

# Preliminary Analysis Results of QCN

IEEE 802.1Qau  
San Francisco Plenary July 2007

M. Gusat, C. Minkenberg, R. Birke, L. Chen and R. Luijten

IBM Research GmbH, Zurich  
July 13<sup>th</sup> 2007

# Outline: Answering the .1au request to analyse 2pt-QCN

1. Validation and simulation results\*
  1. Single Hop: Input Generated and Output Generated
    1. IG: baseline
    2. OG: small system, few flows, medium severity
  2. Fat Tree: Input Generated and Output Generated
    1. Topologies from 3 to 7 levels: 32 to 256 nodes
    2. Traffic: IG and OG of small to medium severity
2. First analytical observations
  1. linearization
  2. statistical analysis with M/M/1
3. On design complexity and protocol issues
  1. RLT aggregation, multipath, CPID
4. Summary and proposal

\* Caveat: Preliminary work on a scheme under development.

# ECM and QCN at a Glance

## ECM

- Probabilistic sampling
  - $P = \text{constant (1\%)}$
  - easy to implement
- Feedback (16bits) range is limited
$$F_b \in [-Q_{eq}(1+2w), Q_{eq}(1+2w)]$$
- Multiplicative rate decrease based on negative feedback
  - No limit
- Rate increase driven by positive feedback

## 2-point QCN

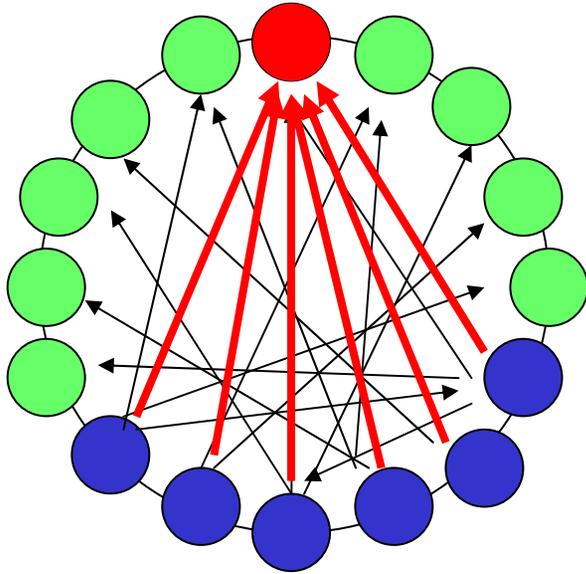
- Probabilistic and adaptive sampling
  - $P = f(\text{feedback}) P$  [1%, 10%]
    - ✓  $P_{\text{reflection}}$  calculated online/queue
  - Dynamic range for  $q, q'$ : no cap
  - in situ Fb reduction:  $\{q(t), q'(t)\} \rightarrow Fb$ 
    - ✓ quantized to 6 MSB
- Multiplicative rate decrease based on negative feedback
- Self-clocked recovery
  - Fb-independent increase based on
    - ✓ Binary increase followed by additive
      - Self-clocked based on tx packets
    - ✓ Multiplicative increase
      - Time based
  - Per phase
    - ✓ extra fast recovery (EFR)
    - ✓ fast rec. (FR)
    - ✓ (hyper)/active increase (AI)
    - ✓ drift (true timer)
  - Not simulated: discount FR and jitter.

# Simulation Results - One Hop

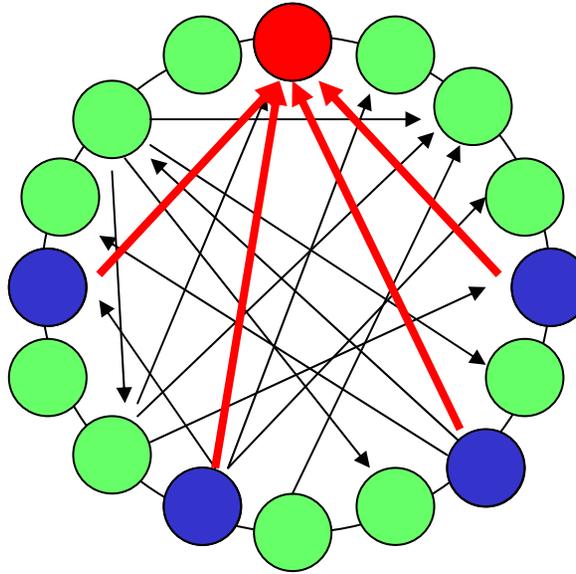
1. Validation and simulation results
  1. Single Hop: Input Generated and Output Generated
    1. IG: baseline
    2. OG: small system, few flows, medium severity

# Single Hop Comparison Scenarios

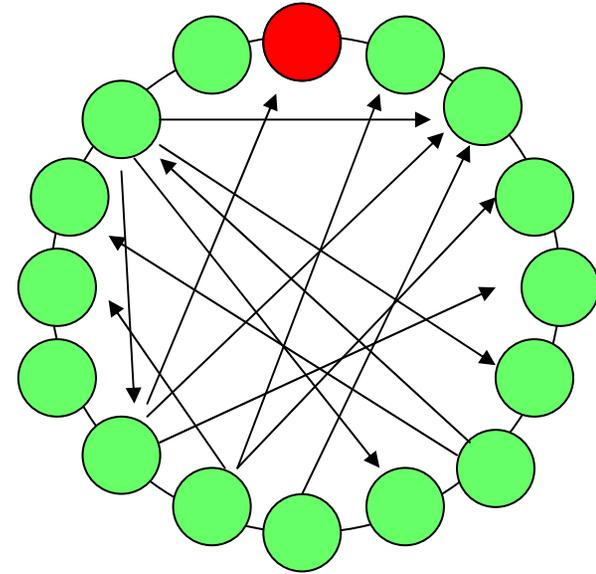
Scenario 1 IG: HSV=3



Scenario 2 IG : HSV ~ 2.5



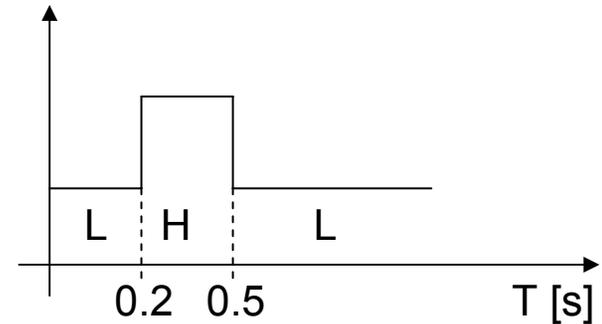
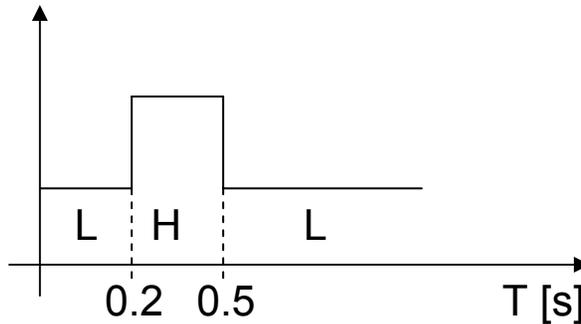
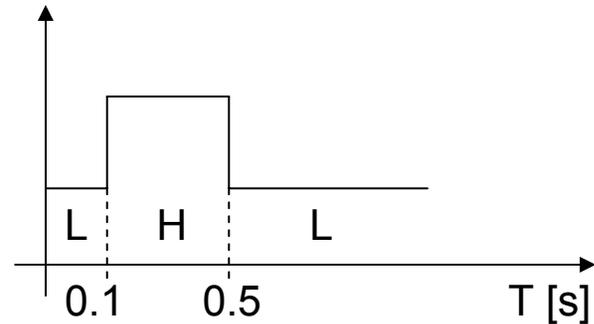
Scenario 3 OG : HSV=5



$L_s = 0.5$   
 L: 6 src send to 15 dst  
 H: 6 src converge on 1 dst

$L_s = 0.5$   
 L: 16 src send to 16 dst  
 H: 4 src -> 1 dst (12 continue)

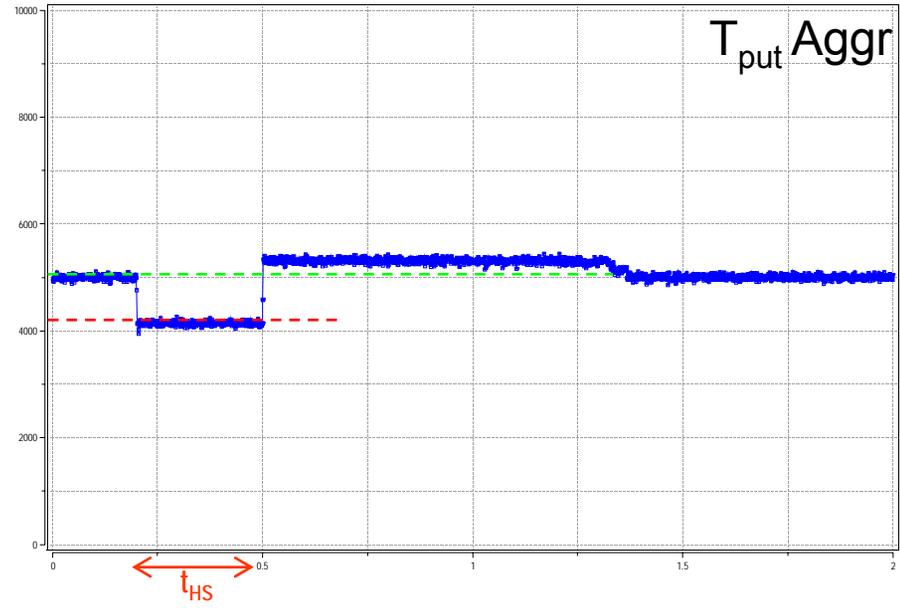
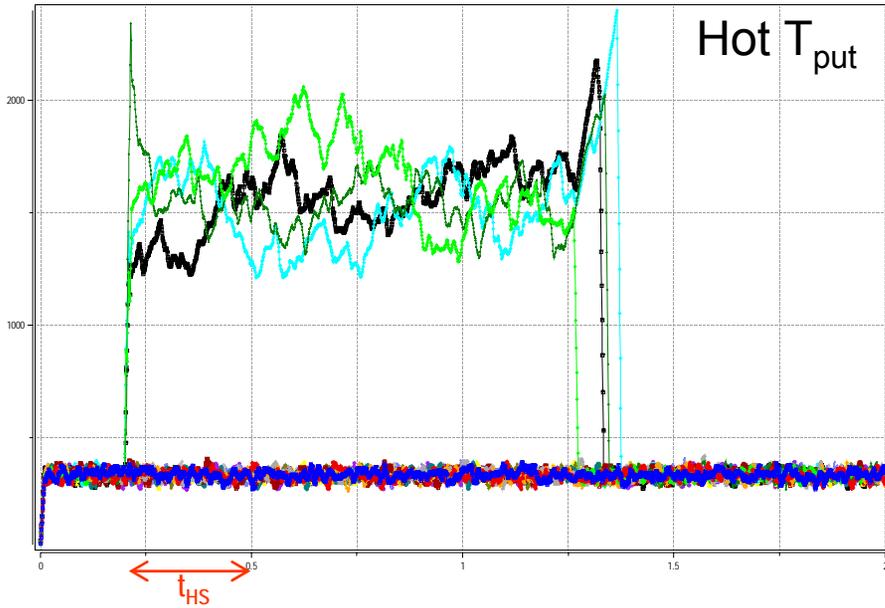
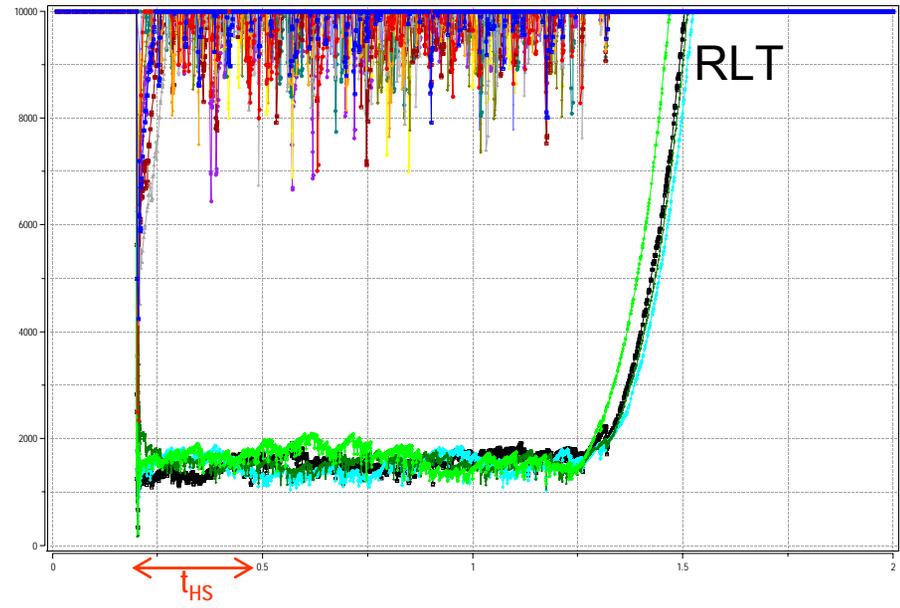
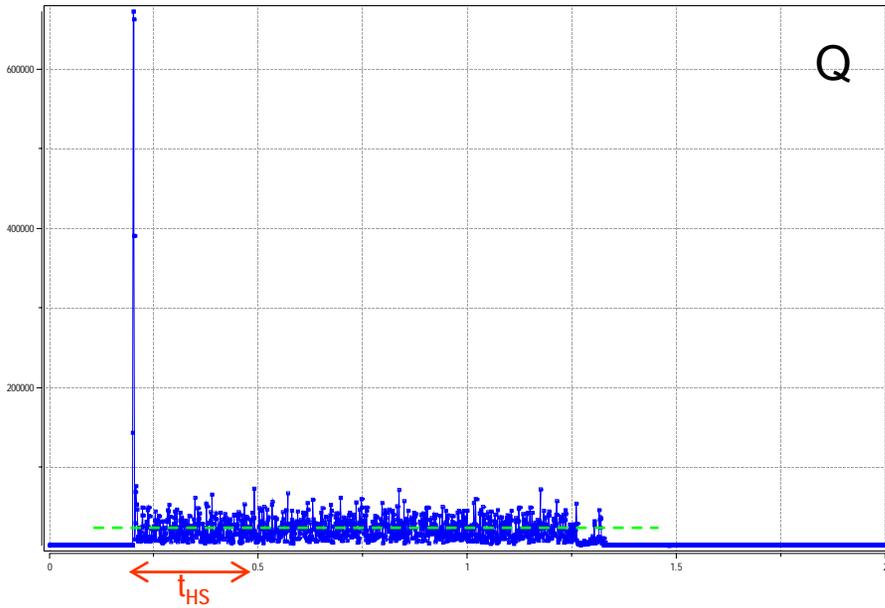
$L_s = 0.5$   
 L: 16 src send each to other 16 dst  
 H: 1 dst drops service rate to 0.1



