

End Station Reaction Points
Which Frames should a Rate Limiter slow?
May 2008 – Rev 2

Caitlin Bestler
Caitlin.bestler@neterion.com

Congestion Notification Message Scope

- **Already limited**
 - Generated based on sampling at CP.
 - Unicast delivery back to a single end station.
- **But the CNM supplies information**
 - It is not a “speeding ticket”
 - Ideally all flows from this end station that reach the congested CP should be throttled
 - But what is realistic?
 - What set of frames should be impacted?

Prior queuing should be Irrelevant

- **End stations have many designs**
 - Specific internal queue structures should neither be rewarded or penalized.
- **Frequently the pre-CNM queue will be too wide**
 - The end station will have had no reason to separate flows based on this destination.
 - Therefore many innocent flows will be slowed.
- **Sometimes the pre-CNM queue will be too narrow**
 - TOE/RDMA per-connection flows that are not the entire output from the end station to the destination.
- **Reaction Points may be created *after* the CNM is received, or it may only identify a *potential* queue.**

Use of Multiple SAs

- **Using Multiple Source Addresses can benefit network utilization when they actually use multiple paths.**
- **But when they hit the same CP, they at best just hog a greater slice of the bandwidth.**
 - The same traffic divided over more flows will be less “dinged” than a single flow would have been.
 - The only escape from this is to make the Source Address irrelevant to the scope of the Rate Limiter created *except* when there is specific reason to believe that Source Address truly will cause the CP to be avoided.
 - We should avoid creating an incentive to use *more* Source Addresses in each NIC.

Multiple Queues Can Be Tightly Coupled

- **Multiple source queues can be tightly coupled and have different Source Addresses**
 - Slowing one source will *instantly* cause other flows to increase their output.
 - Within many end stations the scheduler *pulls* “transmit descriptors” or “work requests” to fill the wire capacity.
 - Not the same as independent sources that “push” frames into a set of queues.
 - *Instantly* replacing the output capacity with frames that could be going to the same CP means that the CP will see *no* relief.

Deliberate Cheating Not Required

- **Many legitimate design trade-offs can result in use of more SAs.**
 - QCN should be neutral on these design trade-offs rather than encouraging or forbidding the use of more Source Addresses.
- **Example: Storage Client**
 - VM's use virtual drives. Parent partition is the sole client of the actual storage service.
 - Each VM acts as its own client.
- **Example: HPC**
 - Each rank uses a different VF in a multi-function NIC.
 - All ranks use a single VF.

Which Frames Should be slowed?

- **Ideal would be all frames that:**
 - Are from this end station
 - Will hit the same Congestion Point.
- **How close to this ideal be achieved with realistic real-time decision making?**
- **Initial assumptions:**
 - Different Priority, probably a different CP
 - Different VID+DA: probably a different CP
 - But maybe not for “next hop” CPs.
 - Different SA: probably the same CPs
 - Unless the SA selects a different egress port.
 - Or there is another reason to expect a different path.

L2 Flows that **SHOULD NOT** be impacted

- **Different Priority**
- **Different Destination End Station**
 - Which should be presumed if VID + DA is unique.
 - Not feasible to know remote VID to FID mapping.
 - Not feasible to know when multiple remote DAs are really the same end station.
 - Different non-aggregated egress port
 - If the first hop is a different non-aggregated port then it is reasonable to assume different CPs will be hit.
 - At least until reaching the final destination.

L2 Flows that SHOULD be impacted

- **Full match on:**

- Egress Port
- Priority
- Destination VID+DA

- **Rationale:**

- Other factors such as SA or L3/L4 headers are unlikely to have an impact on whether the same CP will be hit when they do not impact the egress port on the first hop.
- Merely creating more SAs will *appear* to improve congestion robustness *locally* by *stealing* bandwidth.
- Require actual knowledge of specific multi-pathing to justify NOT including the flows.

Possible special cases

- **When the CP is the last funnel before the destination then multi-pathing will not avoid it.**
 - Could be inferred by comparing CP's MAC Address with Destination.
 - Could be a boolean flag in the CNM.
- **When the CP is on the first hop**
 - End station could learn first hop on each port, and apply the Rate Limiter more broadly.
 - Alternate: CPs could be explicitly allowed to increase sampling rate on ports they know connect directly to end stations.

