```
#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/metaslab.h>
#include <sys/metaslab_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>
```

#ifdef  _KERNEL ?

Yes        No

```
#include <sys/bootprops.h>
#include <sys/callb.h>
#include <sys/cpupart.h>
#include <sys/pool.h>
#include <sys/sysdc.h>
#include <sys/zone.h>
```

_KERNEL

```
#include "zfs_prop.h"
#include "zfs_comutil.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.

#ifdef _KERNEL ?

Yes

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

No

End

```
typedef struct
zio_taskq_info
```

```
enum zti_modes zti_mode;
uint_t zti_value;
```

End

```
static void
spa_prop_add_list(nvlist_
t *nvl, zpool_prop_t
prop, char *strval,
uint64_t intval,
zprop_source_t src)
```

* Add a (source=src,
propname=propval)
list to an nvlist.

```
const char *propname = zpool_prop_to_name(prop);
nvlist_t *propval;
VERIFY(nvlist_alloc(&propval, NV_UNIQUE_NAME, KM_SLEEP) == 0);
VERIFY(nvlist_add_uint64(propval, ZPROP_SOURCE, src) == 0);
```

strval != NULL ?

Yes

No

```
VERIFY(nvlist_add_string
(propval, ZPROP_VALUE,
strval) == 0);
```

```
VERIFY(nvlist_add_uint64
(propval, ZPROP_VALUE,
intval) == 0);
```

```
VERIFY(nvlist_add_nvlist(nvl, propname, propval) == 0);
nvlist_free(propval);
```

End

```
static void
spa_prop_get_config(spa_t
*spa, nvlist_t **nvp)
```

* Get property values
from the spa
configuration.

```
vdev_t *rvd = spa->spa_root_vdev;
dsl_pool_t *pool = spa->spa_dsl_pool;
uint64_t size;
uint64_t alloc;
uint64_t space;
uint64_t cap, version;
zprop_source_t src = ZPROP_SRC_NONE;
spa_config_dirent_t *dp;
ASSERT(MUTEX_HELD(&spa->spa_props_lock));
```

rvd != NULL ?

— No
— Yes

```
alloc = metaslab_class_get_alloc(spa_normal_class(spa));
size = metaslab_class_get_space(spa_normal_class(spa));
spa_prop_add_list(*nvp, ZPOOL_PROP_NAME, spa_name(spa), 0, src);
spa_prop_add_list(*nvp, ZPOOL_PROP_SIZE, NULL, size, src);
spa_prop_add_list(*nvp, ZPOOL_PROP_ALLOCATED, NULL, alloc, src);
spa_prop_add_list(*nvp, ZPOOL_PROP_FREE, NULL, size - alloc, src);
space = 0;
int c = 0
```

c < rvd->vdev_children ?

— No
— Yes

```
vdev_t *tvd = rvd->vdev_child[c];
space += tvd->vdev_max_asize - tvd->vdev_asize;
```

C++

```
spa_prop_add_list(*nvp, ZPOOL_PROP_EXPANDSZ, NULL, space, src);
spa_prop_add_list(*nvp, ZPOOL_PROP_READONLY, NULL, (spa_mode(spa) == FREAD), src);
cap = (size == 0) ? 0 : (alloc * 100 / size);
spa_prop_add_list(*nvp, ZPOOL_PROP_CAPACITY, NULL, cap, src);
spa_prop_add_list(*nvp, ZPOOL_PROP_DEDUPRATIO, NULL, ddt_get_pool_dedup_ratio(spa), src);
spa_prop_add_list(*nvp, ZPOOL_PROP_HEALTH, NULL, rvd->vdev_state, src);
version = spa_version(spa);
```

version ==
zpool_prop_default_
numeric(ZPOOL_
PROP_VERSION) ?

— Yes
— No

src = ZPROP_SRC_DEFAULT;

src = ZPROP_SRC_LOCAL;

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_VERSION,
NULL, version, src);
```

pool != NULL ?

— No
— Yes

```
dsl_dir_t *freedir =
pool->dp_free_dir;
```

* The $FREE directory was introduced in SPA_VERSION_DEADLISTS,
* when opening pools before this version freedir will be NULL.

freedir != NULL ?

— Yes
— No

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_FREEING,
NULL, freedir->dd_phys-
>dd_used_bytes, src);
```

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_FREEING,
NULL, 0, src);
```

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_GUID, NULL,
spa_guid(spa), src);
```

spa->spa_comment
!= NULL ?

— No
— Yes

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_COMMENT,
spa->spa_comment, 0,
ZPROP_SRC_LOCAL);
```

spa->spa_root != NULL ?

— No
— Yes

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_ALTROOT,
spa->spa_root, 0,
ZPROP_SRC_LOCAL);
```

(dp = list_head(&spa-
>spa_config_list))
!= NULL ?

— No
— Yes

dp->scd_path == NULL ?

— Yes
— No

strcmp(dp->scd_path,
spa_config_path) != 0 ?

— No
— Yes

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_CACHEFILE,
"none", 0,
ZPROP_SRC_LOCAL);
```

```
spa_prop_add_list(*nvp,
ZPOOL_PROP_CACHEFILE,
dp->scd_path, 0,
ZPROP_SRC_LOCAL);
```

End

static int
spa_prop_validate(spa_t
*spa, nvlist_t *props)

* Validate the given pool properties nvlist and modify the list
* for the property values to be set.

nvpair_t *elem;
int error = 0, reset_bootfs = 0;
uint64_t objnum;
boolean_t has_feature = B_FALSE;
elem = NULL;

(elem = nvlist_next_
nvpair(props,
elem)) != NULL ?

uint64_t intval;
char *strval, *slash, *check, *fname;
const char *propname = nvpair_name(elem);
zpool_prop_t prop = zpool_name_to_prop(propname);

switch (prop)

ZPROP_INVAL

ZPOOL_PROP_VERSION ?

ZPOOL_PROP_DELEGATION ?

ZPOOL_PROP_AUTOREPLACE ?

ZPOOL_PROP_LISTSNAPS ?

ZPOOL_PROP_AUTOEXPAND ?

ZPOOL_PROP_BOOTFS ?

ZPOOL_PROP_FAILUREMODE ?

ZPOOL_PROP_CACHEFILE ?

ZPOOL_PROP_COMMENT ?

ZPOOL_PROP_DEDUPDITTO ?

error = nvpair_value_
uint64(elem, &intval);

!error && (intval <
spa_version(spa) ||
intval >
SPA_VERSION_BEFORE_
FEATURES ||
has_feature) ?

error = EINVAL;

* If the pool version is less than SPA_VERSION_BOOTFS,
* or the pool is still being created (version == 0),
* the bootfs property cannot be set.

spa_version(spa) <
SPA_VERSION_BOOTFS ?

* Make sure the vdev
config is bootable

!vdev_is_bootable(spa-
>spa_root_vdev) ?

error = ENOTSUP;

error = ENOTSUP;

reset_bootfs = 1;
error = nvpair_value_string(elem, &strval);

error ?

objset_t *os;
uint64_t compress;

error = nvpair_value_
uint64(elem, &intval);

error && intval > 1 ?

error = EINVAL;

error = nvpair_value_
uint64(elem, &intval);

!error && (intval <
ZIO_FAILURE_MODE_WAIT ||
intval >
ZIO_FAILURE_MODE_PANIC) ?

error = EINVAL;

* This is a special case which only occurs when
* the pool has completely failed. This allows
* the user to change the in-core failmode property
* without syncing it out to disk (I/Os might
* currently be blocked). We do this by returning
* EIO to the caller (spa_prop_set) to trick it
* into thinking we encountered a property validation
* error.

strval == NULL ||
strval[0] == '\0' ?

error = dmu_objset_hold(
strval, FTAG, &os) ?

objnum =
zpool_prop_default_
numeric(
ZPOOL_PROP_BOOTFS);

dmu_objset_type(os)
!= DMU_OST_ZFS ?

(error =
dsl_prop_get_integer
(strval,
zfs_prop_to_name(ZFS_
PROP_COMPRESSION),
&compress, NULL)) == 0
&& !BOOTFS_COMPRESS_
VALID(compress) ?

* Must be ZPL and not
gzip compressed.

error = ENOTSUP;

error = ENOTSUP;

objnum =
dmu_objset_id(os);

dmu_objset_rele(os,
FTAG);

(error =
nvpair_value_string(elem,
&strval)) != 0 ?

spa_version(spa) <
SPA_VERSION_DEDUP ?

error = nvpair_value_
uint64(elem, &intval);

error = ENOTSUP;

check = strval;

check != '\0' ?

* The kernel doesn't have an easy isprint()
* check. For this kernel check, we merely
* check ASCII apart from DEL. Fix this if
* there is an easy-to-use kernel isprint().

check >= 0x7f ?

error = EINVAL;

check++;

error = EINVAL;

check++;

error == 0 && intval !=
0 && intval <
ZIO_DEDUPDITTO_MIN ?

error = EINVAL;

strval[0] == '\0' ?

!strcmp(strval,
"none") == 0 ?

strval[0] != '/' ?

slash = strrchr(strval, '/');
ASSERT(slash != NULL);

strlen(strval) >
ZPROP_MAX_COMMENT ?

error = E2BIG;

slash[1] == '\0' ||
strcmp(slash, "/.") == 0
|| strcmp(slash,
"/..") == 0 ?

error = EINVAL;

error = EINVAL;

!error &&
spa_suspended(spa) ?

spa->spa_failmode = intval;
error = EIO;

* Sanitize the input.

nvpair_type(elem) !=
DATA_TYPE_UINT64 ?

zpool_prop_feature
(propname) ?

error = EINVAL;

nvpair_value_uint64(elem,
&intval) != 0 ?

error = EINVAL;

intval != 0 ?

error = EINVAL;

fname = strchr(propname,
'@') + 1;

error = EINVAL;

zfeature_lookup_name
(fname, NULL) != 0 ?

error = EINVAL;

has_feature = B_TRUE;

error

!error && reset_bootfs ?

error = nvlist_remove
(props, zpool_prop_to_
name(ZPOOL_PROP_BOOTFS),
DATA_TYPE_STRING);

error

error = nvlist_add_uint64
(props, zpool_prop_to_
name(ZPOOL_PROP_BOOTFS),
objnum);

return (error);

```
void spa_configfile_set
(spa_t *spa, nvlist_t
*nvp, boolean_t
need_sync)
```

```
char *cachefile;
spa_config_dirent_t *dp;
```

```
nvlist_lookup_string(nvp,
zpool_prop_to_name(ZPOOL_
PROP_CACHEFILE),
&cachefile) != 0 ?
```

No

```
dp = kmem_alloc(sizeof
(spa_config_dirent_t),
KM_SLEEP);
```

```
cachefile[0] == '\0' ?
```

Yes

No

```
strcmp(cachefile,
"none") == 0 ?
```

Yes

Yes

No

```
dp->scd_path =
spa_strdup(spa_
config_path);
```

```
dp->scd_path = NULL;
```

```
dp->scd_path =
spa_strdup(cachefile);
```

```
list_insert_head(&spa-
>spa_config_list, dp);
```

```
need_sync ?
```

No

Yes

```
spa_async_request(spa,
SPA_ASYNC_CONFIG_UPDATE);
```

End

```
                          int spa_prop_set(spa_t
                          *spa, nvlist_t *nvp)


                          int error;
                          nvpair_t *elem = NULL;
                          boolean_t need_sync = B_FALSE;


                              (error =
                          spa_prop_validate(spa,
                               nvp)) != 0 ?

              Yes                              No                     (elem = nvlist_next_
                                                                       nvpair(nvp,
                                                                      elem)) != NULL ?
   return (error);
                                                          Yes    Yes                           No          No

                                                               zpool_prop_t prop =
                                                               zpool_name_to_prop
                                                               (nvpair_name(elem));


                                                              prop == ZPOOL_PROP_
                                                              CACHEFILE || prop ==
                                                              ZPOOL_PROP_ALTROOT ||
                                                                   prop == ZPOOL_
                                                                  PROP_READONLY ?

                                                                         No

                                                              prop == ZPOOL_PROP_
                                                               VERSION || prop
                                                               == ZPROP_INVAL ?

                                                             Yes              No

                                                         uint64_t ver;          need_sync = B_TRUE;


                                                         prop == ZPOOL_PROP_                                need_sync ?
                                                              VERSION ?
                                                       No              Yes                              Yes         No

    Save time if the version
    is already set.                  ASSERT(zpool_prop_feature(nvpair_name(elem)));   VERIFY(nvpair_value_
                                 Yes    ver = SPA_VERSION_FEATURES;               uint64(elem, &ver) == 0);
                                        need_sync = B_TRUE;                                                            return (0);

   * In addition to the pool directory object, we might
   * create the pool properties object, the features for         ver == spa_version(spa) ?
   * read object, the features for write object, or the
   * feature descriptions object.
                                                                         No               return (dsl_sync_task_do
                                                                                           (spa_get_dsl(spa), NULL,
                                                                                            spa_sync_props,
                                                                error = dsl_sync_task_do     spa, nvp, 6));
                                                                (spa_get_dsl(spa), NULL,
                                                                  spa_sync_version,
                                                                   spa, &ver, 6);


                                                                             error ?

                                                                           Yes

                                                                      return (error);                 End
```

```
void spa_prop_clear_
bootfs(spa_t *spa,
uint64_t dsobj,
dmu_tx_t *tx)
```

* If the bootfs property
value is dsobj, clear it.

```
spa->spa_bootfs == dsobj
&& spa->spa_pool_
props_object != 0 ?
```

No

Yes

```
VERIFY(zap_remove(spa->spa_meta_objset,
spa->spa_pool_props_object,
zpool_prop_to_name(ZPOOL_PROP_BOOTFS), tx) == 0);
spa->spa_bootfs = 0;
```

End

```
static int
spa_change_guid_check
(void *arg1, void *arg2,
dmu_tx_t *tx)
```

ARGSUSED

```
spa_t *spa = arg1;
uint64_t *newguid = arg2;
vdev_t *rvd = spa->spa_root_vdev;
uint64_t vdev_state;
spa_config_enter(spa, SCL_STATE, FTAG, RW_READER);
vdev_state = rvd->vdev_state;
spa_config_exit(spa, SCL_STATE, FTAG);
```

vdev_state !=
VDEV_STATE_HEALTHY ?

No

Yes

```
ASSERT3U(spa_guid(spa), !=, *newguid);
return (0);
```

```
return (ENXIO);
```

End

```
static void
spa_change_guid_sync(void
*arg1, void *arg2,
dmu_tx_t *tx)
```

```
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);
```

End

```
int spa_change_guid(spa_t
*spa)
```

* 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);
```

error == 0 ?

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);
```

End

```
          static int
    spa_error_entry_compare
        (const void *a,
         const void *b)
```

* ========================================================================
* 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));
```

ret < 0 ?

Yes → return (-1);

No → ret > 0 ?

Yes → return (1);

No → return (0);

End

```
void spa_get_errlists
(spa_t *spa, avl_tree_t
*last, avl_tree_t *scrub)
```

* 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));
```

End

```
static taskq_t *
spa_taskq_create(spa_t
*spa, const char *name,
enum zti_modes
mode, uint_t value)
```

```
uint_t flags = 0;
boolean_t batch = B_FALSE;
```

switch (mode)

zti_mode_null ?

no taskq needed

Yes → return (NULL);

No ↓

zti_mode_fixed ?

Yes →
```
ASSERT3U(value, >=, 1);
value = MAX(value, 1);
```

No ↓

zti_mode_batch ?

Yes →
```
batch = B_TRUE;
flags |= TASKQ_THREADS_CPU_PCT;
value = zio_taskq_batch_pct;
```

No ↓

zti_mode_online_percent ?

Yes →
```
flags |=
TASKQ_THREADS_CPU_PCT;
```

default →
```
panic("unrecognized mode
for %s taskq (%u:%u) in
" "spa_activate()",
name, mode, value);
```

```
zio_taskq_sysdc &&
spa->spa_proc != &p0 ?
```

No →
```
return (taskq_create_proc
(name, value,
maxclsyspri, 50,
INT_MAX,
spa->spa_proc, flags));
```

Yes →

batch ?

No →
```
return (taskq_create_
sysdc(name, value, 50,
INT_MAX, spa->spa_proc,
zio_taskq_basedc,
flags));
```

Yes → `flags |= TASKQ_DC_BATCH;`

End

```
                    static void
                    spa_create_zio_
                    taskqs(spa_t *spa)

                         int t = 0

                      t < ZIO_TYPES ?
                    No │      │ Yes
                  ┌────┘      └────┐
                 End          int q = 0

                         q < ZIO_TASKQ_TYPES ?
              Yes │                    │ No
          ┌───────┘                    └──────┐
                                              t++

const zio_taskq_info_t *ztip = &zio_taskqs[t][q];
enum zti_modes mode = ztip->zti_mode;
uint_t value = ztip->zti_value;
char name[32];
(void) snprintf(name, sizeof (name), "%s_%s", zio_type_name[t], zio_taskq_types[q]);
spa->spa_zio_taskq[t][q] = spa_taskq_create(spa, name, mode, value);

                                     q++
```

```
static void
spa_thread(void *arg)
```

```
callb_cpr_t cprinfo;
spa_t *spa = arg;
user_t *pu = PTOU(curproc);
CALLB_CPR_INIT(&cprinfo, &spa->spa_proc_lock, callb_generic_cpr, spa->spa_name);
ASSERT(curproc != &p0);
(void) snprintf(pu->u_psargs, sizeof (pu->u_psargs), "zpool-%s", spa->spa_name);
(void) strlcpy(pu->u_comm, pu->u_psargs, sizeof (pu->u_comm));
```

bind this thread to the
requested psrset

**zio_taskq_psrset_bind != PS_NONE ?**

No

Yes

```
pool_lock();
mutex_enter(&cpu_lock);
mutex_enter(&pidlock);
mutex_enter(&curproc->p_lock);
```

**cpupart_bind_thread (curthread, zio_taskq_psrset_bind, 0, NULL, NULL) == 0 ?**

Yes

No

```
curthread->t_bind_pset =
zio_taskq_psrset_bind;
```

```
cmn_err(CE_WARN,
"Couldn't bind process
for zfs pool \"%s\" to "
"pset %d\n",
spa->spa_name,
zio_taskq_psrset_bind);
```

```
mutex_exit(&curproc->p_lock);
mutex_exit(&pidlock);
mutex_exit(&cpu_lock);
pool_unlock();
```

**zio_taskq_sysdc ?**

No

Yes

```
sysdc_thread_enter
(curthread, 100, 0);
```

```
spa->spa_proc = curproc;
spa->spa_did = curthread->t_did;
spa_create_zio_taskqs(spa);
mutex_enter(&spa->spa_proc_lock);
ASSERT(spa->spa_proc_state == SPA_PROC_CREATED);
spa->spa_proc_state = SPA_PROC_ACTIVE;
cv_broadcast(&spa->spa_proc_cv);
CALLB_CPR_SAFE_BEGIN(&cprinfo);
```

**spa->spa_proc_state == SPA_PROC_ACTIVE ?**

No

Yes

```
CALLB_CPR_SAFE_END(&cprinfo, &spa->spa_proc_lock);
ASSERT(spa->spa_proc_state == SPA_PROC_DEACTIVATE);
spa->spa_proc_state = SPA_PROC_GONE;
spa->spa_proc = &p0;
cv_broadcast(&spa->spa_proc_cv);
```

```
cv_wait(&spa->spa_proc_
cv, &spa->spa_proc_lock);
```

drops spa_proc_lock

```
CALLB_CPR_EXIT(&cprinfo);
mutex_enter(&curproc->p_lock);
lwp_exit();
```

End

```
static void
spa_activate(spa_t
*spa, int mode)
```

* Activate an
uninitialized pool.

```
ASSERT(spa->spa_state == POOL_STATE_UNINITIALIZED);
spa->spa_state = POOL_STATE_ACTIVE;
spa->spa_mode = mode;
spa->spa_normal_class = metaslab_class_create(spa, zfs_metaslab_ops);
spa->spa_log_class = metaslab_class_create(spa, zfs_metaslab_ops);
```

Try to create a
covering process

```
mutex_enter(&spa->spa_proc_lock);
ASSERT(spa->spa_proc_state == SPA_PROC_NONE);
ASSERT(spa->spa_proc == &p0);
spa->spa_did = 0;
```

Only create a process if
we're going to be
around a while.

spa_create_process &&
strcmp(spa->spa_name,
TRYIMPORT_NAME) != 0 ?