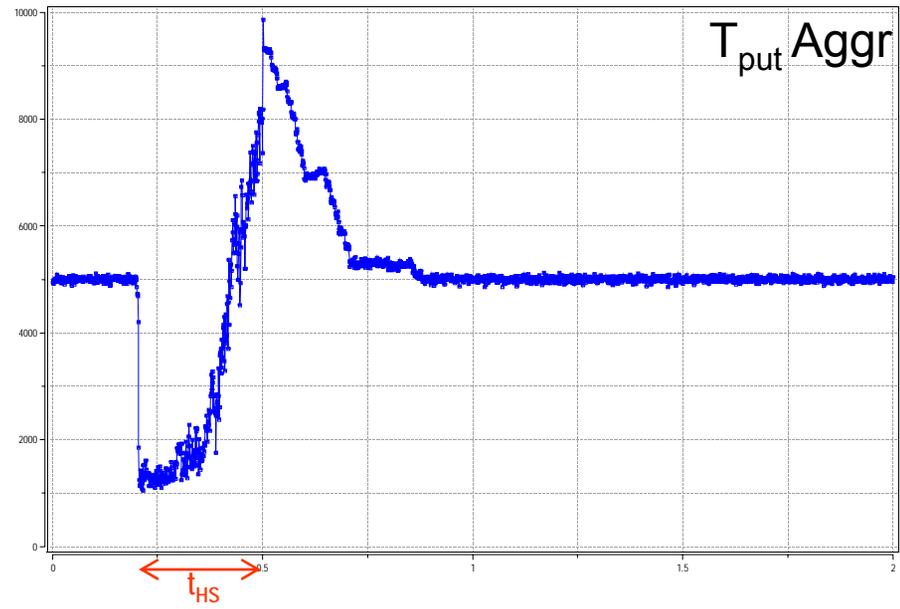
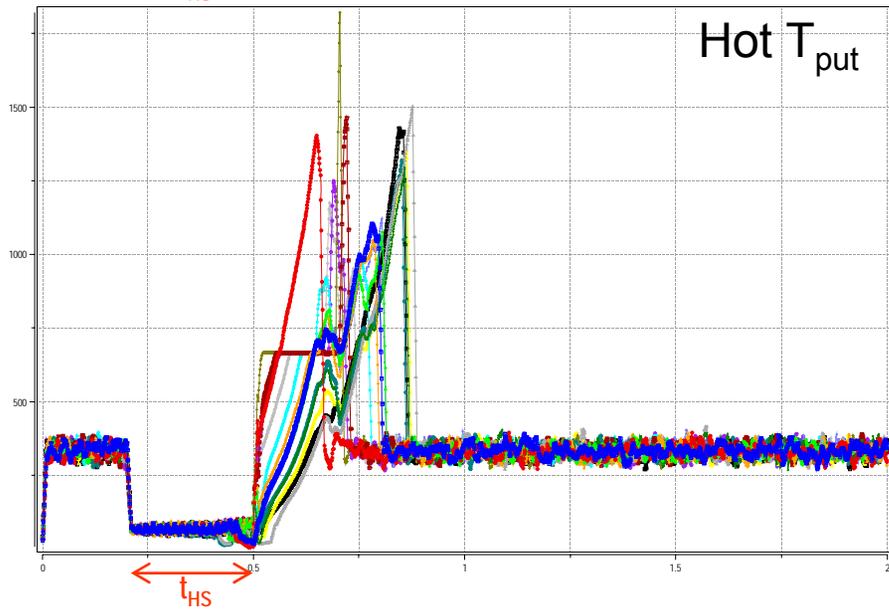
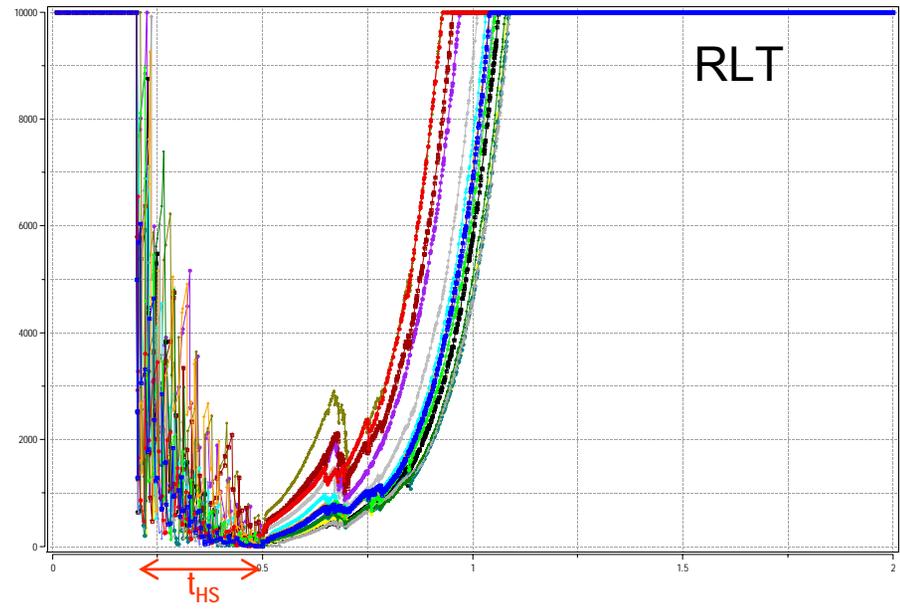
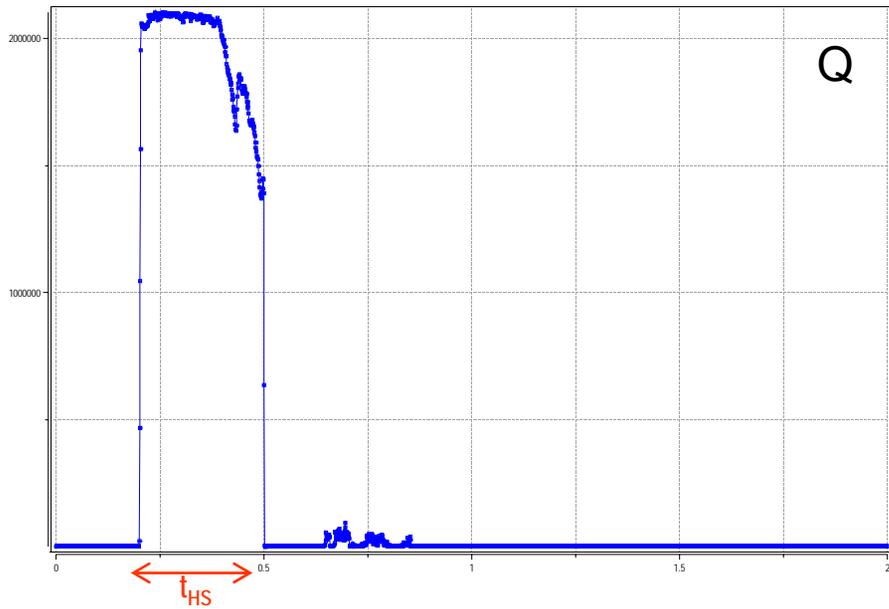
# Simulation Parameters

- Traffic
  - I.i.d. Bernoulli arrivals
  - Uniform destination distribution (to all nodes except self)
  - Fixed frame size = 1500 B
- Switch
  - VOQ with 2.4MB shared mem
  - Partitioned memory per input, shared among all outputs
  - No limit on per-output memory usage
  - PAUSE enabled
    - ✓ Applied on a per input basis based on local high/low watermarks
    - ✓  $\text{watermark}_{\text{high}} = 141.5 \text{ KB}$
    - ✓  $\text{watermark}_{\text{low}} = 131.5 \text{ KB}$
- Adapter
  - RLT: VOQ and single; RR service
  - One rate limiter per destination
  - Egress buffer size = 1500 KB,
  - Ingress buffer size = Unlimited
  - PAUSE enabled
    - ✓  $\text{watermark}_{\text{high}} = 150 - \text{rtt} * \text{bw} \text{ KB}$
  - $\text{watermark}_{\text{low}} = \text{watermark}_{\text{high}} - 10 \text{ KB}$
- ECM
  - $W = 2.0$
  - $Q_{\text{eq}} = 37.5 \text{ KB}$
  - $G_d = 0.5 / ((2*W+1)*Q_{\text{eq}})$
  - $G_{i0} = (R_{\text{link}} / R_{\text{unit}}) * ((2*W+1)*Q_{\text{eq}})$
  - $G_i = 0.1 * G_{i0}$
  - $P_{\text{sample}} = 2\%$  (on average 1 sample every 75 KB)
  - $R_{\text{unit}} = R_{\text{min}} = 1 \text{ Mb/s}$
  - BCN\_MAX enabled, thshld = 150 KB
  - BCN(0,0) dis/enabled, thshld = 300KB
- QCN
  - Drift Factor = 1.005
  - Timer Period Drift = 0.0005 s
  - Extra Fast Recovery enabled
  - EFR MAX disabled.
  - $A = 3 \text{ Mbps}$
  - Fast Recovery Threshold = 5
  - Hyper Active Increase disabled, unless otherwise specified

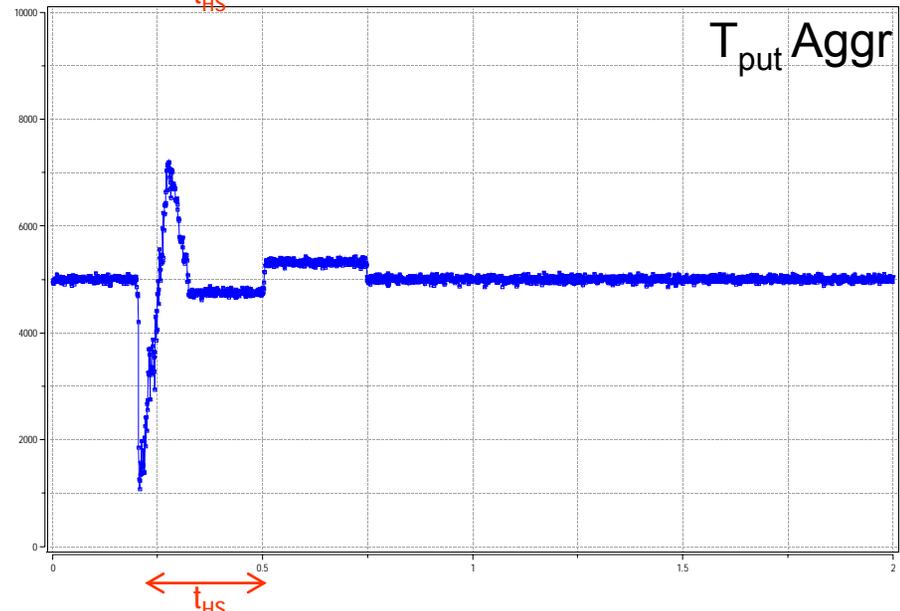
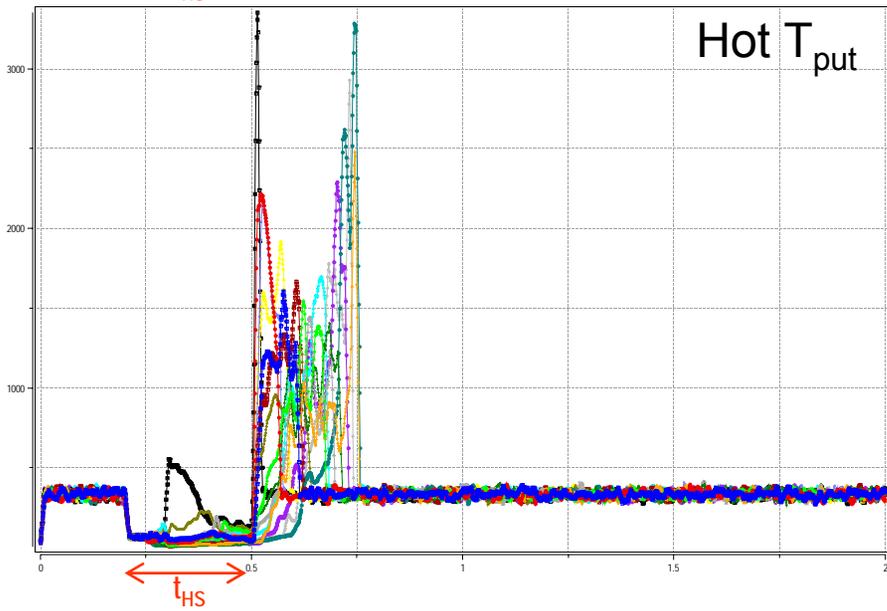
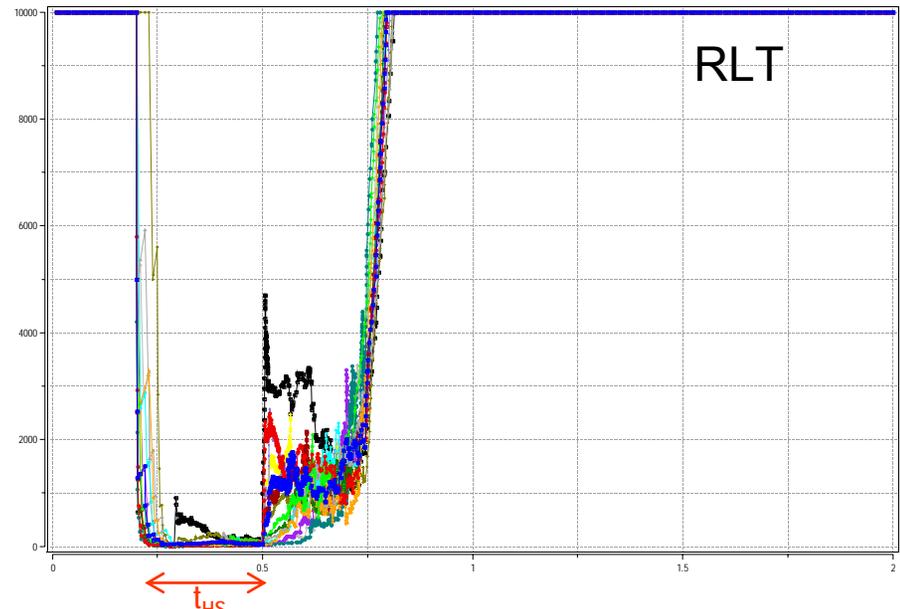
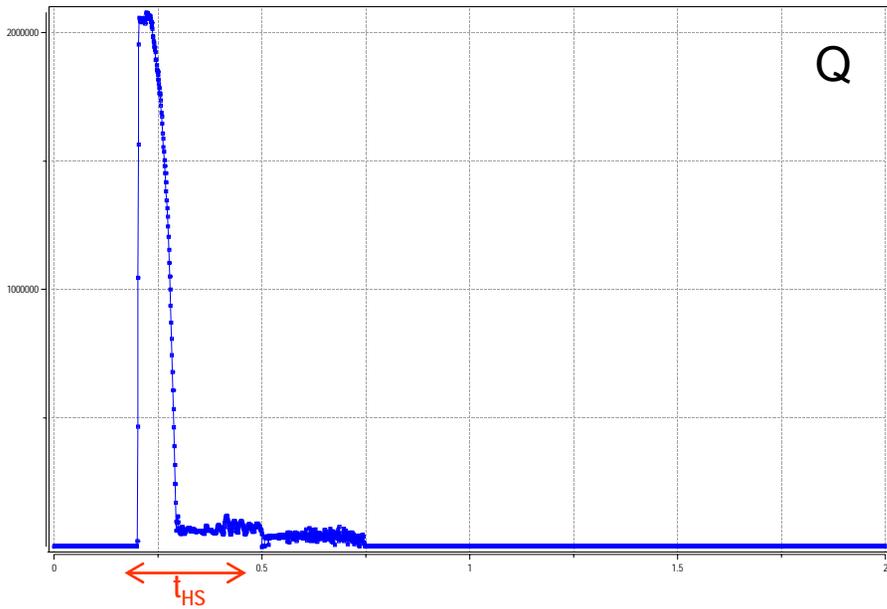
# QCN - IG - Scenario1



