

Master Mini Project

Performance Evaluation of Minimum Delivery Delay Multipath TCP Scheduler Implemented in Linux Kernel

Faridul Alam

Date: 21 June, 2016

Supervisors

Amanpreet Singh, M.Sc

Dr. Andreas Könsgen

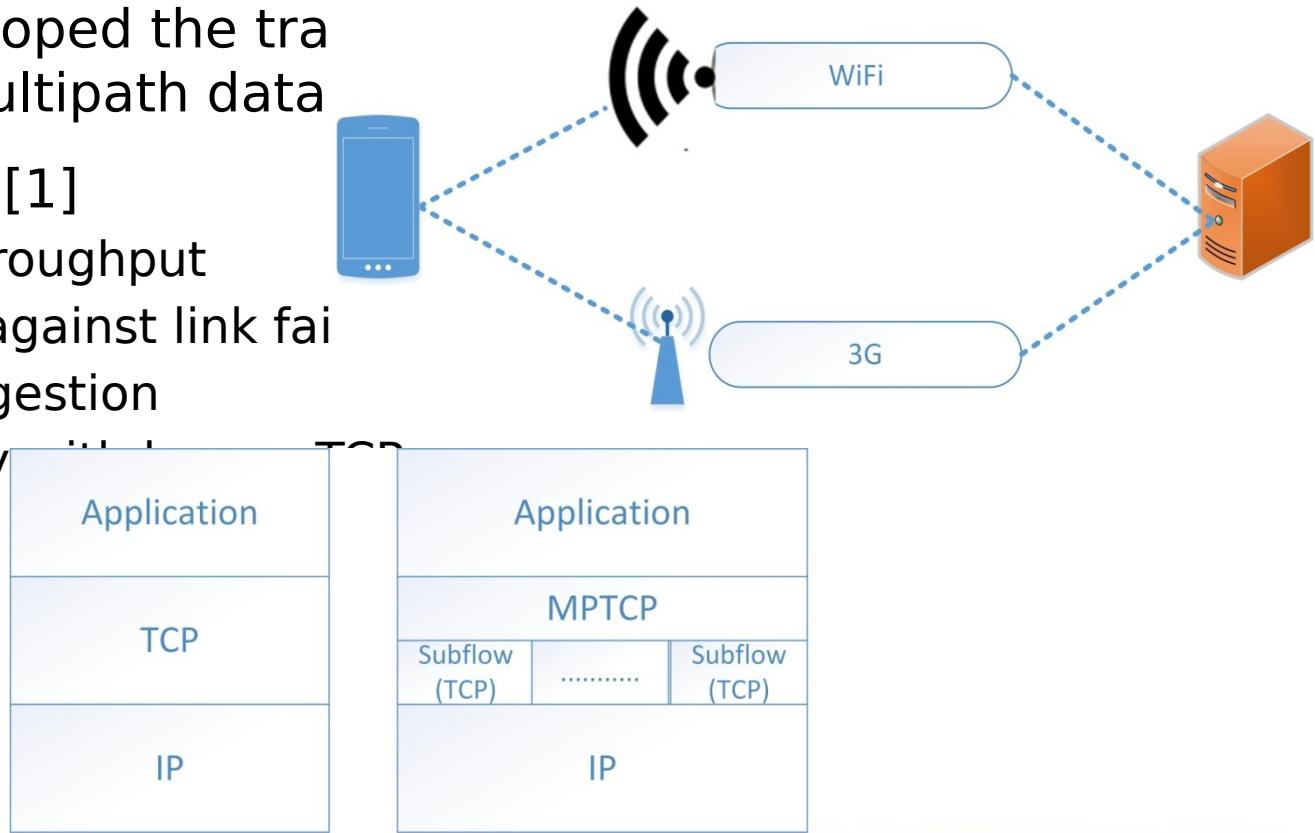
Prof. Dr. Anna Förster

Outline

- Multipath TCP (MPTCP)
- MPTCP Scheduler
 - Existing MPTCP Schedulers on Linux
 - Round Robin Scheduler
 - Shortest Round Trip Time Scheduler
- Implementation of Delivery Delay Scheduler
- Performance Evaluation
- Conclusion & Outlook

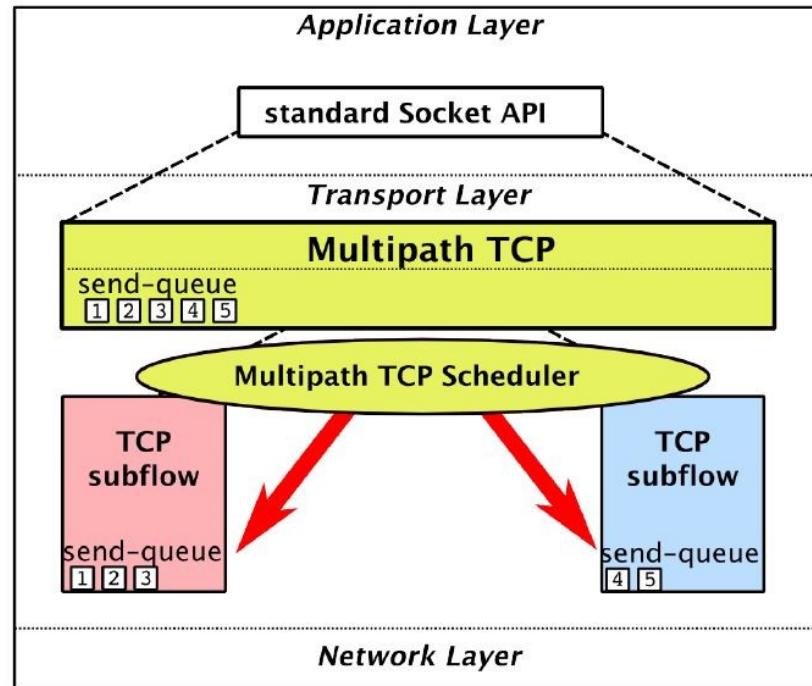
Multipath TCP (MPTCP)

- With the emergence of multihoming capable devices, data transfer through multiple paths for a session become a viable option
- IETF has developed the tra to facilitate multipath data
- MPTCP offers- [1]
 - Increased throughput
 - Robustness against link fai
 - Balance congestion
 - Compatibility



MPTCP Scheduler

- Heterogeneous nature of delay and bandwidth asymmetry of the path can create out-of-order data arrival at the receiver [2]
 - Head-of-Line Blocking
 - Receiver-Window Limitations
- Results higher reordering delay and might give lower throughput



[2]

Round Robin (RR) Scheduler

- Schedule data packets in round robin or cyclic fashion
- Maximum number of consecutive scheduled segments for an instance is configurable
 - Fill up the congestion window of a subflow or
 - Strictly Round Robin
- Disadvantages –
 - May not be able to utilize the whole path capacity when paths have significant delay or capacity differences as the scheduler always tries to distribute equally among all available subflows [5]

Shortest Round Trip Time (sRTT) Scheduler

- Schedule data packets on a subflow with the lowest RTT amongst all available subflows
- Considers delay (RTT) of the paths, which also reflects network condition to some extent
- Disadvantages -
 - High delay assymetry in backward path might affect the scheduling which might influence out-of-order delivery

