



## SELinux support for RDMA

Dan Jurgens

Sept 23, 2015

- **SELinux is a Mandatory Access Control (MAC) scheme for Linux**
  - Central policy is loaded upfront into the kernel
    - Standard policies are typically provided by the Linux distribution
  - Applications cannot override or modify this policy
  
- **Benefits**
  - Differentiate a user from the applications that the user runs
  - Restrict application access only to what is required to perform its task
  - Allow granular policy segregation
  - Example
    - Run 2 instances of a Web Server: “top-secret” and “standard”
    - Each server can only
      - Receive traffic from specific network interfaces
      - Open sockets on specific ports
      - Serve files from specific directories
      - Communicate only with specific peer addresses
  
- **Type enforcement is the main security mechanism used by SELinux**



- Applies to all user-visible kernel entities
  - E.g., processes, files, IPC objects, sockets
- Each entity is associated with:
  - A security descriptor
    - Assigned upon creation
    - May be modified based on policy
  - A class and a set of operations
    - Stems from the type of object
    - E.g., a socket can send() and recv()
- TE defines what a <subject> can do on an <object> based on their security descriptors
  - Specified by a policy of access rules
  - Enforced when accesses are made
- Security descriptors
  - Identify the user, role, type, and optionally security level+class of an object
  - Specified by a variable-length string: “user:role:type[:level]”
- Policy rules
  - Specify which source tag can access which target tag and for what operations
    - E.g., “allow source\_t target\_t:class { [op1] [op2] ... }”
  - Typically, only the ‘type’ (a.k.a ‘Domain’) portion of the tag is mentioned

## ■ Network object labeling

- Interfaces
  - E.g., eth2
  - Used in the past for packet tagging
    - Today packets are tagged by network traffic labeling
- Nodes
  - Label IPv4/6 addresses and network masks
- Ports
  - Label TCP/UDP port numbers
- Sockets
  - Usually inherit the security descriptor of the creating process

## ■ Network traffic labeling

- Internal labeling
  - Tag traffic according to local OS policies
  - The de-facto standard is SECMARK
    - Extends standard iptables/netfilter to mark packets with security descriptors
- Labeled networking
  - Labels on the wire
  - Each local system interprets the label to enforce is MAC policies
  - Supported schemes
    - Labeled IPsec
    - CIPSO

## ■ SELinux network policies define

- What a process can do with network objects
  - For example: allow 'ftp\_t' (the FTP process) to bind a socket to 'ftp\_data\_port' (TCP port 20)
- What traffic a process can send/receive

## ■ Policy enforcement

- Object policies are enforced during system calls
- Traffic policies are enforced per-packet

- Verbs objects (PDs, QPs, CQs, SRQs, MRs, etc.)
  - Do not have well-known names
    - The user doesn't even control them
  - **Observation: granular MAC control policies for Verbs resources doesn't make much sense**
    - **E.g., allow a process to modify QP253 to RTR???**
  
- RDMA-CM objects (RDMA IDs)
  - Similar to sockets
    - ServiceID port-spaces map TCP/UDP ports spaces
    - Similar semantics and operation
  - However
    - Actual data path governed by Verbs objects
    - There are RDMA applications that don't use RDMA-CM at all
  - **Observation: RDMA-CM security policies may provide socket-like MAC control for CM service IDs**
    - **However, it does not provide a stronger security model**
    - **Simple policies are still required for P\_Key enforcement**

- Interface objects (RDMA devices and ports)
  - Associated with every communication end-point (QP)
  - **Observation: may be used for object-based policies**
    - But benefit is not clear
    - Network labeling (see below) is a much better alternative
  
- Node objects (GIDs)
  - May be treated similarly to IPv6 addresses
  - Not always used
  - Performance penalty for UD sends (to verify address handle)
  - Does not cover UD reception
  - **Observation: GIDs do not seem to be suitable for an object-based policy**

- Applications initiate IO directly on HW endpoints (QPs)
  - HW generates packets
    - Kernel is completely by-passed
  - **Observation: arbitrary internal labeling or labeled networking is not an option**
  