# QCN - OG - Scenario3



# ECM - OG - Scenario3

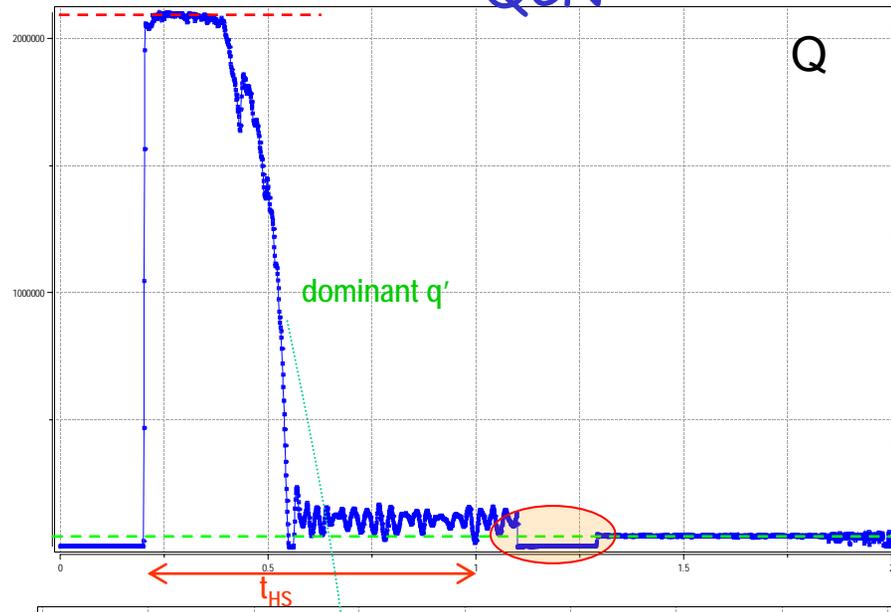


# ECM - QCN - OG - Scenario3 - long HS: Q vs. rate

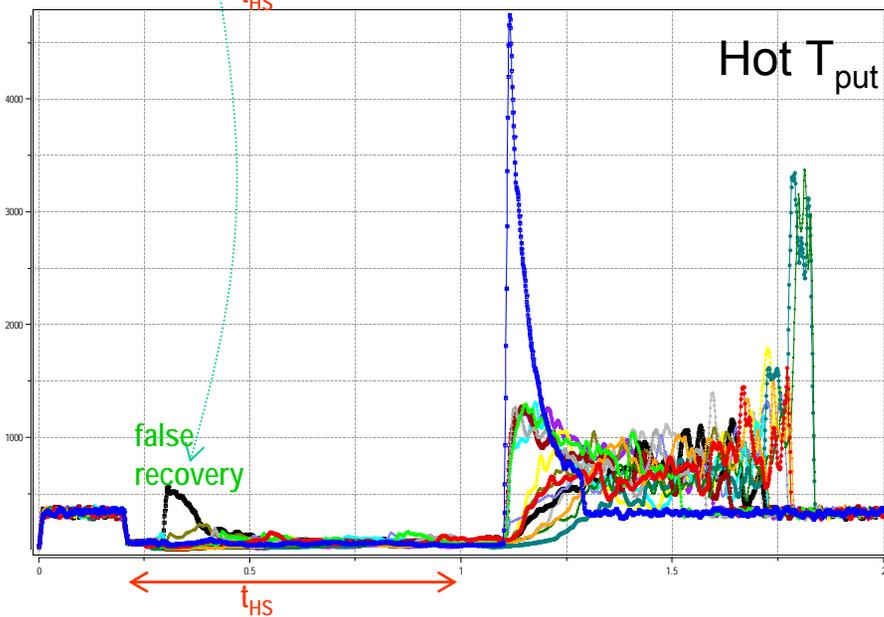
ECM



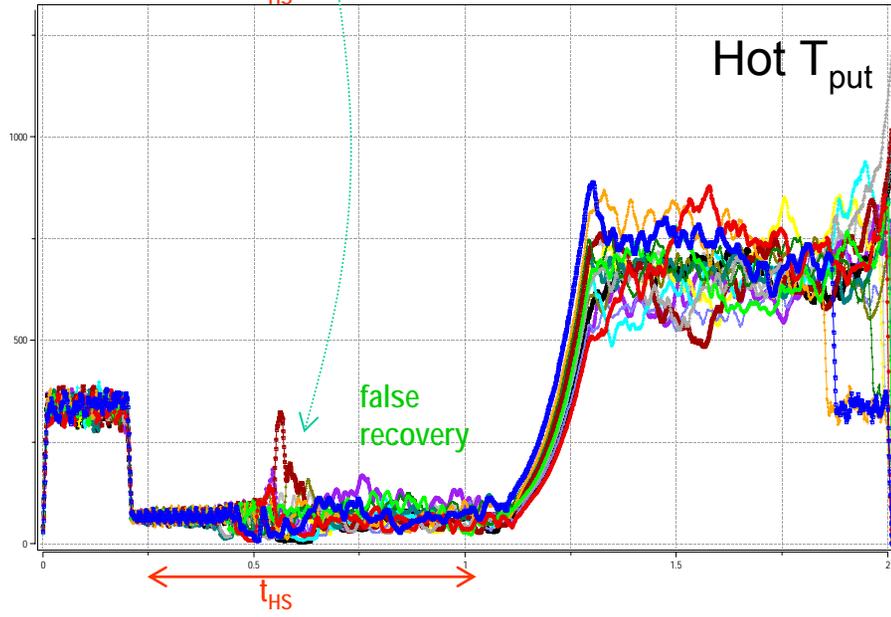
QCN



Hot  $T_{put}$



Hot  $T_{put}$



# Simulation Results in Fat-trees

## 1. Fat Tree: Input Generated and Output Generated

1. Topologies from 3 to 7 levels: 16 to 256 nodes
2. Traffic: IG and OG of small to medium severity
3.  $t_{HS}$ 
  1. short: 100-500ms
  2. long: 200-1100ms

## 2. $N^2$ flows: e.g. for 256 nodes $\Rightarrow$ $\sim$ 64K flows

1. Distributions: uniform traffic without self-traffic. Bernoulli departure times and uniform across destinations. Only the flows going to the HS are recorded (256 nodes  $\rightarrow$  255 flows) and the global  $T_{put}$ .

## 3. IG and OG refer to scenario 2 and 3

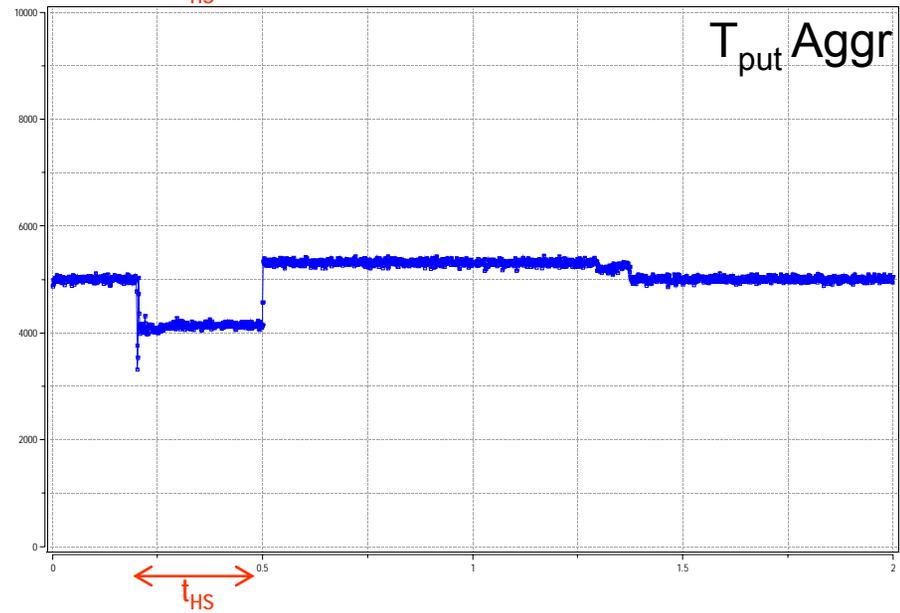
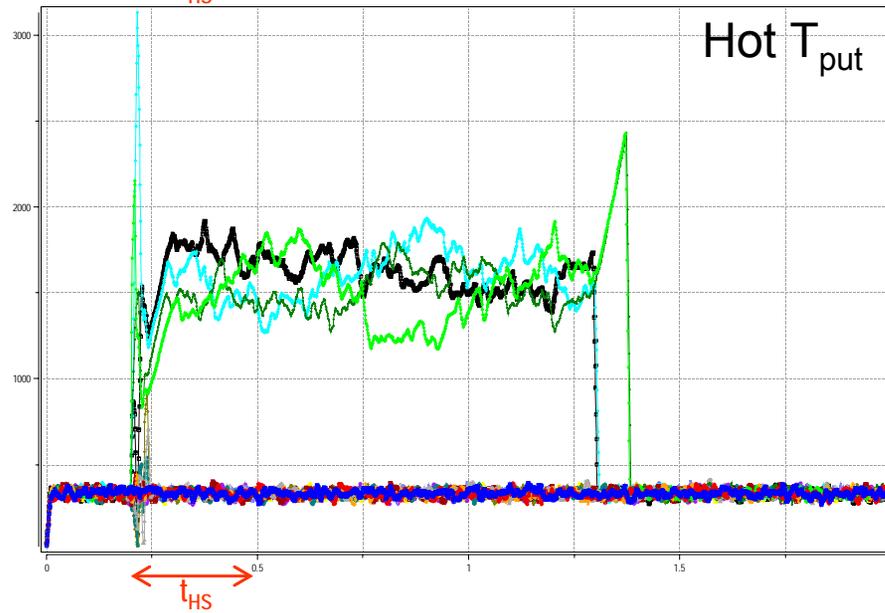
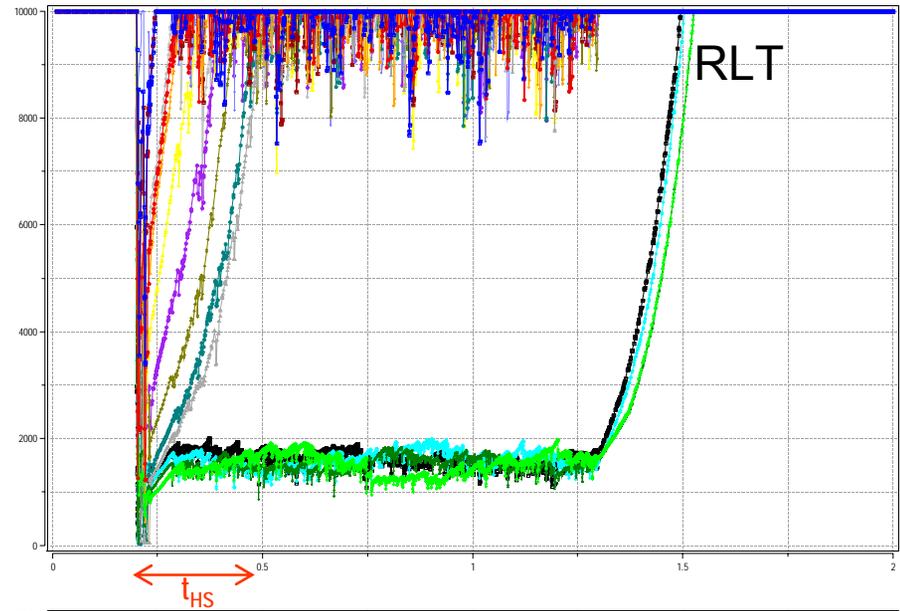
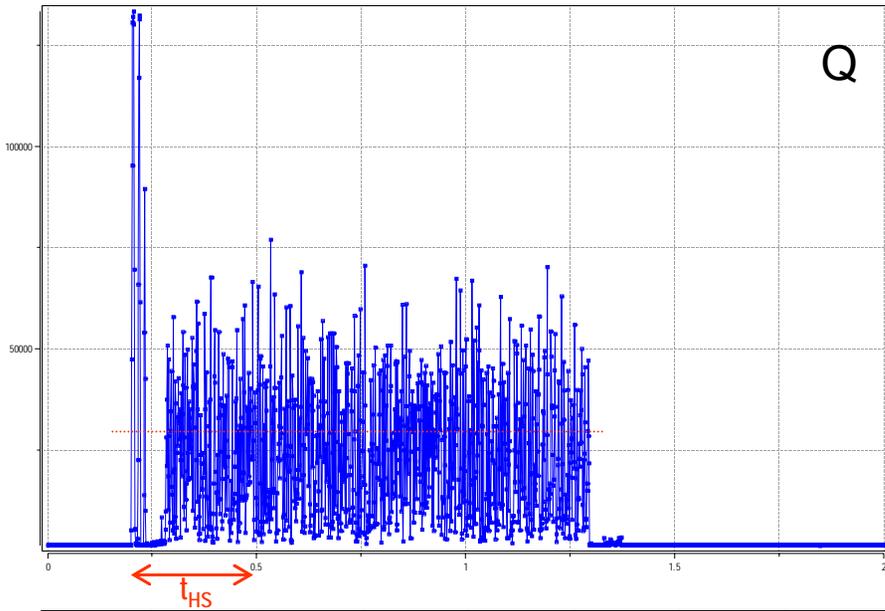
### 1. IG