Delivery Delay (DD) based Scheduler

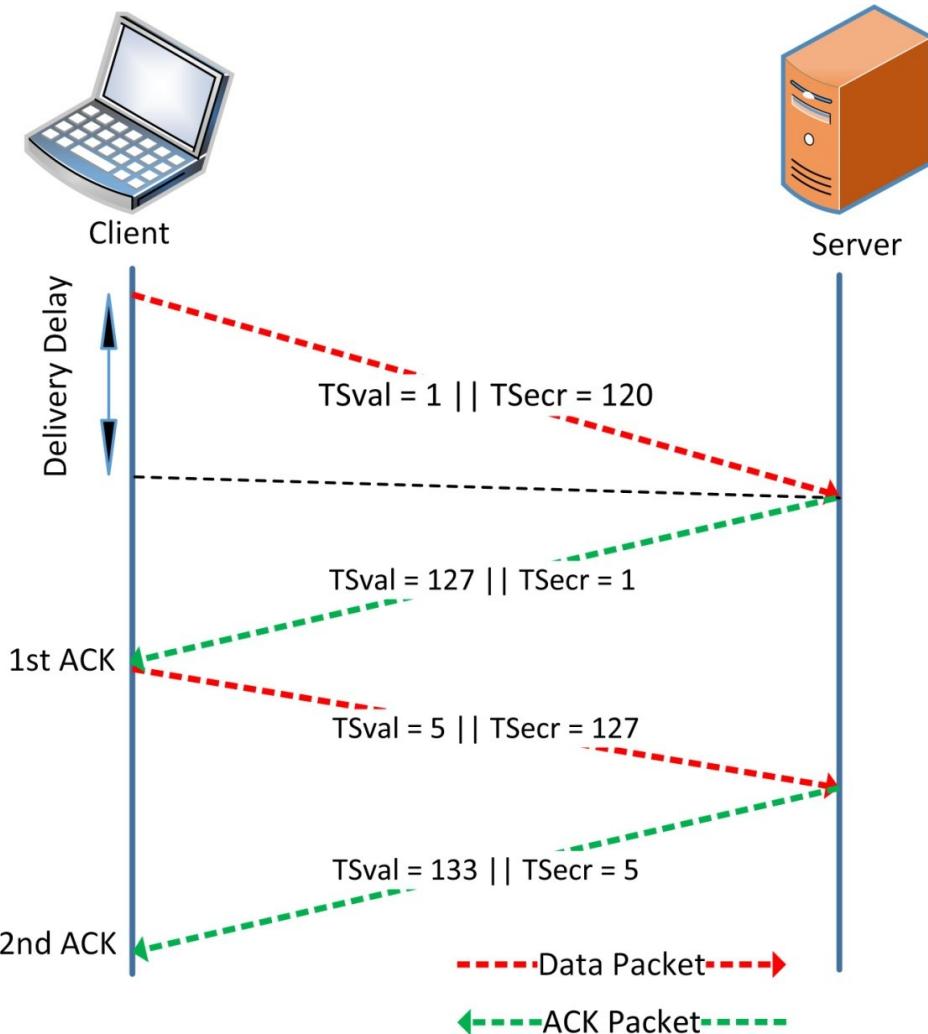
- Schedule data packet on a subflow with the lowest Delivery Delay (DD) amongst all available subflows
- DD scheduler was first proposed in [8] with the title ``Performance Comparison of Scheduling Algorithms for Multipath Transfer''
 - Was implemented in a Simulator and shown that the End-to-End (ENE) delay is reduced
- Focus point –
 - Understand the implemented schedulers (sRTT and RR) in Linux Kernel along with their behaviour
 - Implement the Delivery Delay (DD) based scheduler in Linux Kernel
 - Performance analysis of DD scheduler along with the other two implemented schedulers for different network scenarios

Estimation of Delivery Delay

TCP Timestamps [9]			
Kind = 8	Length	TSval	TSecr
1 Byte	1 Byte	4 Bytes	4 Bytes

- TSval is the current value of the timestamp's clock
- TSecr is the TSval value contained in the received data/ack packet

- Delivery Delay - **1st ACK** = $\text{TSval} - \text{TSecr} = 127 - 1 = 126$
- Delivery Delay - **2st ACK** = $\text{TSval} - \text{TSecr} = 133 - 5 = 128$



Scheduling Algorithm - For each packet

min_delivery_delay = MAX_INT

for each subflow i **do**

if (current packet is a retransmission) and (has been transmitted on the subflow i) **then** **continue**

end if

if (subflow i is **definitely** unavailable) or (subflow i is **temporary** unavailable) **then** **continue**

end if

if (delivery_delay_i < min_delivery_delay) **then**

 min_delivery_delay = delivery_delay_i

selected_subflow = i

end if

end for

Performance Evaluation

Test Topology - 2 Node

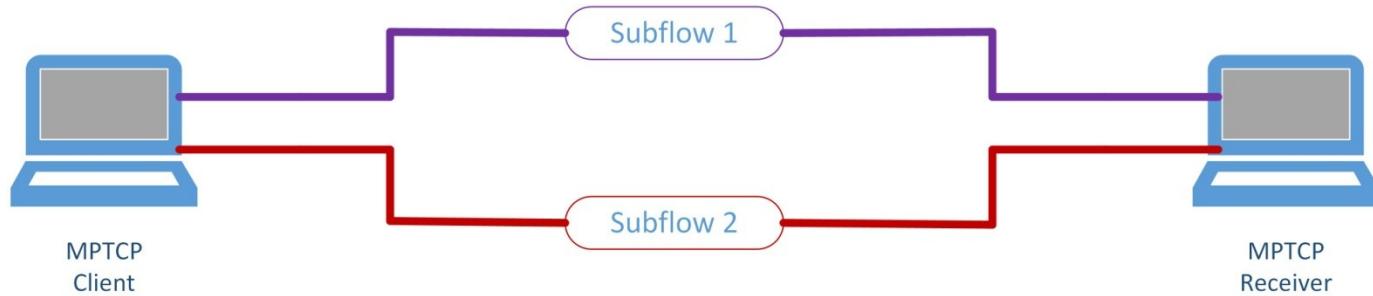


Figure : Two Node Testbed

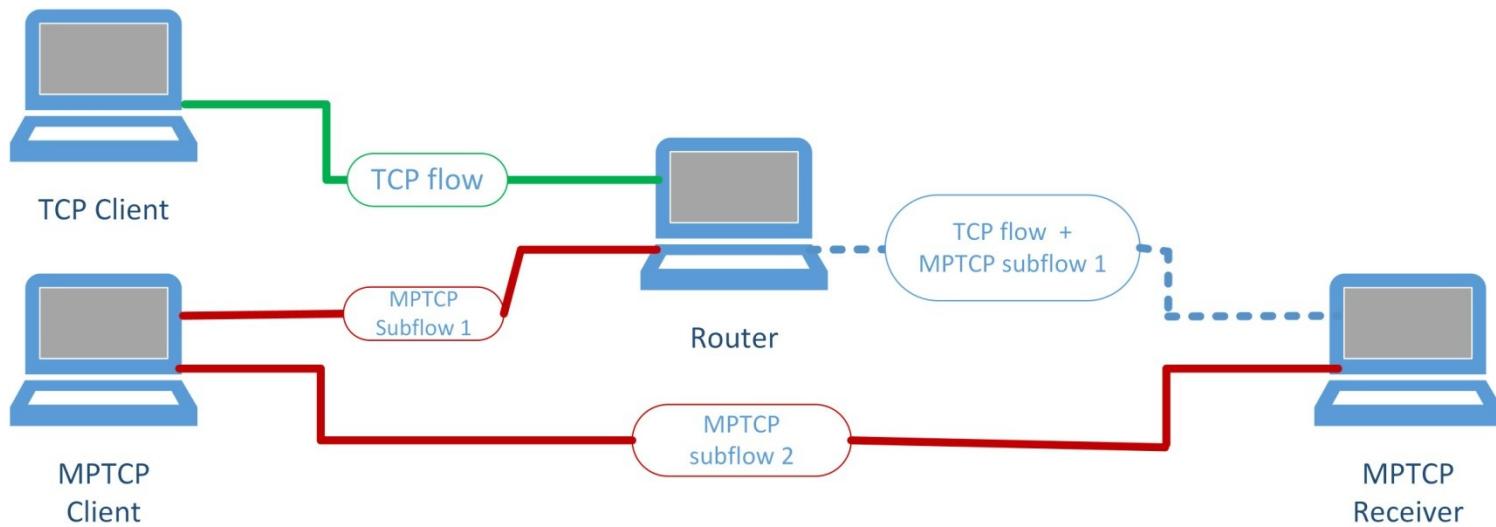


Figure : Extended Testbed

Test Topology - Two Node Testbed with Router

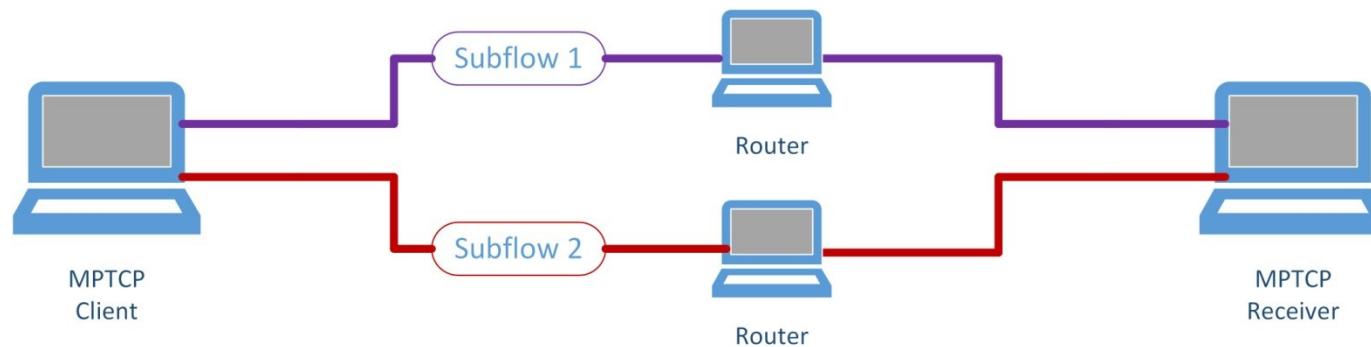


Figure : Two Node Testbed with Router

- Configuration :
 - Congestion Control = Reno
 - `tc qdisc "add" dev "ethX" root netem delay "Yms"` – at router

