

# LNetCtl Utility Improvements

Target release	
Epic	
Document status	DRAFT
Document owner	Amir Shehata
Designer	
Developers	
QA	

## Overview

The DLC project added the ability to dynamically configure LNet. It also added the ability to pull information from LNet and display them in YAML format. The decision was later made to make `lnetctl` the default utility for all LNet operations. This will entail rolling in all LNet related functionality from `lctl` to `lnetctl`.

This document will outline the following:

1. The `sysfs` interface design for the parameters which will be stored in `sysfs` and reported by `lnetctl`
2. The `ioctl` interface design for the operations which will be implemented via `lnetctl`

## Design Principles

There has been several discussions on the mechanics of pulling information from the kernel. What has been agreed on is to use `sysfs` for simple key/value pairs. Primary examples for that are module parameters and statistics. For module parameters the hooks are already in place to add the parameters in `sysfs`. Namely under: `/sys/module/lnet/parameters`. Statistics also fall in the category of simple key/value pairs. These can also be stored under `sysfs`. The current suggested path is: `/sys/module/lnet/statistics/`. The exact hierarchy will be outlined later in the requirements.

This information can then be pulled out of `sysfs` via `lnetctl`. According to an article posted by GregKH and written by Kay Sievers, <https://wn.net/Articles/237664/>, it is not suggested to use `libsysfs` for accessing `sysfs` from userspace, since it doesn't follow the rules he mentions in that article. However, it will be ideal to use an existent library, if one exists, rather than create our own `sysfs` interface. However if one doesn't exist, then the possibility of creating one for lustre should be investigated. The suggestion here is not to have each utility do it's own `sysfs` access, but rather consolidate this in a library which can be utilized by all Lustre/LNet utilities.

For more complex configuration operations, such as configuring network, network interfaces, peers, etc, `IOCTL` will still be used to invoke these operations. `sysfs` is not geared to handle complex information passing to the kernel. And the `IOCTL` mechanism is already in place to handle more complex kernel operations.

## Requirements

### `sysfs` Structure

ID	Status	Class	Description
sysfs-01		MUST HAVE	Create a statistics directory under <code>/sys/modules/lnet</code> <code>/sys/modules/lnet/statistics</code>

sysfs-02		MUST HAVE	<p>Create a directory for each Luster Network type</p> <p>ex:</p> <pre> /sys/modules/lnet/statistics/loLnd/ /sys/modules/lnet/statistics/o2ib/ /sys/modules/lnet/statistics/tcp/ /sys/modules/lnet/statistics/gni/ </pre>
sysfs-03		MUST_HAVE	<p>Create a directory for each separate lustre network under the network type</p> <p>ex:</p> <pre> /sys/modules/lnet/statistics/o2ib/1 </pre>
sysfs-04		MUST HAVE	<p>Create a directory for each configured Network interface under the Lustre Network</p> <p>ex:</p> <pre> /sys/modules/lnet/statistics/o2ib/[NUM]/ib[NUM]/ /sys/modules/lnet/statistics/tcp/[NUM]/eth[NUM]/ </pre>
sysfs-05		MUST HAVE	<p>Place the Network Interface statistics as attribute files under the Network interface directory</p> <p>ex: send_count, recv_count, ...</p>
sysfs-06		MUST HAVE	<p>Create a peer directory under statistics.</p> <pre> /sys/modules/lnet/statistics/peers/ </pre>
sysfs-07		MUST HAVE	<p>Create a separate entry per peer named with the peer's primary_nid</p> <p>ex:</p> <pre> /sys/modules/lnet/statistics/peers/192.168.23.3@o2ib/ </pre>
sysfs-08		MUST HAVE	<p>Create a separate entry per peer network interface (peer_ni) named with the NID of the peer_ni</p> <p>ex:</p> <pre> /sys/modules/lnet/statistics/peers/192.168.23.3@o2ib/192.168.23.4@o2ib2 </pre>
sysfs-09		MUST HAVE	<p>Add peer statistics, each in its own attribute file, under the peer_ni directory.</p> <p>ex:</p> <pre> state, max_ni_tx_credits, etc </pre>
sysfs-10		MUST HAVE	<p>Create a directory for statistics under the ko2iblnd module</p> <pre> /sys/modules/ko2iblnd/statistics </pre>
sysfs-11		MUST HAVE	<p>Create a directory for peers under the statistics directory</p> <pre> /sys/module/ko2iblnd/statistics/peers/ </pre>
sysfs-12		MUST HAVE	<p>Create a directory for each peer_ni under the /sys/module/ko2iblnd/statistics/peers/. Under each peer_ni directory create a directory for each local_ni which has connections established with the peer_ni. This is needed because there could be multiple LND peers for a single LNet peer. One local_ni can talk to multiple peer_nis.</p> <p>ex:</p> <pre> /sys/module/ko2iblnd/statistics/peers/peer_NID/local_NID/ </pre>
sysfs-13		MUST HAVE	<p>Create a link from the LNet peer_ni directory to the specific LND peer-ni directory, which contains all instances of LND peer-nis associated with that LNet peer.</p> <p>ex:</p> <pre> /sys/modules/lnet/statistics/peers/192.168.23.3@o2ib/192.168.23.4@o2ib2/LND_peers/ </pre>

