# SRP - SCSI RDMA Protocol

# Table of contents

# List of Figures

The SCSI RDMA Protocol (SRP) is designed to take full advantage of the protocol offload and RDMA features provided by the InfiniBand architecture. SRP allows a large body of SCSI software to be readily used on InfiniBand architecture. The SRP Initiator controls the connection to an SRP Target in order to provide access to remote storage devices across an InfiniBand fabric. The kSRP Target resides in an IO unit and provides storage services.

# SRP Initiator

This SRP Initiator is based on open source from OpenFabrics (www.openfabrics.org) that implements the SCSI RDMA Protocol-2 (SRP-2). SRP-2 is described in Document # T10/1524-D available from http://www.t10.org.

The SRP Initiator supports

- Basic SCSI Primary Commands -3 (SPC-3)
  (www.t10.org/ftp/t10/drafts/spc3/spc3r21b.pdf)

- Basic SCSI Block Commands -2 (SBC-2)
  (www.t10.org/ftp/t10/drafts/sbc2/sbc2r16.pdf)

- Basic functionality, task management and limited error handling

> (i) **Note**
>
> This package, however, does not include an SRP Target.

# Loading SRP Initiator

➤ *To load the SRP module either:*

- Execute the `modprobe ib_srp` command after the OFED driver is up.

  or

1. Change the value of `SRP_LOAD` in `/etc/infiniband/openib.conf` to "yes".

2. Run /etc/init.d/openibd restart for the changes to take effect.

> (i) **Note**
>
> When loading the ib_srp module, it is possible to set the module parameter srp_sg_tablesize. This is the maximum number of gather/scatter entries per I/O (default: 12).

## SRP Module Parameters

When loading the SRP module, the following parameters can be set (viewable by the "modinfo ib_srp" command):

| | |
|---|---|
| cmd_sg_entries | Default number of gather/scatter entries in the SRP command (default is 12, max 255) |
| allow_ext_sg | Default behavior when there are more than cmd_sg_entries S/G entries after mapping; fails the request when false (default false) |
| topspin_workarounds | Enable workarounds for Topspin/Cisco SRP target bugs |
| reconnect_delay | Time between successive reconnect attempts. Time between successive reconnect attempts of SRP initiator to a disconnected target until dev_loss_tmo timer expires (if enabled), after that the SCSI target will be removed |
| fast_io_fail_tmo | Number of seconds between the observation of a transport layer error and failing all I/O. Increasing this timeout allows more tolerance to transport errors, however, doing so increases the total failover time in case of serious transport failure.<br>Note: fast_io_fail_tmo value must be smaller than the value of reconnect_delay |
| dev_loss_tmo | Maximum number of seconds that the SRP transport should insulate transport layer errors. After this time has been exceeded the SCSI target is |

removed. Normally it is advised to set this to -1 (disabled) which will never remove the scsi_host. In deployments where different SRP targets are connected and disconnected frequently, it may be required to enable this timeout in order to clean old scsi_hosts representing targets that no longer exists

Constraints between parameters:

- dev_loss_tmo, fast_io_fail_tmo, reconnect_delay cannot be all disabled or negative values.

- reconnect_delay must be positive number.

- fast_io_fail_tmo must be smaller than SCSI block device timeout.

- fast_io_fail_tmo must be smaller than dev_loss_tmo.

# SRP Remote Ports Parameters

Several SRP remote ports parameters are modifiable online on existing connection.

➢ *To modify dev_loss_tmo to 600 seconds:*

```
echo 600 > /sys/class/srp_remote_ports/port-xxx/dev_loss_tmo
```

➢ *To modify fast_io_fail_tmo to 15 seconds:*

```
echo 15 > /sys/class/srp_remote_ports/port-xxx/fast_io_fail_tmo
```

➢ *To modify reconnect_delay to 10 seconds:*

```
echo 20 > /sys/class/srp_remote_ports/port-xxx/reconnect_delay
```

# Manually Establishing an SRP Connection

The following steps describe how to manually load an SRP connection between the Initiator and an SRP Target. "Automatic Discovery and Connection to Targets" section explains how to do this automatically.