Special Cases Unlikely to Justify Special Effort

- **Same Egress Port, Same DA, Same Priority –
But interior CPs distribute traffic based on SA or L3/L4 headers.**
 - When this happens then *some* false head-of-line blocking will occur for frames that would really have missed the congested CP.
 - But far more often the SA/L3/L4 will not change the CP, but merely evade the Rate Limiter. Traffic will *instantly* divert to the flows that vary of SA/L3/L4 and the CP will see *no* relief.
- **Different Everything, but same internal CP**
 - Using a link-state databases (from Shortest Path Bridging or TRILL) this case *could* be identified.
 - But even if the data exists it is unlikely to be organized to allow a quick test of “would this frame go to this CP”?
 - Why penalize an end station for having a link-state database available?

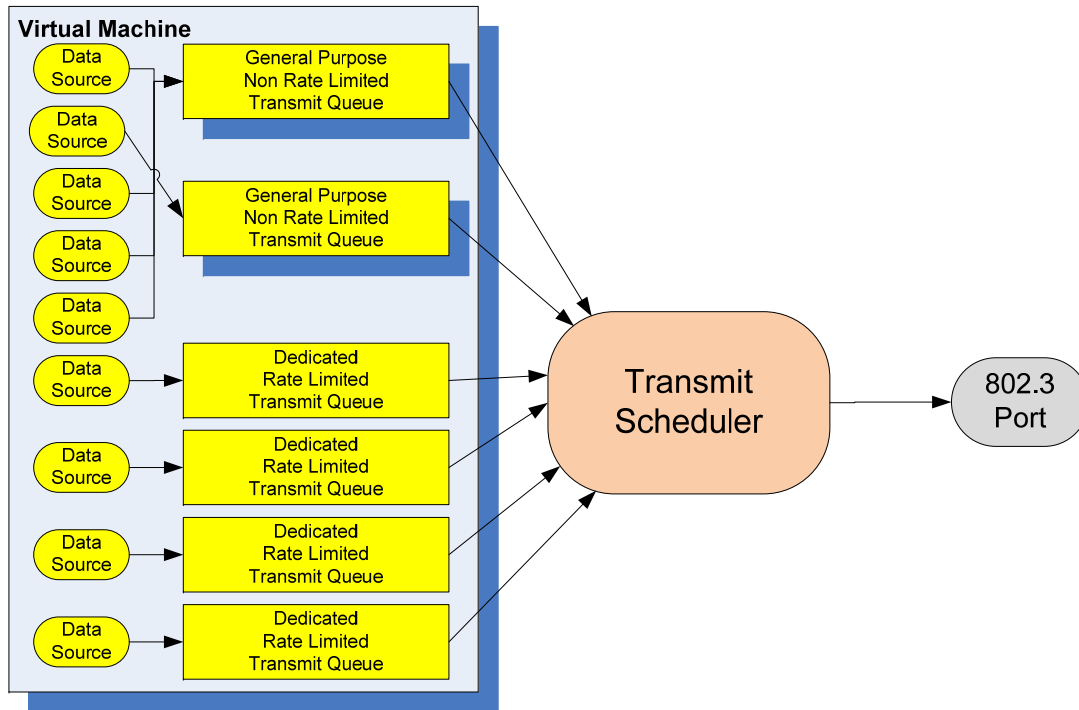
Split Reaction Points

- **End Station may have special purpose Output Queues that have a narrower scope than desired for a Rate Limiter.**
 - Primary example: Send Queues for TOE/RDMA.
- **For some designs the output from these queues would not naturally flow past general purpose Rate Limiters.**
- **Proposed solution: allow “split Rate Limiters” to be created on multiple internal queues in response to a single Congestion Notification Message**