- RDMA traffic labeling must stem from HW architectural attributes
  - Network addresses (GIDs, LIDs) are not suitable
  - Queue keys (Q\_Keys) are not suitable
  - **Observation: partitions are the natural candidate**
    - P\_Key values are held in HCA partition tables
      - Populated by privileged network Subnet-Manager (SM)
    - P\_Keys are carried on the wire of every data packet
      - The only exception is subnet datagram packets (SMPs), which are not accessible to applications
    - Every QPs is associated with a P\_Key value
      - Determined by an index into the partition table
    - Partitions are strictly enforced at all times

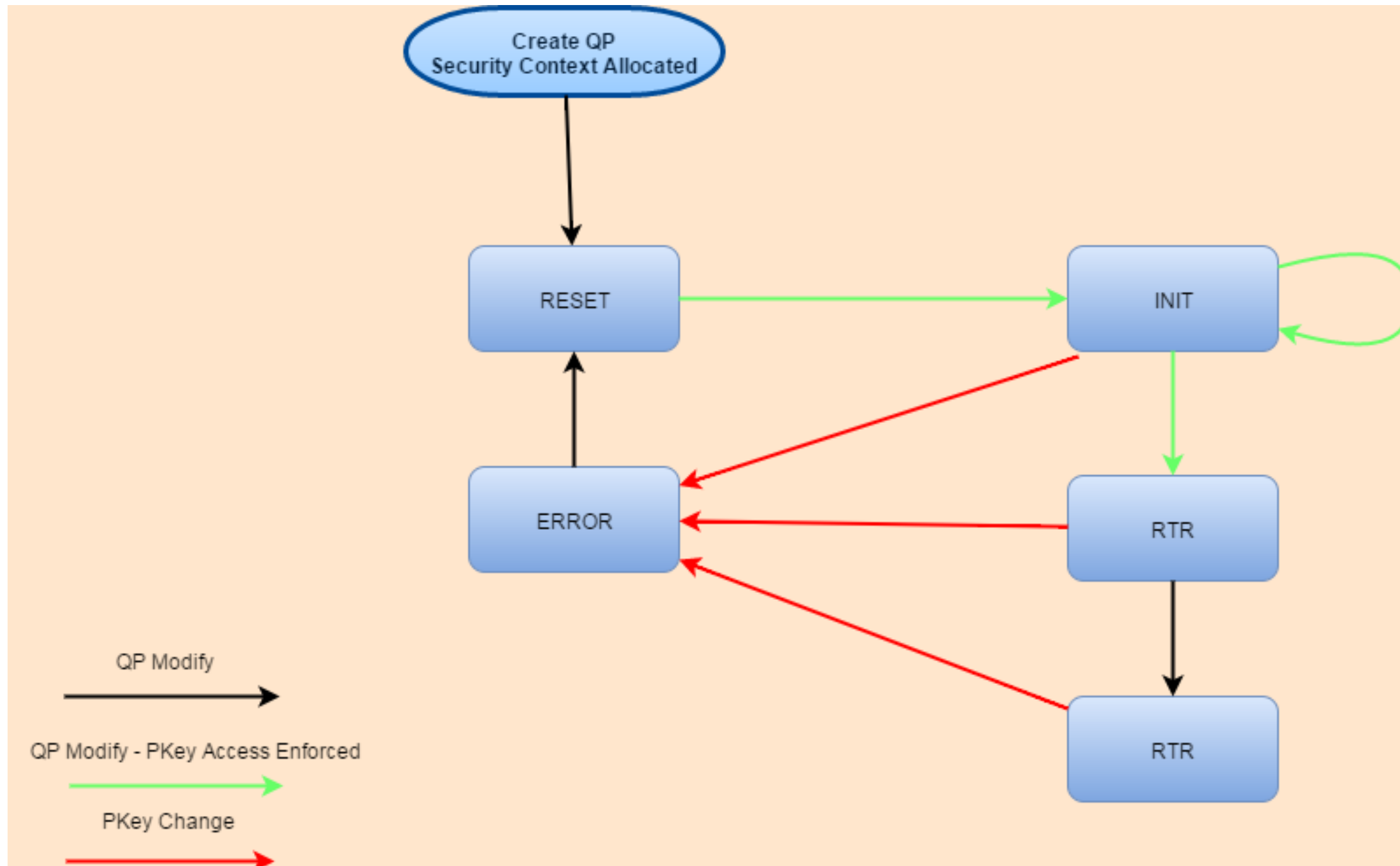
- RDMA network security based on partitioning
  - Host kernels control the association of P\_Key values with security descriptors
  - SM configuration and P\_Key assignment determined by network administrator
    - SELinux policies may be used to control which processes can access the SMI
      - E.g., only SM and tools processes
  
