



IDT RapidIO Driver and Software Training

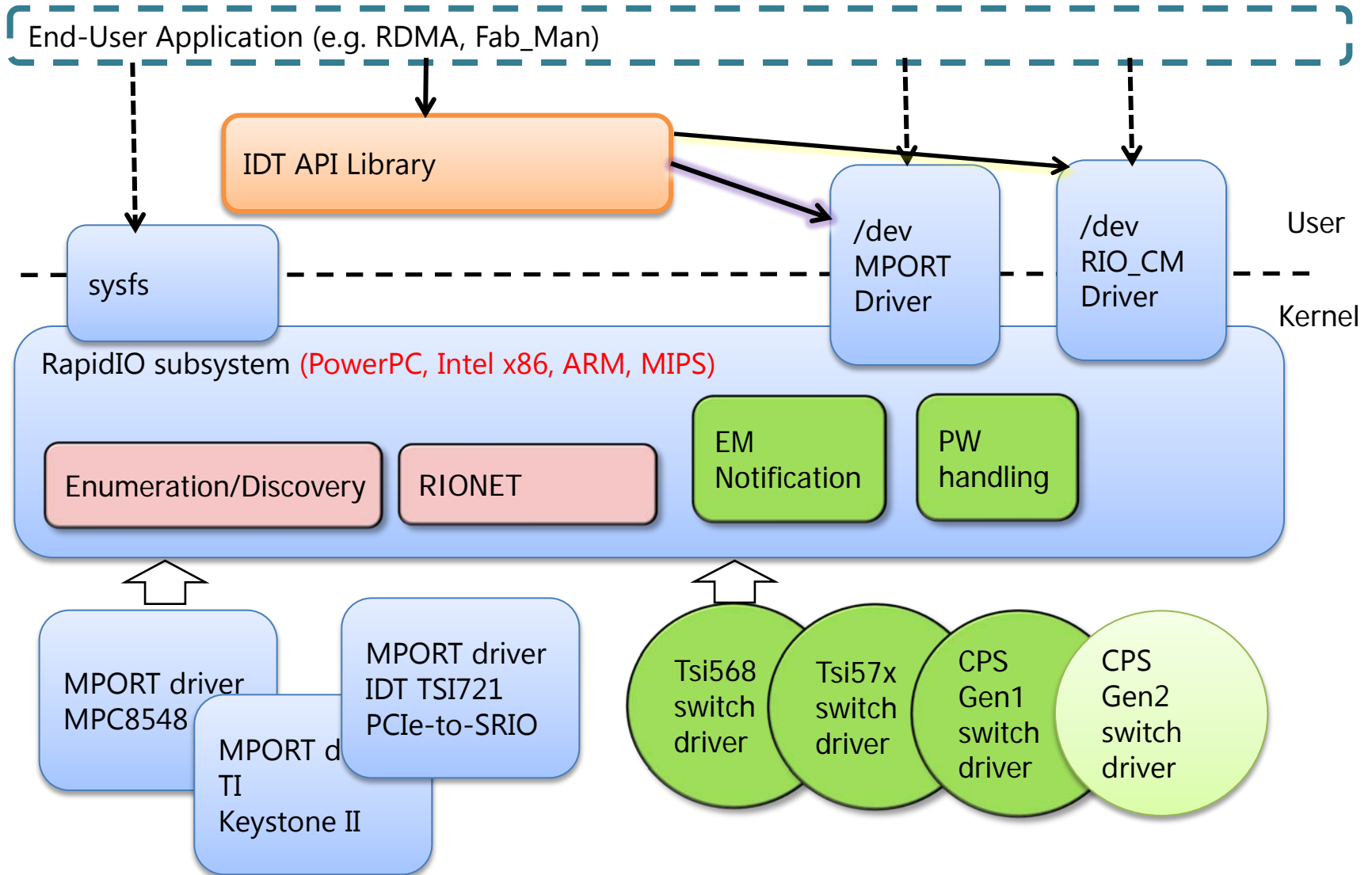
December 23, 2015

Integrated Device Technology, Inc.

Advanced Product and Technology Information
(subject to change)



CORE AND BELOW



- ❑ **Subsystem core** – defines hardware-independent kernel-space interfaces for use by other subsystem components (RIONET, enumeration, mport_cdev driver etc.). This layer allows programming in generic RapidIO terms without knowing details about HW that services requests.
- ❑ **MPORT device drivers** – hardware specific device drivers that provide set of callbacks defined by RapidIO subsystem. Hardware-specific MPORT device drivers are different from commonly known driver types: block and char. Consider them as plug-ins for the core.
- ❑ **RapidIO switch device drivers** – another subsystem/HW-specific elements with their own set of callback routines.

- ❑ For each MPORT device RIO core creates an object (defined as data structure `rio_mport` in `include/linux/rio.h`) which allows unified representation of different HW controllers.
- ❑ One of members of `rio_mport` structure is a pointer to `rio_ops` data structure (see `include /linux/rio.h`) which is created by each specific MPORT driver during registration with RIO core.
- ❑ Core interface functions that need to perform specific RapidIO operation call corresponding callback functions provided for given `rio_mport` object: code fragment from `drivers/rapidio/rio-access.c`

```
int rio_mport_read_config_##size
(struct rio_mport *mport, u16 destid, u8 hopcount, u32 offset, type *value) \
{
    int res;
    unsigned long flags;
    u32 data = 0;
    if (RIO_##size##_BAD) return RIO_BAD_SIZE;
    spin_lock_irqsave(&rio_config_lock, flags);
    res = mport->ops->cread(mport, mport->id, destid, hopcount, offset, len, &data); \
    *value = (type)data;
    spin_unlock_irqrestore(&rio_config_lock, flags);
    return res;
}
```

```

238 * @node: Node in global list of master ports
239 * @nnode: Node in network list of master ports
240 * @iores: I/O mem resource that this master port interface owns
241 * @riores: RIO resources that this master port interfaces owns
242 * @inb_msg: RIO inbound message event descriptors
243 * @outb_msg: RIO outbound message event descriptors
244 * @host_deviceid: Host device ID associated with this master port
245 * @ops: configuration space functions
246 * @id: Port ID, unique among all ports
247 * @index: Port index, unique among all port interfaces of the same type
248 * @sys_size: RapidIO common transport system size
249 * @phy_type: RapidIO phy type
250 * @phys_efptr: RIO port extended features pointer
251 * @name: Port name string
252 * @dev: device structure associated with an mport
253 * @priv: Master port private data
254 * @dma: DMA device associated with mport
255 * @nscan: RapidIO network enumeration/discovery operations
256 */
257 struct rio_mport {
258     struct list_head dbells; /* List of doorbell events */
259     struct list_head node; /* node in global list of ports */
260     struct list_head nnode; /* node in net list of ports */
261     struct resource iores;
262     struct resource riores[RIO_MAX_MPORT_RESOURCES];
263     struct rio_msg inb_msg[RIO_MAX_MBOX];
264     struct rio_msg outb_msg[RIO_MAX_MBOX];
265     int host_deviceid; /* Host device ID */
266     struct rio_ops *ops; /* Low-Level architecture-dependent routines */
267     unsigned char id; /* port ID, unique among all ports */
268     unsigned char index; /* port index, unique among all port
269                          interfaces of the same type */
270     unsigned int sys_size; /* RapidIO common transport system size.
271                          * 0 - Small size. 256 devices.
272                          * 1 - Large size. 65536 devices.
273                          */
274     enum rio_phy_type phy_type; /* RapidIO phy type */
275     u32 phys_efptr;
276     unsigned char name[RIO_MAX_MPORT_NAME];
277     struct device dev;
278     void *priv; /* Master port private data */
279 #ifdef CONFIG_RAPIDIO_DMA_ENGINE
280     struct dma_device dma;
281 #endif
282     struct rio_scan *nscan;
283 };
284
285 /*

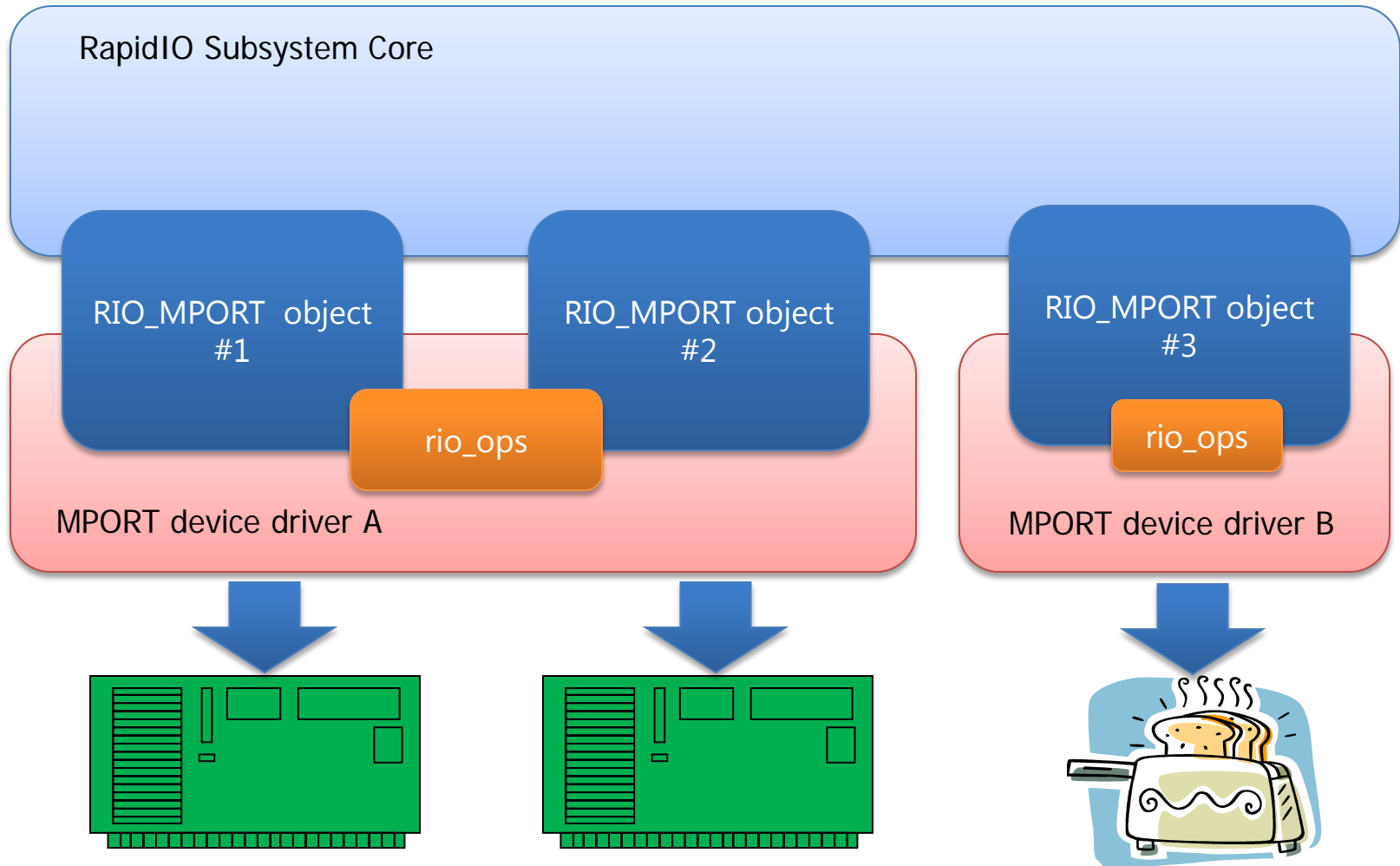
```

```

318
319 /**
320  * struct rio_ops - Low-Level RIO configuration space operations
321  * @lcread: Callback to perform local (master port) read of config space.
322  * @lcwrite: Callback to perform local (master port) write of config space.
323  * @cread: Callback to perform network read of config space.
324  * @cwrite: Callback to perform network write of config space.
325  * @dsend: Callback to send a doorbell message.
326  * @pwenable: Callback to enable/disable port-write message handling.
327  * @open_outb_mbox: Callback to initialize outbound mailbox.
328  * @close_outb_mbox: Callback to shut down outbound mailbox.
329  * @open_inb_mbox: Callback to initialize inbound mailbox.
330  * @close_inb_mbox: Callback to shut down inbound mailbox.
331  * @add_outb_message: Callback to add a message to an outbound mailbox queue.
332  * @add_inb_buffer: Callback to add a buffer to an inbound mailbox queue.
333  * @get_inb_message: Callback to get a message from an inbound mailbox queue.
334  * @map_inb: Callback to map RapidIO address region into local memory space.
335  * @unmap_inb: Callback to unmap RapidIO address region mapped with map_inb().
336  */
337 struct rio_ops {
338     int (*lcread)(struct rio_mport *mport, int index, u32 offset, int len,
339                 u32 *data);
340     int (*lcwrite)(struct rio_mport *mport, int index, u32 offset, int len,
341                  u32 *data);
342     int (*cread)(struct rio_mport *mport, int index, u16 destid,
343                 u8 hopcount, u32 offset, int len, u32 *data);
344     int (*cwrite)(struct rio_mport *mport, int index, u16 destid,
345                  u8 hopcount, u32 offset, int len, u32 *data);
346     int (*dsend)(struct rio_mport *mport, int index, u16 destid, u16 data);
347     int (*pwenable)(struct rio_mport *mport, int enable);
348     int (*open_outb_mbox)(struct rio_mport *mport, void *dev_id,
349                          int mbox, int entries);
350     void (*close_outb_mbox)(struct rio_mport *mport, int mbox);
351     int (*open_inb_mbox)(struct rio_mport *mport, void *dev_id,
352                          int mbox, int entries);
353     void (*close_inb_mbox)(struct rio_mport *mport, int mbox);
354     int (*add_outb_message)(struct rio_mport *mport, struct rio_dev *rdev,
355                            int mbox, void *buffer, size_t len);
356     int (*add_inb_buffer)(struct rio_mport *mport, int mbox, void *buf);
357     void (*get_inb_message)(struct rio_mport *mport, int mbox);
358     int (*map_inb)(struct rio_mport *mport, dma_addr_t lstart,
359                  u64 rstart, u32 size, u32 flags);
360     void (*unmap_inb)(struct rio_mport *mport, dma_addr_t lstart);
361 };
362
363 #define RIO_RESOURCE_MEM      0x00000100
364 #define RIO_RESOURCE_DOORBELL 0x00000200
365 #define RIO_RESOURCE_MAILBOX 0x00000400

```

- ❑ During driver initialization an MPORT driver registers its `rio_mport` object (including supported operations) within the RIO core.
 - ❑ RIO core uses `rio_mport` object as an abstract representation of HW device during its operations.
 - ❑ When RIO core function needs to perform operation controlled by HW (e.g. maintenance read from remote target device) it calls corresponding function supplied through the `rio_mport` object.
- If a RIO operation is not supported by given mport device, pointer to that operation in `rio_ops` must to be set to NULL. Some operations are mandatory for implementation and RIO core does not check for NULL pointer.



```
rio_mport_read_config_32(port, destid, hopcount,  
                        RIO_DEV_ID_CAR,  
                        &result);
```

Upper level kernel module (e.g rio-scan) calls core interface function.

```
res = mport->ops->cread(mport, mport->id,  
                        destid, hopcount, offset, len, &data);
```

RIO core maps the call into corresponding mport operation.

```
static int tsi721_cread_dma(*mport, index,  
                           destid, hcount, offset, len, *data)  
{  
    .....  
    return tsi721_maint_dma(priv, mport->sys_size,  
                           destid, hopcount, offset, len, data, 0);  
}
```

MPORT device driver performs requested SRIO operation and returns a result back through the call chain.

- ❑ Support for DMA data transfers is different from callback scheme described earlier.
- ❑ DMA operations in Linux kernel must be provided through generic DMA engine framework.
- ❑ MPORT device drivers that support DMA implement and register their DMA operations with Linux kernel DMA engine.
- ❑ RapidIO uses SLAVE DMA mode defined by DMA engine with subsystem specific extension.
- ❑ To communicate with DMA engine RIO core provides interface functions:
 - `rio_request_mport_dma(struct rio_mport *mport)`
 - `rio_request_dma(struct rio_dev *rdev)`
 - `rio_release_dma(struct dma_chan *dchan)`
 - `rio_dma_prep_xfer(struct dma_chan *dchan, ...)`

- ❑ Functions listed above address only RIO-specific part of DMA transfers programming.
- ❑ After RIO-capable DMA channel have been allocated standard DMA engine operations should be used.
- ❑ Not every RIO controller provides HW DMA support for data transfers to/from remote RIO target device.
- ❑ Even HW DMA support exists it does not guarantee that DMA engine support have been implemented for that device/mode. E.g. kernel DMA engine code for Freescale's MPC8548 device supports memory-to-memory and regular SLAVE mode, but RIO-compatible SLAVE mode is not implemented (yet).

- ❑ HW controllers supported by kernel RapidIO subsystem can be divided into two groups:
 - Built-in SRIO controllers inside of SoC (Freescale, TI, Cavium, FPGAs, etc.).
 - External bridges to SRIO from busses like PCI or PCI Express (PCI - Tsi620, PCIe – Tsi721).

- ❑ Built-in controllers are specific for given SoC architecture and their mport drivers are located in corresponding arch code branches (this is why we do not see them in “drivers/rapidio” directory).
- ❑ MPORT drivers for external bridges are (mostly) architecture independent and are located in “drivers/rapidio” directory.
- ❑ Both forms of mport drivers register their rio_mport object with RIO subsystem core and therefore their operations are transparent to callers of core interface functions.

SRIO controllers/bridges do not provide an equal level of RapidIO functionality. Currently only Tsi721 provides all set of functions (options) supported by RapidIO subsystem core. Here are examples of some HW or SW limitations:

- Freescale SoCs – have reduced number of RIO messaging MBOXes, do not have DMA engine support for RIO (HW is OK).
- TI Keystone II SoC – does not have inbound window address translation capability.

- ❑ Different enumeration perspective:
 - PCIe is single-root tree (many RCs are possible) with address-based routing.
 - Address-based routing can be a limiting factor in large systems.
 - PCIe eliminated peer-to-peer capability that was available in older PCI implementations.
 - SRIO is fabric with possibility of multipath routing and redundant mport connections on the same node.

ABOVE THE CORE

- ❑ RapidIO subsystem has multiple components that act as users of generic interface functions exported by the core.
- ❑ Those functional components are hardware independent and can work on any platform that support necessary RapidIO operations.
- ❑ Upper layer kernel RapidIO subsystem components:
 - **Enumeration/Discovery (RIO_SCAN)** – implements a basic RapidIO fabric enumeration/discovery process which automatically assigns destination IDs to endpoints and programs switch routing tables.
 - **RIONET** - Ethernet emulation over RapidIO. This driver allows to use RapidIO controllers as common Ethernet NIC with all available TCP/IP stack support.
 - **SYSFS support** – RapidIO subsystem creates its own sysfs entries that allow to analyze state of the subsystem.
 - **MPORT_CDEV** – universal character mode device driver that allows user-space communications with RapidIO fabric.

- ❑ To make RapidIO devices visible/accessible to node users SRIO fabric has to be enumerated.
- ❑ Compared to PCI/PCIe bus enumeration, SRIO fabric enumeration may have multiple methods depending on system configuration and use model.
- ❑ Any RapidIO system can implement its own enumeration method.
- ❑ As an example of enumeration/discovery process Linux kernel RapidIO subsystem includes a basic implementation which automatically assigns destination IDs to endpoints and programs switch routing tables (based on Annex 1 of RapidIO specification):
 - Enumeration is performed by a dedicated node (host). RapidIO enumerating host is defined by kernel boot command line parameter "rapidio.hdid=n,...".
 - All other nodes (targets) in RapidIO network perform the discovery process.

- ❑ The basic implementation can enumerate the RIO fabric only in “fresh after power-on” state (all nodes). It does not support re-enumeration of the fabric or hot-plug node/segment addition/removal.
- ❑ The basic enumeration/discovery process is implemented as a loadable kernel module (rio-scan.c) that can be replaced by an end-user if required.
- ❑ RIO subsystem core offers mechanism to register enumerating process on per-mport basis. Therefore, if required users can implement different enumeration/discovery processes for each mport.
- ❑ With MPORT_CDEV device driver availability users have got ability to implement their own enumeration algorithms as user space applications instead of kernel module.

