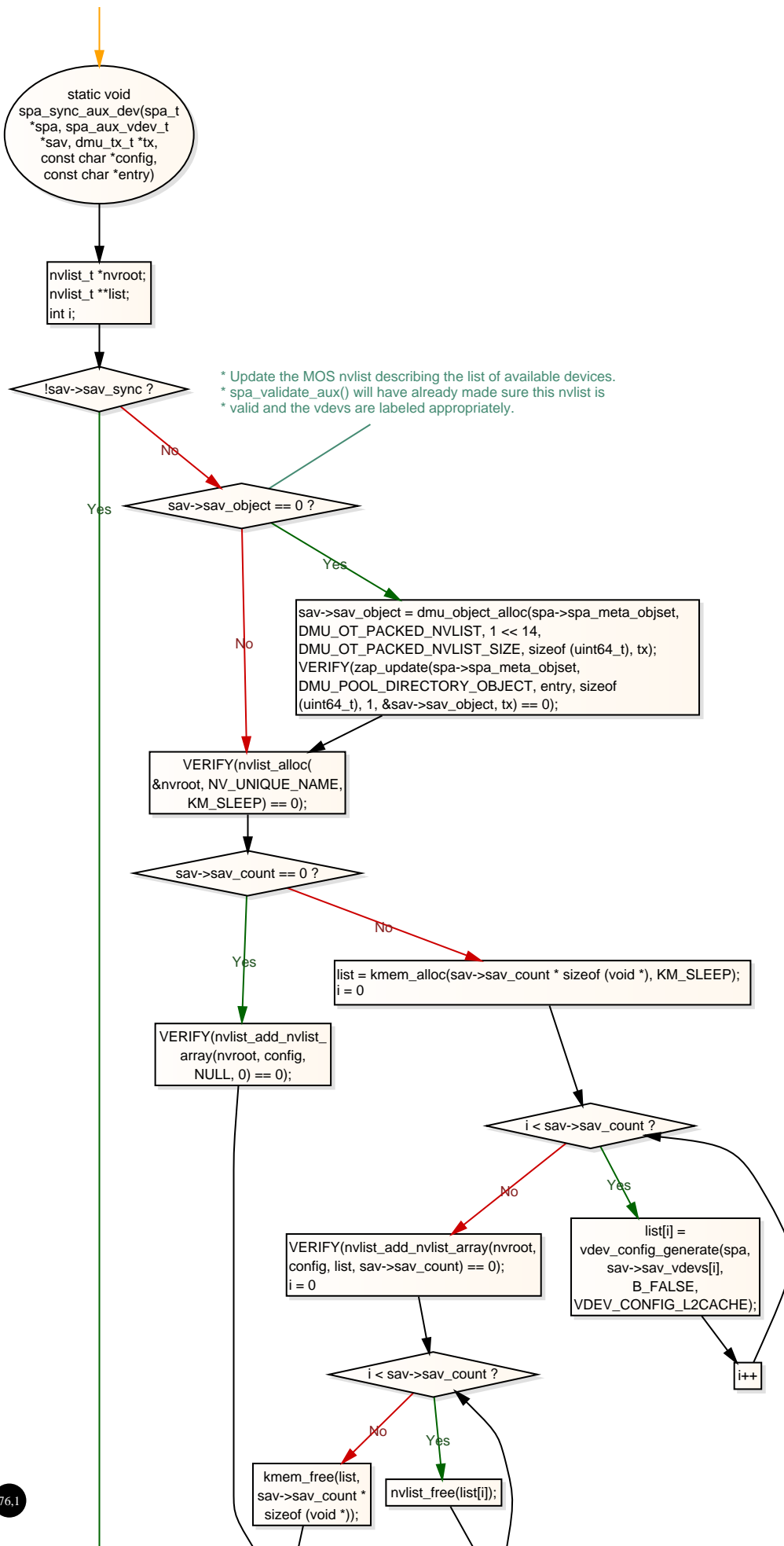
mutex\_exit(&spa->spa\_async\_lock);