Yes

newproc(spa_thread,
(caddr_t)spa, syscid,
maxclsyspri,
NULL, 0) == 0 ?

Yes

No

No

```
spa->spa_proc_state =
SPA_PROC_CREATED;
```

spa->spa_proc_state ==
SPA_PROC_CREATED ?

#ifdef _KERNEL ?

No

Yes

Yes

```
ASSERT(spa->spa_proc_state == SPA_PROC_ACTIVE);
ASSERT(spa->spa_proc != &p0);
ASSERT(spa->spa_did != 0);
```

```
cv_wait(&spa->spa_proc_
cv, &spa->spa_proc_lock);
```

```
cmn_err(CE_WARN,
"Couldn't create process
for zfs pool \"%s\"\n",
spa->spa_name);
```

```
mutex_exit(&spa
->spa_proc_lock);
```

If we didn't create a
process, we need to
create our taskqs.

spa->spa_proc == &p0 ?

No

Yes

```
spa_create_zio_
taskqs(spa);
```

```
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

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

* Opposite of
spa_activate().

```
ASSERT(spa->spa_sync_on == B_FALSE);
ASSERT(spa->spa_dsl_pool == NULL);
ASSERT(spa->spa_root_vdev == NULL);
ASSERT(spa->spa_async_zio_root == NULL);
ASSERT(spa->spa_state != POOL_STATE_UNINITIALIZED);
txg_list_destroy(&spa->spa_vdev_txg_list);
list_destroy(&spa->spa_config_dirty_list);
list_destroy(&spa->spa_state_dirty_list);
int t = 0
```

t < ZIO_TYPES ?

**No**

```
metaslab_class_destroy(spa->spa_normal_class);
spa->spa_normal_class = NULL;
metaslab_class_destroy(spa->spa_log_class);
spa->spa_log_class = NULL;
```

* If this was part of an import or
the open otherwise failed, we may
* still have errors left in the queues.  Empty them just in case.

**Yes**

int q = 0

```
spa_errlog_drain(spa);
avl_destroy(&spa->spa_errlist_scrub);
avl_destroy(&spa->spa_errlist_last);
spa->spa_state = POOL_STATE_UNINITIALIZED;
mutex_enter(&spa->spa_proc_lock);
```

q < ZIO_TASKQ_TYPES ?

**No**

**Yes**

spa->spa_proc_state !=
SPA_PROC_NONE ?

**No**

**Yes**

t++

spa->spa_zio_taskq[t][q]
!= NULL ?

**No**

**Yes**

```
ASSERT(spa->spa_proc_state == SPA_PROC_ACTIVE);
spa->spa_proc_state = SPA_PROC_DEACTIVATE;
cv_broadcast(&spa->spa_proc_cv);
```

```
taskq_destroy(spa->spa_
zio_taskq[t][q]);
```

spa->spa_proc_state ==
SPA_PROC_DEACTIVATE ?

**No**

**Yes**

```
spa->spa_zio_taskq[t][q]
= NULL;
```

```
ASSERT(spa->spa_proc_state == SPA_PROC_GONE);
spa->spa_proc_state = SPA_PROC_NONE;
```

```
ASSERT(spa->spa_proc != &p0);
cv_wait(&spa->spa_proc_cv, &spa->spa_proc_lock);
```

q++

```
ASSERT(spa->spa_proc == &p0);
mutex_exit(&spa->spa_proc_lock);
```

* We want to make sure spa_thread() has actually exited the ZFS
* module, so that the module can't be unloaded out from underneath
* it.

spa->spa_did != 0 ?

**No**

**Yes**

```
thread_join(spa->spa_did);
spa->spa_did = 0;
```

End

```
                              static int
                        spa_config_parse(spa_t
                         *spa, vdev_t **vdp,
                         nvlist_t *nv, vdev_t
                           *parent, uint_t
                            id, int atype)
```

* 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.

```
nvlist_t **child;
uint_t children;
int error;
```

(error = vdev_alloc(spa,
vdp, nv, parent,
id, atype)) != 0 ?

**No** → (*vdp)->vdev_ops-
>vdev_op_leaf ?

**Yes** → return (error);

(*vdp)->vdev_ops->vdev_op_leaf ?

**Yes** → return (0);

**No** →
```
error = nvlist_lookup_
nvlist_array(nv,
ZPOOL_CONFIG_CHILDREN,
&child, &children);
```

error == ENOENT ?

**Yes** → return (0);

**No** → error ?

**Yes** →
```
vdev_free(*vdp);
*vdp = NULL;
return (EINVAL);
```

**No** → `int c = 0`

c < children ?

**No** →
```
ASSERT(*vdp != NULL);
return (0);
```

**Yes** → `vdev_t *vd;`

(error =
spa_config_parse(spa,
&vd, child[c], *vdp, c,
atype)) != 0 ?

**Yes** →
```
vdev_free(*vdp);
*vdp = NULL;
return (error);
```

**No** → `c++`

End

```
                    static void
                 spa_unload(spa_t *spa)

                 * Opposite of spa_load().

                 int i;
                 ASSERT(MUTEX_HELD(&spa_namespace_lock));        * Stop async tasks.

                 spa_async_suspend(spa);        * Stop syncing.

                 spa->spa_sync_on ?

        Yes                          No

 txg_sync_stop(spa->spa_dsl_pool);              * Wait for any
 spa->spa_sync_on = B_FALSE;                    outstanding async
                                                I/O to complete.

                 spa->spa_async_zio_root
                 != NULL ?

                                        Yes

                 No      (void) zio_wait(spa->spa_async_zio_root);
                         spa->spa_async_zio_root = NULL;

                 bpobj_close(&spa->spa_
                 deferred_bpobj);       * Close the dsl pool.

                 spa->spa_dsl_pool ?

                                        Yes

                 No      dsl_pool_close(spa->spa_dsl_pool);
                         spa->spa_dsl_pool = NULL;
                         spa->spa_meta_objset = NULL;

                 ddt_unload(spa);
                 spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);    * Drop and purge
                                                                     level 2 cache

                 spa_l2cache_drop(spa);        * Close all vdevs.

                 spa->spa_root_vdev ?

                                        Yes

                 No      vdev_free(spa-
                         >spa_root_vdev);

                 ASSERT(spa->spa_root_vdev == NULL);
                 i = 0

                         i < spa->spa_spares
                         .sav_count ?

                 No                      Yes

 spa->spa_spares.sav_vdevs       vdev_free(spa->spa_spares
 ?                               .sav_vdevs[i]);

        Yes

 No      kmem_free(spa->spa_spares.sav_vdevs,              i++
         spa->spa_spares.sav_count * sizeof (void *));
         spa->spa_spares.sav_vdevs = NULL;

                 spa->spa_spares
                 .sav_config ?

                                        Yes

                 No      nvlist_free(spa->spa_spares.sav_config);
                         spa->spa_spares.sav_config = NULL;

                 spa->spa_spares.sav_count = 0;
                 i = 0

                         i < spa->spa_
                         l2cache.sav_count ?

                                        Yes

 No      vdev_clear_stats(spa->spa_l2cache.sav_vdevs[i]);
 spa->spa_l2cache        vdev_free(spa->spa_l2cache.sav_vdevs[i]);
 .sav_vdevs ?

        Yes

 No      kmem_free(spa->spa_l2cache.sav_vdevs,              i++
         spa->spa_l2cache.sav_count * sizeof (void *));
         spa->spa_l2cache.sav_vdevs = NULL;

                 spa->spa_l2cache
                 .sav_config ?

                                        Yes

                 No      nvlist_free(spa->spa_l2cache.sav_config);
                         spa->spa_l2cache.sav_config = NULL;

                 spa->spa_l2cache.sav_count = 0;
                 spa->spa_async_suspended = 0;

                 spa->spa_comment
                 != NULL ?

                                        Yes

                 No      spa_strfree(spa->spa_comment);
                         spa->spa_comment = NULL;

                 spa_config_exit(spa,
                 SCL_ALL, FTAG);

                         End
```

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

* 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.

```
nvlist_t **spares;
uint_t nspares;
int i;
vdev_t *vd, *tvd;
ASSERT(spa_config_held(spa,
SCL_ALL, RW_WRITER) == SCL_ALL);
```

* First, close and free
any existing spare vdevs.

```
i = 0
```

```
i < spa->spa_spares
.sav_count ?
```
No / Yes

```
spa->spa_spares.sav_vdevs
?
```
Yes → 
```
kmem_free(spa->spa_spares
.sav_vdevs,
spa->spa_spares.sav_count
* sizeof (void *));
```
No →

```
vd = spa->spa_spares
.sav_vdevs[i];
```

Undo the call to
spa_activate() below

```
(tvd = spa_lookup_by_guid
(spa, vd->vdev_guid,
B_FALSE)) != NULL &&
tvd->vdev_isspare ?
```
No / Yes →
```
spa_spare_remove(tvd);
```

```
vdev_close(vd);
vdev_free(vd);
```

```
i++
```

```
spa->spa_spares.sav_
config == NULL ?
```
Yes →
```
nspares = 0;
```
No →
```
VERIFY(nvlist_lookup_
nvlist_array(spa->spa_
spares.sav_config,
ZPOOL_CONFIG_SPARES,
&spares, &nspares) == 0);
```

```
spa->spa_spares.sav_count = (int)nspares;
spa->spa_spares.sav_vdevs = NULL;
```

* 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.

```
nspares == 0 ?
```
Yes / No →

```
spa->spa_spares.sav_vdevs = kmem_alloc(nspares
* sizeof (void *), KM_SLEEP);
i = 0
```

```
i < spa->spa_spares
.sav_count ?
```
No / Yes →
```
VERIFY(spa_config_parse(spa, &vd, spares[i], NULL, 0, VDEV_ALLOC_SPARE) == 0);
ASSERT(vd != NULL);
spa->spa_spares.sav_vdevs[i] = vd;
```

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

```
VERIFY(nvlist_remove(spa->spa_spares.sav_config, ZPOOL_CONFIG_SPARES, DATA_TYPE_NVLIST_ARRAY) == 0);
spares = kmem_alloc(spa->spa_spares.sav_count * sizeof (void *), KM_SLEEP);
i = 0
```

```
(tvd = spa_lookup_by_guid
(spa, vd->vdev_guid,
B_FALSE)) != NULL ?
```
Yes →
```
!tvd->vdev_isspare ?
```
Yes →
```
spa_spare_add(tvd);
```
No →

* We only mark the spare active if we were successfully
* able to load the vdev. Otherwise, importing a pool
* with a bad active spare would result in strange
* behavior, because multiple pool would think the spare
* is actively in use.
*
* There is a vulnerability here to an equally bizarre
* circumstance, where a dead active spare is later
* brought back to life (onlined or otherwise). Given
* the rarity of this scenario, and the extra complexity
* it adds, we ignore the possibility.

```
i < spa->spa_spares
.sav_count ?
```
No / Yes →
```
spares[i] =
vdev_config_generate(spa,
spa->spa_spares.sav_vdevs
[i], B_TRUE,
VDEV_CONFIG_SPARE);
```

```
VERIFY(nvlist_add_nvlist_array(spa->spa_spares.sav_config,
ZPOOL_CONFIG_SPARES, spares, spa->spa_spares.sav_count) == 0);
i = 0
```

```
i++
```

```
!vdev_is_dead(tvd) ?
```
No / Yes →
```
spa_spare_activate(tvd);
```

```
i < spa->spa_spares
.sav_count ?
```
No / Yes →
```
nvlist_free(spares[i]);
```

```
kmem_free(spares,
spa->spa_spares.sav_count
* sizeof (void *));
```

```
vd->vdev_top = vd;
vd->vdev_aux = &spa->spa_spares;
```

```
End
```

```
i++
```

```
vdev_open(vd) != 0 ?
```
Yes / No →
```
vdev_validate_aux(vd)
== 0 ?
```
Yes →
```
spa_spare_add(vd);
```
No →

```
i++
```

```
                                    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 ?

— Yes → 
```
VERIFY(nvlist_lookup_nvlist_array(sav->sav_config,
ZPOOL_CONFIG_L2CACHE, &l2cache, &nl2cache) == 0);
newvdevs = kmem_alloc(nl2cache * sizeof (void *), KM_SLEEP);
```

— No → `nl2cache = 0;`

```
oldvdevs = sav->sav_vdevs;
oldnvdevs = sav->sav_count;
sav->sav_vdevs = NULL;
sav->sav_count = 0;
```
* Process new
nvlist of vdevs.

`i = 0`

`i < nl2cache ?`

— No → (* Purge vdevs that were dropped) `i = 0`

`i < oldnvdevs ?`

— Yes →
```
uint64_t pool;
vd = oldvdevs[i];
```
`vd != NULL ?`

— Yes → `ASSERT(vd->vdev_isl2cache);`
```
spa_l2cache_exists(vd-
>vdev_guid, &pool) &&
pool != 0ULL &&
l2arc_vdev_present(vd) ?
```
— Yes → `l2arc_remove_vdev(vd);`
— No →
```
vdev_clear_stats(vd);
vdev_free(vd);
```
— No → `i++`

— No → `oldvdevs ?`
— Yes →
```
kmem_free(oldvdevs,
oldnvdevs *
sizeof (void *));
```

`sav->sav_config == NULL ?`
— Yes (down)
— No →
```
sav->sav_vdevs = newvdevs;
sav->sav_count = (int)nl2cache;
```
* Recompute the stashed list of l2cache devices, with status
* information this time.

```
VERIFY(nvlist_remove(sav->sav_config, ZPOOL_CONFIG_L2CACHE, DATA_TYPE_NVLIST_ARRAY) == 0);
l2cache = kmem_alloc(sav->sav_count * sizeof (void *), KM_SLEEP);
i = 0
```

`i < sav->sav_count ?`
— Yes →
```
l2cache[i] =
vdev_config_generate(spa,
sav->sav_vdevs[i],
B_TRUE, VDEV_
CONFIG_L2CACHE);
```
→ `i++`
— No →
```
VERIFY(nvlist_add_nvlist_
array(sav->sav_config,
ZPOOL_CONFIG_L2CACHE,
l2cache,
sav->sav_count) == 0);
```

`out: i = 0`

`i < sav->sav_count ?`
— Yes → `nvlist_free(l2cache[i]);` → `i++`
— No → `sav->sav_count ?`
— Yes →
```
kmem_free(l2cache,
sav->sav_count *
sizeof (void *));
```
— No → `End`

Right branch:
`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 ?
```
— No → `j++`
— Yes → (* Retain previous vdev for add/remove ops.)
```
newvdevs[i] = vd;
oldvdevs[j] = NULL;
```
— No →

`newvdevs[i] == NULL ?`
— Yes → (* Create new vdev)
```
VERIFY(spa_config_parse(spa, &vd, l2cache[i],
NULL, 0, VDEV_ALLOC_L2CACHE) == 0);
ASSERT(vd != NULL);
newvdevs[i] = vd;
```
* Commit this vdev as an l2cache device,
* even if it fails to open.
```
spa_l2cache_add(vd);
vd->vdev_top = vd;
vd->vdev_aux = sav;
spa_l2cache_activate(vd);
```
`vdev_open(vd) != 0 ?`
— No →
```
(void) vdev_
validate_aux(vd);
```
`!vdev_is_dead(vd) ?`
— Yes → `l2arc_add_vdev(spa, vd);`
— No →
— Yes →

`i++`
— No →
```

```
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

```
static void
spa_check_removed(vdev_t
*vd)
```

* 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`

c < vd->vdev_children ?

No

vd->vdev_ops->vdev_op_
leaf && vdev_is_dead(vd)
?

Yes

`spa_check_removed(vd
->vdev_child[c]);`

c++

Yes

No

```
zfs_post_autoreplace(vd->vdev_spa, vd);
spa_event_notify(vd->vdev_spa, vd, ESC_ZFS_VDEV_CHECK);
```

End

```
static boolean_t
spa_config_valid(spa_t
*spa, nvlist_t *config)
```

* Validate the current
config against
the MOS config

```
vdev_t *mrvd, *rvd = spa->spa_root_vdev;
nvlist_t *nv;
VERIFY(nvlist_lookup_nvlist(config, ZPOOL_CONFIG_VDEV_TREE, &nv) == 0);
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
VERIFY(spa_config_parse(spa, &mrvd, nv, NULL, 0, VDEV_ALLOC_LOAD) == 0);
ASSERT3U(rvd->vdev_children, ==, mrvd->vdev_children);
```

* 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.

```
!(spa->spa_import_flags
& ZFS_IMPORT_MISSING_LOG)
?
```

Yes → 

No → 

```
nvlist_t **child, *nv;
uint64_t idx = 0;
child = kmem_alloc(rvd->vdev_children * sizeof (nvlist_t **), KM_SLEEP);
VERIFY(nvlist_alloc(&nv, NV_UNIQUE_NAME, KM_SLEEP) == 0);
int c = 0
```

```
c < rvd->vdev_children ?
```

Yes / No

```
idx ?
```

```
vdev_t *tvd = rvd->vdev_child[c];
vdev_t *mtvd = mrvd->vdev_child[c];
```

```
tvd->vdev_ops ==
&vdev_missing_ops &&
mtvd->vdev_ops !=
&vdev_missing_ops &&
mtvd->vdev_islog ?
```

No / Yes

```
VERIFY(nvlist_add_nvlist_array(nv, ZPOOL_CONFIG_CHILDREN, child, idx) == 0);
VERIFY(nvlist_add_nvlist(spa->spa_load_info, ZPOOL_CONFIG_MISSING_DEVICES, nv) == 0);
int i = 0
```

```
child[idx++] =
vdev_config_generate(spa,
mtvd, B_FALSE, 0);
```

```
i < idx ?
```

Yes / No

c++

```
nvlist_free(child[i]);
```

* Compare the root vdev tree with the information we have
* from the MOS config (mrvd). Check each top-level vdev
* with the corresponding MOS config top-level (mtvd).

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

i++

```
int c = 0
```

```
c < rvd->vdev_children ?
```

No / Yes

* Ensure we were able to
validate the config.

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

* 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.