Validation of Delivery Delay Scheduler - I

```
647.097468] MPTCPcon:mptcp_sch      get_available_subflows_x
647.097469] MPTCPcon:mptcp_sch      subflow_selector_enter
647.097471] MPTCPcon:mptcp_sch      subflow_selector_subflow_enter
647.097472] MPTCPcon:mptcp_sch      subflow_selector_definitely_unavailable
647.097474] MPTCPcon:mptcp_sch      subflow_selector_subflow_enter
647.097476] MPTCPcon:mptcp_sch      subflow_selector_temporary_unavailable
647.097477] MPTCPcon:mptcp_sch      subflow_selector_subflow_enter
647.097478] MPTCPcon:mptcp_sch      subflow_selector_definitely_unavailable
647.097480] MPTCPcon:mptcp_sch      subflow_selector_subflow_enter
647.097484] MPTCPcon:mptcp_sch      subflow_selector_temporary_unavailable
647.097489] MPTCPcon:mptcp_sch      subflow_selector_selected_subflow
647.097494] MPTCPcon:mptcp_sch      subflow_selector_exit
647.097498] MPTCPcon:mptcp_sch      subflow_select_active (null)
647.097504] TCPcong: tcp_input      tcp_rtt_estm
647.097509] TCPcong: tcp_input      d-delay-es tval tsecr
647.097514] TCPcong: tcp_input      smooth fw_delay
647.097519] TCPcong: tcp_input      srtt_f-es tstm tsecr
647.097524] TCPcong: tcp_input      d-delay-es_x srtt_f_c
647.097529] TCPcong: tcp_input      smooth srtt_f
647.097536] TCPcong: tcp_input      srtt_f-es tstm tsval
647.097542] TCPcong: tcp_input      d-delay-es_x bw_delay
647.097547] TCPcong: tcp_input      smooth bw_delay
647.097550] TCPcong: tcp_cong     s_start cwnd
647.097552] TCPcong: tcp_cong     s_start ssthresh
647.097553] TCPcong: tcp_cong     s_start in_flight
647.097555] MPTCPcon:mptcp_sch   get_available_subflows_x
647.097556] MPTCPcon:mptcp_sch   subflow_selector_enter
647.097558] MPTCPcon:mptcp_sch   subflow_selector_subflow_enter
647.097559] MPTCPcon:mptcp_sch   subflow_selector_definitely_unavailable
647.097561] MPTCPcon:mptcp_sch   subflow_selector_subflow_enter
647.097563] MPTCPcon:mptcp_sch   subflow_selector_checking_algorithm
647.097564] MPTCPcon:mptcp_sch   subflow_selector_subflow_exit
647.097565] MPTCPcon:mptcp_sch   subflow_selector_subflow_enter
647.097567] MPTCPcon:mptcp_sch   subflow_selector_definitely_unavailable
647.097568] MPTCPcon:mptcp_sch   subflow_selector_subflow_enter
647.097570] MPTCPcon:mptcp_sch   subflow_selector_temporary_unavailable
647.097572] MPTCPcon:mptcp_sch   subflow_selector_selected_subflow
647.097573] MPTCPcon:mptcp_sch   subflow_selector_exit
647.097574] MPTCPcon:mptcp_sch   subflow_select_active f2351080
```

Both subflows are temporary unavailable

One ACK received on subflow 2, creates free space in cwnd

Subflow 2 is available now

Subflow 2 is selected

Validation of Delivery Delay Scheduler - II

```
647.150192] MPTCPcon:mptcp_sch      get_available_subflows_x
647.150196] MPTCPcon:mptcp_sch      subflow_selector_enter
647.150201] MPTCPcon:mptcp_sch      subflow_selector  definitely_unavailable  0 f2351600
647.150206] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          340200 f2351080
647.150212] MPTCPcon:mptcp_sch      subflow_selector  checking_algorithm    340200 f2351080
647.150220] MPTCPcon:mptcp_sch      subflow_selector  subflow_exit           340200 f2351080
647.150225] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          0 f2350b00
647.150230] MPTCPcon:mptcp_sch      subflow_selector  definitely_unavailable 0 f2350b00
647.150235] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          340199 f3f02c00
647.150240] MPTCPcon:mptcp_sch      subflow_selector  checking_algorithm    340199 f3f02c00
647.150245] MPTCPcon:mptcp_sch      subflow_selector  subflow_exit           340199 f3f02c00
647.150250] MPTCPcon:mptcp_sch      subflow_selector  selected_subflow     f3f02c00
647.150255] MPTCPcon:mptcp_sch      subflow_selector_exit  f3f02c00
647.150260] MPTCPcon:mptcp_sch      subflow_select_active f3f02c00
647.150265] MPTCPcon:mptcp_sch      sel_win        cwnd       506          f3f02c00
647.150271] TCPcong: tcp_out        sel_win        ssthresh   2147483647 f3f02c00
647.150276] TCPcong: tcp_out        sel_win        in_flight  128          f3f02c00
647.150280] TCPcong: tcp_out
647.150287] MPTCPcon:mptcp_sch      get_available_subflows_x
647.150288] MPTCPcon:mptcp_sch      subflow_selector_enter
647.150290] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          0 f2351600
647.150292] MPTCPcon:mptcp_sch      subflow_selector  definitely_unavailable 0 f2351600
647.150293] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          340200 f2351080
647.150295] MPTCPcon:mptcp_sch      subflow_selector  checking_algorithm    340200 f2351080
647.150297] MPTCPcon:mptcp_sch      subflow_selector  subflow_exit           340200 f2351080
647.150298] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          0 f2350b00
647.150300] MPTCPcon:mptcp_sch      subflow_selector  definitely_unavailable 0 f2350b00
647.150302] MPTCPcon:mptcp_sch      subflow_selector  subflow_enter          340199 f3f02c00
647.150303] MPTCPcon:mptcp_sch      subflow_selector  checking_algorithm    340199 f3f02c00
647.150305] MPTCPcon:mptcp_sch      subflow_selector  subflow_exit           340199 f3f02c00
647.150306] MPTCPcon:mptcp_sch      subflow_selector  selected_subflow     f3f02c00
647.150307] MPTCPcon:mptcp_sch      subflow_selector_exit  f3f02c00
647.150309] MPTCPcon:mptcp_sch      subflow_select_active f3f02c00
647.150312] TCPcong: tcp_out        sel_win        cwnd       506          f3f02c00
647.150313] TCPcong: tcp_out        sel_win        ssthresh   2147483647 f3f02c00
647.150315] TCPcong: tcp_out        sel_win        in_flight  129          f3f02c00
647.150320] MPTCPcon:mptcp_sch      sel_win        cwnd       506          f3f02c00
647.150321] MPTCPcon:mptcp_sch      sel_win        ssthresh   2147483647 f3f02c00
647.150322] MPTCPcon:mptcp_sch      sel_win        in_flight  129          f3f02c00
```

Both subflows are available

Subflow 1's delivery delay is lower than that of subflow 2, therefore subflow 1 is selected