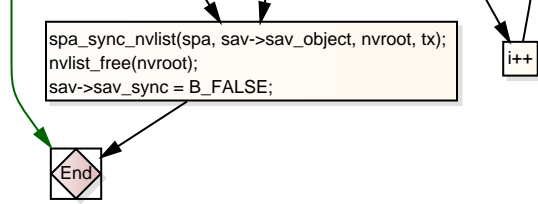
End

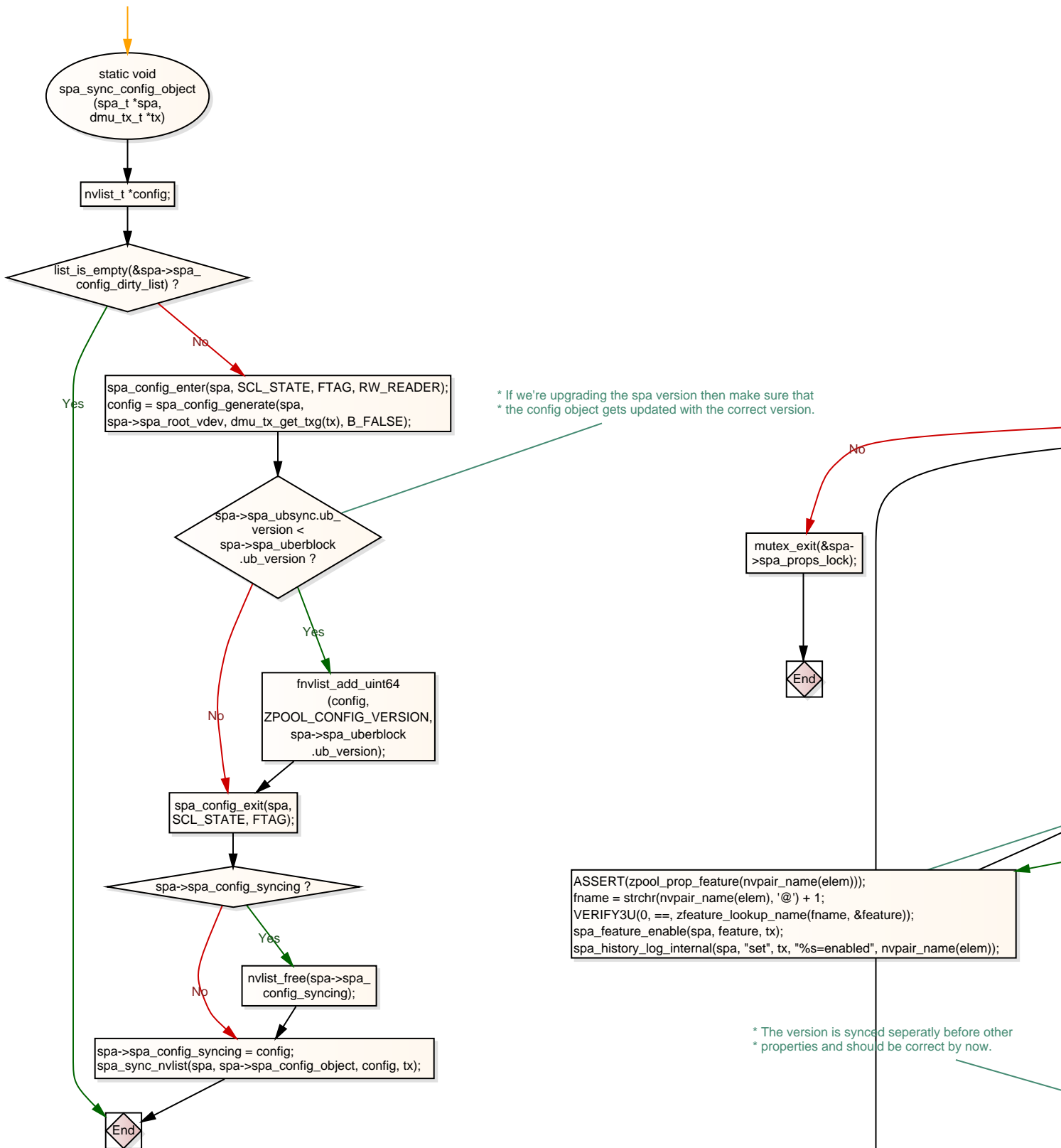










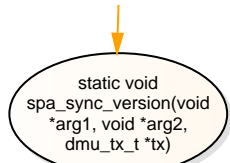


spa->spa\_comment =  
spa\_strdup(strval);

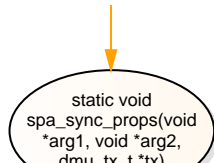
vdev\_config\_dirty(spa-  
>spa\_root\_vdev);

spa\_history\_log\_internal  
(spa, "set", tx,  
"%s=%s",  
nvpair\_name(elem),  
strval);

No



\* Setting the version is special cased when first creating the pool.



\* Set zpool properties.

```

ASSERT(tx->tx_tgx != TXG_INITIAL);
ASSERT(version <= SPA_VERSION);
ASSERT(version >= spa_version(spa));
spa->spa_uberblock.ub_version = version;
vdev_config_dirty(spa->spa_root_vdev);
spa_history_log_internal(spa, "set", tx, "version=%lld", version);
  
```



(elem = nvlist\_next\_nvpair(nvp, elem)) ?

Yes

```

uint64_t intval;
char *strval, *fname;
zpool_prop_t prop;
const char *propname;
zprop_type_t proptype;
zfeature_info_t *feature;
  
```

switch (prop =  
zpool\_name\_to\_prop  
(nvpair\_name(elem)))

ZPROP\_INVALID ?

No

ZPOOL\_PROP\_VERSION ?

No

ZPOOL\_PROP\_ALTROOT ?

No

ZPOOL\_PROP\_READONLY ?

No

ZPOOL\_PROP\_CACHEFILE ?

No

ZPOOL\_PROP\_COMMENT ?

Yes

\* We checked this earlier in spa\_prop\_validate().

Yes

```

VERIFY(nvpair_value_uint64(elem,  
&intval) == 0);
  
```

```

ASSERT3U(spa_version(spa), >=, intval);
  
```

\* 'altroot' is a non-persistent property. It should have been set temporarily at creation or import time.