- Make sure that the ib_srp module is loaded, the SRP Initiator is reachable by the SRP Target, and that an SM is running.

- To establish a connection with an SRP Target and create an SRP (SCSI) device for that target under /dev, use the following command:

```
echo -n id_ext=[GUID value],ioc_guid=[GUID value],dgid=[port GID value],\
pkey=ffff,service_id=[service[0] value] > \
/sys/class/infiniband_srp/srp-mlx[hca number]-[port number]/add_target
```

See "SRP Tools - ibsrpdm, srp_daemon and srpd Service Script" section for instructions on how the parameters in this echo command may be obtained.

**Notes**:

- Execution of the above "echo" command may take some time

- The SM must be running while the command executes

- It is possible to include additional parameters in the echo command:

  - max_cmd_per_lun - Default: 62

  - max_sect (short for max_sectors) - sets the request size of a command

  - io_class - Default: 0x100 as in rev 16A of the specification (In rev 10 the default was 0xff00)

  - tl_retry_count - a number in the range 2..7 specifying the IB RC retry count. Default: 2

  - comp_vector, a number in the range 0..n-1 specifying the MSI-X completion vector. Some HCA's allocate multiple (n) MSI-X vectors per HCA port. If the IRQ affinity masks of these interrupts have been configured such that each MSI-X interrupt is handled by a different CPU then the comp_vector parameter can be used to spread the SRP completion workload over multiple CPU's.

- cmd_sg_entries, a number in the range 1..255 that specifies the maximum number of data buffer descriptors stored in the SRP_CMD information unit itself. With allow_ext_sg=0 the parameter cmd_sg_entries defines the maximum S/G list length for a single SRP_CMD, and commands whose S/G list length exceeds this limit after S/G list collapsing will fail.

- initiator_ext - see "Multiple Connections from Initiator InfiniBand Port to the Target" section.

- To list the new SCSI devices that have been added by the echo command, you may use either of the following two methods:

  - Execute "fdisk -l". This command lists all devices; the new devices are included in this listing.

  - Execute "dmesg" or look at /var/log/messages to find messages with the names of the new devices.

## SRP sysfs Parameters

Interface for making ib_srp connect to a new target. One can request ib_srp to connect to a new target by writing a comma-separated list of login parameters to this sysfs attribute. The supported parameters are:

| | |
|---|---|
| id_ext | A 16-digit hexadecimal number specifying the eight byte identifier extension in the 16-byte SRP target port identifier. The target port identifier is sent by ib_srp to the target in the SRP_LOGIN_REQ request. |
| ioc_guid | A 16-digit hexadecimal number specifying the eight byte I/O controller GUID portion of the 16-byte target port identifier. |
| dgid | A 32-digit hexadecimal number specifying the destination GID. |
| pkey | A four-digit hexadecimal number specifying the InfiniBand partition key. |
| service_id | A 16-digit hexadecimal number specifying the InfiniBand service ID used to establish communication with the SRP target. How to find out the value of the service ID is specified in the documentation of the SRP target. |
| max_sect | A decimal number specifying the maximum number of 512-byte sectors to be transferred via a single SCSI command. |