1. 0.5 background traffic from every host.
2. During HS, 4 hosts redirect their traffic ( $4 * 0.5 = 2$ ) to the HS
3. HSV=2+ background from N-1-4 (self and HS creating hosts)

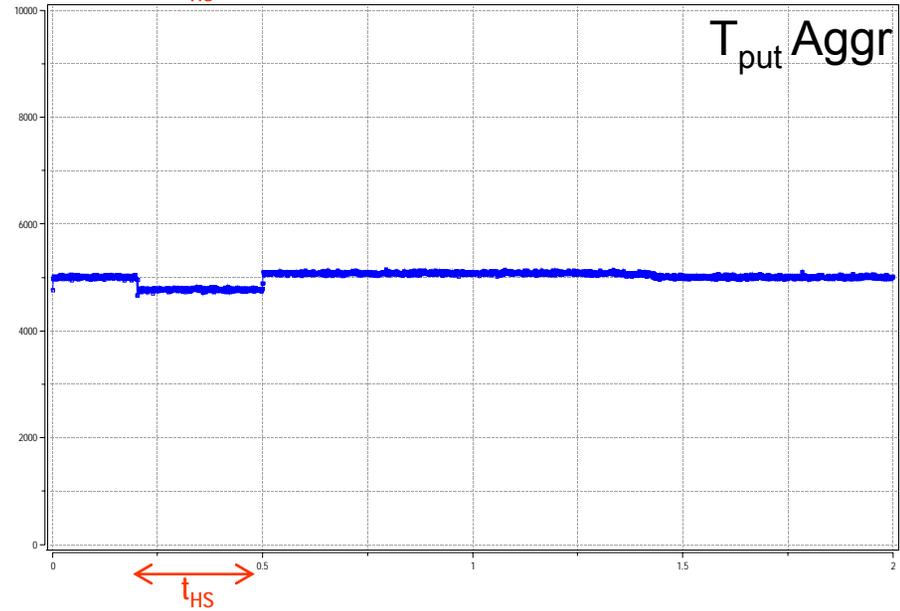
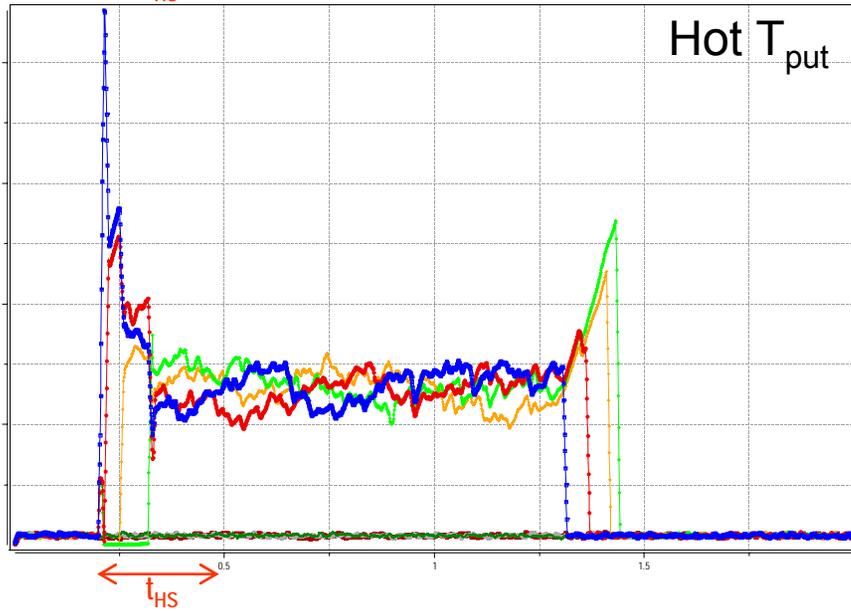
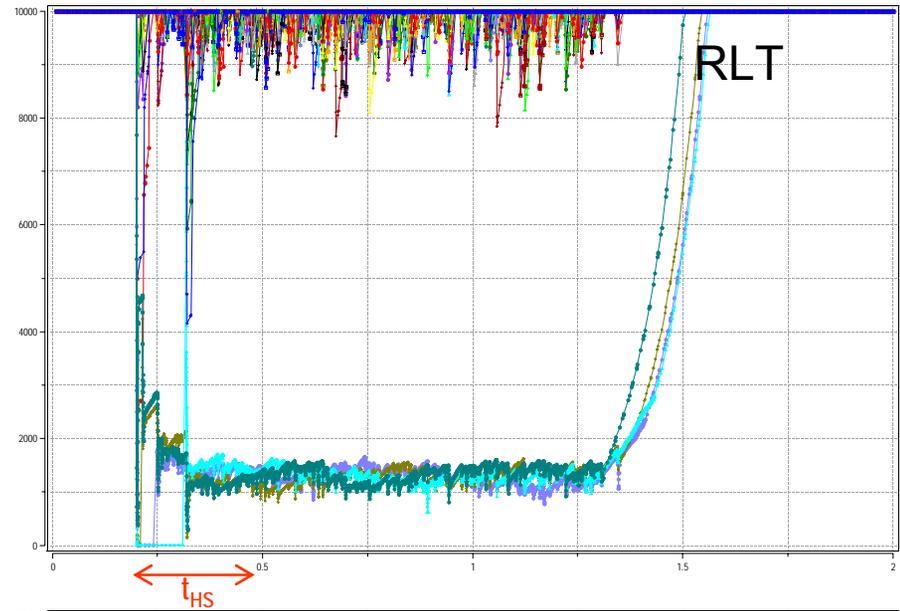
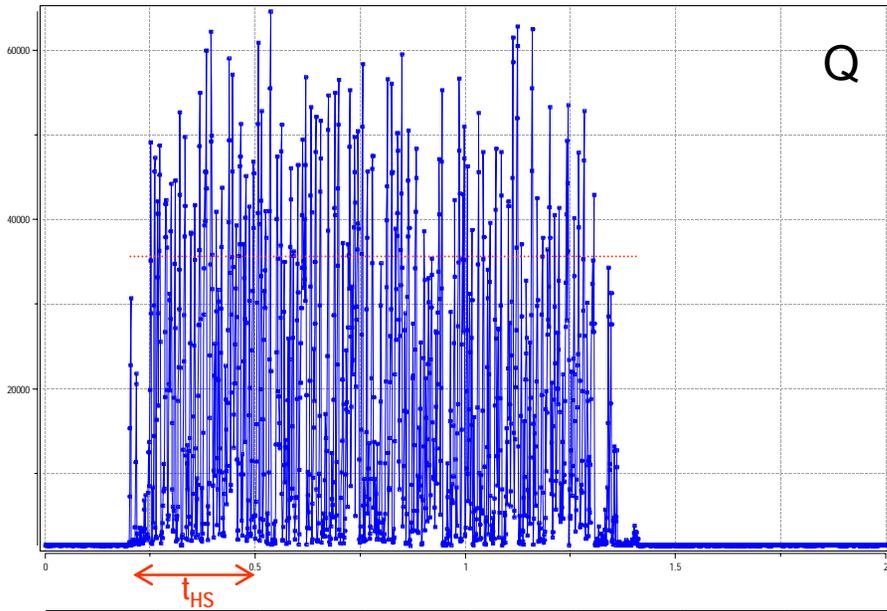
### 2. OG

1. 0.5 background traffic.
2. HS host reduces service rate to 0.1
3. HSV  $0.5 / 0.1 = 5$

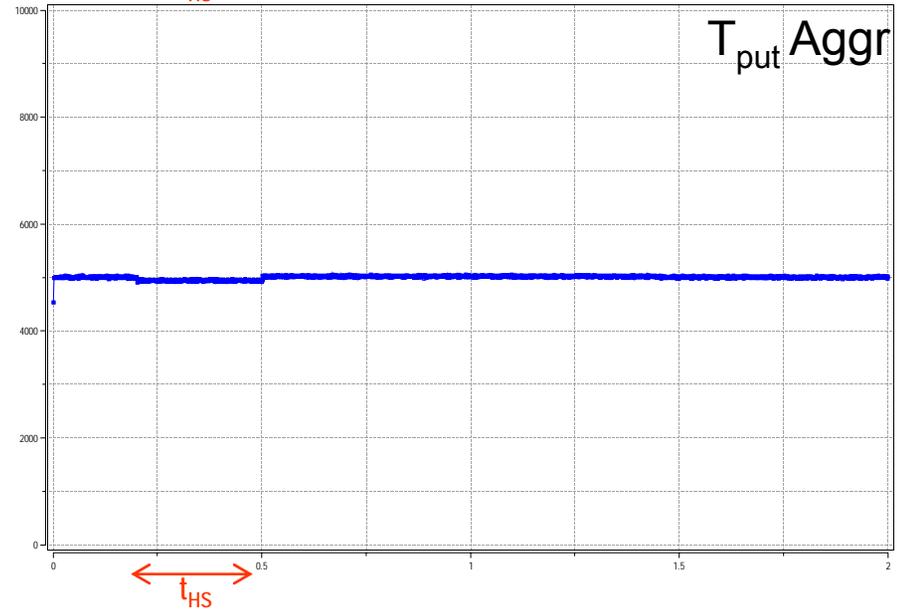
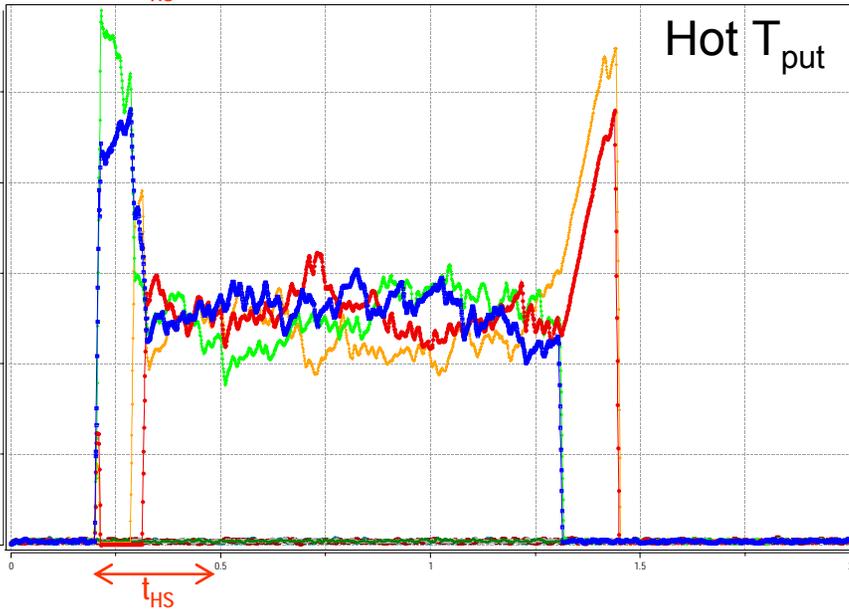
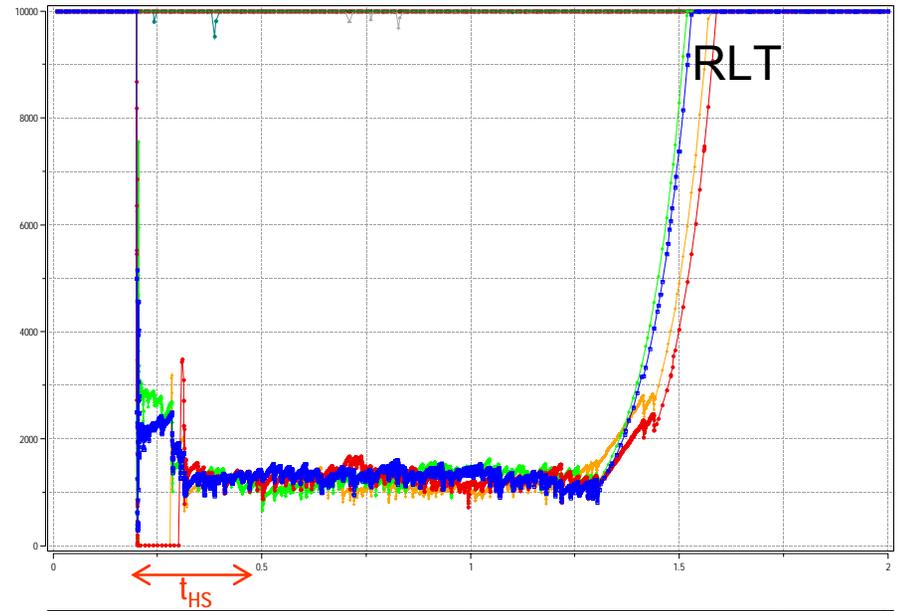
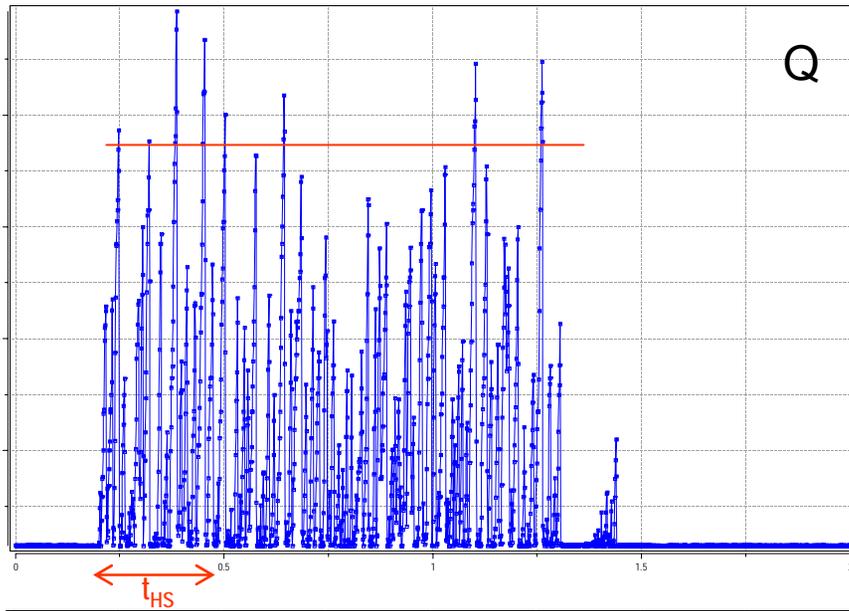
# QCN - IG - FTree3L