Yes

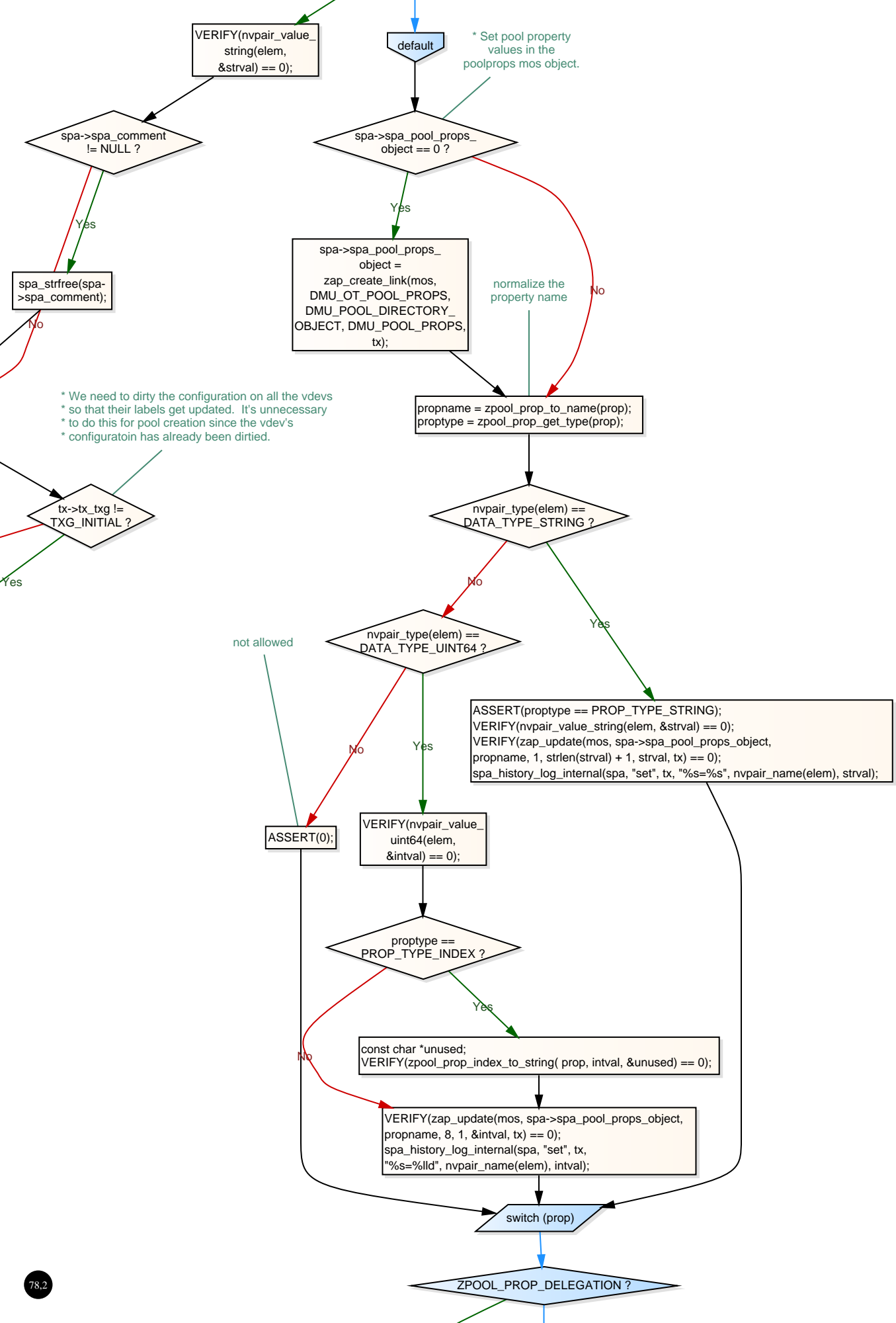
```

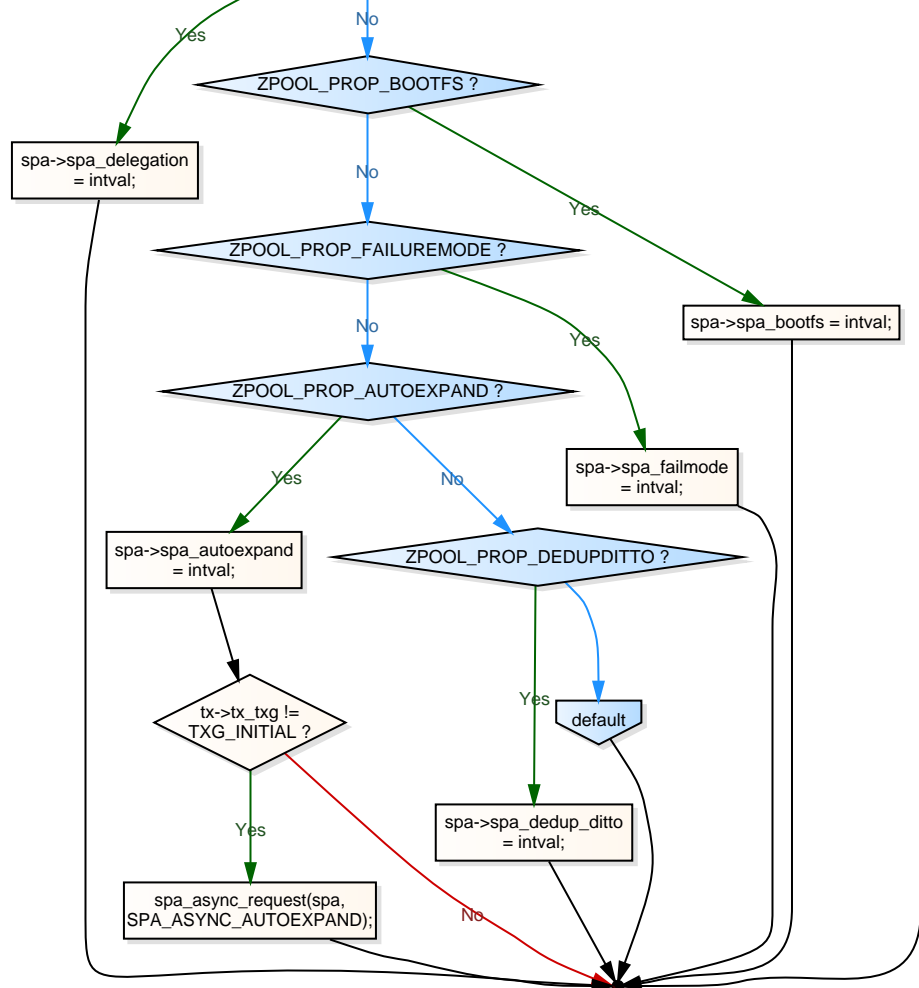
ASSERT(spa->spa_root != NULL);
  
```

Yes

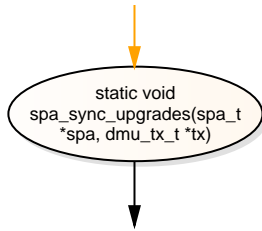
Yes

\* 'readonly' and 'cache' are also non-persistent properties.





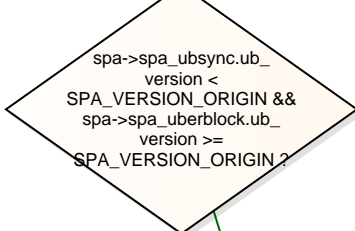




\* Perform one-time upgrade on-disk changes. spa\_version() does not  
 \* reflect the new version this txg, so there must be no changes this  
 \* txg to anything that the upgrade code depends on after it executes.  
 \* Therefore this must be called after dsl\_pool\_sync() does the sync  
 \* tasks.

```

dsl_pool_t *dp = spa->spa_dsl_pool;
ASSERT(spa->spa_sync_pass == 1);
  
```



Yes

```

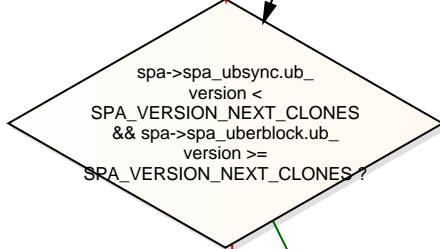
dsl_pool_create_origin
(dp, tx);
  
```

Keeping the origin open  
increases spa\_minref

No

```

spa->spa_minref += 3;
  