```
vdev_t *tvd = rvd->vdev_child[c];
vdev_t *mtvd = mrvd->vdev_child[c];
```

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

```
tvd->vdev_ops ==
&vdev_missing_ops &&
mtvd->vdev_ops !=
&vdev_missing_ops ?
```

Yes / No

<End>

```
!(spa->spa_import_flags
& ZFS_IMPORT_MISSING_LOG)
?
```

* Device
specific actions.

```
mtvd->vdev_islog ?
```

No / Yes

```
mtvd->vdev_islog ?
```

Yes

* Swap the missing vdev with the data we were
* able to obtain from the MOS config.

```
spa_set_log_state(spa,
SPA_LOG_CLEAR);
```

* 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.

* XXX - once we have 'readonly' pool
* support we should be able to handle
* missing data devices by transitioning
* the pool to readonly.

```
vdev_remove_child(rvd, tvd);
vdev_remove_child(mrvd, mtvd);
vdev_add_child(rvd, mtvd);
vdev_add_child(mrvd, tvd);
spa_config_exit(spa, SCL_ALL, FTAG);
vdev_load(mtvd);
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
vdev_reopen(rvd);
```

No

```
vdev_load_log_state(tvd, mtvd);
vdev_reopen(tvd);
```

c++

Yes
```

```
static int
spa_check_logs(spa_t
*spa)
```

* Check for missing
log devices

switch (spa-
>spa_log_state)

SPA_LOG_MISSING ?

No → SPA_LOG_UNKNOWN ?

Yes → dmu_objset_find(spa->spa_
name, zil_check_log_
chain, NULL,
DS_FIND_CHILDREN) ?

need to recheck in case
slog has been restored

default → return (0);

No → return (0);

Yes →
```
spa_set_log_state(spa, SPA_LOG_MISSING);
return (1);
```

End

```
                          static boolean_t
                         spa_passivate_log(spa_t
                                *spa)

                    vdev_t *rvd = spa->spa_root_vdev;
                    boolean_t slog_found = B_FALSE;
                    ASSERT(spa_config_held(spa, SCL_ALLOC, RW_WRITER));

                              !spa_has_slogs(spa) ?
                        No                         Yes

                   int c = 0

                                                  return (B_FALSE);

                    c < rvd->vdev_children ?

              Yes                    No

      vdev_t *tvd = rvd->vdev_child[c];
      metaslab_group_t *mg = tvd->vdev_mg;          return (slog_found);

              tvd->vdev_islog ?

         Yes                No                              End

  metaslab_group_passivate(mg);
  slog_found = B_TRUE;

                                   c++
```

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

```
vdev_t *rvd = spa->spa_root_vdev;
ASSERT(spa_config_held(spa, SCL_ALLOC, RW_WRITER));
int c = 0
```

c < rvd->vdev_children ?

No

Yes

End

```
vdev_t *tvd = rvd->vdev_child[c];
metaslab_group_t *mg = tvd->vdev_mg;
```

tvd->vdev_islog ?

Yes

No

```
metaslab_group_
activate(mg);
```

c++

int spa_offline_log(spa_t
*spa)

int error = 0;

(error =
dmu_objset_find(spa_name
(spa), zil_vdev_offline,
NULL, DS_FIND_CHILDREN))
== 0 ?

* We successfully offlined the log device, sync out the
* current txg so that the "stubby" block can be removed
* by zil_sync().

Yes

No

txg_wait_synced(spa->spa_
dsl_pool, 0);

return (error);

End

```
              static void
         spa_aux_check_removed
         (spa_aux_vdev_t *sav)

                    │
                    ▼
              ┌──────────┐
              │ int i = 0 │
              └──────────┘
                    │
                    ▼
            ◇ i < sav->sav_count ? ◇
           No ╱              ╲ Yes
             ▼                ▼
         ◇ End ◇      ┌──────────────────┐
                      │ spa_check_removed(sav │
                      │   ->sav_vdevs[i]);     │
                      └──────────────────┘
                                │
                                ▼
                            ┌─────┐
                            │ i++ │
                            └─────┘
```

```
                          void spa_claim_
                          notify(zio_t *zio)

                      spa_t *spa = zio->io_spa;

                         zio->io_error ?          any mutex will do

                                          No

                  Yes                        mutex_enter(&spa
                                             ->spa_props_lock);

                                        spa->spa_claim_max_txg <
                                          zio->io_bp->blk_birth ?

                                                            Yes

                                      No
                                               spa->spa_claim_max_txg =
                                                 zio->io_bp->blk_birth;

                                      mutex_exit(&spa-
                                      >spa_props_lock);

                                        End
```

```
typedef struct
spa_load_error
```

```
uint64_t sle_meta_count;
uint64_t sle_data_count;
```

End

```
static void
spa_load_verify_done
(zio_t *zio)
```

```
blkptr_t *bp = zio->io_bp;
spa_load_error_t *sle = zio->io_private;
dmu_object_type_t type = BP_GET_TYPE(bp);
int error = zio->io_error;
```

error ?

No

Yes

(BP_GET_LEVEL(bp) != 0
|| DMU_OT_IS_METADATA
(type)) && type !=
DMU_OT_INTENT_LOG ?

Yes

No

```
atomic_add_64(&sle->sle_
meta_count, 1);
```

```
atomic_add_64(&sle->sle_
data_count, 1);
```

```
zio_data_buf_free(zio-
>io_data, zio->io_size);
```

End

```
static int
spa_load_verify_cb(spa_t
*spa, zilog_t *zilog,
const blkptr_t *bp,
arc_buf_t *pbuf, const
zbookmark_t *zb, const
dnode_phys_t
*dnp, void *arg)
```

ARGSUSED

bp != NULL ?

Yes

No

```
zio_t *rio = arg;
size_t size = BP_GET_PSIZE(bp);
void *data = zio_data_buf_alloc(size);
zio_nowait(zio_read(rio, spa, bp, data, size, spa_load_verify_done, rio->io_private,
ZIO_PRIORITY_SCRUB, ZIO_FLAG_SPECULATIVE | ZIO_FLAG_CANFAIL | ZIO_FLAG_SCRUB | ZIO_FLAG_RAW, zb));
```

return (0);

End

```
static int
spa_load_verify(spa_t
*spa)
```

```
zio_t *rio;
spa_load_error_t sle = { 0 };
zpool_rewind_policy_t policy;
boolean_t verify_ok = B_FALSE;
int error;
zpool_get_rewind_policy(spa->spa_config, &policy);
```

policy.zrp_request &
ZPOOL_NEVER_REWIND ?

No

Yes

```
rio = zio_root(spa, NULL, &sle, ZIO_FLAG_CANFAIL | ZIO_FLAG_SPECULATIVE);
error = traverse_pool(spa, spa->spa_verify_min_txg,
TRAVERSE_PRE | TRAVERSE_PREFETCH, spa_load_verify_cb, rio);
(void) zio_wait(rio);
spa->spa_load_meta_errors = sle.sle_meta_count;
spa->spa_load_data_errors = sle.sle_data_count;
```

return (0);

!error &&
sle.sle_meta_count <=
policy.zrp_maxmeta &&
sle.sle_data_count <=
policy.zrp_maxdata ?

Yes

No

```
int64_t loss = 0;
verify_ok = B_TRUE;
spa->spa_load_txg = spa->spa_uberblock.ub_txg;
spa->spa_load_txg_ts = spa->spa_uberblock.ub_timestamp;
loss = spa->spa_last_ubsync_txg_ts - spa->spa_load_txg_ts;
VERIFY(nvlist_add_uint64(spa->spa_load_info, ZPOOL_CONFIG_LOAD_TIME, spa->spa_load_txg_ts) == 0);
VERIFY(nvlist_add_int64(spa->spa_load_info, ZPOOL_CONFIG_REWIND_TIME, loss) == 0);
VERIFY(nvlist_add_uint64(spa->spa_load_info,
ZPOOL_CONFIG_LOAD_DATA_ERRORS, sle.sle_data_count) == 0);
```

```
spa->spa_load_max_txg =
spa->spa_
uberblock.ub_txg;
```

error ?

No

Yes

error != ENXIO
&& error != EIO ?

No

Yes

```
return (verify_ok
? 0 : EIO);
```

error = EIO;

return (error);

End

static void
spa_prop_find(spa_t
*spa, zpool_prop_t prop,
uint64_t *val)

* Find a value in the
pool props object.

(void) zap_lookup(spa-
>spa_meta_objset,
spa->spa_pool_props_
object, zpool_prop_to_
name(prop), sizeof
(uint64_t), 1, val);

End

```
static int
spa_dir_prop(spa_t *spa,
const char *name,
uint64_t *val)
```

* Find a value in the
pool directory object.

```
return (zap_lookup(spa-
>spa_meta_objset,
DMU_POOL_DIRECTORY_
OBJECT, name, sizeof
(uint64_t), 1, val));
```

End

```
static int
spa_vdev_err(vdev_t
*vdev, vdev_aux_t
aux, int err)
```

```
vdev_set_state(vdev, B_TRUE, VDEV_STATE_CANT_OPEN, aux);
return (err);
```

End

```
static void
spa_try_repair(spa_t
*spa, nvlist_t *config)
```

* 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 guids 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.

```
uint_t extracted;
uint64_t *glist;
uint_t i, gcount;
nvlist_t *nvl;
vdev_t **vd;
boolean_t attempt_reopen;
```

nvlist_lookup_nvlist
(config,
ZPOOL_CONFIG_SPLIT,
&nvl) != 0 ?

check that the config is complete

No

Yes

nvlist_lookup_uint64_
array(nvl,
ZPOOL_CONFIG_SPLIT_LIST,
&glist, &gcount) != 0 ?

Yes

No

```
vd = kmem_zalloc(gcount
* sizeof (vdev_t
*), KM_SLEEP);
```

attempt to online all the vdevs & validate

```
attempt_reopen = B_TRUE;
i = 0
```

i < gcount ?

vdev is hole

Yes

No

glist[i] == 0 ?

attempt_reopen ?

Yes

No

```
vd[i] = spa_lookup_by_
guid(spa, glist[i],
B_FALSE);
```

```
vdev_reopen(spa
->spa_root_vdev);
```

check each device to see what state it's in

vd[i] == NULL ?

attempt to re-online it

```
extracted = 0, i = 0
```

* Don't bother attempting to reopen the disks;
* just do the split.

Yes

No

```
vd[i]->vdev_offline
= B_FALSE;
```

i < gcount ?

```
attempt_reopen = B_FALSE;
```

Yes

Yes

```
vd[i] != NULL &&
vd[i]->vdev_stat.vs_aux
!= VDEV_AUX_SPLIT_POOL ?
```

No

No

i++

* If every disk has been moved to the new pool, or if we never
* even attempted to look at them, then we split them off for
* good.

Yes

++extracted;

```
!attempt_reopen ||
gcount == extracted ?
```

No

Yes

i++

No

i = 0

i < gcount ?

No

Yes

```
vdev_reopen(spa
->spa_root_vdev);
```

vd[i] != NULL ?

```
kmem_free(vd, gcount *
sizeof (vdev_t *));
```

Yes

No

vdev_split(vd[i]);

<End>

i++

```
static int
spa_load(spa_t *spa,
spa_load_state_t state,
spa_import_type_t type,
boolean_t mosconfig)
```

```
nvlist_t *config = spa->spa_config;
char *ereport = FM_EREPORT_ZFS_POOL;
char *comment;
int error;
uint64_t pool_guid;
nvlist_t *nvl;
```

nvlist_lookup_uint64
(config,
ZPOOL_CONFIG_POOL_GUID,
&pool_guid) ?

No → ASSERT(spa->spa_comment == NULL);

Yes → return (EINVAL);

nvlist_lookup_string
(config,
ZPOOL_CONFIG_COMMENT,
&comment) == 0 ?

Yes → spa->spa_comment = spa_strdup(comment);

No →

```
* Versioning wasn't explicitly added
to the label until later, so if
* it's not present treat it as the initial version.
```

nvlist_lookup_uint64
(config,
ZPOOL_CONFIG_VERSION,
&spa->spa_ubsync.ub_
version) != 0 ?

Yes → spa->spa_ubsync.ub_
version =
SPA_VERSION_INITIAL;

No →

```
(void) nvlist_lookup_
uint64(config,
ZPOOL_CONFIG_POOL_TXG,
&spa->spa_config_txg);
```

(state ==
SPA_LOAD_IMPORT || state
== SPA_LOAD_TRYIMPORT)
&& spa_guid_exists
(pool_guid, 0) ?

Yes → error = EEXIST;

No → spa->spa_config_guid = pool_guid;

nvlist_lookup_nvlist
(config,
ZPOOL_CONFIG_SPLIT,
&nvl) == 0 ?

Yes → VERIFY(nvlist_dup(nvl,
&spa->spa_config_
splitting,
KM_SLEEP) == 0);

No →

```
nvlist_free(spa->spa_load_info);
spa->spa_load_info = fnvlist_alloc();
gethrestime(&spa->spa_loaded_ts);
error = spa_load_impl(spa, pool_guid, config, state, type, mosconfig, &ereport);
```

```
spa->spa_minref =
refcount_count(&spa
->spa_refcount);
```

error ?

Yes → error != EEXIST ?

Yes → spa->spa_loaded_ts.tv_sec = 0;
spa->spa_loaded_ts.tv_nsec = 0;

No →

error != EBADF ?

Yes → zfs_ereport_post(ereport,
spa, NULL, NULL, 0, 0);

No →

```
spa->spa_load_state = error ? SPA_LOAD_ERROR : SPA_LOAD_NONE;
spa->spa_ena = 0;
return (error);
```

End

```
static int
spa_load_retry(spa_t
*spa, spa_load_state_t
state, int mosconfig)
```

```
int mode = spa->spa_mode;
spa_unload(spa);
spa_deactivate(spa);
spa->spa_load_max_txg--;
spa_activate(spa, mode);
spa_async_suspend(spa);
return (spa_load(spa, state,
SPA_IMPORT_EXISTING, mosconfig));
```

End

```
static int
spa_load_best(spa_t
*spa, spa_load_state_t
state, int mosconfig,
uint64_t max_request,
int rewind_flags)
```

* 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().

```
nvlist_t *loadinfo = NULL;
nvlist_t *config = NULL;
int load_error, rewind_error;
uint64_t safe_rewind_txg;
uint64_t min_txg;
```

spa->spa_load_txg &&
state ==
SPA_LOAD_RECOVER ?

**Yes** → 
```
spa->spa_load_max_txg = spa->spa_load_txg;
spa_set_log_state(spa, SPA_LOG_CLEAR);
```

**No** →
```
spa->spa_load_max_txg
= max_request;
```

```
load_error =
rewind_error =
spa_load(spa, state,
SPA_IMPORT_EXISTING,
mosconfig);
```

load_error == 0 ?

**Yes** → `return (0);`

**No** → spa->spa_root_vdev
!= NULL ?

**Yes** →
```
config =
spa_config_generate(spa,
NULL, -1ULL, B_TRUE);
```

**No** →
```
spa->spa_last_ubsync_txg = spa->spa_uberblock.ub_txg;
spa->spa_last_ubsync_txg_ts = spa->spa_uberblock.ub_timestamp;
```

rewind_flags &
ZPOOL_NEVER_REWIND ?

**Yes** →
```
nvlist_free(config);
return (load_error);
```

* 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.

**No** → state ==
SPA_LOAD_RECOVER ?

Price of rolling back is
discarding txgs,
including log

**Yes** →
```
spa_set_log_state(spa,
SPA_LOG_CLEAR);
```

**No** →
```
loadinfo = spa->spa_load_info;
spa->spa_load_info = fnvlist_alloc();
```

```
spa->spa_load_max_txg = spa->spa_last_ubsync_txg;
safe_rewind_txg = spa->spa_last_ubsync_txg - TXG_DEFER_SIZE;
min_txg = (rewind_flags & ZPOOL_EXTREME_REWIND) ? TXG_INITIAL : safe_rewind_txg;
```

* Continue as long as we're finding errors, we're still within
* the acceptable rewind range, and we're still finding uberblocks

```
rewind_error &&
spa->spa_uberblock.ub_txg
>= min_txg &&
spa->spa_uberblock.ub_txg
<= spa->spa_load_max_txg
?
```

**No** →
```
spa->spa_extreme_rewind = B_FALSE;
spa->spa_load_max_txg = UINT64_MAX;
```

**Yes** →
spa->spa_load_max_txg <
safe_rewind_txg ?

**Yes** →
```
spa->spa_extreme_rewind
= B_TRUE;
```

**No** →
```
rewind_error =
spa_load_retry(spa,
state, mosconfig);
```

```
config && (rewind_error
|| state !=
SPA_LOAD_RECOVER) ?
```

**Yes** →
```
spa_config_set(spa,
config);
```

Store the rewind info as
part of the
initial load info

**No** → state ==
SPA_LOAD_RECOVER ?

Restore the
initial load info

**No** →
```
fnvlist_add_nvlist
(loadinfo,
ZPOOL_CONFIG_REWIND_INFO,
spa->spa_load_info);
```

**Yes** →
```
ASSERT3P(loadinfo, ==, NULL);
return (rewind_error);
```

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

End

```
                    static int
              spa_open_common(const
              char *pool, spa_t
              **spapp, void *tag,
              nvlist_t *nvpolicy,
              nvlist_t **config)
```

* 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.

```
spa_t *spa;
spa_load_state_t state
= SPA_LOAD_OPEN;
int error;
int locked = B_FALSE;
*spapp = NULL;
```

* 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_owner(&spa_
namespace_lock)
!= curthread ?
```
Yes
```
mutex_enter(&spa_namespace_lock);
locked = B_TRUE;
```
No

```
(spa = spa_lookup(pool))
== NULL ?
```
Yes → 
```
locked ?
```
No →
```
return (ENOENT);
```
Yes →
```
mutex_exit(
&spa_namespace_lock);
```

No →
```
spa->spa_state ==
POOL_STATE_UNINITIALIZED
?
```
Yes →
```
zpool_rewind_policy_t policy;
zpool_get_rewind_policy(nvpolicy ?
nvpolicy : spa->spa_config, &policy);
```

```
policy.zrp_request &
ZPOOL_DO_REWIND ?
```
Yes →
```
state = SPA_LOAD_RECOVER;
```
No →
```
spa_activate(spa,
spa_mode_global);
```

```
state !=
SPA_LOAD_RECOVER ?
```
Yes →
```
spa->spa_last_ubsync_txg
= spa->spa_load_txg = 0;
```
No →
```
error = spa_load_best
(spa, state, B_FALSE,
policy.zrp_txg,
policy.zrp_request);
```

```
error == EBADF ?
```

* If vdev_validate() returns failure (indicated by
* EBADF), it indicates that one of the vdevs indicates
* that the pool has been exported or destroyed.  If
* this is the case, the config cache is out of sync and
* we should remove the pool from the namespace.

Yes →
```
spa_unload(spa);
spa_deactivate(spa);
spa_config_sync(spa, B_TRUE, B_TRUE);
spa_remove(spa);
```
No →
```
error ?
```

* We can't open the pool, but we still have useful
* information: the state of each vdev after the
* attempted vdev_open().  Return this to the user.

Yes →
```
config != NULL &&
spa->spa_config ?
```
Yes →
```
VERIFY(nvlist_dup(spa->spa_config, config, KM_SLEEP) == 0);
VERIFY(nvlist_add_nvlist(*config,
ZPOOL_CONFIG_LOAD_INFO, spa->spa_load_info) == 0);
```
No →
```
spa_unload(spa);
spa_deactivate(spa);
spa->spa_last_open_failed = error;
```

```
locked ?
```
Yes →
```
mutex_exit(
&spa_namespace_lock);
```
No →
```
*spapp = NULL;
return (error);
```

(from error ? No branch)
```
locked ?
```
Yes →
```
mutex_exit(
&spa_namespace_lock);
```
No →
```
return (ENOENT);
```

```
spa_open_ref(spa, tag);
```

```
config != NULL ?
```
Yes →
```
*config =
spa_config_generate(spa,
NULL, -1ULL, B_TRUE);
```

* If we've recovered the pool, pass back any information we
* gathered while doing the load.