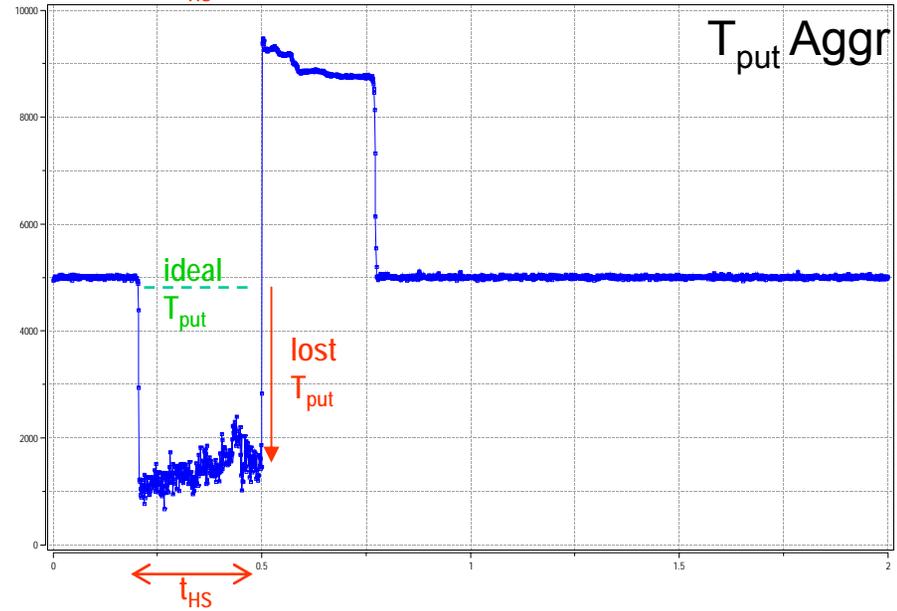
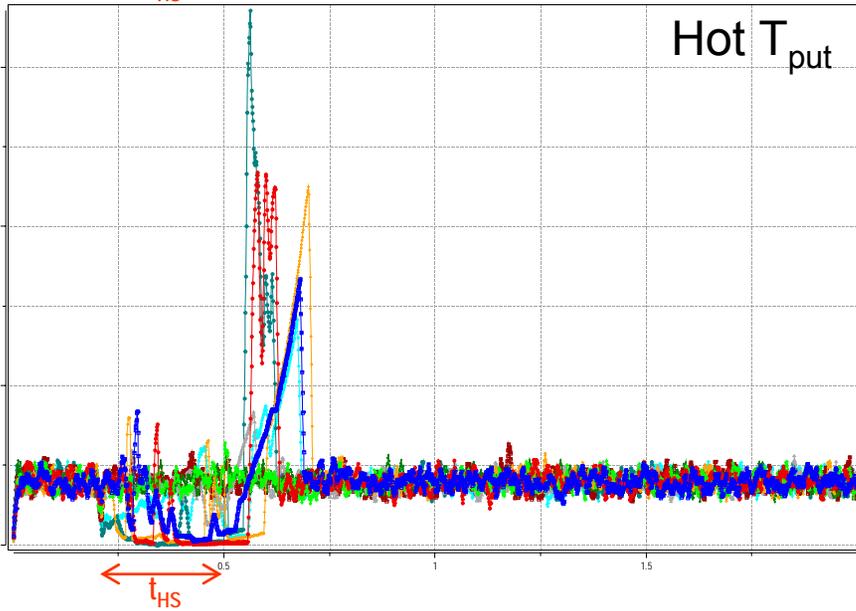
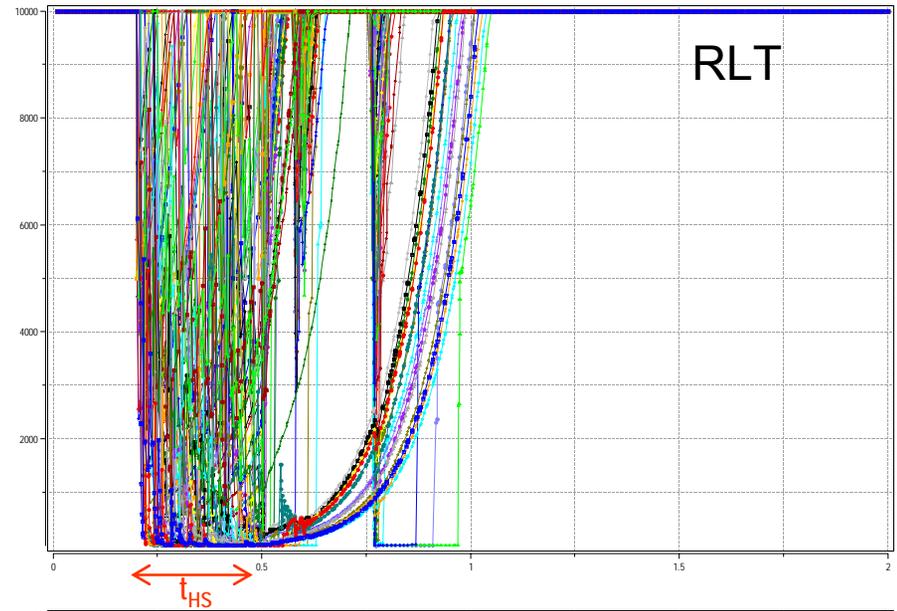
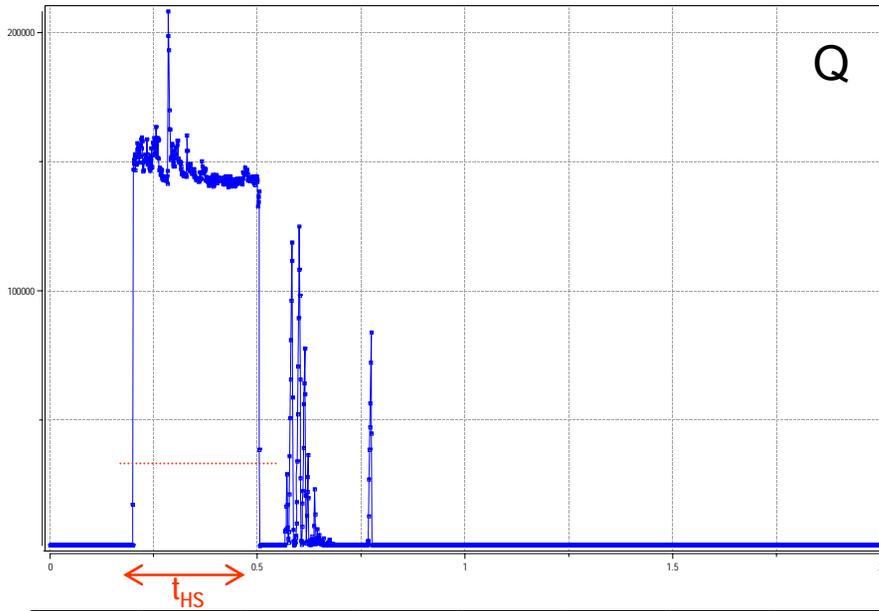
# QCN - IG - FTree5L



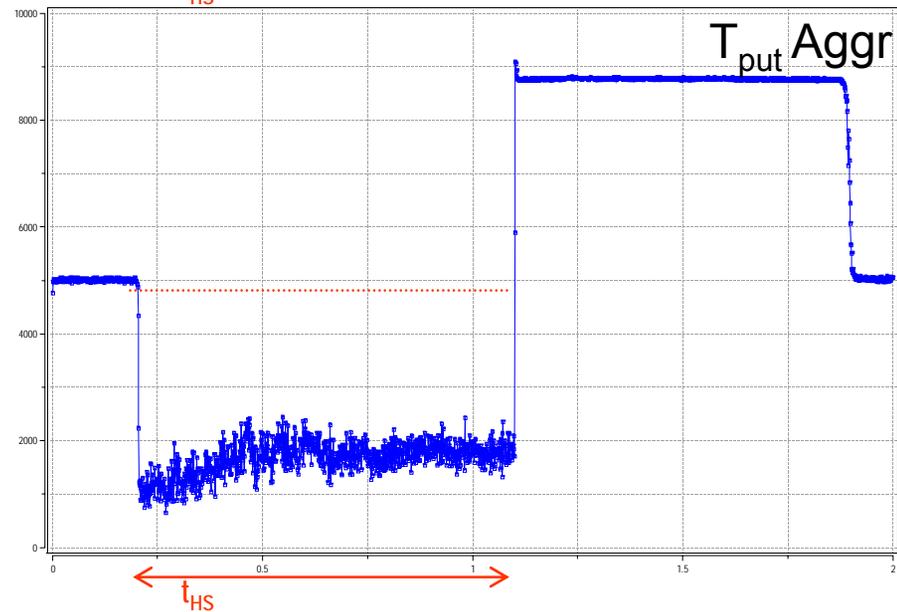
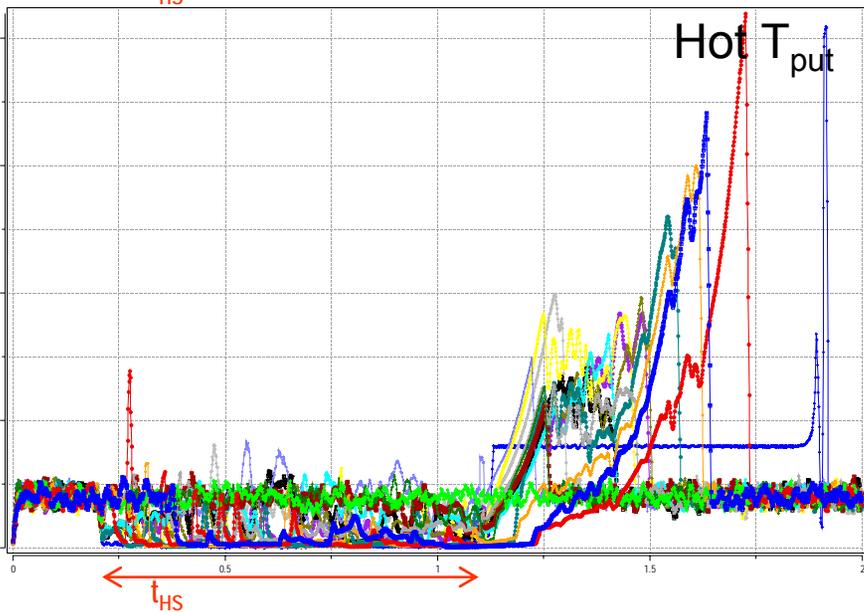
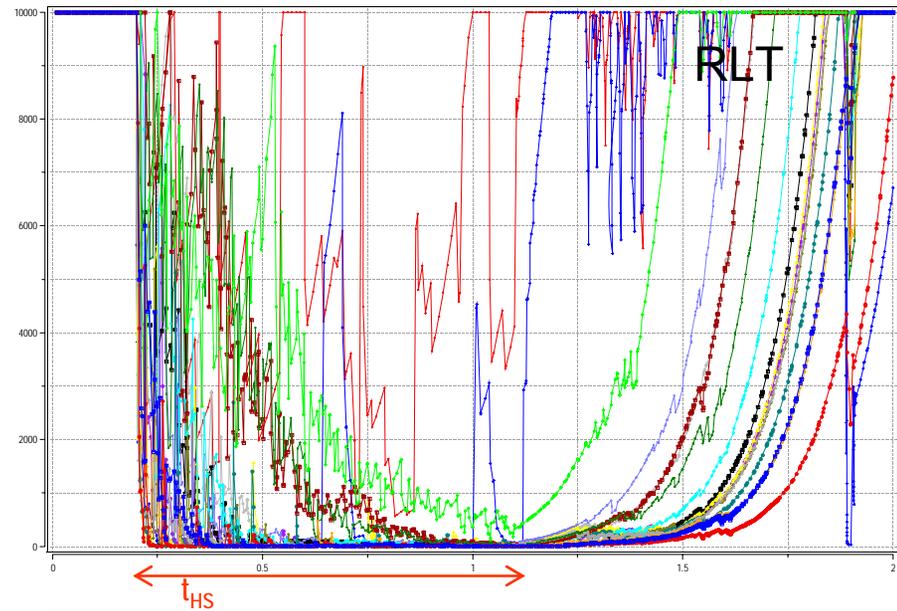
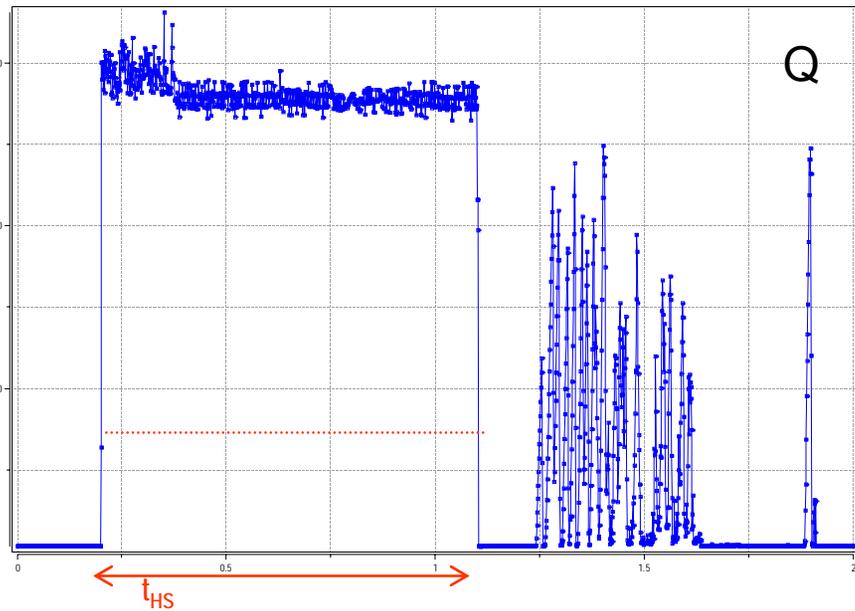
# QCN - IG - FTree7L



