



# Network Adapter Flow Steering

OFA 2012

Author: Tzahi Oved Date: March 2012

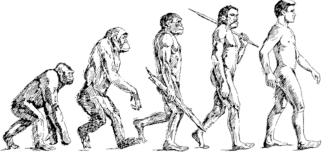
# **Receive Steering Evolution**



- The traditional Single Ring
  - All ingress traffic to land on a single receive ring
- Kernel threads / DPC de-multiplexing
  - Interrupt to de-mux ingress handling to multiple Kernel threads / DPCs
  - Still use single HW ring
  - May be aided with flow based HW hash
- Receive Side Scaling (RSS)
  - Multiple RX rings and interrupts, typically per core
  - Ingress flow hash based distribution to cores (QP->CQ->EQ->MSI->Core)

#### Flow Steering

- Network adapters to support TCP/UDP/IP transports flow steering
- Create as many receive rings as you need
- Adapter to steer ingress traffic to dedicated receive rings according to 5-tuple classification rule
- Rich specification of flow to ring mapping
- And... do it from user mode too!



Receive Steering – Why?

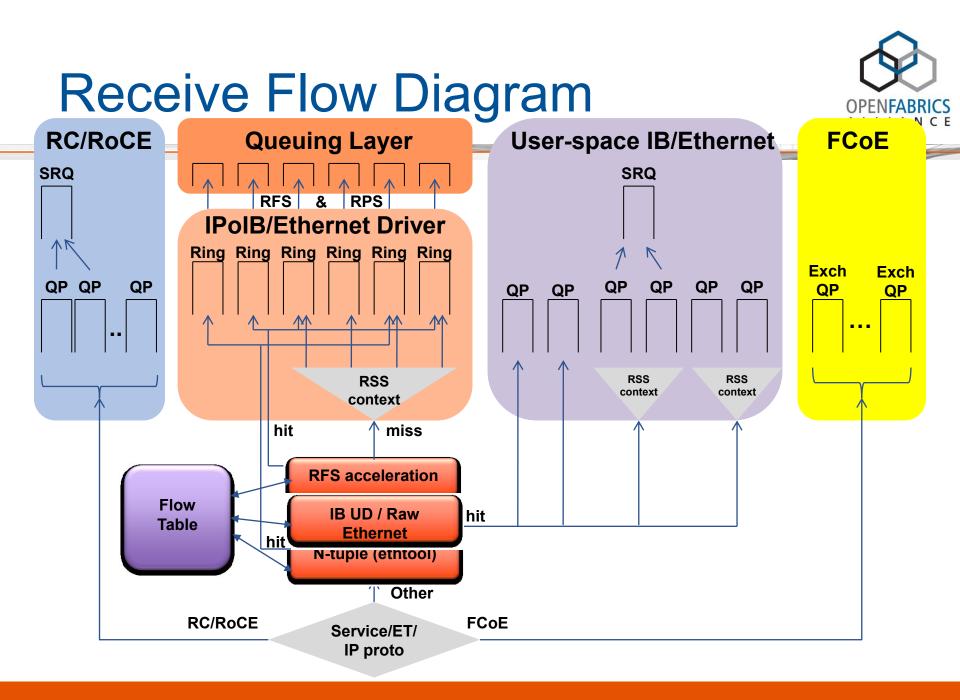


- Network stack acceleration
  - NUMA optimized
    - Receive buffer is allocated close to the flow handling core
  - Granular flow to core processing
    - Ingress flows are directed to the matching receiving application core
      - Don't interrupt other cores
    - OS/bypass, application ingress processing run on the same cache
  - Receive filtering at hardware level
- User mode application acceleration
  - TCP/UDP offload applications can steer its flows directly to user mode
    - Hijacking flows from kernel driver
  - Building block of user mode stack

# **Implementation Overview**



- Hierarchical Steering Capabilities
  - Transport: Ethernet, RoCE, and FCoE / IB-UD
  - Flow
    - MAC/VLAN/priority/ET or d.QP
    - IP protocol
    - TCP/UDP ports
  - RSS contexts, QPs, and SRQs
- Concurrent Support for Multiple steering Interfaces
  - NIC Driver
    - ethtool n-tuple settings
    - RPS, RFS acceleration
    - RSS
  - Verbs Consumers (Kernel and User)
    - RC, RoCE QPs
    - UD, Raw Ethernet QPs
  - FCoE
- Extends Naturally to Virtualized Environment (SR-IOV)
  - Steering rules table is maintained per VF



