



Evaluating Transient Duration in a Multi-hop Output Generated Hotspot Scenario

IEEE Interim (Monterey, CA)
Bruce Kwan, Jin Ding, & Brad Matthews
January 22-25, 2007

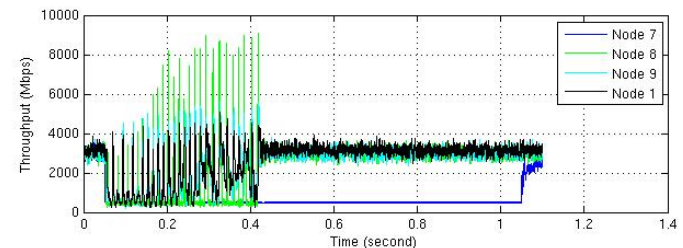
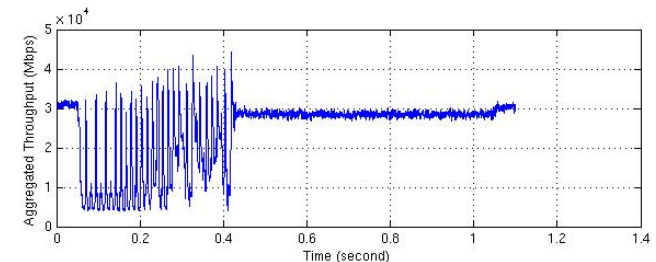
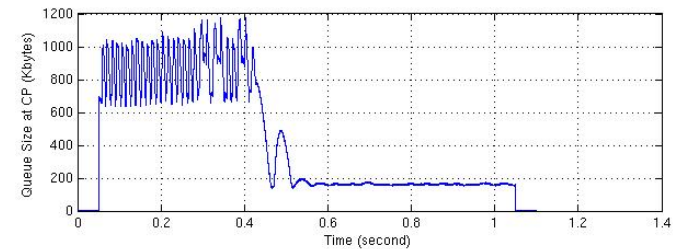
Overview

- Goals
- Observations
- Simulation Setup
- Results
- Summary & Next Steps

Experiment Goals

- Quantify Transient Duration
 - Quantify how long it takes for BCN to resolve a congestion event so that throughput in the system achieves steady state.
 - Similar approach taken in Cisco presentation on Baseline scenario
 - <http://www.ieee802.org/1/files/public/docs2006/au-sim-bergamasco-bcnmax-comparison-110906v2.pdf>
- Trigger Discussions
 - What should the target be?
 - What is required for our target data center applications?
 - What severe-congestion schemes are needed to improve the responsiveness?

Multi Hop Output Generated Scenario (BCN + PAUSE Only)



Overview

- Goals
- Observations
- Simulation Setup
- Results
- Summary & Next Steps

Observations

- Without additional severe congestion schemes, the transient duration for the multi-hop output generated scenario can be on the order of 462 ms (184 x RTT).
- With the use of one of the simple severe congestion management schemes (i.e. BCN-MAX), the transient duration for the multi-hop output generated scenario can be reduced to 52 ms (21 x RTT).
- Small buffer sizes can hamper transient duration (when measured in terms of RTT). Requires additional study.