```
state ==
SPA_LOAD_RECOVER ?
```
Yes →
```
VERIFY(nvlist_add_nvlist(
*config,
ZPOOL_CONFIG_LOAD_INFO,
spa->spa_load_info)
== 0);
```
No →
```
locked ?
```
Yes →
```
spa->spa_last_open_failed = 0;
spa->spa_last_ubsync_txg = 0;
spa->spa_load_txg = 0;
mutex_exit(&spa_namespace_lock);
```
No →
```
*spapp = spa;
return (0);
```

End

int spa_open_rewind(const
char *name, spa_t
**spapp, void *tag,
nvlist_t *policy,
nvlist_t **config)

return (spa_open_common
(name, spapp, tag,
policy, config));

End

```
int spa_open(const char
*name, spa_t
**spapp, void *tag)
```

```
return (spa_open_common
(name, spapp,
tag, NULL, NULL));
```

End

```
spa_t * spa_inject_addref
    (char *name)
```

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

```
spa_t *spa;
mutex_enter(&spa_namespace_lock);
```

```
(spa = spa_lookup(name))
    == NULL ?
```

No

Yes

```
spa->spa_inject_ref++;
mutex_exit(&spa_namespace_lock);
return (spa);
```

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

End

```
void spa_inject_
delref(spa_t *spa)
```

```
mutex_enter(&spa_namespace_lock);
spa->spa_inject_ref--;
mutex_exit(&spa_namespace_lock);
```

End

```
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 ?

No

Yes

```
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 ?

No

Yes

End

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

spa_spare_exists(guid,
&pool, NULL) &&
pool != 0ULL ?

No

Yes

```
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;
```

i++

```
static void
spa_add_l2cache(spa_t
*spa, nvlist_t *config)
```

* Add l2cache device
information to the
nvlist, including
vdev stats.

```
nvlist_t **l2cache;
uint_t i, j, nl2cache;
nvlist_t *nvroot;
uint64_t guid;
vdev_t *vd;
vdev_stat_t *vs;
uint_t vsc;
ASSERT(spa_config_held(spa,
SCL_CONFIG, RW_READER));
```

spa->spa_l2cache
.sav_count == 0 ?

Yes

No

```
VERIFY(nvlist_lookup_nvlist(config,
ZPOOL_CONFIG_VDEV_TREE, &nvroot) == 0);
VERIFY(nvlist_lookup_nvlist_array(spa->spa_l2cache.sav_config,
ZPOOL_CONFIG_L2CACHE, &l2cache, &nl2cache) == 0);
```

nl2cache != 0 ?

Yes

No

```
VERIFY(nvlist_add_nvlist_array(nvroot,
ZPOOL_CONFIG_L2CACHE, l2cache, nl2cache) == 0);
VERIFY(nvlist_lookup_nvlist_array(nvroot,
ZPOOL_CONFIG_L2CACHE, &l2cache, &nl2cache) == 0);
```

* Update level 2
cache device stats.

i = 0

i < nl2cache ?

No

Yes

End

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

j < spa->spa_
l2cache.sav_count ?

No

Yes

guid == spa->spa_l2cache
.sav_vdevs[j]->vdev_guid
?

Yes

No

```
vd = spa->spa_l2cache
.sav_vdevs[j];
```

j++

```
ASSERT(vd != NULL);
VERIFY(nvlist_lookup_uint64_array(l2cache[i],
ZPOOL_CONFIG_VDEV_STATS, (uint64_t **)&vs, &vsc) == 0);
vdev_get_stats(vd, vs);
```

i++

```
                              static void
                           spa_add_feature_stats
                              (spa_t *spa,
                             nvlist_t *config)


  nvlist_t *features;
  zap_cursor_t zc;
  zap_attribute_t za;
  ASSERT(spa_config_held(spa, SCL_CONFIG, RW_READER));
  VERIFY(nvlist_alloc(&features, NV_UNIQUE_NAME, KM_SLEEP) == 0);


                    spa->spa_feat_for_
                       read_obj != 0 ?
                                              Yes
                                                        zap_cursor_init(&zc,
                                                        spa->spa_meta_objset,
              No                                        spa->spa_feat_for_
                                                            read_obj)


                                            zap_cursor_retrieve(&zc,
                                               &za) == 0 ?

                                   No                          Yes

                                                    ASSERT(za.za_integer_length == sizeof
                  zap_cursor_fini(&zc);             (uint64_t) && za.za_num_integers == 1);
                                                    VERIFY3U(0, ==, nvlist_add_uint64(features,
                                                    za.za_name, za.za_first_integer));


                    spa->spa_feat_for_                             zap_cursor_advance(&zc)
                       write_obj != 0 ?
                                              Yes
                                                        zap_cursor_init(&zc,
                                                        spa->spa_meta_objset,
              No                                        spa->spa_feat_for_
                                                            write_obj)


                                            zap_cursor_retrieve(&zc,
                                               &za) == 0 ?

                                   No                          Yes

                                                    ASSERT(za.za_integer_length == sizeof
                  zap_cursor_fini(&zc);             (uint64_t) && za.za_num_integers == 1);
                                                    VERIFY3U(0, ==, nvlist_add_uint64(features,
                                                    za.za_name, za.za_first_integer));


  VERIFY(nvlist_add_nvlist(config,                            zap_cursor_advance(&zc)
  ZPOOL_CONFIG_FEATURE_STATS, features) == 0);
  nvlist_free(features);


                              End
```

```
int spa_get_stats(const
char *name, nvlist_t
**config, char *altroot,
size_t buflen)
```

```
int error;
spa_t *spa;
*config = NULL;
error = spa_open_common(name,
&spa, FTAG, NULL, config);
```

spa != NULL ?

* This still leaves a window of inconsistency where the spares
* or l2cache devices could change and the config would be
* self-inconsistent.

No

Yes

```
spa_config_enter(spa,
SCL_CONFIG, FTAG,
RW_READER);
```

*config != NULL ?

No

Yes

```
uint64_t loadtimes[2];
loadtimes[0] = spa->spa_loaded_ts.tv_sec;
loadtimes[1] = spa->spa_loaded_ts.tv_nsec;
VERIFY(nvlist_add_uint64_array(*config, ZPOOL_CONFIG_LOADED_TIME, loadtimes, 2) == 0);
VERIFY(nvlist_add_uint64(*config, ZPOOL_CONFIG_ERRCOUNT, spa_get_errlog_size(spa)) == 0);
```

spa_suspended(spa) ?

No

Yes

```
VERIFY(nvlist_add_uint64(
*config,
ZPOOL_CONFIG_SUSPENDED,
spa->spa_failmode) == 0);
```

```
spa_add_spares(spa, *config);
spa_add_l2cache(spa, *config);
spa_add_feature_stats(spa, *config);
```

* We want to get the alternate root
even for faulted pools, so we cheat
* and call spa_lookup() directly.

altroot ?

No

Yes

spa == NULL ?

No

Yes

```
mutex_enter(&spa_namespace_lock);
spa = spa_lookup(name);
```

```
spa_altroot(spa,
altroot, buflen);
```

spa ?

Yes

No

```
spa_altroot(spa,
altroot, buflen);
```

```
altroot[0] = '\0';
```

```
spa = NULL;
mutex_exit(&spa_namespace_lock);
```

spa != NULL ?

No

Yes

```
spa_config_exit(spa, SCL_CONFIG, FTAG);
spa_close(spa, FTAG);
```

```
return (error);
```

End

```
                        static int
                   spa_validate_aux_devs
                   (spa_t *spa, nvlist_t
                   *nvroot, uint64_t crtxg,
                   int mode, spa_aux_vdev_t
                      *sav, const char
                      *config, uint64_t
                          version,
                   vdev_labeltype_t label)
```

* 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.

```
nvlist_t **dev;
uint_t i, ndev;
vdev_t *vd;
int error;
ASSERT(spa_config_held(spa,
SCL_ALL, RW_WRITER) == SCL_ALL);
```

* It's acceptable to
have no devs specified.

```
nvlist_lookup_nvlist_
array(nvroot, config,
&dev, &ndev) != 0 ?
```

Yes → return (0);

No →

```
ndev == 0 ?
```

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

Yes → return (EINVAL);

No →

```
spa_version(spa)
< version ?
```

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

Yes → return (ENOTSUP);

No →
```
sav->sav_pending = dev;
sav->sav_npending = ndev;
i = 0
```

```
i < ndev ?
```

Yes →
```
(error =
spa_config_parse(spa,
&vd, dev[i], NULL,
0, mode)) != 0 ?
```

Yes →
```
vdev_free(vd);
error = EINVAL;
```

No →
```
!vd->vdev_ops-
>vdev_op_leaf ?
```

* The L2ARC currently only supports disk devices in
* kernel context.  For user-level testing, we allow it.

Yes →

No →
```
#ifdef _KERNEL ?
```

Yes →
```
(strcmp(config,
ZPOOL_CONFIG_L2CACHE ==
0) && strcmp(vd->vdev_ops
->vdev_op_type,
VDEV_TYPE_DISK) != 0 ?
```

No →

Yes →
```
error = ENOTBLK;
vdev_free(vd);
```

No →
```
vd->vdev_top = vd;
```

```
(error = vdev_open(vd))
== 0 && (error =
vdev_label_init(vd,
crtxg, label)) == 0 ?
```

Yes →
```
VERIFY(nvlist_add_uint64
(dev[i],
ZPOOL_CONFIG_GUID,
vd->vdev_guid) == 0);
```

No →
```
vdev_free(vd);
```

```
error && (mode !=
VDEV_ALLOC_SPARE && mode
!= VDEV_ALLOC_L2CACHE) ?
```

Yes →
```
out: sav->sav_pending = NULL;
sav->sav_npending = 0;
return (error);
```

No →
```
error = 0;
```

```
i++
```

End

```
static int
spa_validate_aux(spa_t
*spa, nvlist_t *nvroot,
uint64_t crtxg, int mode)
```

```
int error;
ASSERT(spa_config_held(spa, SCL_ALL, RW_WRITER) == SCL_ALL);
```

```
(error =
spa_validate_aux_devs
(spa, nvroot, crtxg,
mode, &spa->spa_spares,
ZPOOL_CONFIG_SPARES,
SPA_VERSION_SPARES,
VDEV_LABEL_SPARE)) != 0 ?
```

No

Yes

```
return (spa_validate_aux_
devs(spa, nvroot, crtxg,
mode, &spa->spa_l2cache,
ZPOOL_CONFIG_L2CACHE,
SPA_VERSION_L2CACHE,
VDEV_LABEL_L2CACHE));
```

```
return (error);
```

End

```
                                    static void
                                 spa_set_aux_vdevs(spa_
                                   aux_vdev_t *sav,
                                  nvlist_t **devs, int
                                     ndevs, const
                                    char *config)

                                        int i;

                              sav->sav_config != NULL ?

                     No                              Yes

VERIFY(nvlist_alloc(&sav->sav_config,                              nvlist_t **olddevs;
NV_UNIQUE_NAME, KM_SLEEP) == 0);    * Generate new dev list by concatentating with the    uint_t oldndevs;
VERIFY(nvlist_add_nvlist_array(sav->sav_config,   * current dev list.                 nvlist_t **newdevs;
config, devs, ndevs) == 0);

* Generate a
new dev list.

                    VERIFY(nvlist_lookup_nvlist_array(sav->sav_config, config, &olddevs, &oldndevs) == 0);
                    newdevs = kmem_alloc(sizeof (void *) * (ndevs + oldndevs), KM_SLEEP);
                    i = 0

                                  i < oldndevs ?

                         Yes               No

                    VERIFY(nvlist_dup(olddevs         i = 0
                       [i], &newdevs[i],
                       KM_SLEEP) == 0);

                                           i < ndevs ?

                           i++          Yes                        No

                    VERIFY(nvlist_dup(devs[i]    VERIFY(nvlist_remove(sav->sav_config, config, DATA_TYPE_NVLIST_ARRAY) == 0);
                      , &newdevs[i +          VERIFY(nvlist_add_nvlist_array(sav->sav_config, config, newdevs, ndevs + oldndevs) == 0);
                       oldndevs],            i = 0
                     KM_SLEEP) == 0);

                           i++              i < oldndevs + ndevs ?

                                         No              Yes

                    kmem_free(newdevs,         nvlist_free(newdevs[i]);
                    (oldndevs + ndevs) *
                      sizeof (void *));

                         End                                i++
```

```
void spa_l2cache_drop
(spa_t *spa)
```

* Stop and drop
level 2 ARC devices

```
vdev_t *vd;
int i;
spa_aux_vdev_t *sav
= &spa->spa_l2cache;
i = 0
```

i < sav->sav_count ?

No

Yes

End

```
uint64_t pool;
vd = sav->sav_vdevs[i];
ASSERT(vd != NULL);
```

```
spa_l2cache_exists(vd-
>vdev_guid, &pool) &&
pool != 0ULL &&
l2arc_vdev_present(vd) ?
```

No

Yes

l2arc_remove_vdev(vd);

i++

```
static nvlist_t *
spa_generate_rootconf
(char *devpath, char
*devid, uint64_t *guid)
```

```
nvlist_t *config;
nvlist_t *nvtop, *nvroot;
uint64_t pgid;
```

```
vdev_disk_read_rootlabel
(devpath, devid,
&config) != 0 ?
```

* Add this top-level
vdev to the child array.

No

Yes

```
VERIFY(nvlist_lookup_nvlist(config, ZPOOL_CONFIG_VDEV_TREE, &nvtop) == 0);
VERIFY(nvlist_lookup_uint64(config, ZPOOL_CONFIG_POOL_GUID, &pgid) == 0);
VERIFY(nvlist_lookup_uint64(config, ZPOOL_CONFIG_GUID, guid) == 0);
```

* Put this pool's
top-level vdevs
into a root vdev.

```
return (NULL);
```

```
VERIFY(nvlist_alloc(&nvroot, NV_UNIQUE_NAME, KM_SLEEP) == 0);
VERIFY(nvlist_add_string(nvroot, ZPOOL_CONFIG_TYPE, VDEV_TYPE_ROOT) == 0);
VERIFY(nvlist_add_uint64(nvroot, ZPOOL_CONFIG_ID, 0ULL) == 0);
VERIFY(nvlist_add_uint64(nvroot, ZPOOL_CONFIG_GUID, pgid) == 0);
VERIFY(nvlist_add_nvlist_array(nvroot, ZPOOL_CONFIG_CHILDREN, &nvtop, 1) == 0);
```

* Replace the existing vdev_tree with the new root vdev in
* this pool's configuration (remove the old, add the new).

```
VERIFY(nvlist_add_nvlist(config, ZPOOL_CONFIG_VDEV_TREE, nvroot) == 0);
nvlist_free(nvroot);
return (config);
```

End

```
static void
spa_alt_rootvdev(vdev_t
*vd, vdev_t **avd,
uint64_t *txg)
```

* Walk the vdev tree and see if we can find a device with "better"
* configuration. A configuration is "better" if the label on that
* device has a more recent txg.

```
int c = 0
```

c < vd->vdev_children ?

No → vd->vdev_ops->vdev_op_leaf ?

Yes → 
```
spa_alt_rootvdev(vd->vdev_child[c],
avd, txg);
```
→ c++

vd->vdev_ops->vdev_op_leaf ?

No → End

Yes →
```
nvlist_t *label;
uint64_t label_txg;
```

```
vdev_disk_read_rootlabel
(vd->vdev_physpath,
vd->vdev_devid,
&label) != 0 ?
```

Yes → End

No →
```
VERIFY(nvlist_lookup_
uint64(label,
ZPOOL_CONFIG_POOL_TXG,
&label_txg) == 0);
```

* Do we have a
better boot device?

label_txg > *txg ?

Yes →
```
*txg = label_txg;
*avd = vd;
```

No →
```
nvlist_free(label);
```

→ End

```
int spa_import_rootpool
(char *devpath,
char *devid)
```

* 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"

```
spa_t *spa;
vdev_t *rvd, *bvd, *avd = NULL;
nvlist_t *config, *nvtop;
uint64_t guid, txg;
char *pname;
int error;
```

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

```
config =
spa_generate_rootconf
(devpath, devid, &guid);
```

defined(_OBP) &&
defined(_KERNEL) ?

config == NULL ?

strstr(devpath,
"/iscsi/ssd") != NULL ?

iscsi boot

```
get_iscsi_bootpath_phy(devpath);
config = spa_generate_rootconf(devpath, devid, &guid);
```

config == NULL ?

No

```
VERIFY(nvlist_lookup_string(config, ZPOOL_CONFIG_POOL_NAME, &pname) == 0);
VERIFY(nvlist_lookup_uint64(config, ZPOOL_CONFIG_POOL_TXG, &txg) == 0);
mutex_enter(&spa_namespace_lock);
```

Yes

```
cmn_err(CE_NOTE, "Cannot read the
pool label from '%s'", devpath);
return (EIO);
```

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

(spa = spa_lookup(pname))
!= NULL ?

Yes

spa_remove(spa);

No

* Build up a vdev tree
based on the boot
device's label config.

```
spa = spa_add(pname, config, NULL);
spa->spa_is_root = B_TRUE;
spa->spa_import_flags = ZFS_IMPORT_VERBATIM;
```

```
VERIFY(nvlist_lookup_nvlist(config, ZPOOL_CONFIG_VDEV_TREE, &nvtop) == 0);
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
error = spa_config_parse(spa, &rvd, nvtop, NULL, 0, VDEV_ALLOC_ROOTPOOL);
spa_config_exit(spa, SCL_ALL, FTAG);
```

error ?

* Get the boot vdev.

Yes

No

(bvd = vdev_lookup_by_
guid(rvd,
guid)) == NULL ?

* Determine if there is
a better boot device.