End Station Congestion Points

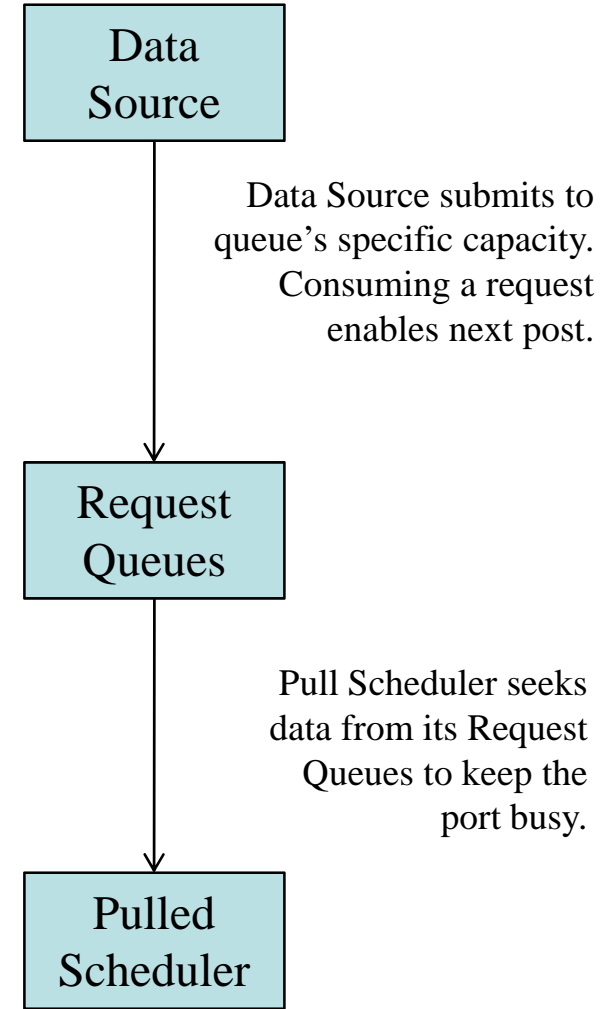
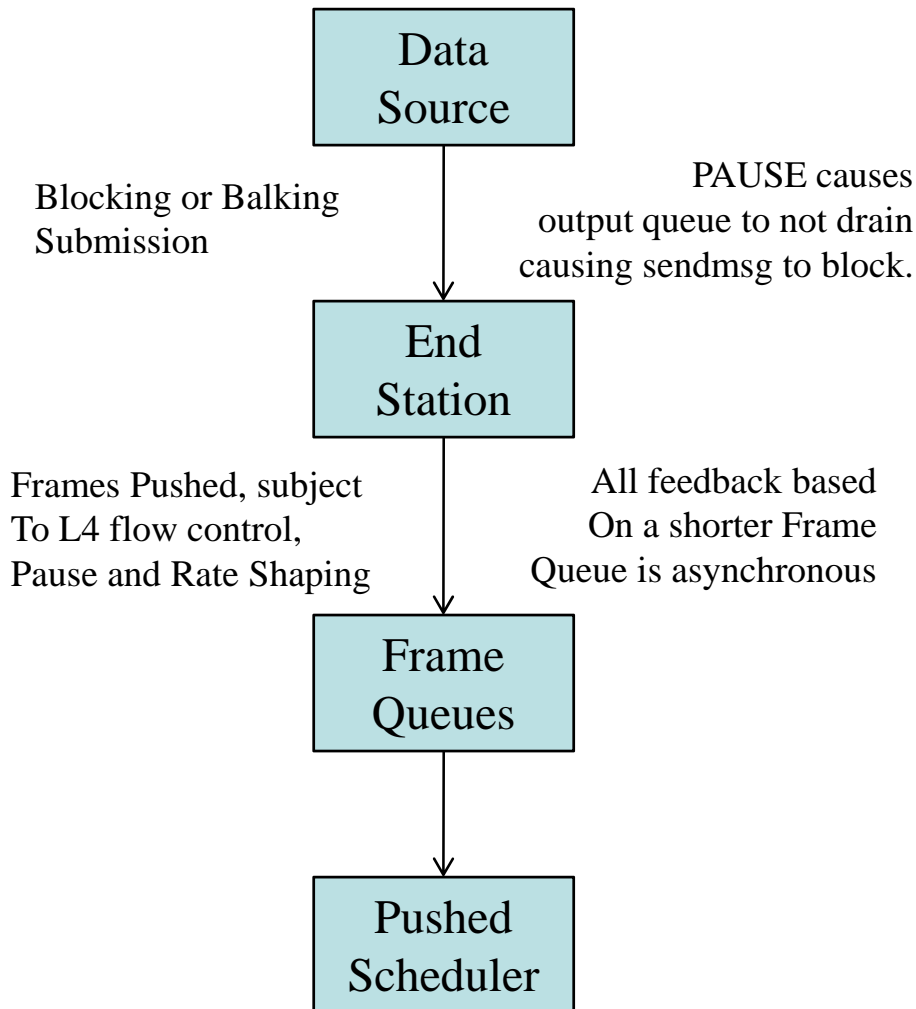
- **Not the topic of this presentation, but...**
- **End Station Congestion Points are *NOT* necessarily the inverse of its Reaction Points.**
- **For multi-function devices, the CPs are likely VF (Virtual Function) dependent.**
 - VID + DA determines VF, but multiple indexes could yield the same VF.
 - This is frequently a “default” VF for unknown addresses.
- **Having VF sensitive QCN triggers is desirable to limit inbound traffic based on VF.**