Overview

- Goals
- Observations
- Simulation Setup
- Results
- Summary & Next Steps

Parameters

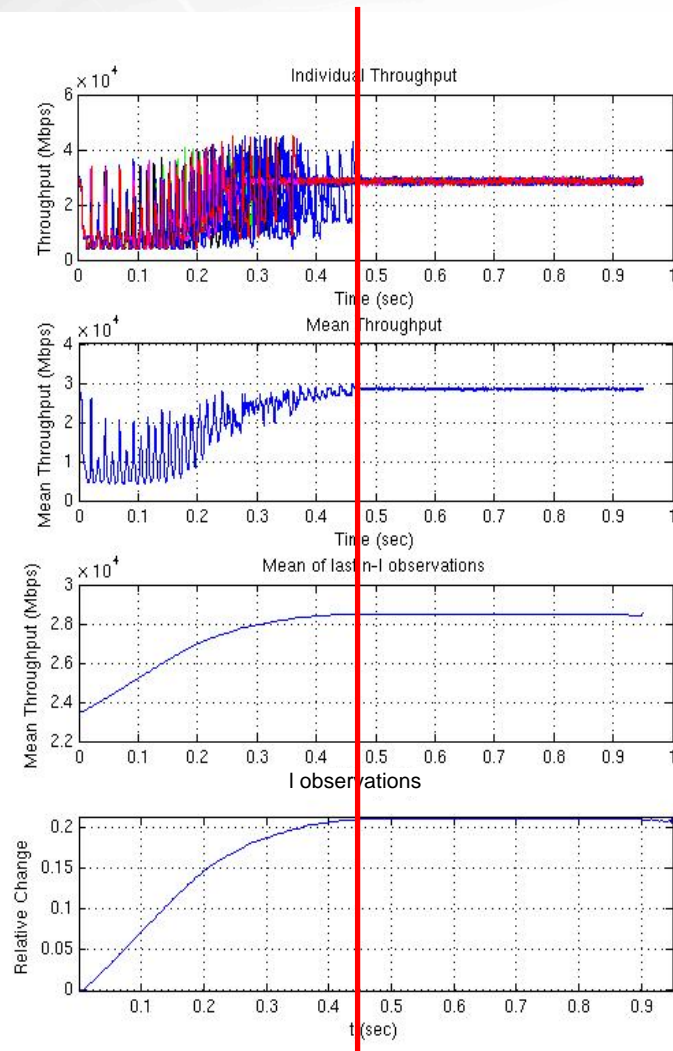
- **Switch Parameters**

- Core switch and edge switches are all 4 port switches
- Buffer Size (B)
 - 600Kbytes/Port
 - 75Kbytes/Port
- Shared Memory Switch Devices, total switch memory size = $4 * B$
- PAUSE Flow Control Settings
 - Applied per ingress port basis based on XON/XOFF thresholds
 - XOFF Threshold = $B - RTT * BW$
 - XON Threshold = $B/2$

- **BCN Parameters**

- Frame Sampling
 - Frames are periodically sampled (on avg) every 75KB (2%)
- $W = 2$
- $Q_{eq} = B/4$
- $R_u = 1\text{Mbps}$
- G_i (Initial)
 - Computed as $(\text{Linerate}/10) * [1/((1+2*W)*Q_{eq})]$
 - Same as in baseline
- G_d (Initial)
 - Computed as $0.5 * 1/((1+2*W)*Q_{eq})$
 - Same as in baseline
- Other BCN Enhancements
 - BCN(Max)
 - OverSampling during severe congestion or consistent oversampling
 - No BCN(0,0)
 - No Self Increase

Defining Transient Duration



- Use *Initial Data Deletion** method to determine when steady state is achieved.

- Study an averaged version across multiple replications.

- Computation of Transient Duration

- Collect multiple simulation replications (i.e. 10) and compute the time series average across the samples.

- Obtain the overall mean throughout the sample run

- Obtain an overall mean from the remaining n-l values as a function of l

- Compute the relative change in the overall mean. The knee of this curve is where the transient duration ends