sysfs-14		MUST HAVE	<p>Add LND peer_ni statistics each in a separate attribute file under the LND peer_ni directory as defined above.</p> <p>The statistics to be added are as follows:</p> <ul style="list-style-type: none"> <li>• number of active connections</li> <li>• number of messages waiting for a connection to be established</li> <li>• number passive connects</li> <li>• number active connects</li> <li>• number of times connection races were triggered</li> <li>• the last time the peer was alive</li> <li>• the number of users for this peer (refcount)</li> </ul>
sysfs-15		MUST HAVE	<p>Add a directory under each LND peer NI director that will contain all the connections established to that peer.</p> <p>ex:</p> <pre>/sys/modules/ko2ibnd/statistics/peers/peer_ni/local_ni/conns</pre>
sysfs-16		MUST HAVE	<p>Add a directory per connection under the connections directory. The name shall be "&lt;num&gt;"</p> <p>ex:</p> <pre>/sys/modules/ko2ibnd/statistics/peers/peer_ni/local_ni/conns/1/</pre>
sysfs-17		MUST HAVE	<p>Add all the connections stats as attribute files under the relevant &lt;num&gt; directory</p> <p>The statistics to be added are as follows:</p> <ul style="list-style-type: none"> <li>• the number of users for this connection (<code>ibc_refcount</code>)</li> <li>• the state of the connection (<code>ibc_state</code>)</li> <li>• number of incomplete sends (<code>ibc_nsend_posted</code>)</li> <li>• number of incomplete no-ops (<code>ibc_noops_posted</code>)</li> <li>• number of available credits (<code>ibc_credits</code>)</li> <li>• number of credits to return (<code>ibc_outstanding_credits</code>)</li> <li>• number of reserved credits (<code>ibc_reserved_credits</code>)</li> <li>• last communication error on that connection (<code>ibc_comms_error</code>) <ul style="list-style-type: none"> <li>• Note if the connection is immediately closed on error then this value will probably never be seen.</li> </ul> </li> <li>• number of queued messages of type IMMEDIATE (<code>ibc_tx_queue</code>)</li> <li>• number of queued messages not requiring credits of type PUT_NAK, PUT_ACK, PUT_DONE, GET_DONE. (<code>ibc_tx_queue_nocred</code>)</li> <li>• number of queued messages of type PUT_REQ, GET_REQ. (<code>ibc_tx_queue_rsrvd</code>)</li> <li>• number of transmits waiting for completion (<code>ibc_active_txs</code>)</li> </ul>
sysfs-18		MUST HAVE	<p>provide a reset attribute file under <code>/sys/modules/lnd/statistics/</code>. When anything is written to it, then all statistics are reset.</p>
sysfs-19		MUST HAVE	<p>provide a reset attribute file under <code>/sys/modules/ko2ibnd/statistics/</code>. When anything is written to it then all statistics are reset.</p>
sysfs-20		MUST HAVE	<p>provide a reset attribute file under <code>/sys/modules/socklnd/statistics/</code>. When anything is written to it then all statistics are reset.</p>