- Object labeling
  - Associate QPs and RDMA IDs with a security descriptor
    - Inherited by the creating process in the absence of a specific policy
  - P\_Key value labeling
    - Associates a P\_Key value with a security descriptor
    - System object descriptors are a good example (like network interfaces or nodes)
      - “system\_u:object\_r:rdma\_partition\_default\_t”
      - “system\_u:object\_r:rdma\_partition\_topsecret\_t”
  - Other objects not labeled



- **Traffic labeling**
  - Network labeling based on P\_Key values
  
- **Policies**
  - Allow a process access to a P\_Key value
    - E.g., “allow hpc\_default\_t rdma\_partition\_default\_t : rdma\_partition { modify }”, where
      - ‘hpc\_default\_t’ is the QP / RDMA\_ID domain (type) inherited from the creating process
      - ‘rdma\_partition\_default\_t’ is a partition security descriptor domain
      - ‘rdma\_partition’ indicates that the subject is of partition type
      - ‘modify’ indicates that the QP is allowed to modify to reference the corresponding partition tag
  - Allow a process access to the SMI
  
- **Partition enforcement**
  - QP partitioning enforced at all times
    - Upon QP creation, a violation shall result in an immediate error
    - During runtime
      - Any runtime violation due to policy changes or P\_Key value changes shall transition the QP into ERROR state
  - RDMA-ID
    - All ingress/egress CM MADs shall be checked according to the partition policy
    - Any violation shall result in an immediate packet drop
  - umad interface for GMPs

- Enforce access to IB partitions with access controls on Pkeys.
  - Label Pkeys in the SELinux policy.
  - When QP partition key settings are changed access to the Pkey is enforced.
  - QP inherits it's creating tasks' security context.
  
- Changes required.
  - New LSM hooks for managing the QP security field and checking Pkey permissions.
  - Add SELinux kernel support for pkey labels and access control.
  - Changes in ib\_core kernel module to enforce QP access to Pkeys and incoming and outgoing management datagrams.
    - Changes are device independent.
  - SELinux user space utilities modified to support Pkey labeling in the policy language.
  - Refpolicy changes to label the Pkeys.

# QP State Transitions



- Added a security field to the `ib_qp` structure.
- Added new LSM and SELinux hooks to allocate the security structure and free it.
  - Allocation happens at QP create, and it's security context is set to that of the calling task.
  - De-allocation happens during QP destroy.
- Added support to SELinux policy language to label `P_Keys` based on their value.
  - Similar to port labeling, allowing for individual `P_Keys` or ranges.
    - `pkeycon <pkey_number> gencontext(<label>)`
    - `pkeycon <low_pkey_number>-<high_pkey_number> gencontext(<label>)`
    - Similar to port labeling.
- Added one new LSM and SELinux hook that takes the `P_Key` value and a QP security context to check for permission.
  - The hook is executed in `ib_modify_qp` whenever `IB_QP_PKEY_INDEX` is set in the attribute mask.
  - It is also run against all QPs on a port when the `P_Key` table changes
    - This requires keeping a list of QP using a particular `P_Key` index on each port.
    - If permission is not allowed for the new `P_Key` in the index the QP is moved to the error state.



```
attribute pkey_type;  
attribute protected_pkey_type;  
attribute unprotected_pkey_type;
```

```
type pkey_t, pkey_type;  
sid pkey gen_context(system_u:object_r:pkey_t,s0)
```

```
type staff_allowed_pkey_t, pkey_type, protected_pkey_type;  
type admin_allowed_pkey_t, pkey_type, protected_pkey_type;  
type default_pkey_t, pkey_type, unprotected_pkey_type;
```

```
pkeycon 65535 gen_context(system_u:object_r:default_pkey_t,s0)  
pkeycon 32769 gen_context(system_u:object_r:staff_allowed_pkey_t,s0)  
pkeycon 32770 gen_context(system_u:object_r:admin_allowed_pkey_t,s0)
```

```
allow sysadm_t default_pkey_t:rdma_pkey modify;  
allow staff_t default_pkey_t:rdma_pkey modify;
```