| | |
|---|---|
| max_cmd_per_lun | A decimal number specifying the maximum number of outstanding commands for a single LUN. |
| io_class | A hexadecimal number specifying the SRP I/O class. Must be either 0xff00 (rev 10) or 0x0100 (rev 16a). The I/O class defines the format of the SRP initiator and target port identifiers. |
| initiator_ext | A 16-digit hexadecimal number specifying the identifier extension portion of the SRP initiator port identifier. This data is sent by the initiator to the target in the SRP_LOGIN_REQ request. |
| cmd_sg_entries | A number in the range 1..255 that specifies the maximum number of data buffer descriptors stored in the SRP_CMD information unit itself. With allow_ext_sg=0 the parameter cmd_sg_entries defines the maxi- mum S/G list length for a single SRP_CMD, and commands whose S/G list length exceeds this limit after S/G list collapsing will fail. |
| allow_ext_sg | Whether ib_srp is allowed to include a partial memory descriptor list in an SRP_CMD instead of the entire list. If a partial memory descriptor list has been included in an SRP_CMD the remaining memory descriptors are communicated from initiator to target via an additional RDMA transfer. Setting allow_ext_sg to 1 increases the maximum amount of data that can be transferred between initiator and target via a single SCSI command. Since not all SRP target implementations support partial memory descriptor lists the default value for this option is 0. |
| sg_tablesize | A number in the range 1..2048 specifying the maximum S/G list length the SCSI layer is allowed to pass to ib_srp. Specifying a value that exceeds cmd_sg_entries is only safe with partial memory descriptor list support enabled (allow_ext_sg=1). |
| comp_vector | A number in the range 0..n-1 specifying the MSI-X completion vector. Some HCA's allocate multiple (n) MSI-X vectors per HCA port. If the IRQ affinity masks of these interrupts have been configured such that each MSI-X interrupt is handled by a different CPU then the comp_vector parameter can be used to spread the SRP completion workload over multiple CPU's. |
| tl_retry_count | A number in the range 2..7 specifying the IB RC retry count. |

# SRP Tools - ibsrpdm, srp_daemon and srpd Service Script

The OFED distribution provides two utilities: ibsrpdm and srp_daemon:

- They detect targets on the fabric reachable by the Initiator (Step 1)

- Output target attributes in a format suitable for use in the above "echo" command (Step 2)

- A service script srpd which may be started at stack startup

The utilities can be found under /usr/sbin/, and are part of the srptools RPM that may be installed using the OFED installation. Detailed information regarding the various options for these utilities are provided by their man pages.
Below, several usage scenarios for these utilities are presented.

## ibsrpdm

ibsrpdm has the following tasks:

1. Detecting reachable targets.

    1. To detect all targets reachable by the SRP initiator via the default umad device (/sys/class/infiniband_mad/umad0), execute the following command:

    ```
    ibsrpdm
    ```

    This command will result into readable output information on each SRP Target detected. Sample:

    ```
    IO Unit Info:
    port LID: 0103
    port GID: fe800000000000000002c90200402bd5
    change ID: 0002
    max controllers: 0x10
    controller[ 1]
    GUID: 0002c90200402bd4
    vendor ID: 0002c9
    device ID: 005a44
    IO class : 0100
    ```

> ID: LSI Storage Systems SRP Driver 200400a0b81146a1
>
> service entries: 1
>
> service[ 0]: 200400a0b81146a1 / SRP.T10:200400A0B81146A1

2. To detect all the SRP Targets reachable by the SRP Initiator via another umad device, use the following command:

```
ibsrpdm -d <umad device>
```

2. Assisting in SRP connection creation.

1. To generate an output suitable for utilization in the "echo" command in "Manually Establishing an SRP Connection" section, add the '-c' option to ibsrpdm:

```
ibsrpdm -c
```

Sample output:

```
id_ext=200400A0B81146A1,ioc_guid=0002c90200402bd4,
dgid=fe800000000000002c90200402bd5,pkey=ffff,service_id=200400a0
```

2. To establish a connection with an SRP Target using the output from the 'ibsrpdm -c' example above, execute the following command:

```
echo -n id_ext=200400A0B81146A1,ioc_guid=0002c90200402bd4,
dgid=fe800000000000002c90200402bd5,pkey=ffff,service_id=200400a0
> /sys/
class/infiniband_srp/srp-mlx5_0-1/add_target
```

The SRP connection should now be up; the newly created SCSI devices should appear in the listing obtained from the 'fdisk -l' command.

3. Discover reachable SRP Targets given an InfiniBand HCA name and port, rather than by just running /sys/class/infiniband_mad/umad<N> where <N> is a digit.

## srpd

The srpd service script allows automatic activation and termination of the srp_daemon utility on all system live InfiniBand ports.