```
mutex_exit(&spa_namespace_lock);
nvlist_free(config);
cmn_err(CE_NOTE, "Can not parse the config for pool '%s'", pname);
return (error);
```

No

Yes

```
avd = bvd;
spa_alt_rootvdev(rvd, &avd, &txg);
```

```
cmn_err(CE_NOTE, "Can not find the boot
vdev for guid %llu", (u_longlong_t)guid);
error = ENOENT;
```

* If the boot device is part of a spare vdev then ensure that
* we're booting off the active spare.

avd != bvd ?

No

bvd->vdev_parent->vdev_
ops == &vdev_spare_ops
&& !bvd->vdev_isspare ?

Yes

```
cmn_err(CE_NOTE, "The boot device is 'degraded'. Please " "try
booting from '%s'", avd->vdev_path);
error = EINVAL;
```

No

Yes

error = 0;

```
cmn_err(CE_NOTE, "The boot device is currently spared. Please "
"try booting from '%s'", bvd->vdev_parent->
vdev_child[bvd->vdev_parent->vdev_children - 1]->vdev_path);
error = EINVAL;
```

```
out: spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
vdev_free(rvd);
spa_config_exit(spa, SCL_ALL, FTAG);
mutex_exit(&spa_namespace_lock);
nvlist_free(config);
return (error);
```

End

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

* Import a non-root pool
  into the system.

```
spa_t *spa;
char *altroot = NULL;
spa_load_state_t state = SPA_LOAD_IMPORT;
zpool_rewind_policy_t policy;
uint64_t mode = spa_mode_global;
uint64_t readonly = B_FALSE;
int error;
nvlist_t *nvroot;
nvlist_t **spares, **l2cache;
uint_t nspares, nl2cache;
```

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

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

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

* 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 ?
```

```
mode = FREAD;
```

* Verbatim import - Take a pool and insert it into the namespace
  as if it had been loaded at boot.

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

```
spa->spa_import_flags &
SPA_IMPORT_VERBATIM ?
```

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

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

```
props != NULL ?
```

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

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

```
policy.zrp_request &
ZPOOL_DO_REWIND ?
```

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

* Pass off the heavy lifting to spa_load().  Pass TRUE for mosconfig
  because the user-supplied config is actually the one to trust when
  doing an import.

```
state = SPA_LOAD_RECOVER;
```

```
state !=
SPA_LOAD_RECOVER ?
```

```
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 sparelist, as it doesn't have any validity
  anymore, and conflicts with spa_has_spares).

```
spa->spa_spares.
sav_config ?
```

```
nvlist_free(spa->spa_spares.sav_config);
spa->spa_spares.sav_config = NULL;
spa_load_spares(spa);
```

```
spa->spa_l2cache.
sav_config ?
```

```
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 == 0 ?
```

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

```
error == 0 ?
```

```
error = spa_validate_aux
(spa, nvroot, -1ULL,
VDEV_ALLOC_L2CACHE);
```

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

```
props != NULL ?
```

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

```
error != 0 || (props &&
spa_writeable(spa) &&
persist =
spa_prop_set(spa,
props)) ?
```

```
spa_async_resume(spa);
```

* Override any spares and level 2 cache devices as specified by
  the user, as these may have correct device names/devids, etc.

```
spa_unload(spa);
spa_deactivate(spa);
spa_remove(spa);
mutex_exit(&spa_namespace_lock);
return (error);
```

```
nvlist_lookup_nvlist_
array(nvroot,
ZPOOL_CONFIG_SPARES,
&spares, &nspares) == 0 ?
```

```
spa->spa_spares.
sav_config ?
```

```
VERIFY(nvlist_alloc(&spa-
>spa_spares.sav_config,
NV_UNIQUE_NAME,
KM_SLEEP) == 0);
```

```
VERIFY(nvlist_remove(spa-
>spa_spares.sav_config,
ZPOOL_CONFIG_SPARES,
DATA_TYPE_NVLIST_ARRAY)
== 0);
```

```
VERIFY(nvlist_add_nvlist_array(spa->spa_spares.sav_config,
ZPOOL_CONFIG_SPARES, spares, nspares) == 0);
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
spa_load_spares(spa);
spa_config_exit(spa, SCL_ALL, FTAG);
spa->spa_spares.sav_sync = B_TRUE;
```

```
nvlist_lookup_nvlist_
array(nvroot,
ZPOOL_CONFIG_L2CACHE,
&l2cache,
&nl2cache) == 0 ?
```

```
spa->spa_l2cache.
sav_config ?
```

```
VERIFY(nvlist_alloc(&spa-
>spa_l2cache.sav_config,
NV_UNIQUE_NAME,
KM_SLEEP) == 0);
```

```
VERIFY(nvlist_remove(spa-
>spa_l2cache.sav_config,
ZPOOL_CONFIG_L2CACHE,
DATA_TYPE_NVLIST_ARRAY)
== 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;
```

* Check for any
  removed devices.

```
spa->spa_autoreplace ?
```

```
spa_aux_check_removed(&spa->spa_spares);
spa_aux_check_removed(&spa->spa_l2cache);
```

* Update the config
  cache to include the
  newly-imported pool.

```
spa_writeable(spa) ?
```

```
spa_config_update(spa,
SPA_CONFIG_UPDATE_POOL);
```

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

```
spa_async_request(spa, SPA_ASYNC_AUTOEXPAND);
mutex_exit(&spa_namespace_lock);
spa_history_log_version(spa, "import");
return (0);
```

```
                    nvlist_t *
               spa_tryimport(nvlist_t
                    *tryconfig)
                          │
                          ▼
            ┌──────────────────────────┐
            │ nvlist_t *config = NULL;  │
            │ char *poolname;           │
            │ spa_t *spa;               │
            │ uint64_t state;           │
            │ int error;                │
            └──────────────────────────┘
                          │
                          ▼
            ╱───────────────────────╲
           ╱ nvlist_lookup_string     ╲
          ╱   (tryconfig,              ╲
          ╲ ZPOOL_CONFIG_POOL_NAME,    ╱
           ╲   &poolname) ?           ╱
            ╲───────────────────────╱
          Yes │              │ No
              │              ▼
              │      ╱───────────────────────╲
              │     ╱ nvlist_lookup_uint64    ╲        * Create and initialize
              │    ╱   (tryconfig,             ╲        the spa structure.
              │    ╲ ZPOOL_CONFIG_POOL_STATE,  ╱
              │     ╲   &state) ?             ╱
              │      ╲───────────────────────╱
              │     Yes │          │ No
    ┌──────────────┐   │          ▼
    │ return (NULL);│   │   ┌────────────────────────────────────────────┐
    └──────────────┘   │   │ mutex_enter(&spa_namespace_lock);           │   * Pass off the heavy lifting to spa_load().
          │            │   │ spa = spa_add(TRYIMPORT_NAME, tryconfig, NULL);│  * Pass TRUE for mosconfig because the user-supplied config
          │            ▼   │ spa_activate(spa, FREAD);                   │   * is actually the one to trust when doing an import.
          │   ┌──────────────┐└────────────────────────────────────────────┘
          │   │ return (NULL);│          │
          │   └──────────────┘          ▼
          │          │         ┌────────────────────┐
          │          │         │ error = spa_load(spa,│      * If 'tryconfig' was at
          │          │         │ SPA_LOAD_TRYIMPORT,  │       least parsable, return
          │          │         │ SPA_IMPORT_EXISTING, │       the current config.
          │          │         │ B_TRUE);             │
          │          │         └────────────────────┘
          │          │                  │
          │          │                  ▼
          │          │         ╱──────────────────╲
          │          │        ╱ spa->spa_root_vdev ╲
          │          │        ╲    != NULL ?        ╱
          │          │         ╲──────────────────╱
          │          │        No │         │ Yes
          │          │           │         ▼
          │          │           │  ┌──────────────────────────────────────────────────────────────────────────────────┐
          │          │           │  │ config = spa_config_generate(spa, NULL, -1ULL, B_TRUE);                           │
          │          │           │  │ VERIFY(nvlist_add_string(config, ZPOOL_CONFIG_POOL_NAME, poolname) == 0);         │
          │          │           │  │ VERIFY(nvlist_add_uint64(config, ZPOOL_CONFIG_POOL_STATE, state) == 0);           │
          │  * If the bootfs property│ │ VERIFY(nvlist_add_uint64(config, ZPOOL_CONFIG_TIMESTAMP, spa->spa_uberblock.ub_timestamp) == 0);│
          │  exists on this pool then we│ │ VERIFY(nvlist_add_nvlist(config, ZPOOL_CONFIG_LOAD_INFO, spa->spa_load_info) == 0);│
          │  * copy it out so that external consumers can tell which└──────────────────────────────────────────────────────────────────────────────────┘
          │  * pools are bootable.                    │
          │          │                                ▼
          │          │                      ╱──────────────────╲
          │          │                     ╱ (!error || error ==╲
          │          │                     ╲    EEXIST) &&       ╱
          │          │                      ╲ spa->spa_bootfs ? ╱
          │          │                       ╲───────────────╱
          │          │                      No │        │ Yes
          │          │                         │        ▼
          │          │                         │ ┌──────────────────┐
          │          │                         │ │ char *tmpname =   │    * We have to play games with the name since the
          │          │                         │ │ kmem_alloc(MAXPATHLEN,│  * pool was opened as TRYIMPORT_NAME.
          │          │                         │ │ KM_SLEEP);        │
          │          │                         │ └──────────────────┘
          │          │                         │        │
          │          │                         │        ▼
          │          │                         │  ╱──────────────────╲
          │          │                         │ ╱ dsl_dsobj_to_dsname(spa╲
          │          │                         │ ╲    name(spa),          ╱
          │          │                         │  ╲ spa->spa_bootfs,     ╱
          │          │                         │   ╲ tmpname) == 0 ?    ╱
          │          │                         │    ╲────────────────╱
          │          │                         │   No │          │ Yes
          │          │                         │      │          ▼
          │          │                         │      │  ┌─────────────────────────────┐
          │          │                         │      │  │ char *cp;                   │
          │          │                         │      │  │ char *dsname = kmem_alloc(MAXPATHLEN, KM_SLEEP);│
          │          │                         │      │  │ cp = strchr(tmpname, '/');  │
          │          │                         │      │  └─────────────────────────────┘
          │          │                         │      │          │
          │          │                         │      │          ▼
          │          │                         │      │    ╱──────────╲
          │          │                         │      │   ╱ cp == NULL ?╲
          │          │                         │      │    ╲──────────╱
          │          │                         │      │   Yes │     │ No
          │          │                         │      │       ▼     ▼
          │          │                         │      │  ┌──────────┐ ┌────────────────────┐
          │          │                         │      │  │ (void) strlcpy(dsname,│ │ (void) snprintf(dsname,│
          │          │                         │      │  │ tmpname, MAXPATHLEN);│  │ MAXPATHLEN, "%s/%s",│
          │          │                         │      │  └──────────┘ │ poolname, ++cp);   │
          │          │                         │      │       │       └────────────────────┘
          │          │                         │      │       └───────┬─────┘
          │          │                         │      │               ▼
          │          │                         │      │  ┌─────────────────────────────┐
          │          │                         │      │  │ VERIFY(nvlist_add_string(config,│
          │          │                         │      │  │ ZPOOL_CONFIG_BOOTFS, dsname) == 0);│
          │          │                         │      │  │ kmem_free(dsname, MAXPATHLEN);│
          │          │                         │      │  └─────────────────────────────┘
          │          │                         │      │               │
          │          │                         │      ▼               ▼
          │          │                         │  ┌──────────────┐   * Add the list of hot
          │          │                         │  │ kmem_free(tmpname,│   spares and level
          │          │                         │  │ MAXPATHLEN);  │    2 cache devices.
          │          │                         │  └──────────────┘
          │          │                         │          │
          │          │                         ▼          ▼
          │          │                ┌──────────────────────────────────────────┐
          │          │                │ spa_config_enter(spa, SCL_CONFIG, FTAG, RW_READER);│
          │          │                │ spa_add_spares(spa, config);              │
          │          │                │ spa_add_l2cache(spa, config);             │
          │          │                │ spa_config_exit(spa, SCL_CONFIG, FTAG);   │
          │          │                └──────────────────────────────────────────┘
          │          │                         │
          │          ▼                         ▼
          │   ┌──────────────────────────────────┐
          │   │ spa_unload(spa);                  │
          │   │ spa_deactivate(spa);              │
          │   │ spa_remove(spa);                  │
          │   │ mutex_exit(&spa_namespace_lock);  │
          │   │ return (config);                  │
          │   └──────────────────────────────────┘
          │          │
          ▼          ▼
        ◇ End ◇
```

```
static int
spa_export_common(char
*pool, int new_state,
nvlist_t **oldconfig,
boolean_t force,
boolean_t hardforce)
```

* 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.

`spa_t *spa;`

**oldconfig ?**

→ Yes → `*oldconfig = NULL;`

→ No

**!(spa_mode_global & FWRITE) ?**

→ Yes → `return (EROFS);`

→ No → `mutex_enter(&spa_namespace_lock);`

**(spa = spa_lookup(pool) == NULL ?**

* Put a hold on the pool, drop the
  namespace lock, stop async tasks,
* reacquire the namespace lock, and see if we can export.

→ Yes → `mutex_exit(&spa_namespace_lock); return (ENOENT);`

→ No →
```
spa_open_ref(spa, FTAG);
mutex_exit(&spa_namespace_lock);
spa_async_suspend(spa);
mutex_enter(&spa_namespace_lock);
spa_close(spa, FTAG);
```

* The pool will be in core if it's openable,
* in which case we can modify its state.

**spa->spa_state != POOL_STATE_UNINITIALIZED && spa->spa_sync_on ?**

* Objsets may be open only because they're dirty, so we
* have to force it to sync before checking spa_refcnt.

→ Yes → `txg_wait_synced(spa->spa_dsl_pool, 0);`

→ No

* A pool cannot be exported or destroyed if there are active
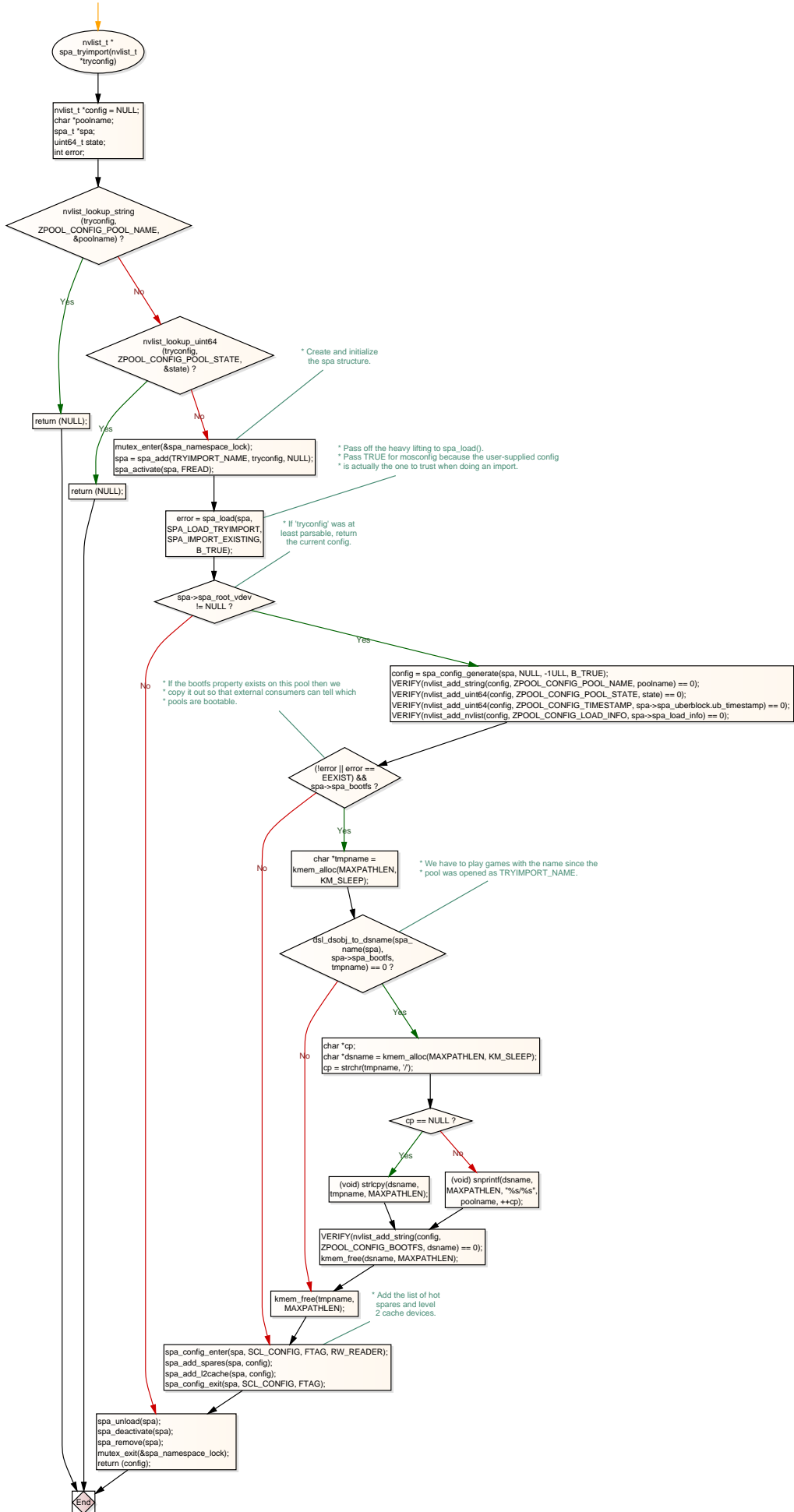* references.  If we are resetting a pool, allow references by
* fault injection handlers.

```
!spa_refcount_zero(spa)
|| (spa->spa_inject_ref
!= 0 && new_state !=
POOL_STATE_UNINITIALIZED)
?
```

* 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 forcedly exported.

→ Yes → `spa_async_resume(spa); mutex_exit(&spa_namespace_lock); return (EBUSY);`

→ No →

```
!force && new_state ==
POOL_STATE_EXPORTED &&
spa_has_active_
shared_spare(spa) ?
```

* 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.

→ Yes → `spa_async_resume(spa); mutex_exit(&spa_namespace_lock); return (EXDEV);`

→ No →

**new_state != POOL_STATE_UNINITIALIZED && !hardforce ?**

→ Yes →
```
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);
```

→ No →

`spa_event_notify(spa, NULL, ESC_ZFS_POOL_DESTROY);`

**spa->spa_state != POOL_STATE_UNINITIALIZED ?**

→ Yes → `spa_unload(spa); spa_deactivate(spa);`

→ No →

**oldconfig && spa->spa_config ?**

→ Yes → `VERIFY(nvlist_dup(spa->spa_config, oldconfig, 0) == 0);`

→ No →

**new_state != POOL_STATE_UNINITIALIZED ?**

→ Yes → **!hardforce ?**

→ No

→ Yes → `spa_config_sync(spa, B_TRUE, B_TRUE);`

→ No → `spa_remove(spa);`

`mutex_exit(&spa_namespace_lock); return (0);`

**End**

```
int spa_destroy(char
      *pool)
```

* Destroy a storage pool.

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

End

int spa_export(char *pool, nvlist_t **oldconfig, boolean_t force, boolean_t hardforce)

* Export a storage pool.

return (spa_export_common (pool, POOL_STATE_ EXPORTED, oldconfig, force, hardforce));

End

int spa_reset(char *pool)

* 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));
```

End

```
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 *rvd = 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);
```

(error =
spa_config_parse(spa,
&vd, nvroot, NULL, 0,
VDEV_ALLOC_ADD)) != 0 ?

spa_vdev_exit()
will clear this

No → spa->spa_pending_vdev
= vd;

Yes

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

nvlist_lookup_nvlist_
array(nvroot,
ZPOOL_CONFIG_SPARES,
&spares, &nspares) != 0 ?