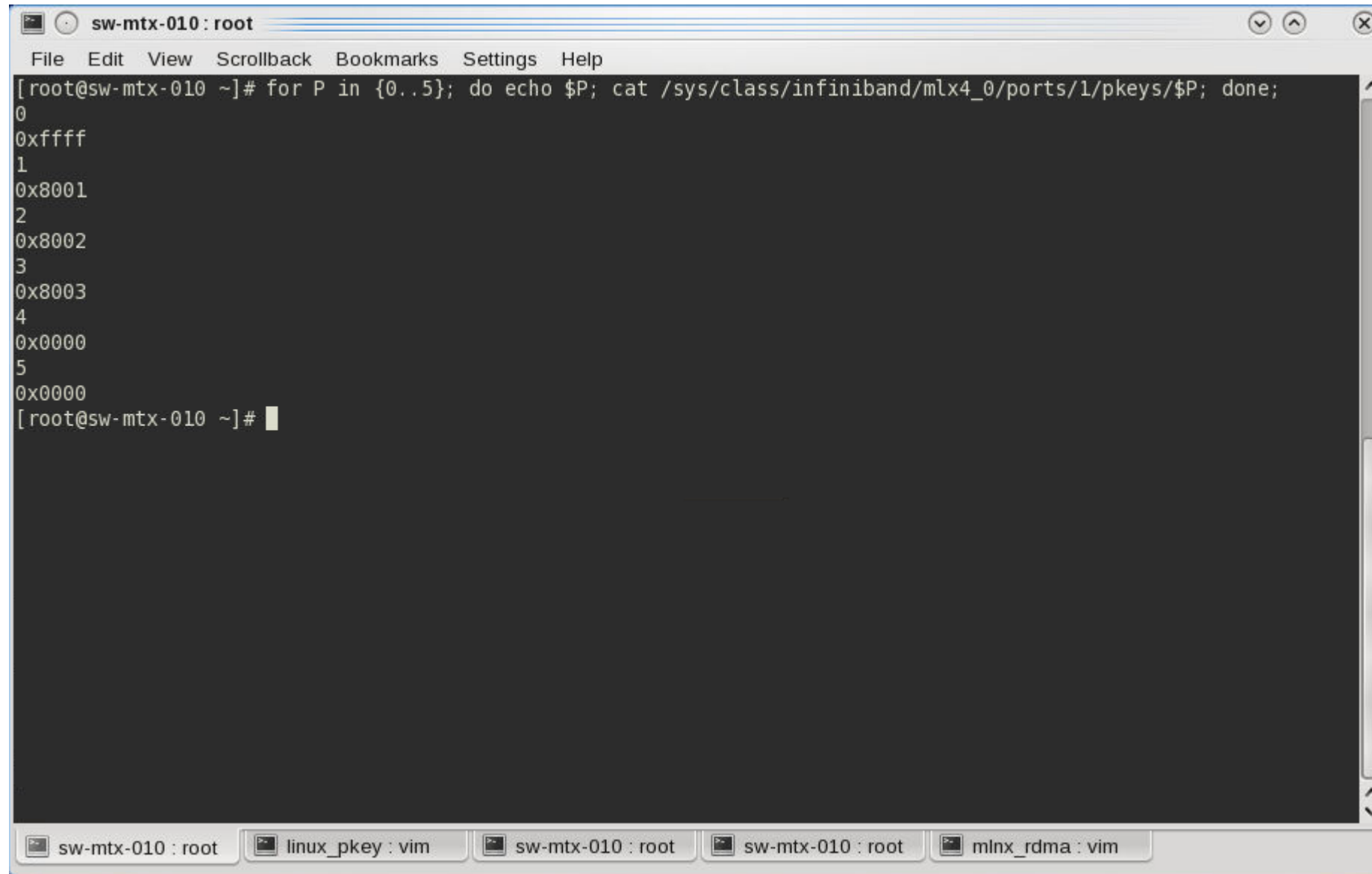
```
allow sysadm_t admin_allowed_pkey_t:rdma_pkey modify;  
allow staff_t staff_allowed_pkey_t:rdma_pkey modify;
```

- First new LSM hooks must be added to the upstream Linux kernel.
- After that SELinux kernel and users space changes can be submitted.
- After that the reference policy changes to provide default rules for Pkey access can be submitted.