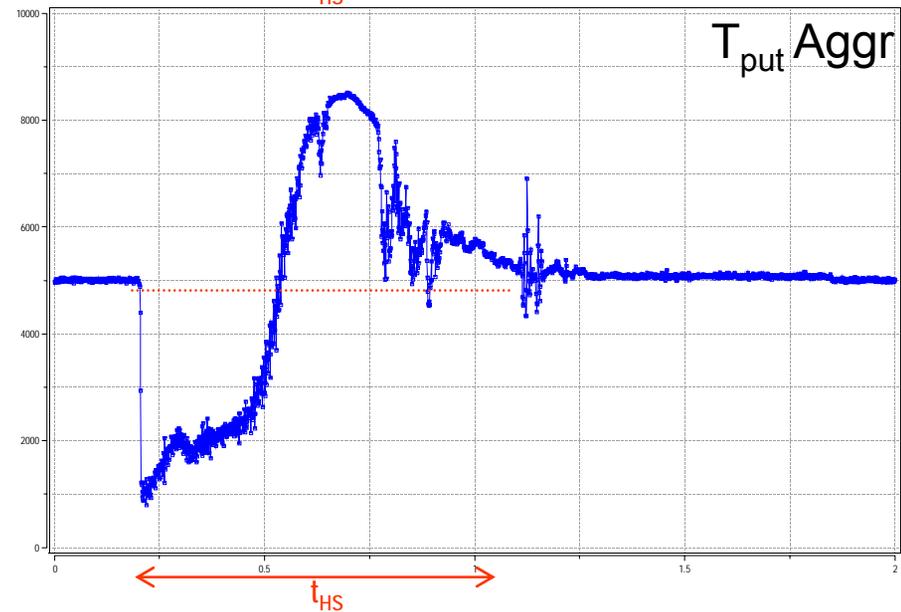
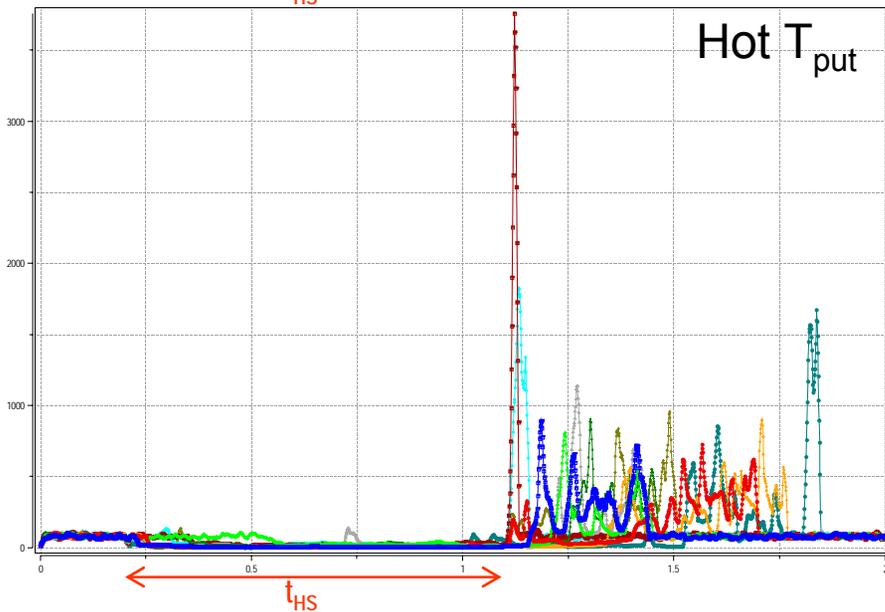
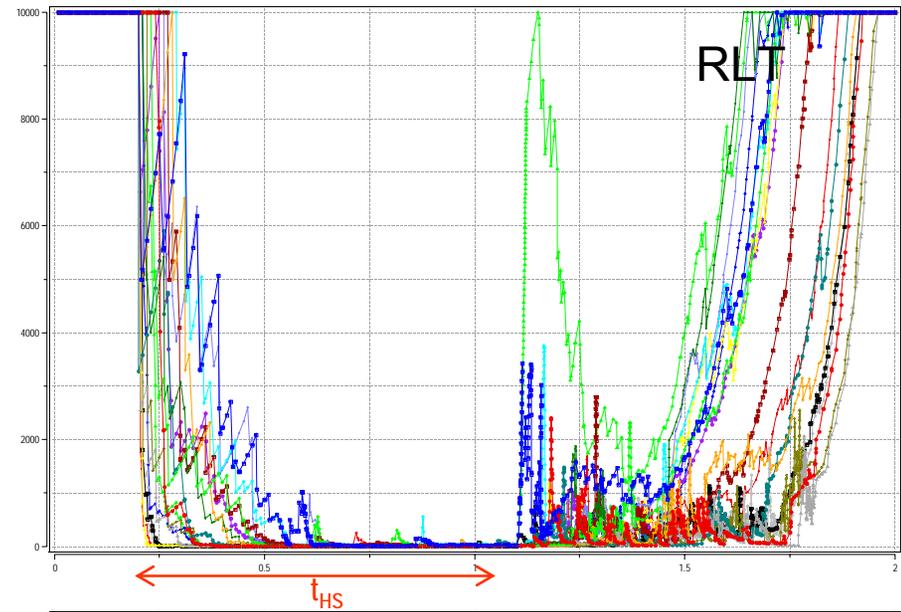
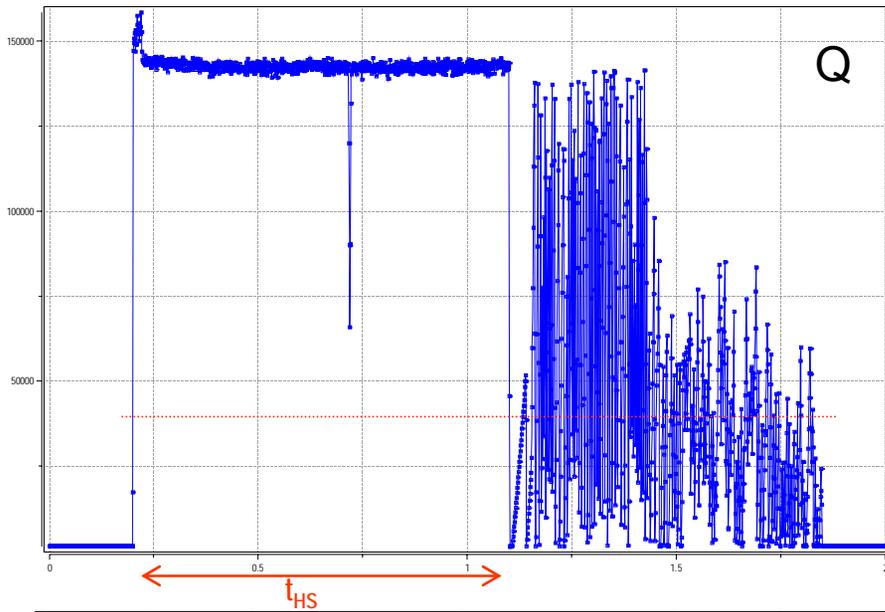
# QCN - OG - FTree5L - low e2e delay



# QCN - OG - FTree5L - low e2e delay - long HS



# ECM - OG - FTree5L - low e2e delay - long HS



# First analytical observations in single hop QCN

1. Assume steady state w/ IG load
  1. preliminary partial linearization of QCN
  2. statistical analysis with M/M/1

# 2pt QCN Linearization

Coarse ODE model of 2-pt QCN

1. Conservation:  $dq/dt = \text{HSD} * \lambda(t) - \mu_{\text{HS}}$
2.  $q(s) = \text{HSD} * \lambda(s) / s$
3. Feedback: ECM's 2 state vars  $q_{\text{off}}$  and  $q'$  reduced 2D  $\rightarrow$  1D
4. Neg. FB:  $\text{Fb}_-(t) = -(q(t) - Q_{\text{eq}}) + w * (dq/dt) / (\mu_{\text{HS}} * p_s)$ 
  1.  $\text{Fb}_-(t) < 0$
  2.  $(q(t) - Q_{\text{eq}}) + w * (dq/dt)$  : calculated in situ (per switch queue) and 6b quantized as a single variable
5.  $\text{Fb}_-(s) \approx G * [1 + w * s / (\mu_{\text{HS}} * p_s)]$
6. Rate decrease control (RDC):  $d\lambda(t)/dt = G_d * \lambda(t) * \lambda(t-\tau) * p_s * \text{Fb}_-(t-\tau)$
7.  $\delta \text{MD}(t) / \delta \text{Fb}_-(t-\tau) \approx G_d * p_s * (\mu_{\text{HS}} / \text{HSD})^2 \Rightarrow$
8. Sensitivity of  $G_d = \delta \text{MD}(t) / \delta \text{Fb}_-(t-\tau) * (\mu_{\text{HS}} / \text{HSD})^{-2} / p_s$ .

QCN's RDC is similar to ECM; however, the switching function decr  $\leftrightarrow$  incr differs.

$q(t)$  =queue occupancy; HSD=no. of hot flows, each with rate  $\lambda(t)$ , at hotspot (HS) served w/ rate  $\mu_{\text{HS}}$

# QCN recovery: Ongoing...

Rate Increase Control (RIC) in 3 concurrent/sequential phases:

1) Extra / Fast recovery: Reclaim the previous  $R_d$  by binary increase

$$t_{Fb-} < t \leq t_{FR} \Rightarrow r_{new} + \sum_{i=0}^5 \sum_{j=0}^{25..100T_f} f(R_0, R_d, t_{Fb-}) \rightarrow r_j(t)dt + \frac{R_d}{2^i} \approx 1 - e^{-kt}$$

\* double integrator w/ (a) initial condition  $R_d$ ; (b) enable  $t_{Fb-}$ ; (c) reset. Executes only once after enable. Byte-based counters, possibly enhanced w/ timer (switch condition?).

2) Active (AI) or hyperactive increase (MI): Probing for the previous equil.

$$t_{FR} \leq t < t_{AI/MI} \Rightarrow r_{new} \approx e^{xt}$$

\* the choice of AI vs. MI depends on traffic and CM target

3) Drift: MI to claim excess  $C$  (newly available Bw)

$$t_{AI/MI} \leq t \Rightarrow \text{multiplicative increase}$$

The improving RIC(t) algorithms and analogy w/ CUBIC are subject of further investigation...

# Optimal admission control rate

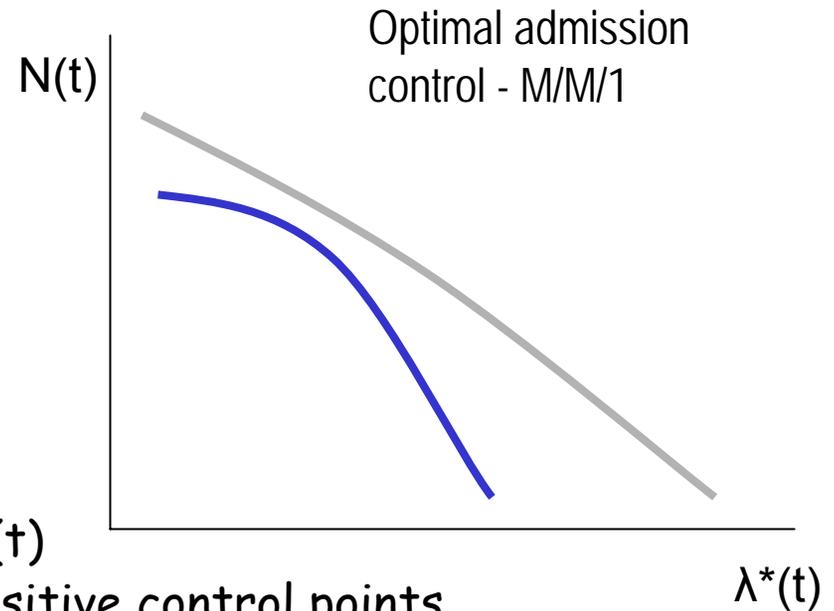
- **System specs**

- Objective: minimize the cost of holding packets in the system (M/M/1)
- Control variable: rate
- Observable system variable: number of packets  $N(t)$  in  $Q$
- Control epoch: every unit time (e.g., 1 MTU = 1us)
- Methodology: dynamic programming

- **Optimal policy:** Arrival rate  $\lambda^*(t)$  is a non increasing function of the load  $N(t)$ .

→ exact functional form depends on the system calibration.

# What is the correlation coefficient of the controlled rate vs. system load ?



## 1. Correlation across all control points

- Arrival rate at time  $t$ ,  $A(t)$
- Number of packets at time  $t$ ,  $N(t)$

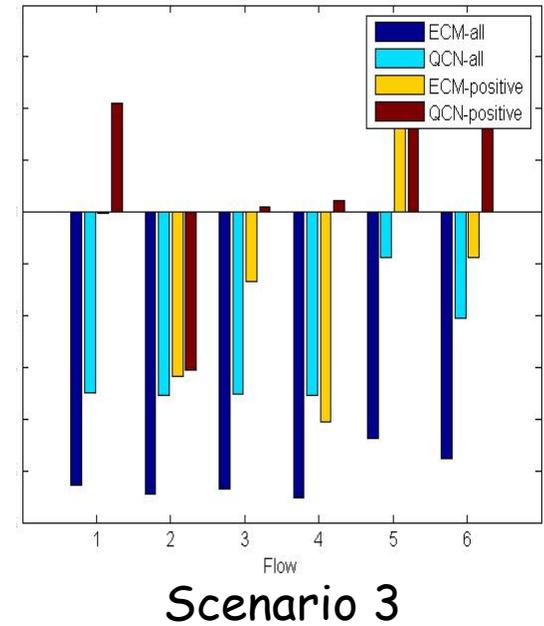
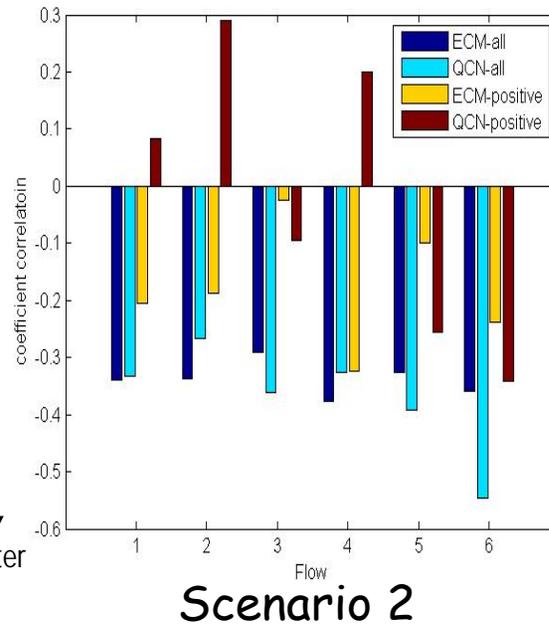
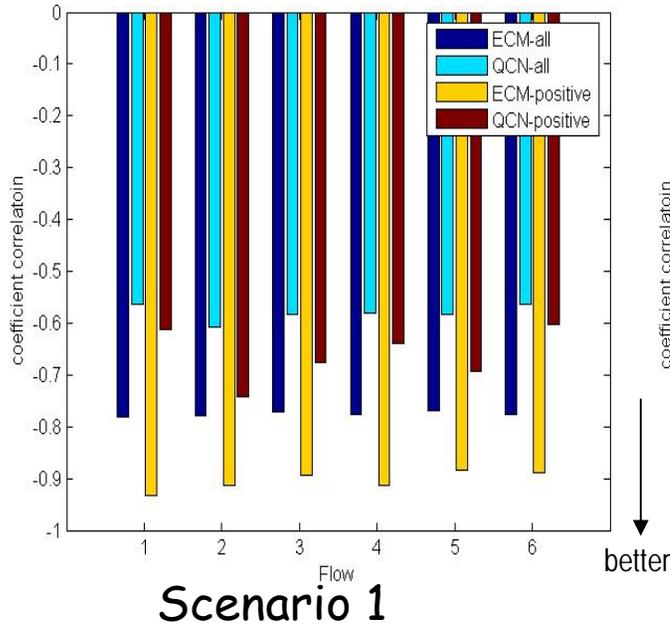
## 2. Coefficient of correlation across positive control points

- Positive increment of arrival rate,  $A+(t)$
- Number of packet changes,  $N+(t)$

## ➔ Our expectation for the coefficient of correlation

- o Negative values should be for (1) and (2).
- o The **closer to -1** → **the more causal/accurate** the control action (RLT).
- o Deviations from -1 may result from:
  - o (a) lags and control delay;
  - o (b) inaccurate calibration of system dynamics.