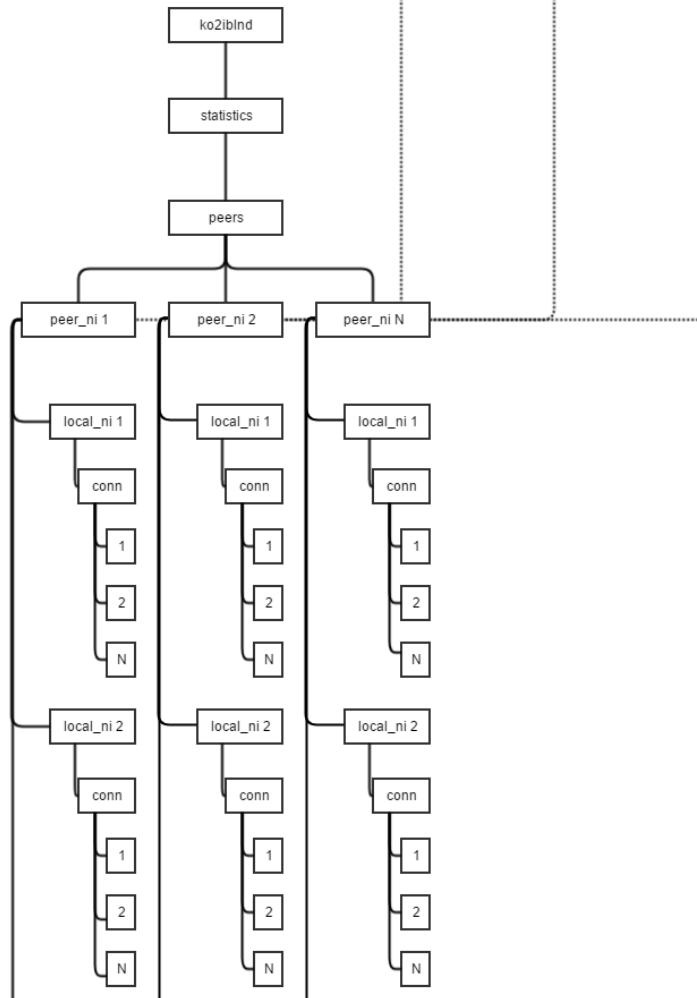
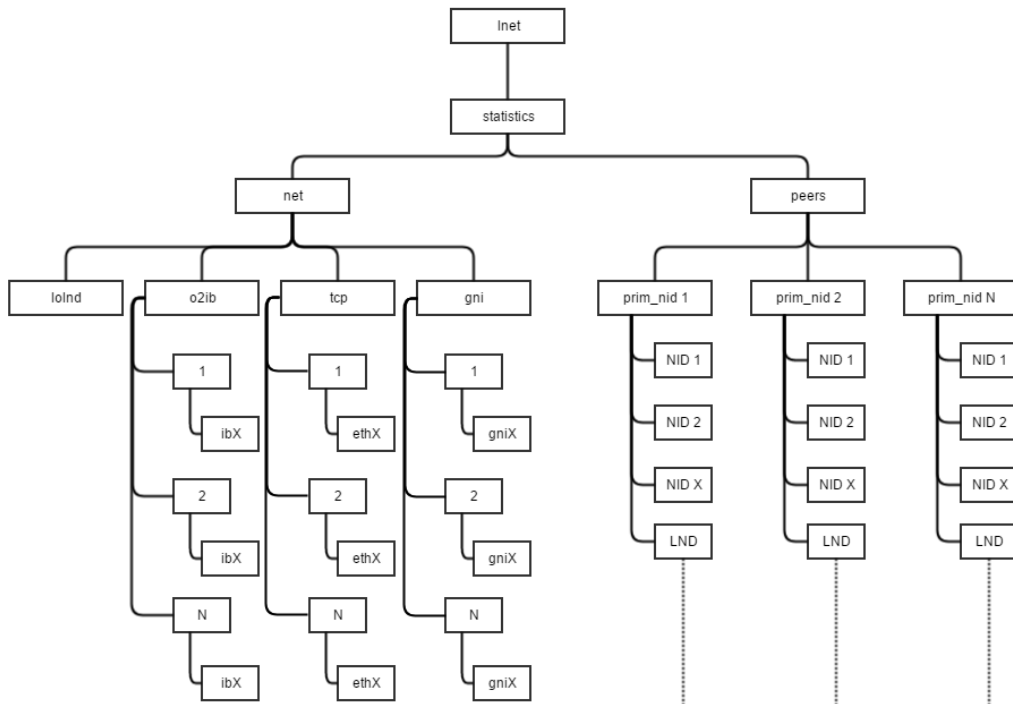
## Inetctl Requirements