RIO_MPORT_CDEV DEVICE DRIVER

- ❑ This driver allows to perform RIO operations from user space applications.
- ❑ It implements common IOCTL interfaces between kernel RIO subsystem core and DMA engine and user space.
- ❑ Next slide demonstrates request path from user space to SRIO bus.

```
ioctl(hnd->fd, IO_MPORT_MAINT_READ_REMOTE, &mt)
```

User-space application or library issues IOCTL request to mport_cdev device driver.

```
rio_mport_read_config_32(port, destid, hopcount,  
RIO_DEV_ID_CAR,  
&result;);
```

mport_cdev processes IOCTL request and calls core interface function.

```
res = mport->ops->cread(mport, mport->id,  
destid, hopcount, offset, len, &data;);
```

RIO core maps the call into corresponding mport operation.

MPORT device driver performs requested SRIO operation and returns a result back through the call chain (see slide 11).

- ❑ RapidIO subsystem offers several debug tracing options that allow end user to analyze component operations for better understanding or issue resolution.
- ❑ At the top level debug traces for RapidIO subsystem are controlled by kernel configuration option `CONFIG_RAPIDIO_DEBUG`.
- ❑ `CONFIG_RAPIDIO_DEBUG` enables debug output from all subsystem components.
- ❑ To reduce amount of produced debug information some newer components implement filtered debug output options controlled through kernel module parameters.
- ❑ Currently this option available for Tsi721 mport driver, `rio_mport_cdev` driver and `rio_cm` driver.
- ❑ Older subsystem components will be converted to support this option as time permits.

- ❑ Debugging output from this driver is controlled by kernel module parameter "dbg_level".
- ❑ By default this parameter is set to 0, what means "disabled".
- ❑ Two ways to set this parameter (values are for example only):
 - When loading the module: "modprobe tsi721_mport dbg_level=0x23", or
 - at runtime using sysfs (with root privileges):
"echo 0x1f > /sys/module/tsi721_mport/parameters/dbg_level"
- ❑ Current value of dbg_level can be displayed at any time by entering the following command: "cat /sys/module/tsi721_mport/parameters/dbg_level". NOTE: value is returned in decimal format, convert to hex to match the mask.
- ❑ Debug output bit masks are listed below (see drivers/rapidio/devices/tsi721.h for updates)

```
DBG_INIT= 0x01    /* driver initialization */
DBG_EXIT= 0x02    /* driver exit/removal */
DBG_MPORT= 0x04   /* mport add/remove */
DBG_MAINT= 0x08   /* maintenance ops messages */
DBG_DMA= 0x10     /* DMA transfer messages */
DBG_DMAV= 0x20    /* verbose DMA transfer messages */
DBG_IBW= 0x40     /* inbound window messages */
DBG_EVENT= 0x80   /* event handling messages */
DBG_OBW= 0x100    /* outbound window messages */
DBG_DBELL=0x200   /* doorbell configuration/handling messages */
DBG_OMSG=0x400   /* outbound MBOX messages */
DBG_IMSG= 0x800   /* inbound MBOX messages */
```

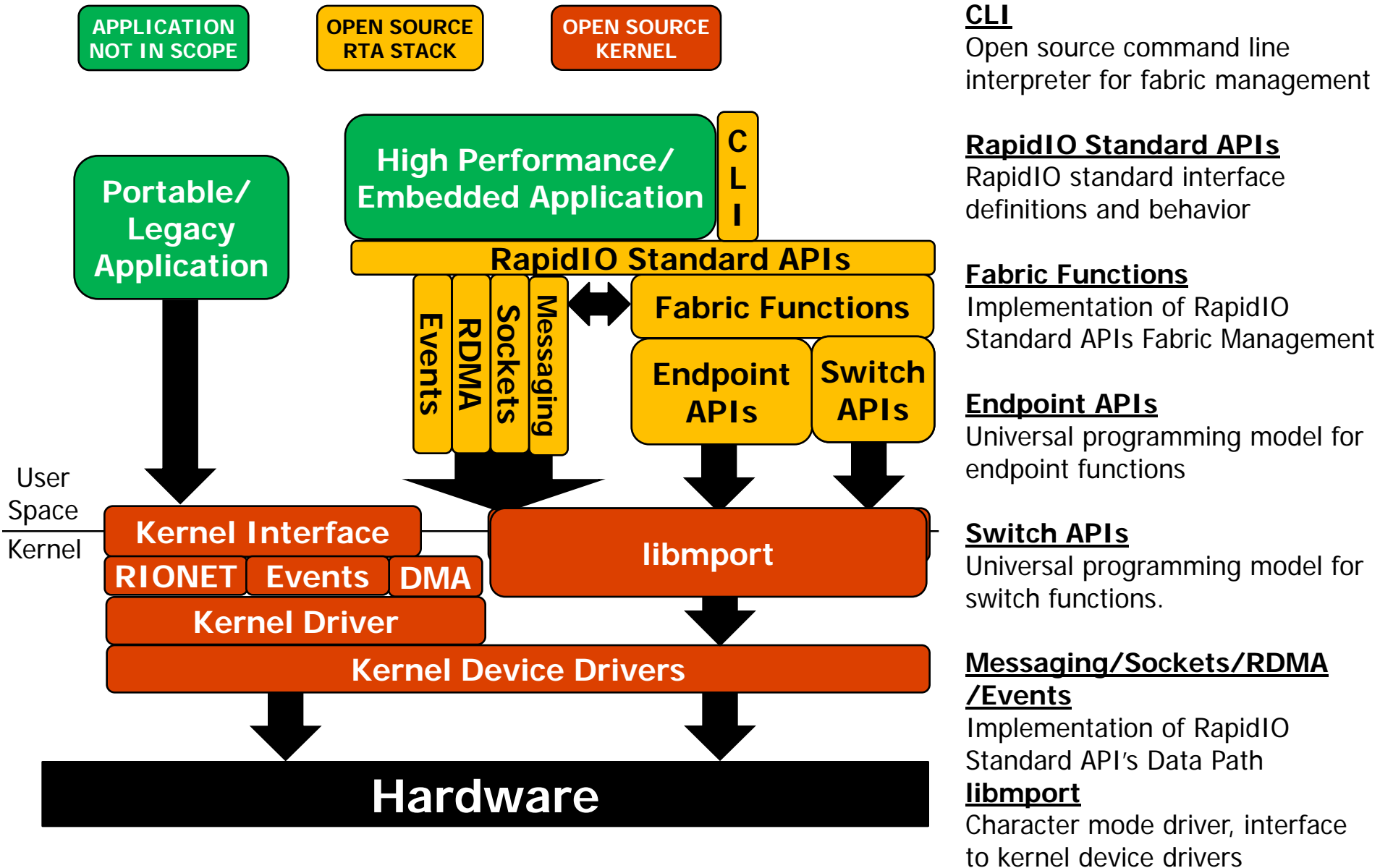

- ❑ Debugging output from this driver is controlled by kernel module parameter "dbg_level".
- ❑ By default this parameter is set to 0, what means "disabled".
- ❑ Two ways to set this parameter (values are for example only):
 - When loading the module: "modprobe rio_mport_cdev dbg_level=0x23", or
 - at runtime using sysfs (with root privileges):
"echo 0x1f > /sys/module/rio_mport_cdev/parameters/dbg_level"
- ❑ Current value of dbg_level can be displayed at any time by entering the following command: "cat /sys/module/rio_mport_cdev/parameters/dbg_level". NOTE: value is returned in decimal format, convert to hex to match the mask.
- ❑ Debug output bit masks are listed below (see drivers/rapidio/devices/rio_mport_cdev.c)

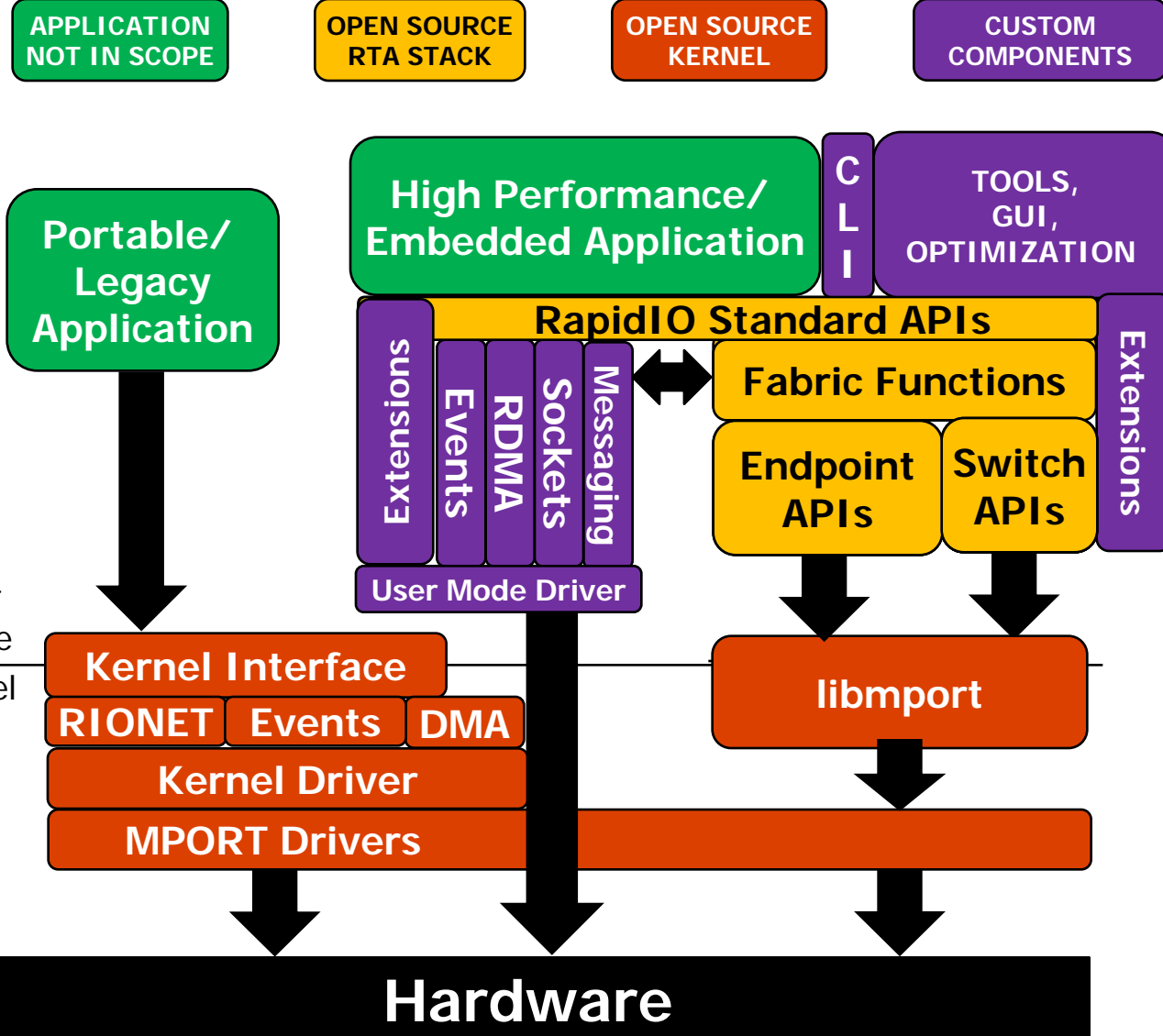
```
DBG_INIT= 0x01    /* driver init */
DBG_EXIT= 0x02    /* driver exit */
DBG_MPORT= 0x04   /* mport add/remove */
DBG_RDEV=0x08     /* RapidIO device add/remove */
DBG_DMA= 0x10     /* DMA transfer messages */
DBG_MMAP=0x20     /* mapping messages */
DBG_IBW=0x40      /* inbound window */
DBG_EVENT= 0x80   /* event handling messages */
DBG_OBW=0x100     /* outbound window messages */
DBG_DBELL=0x200   /* doorbell messages */
```

- ❑ Masks for both drivers seems to be similar but because they are used at different layers sequential approach can be useful to reduce amount of produced messages.
- ❑ When debugging user space application start with traces from `rio_mport_cdev` and decide if an error is reported by this driver or reported by `tsi721_mport` driver.
- ❑ If needed add traces from `tsi721_mport` driver.

RAPIDIO SOFTWARE BACKGROUND

Conceptual - Open Source Stack





TOOLS, GUI, OPTIMIZATION

Value adds for debug/monitoring, system visualization, data interpretation, and topology specific functions

User Mode Driver

User mode drivers optimized for executing hardware operations in user mode.

Extensions

Additional data path and/or fabric management services provided to applications

Minimum functionality to demonstrate RapidIO value:

- High throughput
- Low latency
- Fault tolerant

Portable between operating systems, hardware configurations, and applications

Interoperable between any endpoint and/or switch compliant to the standard

Extensible to enable innovation by all members of the RapidIO community

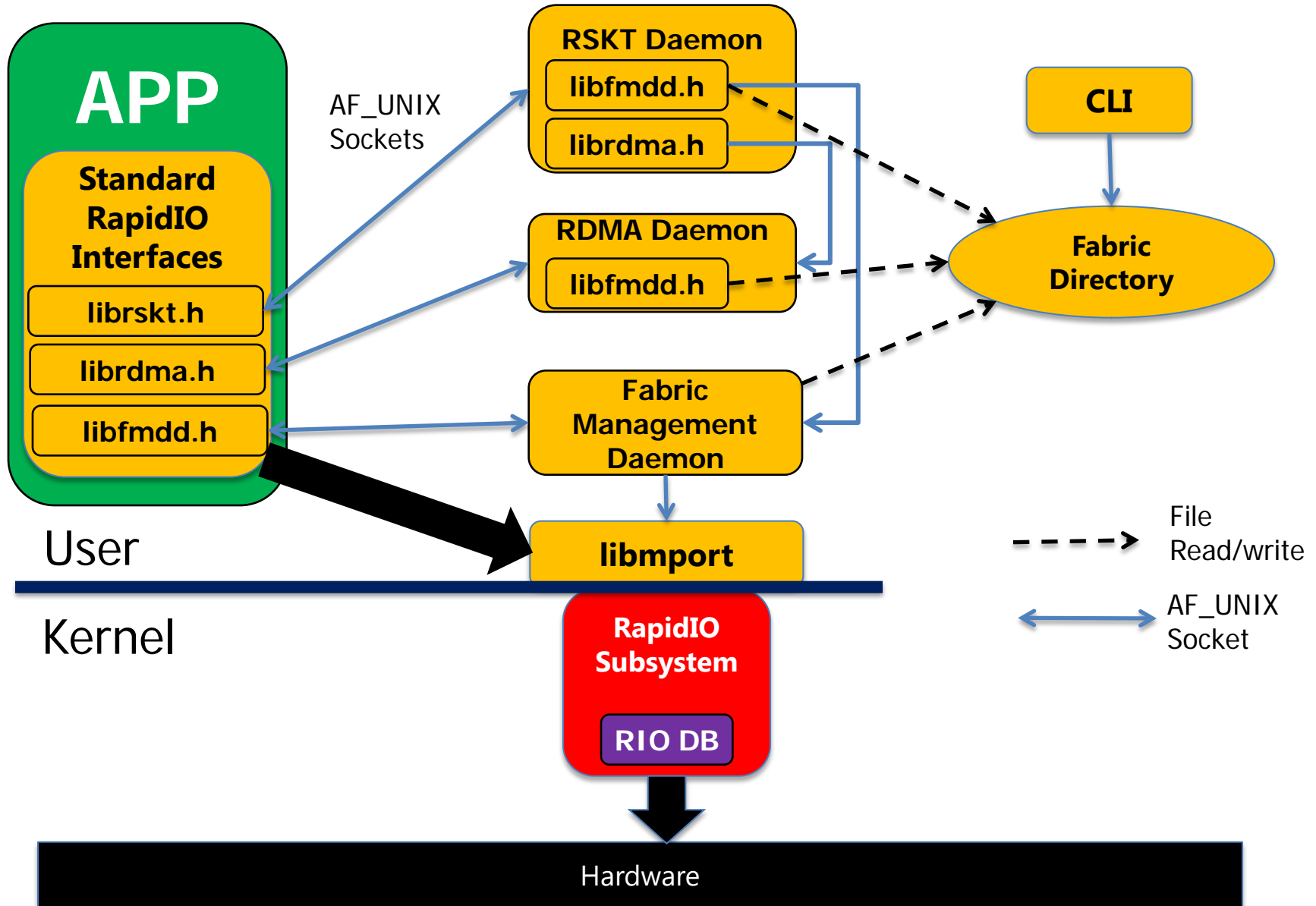
RAPIDIO SOFTWARE STACK PROCESS STRUCTURE

- Clone the rapidio_sw repository from www.github.com/RapidIO into /opt/rapidio/
 - Creates /opt/rapidio/rapidio_sw directory
- in the “rapidio_sw” directory, type “make all”
- create /etc/rapidio/rio.conf file to control master and slave Fabric Management Daemons
 - Examples: fabric_management/daemon/cfg
 - GRY05: Master
 - GRY06: Slave

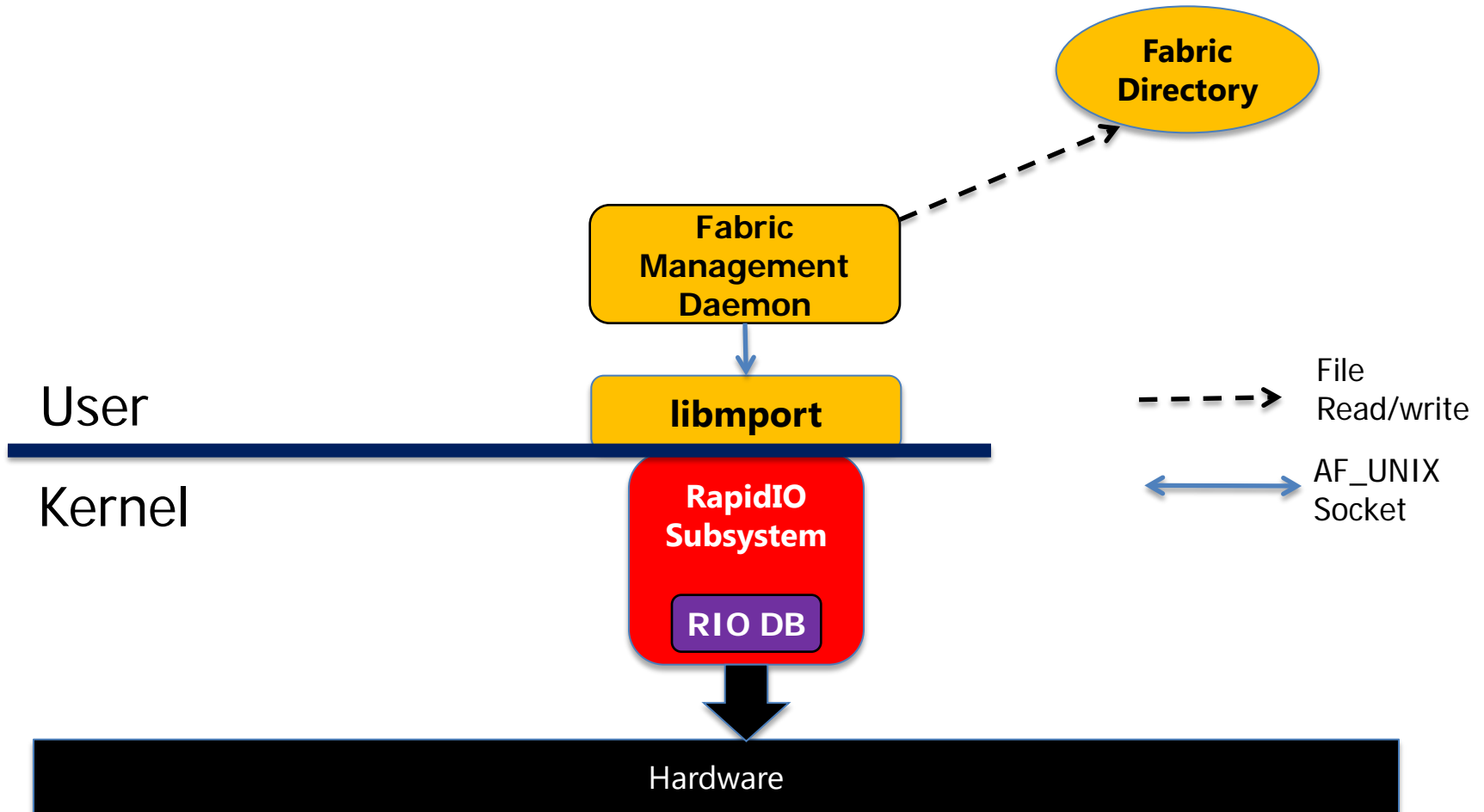
- the scripts in `/opt/rapidio/` will use the executables in `/opt/rapidio/rapidio_sw` to start the RapidIO Software stack
- Similar to previous installation:
 - `./all_down.sh` – powers down all nodes
 - `./all_start.sh` – starts kernel RapidIO subsystem with kernel enumeration
 - `./ibox_start_all` – starts user mode RapidIO subsystem with user mode enumeration