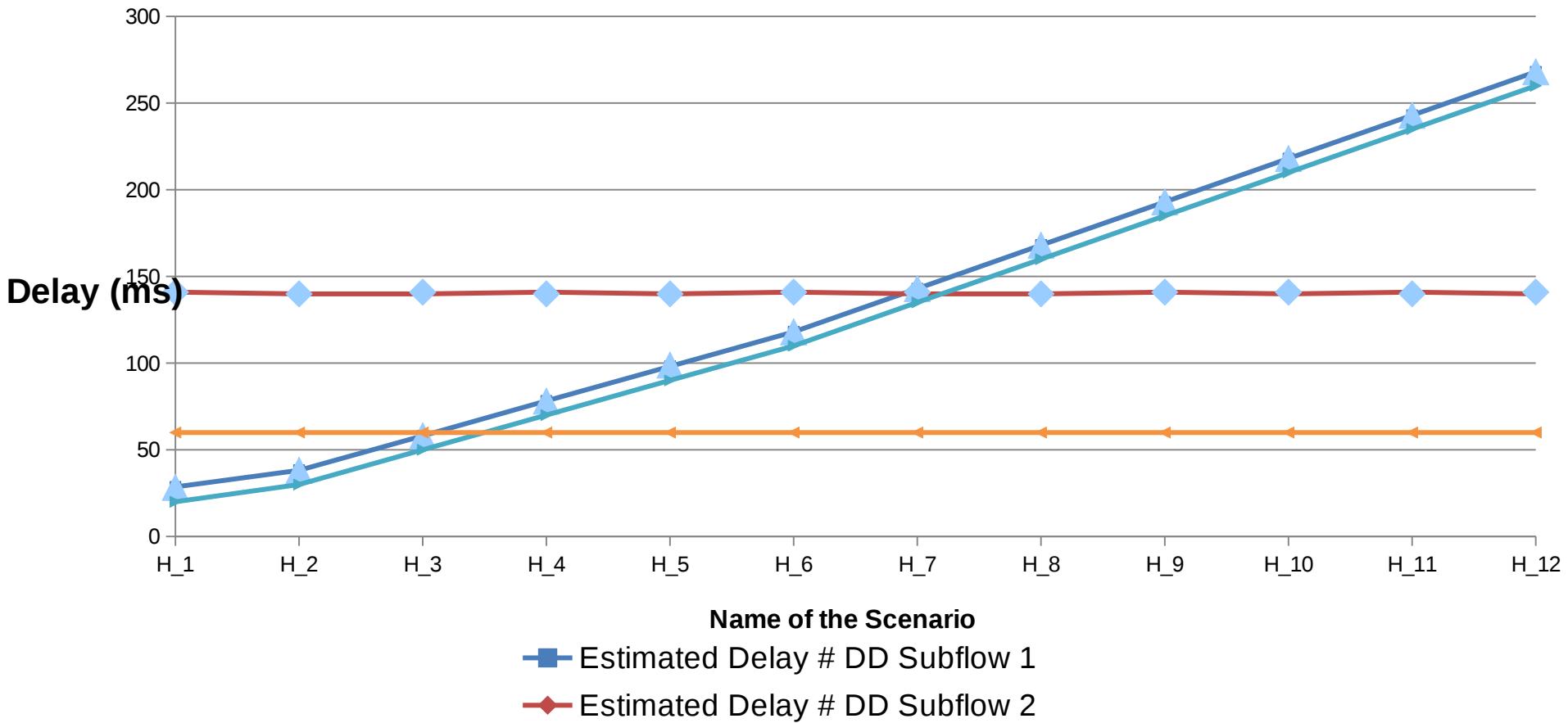
For the next packet, subflow 1 is also selected

Performance Evaluation - Two Node Testbed with Router

Name of the Scenario	Subflow 1		Subflow 2	
	RTT (ms)	FPD+BPD (ms)	RTT (ms)	FPD+BPD (ms)
H_1	20	10 + 10	60	50 + 10
H_2	30	10 + 20	60	50 + 10
H_3	50	10 + 40	60	50 + 10
H_4	70	10 + 60	60	50 + 10
H_5	90	10 + 80	60	50 + 10
H_6	110	10 + 100	60	50 + 10
H_7	135	10 + 125	60	50 + 10
H_8	160	10 + 150	60	50 + 10
H_9	185	10 + 175	60	50 + 10
H_10	210	10 + 200	60	50 + 10
H_11	235	10 + 225	60	50 + 10
H_12	260	10 + 250	60	50 + 10

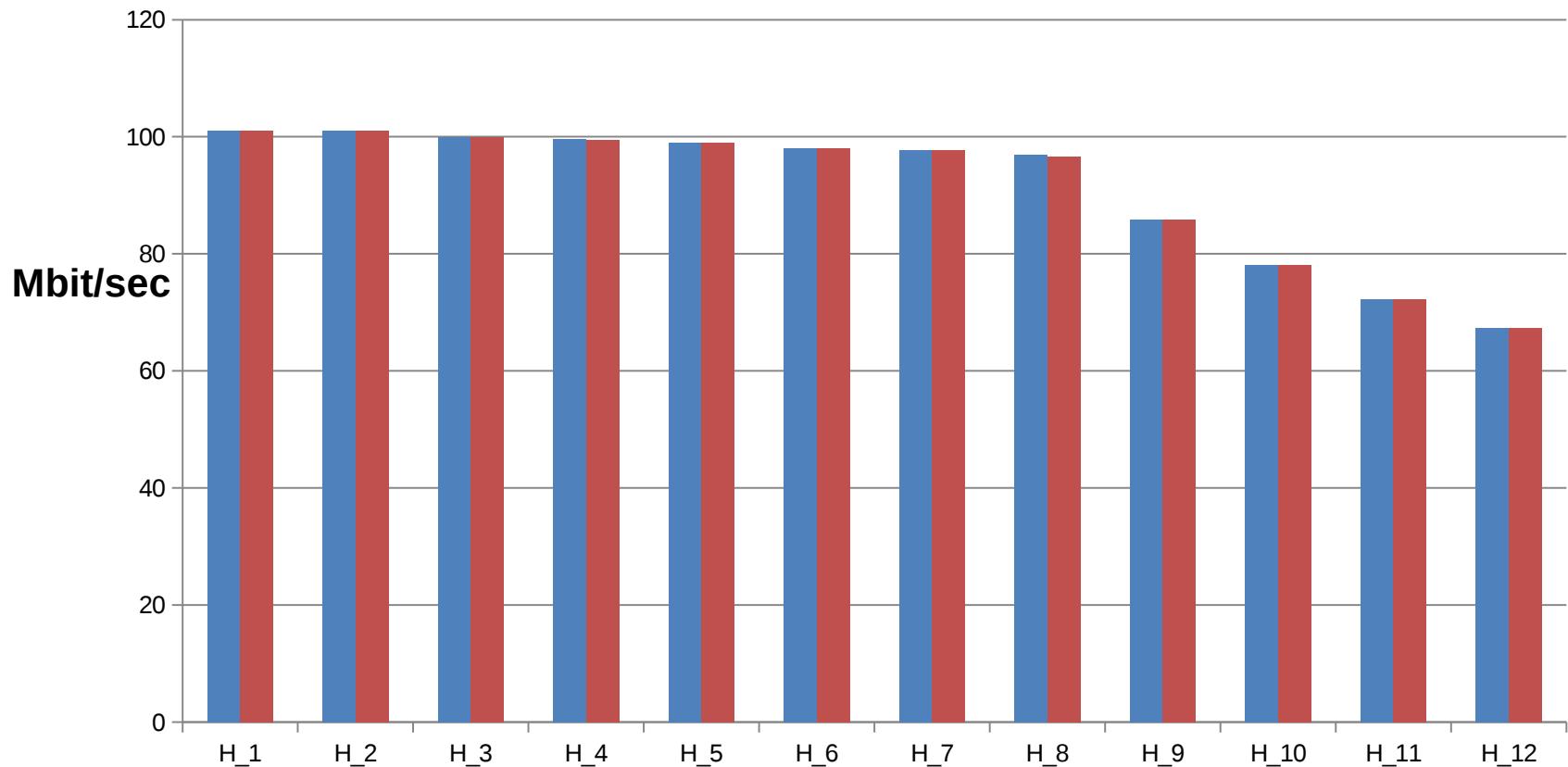
Link Capacity - 100/10 (Mbit/sec)

Measured Delay Vs Configured Delay



Link Capacity 100/10 (Mbit/sec)