```

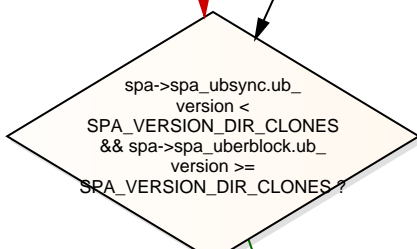


Yes

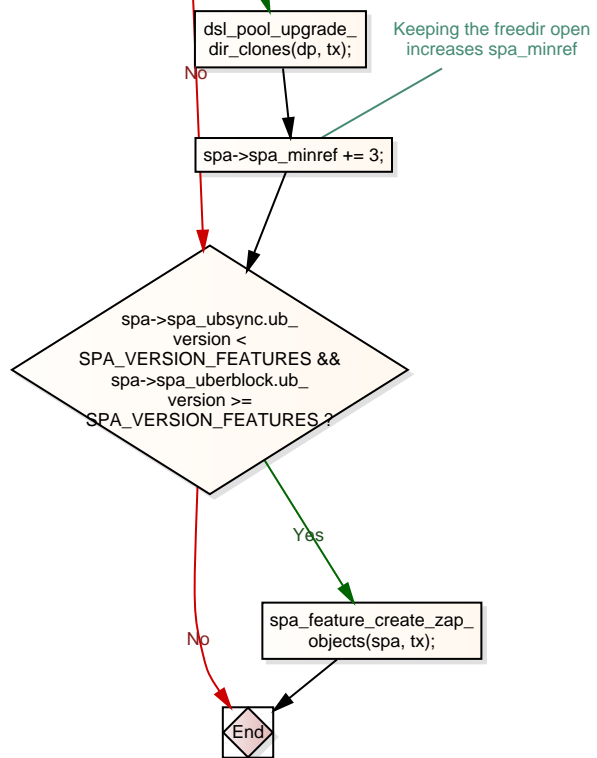
No

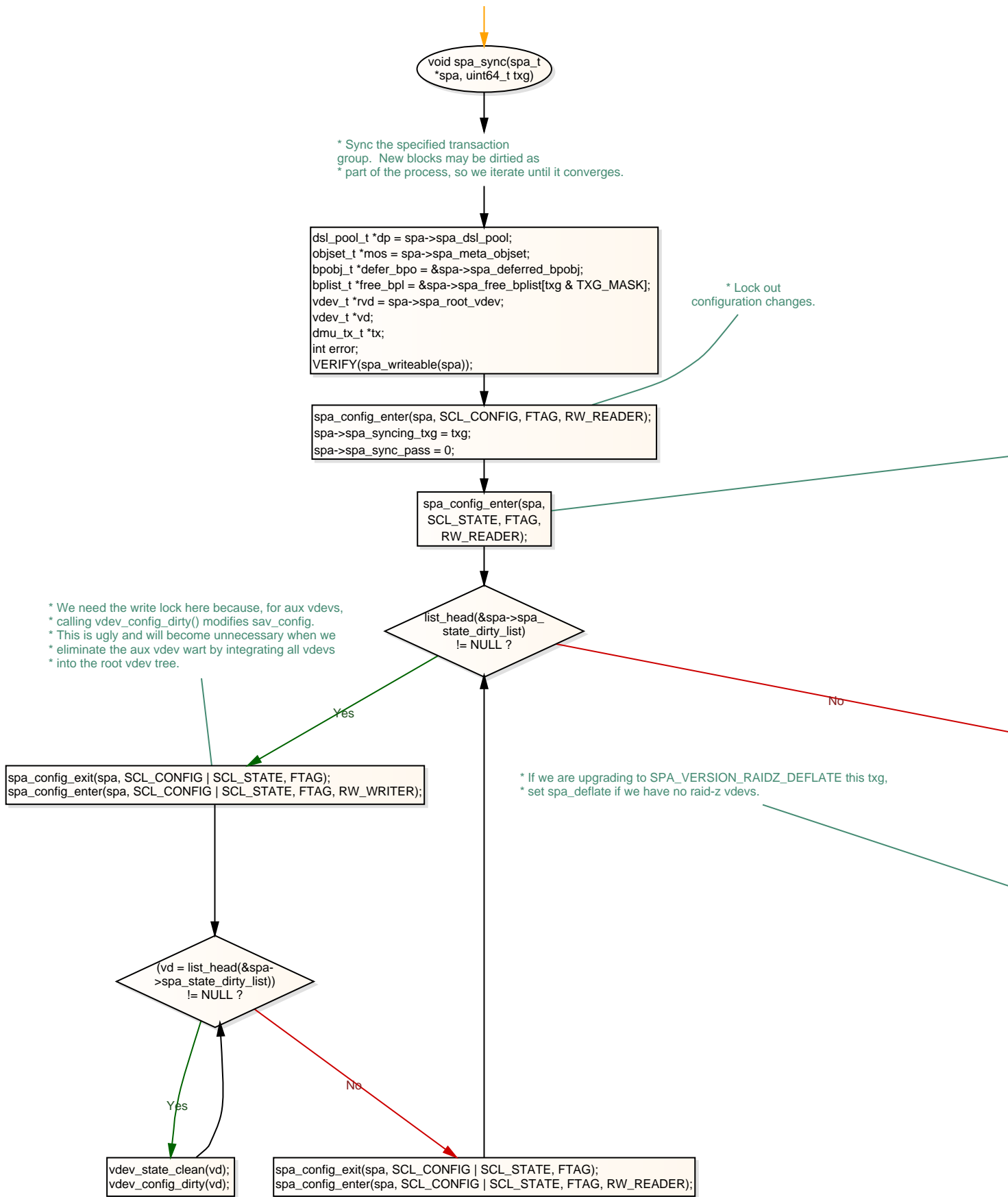
```

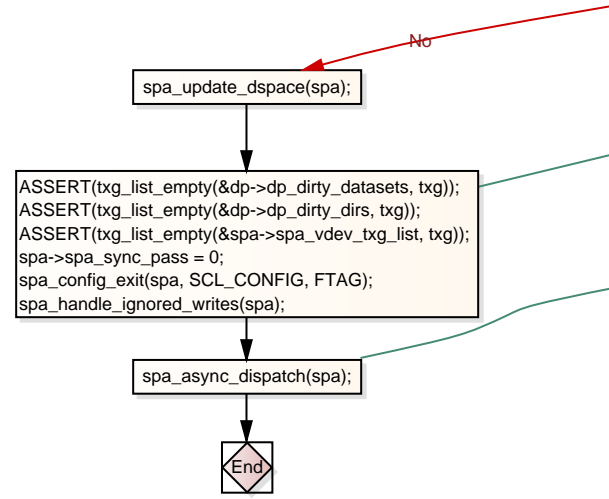
dsl_pool_upgrade_
clones(dp, tx);
  
```



Yes







void spa\_  
sync\_allpools(void)

\* Sync all pools. We don't want to hold the namespace lock across these  
\* operations, so we take a reference on the spa\_t and drop the lock during the  
\* sync.

spa\_t \*spa = NULL;  
mutex\_enter(&spa\_namespace\_lock);

(spa = spa\_next(spa))  
!= NULL ?

Yes

Yes

spa\_state(spa) !=  
POOL\_STATE\_ACTIVE ||  
!spa\_writable(spa) ||  
spa\_suspended(spa) ?

No

spa\_open\_ref(spa, FTAG);  
mutex\_exit(&spa\_namespace\_lock);  
txg\_wait\_synced(spa\_get\_dsl(spa), 0);  
mutex\_enter(&spa\_namespace\_lock);  
spa\_close(spa, FTAG);

\* If there are any pending vdev state changes, convert them  
\* into config changes that go out with this transaction group.

mutex\_exit(  
&spa\_namespace\_lock);

End

spa\_config\_exit(spa, SCL\_STATE, FTAG);  
tx = dmu\_tx\_create\_assigned(dp, txg);  
spa->spa\_sync\_starttime = gethrtime();  
VERIFY(cyclic\_reprogram(spa->spa\_deadman\_cycid,  
spa->spa\_sync\_starttime + spa->spa\_deadman\_syntime));

spa->spa\_ubsync.ub\_  
version <  
SPA\_VERSION\_RAIDZ\_DEFLATE  
&& spa->spa\_uberblock.ub\_  
version >=  
SPA\_VERSION\_RAIDZ\_DEFLATE  
?