## srp_daemon

srp_daemon utility is based on ibsrpdm and extends its functionality. In addition to the ibsrpdm functionality described above, srp_daemon can:

- Establish an SRP connection by itself (without the need to issue the "echo" command described in "Manually Establishing an SRP Connection" section)

- Continue running in background, detecting new targets and establishing SRP connections with them (daemon mode)

- Discover reachable SRP Targets given an infiniband HCA name and port, rather than just by /dev/umad<N> where <N> is a digit

- Enable High Availability operation (together with Device-Mapper Multipath)

- Have a configuration file that determines the targets to connect to:

1. srp_daemon commands equivalent to ibsrpdm:

> "srp_daemon -a -o" is equivalent to "ibsrpdm"
>
> "srp_daemon -c -a -o" is equivalent to "ibsrpdm -c"

   **Note**: These srp_daemon commands can behave differently than the equivalent ibsrpdm command when /etc/srp_daemon.conf is not empty.

2. srp_daemon extensions to ibsrpdm.

- To discover SRP Targets reachable from the HCA device <InfiniBand HCA name> and the port <port num>, (and to generate output suitable for 'echo'), execute:

```
host1# srp_daemon -c -a -o -i <InfiniBand HCA name> -p <port number>
```

**Note**: To obtain the list of InfiniBand HCA device names, you can either use the ibstat tool or run 'ls /sys/class/infiniband'.

- To both discover the SRP Targets and establish connections with them, just add the -e option to the above command.

- Executing srp_daemon over a port without the -a option will only display the reachable targets via the port and to which the initiator is not connected. If executing with the -e option it is better to omit -a.

- It is recommended to use the -n option. This option adds the initiator_ext to the connecting string (see "Multiple Connections from Initiator InfiniBand Port to the Target" section).

- srp_daemon has a configuration file that can be set, where the default is /etc/srp_daemon.conf. Use the -f to supply a different configuration file that configures the targets srp_daemon is allowed to connect to. The configuration file can also be used to set values for additional parameters (e.g., max_cmd_per_lun, max_sect).

- A continuous background (daemon) operation, providing an automatic ongoing detection and connection capability. See "Automatic Discovery and Connection to Targets" section.

## Automatic Discovery and Connection to Targets

- Make sure the ib_srp module is loaded, the SRP Initiator can reach an SRP Target, and that an SM is running.

- To connect to all the existing Targets in the fabric, run "srp_daemon -e -o". This utility will scan the fabric once, connect to every Target it detects, and then exit.

> ⓘ **Note**

> srp_daemon will follow the configuration it finds in /etc/srp_daemon.conf. Thus, it will ignore a target that is disallowed in the configuration file.

- To connect to all the existing Targets in the fabric and to connect to new targets that will join the fabric, execute srp_daemon -e. This utility continues to execute until it is either killed by the user or encounters connection errors (such as no SM in the fabric).

- To execute SRP daemon as a daemon on all the ports:

    - srp_daemon.sh (found under /usr/sbin/). srp_daemon.sh sends its log to /var/log/srp_daemon.log.

    - Start the srpd service script, run service srpd start

For the changes in openib.conf to take effect, run:

```
/etc/init.d/openibd restart
```

# Multiple Connections from Initiator InfiniBand Port to the Target

Some system configurations may need multiple SRP connections from the SRP Initiator to the same SRP Target: to the same Target IB port, or to different IB ports on the same Target HCA.
In case of a single Target IB port, i.e., SRP connections use the same path, the configuration is enabled using a different initiator_ext value for each SRP connection. The initiator_ext value is a 16-hexadecimal-digit value specified in the connection command.

Also in case of two physical connections (i.e., network paths) from a single initiator IB port to two different IB ports on the same Target HCA, there is need for a different initiator_ext value on each path. The conventions is to use the Target port GUID as the initiator_ext value for the relevant path.