Function	Daemon	Implementation rapidio_sw/	Library in /include	Implementation rapidio_sw/
Fabric Management	FMD	fabric_management/daemon	libfmdd.h	fabric_management/libfmdd.h
RDMA	RDMAD	rdma/daemon	librdma.h	rdma/lib
RDMA Sockets	RSKTD	rdma/rskt/daemon	librskt.h	rdma/rskt/libr
Command Line Interpreter	CLI	TBD	N/A	N/A
Fabric Directory	N/A	Posix Shared Memory Region	N/A	N/A

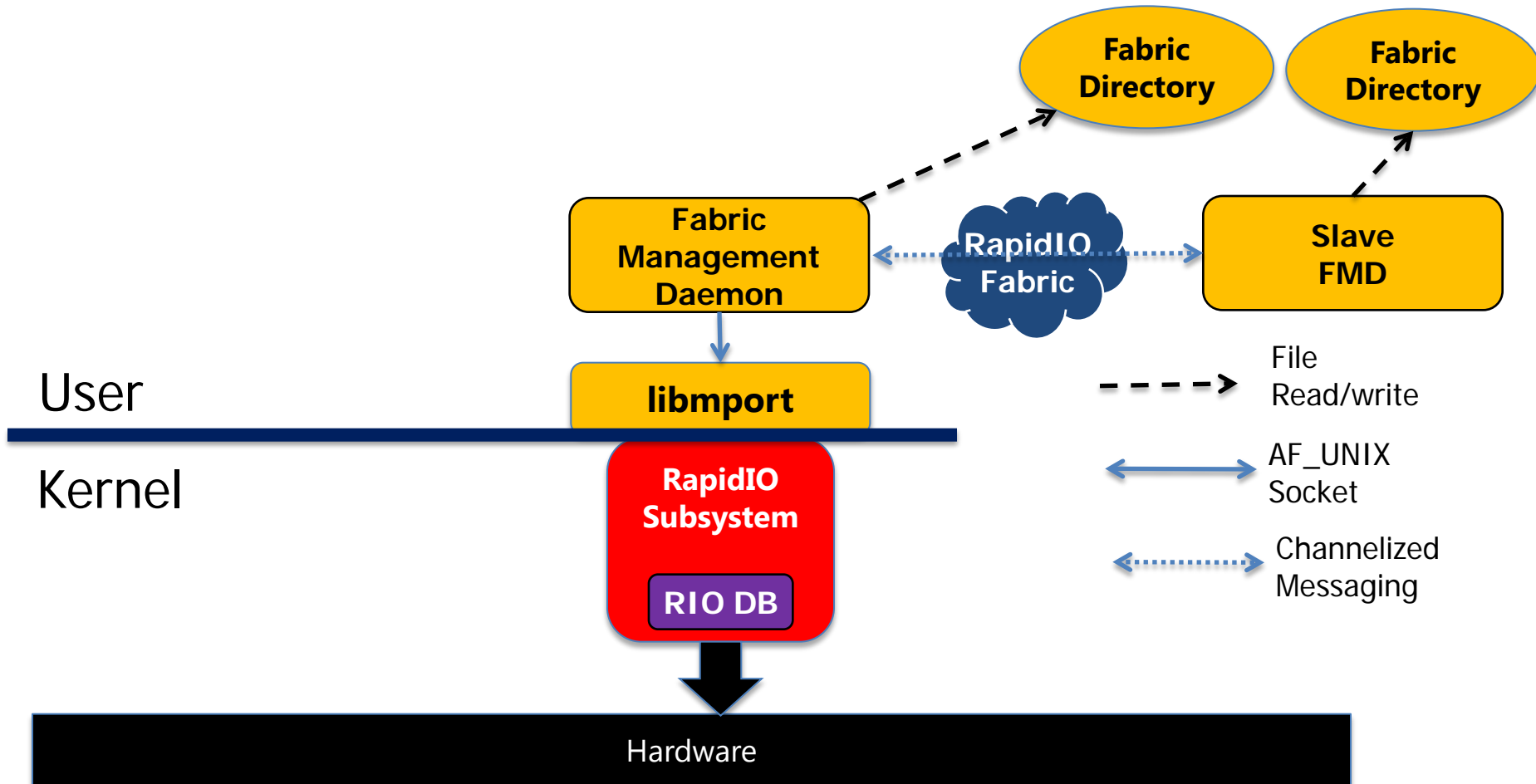
Library & Daemon Relationships



Startup – Master FMD First

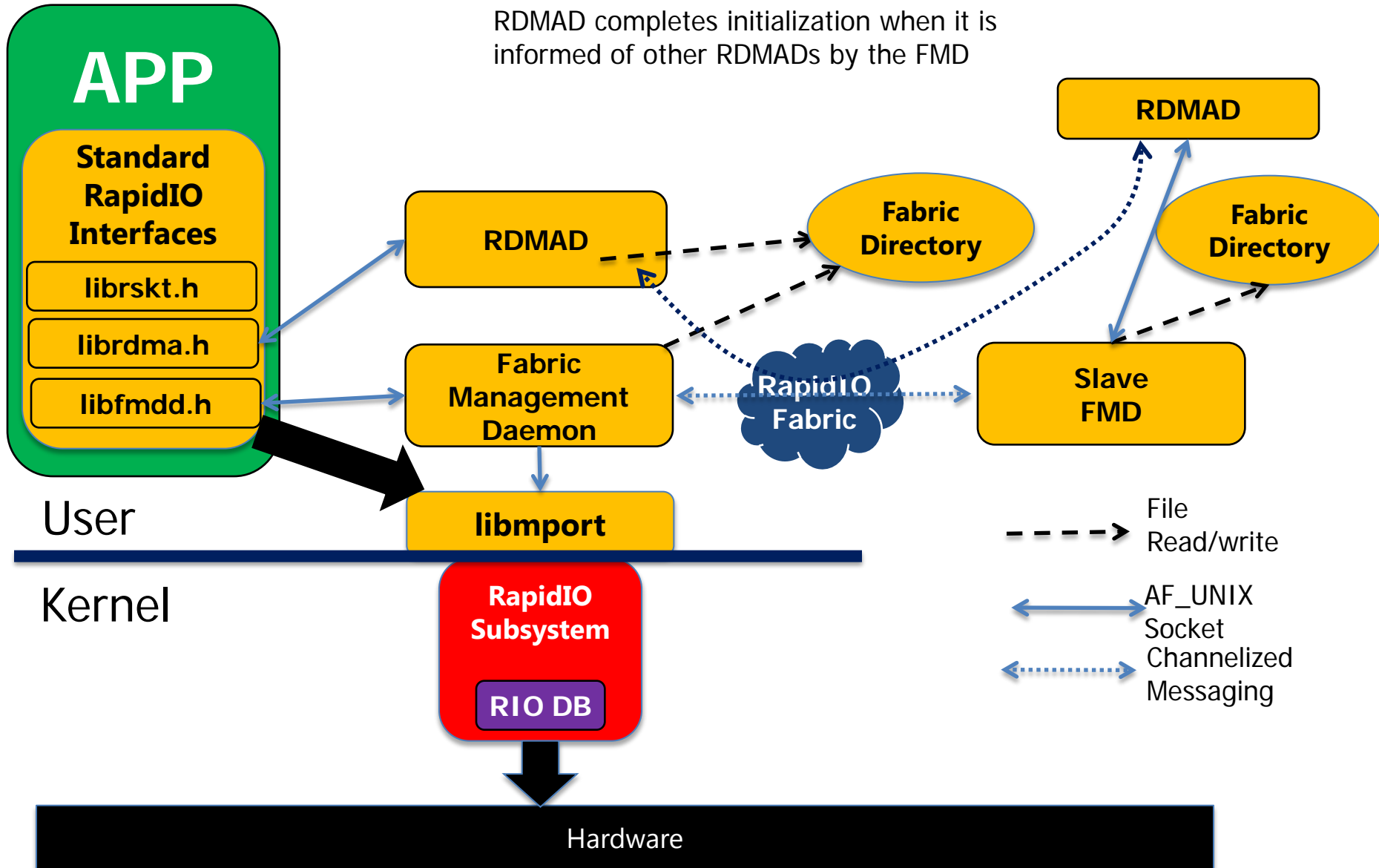


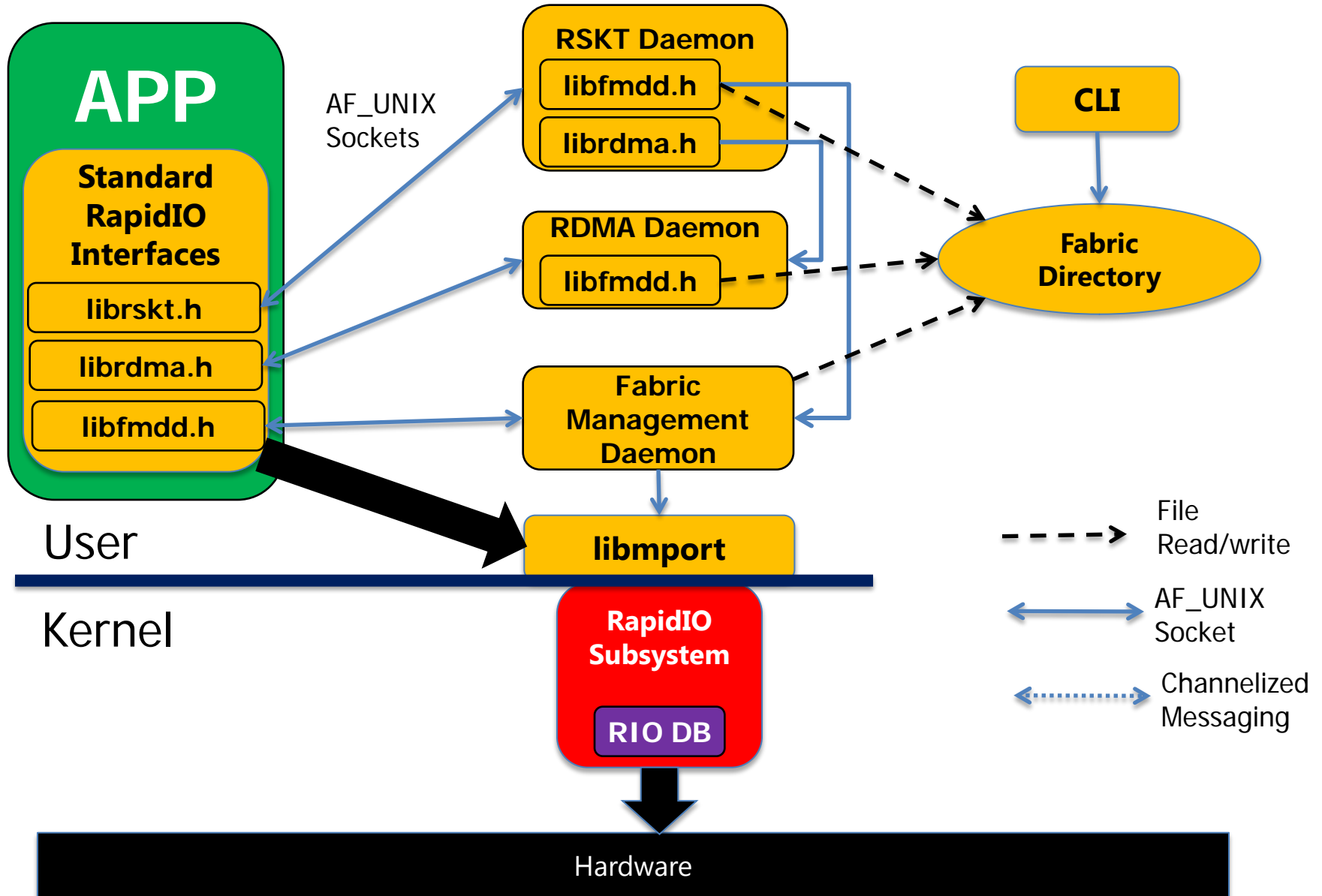
Startup – Slave FMDs Next



Startup – RDMADs, in any order

RDMAD completes initialization when it is informed of other RDMADs by the FMD





RAPIDIO FABRIC MANAGEMENT OVERVIEW

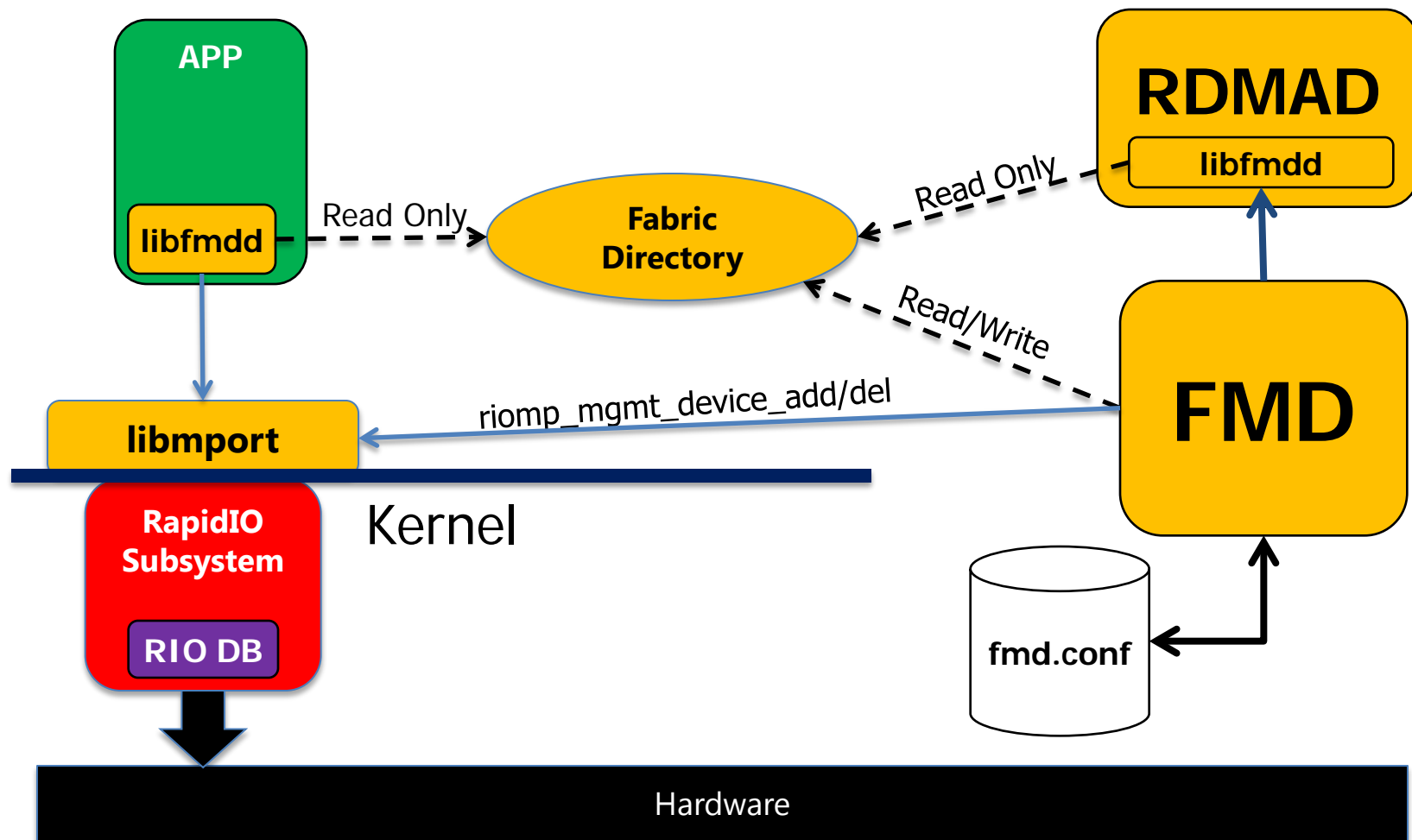
System enumeration

- uses configuration file

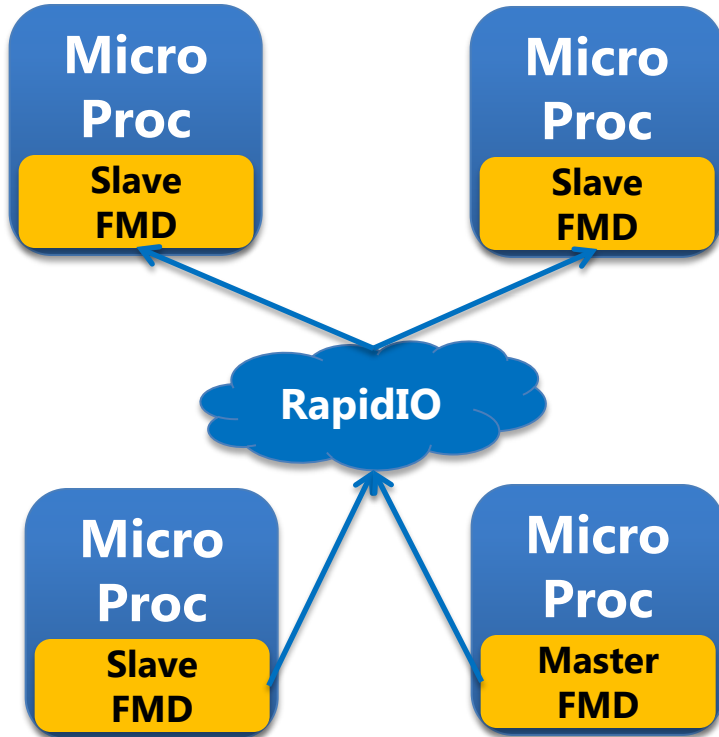
Notification of application availability

- Tells all nodes when RDMA, RSKT applications are available

- Fault tolerance
- Supports hot swap
- hardware and software robustness



RDMA Library gets status information from the Fabric Directory
 Fabric Management Daemon receives all port writes
 Fabric Management Daemon uses libmport to manage Kernel Devices

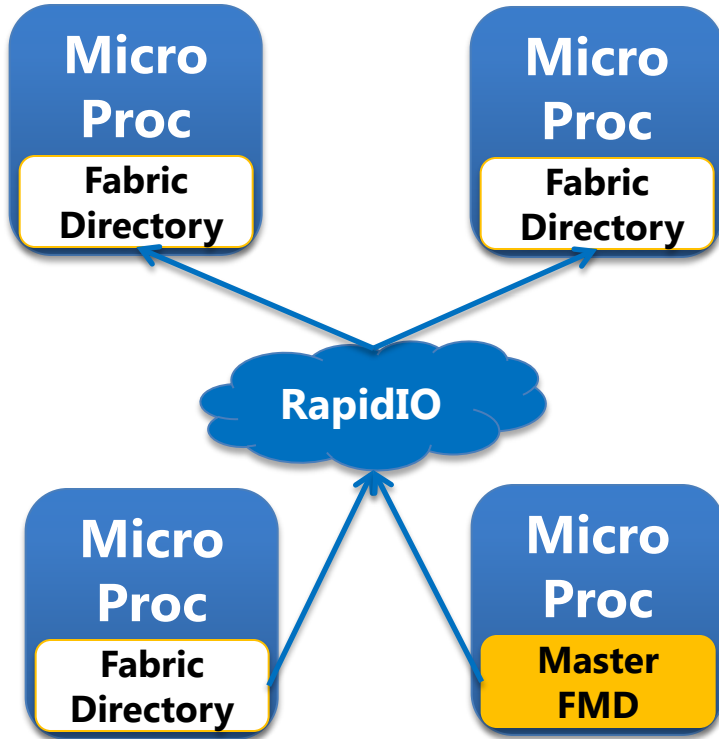


Master Fabric Management Daemon (FMD) enumerates system, then informs Slave FMDs to read configuration of the system..