Yes

int i;  
i = 0

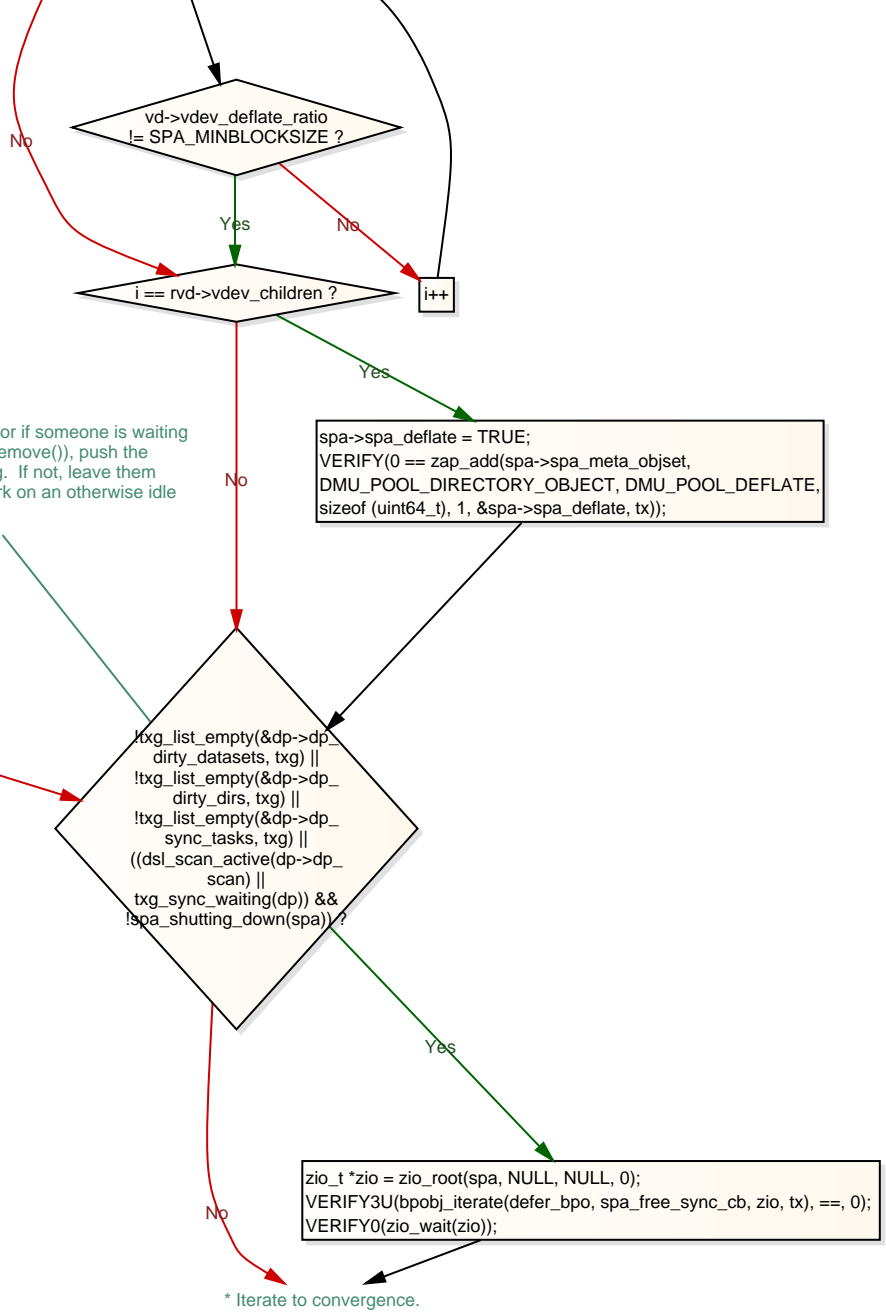
i < rvd->vdev\_children ?

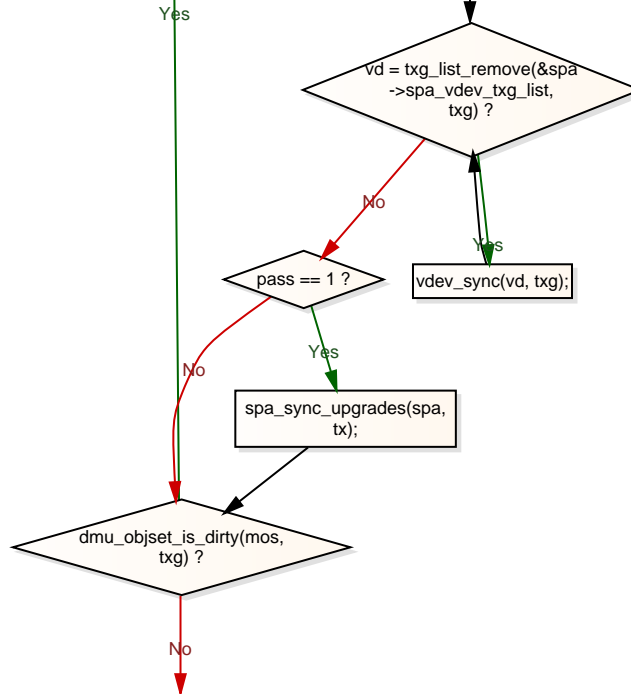
Yes

vd = rvd->vdev\_child[i];

