

"Denis V. Lunev" <den@sw.ru> writes:

> Hello, Eric!  
>  
> Current locking in mainstream seems broken to me.

Thanks. After looking at this I concur.

> 1. struct net->list is manipulated under double net\_mutex/net\_list\_mutex

Yes. Making iteration safe if we hold only one of those.

> 2. net\_list\_mutex has been taken only in cleanup\_net/net\_ns\_init inside  
> net\_mutes and seems pointless now

And in rtnl\_unlock (although that isn't upstream just yet).  
It looks like I forgot to call net\_lock in some of my later  
insertions of for\_each\_net.

Certainly it looks like too many locks.

Thinking.

net\_mutex appears to be there to serial the addition/removal of  
subsystems/modules and the creation/destruction of network namespaces.

net\_list\_mutex is just there to serialize operations on the list of  
namespaces.

I'm trying to see if there is something that implies a nesting of:  
net\_mutex, rtnl, net\_list\_mutex.

Although it is no longer an issue now that I am making fewer locks  
per network namespace.

I am remembering that there was something keeping from using the rtnl.

> 3. for\_each\_net (iterating against net\_namespace\_list) is called from  
> a) register\_netdevice\_notifier/\_\_rtnl\_link\_unregister

Yes this is fishy, and probably needs to be fixed.

> b) register\_pernet\_operations/unregister\_pernet\_operations  
> In the case b) the situation is sane, i.e. net\_mutex is held while in

> the case b) we held rtnl\_only  
>  
> So, this does not look good to me for now.  
> How to cure this situation? I think that we can drop all locks for now  
> and perform all operations under rtnl only. In the other case we must  
> decide now should we make rtnl inner or outer for net\_mutex.

Ok. I have found an important case. loopback.

We must hold net\_mutex when we are calling all of the .init routines.  
The loopback code calls register\_netdev which grabs rtnl.

- So we have net\_mutex must be outside of rtnl.

We have to do for\_each\_net in rtnl\_unlock so we can find all of the  
rtnl netlink sockets and sk\_data\_ready aka rtnetlink\_rcv which takes  
the rtnl\_lock.

- So net\_list\_lock should be taken outside of rtnl\_lock.

We take net\_list\_mutex in rtnl\_unlock() but not under rtnl\_mutex. And  
rtnl\_unlock is called inside of net\_mutex, so we can't use net\_mutex.

- So we need both net\_list\_lock and net\_mutex.

Therefore it looks like we just need to take net\_lock() outside of  
rtnl\_lock() in register\_netdevice\_notifier.

>>From my point of view net\_mutex should be taken inside rtnl lock and we  
> must add it now into list manipulation routines.

I think that is where I started and I failed miserably. The per  
network namespace instances of the rtnl socket look to make that  
impossible.

> Plz point me to my mistake in logic :)

Does what I said sound reasonable now.

Thanks for spotting the missing lock by the way.

You want to cook up the patch to fix register\_netdevice\_notifier?

Eric

---