# Correlation Results

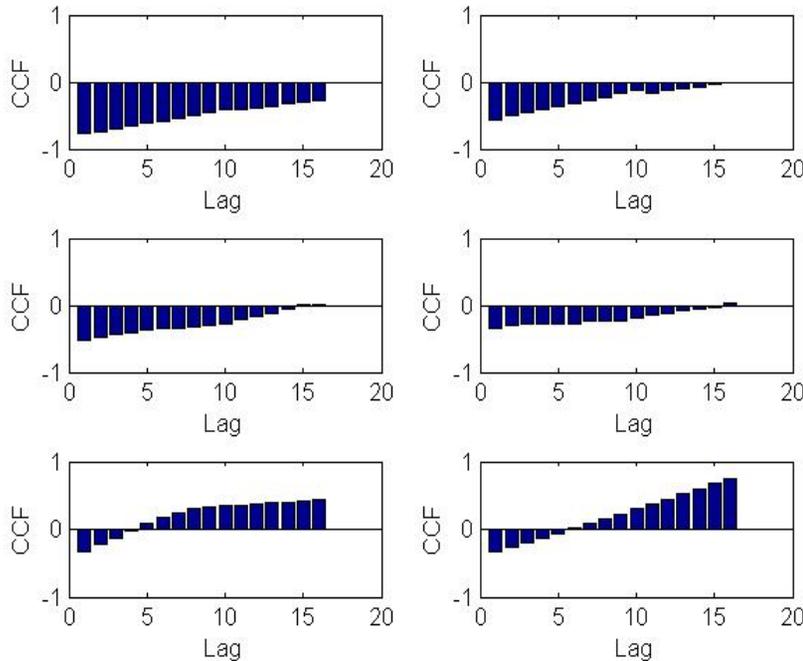


Mean no. pkts/Q	Scenario 1	Scenario 2	Scenario 3
ECM	22963	23747	85280
QCN	21429	14179	279940

QCN's open loop recovery is self-healed through negative feedback (switched loop). This increases its sensitivity to load changes, as in e.g. OG.

# Controlled arrival rate ( $\lambda(t)$ ) of flow A v.s $N(t) \rightarrow$ CCF

ECM-scenario 1	QCN-scenario 1
ECM-scenario 2	QCN-scenario 2
ECM-scenario 3	QCN-scenario 3

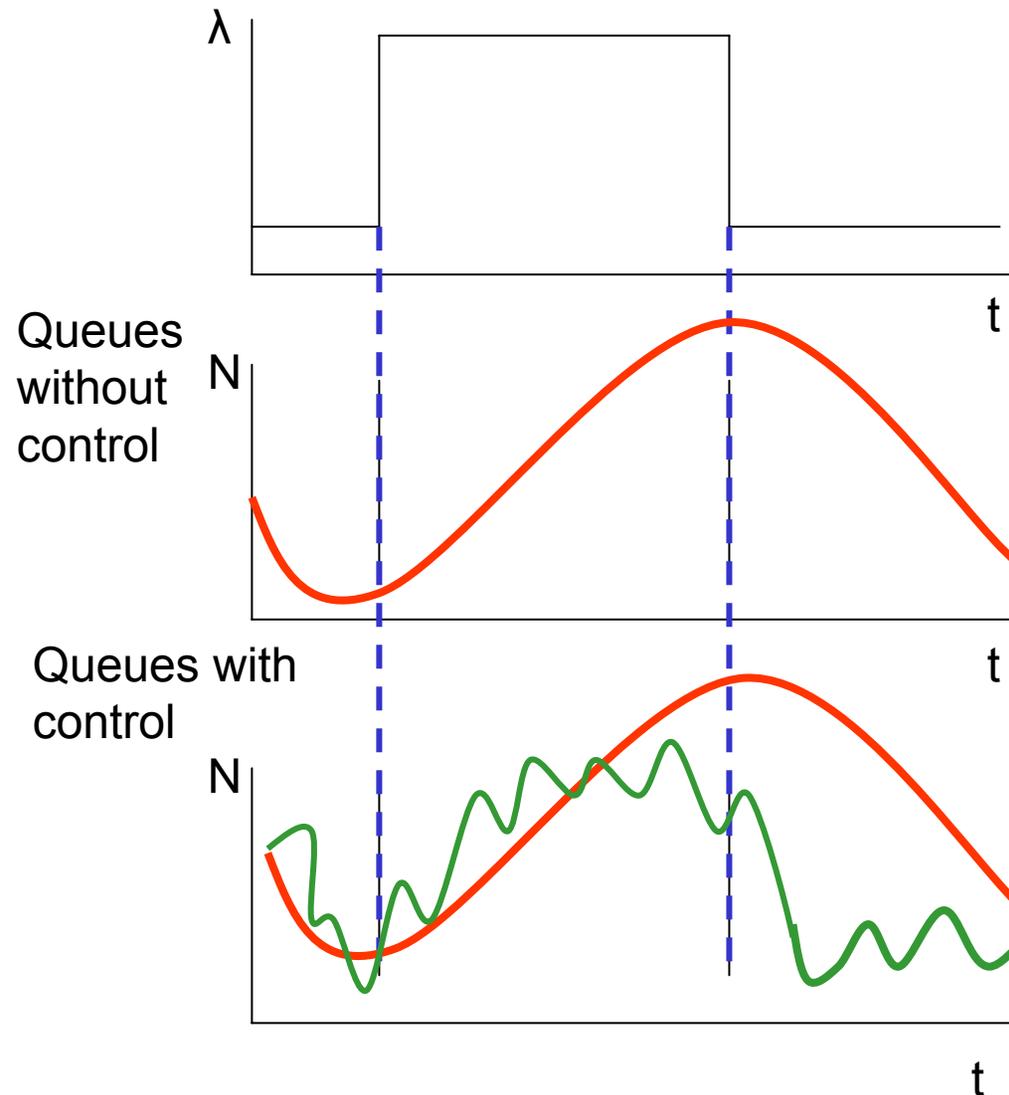


- Each lag corresponds to  $T_{\text{sample}} = 10$  ms. That's how the preliminary statistics were collected.
- The transition from negative to positive CCF can be interpreted as following:
  - Total lags of negative values correspond to the fluctuated period due to the controlled action. (Since we know our load is generated in certain way. Open loop should show smooth transition.)
  - The longer the period is, the less fluctuation is brought into the system. (More like off-line analysis for system calibration.)
- The above holds based on this specific workload generator.

# Cross-correlation factor: Open vs. Closed Loop

Note: if the controlled rate is highly correlated with the queue status, their CCF will be similar to the ACF of an uncontrolled queue. In our case here,  $T=2-300\text{ms}$ .

The control loop brings in queue fluctuations (green at bottom). CCF of queues and controlled arrival rate will reflect length of control cycle / lag. Ideally a shorter lag (1-3) is preferable.

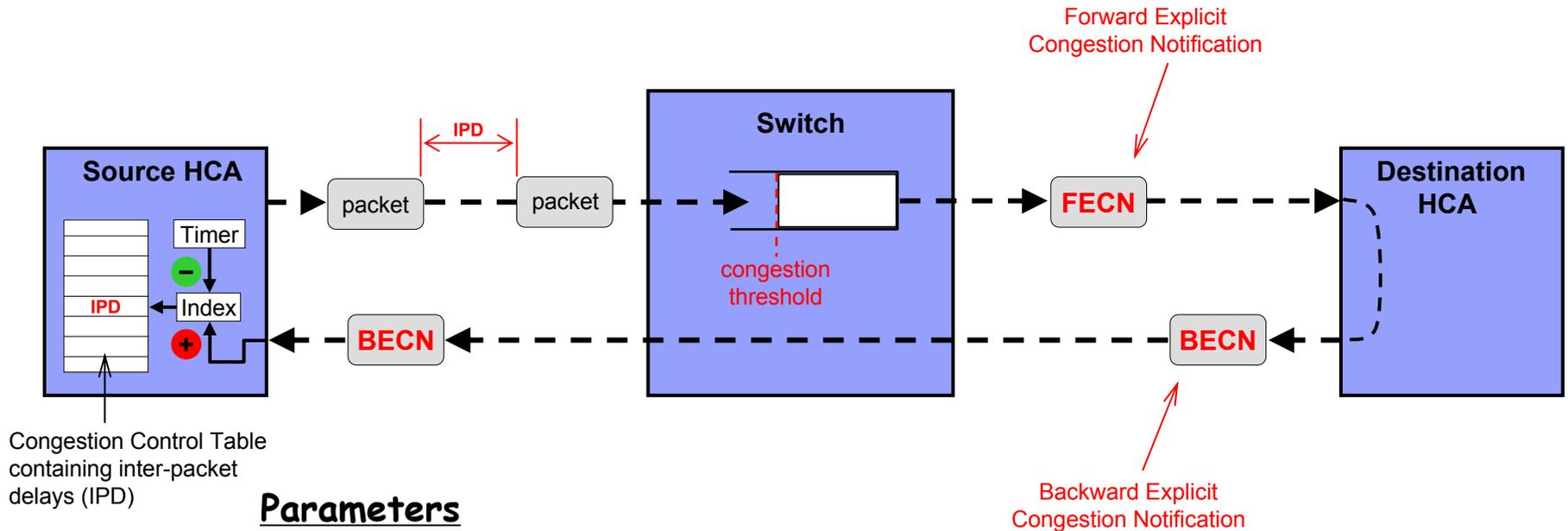


Change in  $N(t)$   
due to Fb\_-based  
control.

# Summary Correlation Analysis

- Instead of utilizing positive feedback directly, QCN *autonomically* recovers the arrival rate in quasi-open loop, disregarding the load sensors (CPID)
  - pro: w/ proper dither, may avoid resonance (synchro)
  - "overshoot" feedback: self-heals on (strongly) negative feedback
- Hence higher mean control delay/lags.
  - the recovery algorithms packed in RIC(t) elicit further study... we'll help the proponents
  - Next: Comparison w/ IBA's CM

# QCN vs. IBA CM : Similarity or Contrast?



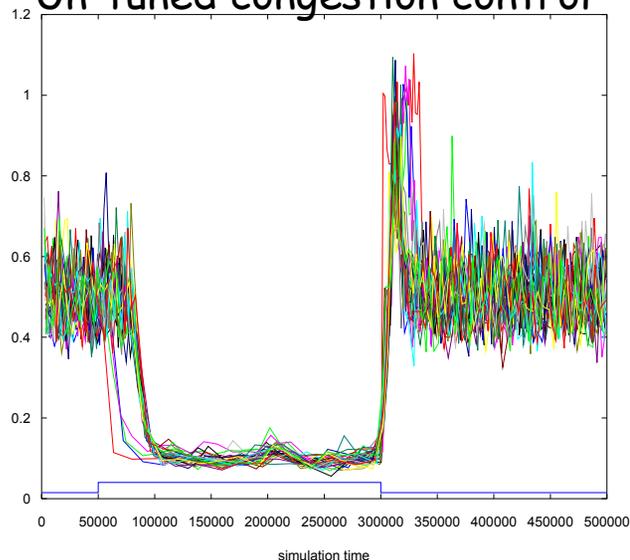
## Parameters

- Switch queue threshold ( $\sim Q_{eq}$ ): *sw\_th*
  - congestion detection and FECN marking
- IPD table size = 64-256 entries
  - dynamic range of per flow rate control
- highest IPD entry: *max\_ipd*
  - largest inter-packet gap =  $1/R_{min}$
- IPD index increment: *ipd\_idx\_incr*
  - step on each new BECN  $\Rightarrow$  rate control granularity
- IPD recovery timer: *rec\_time*
  - timeout of the autonomic rate increase timer

# IBA CCA/CM Stable?

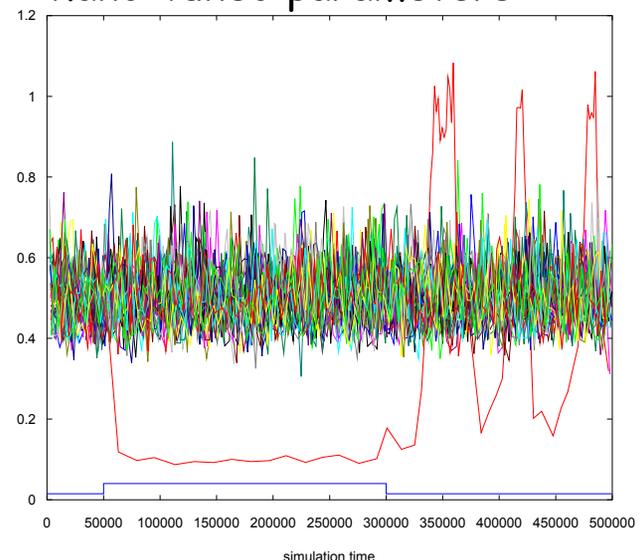
- Qualified "yes" => needs online tuning
  - ❖ easy for small fabrics w/ simple traffic, hard for others...
- Param *tuning* required per (1) fabric architecture and (2) traffic pattern
- Conditional stability: CCA has high sensitivity to traffic...