- Target Transient Duration

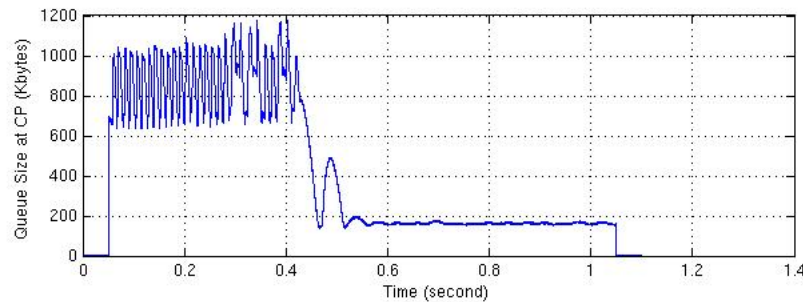
- 10 x RTT

*Raj Jain "The Art of Computer Systems Performance Analysis", 1991, Pg 424

Overview

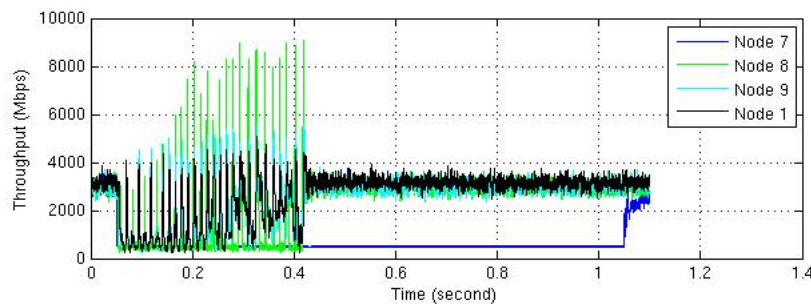
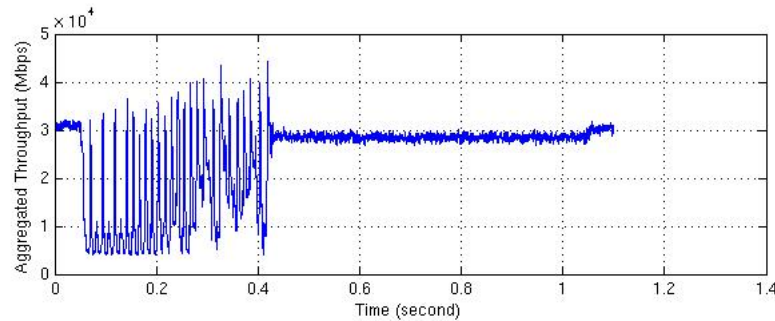
- Goals
- Observations
- Simulation Setup
- Results
- Summary & Next Steps

Basic BCN + PAUSE (No BCN-Max, No Oversampling)



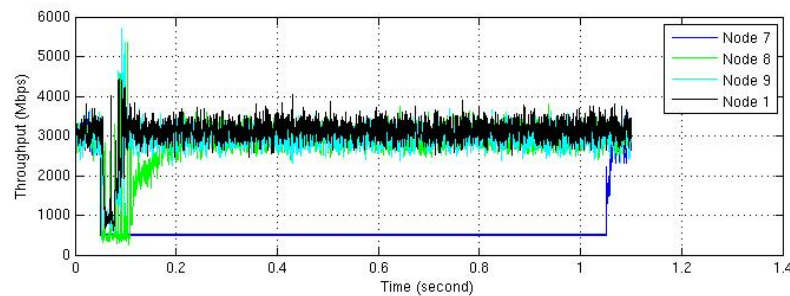
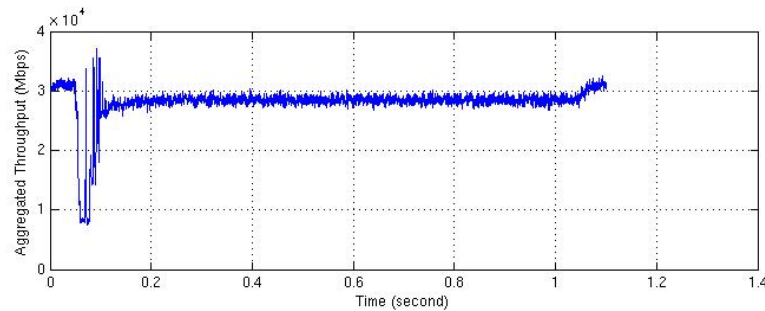
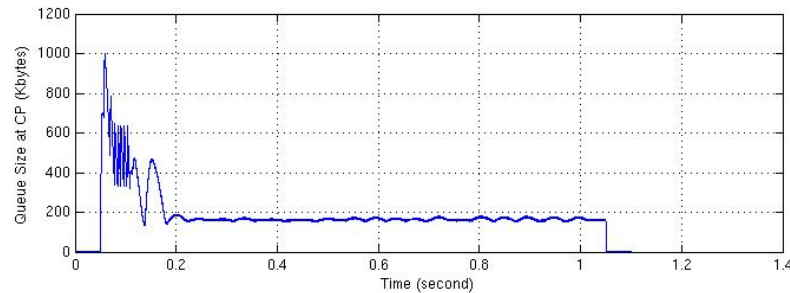
- Transient Duration
 - 462 ms (184 x RTT)

- Next: Consider BCN-MAX



**Multi Hop Output Generated Scenario
(Single Replication)**

BCN w/BCN-MAX + PAUSE (No Oversampling)



*Multi Hop Output Generated Scenario
(Single Replication)*

- Transient Duration
 - 52 ms (21 x RTT)
- Next: Consider additional severe congestion enhancements