# Flows and Priorities



			Domain	Priority	Pattern	QP	ICE
	Steering rules ordered in a layered hierarchy <ul> <li>Domains</li> <li>ethtool</li> <li>Verbs (offloaded traffic)</li> <li>RFS, Netdev (network stack)</li> </ul> <li>Priority</li>		ethtool	1	tcp flow f	ring1	
		etht		2	src+dst ip g	ring2	
		rule		3	dst ip h	drop	
			Verbs & RDMA- CM	1	tcp flow a	qp1	
				1	tcp flow b	qp2	
	Multiple rules can exist in the same domain+priority		ows	1	udp mcast c	qp2	
	<ul> <li>For different flows</li> <li>Insertion order applies, last wins</li> </ul>	Lis	teners	2	tcp listener d	qp3	
	<ul> <li>For same flows</li> </ul>			2	tcp listener e	qp3	
	<ul> <li>Adapter to distribute ingress to multiple rings</li> </ul>		Linux RFS	999	dst_ip i	ring3	
٠	Rule processing order			999	tcp flow j	ring4	
	<ul> <li>Scan:</li> <li>(1) domain, (2) priority, (3) location</li> </ul>	RFS accel	eration	999	tcp flow k	ring4	
•	<ul> <li>First match takes</li> <li>Rule actions</li> </ul>		Net_dev	999	tcp flow l	ring4	
·	– Steer to QP att	ach_	mcast	999	tcp flow m	ring5	
	Multiple QPs are also allowed     Drop			999	udp flow n	ring6	

#### The API - Kernel Verbs



- flow\_id ib\_attach\_flow(target\_qp, domain, prio, flow\_spec\*)
  - Attach QP target\_qp to flow flow\_spec in
    - Domain *domain*
    - Priority *prio*
  - Attaching multiple QPs to the same flow+domain+prio allowed
    - Adds the QPs to the same rule
    - Replicates traffic to all QPs
- ib\_detach\_flow(flow\_id)
  - Detach QP target\_qp from flow flow\_spec in
    - Domain *domain*
    - Priority *prio*
  - Rule is actually removed once all QPs are removed
- Flow specification: IPv4, IPv6, IB-UD/RoCE-UD
  - Built from L2, L3, L4 protocol headers ordered filters and masks
  - L2: DMAC, VID, ET
  - L3: src IP, dst IP, protocol
  - L4: src port, dst port
  - IB/RoCE: dst QP, dst GID

# **Flow Filter Specification**



- Extendible flow specification structure
- Preceded by a flow specification header
- struct flow\_spec
  - uint filters\_num; // Number of filters describing the flow
  - void \*next\_flow\_spec; // pointer to next flow spec filter
- General flow specification structure:
  - struct flow\_spec\_xx
    - enum flow\_spec\_type
    - struct flow\_spec\_xx
    - void\* next\_flow\_spec; characteristics
- // eth/arp/ipv4/ipv6/udp/tcp/ud
- // flow filter characteristics + mask
- // header specific flow

# Flow Filter Specification – An Example

- struct flow\_spec\_eth
  - enum flow\_spec\_type type;
  - struct flow\_spec\_eth\_filter filter;
  - struct flow\_spec\_eth\_filter filter\_mask;
  - void \* next\_flow\_spec;
- struct flow\_spec\_eth\_filter
  - u8 d.addr[ETH\_ALEN];
  - u8 s.addr[ETH\_ALEN];
  - u16 protocol;
  - u16 vid;
  - u8 priority;
- The same for rest of header types:
  - struct flow\_spec\_arp, struct flow\_spec\_arp\_filter
  - struct flow\_spec\_ipv4, struct flow\_spec\_ipv4\_filter
  - struct flow\_spec\_ipv6, struct flow\_spec\_ipv6\_filter
  - struct flow\_spec\_udp, struct flow\_spec\_udp\_filter
  - struct flow\_spec\_tcp, struct\_flow\_spec\_tcp\_filter
  - struct flow\_spec\_ud, struct flow\_spec\_ud\_filter