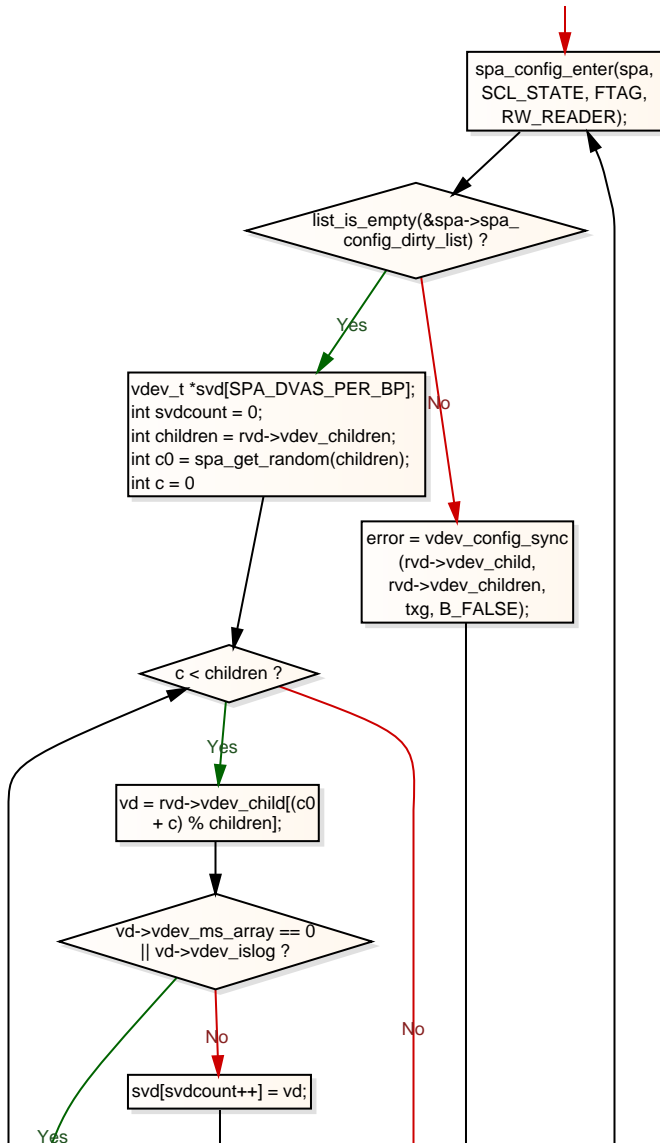
No

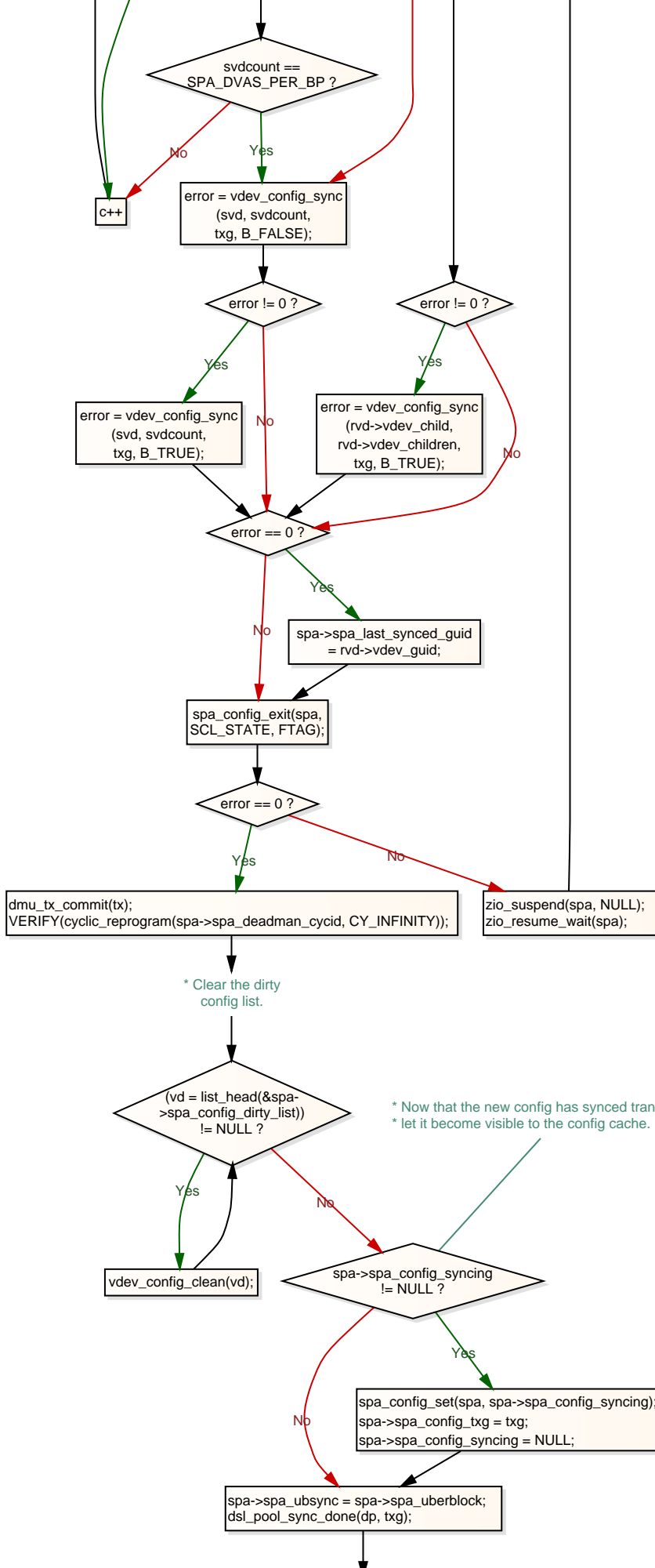
\* If anything has changed in this txg, or if someone is waiting  
 \* for this txg to sync (eg, spa\_vdev\_remove()), push the  
 \* deferred frees from the previous txg. If not, leave them  
 \* alone so that we don't generate work on an otherwise idle  
 \* system.





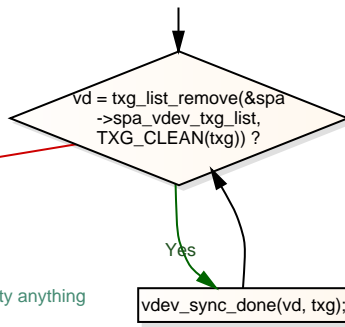
\* Rewrite the vdev configuration (which includes the uberblock)  
 \* to commit the transaction group.  
 \*  
 \* If there are no dirty vdevs, we sync the uberblock to a few  
 \* random top-level vdevs that are known to be visible in the  
 \* config cache (see spa\_vdev\_add() for a complete description).  
 \* If there "are" dirty vdevs, sync the uberblock to all vdevs.