Throughput Comparison Table



Graphical Analysis

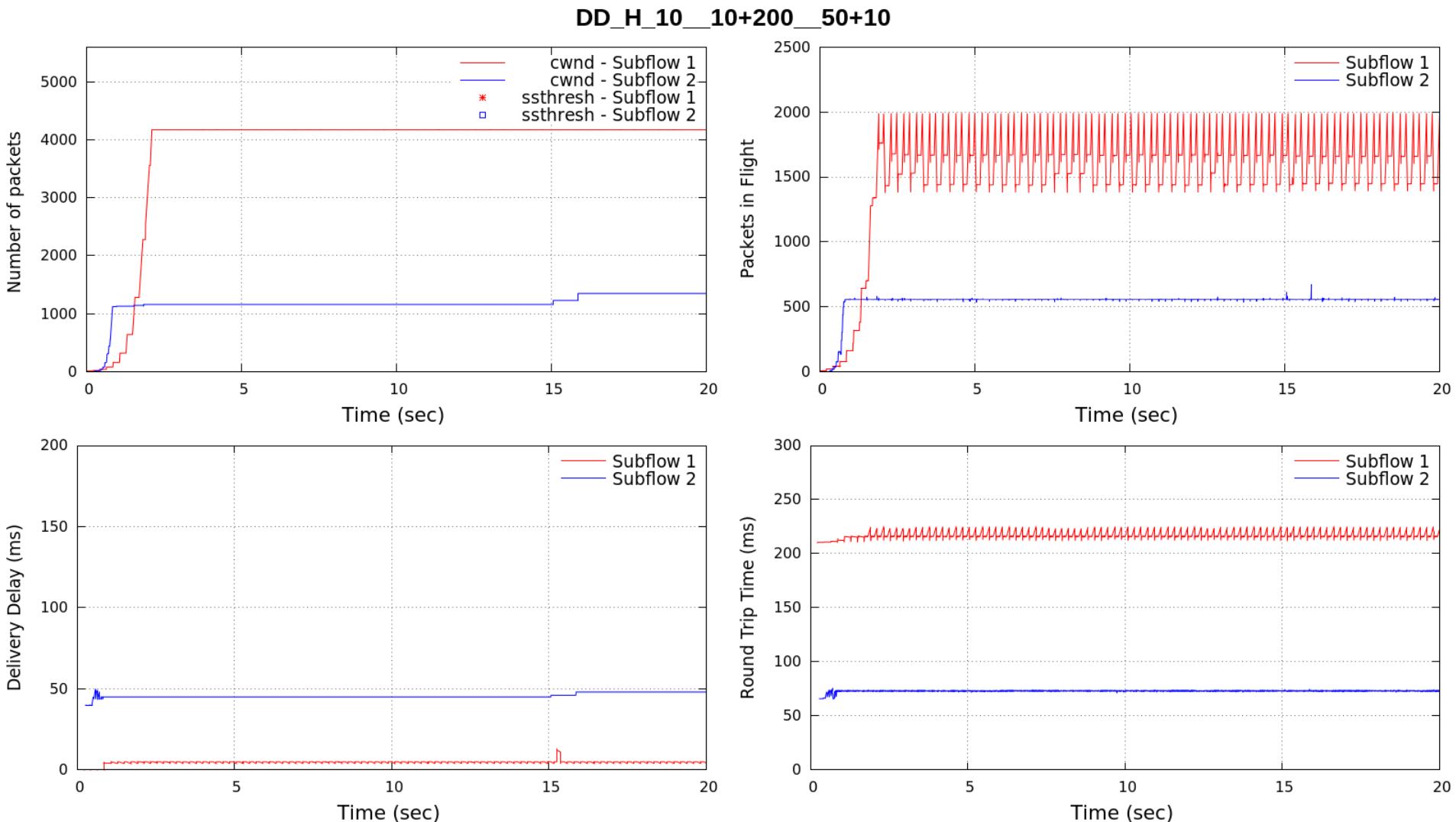
Link Capacity 100/100 (Mbit/sec)

Scenario H_10

Subflow 1 RTT - 210 ms (10 ms + 200 ms)

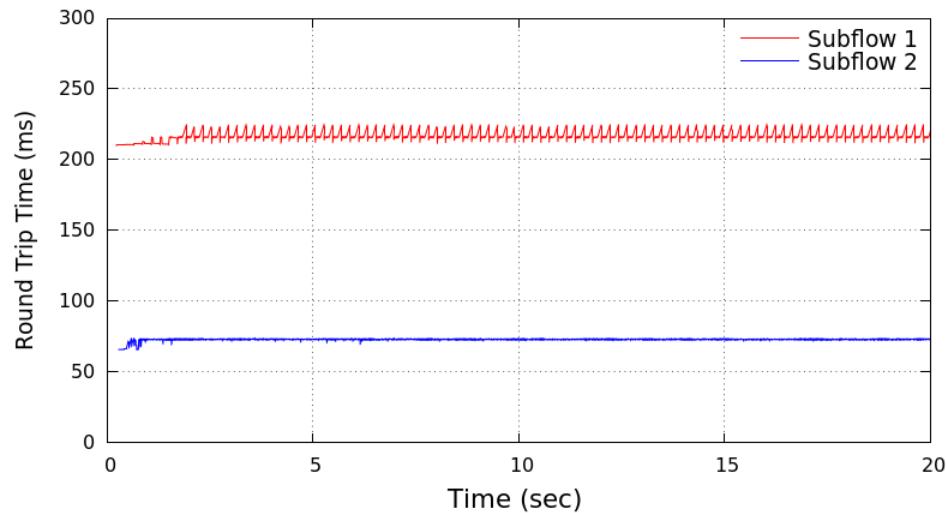
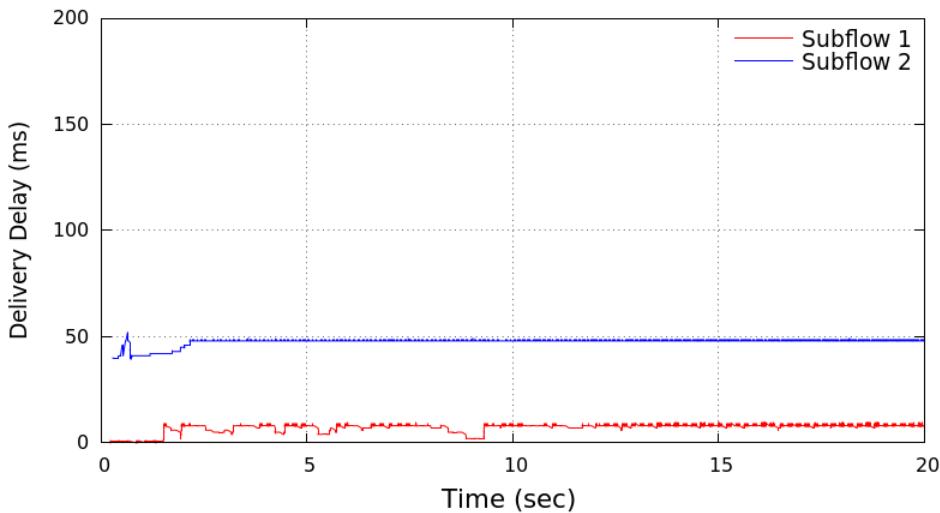
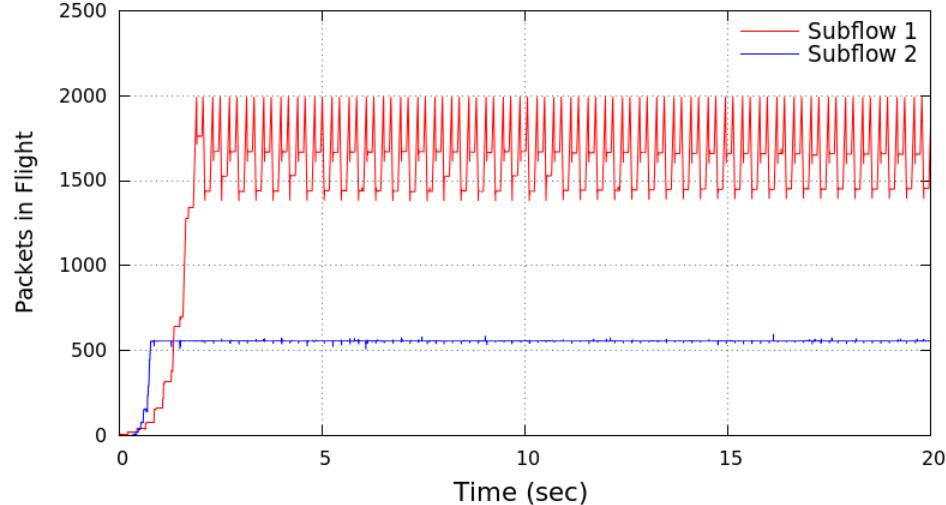
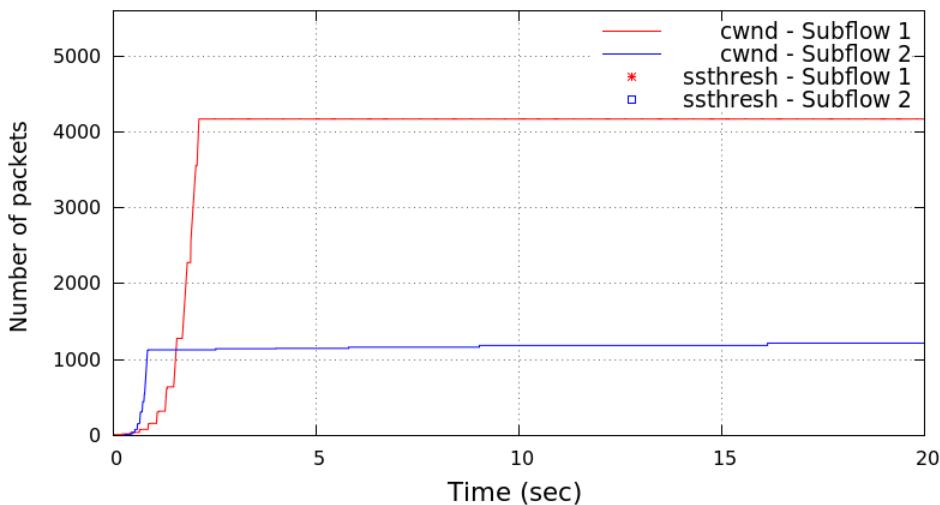
Subflow 2 RTT - 60 ms (50 ms + 10 ms)