- Two roles
  - Staff\_r
  - Admin\_r
  
- Four available partitions
  - Default (0xFFFF) – both allowed
  - Staff allowed (0x8001)
  - Admin allowed (0x8002)
  - Neither allowed (0x8003)



# Pkey Table



```
sw-mtx-010 : root
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# for P in {0..5}; do echo $P; cat /sys/class/infiniband/mlx4_0/ports/1/pkeys/$P; done;
0
0xffff
1
0x8001
2
0x8002
3
0x8003
4
0x0000
5
0x0000
[root@sw-mtx-010 ~]#
```

The terminal window displays the output of a shell script that iterates over the values 0 through 5. For each value, it prints the value and the contents of the file `/sys/class/infiniband/mlx4_0/ports/1/pkeys/$P`. The output shows that for values 0, 1, 2, and 3, the pkey values are 0xffff, 0x8001, 0x8002, and 0x8003 respectively. For values 4 and 5, the pkey values are 0x0000.

# Admin to Staff on Default Partition



```
sw-mtx-010 : root <2>
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:staff_r:staff_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=0 --ib-port=1 sw-mtx-010
-----
RDMA_Write BW Test
Dual-port      : OFF      Device      : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ   : OFF
TX depth       : 128
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs   : OFF
Data ex. method : Ethernet
-----
local address: LID 0x01 QPN 0x0211 PSN 0xf76897 RKey 0x30011a00 VAddr 0x007f4f50d06000
remote address: LID 0x02 QPN 0x0210 PSN 0x5fdcaf RKey 0x30011900 VAddr 0x007fddd781f000
-----
#bytes  #iterations  BW peak[MB/sec]  BW average[MB/sec]  MsgRate[Mpps]
65536   190700       0.00             5966.21             0.095459
-----
[root@sw-mtx-010 ~]#

sw-mtx-010 : root
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=0 --ib-port=2
*****
* Waiting for client to connect... *
*****
-----
RDMA_Write BW Test
Dual-port      : OFF      Device      : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ   : OFF
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs   : OFF
Data ex. method : Ethernet
-----
local address: LID 0x02 QPN 0x0210 PSN 0x5fdcaf RKey 0x30011900 VAddr 0x007fddd781f000
remote address: LID 0x01 QPN 0x0211 PSN 0xf76897 RKey 0x30011a00 VAddr 0x007f4f50d06000
-----
#bytes  #iterations  BW peak[MB/sec]  BW average[MB/sec]  MsgRate[Mpps]
65536   190700       0.00             5966.21             0.095459
-----
[root@sw-mtx-010 ~]#
```