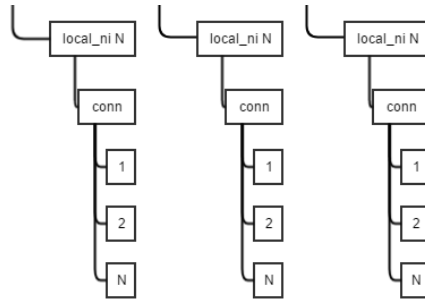
ID	Status	Class	Description
ctl-01		MUST HAVE	Inetctl shall traverse the <code>/sys/lnd/statistics</code> and <code>/sys/&lt;LND&gt;/statistics</code> directory and gather the stats and displays them in YAML format

ctl-02	MUST HAVE	<p>Inetctl shall display the peer stats using the following hierarchy:</p> <pre> peer: - primary nid: 192.168.1.1@o2ib Multi-Rail: True peer ni: - nid: 192.168.1.1@o2ib &lt;list of stats&gt; LND peer ni: - ni: &lt;local ni&gt; &lt;list of stats&gt; connections: - connection: &lt;num&gt; &lt;list of stats&gt; </pre>
ctl-03	MUST HAVE	<p>Inetctl shall show the LND peer only when a the --lnd_peer option is provided. the LND peer information will be displayed as an addition as represented the hierarchy described above</p> <p>ex:</p> <pre>lnetctl peer show --lnd_peer</pre>
ctl-04	MUST HAVE	<p>Inetctl shall allow peers to be filtered on their primary NID</p> <p>ex:</p> <pre>lnetctl peer show --prim_nid &lt;nid&gt;</pre>
ctl-05	MUST HAVE	<p>Inetctl shall allow peers to be filtered on their networks. The semantics are "Show all peers on &lt;network&gt;."</p> <p>ex:</p> <pre>lnetctl peer show --net &lt;net&gt;</pre>
ctl-06	MUST HAVE	<p>Inetctl shall use the function <code>libcfs_str2nid()</code> to resolve a string to a NID. This function is able to resolve an IP address or a hostname into a NID, for socklnd and o2iblnd networks. It can also handle gni networks</p>
ctl-07	MUST HAVE	<p>Inetctl shall show the LND connections per peer only when the --connections option is provided. The LND connection information will be displayed per peer as represented in the hierarchy above.</p> <p>By providing the --connection option the --lnd_peer option is activated automatically as it is necessary to display the connection information per lnd peer.</p> <p>ex:</p> <pre>lnetctl peer show --prim_nid &lt;nid&gt; --connections</pre>
ctl-08	MUST HAVE	<p>LNet shall extend the del peer command to shutdown all active connection to a peer.</p>
ctl-09	NICE-TO-HAVE	<p>LNet shall allow the termination of a single connection to a peer in the case when multiple connections are opened to the same peer.</p> <pre>lnetctl peer del connection --nid &lt;nid&gt; --id &lt;connection id as displayed in the show&gt;</pre>
ctl-10	MUST HAVE	<p>Inetctl shall provide a command to reset all statistics</p> <pre>lnetctl statistics reset</pre>

## sysfs structure

### Overview





## Example

For the sake of demonstrating the sysfs structure, lets assume the following system:

Node A with the following NIDs:

- 192.168.1.2@o2ib
- 192.168.1.3@o2ib
- 192.168.1.4@o2ib

Peer B with the following NIDs

- 192.168.2.2@o2ib
- 192.168.2.3@o2ib
- 172.168.2.2@o2ib1

Currently the following connections are established

- 2 connections between 192.168.1.2@o2ib and 192.168.2.2@o2ib
- 2 connections between 192.168.1.2@o2ib and 192.168.2.3@o2ib
- 1 connection between 192.168.1.3@o2ib and 192.168.2.2@o2ib
- 1 connection between 192.168.1.4@o2ib and 192.168.2.3@o2ib