Transient Duration Results (Effects of Qoffset/Qdelta Range)

- Setup

- Buffer Size = 600KB/port
- Sampling Rate = 2%

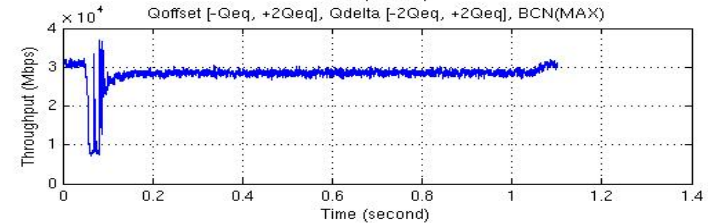
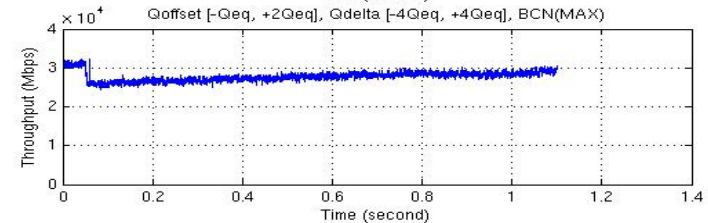
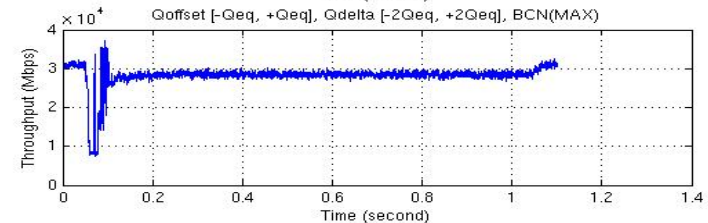
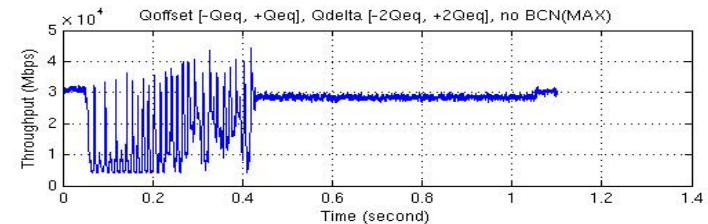
Case	BCN (Max)	Qoffset & Qdelta Range	Transient Duration
1	No	Qoffset: (-Qeq,+Qeq)	462 ms (184 x RTT)
2	Yes	Qdelta: (-2Qeq, +2Qeq)	52ms (21 x RTT)
3	No	Qoffset: (-Qeq,+2Qeq)	170 ms (68 x RTT)
4	Yes	Qdelta: (-4Qeq, +4Qeq)	274 ms (109 x RTT)
5	No	Qoffset: (-Qeq,+2Qeq)	108ms(43 x RTT)
6	Yes	Qdelta: (-2Qeq, +2Qeq)	62ms(25 x RTT)

Transient duration is not further improved by increasing the range of Qoffset (nor Qdelta).

Effects of Qoffset/Qdelta Range (Aggregate Throughput)

- Set up
 - Buffer Size 600KB/port
 - Sampling Rate = 2%

Case	BCN (Max)	Qoffset & Qdelta Range	Transient Duration
1	No	Qoffset: (-Qeq,+Qeq) Qdelta: (-2Qeq, +2Qeq)	462 ms (184 x RTT)
2	Yes	Qoffset: (-Qeq,+Qeq) Qdelta: (-2Qeq, +2Qeq)	52ms (21 x RTT)
4	Yes	Qoffset: (-Qeq,+2Qeq) Qdelta: (-4Qeq, +4Qeq)	274 ms (109 x RTT)
6	Yes	Qoffset: (-Qeq,+2Qeq) Qdelta: (-2Qeq, +2Qeq)	62ms(25 x RTT)