Hot swap events cause port-writes to be sent to all FMDs.

Each FMD updates its local database and applications.

DISTRIBUTED FABRIC MANAGEMENT



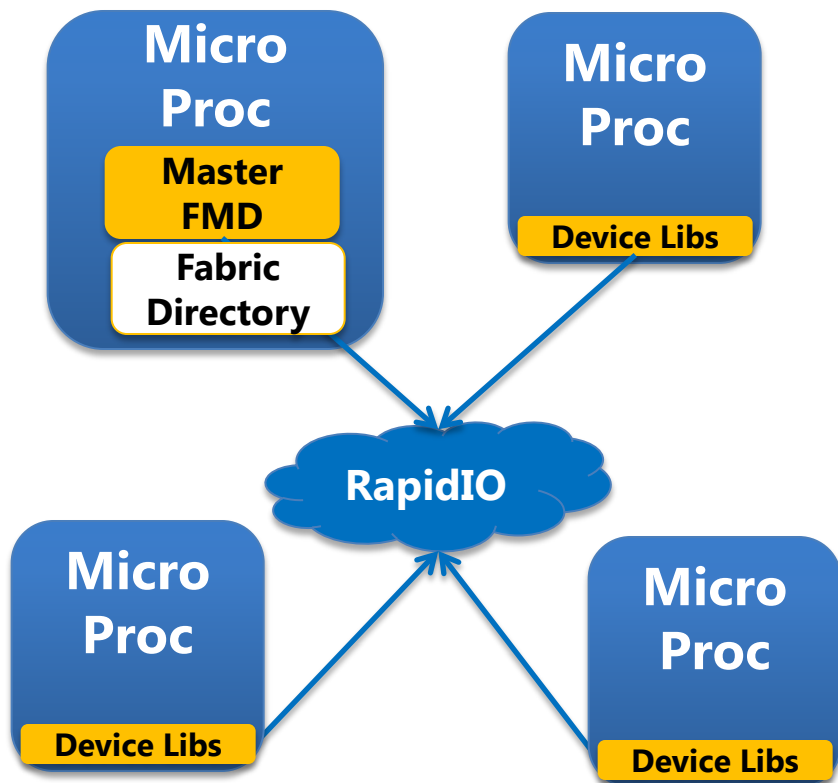
Master Fabric Management Daemon (FMD) pushes/writes system configuration information to RDMA memory on remote devices.

Centralizes fabric event management

Redundant Master FMD's possible

Simplifies system partitioning
 - Master FMD "knows" what switches, endpoints, destination IDs it controls
 - Services outside the FMD only know the devices in the Fabric Directory

CENTRALIZE FABRIC MANAGEMENT



Services (i.e. RDMA) running on remote devices pull/read system configuration information from Master RDMA read-only accessible region

Remote access using source/target enforced read only RDMA

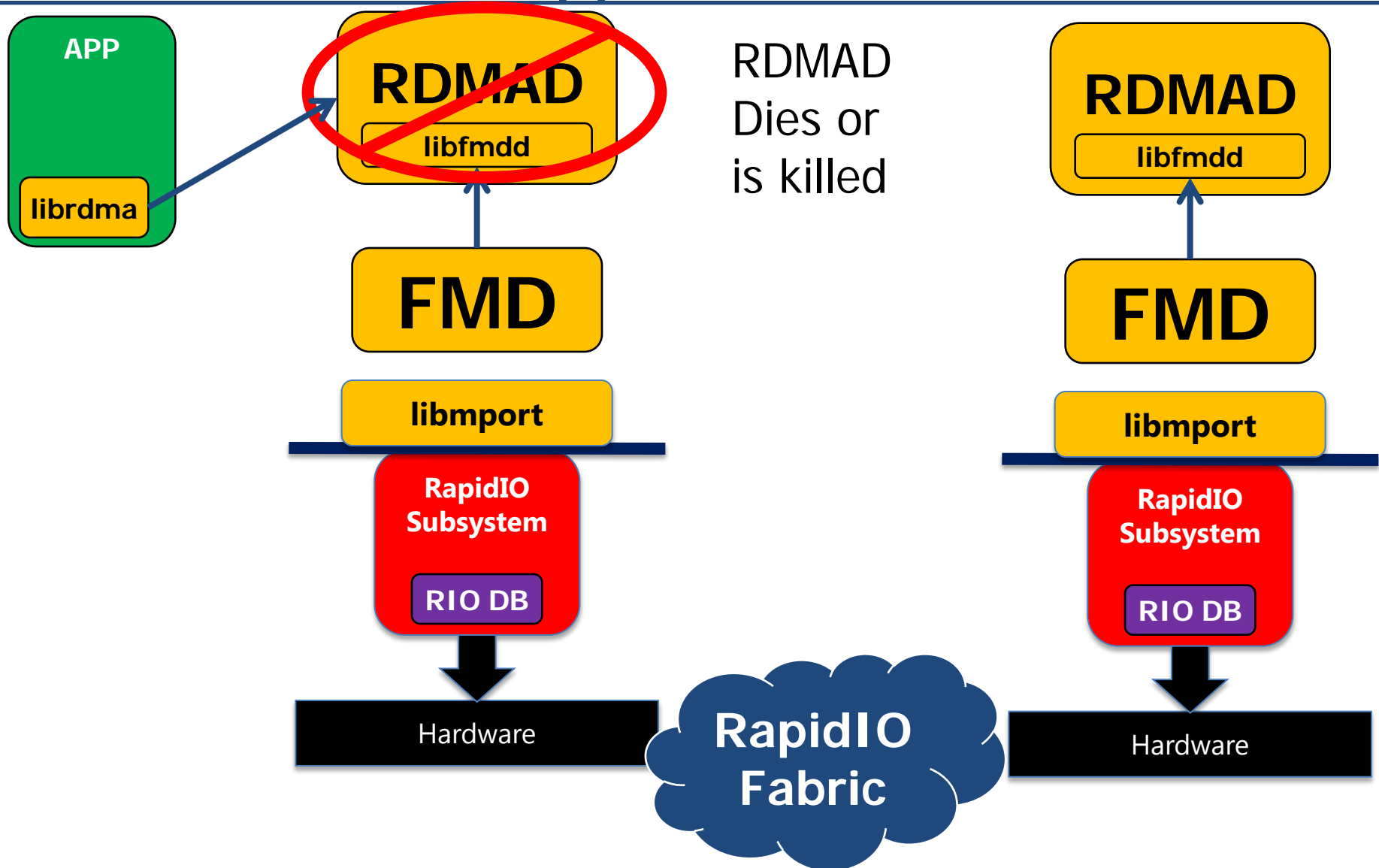
Could also use a cache coherency protocol for synchronization

Centralized, secure, scalable network management

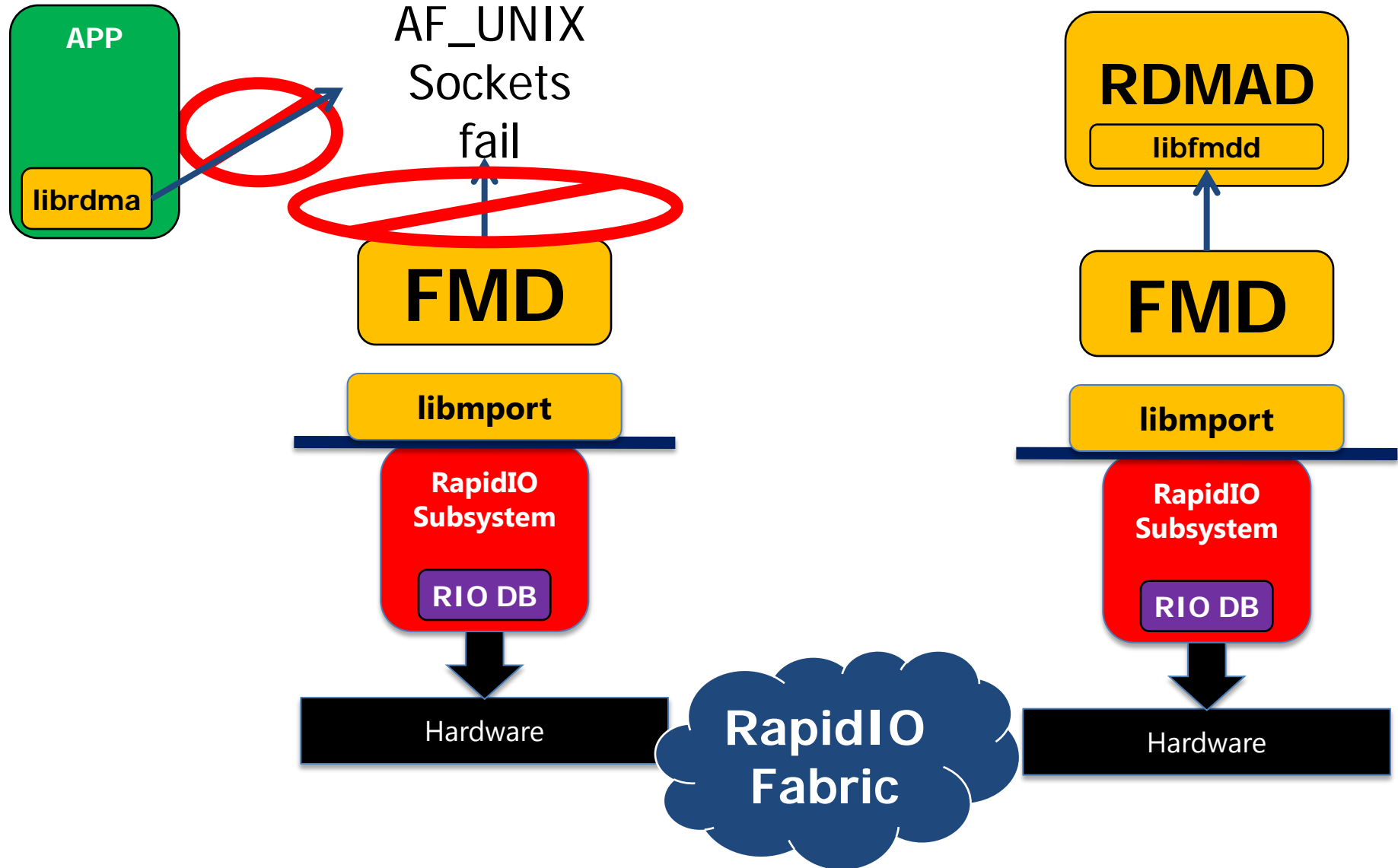
Master FMD can manage multiple different partitions, publish to different Fabric Directories