Delivery Delay Scheduler

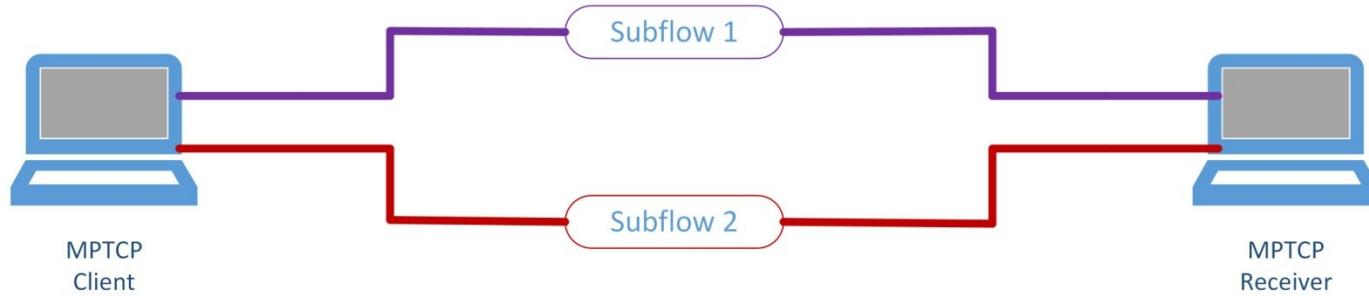


Shortest Round Trip Time Scheduler

sRTT_H_10_10+200_50+10



Graphical Analysis- Two Node Testbed



Configuration

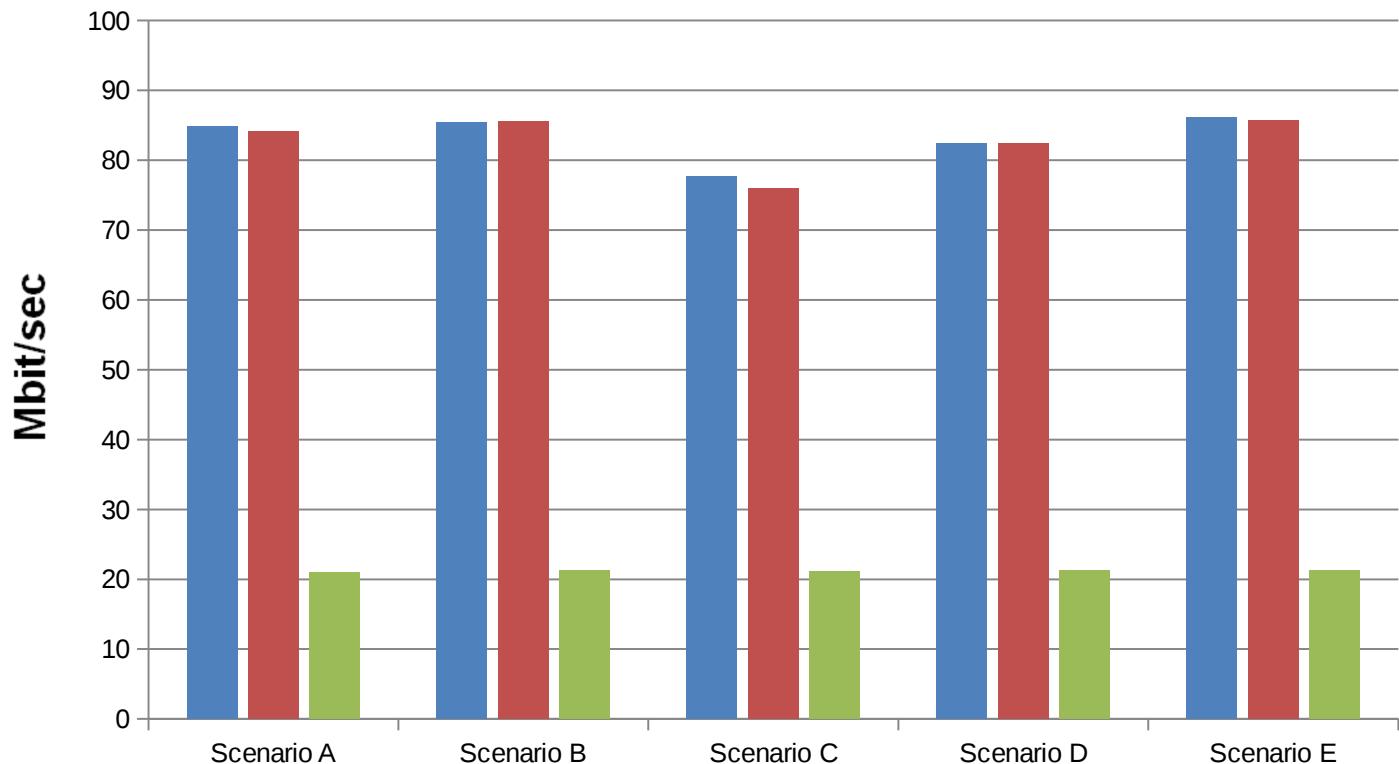
- Congestion Control = Reno
- `tc qdisc "add" dev "ethX" root netem delay "Yms"` – at Client and Server

Scenarios

Capacity (Mbit/sec)	Name of Scenario	Subflow 1			Subflow 2		
		FPD (ms)	BPD (ms)	RTT (ms)	FPD (ms)	BPD (ms)	RTT (ms)
100/100 & 100/10 & 10/10	A	10	20	30	20	10	30
	B	10	10	20	20	10	30
	C	10	40	50	20	20	40
	D	10	20	30	20	40	60
	E	10	10	20	20	20	40
	F	10	200	210	20	400	420

Link Capacity 100/10 (Mbit/sec)