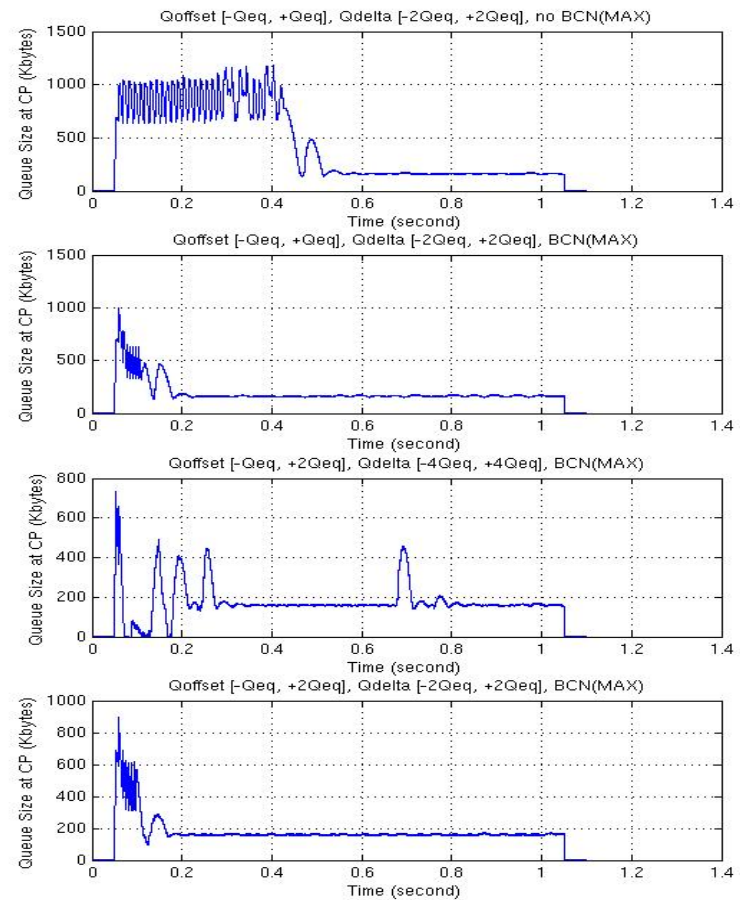


Effects of Qoffset/Qdelta Range (Queue Behavior @ CP)

- Set up

- Buffer Size 600KB/port
- Sampling Rate = 2%

Case	BCN (Max)	Qoffset & Qdelta Range	Transient Duration
1	No	Qoffset: (-Qeq,+Qeq) Qdelta: (-2Qeq, +2Qeq)	462 ms (184 x RTT)
2	Yes	Qoffset: (-Qeq,+Qeq) Qdelta: (-2Qeq, +2Qeq)	52ms (21 x RTT)
4	Yes	Qoffset: (-Qeq,+2Qeq) Qdelta: (-4Qeq, +4Qeq)	274 ms (109 x RTT)
6	Yes	Qoffset: (-Qeq,+2Qeq) Qdelta: (-2Qeq, +2Qeq)	62ms(25 x RTT)



Transient Duration Results (Effects of Sampling Behavior)

- Setup

- Buffer Size = 600KB/port
- Qoffset: (-Qeq, +Qeq)
- Qdelta: (-2Qeq, +2Qeq)

Case	BCN (Max)	Sampling	Transient Duration
1	No	2%	462 ms (184 x RTT)
2	Yes	2%	52ms (or 21 x RTT)
3	Yes	2% (Qlen < Qsc) 10% (Qlen >= Qsc)	86 ms (34 x RTT)
4	Yes	2% (Qlen < Qsc) 20% (Qlen >= Qsc)	144 ms (58 x RTT)
5	Yes	2% (Qlen < Qsc) 30% (Qlen >= Qsc)	96 ms (38 x RTT)
6	Yes	5%	52ms (21 x RTT)
7	Yes	10%	38 ms (15 x RTT)