SECURE, SCALABLE, PARTIONABLE FABRIC MANAGEMENT

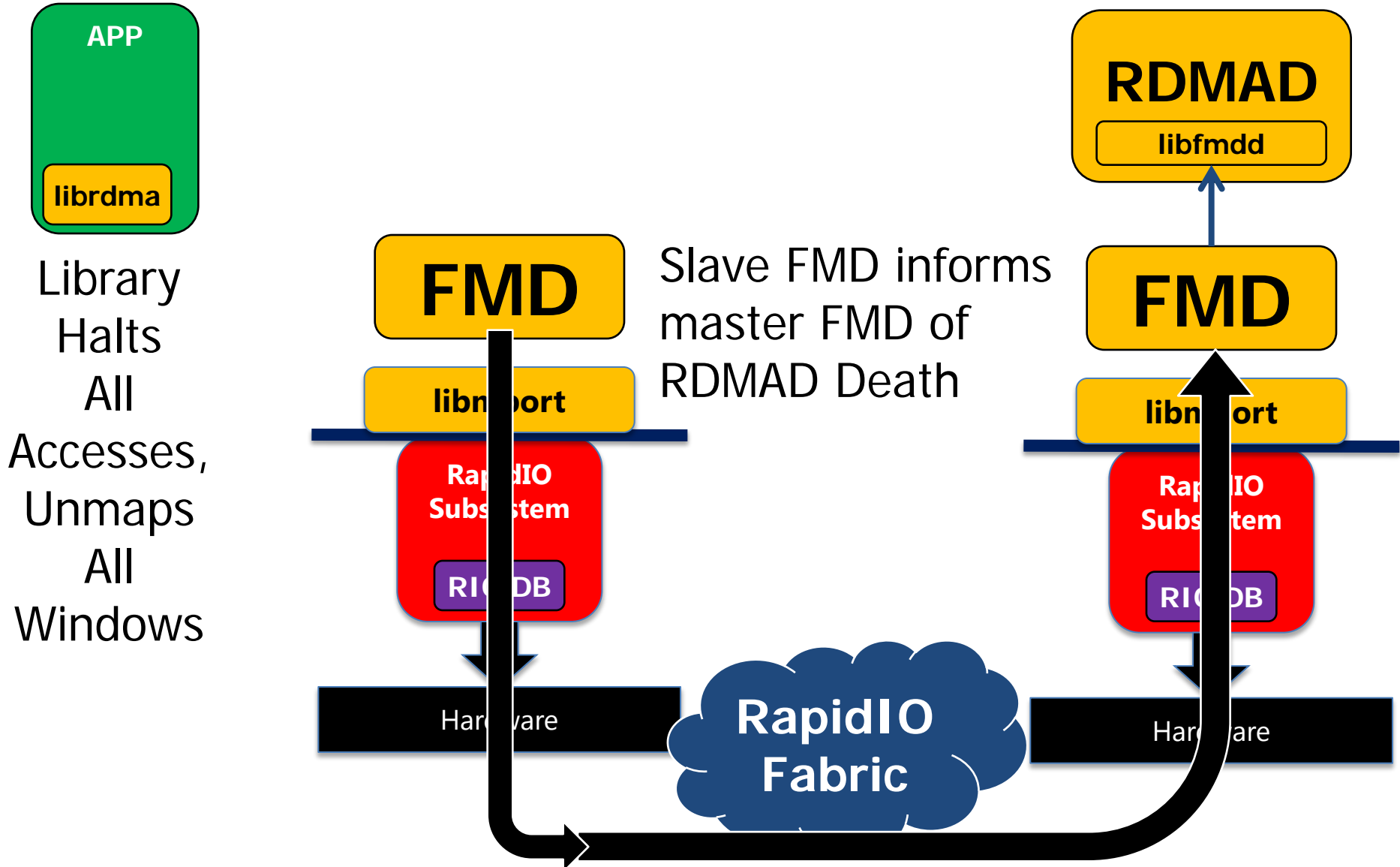
FMD Application Fault Tolerance



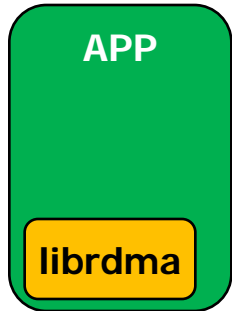
FMD Application Fault Tolerance



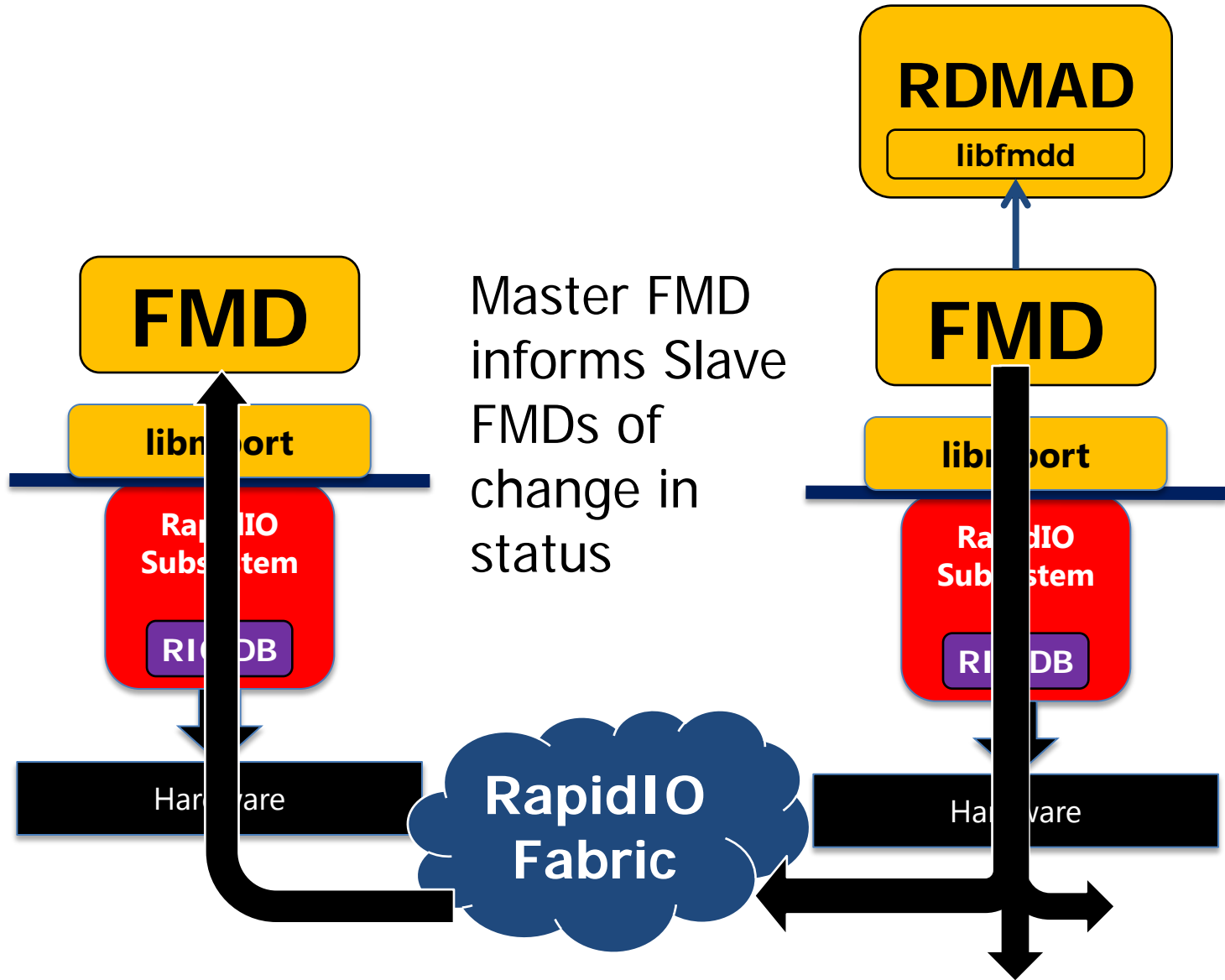
FMD Application Fault Tolerance



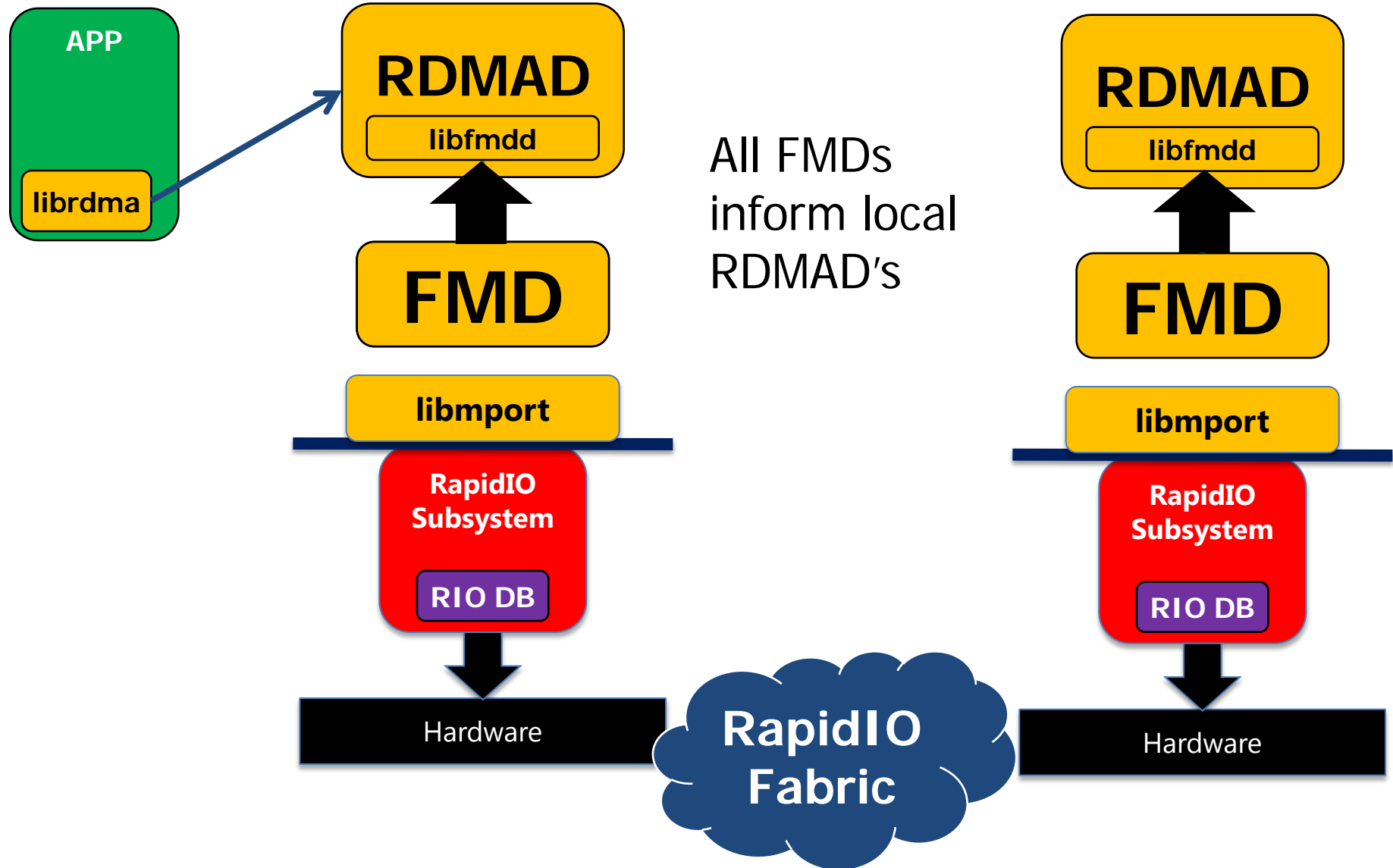
FMD Application Fault Tolerance



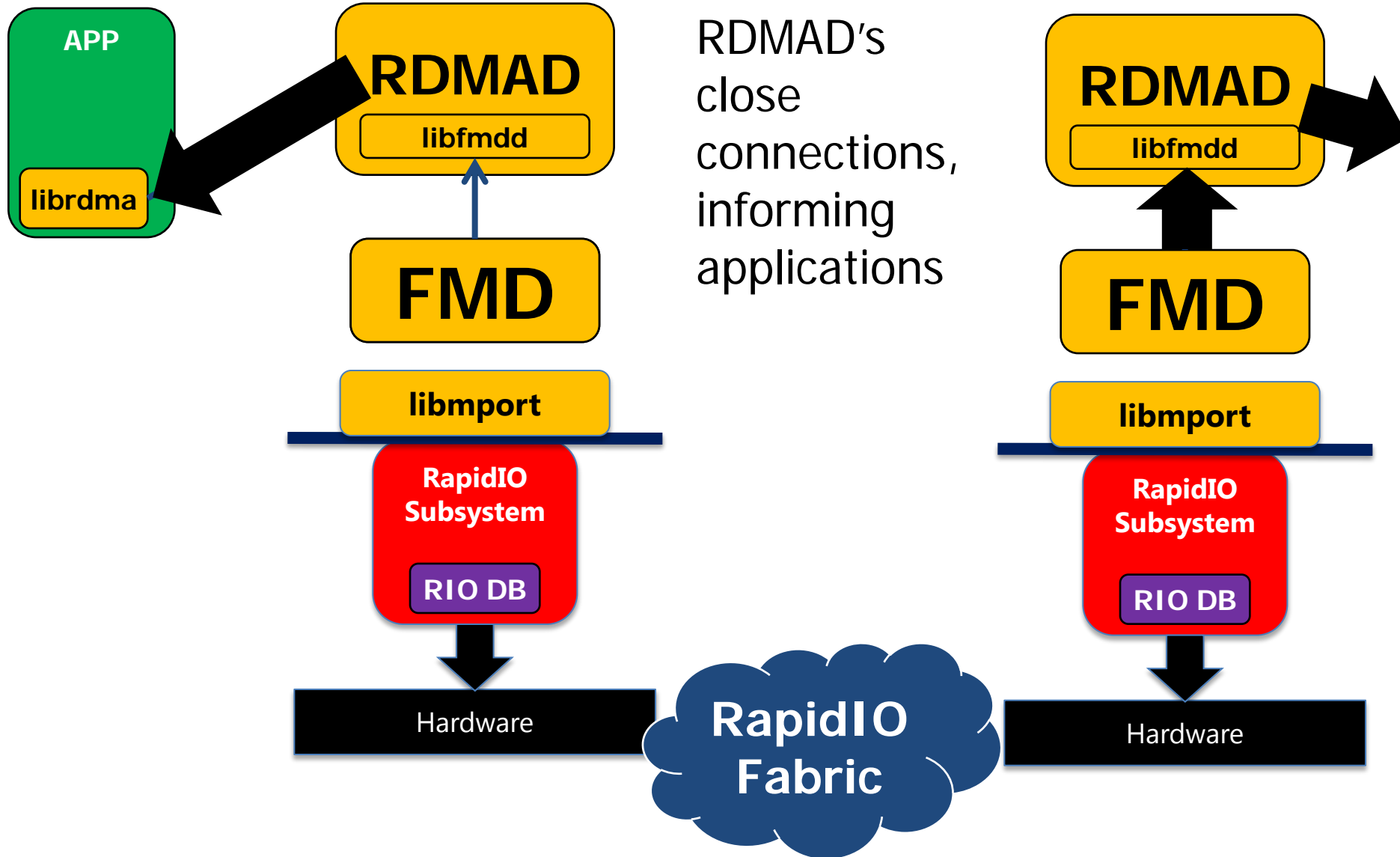
Kernel
reclaims
inbound
windows
&
memory



FMD Application Fault Tolerance



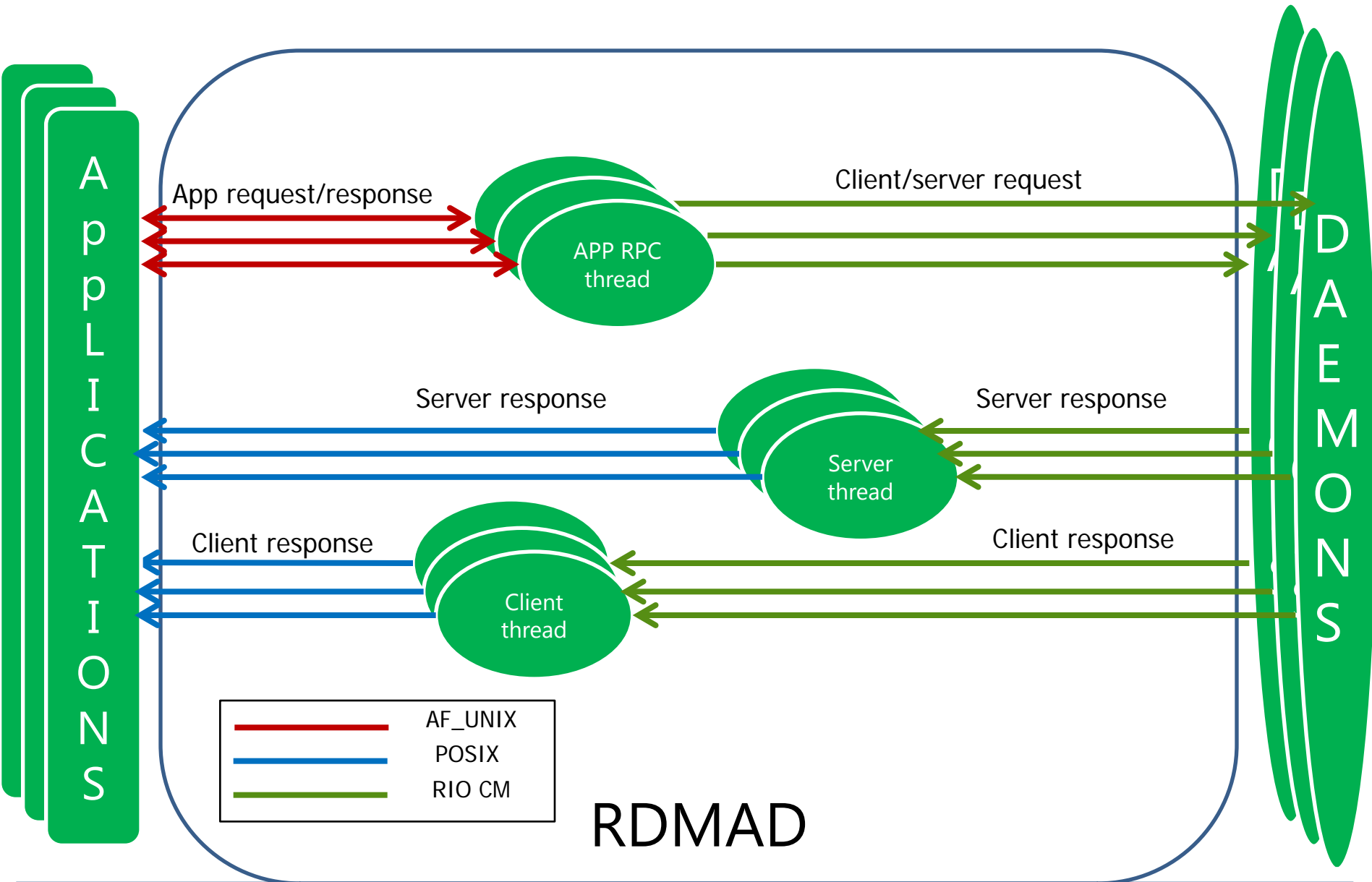
FMD Application Fault Tolerance



RAPIDIO RDMA OVERVIEW

Routine	Description
rdma_create_mso_h	Owner: Create Memory Space Owner handle
rdma_open_mso_h	User: Open memory space owner handle
rdma_close_mso_h	User: Close memory space owner handle
rdma_destroy_mso_h	Owner: Destroy Memory Space Owner handle
rdma_create_ms_h	Owner: Create a Memory Space
rdma_open_ms_h	User: Open memory space
rdma_close_ms_h	User: Close memory space
rdma_destroy_ms_h	Owner: Destroy Memory Space
rdma_create_msub_h	User: Create a subspace handle for a Memory Space
rdma_destroy_msub_h	User: Destroy a subspace handle
rdma_mmap_msub	Memory map a subspace
rdma_munmap_msub	Destroy memory mapping of a subspace

Routine	Description
<code>rdma_conn_ms_h</code>	Connect to a memory subspace on another device
<code>rdma_accept_ms_h</code>	Accept a connection request for a memory subspace
<code>rdma_disc_ms_h</code>	Destroy connection between remote and local subspaces
<code>rdma_push_msub</code>	Write data from local subspace to remote subspace
<code>rdma_pull_msub</code>	Read data from remote subspace to local subspace
<code>rdma_push_buf</code>	Write data from local buffer to remote subspace
<code>rdma_pull_buf</code>	Read data from remote subspace to local buffer
<code>rdma_sync_chk_push_pull</code>	<p>Push and pull routines support DMA completions:</p> <ul style="list-style-type: none"> • Blocking: Return when DMA is complete • Asynchronous: Continue, check for completion later • Fire and forget: Don't tell anyone <p><code>rdma_sync_chk_push_pull</code> is the "check for completion later" routine.</p>



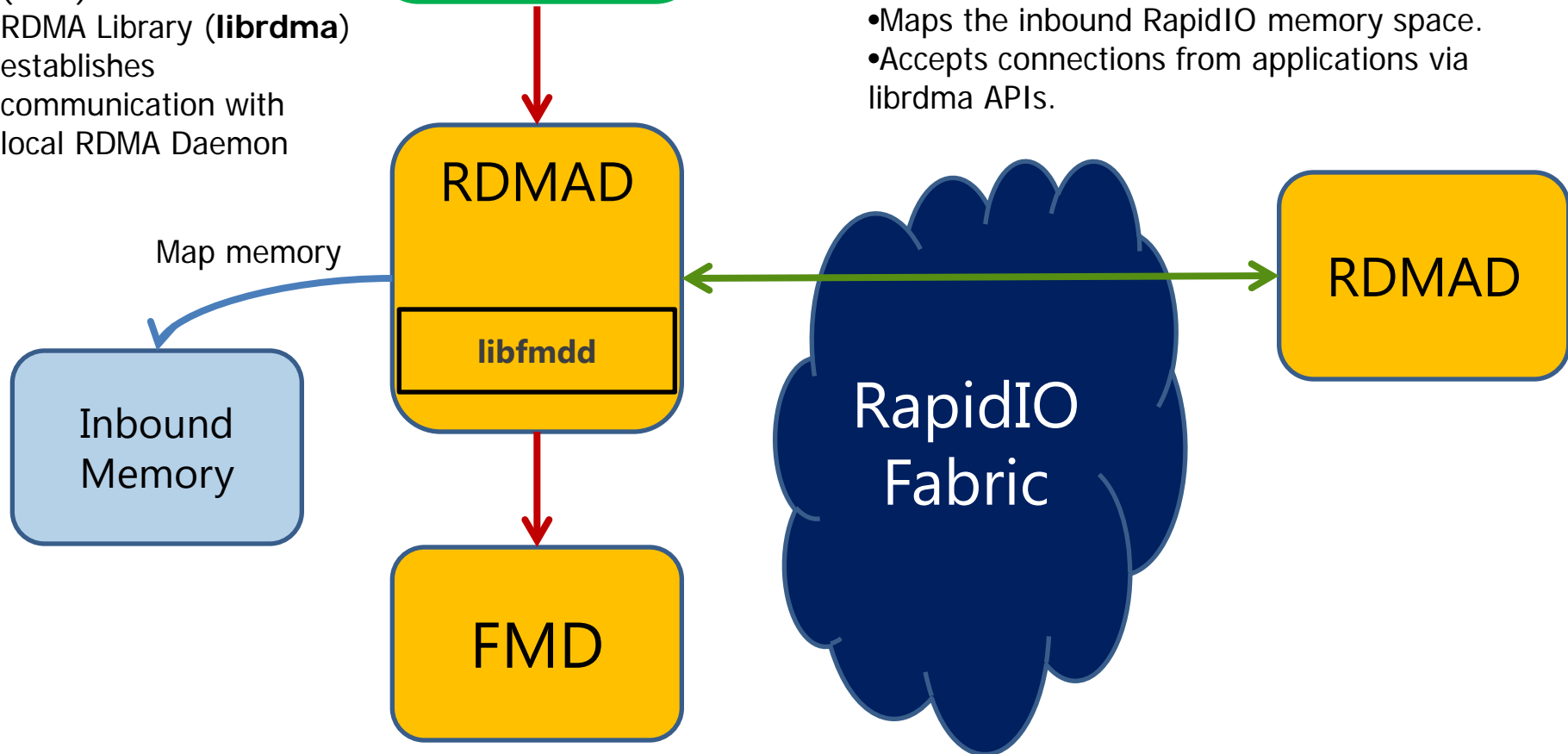
RDMA Daemon – Initialization

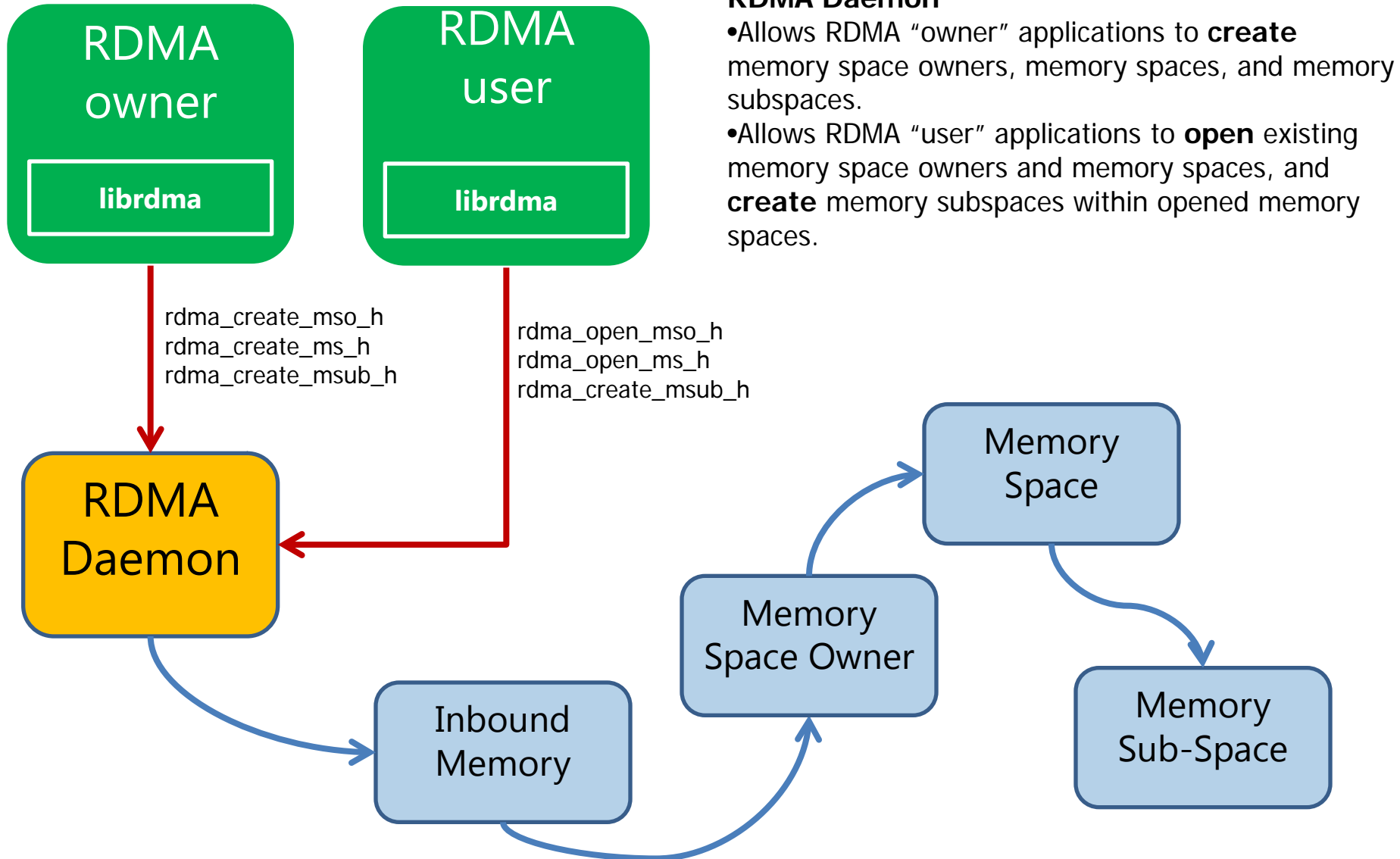
The Fabric Daemon Library (**libfmdd**) establishes communication with Fabric Management (FMD) daemon.
RDMA Library (**librdma**) establishes communication with local RDMA Daemon



RDMA Daemon

- Queries FMD to get a list of remote RDMA daemons.
- Establishes connections with remote RDMA daemons.
- Maps the inbound RapidIO memory space.
- Accepts connections from applications via librdma APIs.