Yes → nspares = 0;

No

nvlist_lookup_nvlist_
array(nvroot,
ZPOOL_CONFIG_L2CACHE,
&l2cache,
&nl2cache) != 0 ?

Yes → nl2cache = 0;

No

vd->vdev_children == 0
&& nspares == 0
&& nl2cache == 0 ?

Yes

No

vd->vdev_children != 0
&& (error =
vdev_create(vd, txg,
B_FALSE)) != 0 ?

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

```
return (spa_vdev_exit
(spa, vd, txg, EINVAL));
```

No

Yes

(error =
spa_validate_aux(spa,
nvroot, txg,
VDEV_ALLOC_ADD)) != 0 ?

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

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

Yes

No → int c = 0

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

c < vd->vdev_children ?

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

No → nspares != 0 ?

Yes → id = 0

Yes

```
spa_set_aux_vdevs(&spa->spa_spares, spares, nspares, ZPOOL_CONFIG_SPARES);
spa_load_spares(spa);
spa->spa_spares.sav_sync = B_TRUE;
```

No

id < rvd->vdev_children ?

nl2cache != 0 ?

Yes

rvd->vdev_child[id]
->vdev_ishole ?

No

* We have to be careful when adding new vdevs to an existing pool.
* If other threads start allocating from these vdevs before we
* sync the config cache, and we lose power, then upon reboot we may
* fail to open the pool because there are DVAs that the config cache
* can't translate.  Therefore, we first add the vdevs without
* initializing metaslabs; sync the config cache (via spa_vdev_exit());
* and then let spa_config_update() initialize the new metaslabs.
*
* spa_load() checks for added-but-not-initialized vdevs, so that
* if we lose power at any point in this sequence, the remaining
* steps will be completed the next time we load the pool.

Yes

Yes → vdev_free(rvd-
>vdev_child[id]);

No → id++

```
spa_set_aux_vdevs(&spa->spa_l2cache, l2cache, nl2cache, ZPOOL_CONFIG_L2CACHE);
spa_load_l2cache(spa);
spa->spa_l2cache.sav_sync = B_TRUE;
```

No

```
(void) spa_vdev_exit(spa, vd, txg, 0);
mutex_enter(&spa_namespace_lock);
spa_config_update(spa, SPA_CONFIG_UPDATE_POOL);
mutex_exit(&spa_namespace_lock);
return (0);
```

```
tvd = vd->vdev_child[c];
vdev_remove_child(vd, tvd);
tvd->vdev_id = id;
vdev_add_child(rvd, tvd);
vdev_config_dirty(tvd);
```

c++

End

spa_vdev_attach
(spa_uint64_t guid,
nvlist_t *nvroot,
int replacing)

* 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. In actuality reuse all the same ops() but has a few
* extra rules: you can't attach to it after the vdev has 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 *pvops;
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);

oldvd == NULL ?

⟶ spa->spa_root_
vdev_op_leaf ?

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

pvd = oldvd->vdev_parent

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

(error =
spa_config_parse(spa,
&newrootvd, nvroot,
NULL, 0,
VDEV_ALLOC_ATTACH))
!= 0 ?

newrootvd->vdev_children
!= 1 ?

newvd = newrootvd
->vdev_child[0]

return (spa_vdev_exit
(spa, NULL,
txg, EINVAL));

return (spa_vdev_exit
(spa, newrootvd,
txg, EINVAL));

!newvd->vdev_ops->
vdev_op_leaf ?

return (spa_vdev_exit
(spa, newrootvd,
txg, EINVAL));

(error =
vdev_create(newrootvd,
txg, replacing)) != 0 ?

* Spares can't
replace logs

!oldvd->vdev_top->vdev
_islog &&
newvd->vdev_isspare

return (spa_vdev_exit
(spa, newrootvd,
txg, error));

return (spa_vdev_exit
(spa, newrootvd,
txg, ENOTSUP));

!replacing

* Active hot spares can only be replaced by inactive hot
* spares.

pvd->vdev_ops ==
&vdev_spare_ops &&
oldvd->vdev_isspare &&
!spa_has_spares(spa,
newvd->vdev_guid)

* If the source is a hot spare, and the pooled 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 'replace' state is
* the same (spare replaces spare, non-spare replaces
* non-spare).

pvd->vdev_ops ==
&vdev_replacing_ops &&
spa_version(spa) <
SPA_VERSION_MULTI_REPLACE
?

pvd->vdev_ops ==
&vdev_spare_ops &&
newvd->vdev_isspare !=
oldvd->vdev_isspare

return (spa_vdev_exit
(spa, newrootvd,
txg, ENOTSUP));

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

pvd->vdev_ops !=
&vdev_mirror_ops &&
pvd->vdev_ops !=
&vdev_root_ops ?

return (spa_vdev_exit
(spa, newrootvd,
txg, ENOTSUP));

* Make sure the new
device is big enough.

!newvd->vdev_isspare

pvops = &vdev_mirror_ops

pvops = &vdev_spare_ops

pvops = &vdev_replacing_
ops;

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

!!newvd->vdev_asize
vdev_get_min_asize(oldvd)
?

return (spa_vdev_exit
(spa, newrootvd,
txg, EOVERFLOW));

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

strcmp(oldvd->vdev_path,
newvd->vdev_path) == 0 ?

return (spa_vdev_exit
(spa, newrootvd,
txg, EDOM));

spa_strfree(oldvd->vdev_path);
oldvd->vdev_path = kmem_alloc(strlen(newvd->vdev_p)ath) + 5, KM_SLEEP);
(void) sprintf(oldvd->vdev_path, "%s/%s", newvd->vdev_path, "old");

oldvd->vdev_devid
!= NULL ?

* mark the device
being replaced

spa_strfree(oldvd->vdev_devid);
oldvd->vdev_devid = NULL;

newvd->vdev_replacing
= B_TRUE;

* 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

pvd = vdev_add_parent
(oldvd, pvops);

ASSERT(pvd->vdev_top->vdev_parent == rvd);
ASSERT(pvd->vdev_ops == pvops);
ASSERT(oldvd->vdev_parent == pvd);

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

vdev_remove_child(newrootvd, newvd);
newvd->vdev_id = pvd->vdev_children;
newvd->vdev_crtxg = oldvd->vdev_crtxg;
tvd = newvd->vdev_top;
ASSERT(pvd->vdev_top == tvd);
ASSERT(tvd->vdev_parent == rvd);
vdev_config_dirty(tvd);

* 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().

dtl_max_txg = txg + TXG_CONCURRENT_STATES;
vdev_dtl_dirty(newvd, DTL_MISSING,
TXG_INITIAL, dtl_max_txg - TXG_INITIAL);

newvd->vdev_isspare

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

oldvdpath = spa_strdup(oldvd->vdev_path);
newvdpath = spa_strdup(newvd->vdev_path);
newvd_isspare = newvd->vdev_isspare;

* Mark newvd's DTL
dirty in this txg.

vdev_dirty(tvd,
VDD_DTL, newvd, txg);

* Restart the resilver

dsl_resilver_restart(spa-
>spa_dsl_pool);

* Commit the config

(void) spa_vdev_exit(spa, newrootvd, dtl_max_txg, 0);
spa_history_log_internal(spa, "vdev attach", NULL, "%s vdev=%s %s vdev=%s", replacing &&
newvd_isspare ? "spare in" : replacing ? "replace" :
"attach", newvdpath, replacing ? "for" : "to", oldvdpath);
spa_strfree(newvdpath);

spa->vdev_bcache

spa_event_notify(spa,
newvd, ESC_ZFS_
BOOTFS_VDEV_ATTACH)

return (0)

spa_vdev_detach(spa, *spa_uint64_t guid, uint64_t cguid, int replace_done)

*Detach a device from a mirror or replacing vdev.
* If 'replace_done' is specified, only detach it if that parent
* is a replacing vdev.

```
uint64_t txg;
int error;
vdev_t *vd = 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

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

pvd = vd->vdev_parent

*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 R(B,C) involving devices B and C. The user's intent is replacing
* B to go from M(A,R(B,C)) 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 decide to detach C,
* the replacement of B completes. We would have R(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.

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

pvd->vdev_guid != pguid &&
(pguid != 0)

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

replace_done &&
(vdev_replacing_ops &&
pvd->vdev_ops !=
&vdev_spare_ops

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

```
ASSERT(pvd->vdev_ops !=
&vdev_spare_ops ||
spa_version(spa) >=
SPA_VERSION_SPARES);
```

*Only mirror,
* replacing, and spare
* vdevs support detach.

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

pvd->vdev_ops ==
&vdev_replacing_ops &&
pvd->vdev_ops !=
&vdev_mirror_ops &&
pvd->vdev_ops !=
&vdev_spare_ops

*If this device has the only valid copy of some data,
* we cannot safely detach it.

vdev_dtl_required(vd)

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

```
ASSERT(pvd->vdev_children
>= 2);
```

*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.

pvd->vdev_ops ==
&vdev_replacing_ops &&
vd->vdev_id > 0 &&
vd->vdev_path != NULL

```
spa_t len = strlen(vd->vdev_path);
int c = 0
```

pvd->vdev_children

```
cvd = pvd->vdev_child[c];
```

int c == id ||
cvd->vdev_path == NULL

```
strlen(cvd->vdev_path)
== len &&
strncmp(cvd->vdev_path,
vd->vdev_path, len) == 0 &&
memcmp(vd->vdev_path + len, "/old", 4) == 0
```

*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.

```
spa_strfree(vd->vdev_path);
vd->vdev_path = spa_strdup(cvd->vdev_path);
```

pvd->vdev_ops ==
&vdev_spare_ops &&
vd->vdev_id == 0 &&
pvd->vdev_child[pvd
->vdev_children -
1]->vdev_isspare

*Erase the disk labels so the disk can be used for other things.
* This must be done after all other error cases are handled,
* but before we disembowel vd (so we can still do I/O to it).
* But if we can't do it, don't treat the error as fatal --
* it may be that the unavailability of the disk is the reason
* it's being detached.

unspare = B_TRUE;

*Remove vd from its
* parent and compact the
* parent's children.

```
error = vdev_label_init
(vd, 0, VDEV_
LABEL_REMOVE);
```

*Remember one of the
* remaining children so we
* can get tvd below.

```
vdev_remove_child(pvd, vd);
vdev_compact_children(pvd);
```

```
cvd = pvd->vdev_child[pvd
->vdev_children - 1];
```

*If we need to remove the remaining child from the list of top-level vdevs,
* do it now, marking the vdev as no longer a spare in the process.
* We must do this before vdev_remove_parent(), because that can
* change the GUID if it creates a new toplevel GUID. For a similar
* reason, we must remove the spare now, or the same tvg as the detach,
* otherwise someone could attach a new spare, change the GUID, and
* the subsequent attempt to spa_vdev_remove(unspare_guid) would fail.

```
ASSERT(cvd->vdev_isspare);
spa_spare_remove(cvd);
unspare_guid = cvd->vdev_guid;
(void) spa_vdev_remove(spa, unspare_guid, B_TRUE);
cvd->vdev_unspare = B_TRUE;
```

*If the parent mirror/replacing vdev only has one child,
* the parent is no longer needed. Remove it from the tree.

cvd->vdev_children

pvd->vdev_ops ==

```
cvd->vdev_unspare
= B_FALSE;
```

```
vdev_remove_parent(cvd);
pvd->vdev_reopening = B_FALSE;
```

*The dtl isn't set (old now because the parent we just removed
* may have been the previous top-level vdev.

```
tvd = cvd->vdev_top;
ASSERT(tvd->vdev_parent == tvd);
```

*Reevaluate the
* parent vdev state.

```
vdev_propagate_
state(cvd);
```

*If the 'autoexpand' property is set on the pool then automatically
* try to expand the size of the pool. For example if the device we
* just detached was smaller than the others, it may be possible to
* add metaslabs (i.e. grow the pool). We need to reopen the vdev
* first so that we can obtain the updated sizes of the leaf vdevs.

spa_vdev_autoexpand(spa)

```
vdev_reopen(tvd);
vdev_expand(tvd, txg);
```

*Mark vd's DTL as dirty in this txg. vdev_dtl_sync() will see that
* vd->vdev_detached is set and free vd's DTL object in syncing context.
* But first make sure we're not on any "other" txg's DTL list, to
* prevent vd from being accessed after it's freed.

```
vdev_dirty(tvd);
```

```
vdpath = spa_strdup(vd->vdev_path);
int c = 0
```

VDEV_RAID

vd->vdev_detached = B_TRUE;
vdev_dirty(vd, VDD_DTL, vd, txg);

```
(void) txg_list_remove_
this(&vd->vdev_dtl_list,
vd, t);
```

*hang on to the spa
* before we
* release this lock

```
spa_event_notify(spa, vd, ESC_ZFS_VDEV_REMOVE);
```

*If this was the removal of the original device in a hot spare vdev,
* then we want to go through and remove the device from the hot spare
* list of every other pool.

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

spa_t *altspa = NULL;
&vdev_enter(&spa_namespace_lock);

altspa =
spa_next(altspa)
!= NULL

altspa->spa_state !=
POOL_STATE_ACTIVE
altspa == spa

*search the rest of the
* spares to remove

```
mutex_exit(
&spa_namespace_lock);
```

```
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);
```

*all done with the
* spa; OK to release

```
spa_vdev_resilver_
done(spa);
```

```
mutex_enter(&spa_namespace_lock);
spa_close(spa, FTAG);
mutex_exit(&spa_namespace_lock);
return (error);
```

```
static nvlist_t *
spa_nvlist_lookup_by_guid
(nvlist_t **nvpp, int
count, uint64_t
target_guid)
```

int i = 0

i < count ?

**No**

**Yes**

```
uint64_t guid;
VERIFY(nvlist_lookup_uint64(nvpp[i],
ZPOOL_CONFIG_GUID, &guid) == 0);
```

return (NULL);

guid == target_guid ?

**Yes**

**No**

return (nvpp[i]);

i++

End

```
                                        static void
                                   spa_vdev_remove_aux
                                  (nvlist_t *config, char
                                   *name, nvlist_t **dev,
                                    int count, nvlist_t
                                    *dev_to_remove)

                                  nvlist_t **newdev = NULL;

                                        count > 1 ?
                                                        Yes
                                     No                         newdev =
                                                           kmem_alloc((count - 1) *
                                                                sizeof (void
                                                             *), KM_SLEEP);

                                       int i = 0, j = 0

                                         i < count ?
                            No                              Yes

  VERIFY(nvlist_remove(config, name, DATA_TYPE_NVLIST_ARRAY) == 0);
  VERIFY(nvlist_add_nvlist_array(config, name, newdev, count - 1) == 0);
  int i = 0                                      dev[i] == dev_to_remove ?
                                                        Yes        No
                                                                        VERIFY(nvlist_dup(dev[i],
                  i < count - 1 ?                                          &newdev[j++],
          No                     Yes                                      KM_SLEEP) == 0);

      count > 1 ?           nvlist_free(newdev[i]);                              i++
                 Yes
   No
           kmem_free(newdev, (count          i++
            - 1) * sizeof (void *));

         End
```

```
                    static int
             spa_vdev_remove_evacuate
               (spa_t *spa, vdev_t *vd)
```

* Evacuate the device.

```
uint64_t txg;
int error = 0;
ASSERT(MUTEX_HELD(&spa_namespace_lock));
ASSERT(spa_config_held(spa, SCL_ALL, RW_WRITER) == 0);
ASSERT(vd == vd->vdev_top);
```

* Evacuate the device.  We don't hold the config lock as writer
* since we need to do I/O but we do keep the
* spa_namespace_lock held.  Once this completes the device
* should no longer have any blocks allocated on it.

vd->vdev_islog ?

Yes

No

vd->vdev_stat.vs_alloc
!= 0 ?

Yes

No

```
error = ENOTSUP;
```

```
error = spa_
offline_log(spa);
```

error ?

* The evacuation succeeded.  Remove any remaining MOS metadata
* associated with this vdev, and wait for these changes to sync.

Yes

No

```
return (error);
```

```
ASSERT0(vd->vdev_stat.vs_alloc);
txg = spa_vdev_config_enter(spa);
vd->vdev_removing = B_TRUE;
vdev_dirty(vd, 0, NULL, txg);
vdev_config_dirty(vd);
spa_vdev_config_exit(spa, NULL, txg, 0, FTAG);
return (0);
```

End

```
static void
spa_vdev_remove_from_
namespace(spa_t
*spa, vdev_t *vd)
```

* Complete the removal
by cleaning up
the namespace.

```
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);
```

* Only remove any
devices which are empty.

vd->vdev_stat.vs_alloc
!= 0 ?

No

Yes

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

list_link_active(&vd-
>vdev_state_dirty_node) ?

No

Yes

vdev_state_clean(vd);

list_link_active(&vd-
>vdev_config_dirty_node)
?

No

Yes

vdev_config_clean(vd);

vdev_free(vd);

last_vdev ?

Yes

No

```
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);

End
```

```
int spa_vdev_remove(spa_t
*spa, uint64_t guid,
boolean_t unspare)
```

* 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.

```
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));
```

!locked ?

Yes → `txg = spa_vdev_enter(spa);`

No

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

```
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->spa_l2cache.sav_
vdevs != NULL &&
nvlist_lookup_nvlist_
array(spa->spa_l2cache
.sav_config,
ZPOOL_CONFIG_L2CACHE,
&l2cache, &nl2cache) ==
0 && (nv =
spa_nvlist_lookup_by_guid
(l2cache, nl2cache,
guid)) != NULL ?
```

* Cache devices can
always be removed.

Yes →
```
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;
```

No → vd != NULL &&
vd->vdev_islog ?

Yes →
```
ASSERT(!locked);
ASSERT(vd == vd->vdev_top);
```

* XXX - Once we have bp-rewrite this should
* become the common case.

* Stop allocating
from this vdev.

`mg = vd->vdev_mg;`

* Wait for the youngest allocations and frees to sync,
* and then wait for the deferral of those frees to finish.

```
metaslab_
group_passivate(mg);
```

* Attempt to
evacuate the vdev.

```
spa_vdev_config_exit(spa,
NULL, txg +
TXG_CONCURRENT_STATES +
TXG_DEFER_SIZE, 0, FTAG);
```

```
error = spa_vdev_remove_evacuate(spa, vd);
txg = spa_vdev_config_enter(spa);
```

* If we couldn't
evacuate the
vdev, unwind.

error ?

* Only remove the hot spare if it's not currently in use
* in this pool.

Yes (spares)

vd == NULL || unspare ?

Yes →
```
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;
```

No → error = EBUSY;

* There is no vdev of
any kind with the
specified guid.

error = ENOENT;

vd != NULL ?

No → error = ENOENT;

Yes →

* Normal vdevs cannot be
removed (yet).

error = ENOTSUP;

* Clean up the
vdev namespace.

* Normal vdevs cannot be

```
spa_vdev_remove_from_
namespace(spa, vd);
```

No

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

!locked ?

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

No → return (error);

End

```
static vdev_t *
spa_vdev_resilver_done_
hunt(vdev_t *vd)
```

* Find any device that's done replacing,
or a vdev marked 'unspare' that's
* current spared, so we can detach it.