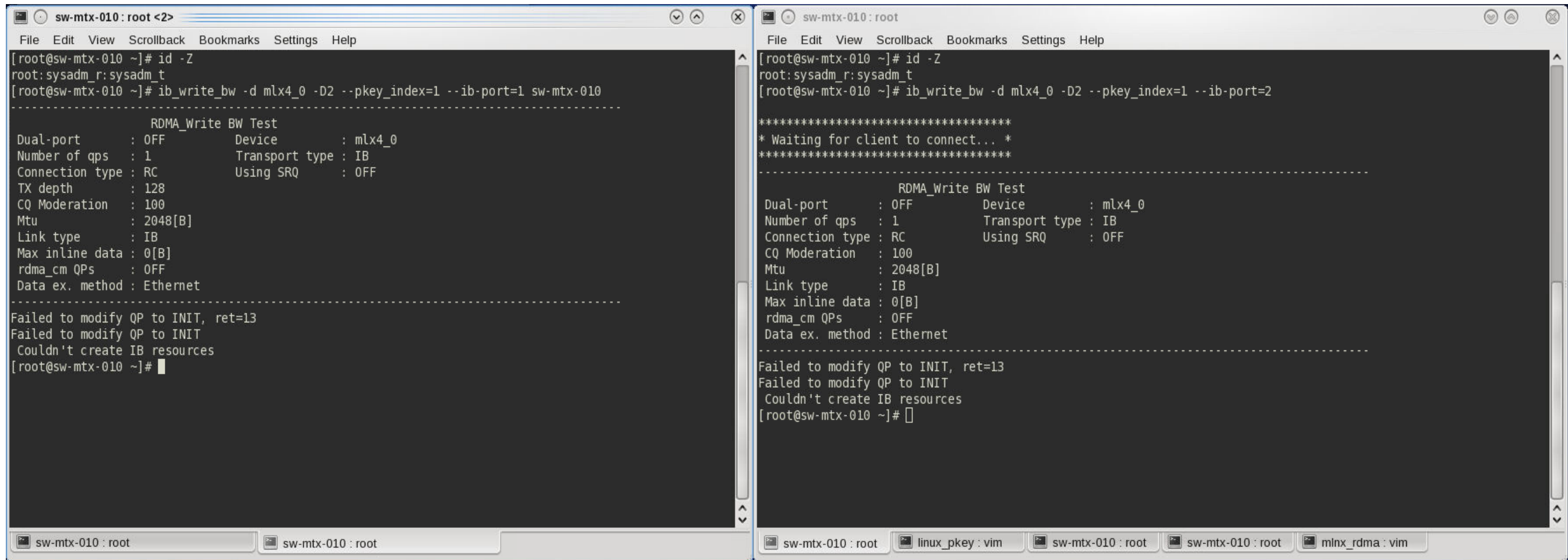
# Admin to Admin on the Admin Partition



```
sw-mtx-010: root <2>
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=2 --ib-port=1 sw-mtx-010
-----
RDMA_Write BW Test
Dual-port      : OFF      Device      : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
TX depth       : 128
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----
local address: LID 0x01 QPN 0x0213 PSN 0x17efaa RKey 0x40011a00 VAddr 0x007f8edecde000
remote address: LID 0x02 QPN 0x0212 PSN 0x7020a8 RKey 0x40011900 VAddr 0x007fab0ebaf000
-----
#bytes  #iterations  BW peak[MB/sec]  BW average[MB/sec]  MsgRate[Mpps]
65536   190800       0.00             5969.34              0.095509
-----
[root@sw-mtx-010 ~]#

sw-mtx-010: root
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=2 --ib-port=2
*****
* Waiting for client to connect... *
*****
-----
RDMA_Write BW Test
Dual-port      : OFF      Device      : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----
local address: LID 0x02 QPN 0x0212 PSN 0x7020a8 RKey 0x40011900 VAddr 0x007fab0ebaf000
remote address: LID 0x01 QPN 0x0213 PSN 0x17efaa RKey 0x40011a00 VAddr 0x007f8edecde000
-----
#bytes  #iterations  BW peak[MB/sec]  BW average[MB/sec]  MsgRate[Mpps]
65536   190800       0.00             5969.34              0.095509
-----
[root@sw-mtx-010 ~]#
```