Hypothetical Multi-function NIC with all Qau, Qaz and Qbb support

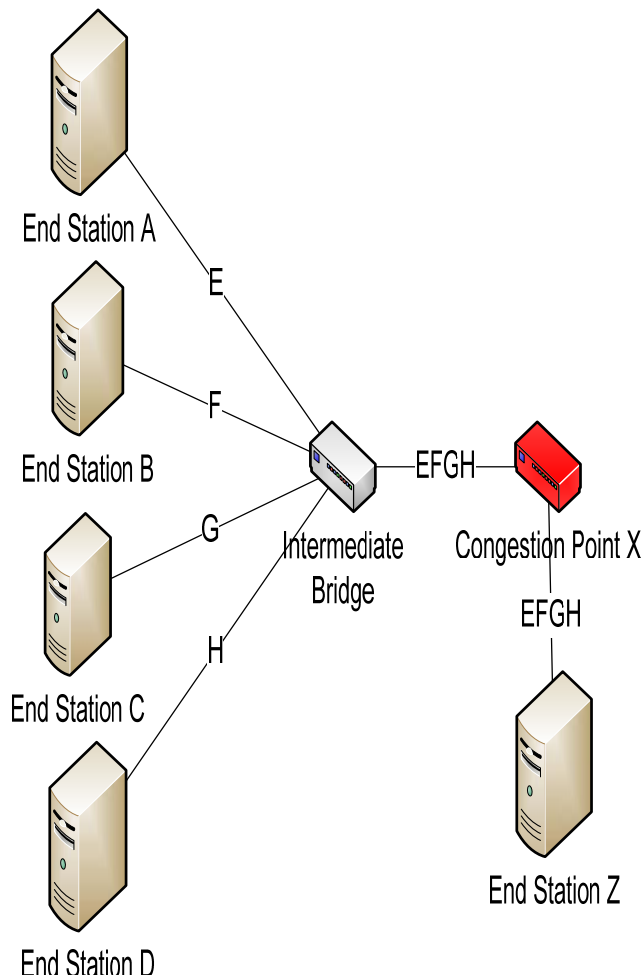


- Most Data Sources feed general purpose transmit queues that are not rate limited.
- Data Sources may be diverted to dynamically allocated rate limited transmit queues
- Data Sources may have dedicated Transmit Queues which are optionally Rate Limited (RDMA/TOE/iSCSI).
- Each Transmit Queue is for
 - Single Virtual NIC
 - Single Traffic Class
- Each PCB priority applies to set of transmit queues.
- Each Transmit Queue is accounted for by one ETS priority.
- Additional weighted round robin likely applies to each VNIC.