Throughput Comparison Table



Graphical Analysis

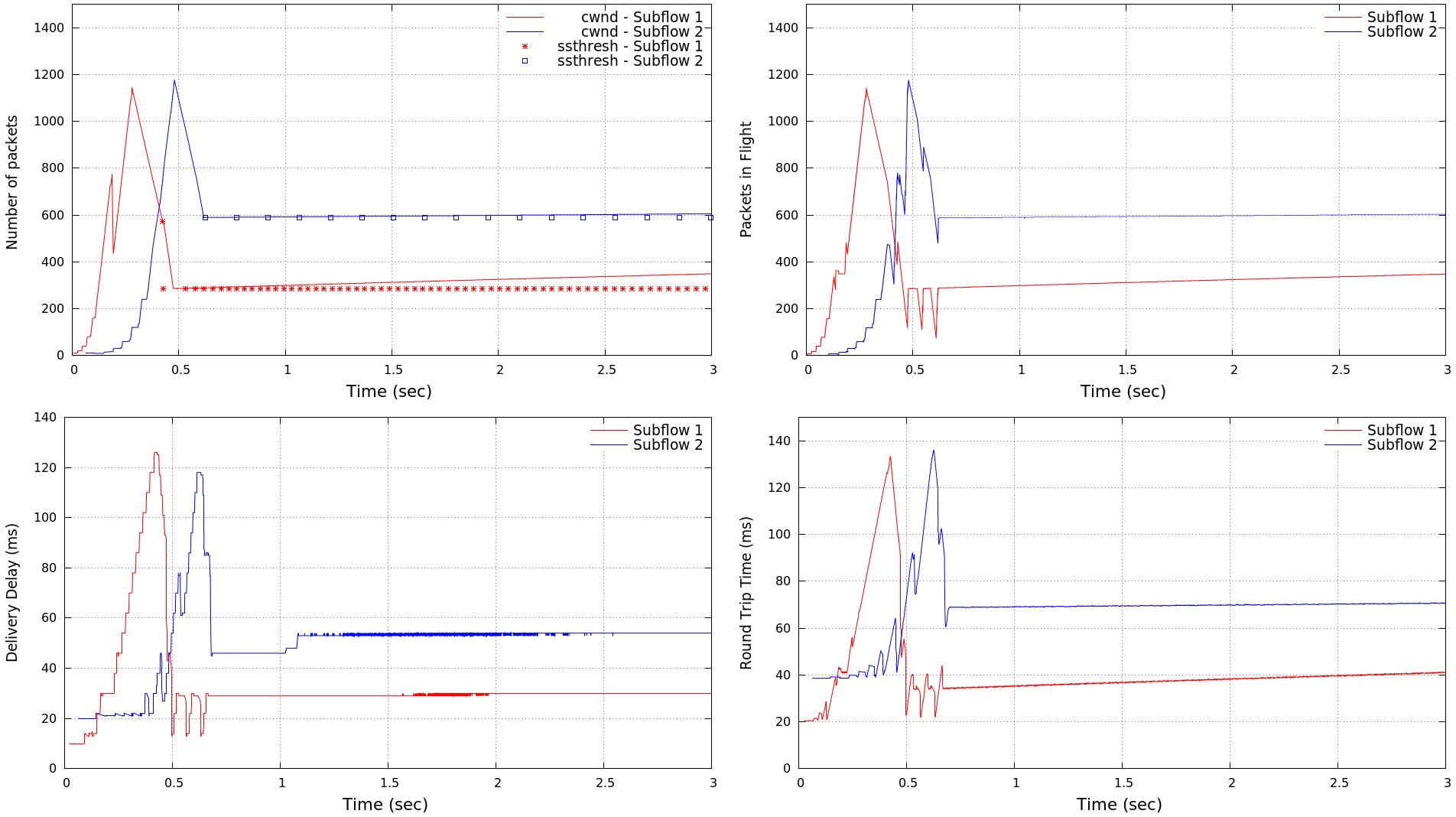
Link Capacity 100/100 (Mbit/sec)

Scenario E

Subflow 1 RTT - 20 ms (10 ms + 20 ms)

Subflow 2 RTT - 40 ms (20 ms + 20 ms)

Delivery Delay Based Scheduler



Receiver Buffer Occupancy

Scenario			
Link Capacity = 1.2 Mbit /sec	Transmission Delay	RTT (ms)	FPD+BPD (ms)
Subflow 1	10 ms	110	10 + 100
Subflow 2	10 ms	60	50 + 10

Receiver Buffer (sRTT Scheduler)

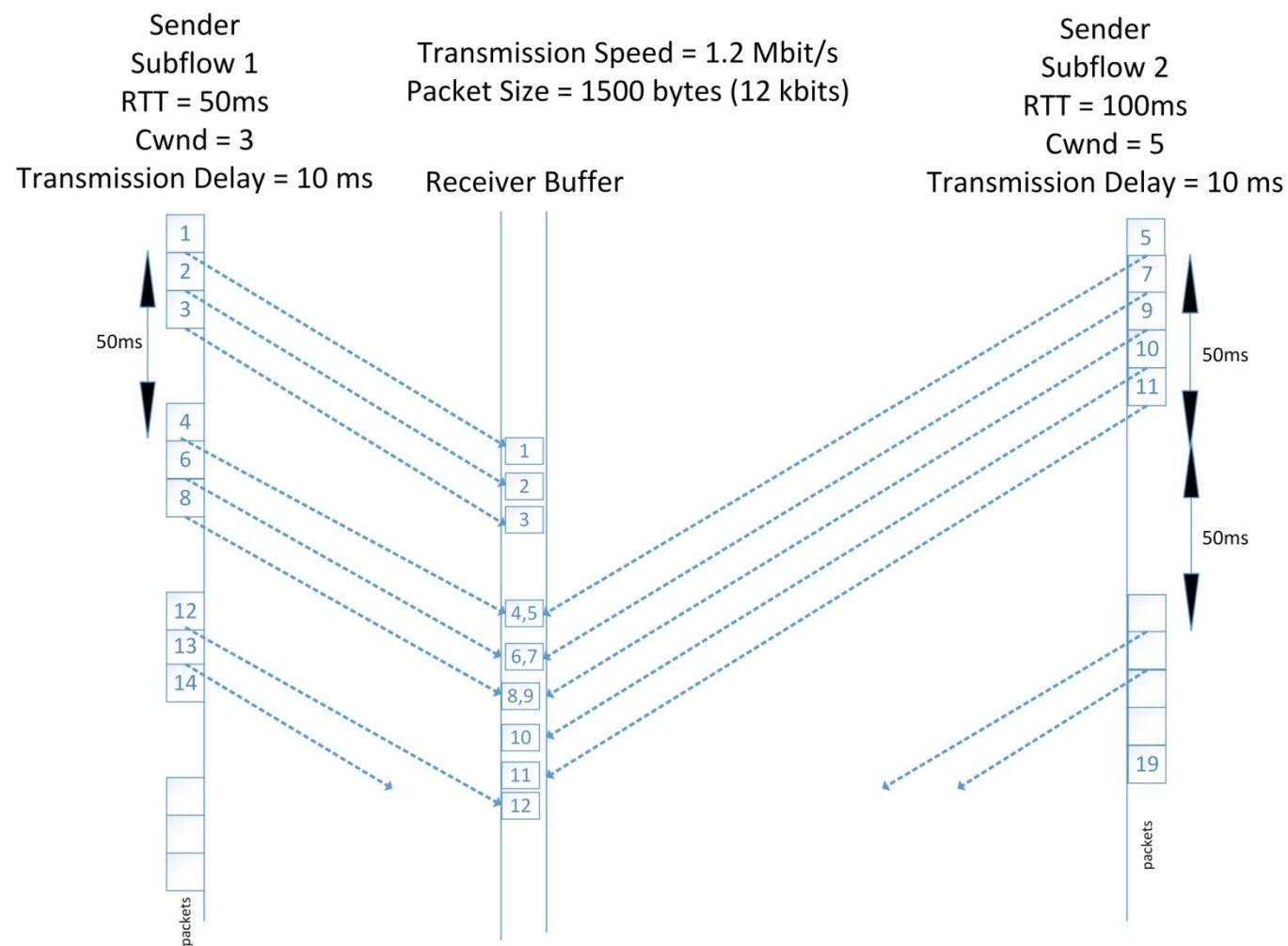
4	5	6	1	2	3	7	10	8	9	12	14	11	13	15	19
20	16	21	17	18	22	23	24	28	29	30	25	26	27	31	32
33	36	38	39	34	35	37	40	41	44	42	46	48	43	45	47

Receiver Buffer (DD Scheduler)

1	2	3	4	5	6	7	10	8	11	9	13	12	14	15	19
16	20	17	21	18	22	23	24	28	29	25	30	26	27	31	32
33	36	38	39	34	35	37	40	41	42	44	46	48	43	45	47

Outlook (Intelligent DD Scheduler)

Packet no	Reordering Delay (ms)	End-to-end Delay (ms)
1	0	50
2	0	50
3	0	50
4	0	50
5	0	100
6	0	50
7	0	100
8	0	50
9	0	100
10	0	100
11	0	100
12	0	50



Conclusion

- The performance of Round Robin scheduler is poor compared to Delivery Delay and sRTT scheduler in case of high link capacity and delay asymmetry
- Throughput performance of Delivery Delay scheduler and sRTT scheduler are similar (ack-coded)
 - Scheduling is considered only among the subflows which have an open congestion window
 - Congestion window's increase is dependent on the ACK packets and hence also controls the scheduler's performance
 - When the backward delay (delivery delay for the ACK packet) on the subflows remains constant throughout the simulation run, the sRTT scheduler almost acts as the delivery delay scheduler
- Artificial delay configuration at the sender nodes created performance issues in the experimental setup
 - Delay should be configured on the routers