# Admin to Admin on Staff Partition



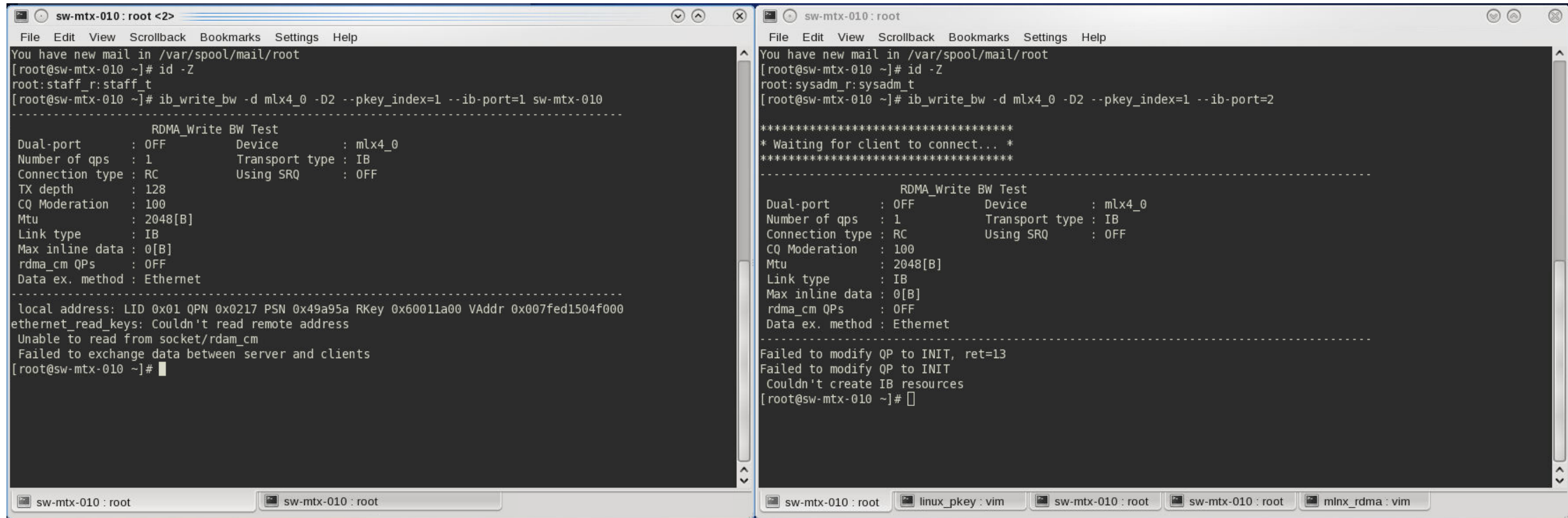
```
sw-mtx-010: root <2>
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=1 --ib-port=1 sw-mtx-010
-----
RDMA_Write BW Test
Dual-port      : OFF      Device       : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
TX depth       : 128
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----
Failed to modify QP to INIT, ret=13
Failed to modify QP to INIT
Couldn't create IB resources
[root@sw-mtx-010 ~]#

sw-mtx-010: root
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=1 --ib-port=2
-----
RDMA_Write BW Test
Dual-port      : OFF      Device       : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----
Failed to modify QP to INIT, ret=13
Failed to modify QP to INIT
Couldn't create IB resources
[root@sw-mtx-010 ~]#
```

- Both sides of the connection encountered an EACCESS error.



# Staff to Admin on Staff Partition

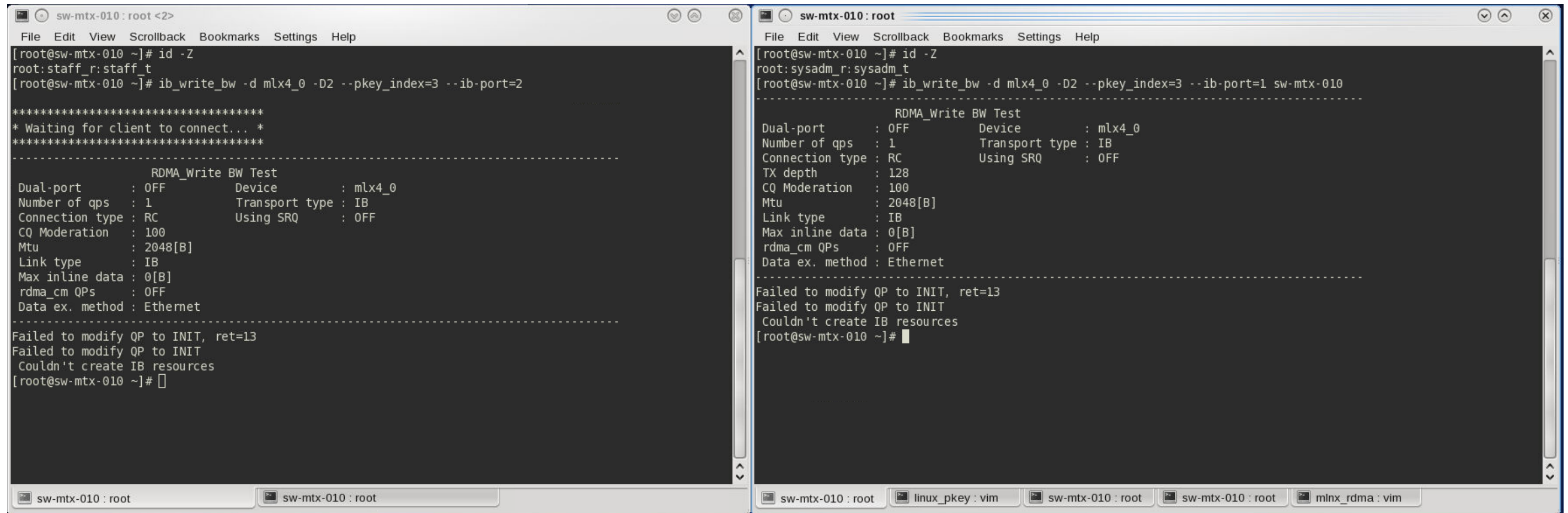


```
sw-mtx-010: root <2>
File Edit View Scrollback Bookmarks Settings Help
You have new mail in /var/spool/mail/root
[root@sw-mtx-010 ~]# id -Z
root:staff_r:staff_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=1 --ib-port=1 sw-mtx-010
-----
RDMA_Write BW Test
Dual-port      : OFF      Device       : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
TX depth       : 128
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----
local address: LID 0x01 QPN 0x0217 PSN 0x49a95a RKey 0x60011a00 VAddr 0x007fed1504f000
ethernet_read_keys: Couldn't read remote address
Unable to read from socket/rdam_cm
Failed to exchange data between server and clients
[root@sw-mtx-010 ~]#

sw-mtx-010: root
File Edit View Scrollback Bookmarks Settings Help
You have new mail in /var/spool/mail/root
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=1 --ib-port=2
-----
*****
* Waiting for client to connect... *
*****
RDMA_Write BW Test
Dual-port      : OFF      Device       : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----
Failed to modify QP to INIT, ret=13
Failed to modify QP to INIT
Couldn't create IB resources
[root@sw-mtx-010 ~]#
```

