Adding reference counters (krefs) to kernel objects

Author: Corey Minyard < minyard@acm.org>
Author: Thomas Hellstrom@vmware.com>

A lot of this was lifted from Greg Kroah-Hartman's 2004 OLS paper and presentation on krefs, which can be found at:

- http://www.kroah.com/linux/talks/ols 2004 kref paper/Reprint-Kroah-Hartman-OLS2004.pdf
- http://www.kroah.com/linux/talks/ols 2004 kref talk/

Introduction

krefs allow you to add reference counters to your objects. If you have objects that are used in multiple places and passed around, and you don't have refcounts, your code is almost certainly broken. If you want refcounts, krefs are the way to go.

To use a kref, add one to your data structures like:

```
struct my_data
{
    .
    .
    struct kref refcount;
    .
};
```

The kref can occur anywhere within the data structure.

Initialization

You must initialize the kref after you allocate it. To do this, call kref init as so:

```
struct my_data *data;
data = kmalloc(sizeof(*data), GFP_KERNEL);
if (!data)
        return -ENOMEM;
kref_init(&data->refcount);
```

This sets the refcount in the kref to 1.

Kref rules

Once you have an initialized kref, you must follow the following rules:

1. If you make a non-temporary copy of a pointer, especially if it can be passed to another thread of execution, you must increment the refcount with kref get() before passing it off:

```
kref_get(&data->refcount);
```

If you already have a valid pointer to a kref-ed structure (the refcount cannot go to zero) you may do this without a lock.

2. When you are done with a pointer, you must call kref put():

```
kref put(&data->refcount, data release);
```

If this is the last reference to the pointer, the release routine will be called. If the code never tries to get a valid pointer to a kref-ed structure without already holding a valid pointer, it is safe to do this without a lock.

3. If the code attempts to gain a reference to a kref-ed structure without already holding a valid pointer, it must serialize access where a kref_put() cannot occur during the kref_get(), and the structure must remain valid during the kref_get().

For example, if you allocate some data and then pass it to another thread to process:

```
void data_release(struct kref *ref)
{
    struct my_data *data = container_of(ref, struct my_data, refcount);
    kfree(data);
}

void more_data_handling(void *cb_data)
{
    struct my_data *data = cb_data;
    .
    . do stuff with data here
```

```
kref put(&data->refcount, data release);
}
int my_data_handler(void)
{
   int rv = 0;
    struct my data *data;
    struct task struct *task;
    data = kmalloc(sizeof(*data), GFP KERNEL);
    if (!data)
           return -ENOMEM;
    kref init(&data->refcount);
    kref get(&data->refcount);
    task = kthread run(more data handling, data, "more data handling");
    if (task == ERR PTR(-ENOMEM)) {
           rv = -ENOMEM:
           kref put(&data->refcount, data release);
           goto out;
    }
    . do stuff with data here
011†:
    kref put(&data->refcount, data release);
    return rv:
```

This way, it doesn't matter what order the two threads handle the data, the kref_put() handles knowing when the data is not referenced any more and releasing it. The kref_get() does not require a lock, since we already have a valid pointer that we own a refeount for. The put needs no lock because nothing tries to get the data without already holding a pointer.

In the above example, kref_put() will be called 2 times in both success and error paths. This is necessary because the reference count got incremented 2 times by kref_init() and kref_get().

Note that the "before" in rule 1 is very important. You should never do something like:

```
task = kthread_run(more_data_handling, data, "more_data_handling");
if (task == ERR_PTR(-ENOMEM)) {
    rv = -ENOMEM;
    goto out;
} else
    /* BAD BAD BAD - get is after the handoff */
    kref get(&data->refcount);
```

Don't assume you know what you are doing and use the above construct. First of all, you may not know what you are doing. Second, you may know what you are doing (there are some situations where locking is involved where the above may be legal) but someone else who doesn't know what they are doing may change the code or copy the code. It's bad style. Don't do it.

There are some situations where you can optimize the gets and puts. For instance, if you are done with an object and enqueuing it for something else or passing it off to something else, there is no reason to do a get then a put:

```
/* Silly extra get and put */
kref_get(&obj->ref);
enqueue(obj);
kref_put(&obj->ref, obj_cleanup);
```

Just do the enqueue. A comment about this is always welcome:

```
enqueue(obj);
/* We are done with obj, so we pass our refcount off
  to the queue. DON'T TOUCH obj AFTER HERE! */
```

The last rule (rule 3) is the nastiest one to handle. Say, for instance, you have a list of items that are each kref-ed, and you wish to get the first one. You can't just pull the first item off the list and kref_get() it. That violates rule 3 because you are not already holding a valid pointer. You must add a mutex (or some other lock). For instance:

```
static DEFINE_MUTEX(mutex);
static LIST_HEAD(q);
struct my_data
{
        struct kref refcount;
        struct list_head link;
};
static struct my_data *get_entry()
{
        struct my_data *entry = NULL;
        mutex lock(&mutex);
```

The kref_put() return value is useful if you do not want to hold the lock during the whole release operation. Say you didn't want to call kfree() with the lock held in the example above (since it is kind of pointless to do so). You could use kref_put() as follows:

This is really more useful if you have to call other routines as part of the free operations that could take a long time or might claim the same lock. Note that doing everything in the release routine is still preferred as it is a little neater.

The above example could also be optimized using kref_get_unless_zero() in the following way:

```
static struct my_data *get_entry()
{
        struct my_data *entry = NULL;
        mutex lock(&mutex);
        if (!list empty(&q)) {
                entry = container of(q.next, struct my data, link);
                if (!kref get unless zero(&entry->refcount))
                        entry = NULL;
        mutex_unlock(&mutex);
        return entry;
}
static void release entry(struct kref *ref)
        struct my data *entry = container of(ref, struct my data, refcount);
        mutex lock(&mutex);
        list del(&entry->link);
        mutex unlock(&mutex);
        kfree (entry);
}
static void put_entry(struct my_data *entry)
{
        kref put(&entry->refcount, release entry);
```

Which is useful to remove the mutex lock around kref_put() in put_entry(), but it's important that kref_get_unless_zero is enclosed in the same critical section that finds the entry in the lookup table, otherwise kref_get_unless_zero may reference already freed memory. Note that it is illegal to use kref_get_unless_zero without checking its return value. If you are sure (by already having a valid pointer) that kref_get_unless_zero() will return true, then use kref_get() instead.

17....C. 1 DAT

Kreis and KUU

The function kref get unless zero also makes it possible to use rcu locking for lookups in the above example:

```
struct my_data
        struct rcu_head rhead;
        struct kref refcount;
};
static struct my_data *get_entry_rcu()
        struct my data *entry = NULL;
        rcu_read_lock();
        if (!list empty(&q)) {
                entry = container_of(q.next, struct my_data, link);
                if (!kref get unless zero(&entry->refcount))
                        entry = NULL;
        }
        rcu read unlock();
        return entry;
static void release entry rcu(struct kref *ref)
{
        struct my_data *entry = container_of(ref, struct my_data, refcount);
       mutex_lock(&mutex);
        list del rcu(&entry->link);
       mutex_unlock(&mutex);
        kfree rcu(entry, rhead);
}
static void put entry(struct my data *entry)
        kref put(&entry->refcount, release entry rcu);
```

But note that the struct kref member needs to remain in valid memory for a rcu grace period after release_entry_rcu was called. That can be accomplished by using kfree_rcu(entry, rhead) as done above, or by calling synchronize_rcu() before using kfree, but note that synchronize_rcu() may sleep for a substantial amount of time.