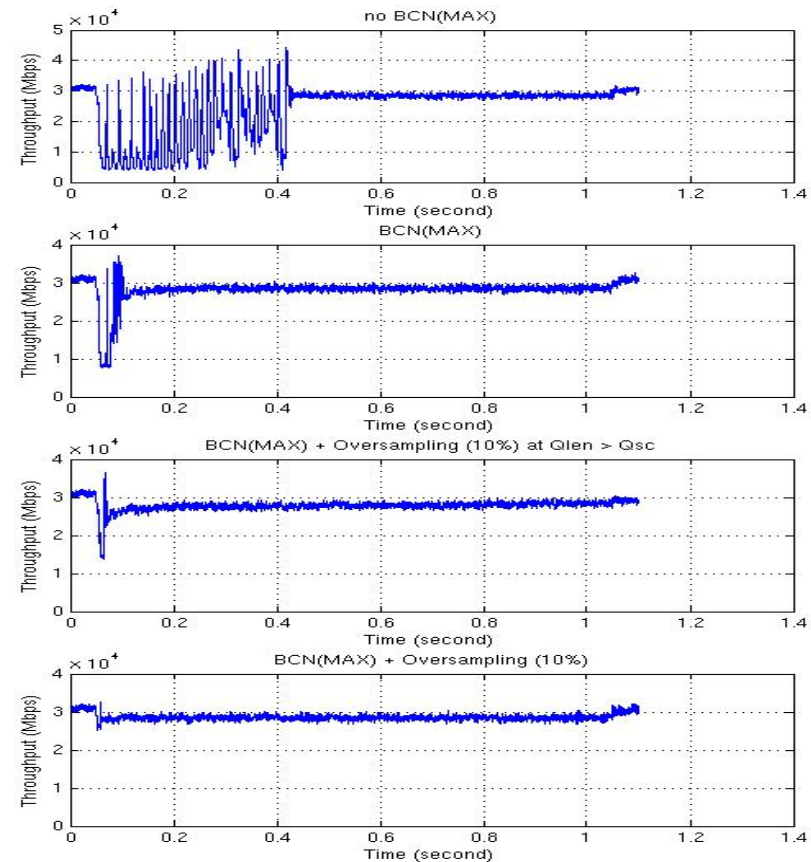
An overall increase in the sampling rate aids in reducing the transient duration down to 38ms. Need to consider also oversampling when Qlen = 0 to improve the adaptive oversampling scheme.

Effects of Sampling Behavior (Aggregate Throughput)

- Set up

- Buffer Size = 600KB/Port
- Qoffset = (-Qeq, +Qeq)
- Qdelta = (-2Qeq, +2Qeq)

Case	BCN (Max)	Sampling	Transient Duration
1	No	2%	462 ms (184 x RTT)
2	Yes	2%	52ms (or 21 x RTT)
3	Yes	2% (Qlen < Qsc) 10% (Qlen >= Qsc)	86 ms (34 x RTT)
7	Yes	10%	38 ms (15 x RTT)

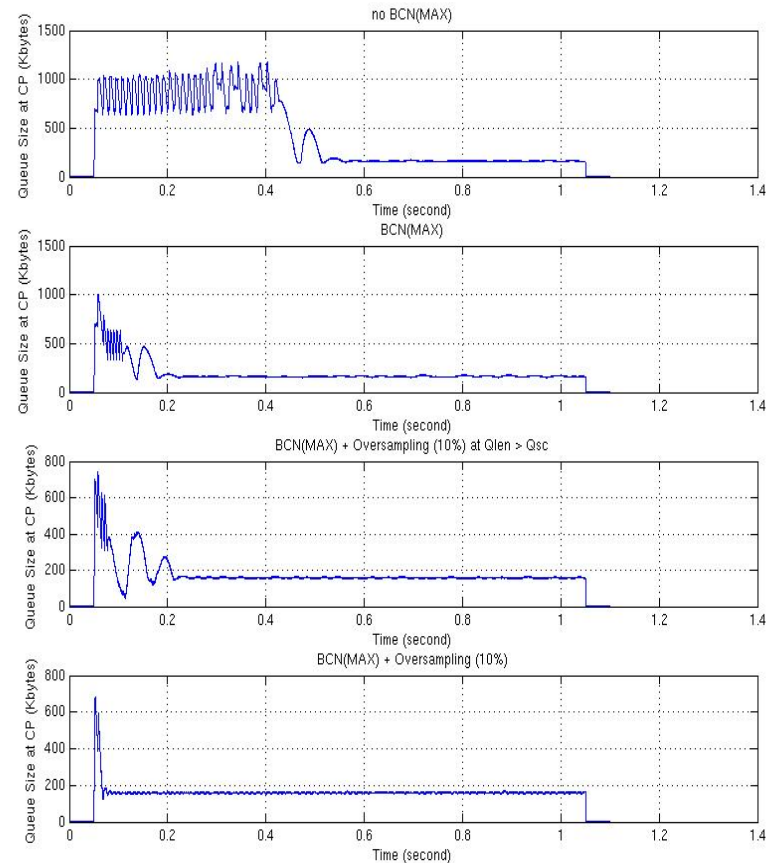


Effects of Sampling Behavior (Queue Behavior @ CP)