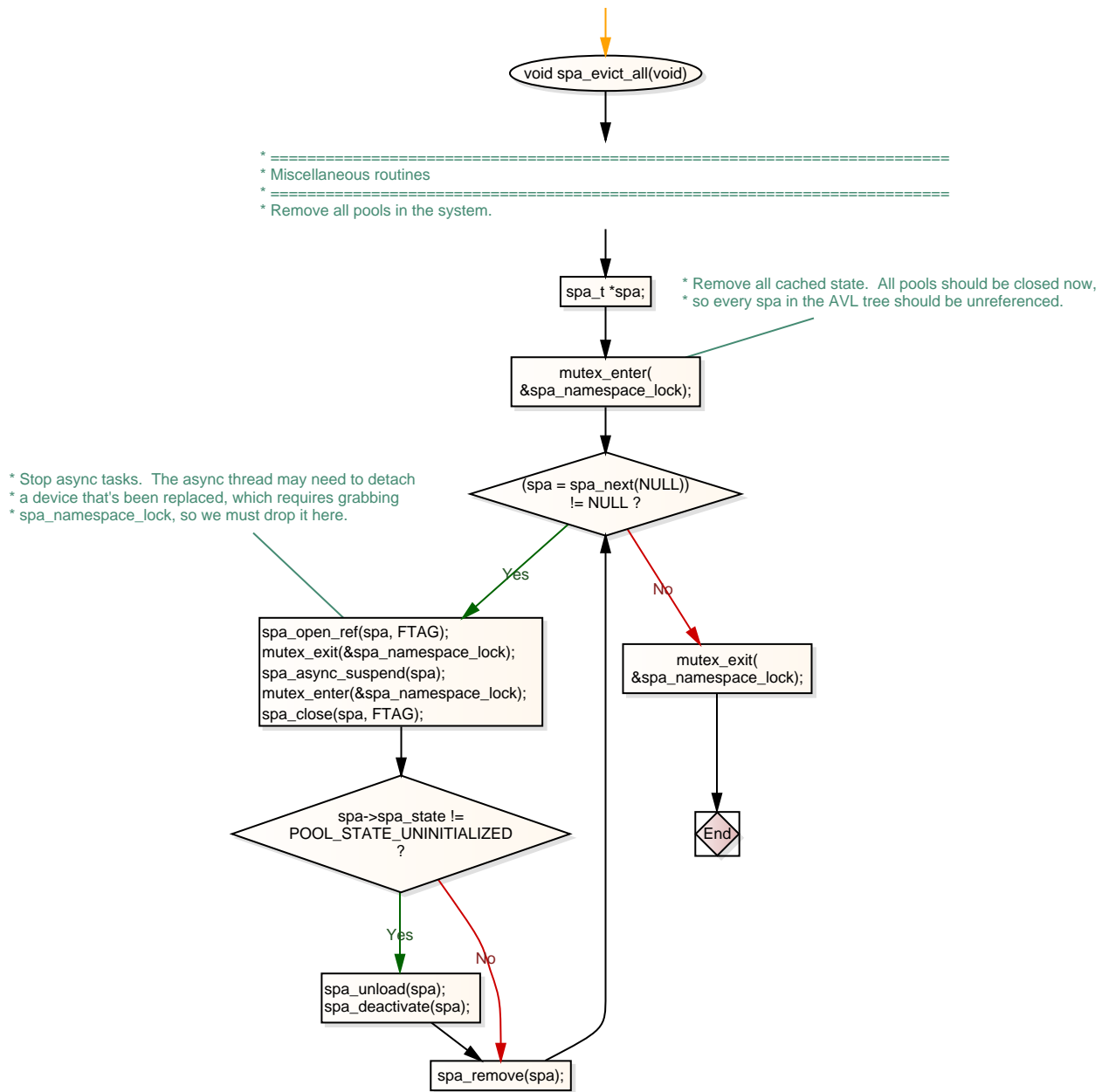


\* Update usable  
space statistics.