RDMA Daemon

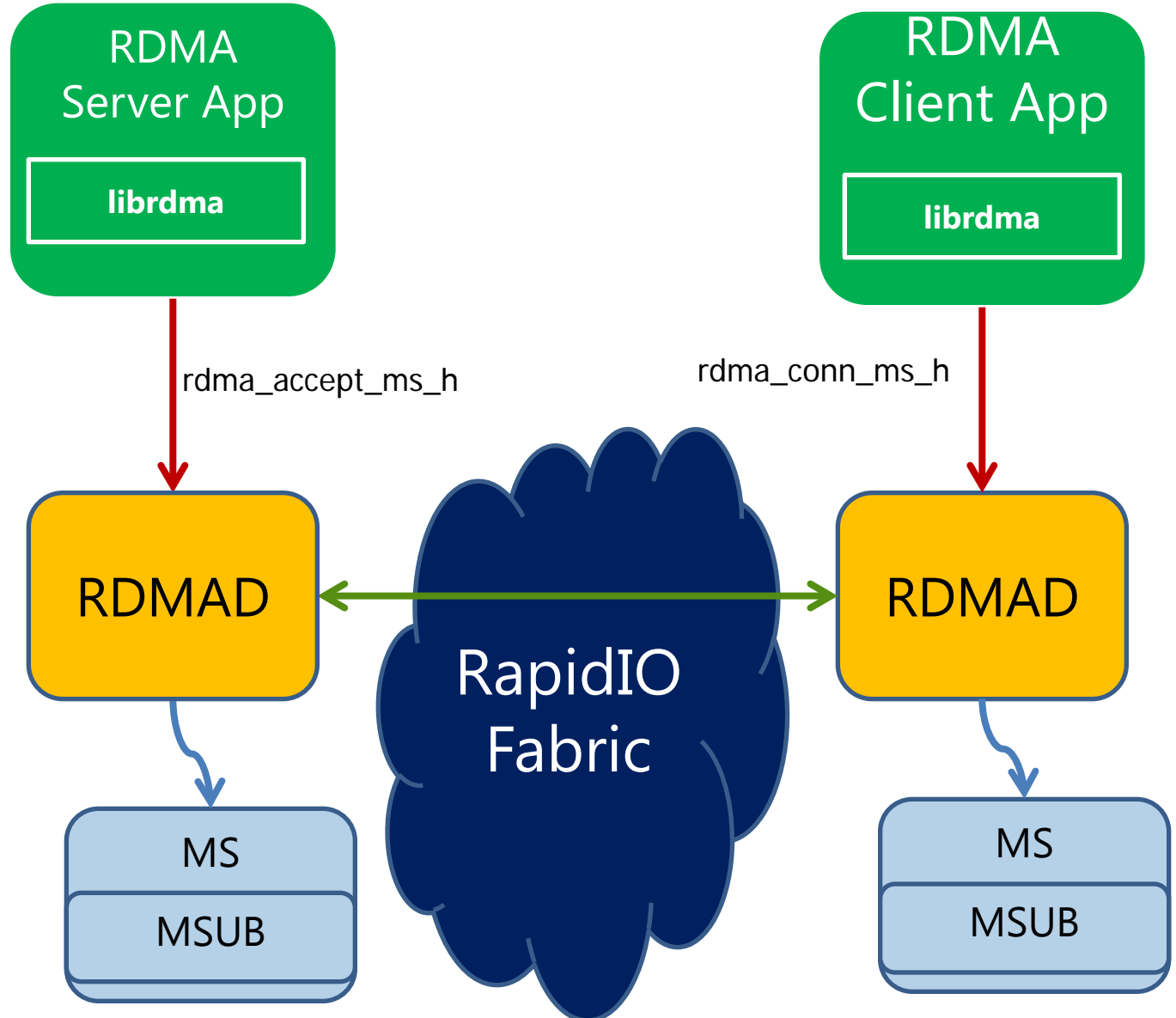
- Allows RDMA “owner” applications to **create** memory space owners, memory spaces, and memory subspaces.
- Allows RDMA “user” applications to **open** existing memory space owners and memory spaces, and **create** memory subspaces within opened memory spaces.

Server application calls `rdma_accept_ms_h()` on a memory space, and provides a memory subspace (msub) to client applications connecting to the memory space.

Client application calls `rdma_conn_ms_h()` and specifies the name of the remote memory space it wishes to connect to. It also provides an msub to the server.

`rdma_accept_ms_h()` is blocking.

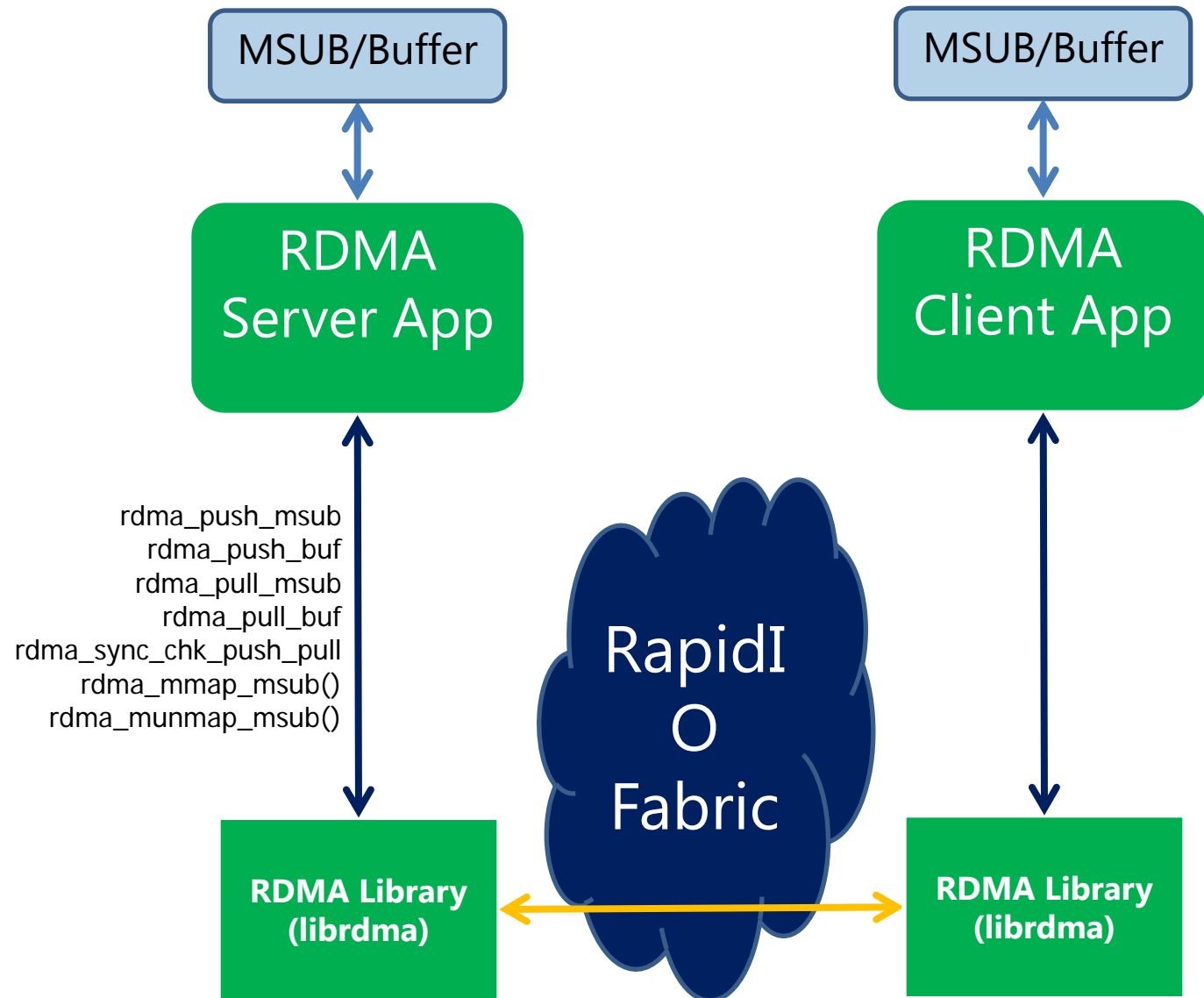
`rdma_conn_ms_h()` has an optional timeout.



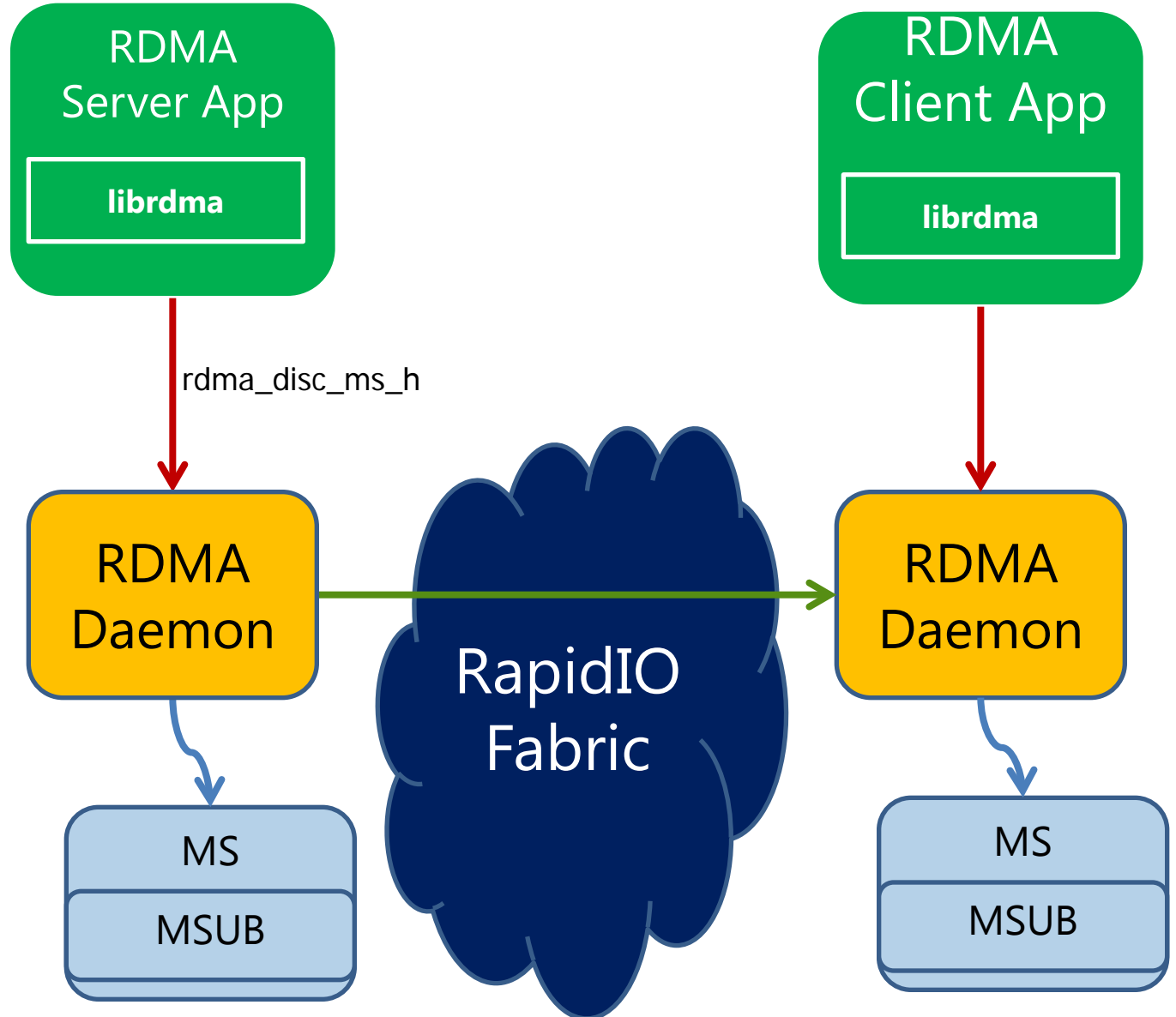
Once connection is established between a server app and a client appl, DMA transfers can go directly between the two sides using only routines in the RDMA library.

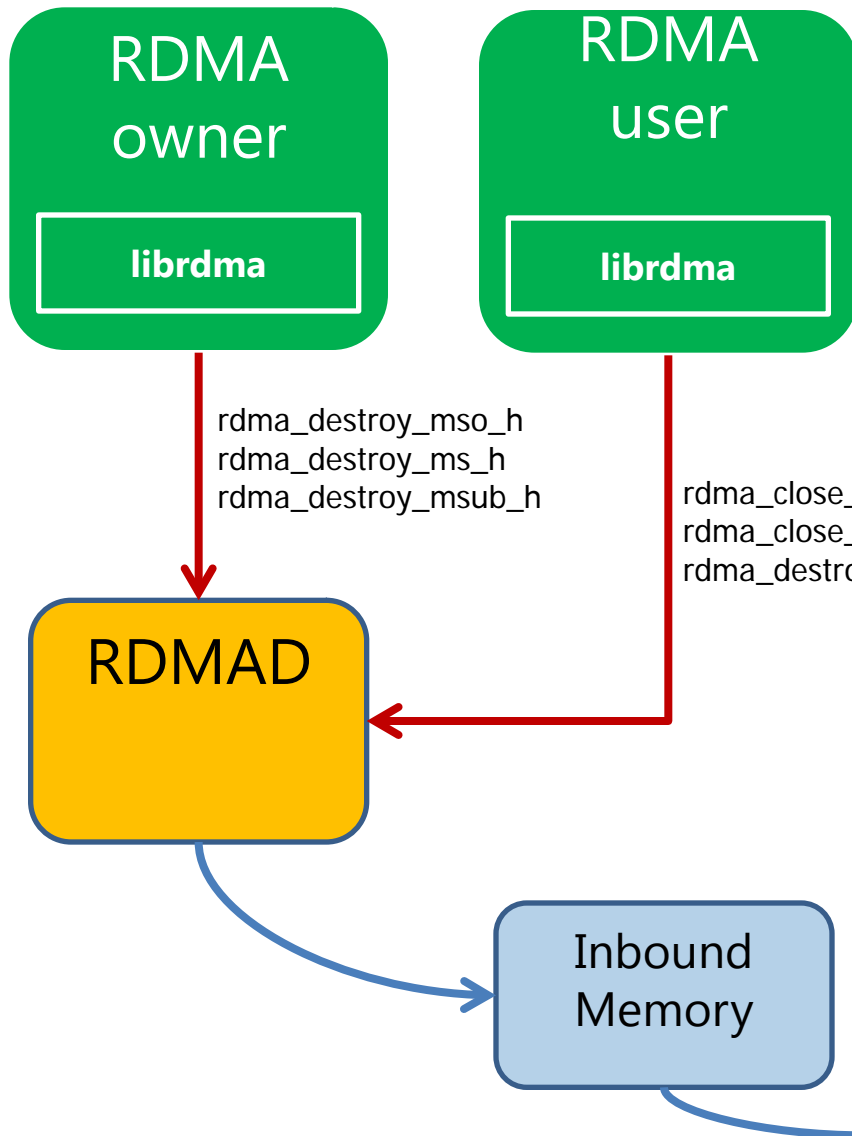
RDMA transfer routines provide push (write), and pull (read) capability in various modes (synchronous, fire-and-forget, and asynchronous).

Data can be exchanged using a memory subspace or a user buffer. For subspaces the space must first be mapped using `rdma_map_msub()`.



Client application calls `rdma_disc_ms_h()` and specifies the handle of the remote memory space it wishes to disconnect from.