**User Verbs Provider Interface** 



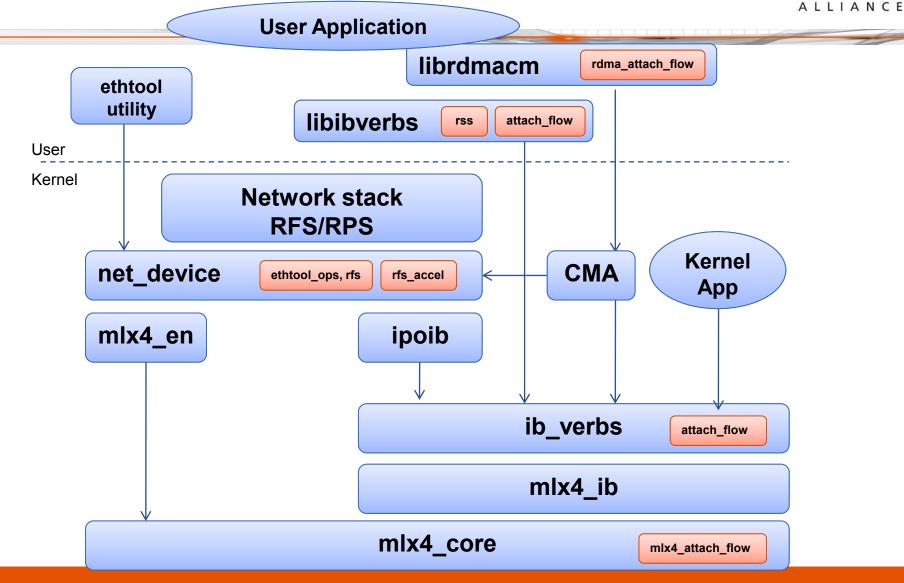
- QP type: IBV\_QPT\_RAW\_ETH, IBV\_QPT\_UD
- Similar to kernel (except for fixed domain = UVERBS\_DOMAIN)
  - flow\_handle ibv\_attach\_flow(target\_qp, prio, flow\_spec\*)
  - ibv\_detach\_flow(flow\_handle)
- Requires CAP\_NET\_RAW privileges
- Resources tracked by uverbs module
  - Flows are tracked as process resource
  - When terminates, all it's remaining allocated flows are detached

### rdmacm Abstraction



- RDMA-CM will use rdma\_cm\_id identifier for flow characteristic
  - rdma\_id is the "handle" to the target QP
- API calls:
  - flow\_handle rdma\_attach\_flow(flow\_rdma\_id, dest\_rdma\_id)
    - flow\_rdma\_id includes enough info to build the flow\_spec
    - For L2 params:
      - Extract local link layer params from local interface (IP address to MAC+VLAN / GID+QP)
      - Depending on rule type, extract remote link layer address through neigh\_lookup()
    - Will use 3 rule priorities according to flow\_rdma\_id
      - 1<sup>st</sup> (highest) For 5-tuple rule (connect)
      - 2<sup>nd</sup> For 3-tuple rule (listen)
      - 3<sup>rd</sup> For 2-tuple rule (fragment receive)
    - Use kernel verb ib\_attach\_flow() call with UVERBS\_DOMAIN
  - rdma\_detach\_flow(flow\_handle)

### Solution block diagram



### Quick Example



rdma\_create\_id(target\_id, PS\_UDP)

rdma\_create\_qp(target\_id)

// Create target id->qp

rdma\_create\_id(source\_id, PS\_UDP)

// Create source flow spec, set protocol

rdma\_bind\_addr(source\_id, 2-tuple of local IP address + port)

// Set rx flow spec characteristics through sockaddr

Or rdma\_resolve\_addr(source\_id, 2 or 4 tuple)

// src sockaddr to point to local interface+port, dst sockaddr to remote, protocol was set in create\_id()

rdma\_attach\_flow(source\_id, target\_id)

// Insert our steering rule

### Status and What's Next



- Kernel API done
- User API under definition, implementation to follow
- IPv4 support only for now
- Next
  - Bitwise flow spec mask support
  - IPv6
  - Extend flow spec filters for L7 application level data inspection
  - TCP/UDP payload based steering