Un-tuned congestion control



output[31]  
output[30]  
output[29]  
output[28]  
output[27]  
output[26]  
output[25]  
output[24]  
output[23]  
output[22]  
output[21]  
output[20]  
output[19]  
output[18]  
output[17]  
output[16]  
output[15]  
output[14]  
output[13]  
output[12]  
output[11]  
output[10]  
output[9]  
output[8]  
output[7]  
output[6]  
output[5]  
output[4]  
output[3]  
output[2]  
output[1]  
output[0]  
hotspot

hand-tuned parameters



output[31]  
output[30]  
output[29]  
output[28]  
output[27]  
output[26]  
output[25]  
output[24]  
output[23]  
output[22]  
output[21]  
output[20]  
output[19]  
output[18]  
output[17]  
output[16]  
output[15]  
output[14]  
output[13]  
output[12]  
output[11]  
output[10]  
output[9]  
output[8]  
output[7]  
output[6]  
output[5]  
output[4]  
output[3]  
output[2]  
output[1]  
output[0]  
hotspot

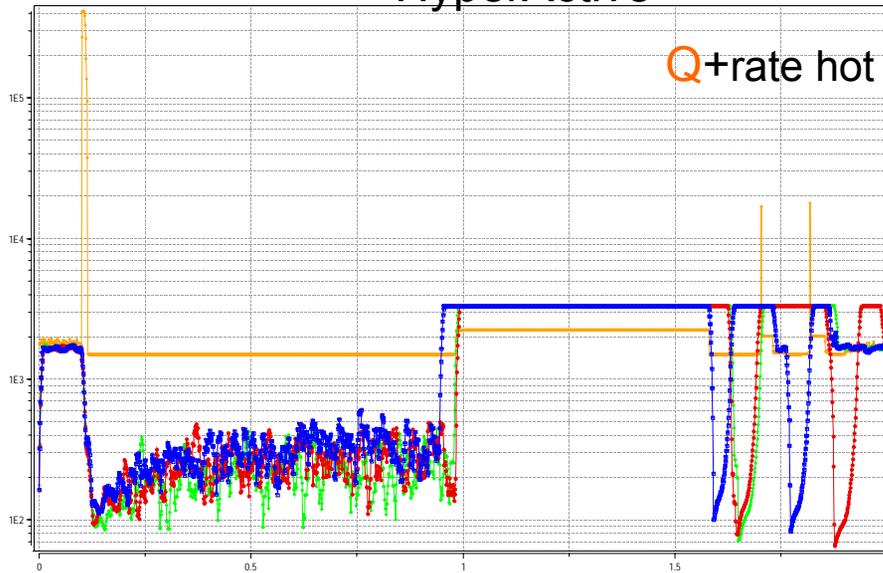
# Robust CM Design vs. Implementation - RLT

## II) Rate Limiter Aliasing

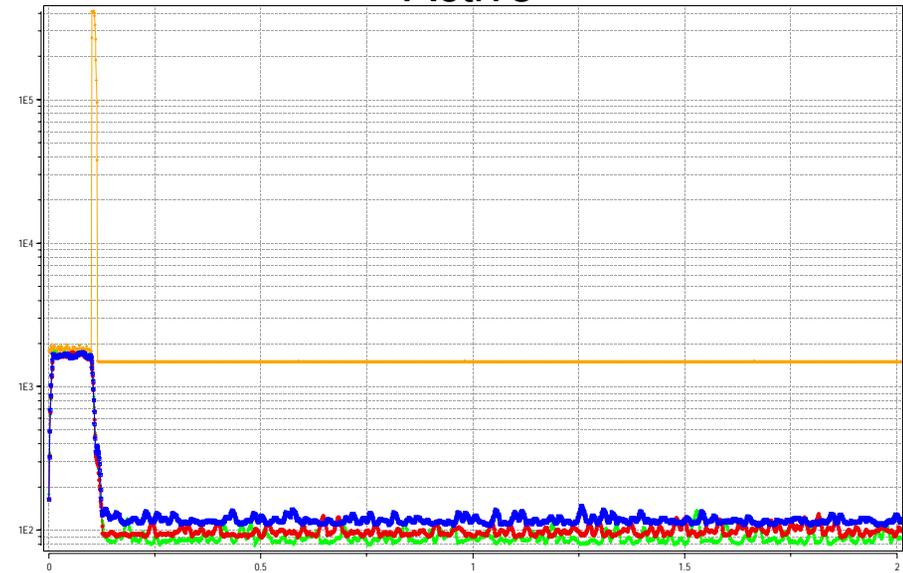
- Scenario
  - Single-path flows
  - Multiple flows sharing one rate limiter (RL)
  - One rate applied to all flows in same RL ... see results next page
- Observations
  - Multi-CP considerations apply (see CPID)
    - ✓ Negative feedback should be applied if generated by *any* CP
    - ✓ Positive feedback should be applied if generated by *all* CPs; simplifications are possible
  - Why not keep track of rate per flow?
    - ✓ Main RL complexity is in buffers
    - ✓ Maintain small table mapping flow ID to rate
    - ✓ Each rate is updated and applied per flow *without* per-flow queuing
      - Hence, flows can still not move faster than the slowest one
    - ✓ Does not require CPID (for this subproblem per se)

# Coalesced RLT in QCN: 4 flows to 1-RLT, OG

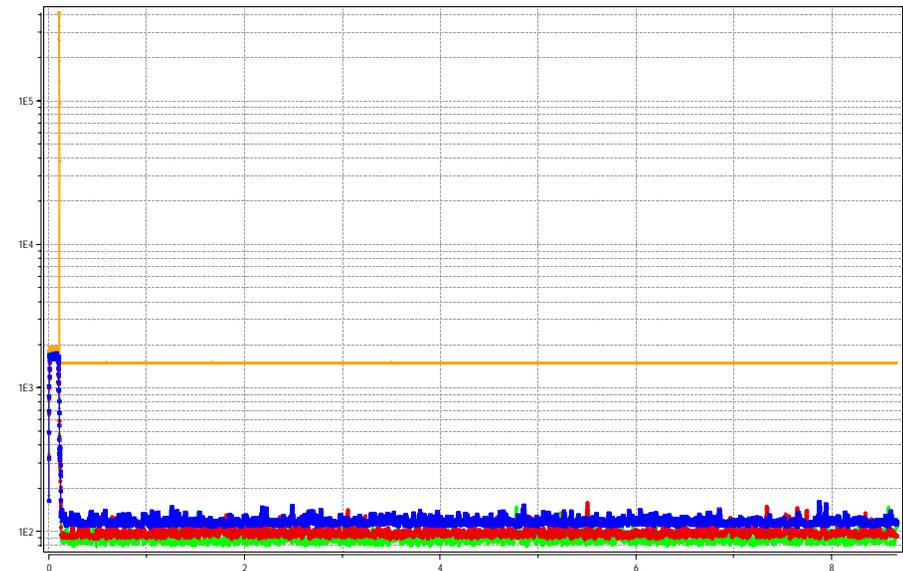
HyperActive



Active



- Step impulse response: HS
  - Start 0.1
  - Stop 0.2
- logY scale
  - Q and per flow Tput plotted together
    - ✓ enables qualitative observation
- one rate to 4 flows => Active Incr. QCN trapped in a local minimum



# Robust CM Design vs. Implementation - CPID

## III) Stale CPIDs

- Rate limiter is associated with specific CPID at some (low) rate
  - a) All traffic to that CPID ceases or
  - b) the CPID disappears altogether (network failure, equipment replacement)
- How do we ensure full rate recovery?
  - a) Closed Loop: Query feedback from specific CPID (e.g. using probing)
  - b) Open Loop: Use conservative timer-based recovery

# Robust CM Design vs. Implementation - Positive Fb

## IV) On **Positive Feedback Signals**

- One flow may receive feedback from multiple congestion points (CPs)
  - May be on the same path (single-path case) or on different paths (multi-path case)
  - Basic conundrum
    - ✓ Negative feedback should be applied if generated by *any* CP
    - ✓ Positive feedback should be applied if generated by *all* CPs
- **Q:** How to obtain simultaneous feedback from all CPs?
- **A1** (Single-path case): Probe the entire path - no CPID association and no tags required
- **A2** (Multi-path case): Disambiguate positive feedback signals by associating CPID with RL
  - With network-triggered sampling (BCN) this requires tagging all frames
  - With NIC-triggered sampling (probes) this requires tagging probes

# Robust CM Design vs. Implementation - MP

## I) Multi-path Congestion

- Scenario
  - Flow  $f$  traverses  $P$  paths
  - Total spare capacity on these  $P$  paths without flow  $f$  equals  $C$
  - Ideal load balancing
- Observations
  - No congestion will occur as long as  $\text{rate}(f) < C$
  - As soon as  $\text{rate}(f) > C$ , all  $P$  paths will be congested
  - Hence, excepting naive hashing, the issue of CP disambiguation is arguably moot
- Corollary
  - If congestion notifications are not coming from all paths, the load balancer is not doing a great job

# Summary

- Both ECM and QCN have clear merits
  - no judgement on complexity and performance can be imparted at this early stage
- QCN
  - works in many typical input-generated hotspots
    - ✓ in both SS and trees
    - ✓ outperforms IBA's CCA (despite conceptual similarity, Fb-oblivious NIC state-based recovery)
  - QCN's aggressive-gentle-aggressive recovery is tantalizing
    - ✓ however, its linearization is not a trivial exercise (not CUBIC-like)
  - Questionable
    1. Open positive feedback loop in non-linear datacenter networks
    2. Feedback degree reduction and calculation in switch
    3. Complexity of a better recovery scheme
      1. despite algorithmical cleverness, QCN can still get "stuck at low"; fixes may entail stability issues
      2. despite appearances, 2pt QCN's nearest kin is IBA CM, not TCP-CUBIC
- 2-pt. QCN and ECM are located at different stability/cost points
  - CM robustness and design issues elicit more attention
- Node- and (sub)path-centric CM are not exclusive when traffic is unknown and overload persistence is a heavy-tailed distribution.
  - Convergence time in seconds is not acceptable, even if fairness is not an issue
  - However, convergence time should not unduly affect the stability margins

# Proposal: Close the "+ Loop"

- Plan A) Use ECM as base CM
  - most tested scheme in last 2 yrs.
  - predictable and *improvable* performance
    - ✓ ample room above (perf. optimizations not squeezed)
      - E2CM and QCN elements
  - improve it w/ the best QCN and E2CM elements (TBD)
  - Tuning with sub-path probing: RP -> CP
- Plan B) QCN in a Fb<sub>+</sub>-driven scheme
  - see Guenter's QCN-P/R proposal

# Backups

## Switch

- Shared-memory output-queued switch
- PAUSE enabled
  - Global high- and low-watermark memory threshold trigger pause and unpauses
  - High watermark  $T_h = M - N * (RTT * B + L_{max})$
  - Low watermark  $T_l = T_h / 2$
  - PAUSE renewed before expiry (take into account RTT)

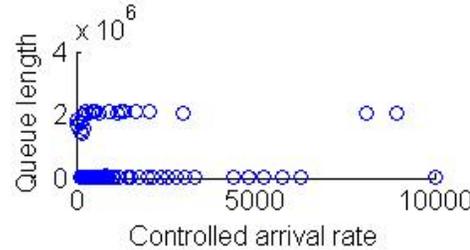
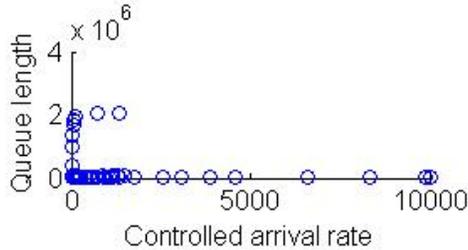
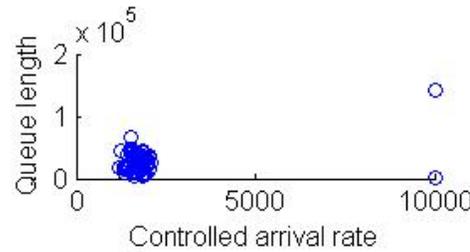
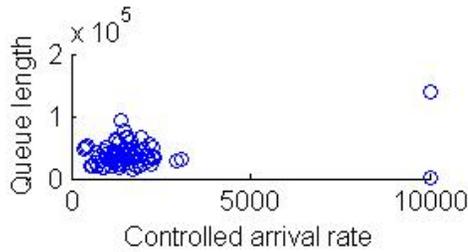
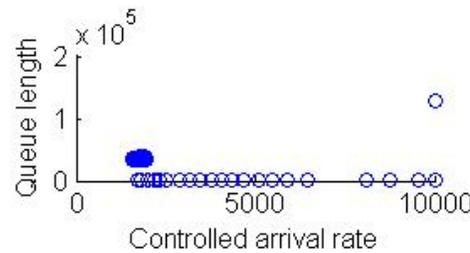
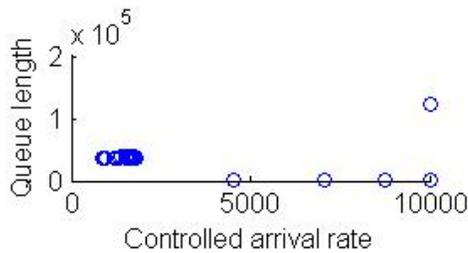
## Adapter

- VOQ-ed per end node
- Round-robin service discipline
- Number of rate limiters unlimited
- Egress buffer flow-controlled using PAUSE (high/low watermarks)

**Lossless operation:  
No frame drops due to buffer overflows!**

# Scatter Plots: Controlled arrival rate of flow A v.s N(t)

ECM-scenario 1	QCN-scenario 1
ECM-scenario 2	QCN-scenario 2
ECM-scenario 3	QCN-scenario 3



- From M/M/1 a downslope is expected. Such trend, however, vanishes from IG to OG.
- Specifically QCN is more sensitive to load changes. This is also reflected by the values of the correlation coefficient above.

# Appendix

- Derivation of optimal admission control
  - Dynamic programming and optimal control ~Dimitri P. Bertsekas.