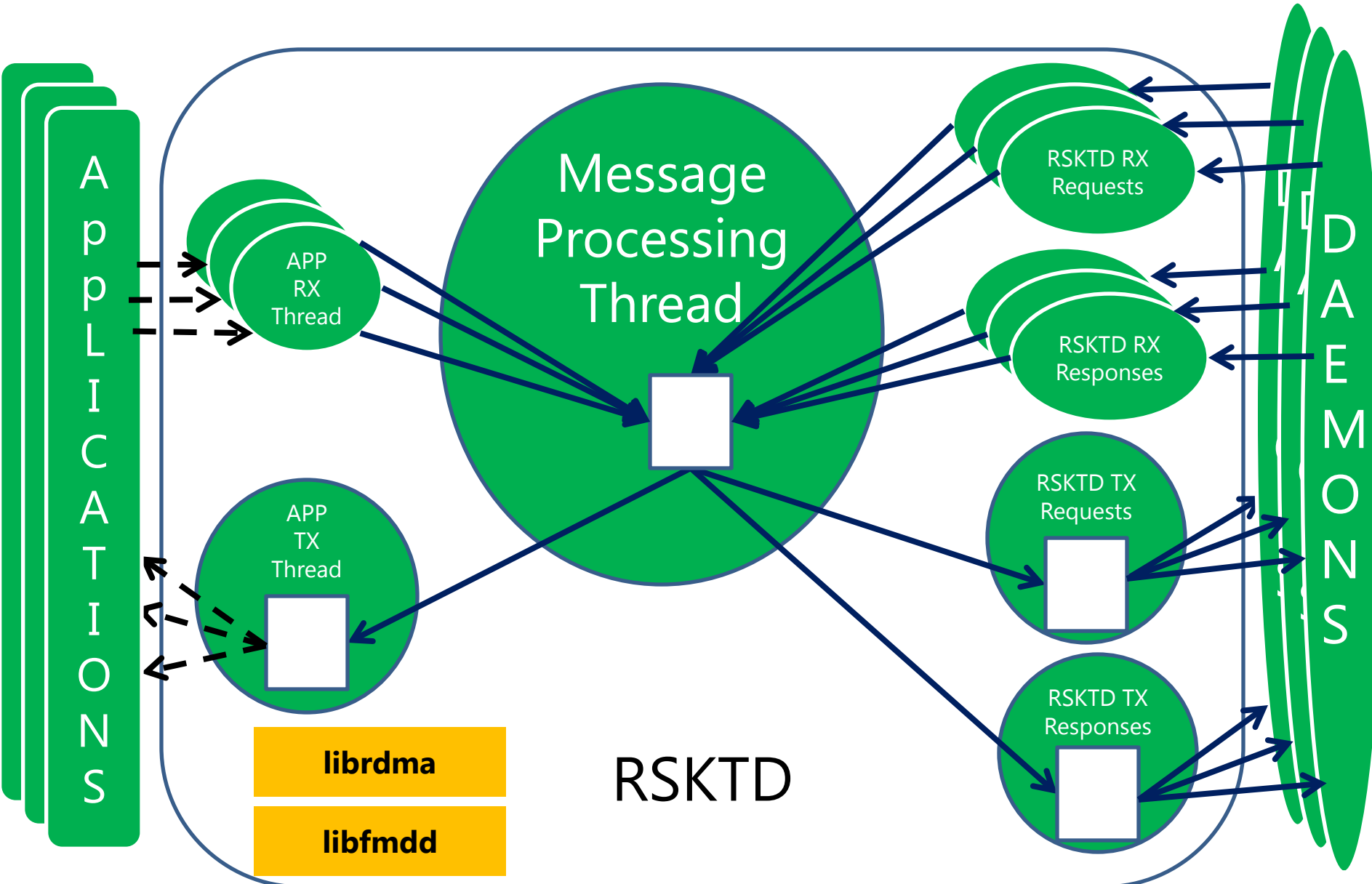


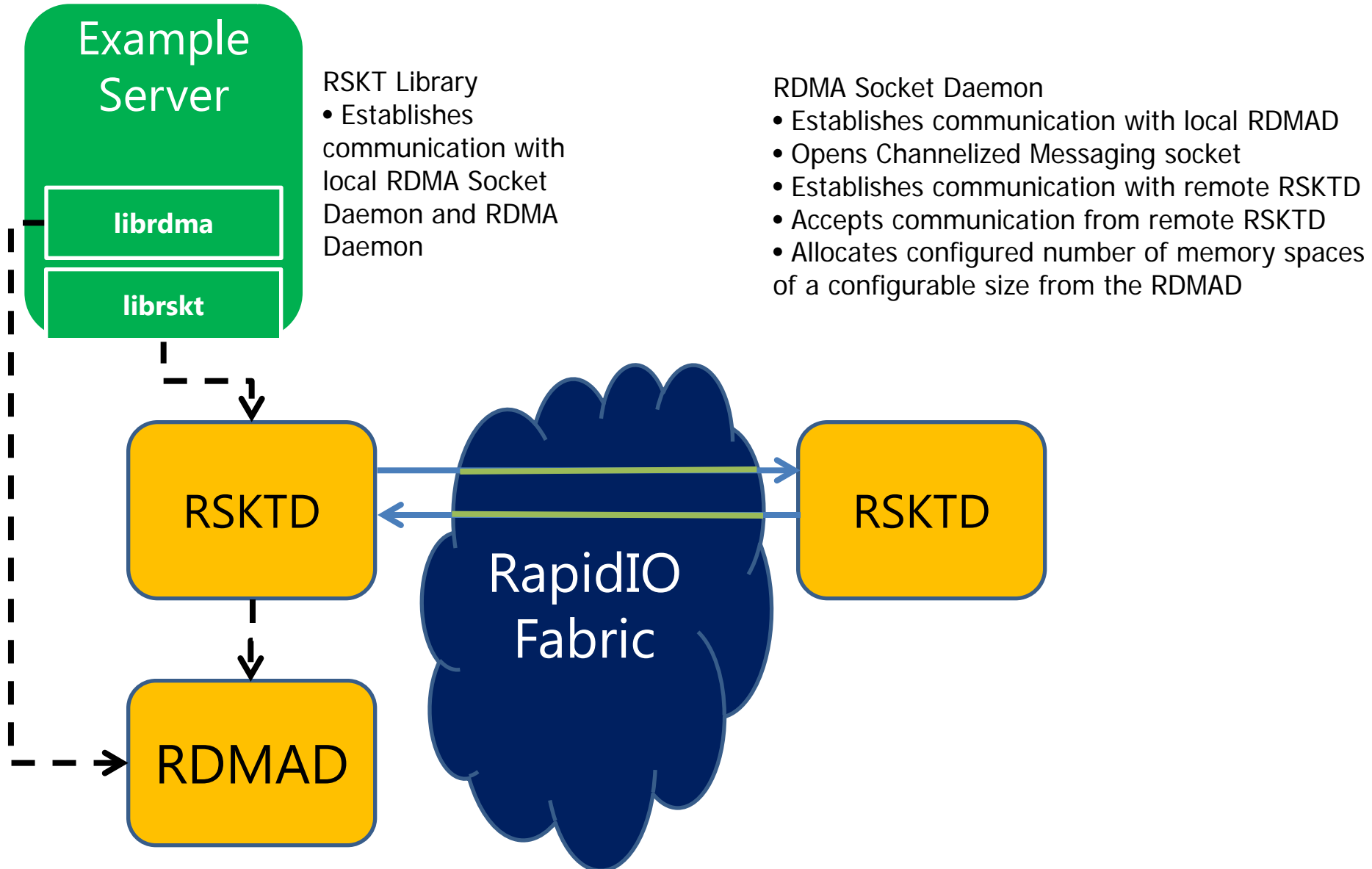
RDMA Daemon

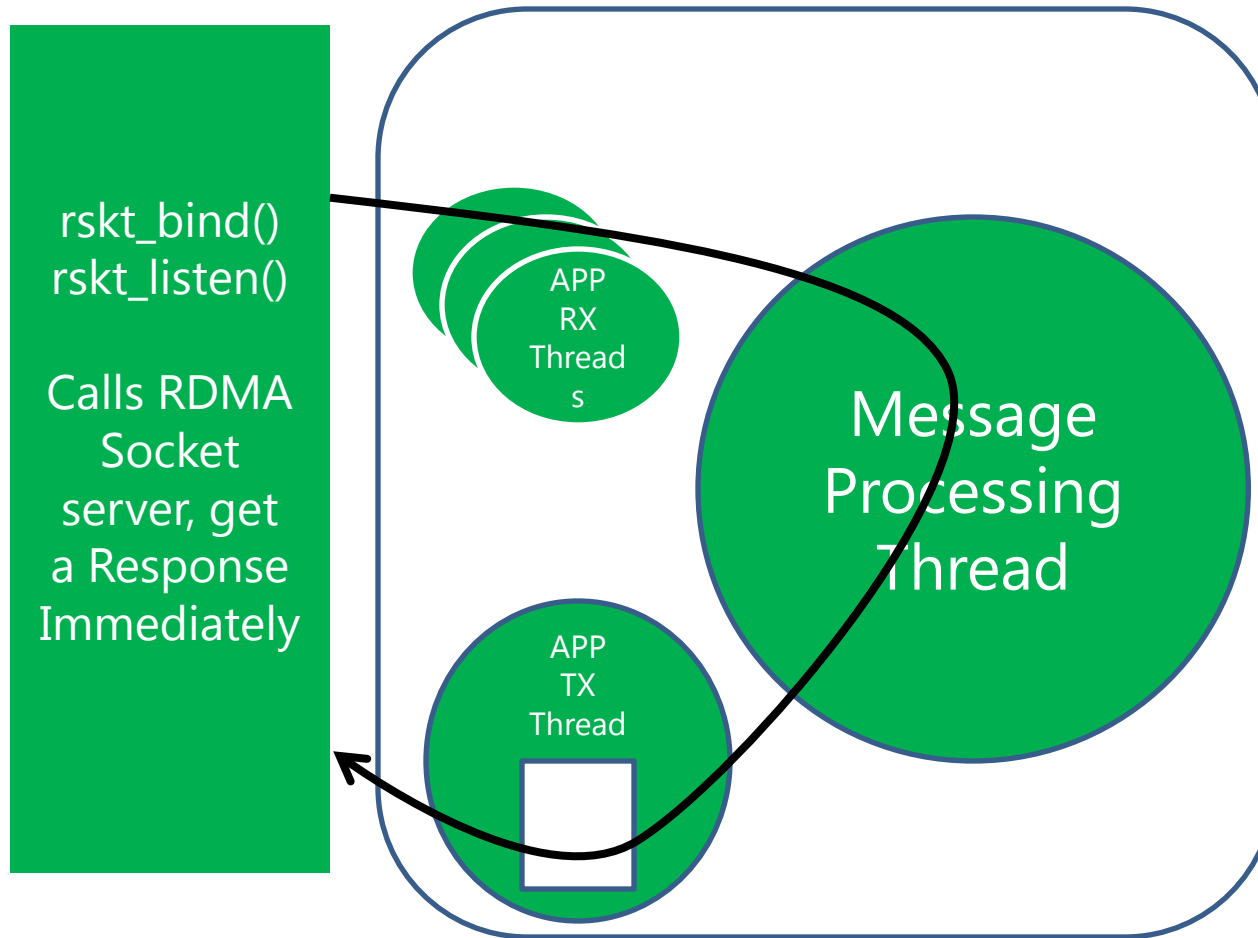
- Allows RDMA “owner” applications to **destroy** memory space owners, memory spaces, and memory subspaces.
- Allows RDMA “user” applications to **close** existing memory space owners and memory spaces, and **destroy** memory subspaces within opened memory spaces.
- Closing and destroying an mso are hierarchical operations.

RDMA SOCKETS (RSKTS) OVERVIEW

Routine	Description
<code>rskt_socket</code>	Create socket data structure and opaque handle
<code>rskt_bind</code>	Bind socket to socket number
<code>rskt_listen</code>	Prepare to accept connect requests on this socket
<code>rskt_accept</code>	Accept a connect request to this socket
<code>rskt_connect</code>	Connect to an RDMA socket on another RapidIO node
<code>rskt_write</code>	Write data to socket connection
<code>rskt_read</code>	Read data from socket connection, stream of bytes
<code>rskt_recv</code>	Read next record from socket connection
<code>rskt_flush</code>	Flush all writes for socket connection
<code>rskt_shutdown</code>	Flush and then close the socket connection
<code>rskt_close</code>	Close the socket connection, may or may not transfer all data before closure.