```
vdev_t *newvd, *oldvd;
int c = 0
```

c < vd->vdev_children ?

* Check for a completed replacement.  We always consider the first
* vdev in the list to be the oldest vdev, and the last one to be
* the newest (see spa_vdev_attach() for how that works).  In
* the case where the newest vdev is faulted, we will not automatically
* remove it after a resilver completes.  This is OK as it will require
* user intervention to determine which disk the admin wishes to keep.

Yes

No

```
oldvd = spa_vdev_
resilver_done_hunt(vd-
>vdev_child[c]);
```

vd->vdev_ops ==
&vdev_replacing_ops ?

Yes

oldvd != NULL ?

No

```
ASSERT(vd->vdev_children > 1);
newvd = vd->vdev_child[vd->vdev_children - 1];
oldvd = vd->vdev_child[0];
```

Yes

c++

* Check for a completed
resilver with the
'unspare' flag set.

vdev_dtl_empty(newvd,
DTL_MISSING) &&
vdev_dtl_empty(newvd,
DTL_OUTAGE) &&
!vdev_dtl_required(oldvd)
?

No

No

return (oldvd);

vd->vdev_ops ==
&vdev_spare_ops ?

Yes

return (oldvd);

```
vdev_t *first = vd->vdev_child[0];
vdev_t *last = vd->vdev_child[vd->vdev_children - 1];
```

No

last->vdev_unspare ?

No

Yes

first->vdev_unspare ?

No

Yes

```
oldvd = first;
newvd = last;
```

```
oldvd = NULL;
```

```
oldvd = last;
newvd = first;
```

oldvd != NULL &&
vdev_dtl_empty(newvd,
DTL_MISSING) &&
vdev_dtl_empty(newvd,
DTL_OUTAGE) &&
!vdev_dtl_required(oldvd)
?

* If there are more than two spares attached to a disk,
* and those spares are not required, then we want to
* attempt to free them up now so that they can be used
* by other pools.  Once we're back down to a single
* disk+spare, we stop removing them.

No

Yes

return (oldvd);

vd->vdev_children > 2 ?

Yes

```
newvd = vd-
>vdev_child[1];
```

No

newvd->vdev_isspare &&
last->vdev_isspare &&
vdev_dtl_empty(last,
DTL_MISSING) &&
vdev_dtl_empty(last,
DTL_OUTAGE) &&
!vdev_dtl_required(newvd)
?

No

Yes

return (newvd);

return (NULL);

End

```
                         static void
                     spa_vdev_resilver_
                       done(spa_t *spa)


              vdev_t *vd, *pvd, *ppvd;
              uint64_t guid, sguid, pguid, ppguid;
              spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);


                    (vd = spa_vdev_resilver_
                     done_hunt(spa->spa_root_
                         vdev)) != NULL ?

        No                      │ Yes

              pvd = vd->vdev_parent;
              ppvd = pvd->vdev_parent;
              guid = vd->vdev_guid;            * If we have just finished replacing a hot spared device, then
              pguid = pvd->vdev_guid;          * we need to detach the parent's first child (the original hot
              ppguid = ppvd->vdev_guid;        * spare) as well.
              sguid = 0;

   spa_config_exit(spa,
   SCL_ALL, FTAG);

                         ppvd->vdev_ops ==
                         &vdev_spare_ops &&
                         pvd->vdev_id == 0 &&
                         ppvd->vdev_children
                            == 2 ?

                                      │ Yes

                         ASSERT(pvd->vdev_ops == &vdev_replacing_ops);
        No               sguid = ppvd->vdev_child[1]->vdev_guid;

                         spa_config_exit(spa,
                         SCL_ALL, FTAG);

                         spa_vdev_detach(spa,
                         guid, pguid,
                         B_TRUE) != 0 ?

                    Yes                    No

                                              sguid &&
                                              spa_vdev_detach(spa,
                                              sguid, ppguid,
                                              B_TRUE) != 0 ?

                                         Yes              No

                    End                        spa_config_enter(spa,
                                               SCL_ALL, FTAG,
                                               RW_WRITER);
```

```
int spa_vdev_set_common
(spa_t *spa, uint64_t
guid, const char *value,
boolean_t ispath)
```

* Update the stored path
or FRU for this vdev.

```
vdev_t *vd;
boolean_t sync = B_FALSE;
ASSERT(spa_writeable(spa));
spa_vdev_state_enter(spa, SCL_ALL);
```

(vd = spa_lookup_by_guid
(spa, guid,
B_TRUE)) == NULL ?

No

Yes

!vd->vdev_ops-
>vdev_op_leaf ?

No

Yes

return (spa_vdev_state_
exit(spa, NULL, ENOENT));

ispath ?

Yes

No

return (spa_vdev_state_
exit(spa, NULL,
ENOTSUP));

vd->vdev_fru == NULL ?

Yes

No

strcmp(value,
vd->vdev_path) != 0 ?

Yes

No

strcmp(value,
vd->vdev_fru) != 0 ?

No

Yes

```
vd->vdev_fru = spa_strdup(value);
sync = B_TRUE;
```

```
spa_strfree(vd->vdev_path);
vd->vdev_path = spa_strdup(value);
sync = B_TRUE;
```

```
spa_strfree(vd->vdev_fru);
vd->vdev_fru = spa_strdup(value);
sync = B_TRUE;
```

return (spa_vdev_state_
exit(spa, sync ?
vd : NULL, 0));

End

int spa_vdev_setpath(spa_t *spa, uint64_t guid, const char *newpath)

return (spa_vdev_set_common(spa, guid, newpath, B_TRUE));

End

```
int spa_vdev_setfru(spa_t
    *spa, uint64_t guid,
   const char *newfru)
```

```
return (spa_vdev_set_
   common(spa, guid,
   newfru, B_FALSE));
```

End

```
int spa_scan_stop(spa_t
         *spa)
```

```
* ====================================================================
* SPA Scanning
* ====================================================================
```

```
ASSERT(spa_config_held
      (spa, SCL_ALL,
    RW_WRITER) == 0);
```

```
dsl_scan_resilvering(spa-
      >spa_dsl_pool) ?
```

No      Yes

```
return (dsl_scan_cancel
   (spa->spa_dsl_pool));
```

```
return (EBUSY);
```

End

```
int spa_scan(spa_t *spa,
pool_scan_func_t func)
```

```
ASSERT(spa_config_held
(spa, SCL_ALL,
RW_WRITER) == 0);
```

```
func >= POOL_SCAN_FUNCS
|| func ==
POOL_SCAN_NONE ?
```

* If a resilver was requested, but there is no DTL on a
* writeable leaf device, we have nothing to do.

No

Yes

```
func == POOL_SCAN_
RESILVER &&
!vdev_resilver_needed(spa
->spa_root_vdev,
NULL, NULL) ?
```

Yes

No

```
return (ENOTSUP);
```

```
return (dsl_scan(spa-
>spa_dsl_pool, func));
```

```
spa_async_request(spa, SPA_ASYNC_RESILVER_DONE);
return (0);
```

End

```
                          static void
                     spa_async_remove(spa_t
                       *spa, vdev_t *vd)
```

* ===========================================================================
* SPA async task processing
* ===========================================================================

```
                    vd->vdev_remove_wanted ?
```

No          Yes

```
vd->vdev_remove_wanted = B_FALSE;
vd->vdev_delayed_close = B_FALSE;
vdev_set_state(vd, B_FALSE, VDEV_STATE_REMOVED, VDEV_AUX_NONE);
```

* 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.

```
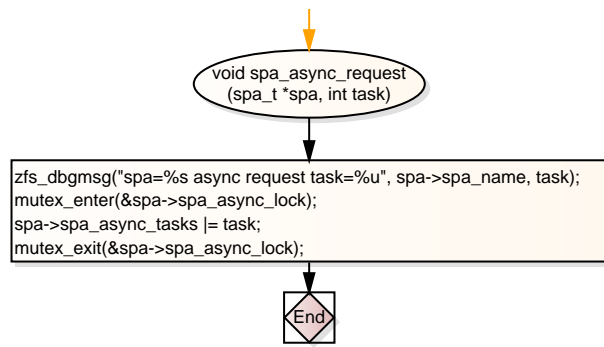vd->vdev_stat.vs_read_errors = 0;
vd->vdev_stat.vs_write_errors = 0;
vd->vdev_stat.vs_checksum_errors = 0;
vdev_state_dirty(vd->vdev_top);
```

```
int c = 0
```

```
c < vd->vdev_children ?
```

No          Yes

```
End
```

```
spa_async_remove(spa,
  vd->vdev_child[c]);
```

```
c++
```

```
                  static void
              spa_async_probe(spa_t
               *spa, vdev_t *vd)


            vd->vdev_probe_wanted ?
         No                  Yes

                        vd->vdev_probe_wanted
                          = B_FALSE;              vdev_open() does
                                                  the actual probe

                        vdev_reopen(vd);

            int c = 0


            c < vd->vdev_children ?

         No            Yes

        End      spa_async_probe(spa,
                 vd->vdev_child[c]);


                              c++
```

```
              static void
        spa_async_autoexpand(spa_
          t *spa, vdev_t *vd)
```

```
sysevent_id_t eid;
nvlist_t *attr;
char *physpath;
```

!spa->spa_autoexpand ?

— No → `int c = 0`

Yes

`c < vd->vdev_children ?`

— No → !vd->vdev_ops->vdev_op_leaf || vd->vdev_physpath == NULL ?

— Yes → `vdev_t *cvd = vd->vdev_child[c];`
`spa_async_autoexpand(spa, cvd);`

`c++`

!vd->vdev_ops->vdev_op_leaf || vd->vdev_physpath == NULL ?

— Yes → End

— No →

```
physpath = kmem_zalloc(MAXPATHLEN, KM_SLEEP);
(void) snprintf(physpath, MAXPATHLEN, "/devices%s", vd->vdev_physpath);
VERIFY(nvlist_alloc(&attr, NV_UNIQUE_NAME, KM_SLEEP) == 0);
VERIFY(nvlist_add_string(attr, DEV_PHYS_PATH, physpath) == 0);
(void) ddi_log_sysevent(zfs_dip, SUNW_VENDOR, EC_DEV_STATUS, ESC_DEV_DLE, attr, &eid, DDI_SLEEP);
nvlist_free(attr);
kmem_free(physpath, MAXPATHLEN);
```

End

```
                              static void
                           spa_async_thread(spa_t
                                  *spa)

                           int tasks;
                           ASSERT(spa->spa_sync_on);           * See if the config
                           mutex_enter(&spa->spa_async_lock);   needs to be updated.
                           tasks = spa->spa_async_tasks;
                           spa->spa_async_tasks = 0;
                           mutex_exit(&spa->spa_async_lock);

                                tasks & SPA_ASYNC_
                                CONFIG_UPDATE ?

                 Yes                                              No

  uint64_t old_space, new_space;
  mutex_enter(&spa_namespace_lock);
  old_space = metaslab_class_get_space(spa_normal_class(spa));   * If the pool grew as a result of the config update,
  spa_config_update(spa, SPA_CONFIG_UPDATE_POOL);                * then log an internal history event.
  new_space = metaslab_class_get_space(spa_normal_class(spa));
  mutex_exit(&spa_namespace_lock);

                           new_space != old_space ?

                   Yes

                        spa_history_log_internal
                        (spa, "vdev online",
                        NULL, "pool '%s' size:      * See if any devices
                             %llu(+%llu)",           need to be
                        spa_name(spa),               marked REMOVED.
                        new_space, new_space
                        - old_space);
            No

                                tasks & SPA_ASYNC_REMOVE
                                         ?

                 No                              Yes

                             spa_vdev_state_enter(spa, SCL_NONE);
                             spa_async_remove(spa, spa->spa_root_vdev);
                             int i = 0

                                   i < spa->spa_
                                   l2cache.sav_count ?

                            Yes                        No

                       spa_async_remove(spa,          int i = 0
                       spa->spa_l2cache.sav_
                       vdevs[i]);

                                                        i < spa->spa_spares
                             i++                        .sav_count ?

                                                No              Yes

                                                          spa_async_remove(spa,
                       (void) spa_vdev_state_             spa->spa_spares.sav_
                       exit(spa, NULL, 0);               vdevs[i]);

                                                              i++
                                (tasks &
                             SPA_ASYNC_AUTOEXPAND) &&
                             !spa_suspended(spa) ?

                 Yes                              No

  spa_config_enter(spa, SCL_CONFIG, FTAG, RW_READER);    * See if any devices
  spa_async_autoexpand(spa, spa->spa_root_vdev);          need to be probed.
  spa_config_exit(spa, SCL_CONFIG, FTAG);

                                tasks & SPA_ASYNC_PROBE ?

                 Yes                              No

  spa_vdev_state_enter(spa, SCL_NONE);                  * If any devices are
  spa_async_probe(spa, spa->spa_root_vdev);              done replacing,
  (void) spa_vdev_state_exit(spa, NULL, 0);              detach them.

                                tasks & SPA_ASYNC_
                                RESILVER_DONE ?

                 Yes                              No

                       spa_vdev_resilver_
                       done(spa);                        * Kick off a resilver.

                                tasks & SPA_
                                ASYNC_RESILVER ?

                 Yes                              No

                       dsl_resilver_restart(spa-        * Let the world know
                       >spa_dsl_pool, 0);                that we're done.

                           mutex_enter(&spa->spa_async_lock);
                           spa->spa_async_thread = NULL;
                           cv_broadcast(&spa->spa_async_cv);
                           mutex_exit(&spa->spa_async_lock);
                           thread_exit();

                                    End
```

```
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);
```

```
cv_wait(&spa->spa_async_
cv, &spa->spa_
async_lock);
```

End

```
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) ?
```

No

Yes

```
              spa->spa_async_thread =
              thread_create(NULL, 0,
              spa_async_thread, spa,
              0, &p0, TS_RUN,
              maxclsyspri);

              mutex_exit(&spa-
              >spa_async_lock);
```

End

```
void spa_async_request
(spa_t *spa, int task)
```

```c
zfs_dbgmsg("spa=%s async request task=%u", spa->spa_name, task);
mutex_enter(&spa->spa_async_lock);
spa->spa_async_tasks |= task;
mutex_exit(&spa->spa_async_lock);
```

End

```
static int
bpobj_enqueue_cb(void
*arg, const blkptr_t
*bp, dmu_tx_t *tx)
```

* ======================================================================
* SPA syncing routines
* ======================================================================

```
bpobj_t *bpo = arg;
bpobj_enqueue(bpo, bp, tx);
return (0);
```

End

```
static int
spa_free_sync_cb(void
*arg, const blkptr_t
*bp, dmu_tx_t *tx)
```

```
zio_t *zio = arg;
zio_nowait(zio_free_sync(zio, zio->io_spa,
dmu_tx_get_txg(tx), bp, zio->io_flags));
return (0);
```

End

```
static void
spa_sync_nvlist(spa_t
*spa, uint64_t obj,
nvlist_t *nv,
dmu_tx_t *tx)
```

```
char *packed = NULL;
size_t bufsize;
size_t nvsize = 0;
dmu_buf_t *db;
VERIFY(nvlist_size(nv,
&nvsize, NV_ENCODE_XDR) == 0);
```

* Write full (SPA_CONFIG_BLOCKSIZE) blocks of configuration
* information.  This avoids the dbuf_will_dirty() path and
* saves us a pre-read to get data we don't actually care about.

```
bufsize = P2ROUNDUP((uint64_t)nvsize, SPA_CONFIG_BLOCKSIZE);
packed = kmem_alloc(bufsize, KM_SLEEP);
VERIFY(nvlist_pack(nv, &packed, &nvsize, NV_ENCODE_XDR, KM_SLEEP) == 0);
bzero(packed + nvsize, bufsize - nvsize);
dmu_write(spa->spa_meta_objset, obj, 0, bufsize, packed, tx);
kmem_free(packed, bufsize);
VERIFY(0 == dmu_bonus_hold(spa->spa_meta_objset, obj, FTAG, &db));
dmu_buf_will_dirty(db, tx);
*(uint64_t *)db->db_data = nvsize;
dmu_buf_rele(db, FTAG);
```

End

```
static void
spa_sync_aux_dev(spa_t
*spa, spa_aux_vdev_t
*sav, dmu_tx_t *tx,
const char *config,
const char *entry)
```

```
nvlist_t *nvroot;
nvlist_t **list;
int i;
```

!sav->sav_sync ?

* Update the MOS nvlist describing the list of available devices.
* spa_validate_aux() will have already made sure this nvlist is
* valid and the vdevs are labeled appropriately.

No → sav->sav_object == 0 ?

Yes

Yes →
```
sav->sav_object = dmu_object_alloc(spa->spa_meta_objset,
DMU_OT_PACKED_NVLIST, 1 << 14,
DMU_OT_PACKED_NVLIST_SIZE, sizeof (uint64_t), tx);
VERIFY(zap_update(spa->spa_meta_objset,
DMU_POOL_DIRECTORY_OBJECT, entry, sizeof
(uint64_t), 1, &sav->sav_object, tx) == 0);
```

No

```
VERIFY(nvlist_alloc(
&nvroot, NV_UNIQUE_NAME,
KM_SLEEP) == 0);
```

sav->sav_count == 0 ?

Yes

No →
```
list = kmem_alloc(sav->sav_count * sizeof (void *), KM_SLEEP);
i = 0
```

```
VERIFY(nvlist_add_nvlist_
array(nvroot, config,
NULL, 0) == 0);
```

i < sav->sav_count ?

No

Yes →
```
list[i] =
vdev_config_generate(spa,
sav->sav_vdevs[i],
B_FALSE,
VDEV_CONFIG_L2CACHE);
```

```
VERIFY(nvlist_add_nvlist_array(nvroot,
config, list, sav->sav_count) == 0);
i = 0
```

i++

i < sav->sav_count ?

No

Yes →
nvlist_free(list[i]);

```
kmem_free(list,
sav->sav_count *
sizeof (void *));
```

i++

```
spa_sync_nvlist(spa, sav->sav_object, nvroot, tx);
nvlist_free(nvroot);
sav->sav_sync = B_FALSE;
```

End

```
            static void
         spa_sync_config_object
              (spa_t *spa,
             dmu_tx_t *tx)


         nvlist_t *config;


     list_is_empty(&spa->spa_
        config_dirty_list) ?
                                          No

                                  spa_config_enter(spa, SCL_STATE, FTAG, RW_READER);      * If we're upgrading the spa version then make sure that
                                  config = spa_config_generate(spa,                       * the config object gets updated with the correct version.
                                  spa->spa_root_vdev, dmu_tx_get_txg(tx), B_FALSE);

  Yes

                                           spa->spa_ubsync.ub_
                                                version <
                                            spa->spa_uberblock
                                               .ub_version ?