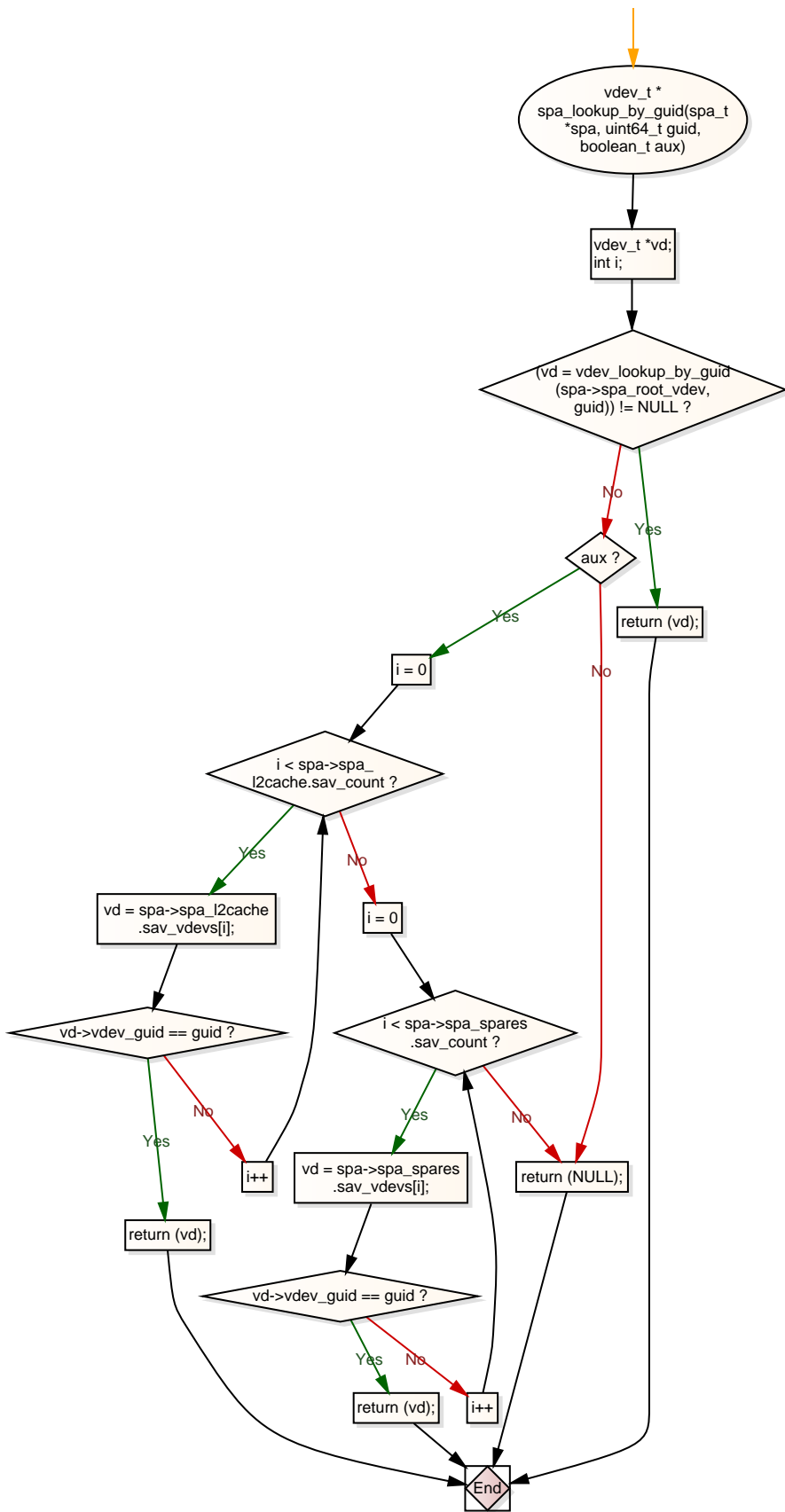


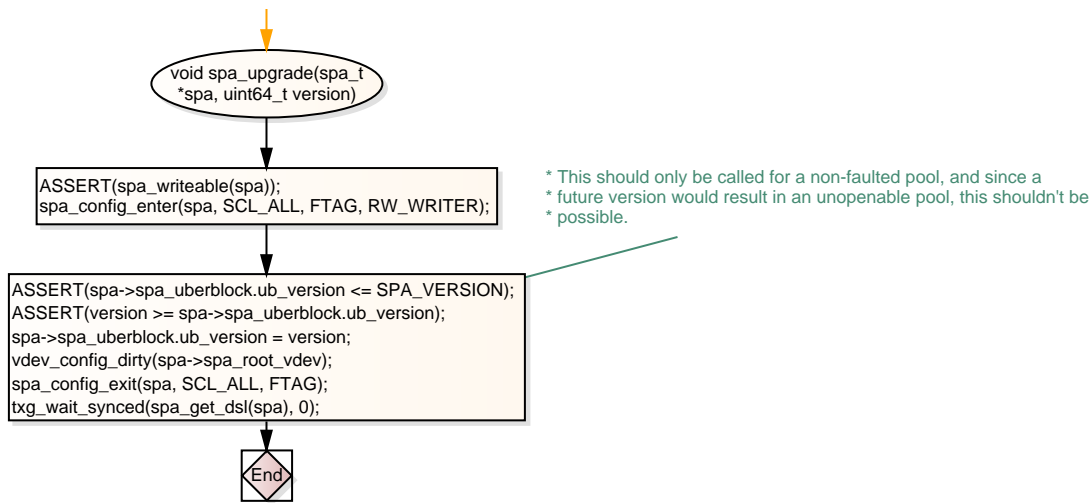
\* It had better be the case that we didn't dirty anything  
\* since `vdev_config_sync()`.

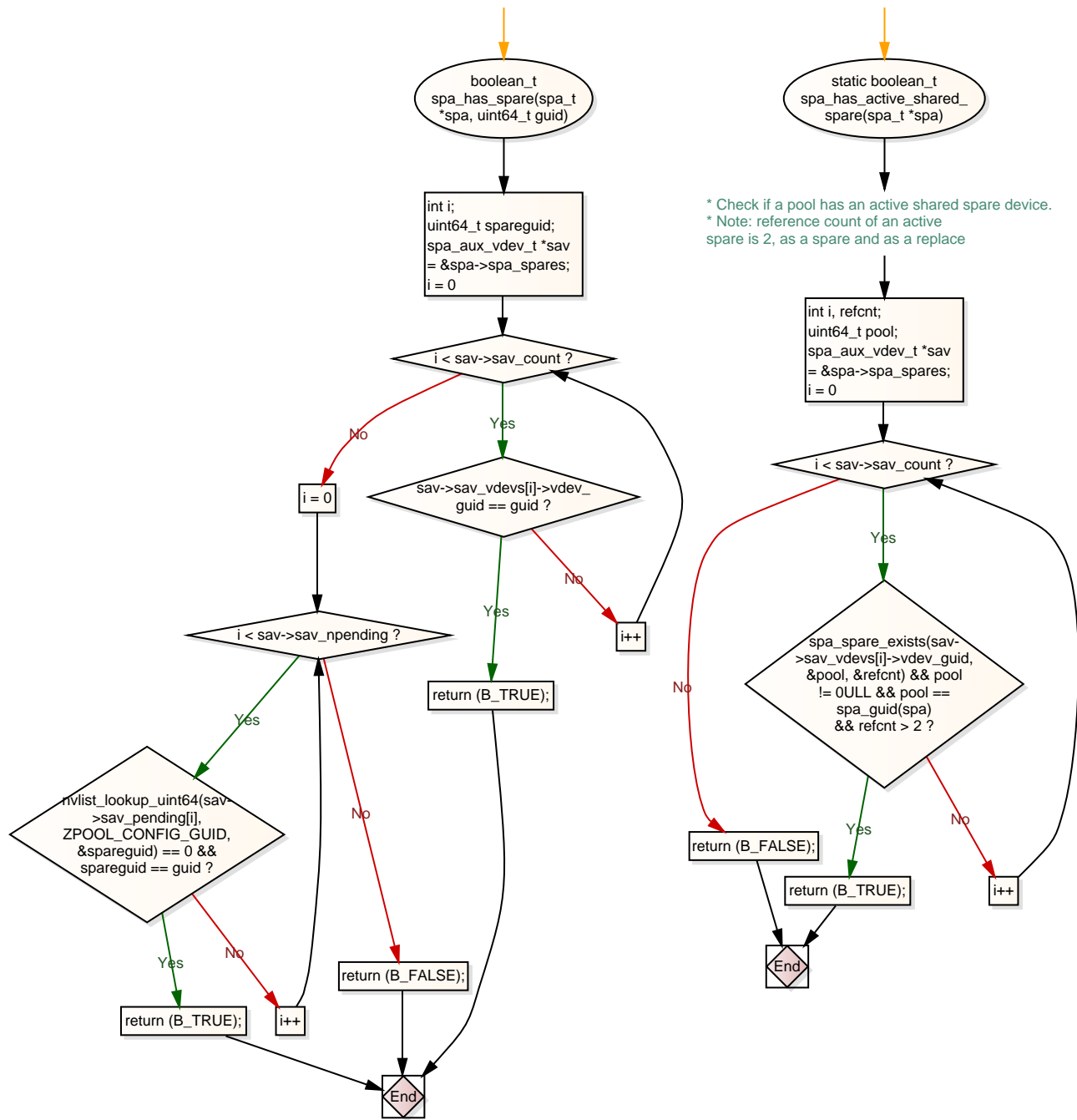
\* If any async tasks  
have been requested,  
kick them off.

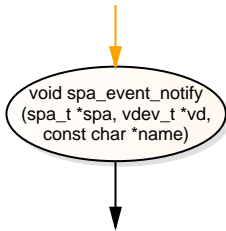


\* We hold SCL\_STATE to prevent vdev open/close/etc.  
\* while we're attempting to write the vdev labels.









\* Post a sysevent corresponding to the given event. The 'name' must be one of  
 \* the event definitions in sys/sysevent/eventdefs.h. The payload will be  
 \* filled in from the spa and (optionally) the vdev. This doesn't do anything  
 \* in the userland libzpool, as we don't want consumers to misinterpret ztest  
 \* or zdb as real changes.