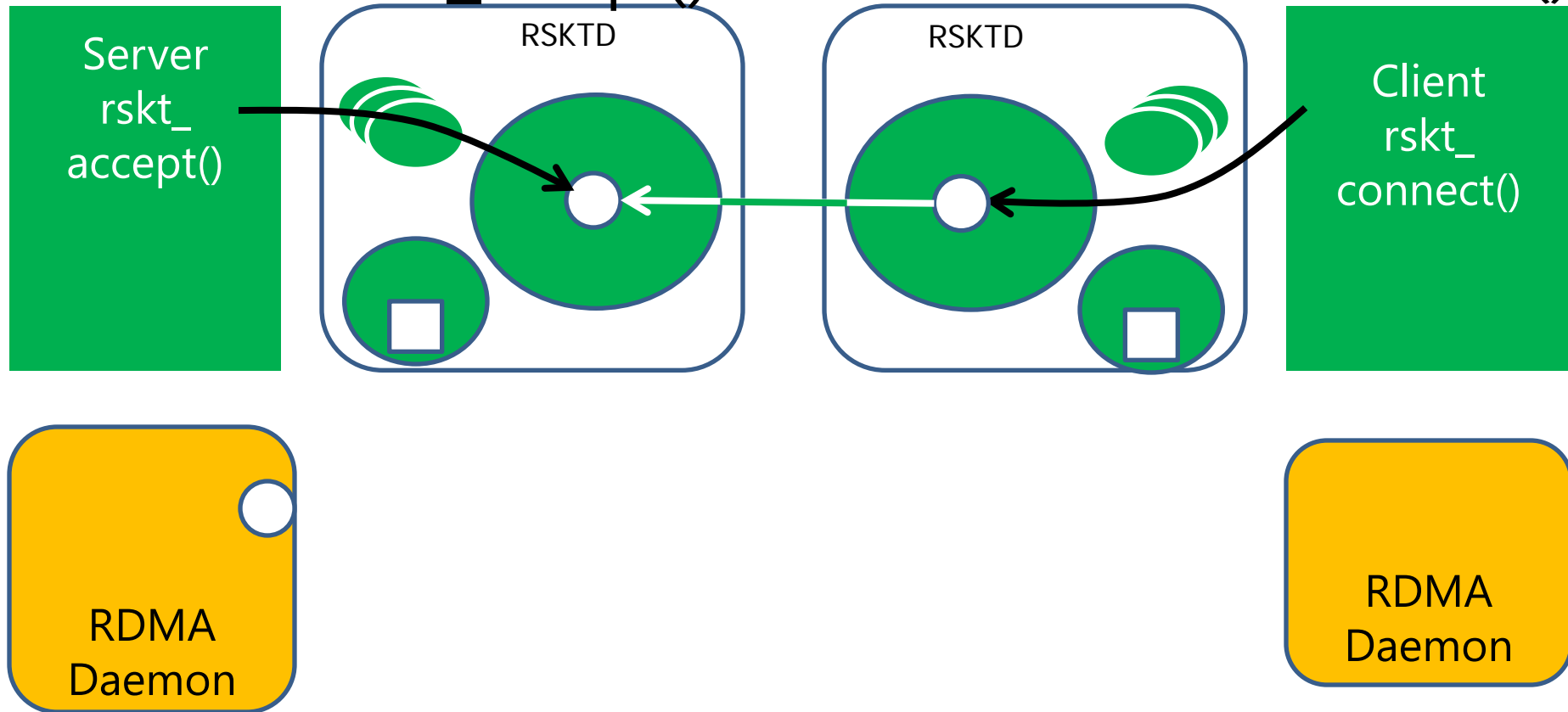


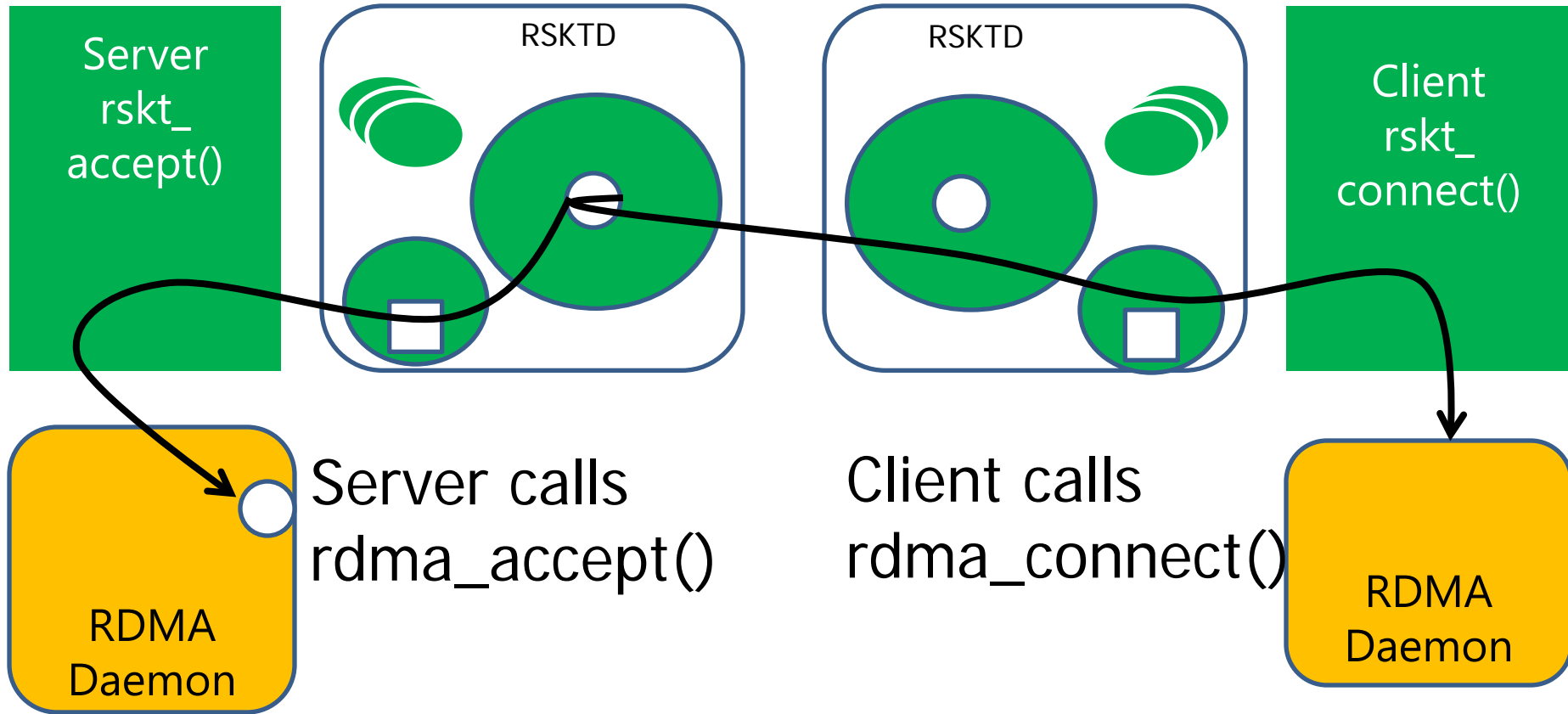


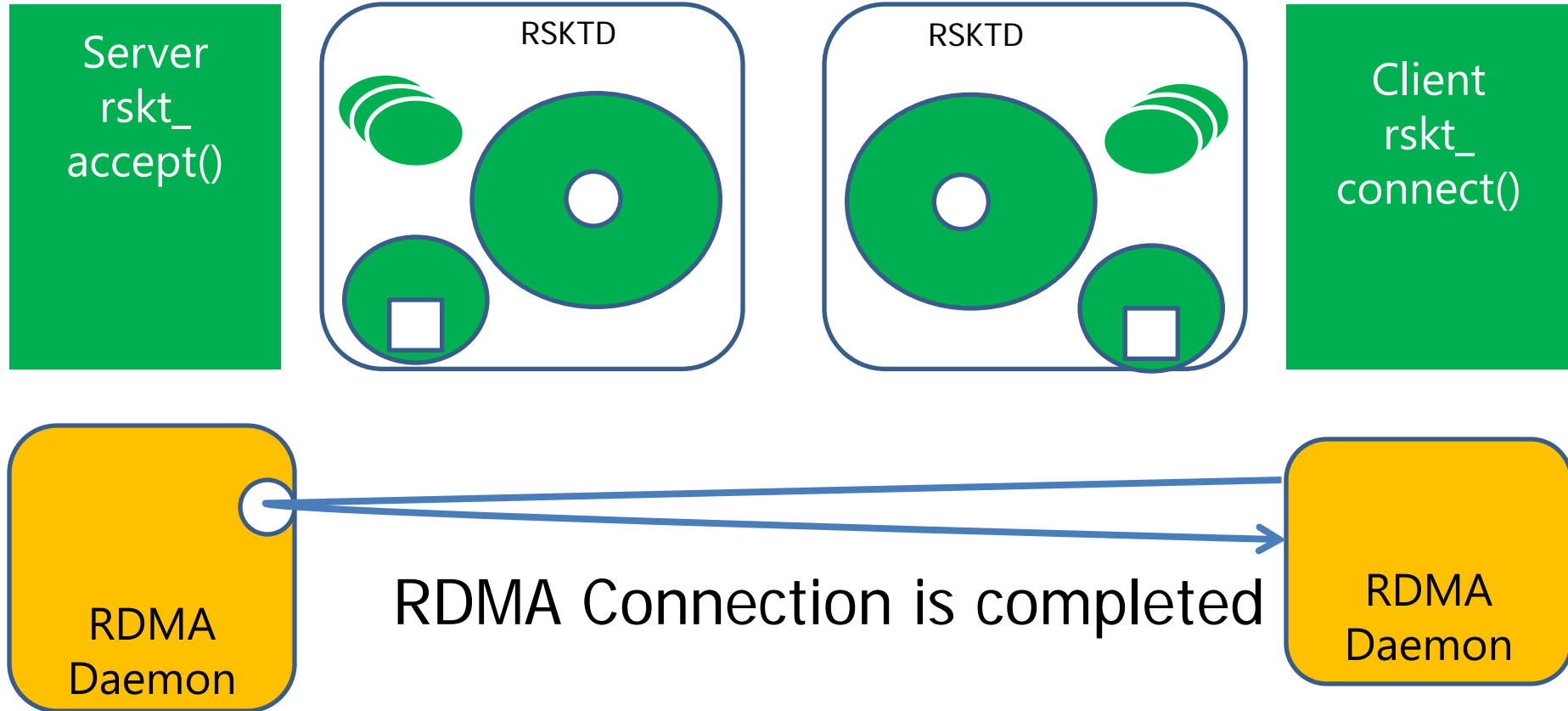


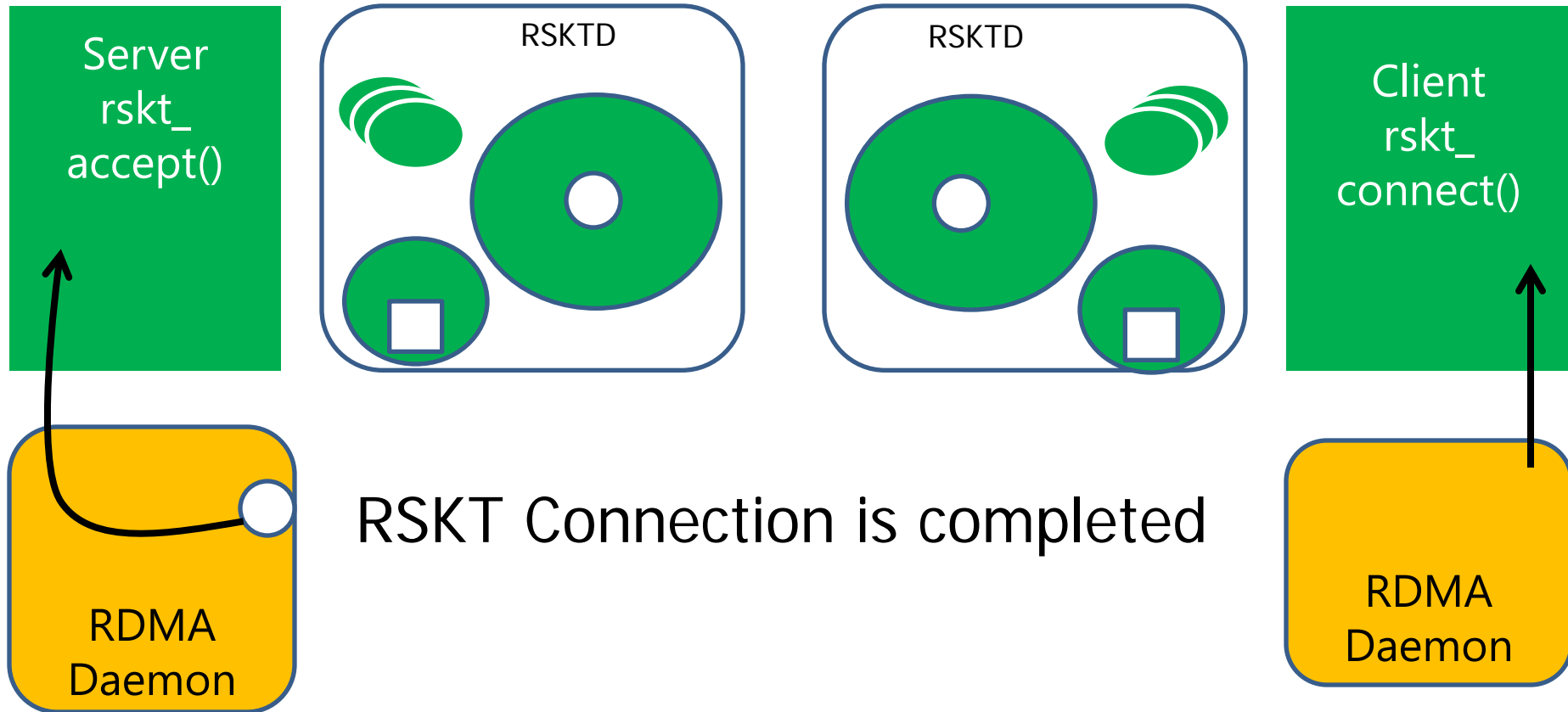
Server calls `rskt_accept()`

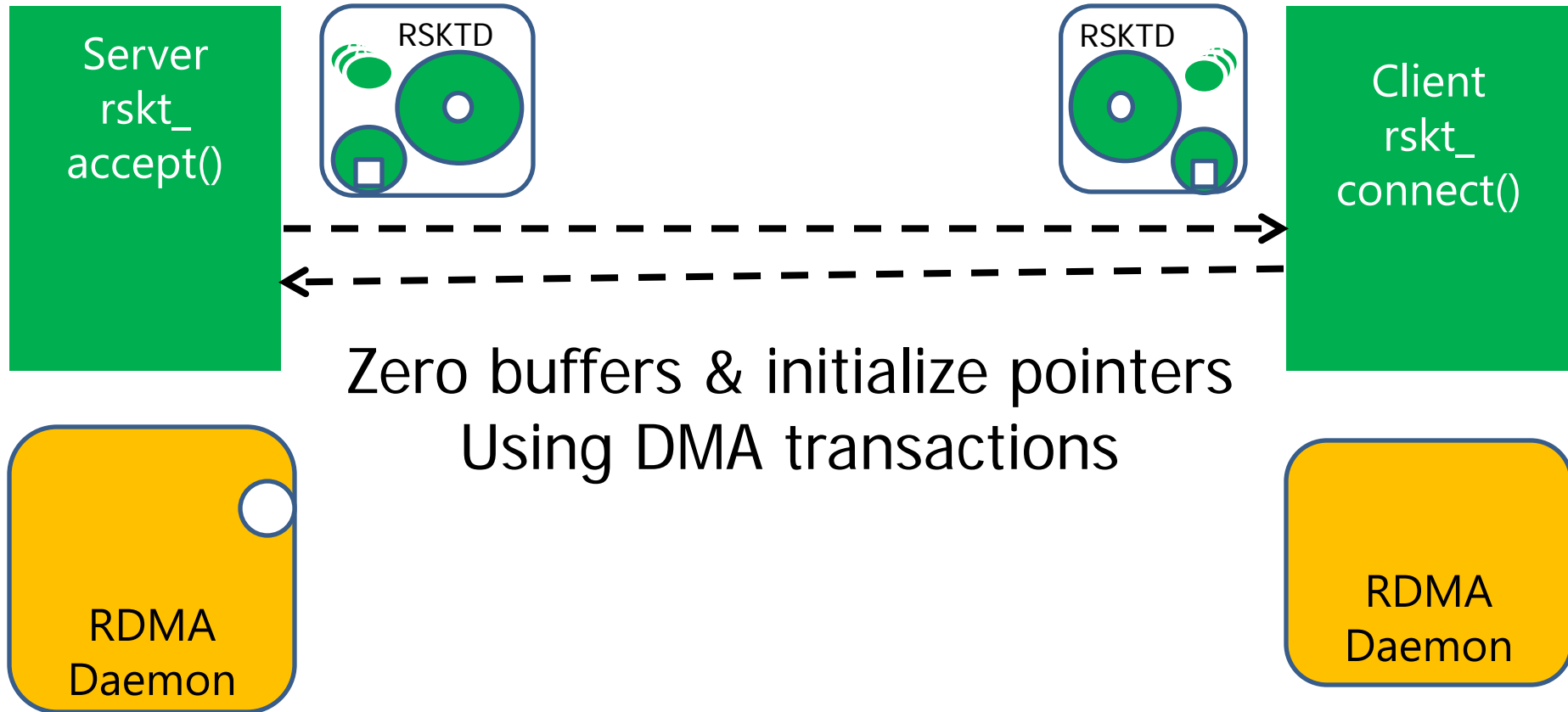
Client calls `rskt_connect()`

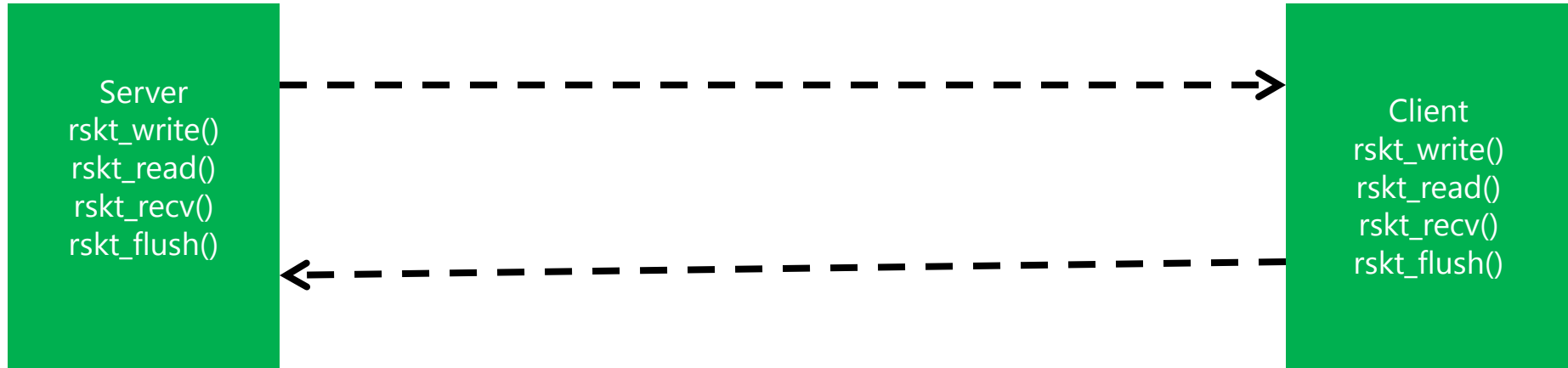






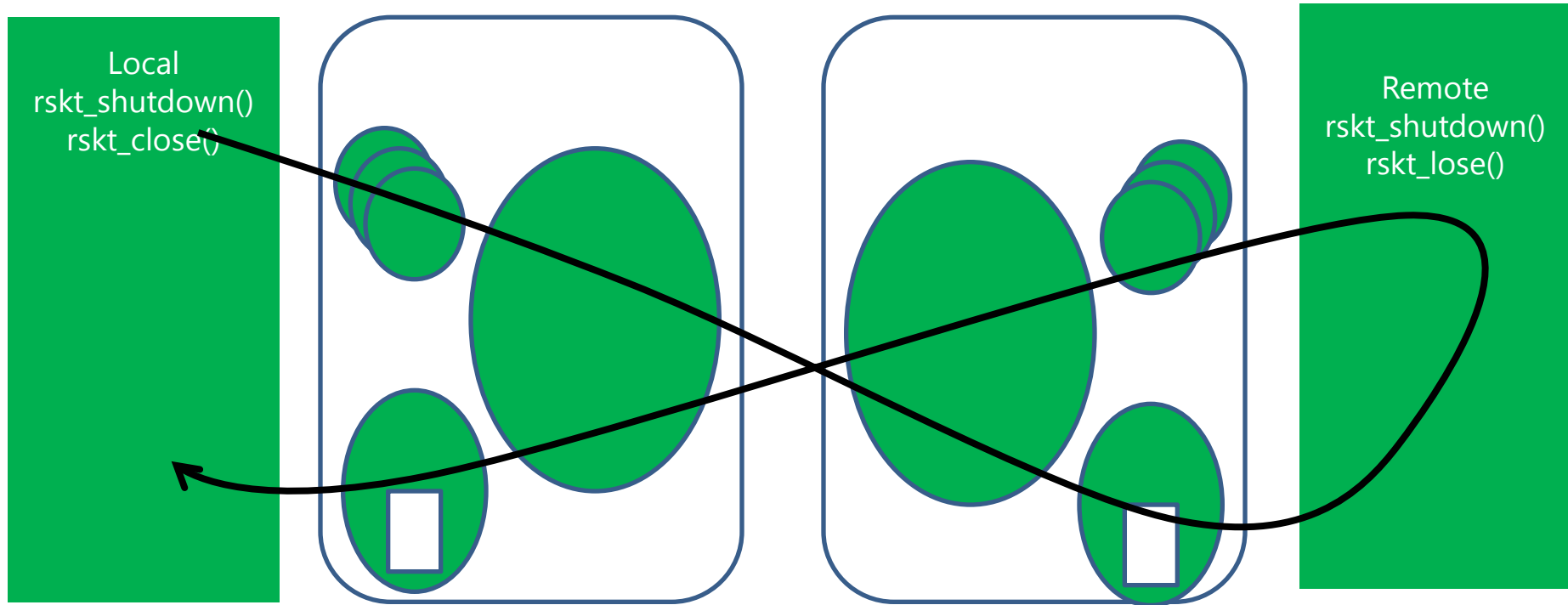






Write/Read/Recv/Flush all use
RDMA library routines and DMA.

No messaging required.



Shutdown and close are initiated

- By procedure call on a local client/server
- Automatically by failure of a DMA write/read/recv/flush request
- Automatically by closure of a POSIX socket with the local client/server
- Automatically by closure of the rio_cm_cdev connection with a remote RSKTD
- By request on a remote client/server

RAPIDIO SOFTWARE STACK DEBUG

Each API routine uses the “errno” global variable to report standard LINUX error codes on a failure.

In addition, each process and library reports status using the `rapidio_sw/common/include/liblog.h`

All log files are found in `/var/log/rdma`

- `fmd.log`
- `rdmad.log`
- `rsktd.log`

Logs can be displayed on each daemon using the “dlog” command.

The level of detail for all logs is controlled using the “levelog” command, available from the CLI for each daemon

- To access the CLI prompt for each daemon, use “screen –r <app>”
- For example, to connect to the RDMAD for debug, use “screen –r rdmad”

Additional debug commands are available from the CLI for process debug

The first step in reporting a problem is examining/transferring the log files found in /var/log/rdma/...

- fmd.log-<date>
- rdmad.log-<date>
- librdma.log-<date>
- rsktd_log.txt-<date>

RAPIDIO FILE TRANSFER UTILITY CHARTS SEPTEMBER 2015

Server

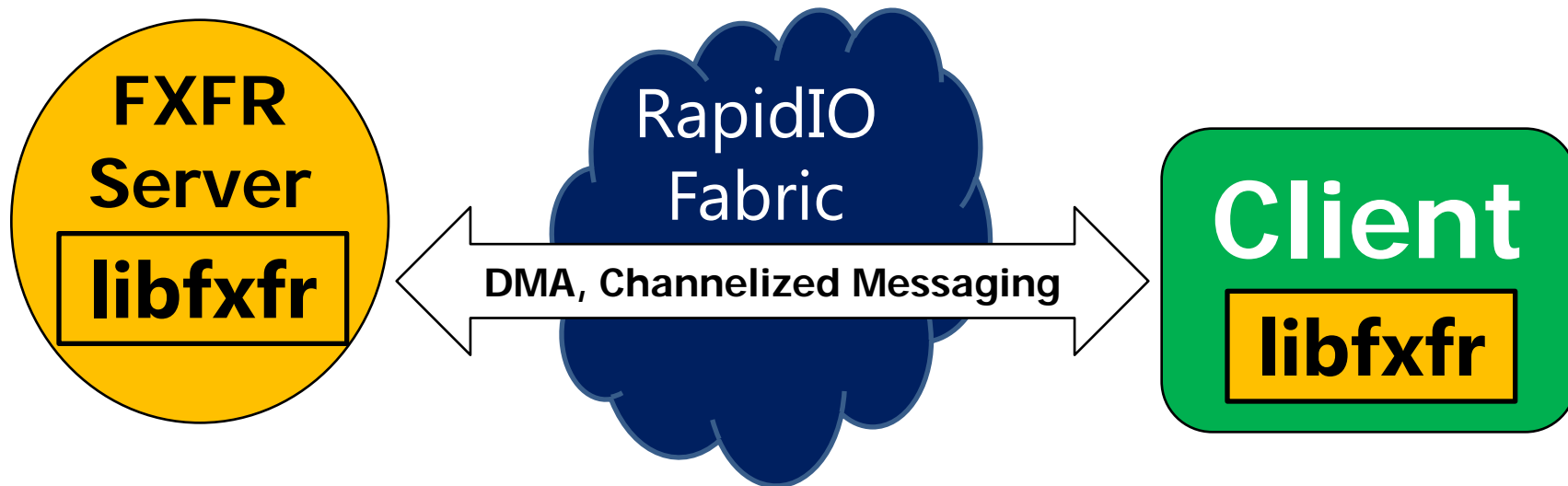
- Acts as a target for “send file” requests
- Manages externally available memory
- Manages multiple large (~2 MB) memory segments as “ping pong” buffers

Library (libxfr)

- exposes “send_file” routine to send a file to the server

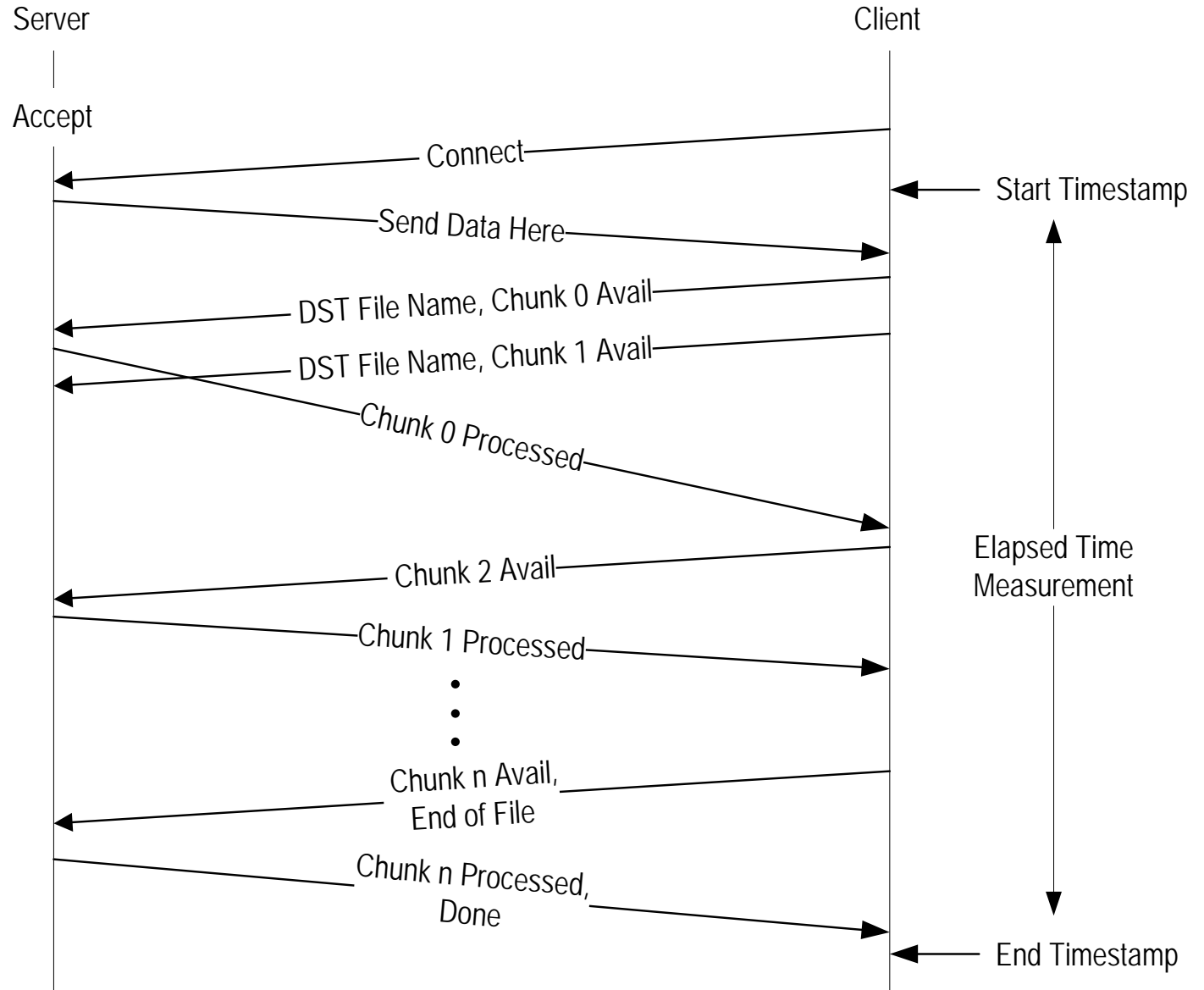
Performance

- Much better for large (multi megabyte) files, rather than small files
- RapidIO transfers can occur at line rate (1700 MBps), even with Kernel Driver
- Limited by the speed of the media
 - SAS can't handle 1700 MBps...



File	Description	Server	libfxfr
server.c	FXFR Server Main Loop	Yes	
libfxfr.h	Client interfaces for libfxfr		Yes
libfxfr_private.h	Server interfaces for libfxfr	Yes	Yes
fxfr_msg.h	Channelized messaging formats for libfxfr	Yes	Yes
fxfr_rx.c	Implementation of file receive	Yes	
fxfr_tx.c	Implementation of send_file		Yes

File Transfer Message Sequence Chart



GOODPUT TOOL OVERVIEW

Measurements

- Goodput
- Latency

Transaction types

- Direct I/O
- DMA
 - Sync, async, fire-and-forget
- Messaging

Drivers

- libmport
- “Demo” user mode driver

Single process, multiple threads

- Command line interpreter (CLI) thread
- 12 worker threads

CLI functions

- start/stop/configuration of worker threads
- dispatch tasks to workers
- measurement display

Bash Scripts

- generate CLI scripts for automated measurements
- summarize logged results into simple tabular format

To compile goodput/ugoodput:

Type "make all" in rapidio_sw/utils/goodput directory

Goodput documentation

- Type "doxygen doxyconfig" in rapidio_sw/utils/goodput
- Open rapidio_sw/utils/goodput/html/index.html file with browser OR
- Refer to goodput/inc/goodput_intro.h