- Note in this case only the Admin side encounters an EACCESS error. The Staff side just has an error connecting.

# Admin to Staff on a Neither Allowed Partition



```
sw-mtx-010:root <2>
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:staff_r:staff_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=3 --ib-port=2

*****
* Waiting for client to connect... *
*****

-----
RDMA_Write BW Test
Dual-port      : OFF      Device      : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
TX depth       : 128
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----

Failed to modify QP to INIT, ret=13
Failed to modify QP to INIT
Couldn't create IB resources
[root@sw-mtx-010 ~]#

sw-mtx-010:root
File Edit View Scrollback Bookmarks Settings Help
[root@sw-mtx-010 ~]# id -Z
root:sysadm_r:sysadm_t
[root@sw-mtx-010 ~]# ib_write_bw -d mlx4_0 -D2 --pkey_index=3 --ib-port=1 sw-mtx-010

-----
RDMA_Write BW Test
Dual-port      : OFF      Device      : mlx4_0
Number of qps  : 1        Transport type : IB
Connection type : RC      Using SRQ    : OFF
TX depth       : 128
CQ Moderation  : 100
Mtu            : 2048[B]
Link type      : IB
Max inline data : 0[B]
rdma_cm QPs    : OFF
Data ex. method : Ethernet
-----

Failed to modify QP to INIT, ret=13
Failed to modify QP to INIT
Couldn't create IB resources
[root@sw-mtx-010 ~]#
```

```

sw-mtx-010 : root
File Edit View Scrollback Bookmarks Settings Help
You have new mail in /var/spool/mail/root
[root@sw-mtx-010 ~]# audit2allow -b -p /etc/selinux/refpolicy-pkey-2/policy/policy.30
DTJ - In policydb_read

#===== newrole_t =====
allow newrole_t kernel_t:unix_dgram_socket sendto;

#===== staff_t =====
allow staff_t init_t:unix_stream_socket connectto;
allow staff_t kmsg_device_t:chr_file { read open };
allow staff_t unlabeled_t:rdma_pkey modify;

#===== sysadm_t =====
allow sysadm_t staff_allowed_pkey_t:rdma_pkey modify;
allow sysadm_t unlabeled_t:filesystem mount;
allow sysadm_t unlabeled_t:rdma_pkey modify;
[root@sw-mtx-010 ~]#

```

- This tool generates policy code to allow violations in the audit log.
- If we added the three allow lines for “rdma\_pkey” the access errors in the demo would be allowed.

- Should we control access to partition values or <partition, port, device> tuples
- Hex format for P\_Key values in policy language.
- Subject of partition rules:
  - Pkeys
  - QPs/RDMA IDs/umad FDs
- “Action” word, currently modify.
- Event on QP modify to ERR.
- SMI interface control mechanism.





Thank You