                                       No                    Yes

                                                      fnvlist_add_uint64
                                                           (config,
                                                    ZPOOL_CONFIG_VERSION,
                                                      spa->spa_uberblock
                                                          .ub_version);

                                   spa_config_exit(spa,
                                   SCL_STATE, FTAG);


                              spa->spa_config_syncing ?

                                                    Yes

                                         No        nvlist_free(spa->spa_
                                                      config_syncing);

                    spa->spa_config_syncing = config;
                    spa_sync_nvlist(spa, spa->spa_config_object, config, tx);


                              End
```

```
static void
spa_sync_version(void
*arg1, void *arg2,
dmu_tx_t *tx)
```

```
spa_t *spa = arg1;
uint64_t version = *(uint64_t *)arg2;
```

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

```
ASSERT(tx->tx_txg != 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);
```

End

```
                          static void
                          spa_sync_props(void
                          *arg1, void *arg2,
                          dmu_tx_t *tx)
```

* Set zpool properties.

```
spa_t *spa = arg1;
objset_t *mos = spa->spa_meta_objset;
nvlist_t *nvp = arg2;
nvpair_t *elem = NULL;
mutex_enter(&spa->spa_props_lock);
```

(elem = nvlist_next_nvpair(nvp, elem) )

No → mutex_exit(&spa->spa_props_lock);

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_INVAL ?

Yes → 
* We checked this earlier in spa_prop_validate().

```
ASSERT(zpool_prop_feature(nvpair_name(elem)));
fname = strchr(nvpair_name(elem), '@') + 1;
VERIFY3U(0, ==, zfeature_lookup_name(fname, &feature));
spa_feature_enable(spa, feature, tx);
spa_history_log_internal(spa, "set", tx, "%s=enabled", nvpair_name(elem));
```

No

ZPOOL_PROP_VERSION ?

Yes →
```
VERIFY(nvpair_value_
uint64(elem,
&intval) == 0);
```

* The version is synced seperatly before other
* properties and should be correct by now.

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

No

ZPOOL_PROP_ALTROOT ?

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

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

No

ZPOOL_PROP_READONLY ?

Yes

No

ZPOOL_PROP_CACHEFILE ?

Yes

No

* 'readonly' and 'cachefile' are also non-persistent
* properties.

ZPOOL_PROP_COMMENT ?

Yes →
```
VERIFY(nvpair_value_
string(elem,
&strval) == 0);
```

```
spa->spa_comment
!= NULL ?
```
Yes →
```
spa_strfree(spa-
>spa_comment);
```

No

```
spa->spa_comment =
spa_strdup(strval);
```

* We need to dirty the configuration on all the vdevs
* so that their labels get updated.  It's unnecessary
* to do this for pool creation since the vdev's
* configuratoin has already been dirtied.

```
tx->tx_txg !=
TXG_INITIAL ?
```
Yes →
```
vdev_config_dirty(spa-
>spa_root_vdev);
```

No

```
spa_history_log_internal
(spa, "set", tx,
"%s=%s",
nvpair_name(elem),
strval);
```

default

* Set pool property values in the poolprops mos object.

```
spa->spa_pool_props_
object == 0 ?
```
Yes →
```
spa->spa_pool_props_
object =
zap_create_link(mos,
DMU_OT_POOL_PROPS,
DMU_POOL_DIRECTORY_
OBJECT, DMU_POOL_PROPS,
tx);
```

No → normalize the property name

```
propname = zpool_prop_to_name(prop);
proptype = zpool_prop_get_type(prop);
```

```
nvpair_type(elem) ==
DATA_TYPE_STRING ?
```
Yes →
```
ASSERT(proptype == PROP_TYPE_STRING);
VERIFY(nvpair_value_string(elem, &strval) == 0);
VERIFY(zap_update(mos, spa->spa_pool_props_object,
propname, 1, strlen(strval) + 1, strval, tx) == 0);
spa_history_log_internal(spa, "set", tx, "%s=%s", nvpair_name(elem), strval);
```

No

```
nvpair_type(elem) ==
DATA_TYPE_UINT64 ?
```
Yes →
```
VERIFY(nvpair_value_
uint64(elem,
&intval) == 0);
```

No → not allowed → ASSERT(0);

```
proptype ==
PROP_TYPE_INDEX ?
```
Yes →
```
const char *unused;
VERIFY(zpool_prop_index_to_string( prop, intval, &unused) == 0);
```

No

```
VERIFY(zap_update(mos, spa->spa_pool_props_object,
propname, 8, 1, &intval, tx) == 0);
spa_history_log_internal(spa, "set", tx,
"%s=%lld", nvpair_name(elem), intval);
```

switch (prop)

ZPOOL_PROP_DELEGATION ?

Yes →
```
spa->spa_delegation
= intval;
```

No

ZPOOL_PROP_BOOTFS ?

Yes → spa->spa_bootfs = intval;

No

ZPOOL_PROP_FAILUREMODE ?

Yes →
```
spa->spa_failmode
= intval;
```

No

ZPOOL_PROP_AUTOEXPAND ?

Yes →
```
spa->spa_autoexpand
= intval;
```

No

ZPOOL_PROP_DEDUPDITTO ?

Yes

default →
```
spa->spa_dedup_ditto
= intval;
```

```
tx->tx_txg !=
TXG_INITIAL ?
```
Yes →
```
spa_async_request(spa,
SPA_ASYNC_AUTOEXPAND);
```

No
```

```
                    static void
              spa_sync_upgrades(spa_t
                *spa, dmu_tx_t *tx)
```

* 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);
```

spa->spa_ubsync.ub_
version <
SPA_VERSION_ORIGIN &&
spa->spa_uberblock.ub_
version >=
SPA_VERSION_ORIGIN ?

Yes

```
dsl_pool_create_origin
        (dp, tx);
```

No

Keeping the origin open
increases spa_minref

```
spa->spa_minref += 3;
```

spa->spa_ubsync.ub_
version <
SPA_VERSION_NEXT_CLONES
&& spa->spa_uberblock.ub_
version >=
SPA_VERSION_NEXT_CLONES ?

Yes

No

```
dsl_pool_upgrade_
   clones(dp, tx);
```

spa->spa_ubsync.ub_
version <
SPA_VERSION_DIR_CLONES
&& spa->spa_uberblock.ub_
version >=
SPA_VERSION_DIR_CLONES ?

Yes

```
dsl_pool_upgrade_
  dir_clones(dp, tx);
```

No

Keeping the freedir open
increases spa_minref

```
spa->spa_minref += 3;
```

spa->spa_ubsync.ub_
version <
SPA_VERSION_FEATURES &&
spa->spa_uberblock.ub_
version >=
SPA_VERSION_FEATURES ?

Yes

No

```
spa_feature_create_zap_
    objects(spa, tx);
```

End
```

*flowchart diagram*

spa_sync_config_object(spa, tx)

ASSERT(!list_empty(&dp->dp_dirty_datasets, txg))
ASSERT(!list_empty(&dp->dp_dirty_dirs, txg))
ASSERT(!txg_list_empty(&spa->spa_vdev_txg_list, txg))
spa_sync_rewrite_vdev_config(spa, tx)
vdev_config_sync(rvd, svdcount, txg)
spa_handle_ignored_writes(spa)

```
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 ?

**No**

**Yes**   **Yes**

```
mutex_exit(
&spa_namespace_lock);
```

spa_state(spa) !=
POOL_STATE_ACTIVE ||
!spa_writeable(spa) ||
spa_suspended(spa) ?

**No**

End

```
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);
```

```
                              ┌──────────────────────┐
                              │ void spa_evict_all(void) │
                              └──────────────────────┘
                                        │
                                        ▼
  * =============================================================================
  * Miscellaneous routines
  * =============================================================================
  * Remove all pools in the system.

                              ┌──────────────┐       * Remove all cached state.  All pools should be closed now,
                              │ spa_t *spa;  │       * so every spa in the AVL tree should be unreferenced.
                              └──────────────┘
                                        │
                                        ▼
                              ┌──────────────────────┐
                              │ mutex_enter(          │
                              │ &spa_namespace_lock); │
                              └──────────────────────┘
                                        │
                                        ▼
  * Stop async tasks.  The async thread may need to detach         ◇ (spa = spa_next(NULL))
  * a device that's been replaced, which requires grabbing            != NULL ?
  * spa_namespace_lock, so we must drop it here.
                                             Yes │                    │ No
                                                 ▼                    ▼
              ┌──────────────────────────────────┐      ┌──────────────────────┐
              │ spa_open_ref(spa, FTAG);          │      │ mutex_exit(           │
              │ mutex_exit(&spa_namespace_lock);  │      │ &spa_namespace_lock); │
              │ spa_async_suspend(spa);           │      └──────────────────────┘
              │ mutex_enter(&spa_namespace_lock); │                 │
              │ spa_close(spa, FTAG);             │                 ▼
              └──────────────────────────────────┘            ◇ End ◇
                                        │
                                        ▼
                    ◇ spa->spa_state !=
                      POOL_STATE_UNINITIALIZED
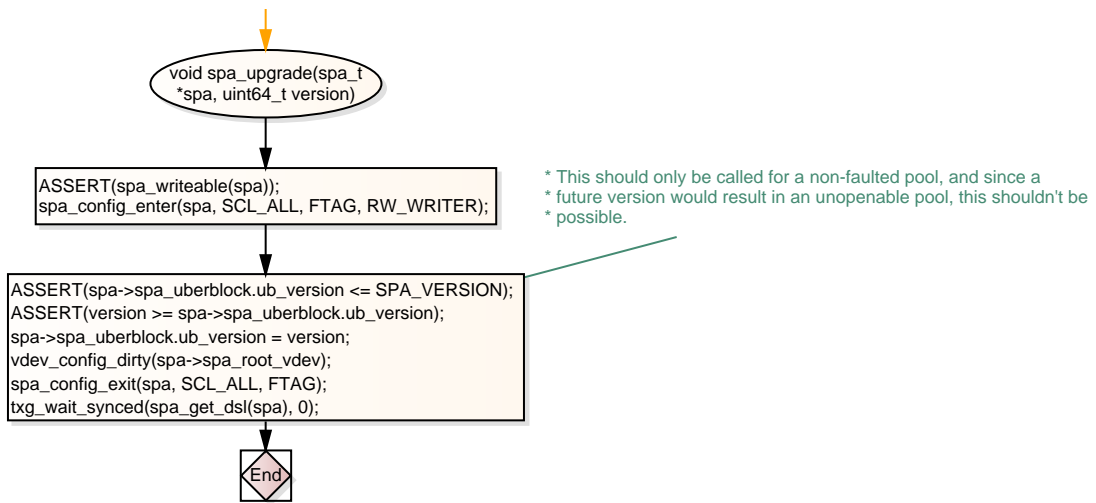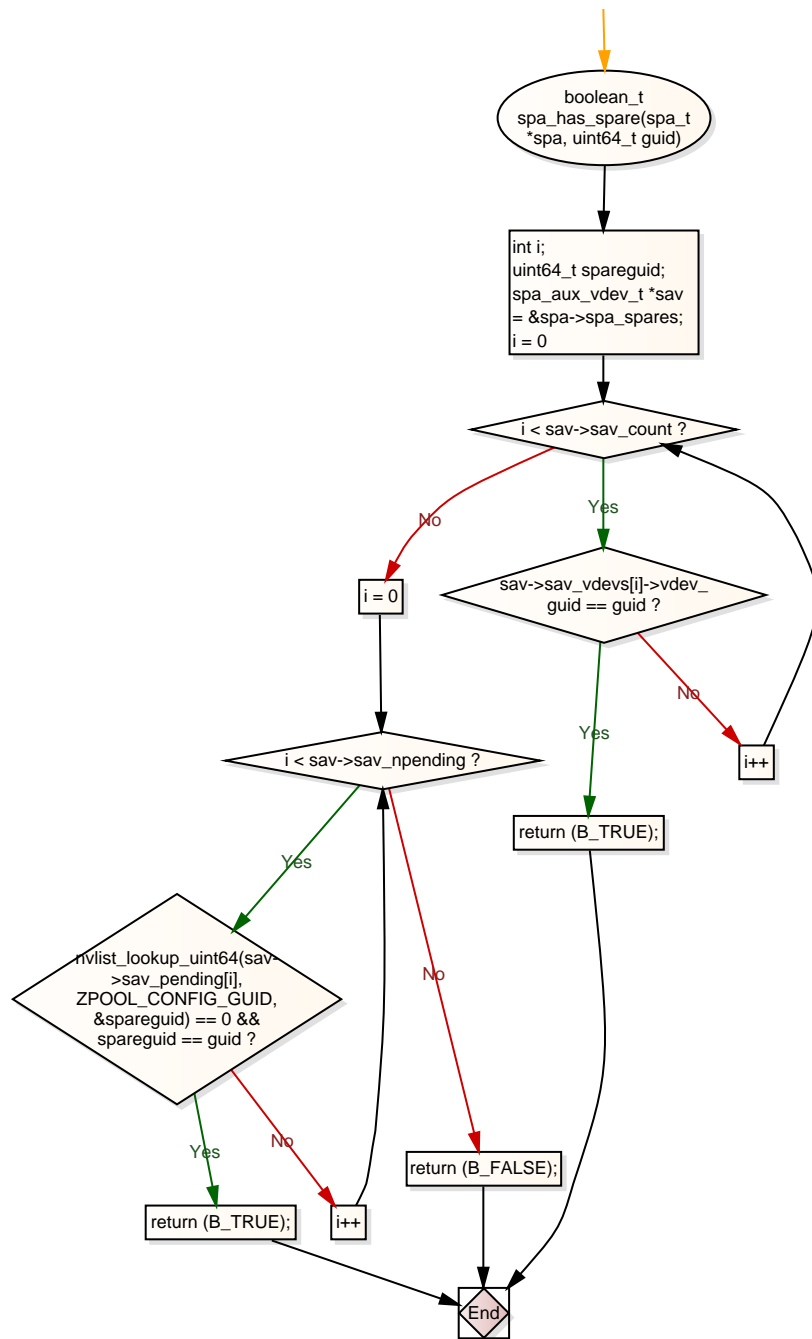                                 ?
                  Yes │                     │ No
                      ▼                     │
         ┌──────────────────────┐           │
         │ spa_unload(spa);      │           │
         │ spa_deactivate(spa);  │           │
         └──────────────────────┘           │
                      │                      ▼
                      └──────────► ┌──────────────────┐
                                   │ spa_remove(spa); │
                                   └──────────────────┘
```

```
vdev_t *
spa_lookup_by_guid(spa_t
*spa, uint64_t guid,
boolean_t aux)
```

```
vdev_t *vd;
int i;
```

```
(vd = vdev_lookup_by_guid
(spa->spa_root_vdev,
guid)) != NULL ?
```

No — aux ?

Yes — return (vd);

aux ?

Yes — i = 0

No — return (NULL);

```
i < spa->spa_
l2cache.sav_count ?
```

Yes —
```
vd = spa->spa_l2cache
.sav_vdevs[i];
```

No — i = 0

```
vd->vdev_guid == guid ?
```

Yes — return (vd);

No — i++

```
i < spa->spa_spares
.sav_count ?
```

Yes —
```
vd = spa->spa_spares
.sav_vdevs[i];
```

No — return (NULL);

```
vd->vdev_guid == guid ?
```

Yes — return (vd);

No — i++

End

```
        void spa_upgrade(spa_t
        *spa, uint64_t version)
```

```
ASSERT(spa_writeable(spa));
spa_config_enter(spa, SCL_ALL, FTAG, RW_WRITER);
```

* This should only be called for a non-faulted pool, and since a
* future version would result in an unopenable pool, this shouldn't be
* possible.

```
ASSERT(spa->spa_uberblock.ub_version <= SPA_VERSION);
ASSERT(version >= spa->spa_uberblock.ub_version);
spa->spa_uberblock.ub_version = version;
vdev_config_dirty(spa->spa_root_vdev);
spa_config_exit(spa, SCL_ALL, FTAG);
txg_wait_synced(spa_get_dsl(spa), 0);
```

End

```
boolean_t
spa_has_spare(spa_t
*spa, uint64_t guid)
```

```
int i;
uint64_t spareguid;
spa_aux_vdev_t *sav
= &spa->spa_spares;
i = 0
```

i < sav->sav_count ?

No → i = 0

Yes → sav->sav_vdevs[i]->vdev_guid == guid ?

No → i++

Yes → return (B_TRUE);

i < sav->sav_npending ?

Yes → nvlist_lookup_uint64(sav->sav_pending[i], ZPOOL_CONFIG_GUID, &spareguid) == 0 && spareguid == guid ?

No → return (B_FALSE);

Yes → return (B_TRUE);

No → i++

End

```
static boolean_t
spa_has_active_shared_
spare(spa_t *spa)
```

* Check if a pool has an active shared spare device.
* Note: reference count of an active
spare is 2, as a spare and as a replace

```
int i, refcnt;
uint64_t pool;
spa_aux_vdev_t *sav
= &spa->spa_spares;
i = 0
```

i < sav->sav_count ?

Yes

```
spa_spare_exists(sav-
>sav_vdevs[i]->vdev_guid,
&pool, &refcnt) && pool
!= 0ULL && pool ==
spa_guid(spa)
&& refcnt > 2 ?
```

No

return (B_FALSE);

Yes

return (B_TRUE);

No

i++

End

void spa_event_notify
(spa_t *spa, vdev_t *vd,
const char *name)

* 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.

#ifdef _KERNEL ?

No

Yes

sysevent_t *ev;
sysevent_attr_list_t *attr = NULL;
sysevent_value_t value;
sysevent_id_t eid;
ev = sysevent_alloc(EC_ZFS, (char
*)name, SUNW_KERN_PUB "zfs", SE_SLEEP);
value.value_type = SE_DATA_TYPE_STRING;
value.value.sv_string = spa_name(spa);

sysevent_add_attr(&attr,
ZFS_EV_POOL_NAME,
&value, SE_SLEEP) != 0 ?

No

Yes

value.value_type = SE_DATA_TYPE_UINT64;
value.value.sv_uint64 = spa_guid(spa);

sysevent_add_attr(&attr,
ZFS_EV_POOL_GUID,
&value, SE_SLEEP) != 0 ?

No

Yes

vd ?

No

Yes

value.value_type = SE_DATA_TYPE_UINT64;
value.value.sv_uint64 = vd->vdev_guid;

sysevent_add_attr(&attr,
ZFS_EV_VDEV_GUID,
&value, SE_SLEEP) != 0 ?

No

Yes

vd->vdev_path ?

Yes

value.value_type = SE_DATA_TYPE_STRING;
value.value.sv_string = vd->vdev_path;

sysevent_add_attr(&attr,
ZFS_EV_VDEV_PATH,
&value, SE_SLEEP) != 0 ?

No

Yes

No

sysevent_attach_
attributes(ev,
attr) != 0 ?

No

Yes

attr = NULL;
(void) log_sysevent(ev, SE_SLEEP, &eid);

done: attr ?

No

Yes

sysevent_free_attr(attr);

sysevent_free(ev);

End