The sysfs directory structure would be as follows:

```

/sys/modules/lnet/statistics/o2ib/
/sys/modules/lnet/statistics/o2ib/ib0/
/sys/modules/lnet/statistics/o2ib/ib0/send_count
/sys/modules/lnet/statistics/o2ib/ib0/rcv_count
/sys/modules/lnet/statistics/o2ib/ib0/drop_count
/sys/modules/lnet/statistics/o2ib/ib0/<other stats>
/sys/modules/lnet/statistics/o2ib/ib1/
/sys/modules/lnet/statistics/o2ib/ib1/send_count
/sys/modules/lnet/statistics/o2ib/ib1/rcv_count
/sys/modules/lnet/statistics/o2ib/ib1/drop_count
/sys/modules/lnet/statistics/o2ib/ib1/<other stats>
/sys/modules/lnet/statistics/o2ib/ib2/
/sys/modules/lnet/statistics/o2ib/ib2/send_count
/sys/modules/lnet/statistics/o2ib/ib2/rcv_count
/sys/modules/lnet/statistics/o2ib/ib2/drop_count
/sys/modules/lnet/statistics/o2ib/ib2/<other stats>

/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.2@o2ib/
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.2@o2ib/state
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.2@o2ib/max_n
i_tx_credits
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.2@o2ib/avail

```

```
able_tx_credits
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.2@o2ib/<other
stats>
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.3@o2ib/
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.3@o2ib/state
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.3@o2ib/max_n
i_tx_credits
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.3@o2ib/avail
able_tx_credits
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/192.168.2.3@o2ib/<other
stats>
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/172.168.2.2@o2ib/
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/172.168.2.2@o2ib/state
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/172.168.2.2@o2ib/max_n
i_tx_credits
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/172.168.2.2@o2ib/avail
able_tx_credits
/sys/modules/lnet/statistics/peers/192.168.2.2@o2ib/172.168.2.2@o2ib/<other
stats>

/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.2@o2ib/
/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.2@o2ib/tx
_queue_noop
/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.2@o2ib/tx
_queue_cr
/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.2@o2ib/<o
ther statistics>

/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.3@o2ib/
/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.3@o2ib/tx
_queue_noop
/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.3@o2ib/tx
_queue_cr
/sys/module/ko2iblnd/statistics/peers/192.168.2.2@o2ib/192.168.1.3@o2ib/<o
ther statistics>

/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.2@o2ib/
/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.2@o2ib/tx
_queue_noop
/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.2@o2ib/tx
_queue_cr
/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.2@o2ib/<o
ther statistics>

/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.4@o2ib/
/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.4@o2ib/tx
_queue_noop
/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.4@o2ib/tx
_queue_cr
```

```
/sys/module/ko2iblnd/statistics/peers/192.168.2.3@o2ib/192.168.1.4@o2ib/other statistics>
```

## Phases

All statistics are currently pulled from the kernel via the `ioctl` interface. As indicated above for simple key/value pairs the `sysfs` interface is preferred. However, the work will be divided in multiple phases or patches in order to concentrate more on adding new functionality rather than rewriting existing ones. These phases are equivalent to gerrit patches.

### Phase 1

Add the necessary infrastructure to create `sysfs` directories and files. This might just boil down to using the provided APIs, but we'd probably want to create `_sysfs.c` file that will contain all the necessary code to handle `sysfs` for `Inet`.

### Phase 2

Create the `lnd peer sysfs` structure and all associated stats.

Add the hooks necessary for `sysfs` to be updated when a read is performed on the exported attributes.

### Phase 3

Create the user space control to pull up and display the `lnd peer sysfs`

### Phase 4

Create the `lnd peer connections sysfs` structure and all associated stats.

Add the hooks necessary for `sysfs` to be updated when a read is performed on the exported attributes

### Phase 5

Create the user space control to pull up and display the `lnd peer connection sysfs stats`

### Phase 6

Create the `LNet peer sysfs` structure and associated stats.

### Phase 7

Update the user space control to traverse the `sysfs` directory tree rather than use `ioctl`.

### Phase 8

Create the `LNet network sysfs` structure and associated stats.

### Phase 9

Update the user space control to traverse the `sysfs` directory tree rather than use `ioctl`