Push vs Pull

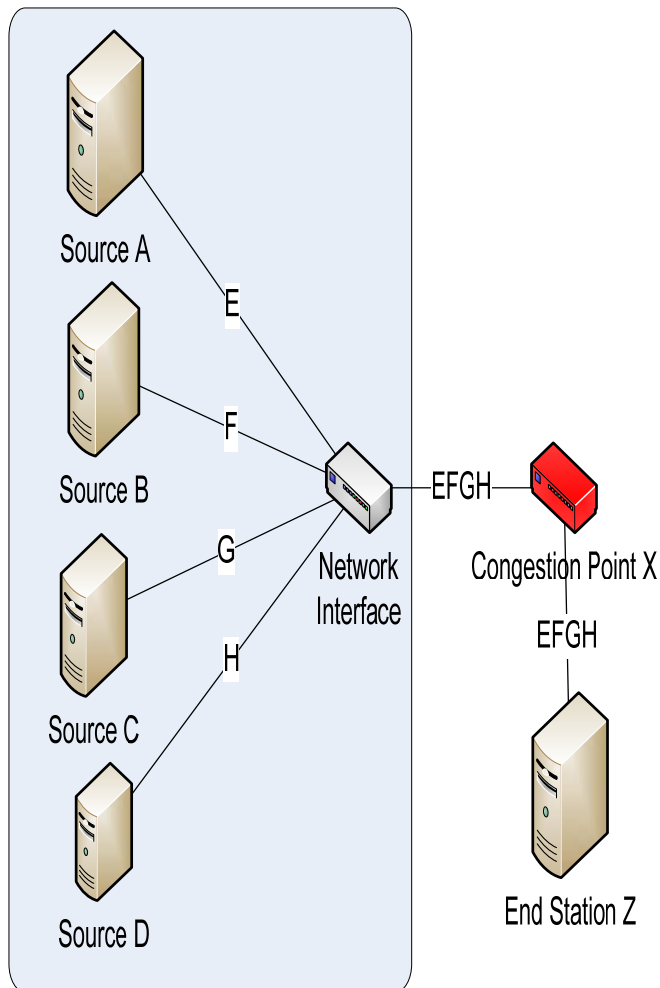


Data Sources on Separate End Stations



- **Flows E,F,G and H to End Station Z all reach CP X in the Red Bridge.**
- **CP X has sent CNM for Flow E to End Station A.**
- **A reduce E's rate.**
- **Queues in the Intermediate Bridge are drained more rapidly because E's rate is reduced.**
 - Immediate reduction in aggregate flow to CP X is unlikely, but there is an immediate drop in the ingress rate (because E is reduced).
 - Draining of queues on the Intermediate Bridge will result in fewer PAUSES to End Stations B, C and D.
 - Eventually this will cause F, G and H to speed up unless they get a CNM. Reducing the ingress rate reduction.
 - But it will not be immediate.

Data Sources on Single End Stations



- **Flows E,F,G and H to End Station Z all reach CP X in the Red Bridge.**
- **CP X has sent CNM for Flow E to Source Address used for flow E.**
- **Minimally scoped Rate Limiter:**
 - Only Source A reduces its rate.
 - Network Interface, seeking to feed a hungry port, increases the rate at which it transmits from B, C and D.
 - There is no immediate reduction in the aggregate flow to CP X.
 - There is no reduction in the ingress to the network of frames destined for CP X.
 - There will be no reduction until all sources on the End Station have received a CNM.
- **End Station scoped Rate Limiter**
 - E,F,G and H are all reduced in response to the first Rate Limiter.

End Station Stack Must Participate

- **When a flow is rate limited the source must ultimately be slowed to match.**
- **With connection-specific RDMA style interfaces this is just a matter of not completing Send Work Requests.**
- **But existing IP stacks generally use a limited number of queues into a given L2 device.**
- **Possible results:**
 - Head of line blocking: a pause on one L2 flow will impact all traffic for the same Priority, whether to the same destination or not.
 - Buffer Drain: to avoid head-of-line blocking the driver will attempt to put rate limited frames in a side-queue.
 - Even if stack supports out-of-order completion, it will result in memory pressure.
 - Worst case: memory pressure causes swap out – to network storage that is reached via the problem Congestion Point.

Method of Participation may vary

- **QCN feedback to L4**
 - Any L4 socket that is impacted by a Rate Limiter is told of the rate limit in L3/L4 terms. It adjusts it's L4 congestion window accordingly.
- **Directed Queuing**
 - L2 driver informs its client that a specific flow should be placed in a distinct input queue.
- **Directed Pausing**
 - L2 driver informs its client that a specific submission cannot not be accepted at this time. The same frame should not be retried until a specified time (or callback). The source socket should block, but not any others.