- Set up

- Buffer Size = 600KB/Port
- Qoffset = (-Qeq, +Qeq)
- Qdelta = (-2Qeq, +2Qeq)

Case	BCN (Max)	Sampling	Transient Duration
1	No	2%	462 ms (184 x RTT)
2	Yes	2%	52ms (or 21 x RTT)
3	Yes	2% (Qlen < Qsc) 10% (Qlen >= Qsc)	86 ms (34 x RTT)
7	Yes	10%	38 ms (15 x RTT)



Effects of Buffer Size

- Setup
 - $Q_{eq} = \text{Queue Size} * 1/4$
 - Qoffset: (- Q_{eq} , + Q_{eq})
 - Qdelta: (-2 Q_{eq} , +2 Q_{eq})

Case	BCN (Max)	Sampling	Transient Duration	
			Queue Size = 600 KBytes ($Q_{eq} = 150\text{KBytes}$)	Queue Size = 75 Kbytes ($Q_{eq} = 18.75\text{KBytes}$)
1	No	2%	462 ms (184 x RTT)	264ms (880 x RTT)
2	Yes	2%	52ms (or 21 x RTT)	56ms (187 x RTT)
3	Yes	2% ($Q_{len} < Q_{sc}$) 10% ($Q_{len} \geq Q_{sc}$)	86 ms (34 x RTT)	42ms (140 x RTT)
4	Yes	2% ($Q_{len} < Q_{sc}$) 20% ($Q_{len} \geq Q_{sc}$)	144 ms (58 x RTT)	52ms (173 x RTT)
5	Yes	2% ($Q_{len} < Q_{sc}$) 30% ($Q_{len} \geq Q_{sc}$)	96 ms (38 x RTT)	92ms (306 x RTT)
6	Yes	5%	52ms (21 x RTT)	26 ms (87 x RTT)
7	Yes	10%	38 ms (15 x RTT)	20 ms (67 x RTT)

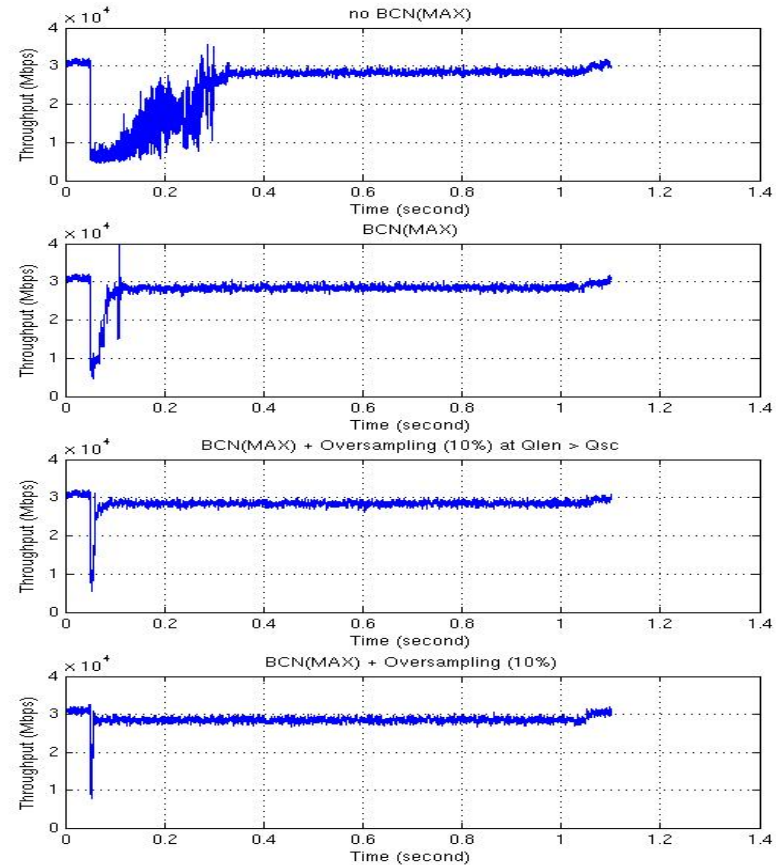
With small buffers, more challenging to minimize transient duration