If you use srp_daemon with -n flag, it automatically assigns initiator_ext values according to this convention. For example:

```
id_ext=200500A0B81146A1,ioc_guid=0002c90200402bec,\
dgid=fe80000000000000002c90200402bed,pkey=ffff,\
service_id=200500a0b81146a1,initiator_ext=ed2b400002c90200
```

**Notes**:

- It is recommended to use the -n flag for all srp_daemon invocations.

- ibsrpdm does not have a corresponding option.

- srp_daemon.sh always uses the -n option (whether invoked manually by the user, or automatically at startup by setting SRP_DAEMON_ENABLE to yes).

# High Availability (HA)

High Availability works using the Device-Mapper (DM) multipath and the SRP daemon. Each initiator is connected to the same target from several ports/HCAs. The DM multipath is responsible for joining together different paths to the same target and for failover between paths when one of them goes offline. Multipath will be executed on newly joined SCSI devices.
Each initiator should execute several instances of the SRP daemon, one for each port. At startup, each SRP daemon detects the SRP Targets in the fabric and sends requests to the ib_srp module to connect to each of them. These SRP daemons also detect targets that subsequently join the fabric, and send the ib_srp module requests to connect to them as well.

## Operation

When a path (from port1) to a target fails, the ib_srp module starts an error recovery process. If this process gets to the reset_host stage and there is no path to the target from this port, ib_srp will remove this scsi_host. After the scsi_host is removed, multipath switches to another path to this target (from another port/HCA).
When the failed path recovers, it will be detected by the SRP daemon. The SRP daemon will then request ib_srp to connect to this target. Once the connection is up, there will be

a new scsi_host for this target. Multipath will be executed on the devices of this host, returning to the original state (prior to the failed path).

## Manual Activation of High Availability

Initialization - execute after each boot of the driver:

1. Execute modprobe dm-multipath

2. Execute modprobe ib-srp

3. Make sure you have created file /etc/udev/rules.d/91-srp.rules as described above

4. Execute for each port and each HCA:

   ```
   srp_daemon -c -e -R 300 -i <InfiniBand HCA name> -p <port number>
   ```

   This step can be performed by executing srp_daemon.sh, which sends its log to /var/log/srp_daemon.log.

Now it is possible to access the SRP LUNs on /dev/mapper/.

> (i) **Note**
>
> It is possible for regular (non-SRP) LUNs to also be present; the SRP LUNs may be identified by their names. You can configure the /etc/multipath.conf file to change multipath behavior.

> (i) **Note**
>
> It is also possible that the SRP LUNs will not appear under /dev/mapper/. This can occur if the SRP LUNs are in the black-list of

multipath. Edit the 'blacklist' section in /etc/multipath.conf and make sure the SRP LUNs are not blacklisted.

## Automatic Activation of High Availability

- Start srpd service, run:

```
service srpd start
```

- From the next loading of the driver it will be possible to access the SRP LUNs on /dev/mapper/

> **ⓘ Note**
>
> It is possible that regular (not SRP) LUNs are also present. SRP LUNs may be identified by their name.

- It is possible to see the output of the SRP daemon in /var/log/srp_daemon.log

# Shutting Down SRP

SRP can be shutdown by using "rmmod ib_srp", or by stopping the OFED driver ("/etc/init.d/openibd stop"), or as a by-product of a complete system shutdown. Prior to shutting down SRP, remove all references to it. The actions you need to take depend on the way SRP was loaded. There are three cases:

1. Without High Availability
   When working without High Availability, you should unmount the SRP partitions that were mounted prior to shutting down SRP.

2. After Manual Activation of High Availability
   If you manually activated SRP High Availability, perform the following steps:

   1. Unmount all SRP partitions that were mounted.

   2. Stop service srpd (Kill the SRP daemon instances).

   3. Make sure there are no multipath instances running. If there are multiple instances, wait for them to end or kill them.

   4. Run: multipath -F

3. After Automatic Activation of High Availability
   If SRP High Availability was automatically activated, SRP shutdown must be part of the driver shutdown ("/etc/init.d/openibd stop") which performs Steps 2-4 of case b above. However, you still have to unmount all SRP partitions that were mounted before driver shutdown.