Effects of Buffer Size (Aggregate Throughput)

- Set up

- Buffer Size = 75KB
- Qoffset: (-Qeq, +Qeq)
- Qdelta: (-2Qeq, +2Qeq)

Case	BCN (Max)	Sampling	Transient Duration
1	No	2%	264ms (880 x RTT)
2	Yes	2%	56ms (187 x RTT)
3	Yes	2% ($Q_{len} < Q_{sc}$) 10% ($Q_{len} \geq Q_{sc}$)	42ms (140 x RTT)
7	Yes	10%	20 ms (67 x RTT)

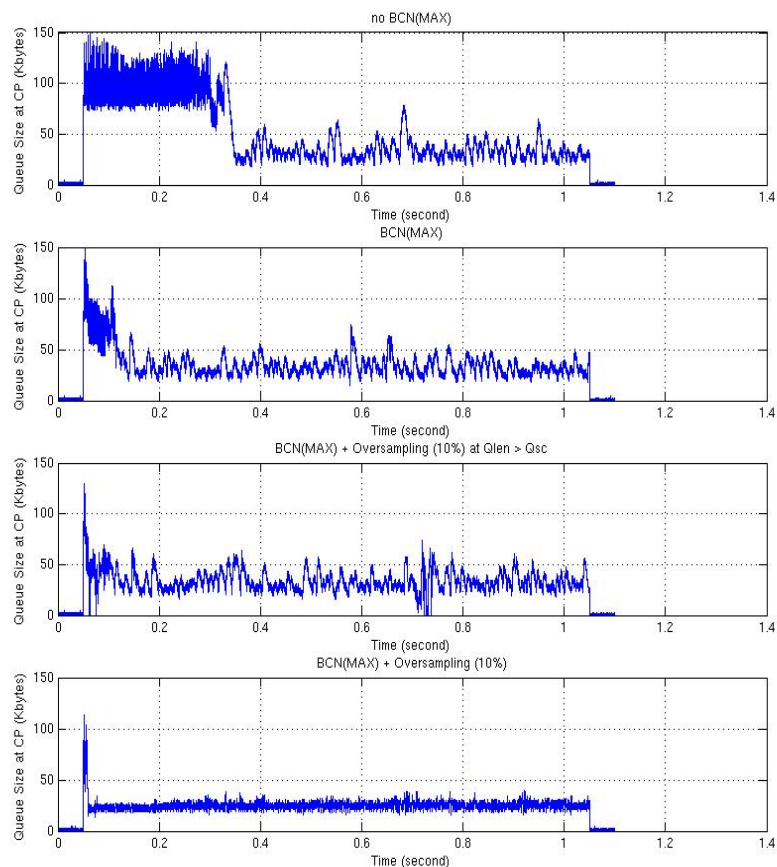


Effects of Buffer Size (Queue Behavior @ CP)

- Set up

- Buffer Size = 75KB
- Qoffset: (-Qeq, +Qeq)
- Qdelta: (-2Qeq, +2Qeq)

Case	BCN (Max)	Sampling	Transient Duration
1	No	2%	264ms (880 x RTT) <i>(Heavy Oscillations)</i>
2	Yes	2%	56ms (187 x RTT)
3	Yes	2% (Qlen < Qsc) 10% (Qlen >= Qsc)	42ms (140 x RTT)
7	Yes	10%	20 ms (67 x RTT)



Overview

- Goals
- Observations
- Simulation Setup
- Results
- Summary & Next Steps

Summary

- Without additional severe congestion schemes, the transient duration for the multi-hop output generated scenario can be on the order of 462 ms (184 x RTT).
- With the use of one of the simple severe congestion management schemes (i.e. BCN-MAX), the transient duration for the multi-hop output generated scenario can be reduced to 52 ms (21 x RTT).
- Small buffer sizes can hamper transient duration (when measured in terms of RTT). Requires additional study.

Next Steps

- Refine metric to quantify the response time of the CN mechanism.
- Collect input from the group on identifying a target response time for the traffic of interest.
- Quantify benefits of oversampling when $Q_{len} = 0$. Measure actual amount of BCN messaging control bandwidth.
- Additional study required of how BCN operates under small buffer conditions at the switch.