Outlook

- The performance comparison of different schedulers can be evaluated in different network conditions
 - In wireless networks
 - With bidirectional data traffic
- Intelligent Delivery Delay scheduler can be implemented to minimize the out-of-order delivery
 - Advance packet scheduling by estimating the forward paths behavior

References

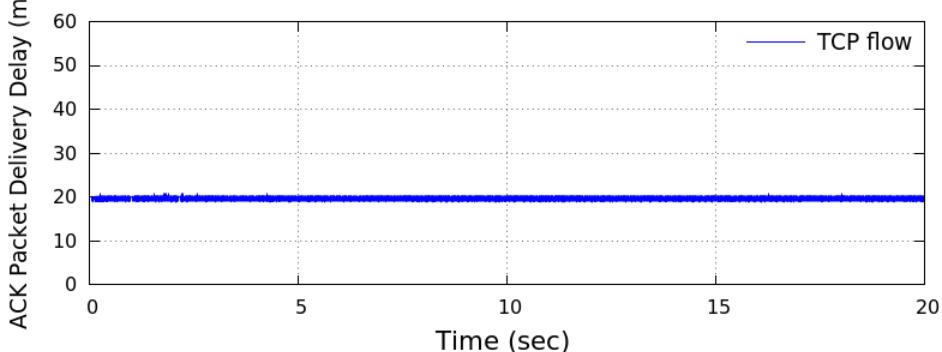
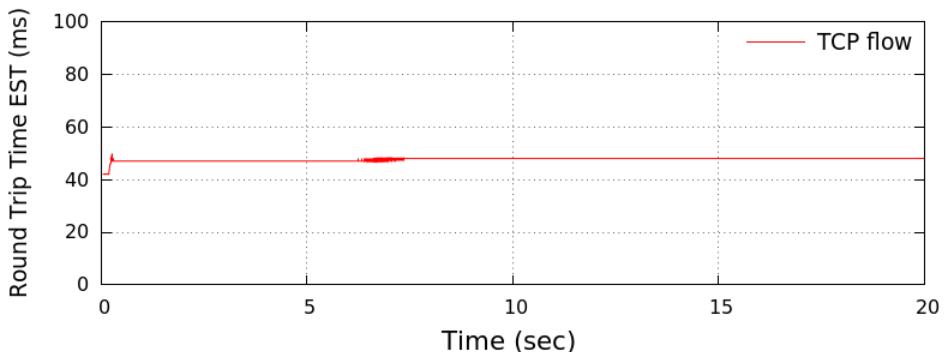
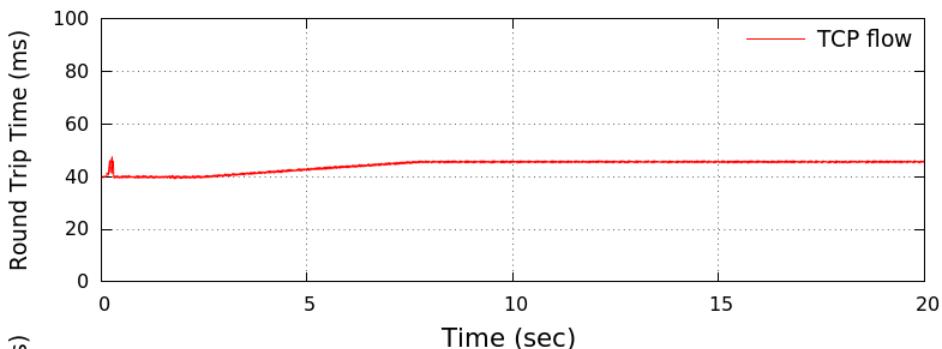
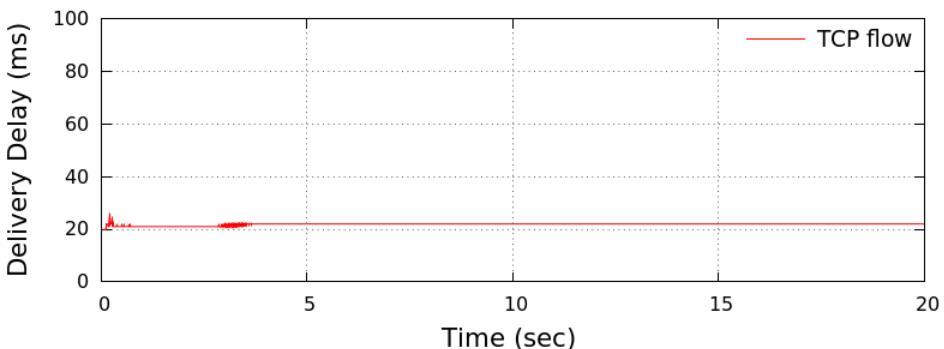
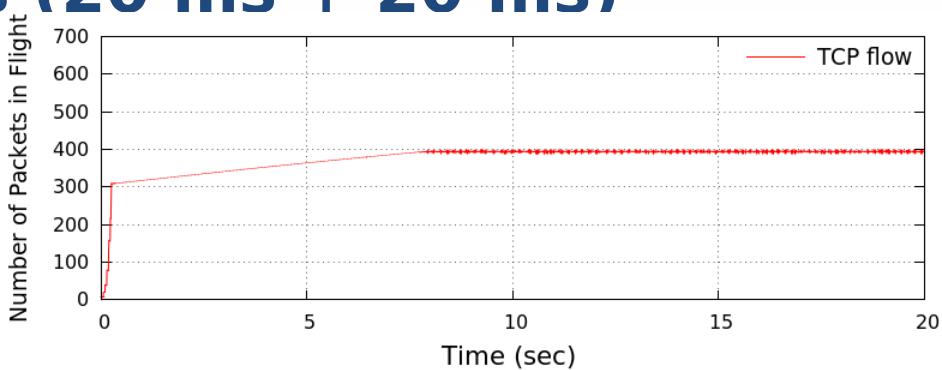
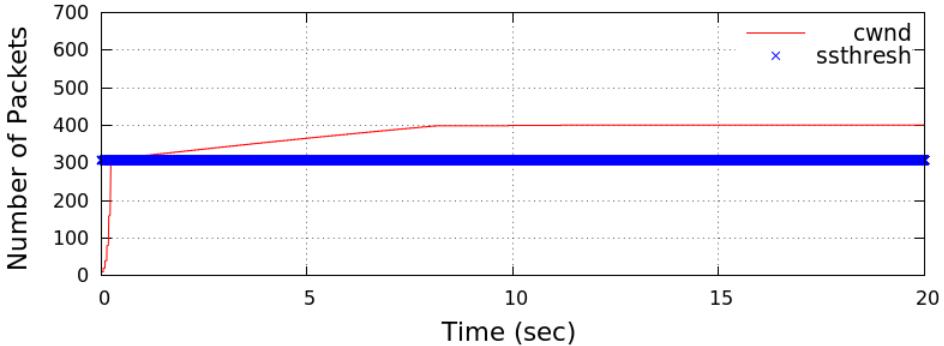
1. Costin Raiciu, Mark Handley, and Damon Wischik. "Coupled congestion control for multipath transport protocols ". RFC6356, October, 2011.
2. Christoph Paasch, Simone Ferlin, Ozgu Alay, and Olivier Bonaventure. "Experimental evaluation of multipath tcp schedulers". In Proceedings of the 2014 ACM SIGCOMM workshop on Capacity sharing workshop, pages 27-32. ACM, 2014.
3. Tuan-Anh Le and Loc X Bui. "Forward delay-based packet scheduling algorithm for multipath tcp". arXiv preprint arXiv:1501.03196, 2015.
4. Multipath TCP, Multipath TCP in the stack. <http://cacm.acm.org/magazines/2014/4/173230-multipath-tcp/fulltext> , Last visited: 02.03.2016.
5. J Popat Khushi, Jigar A Raval, Samuel Johnson, and Bhavesh Patel. ``An efficient scheduling scheme of multipath tcp for mpi''.
6. Costin Raiciu, Mark Handley, and Damon Wischik. ``Coupled congestion control for multipath transport protocols''. RFC6356, October, 2011.
7. Sébastien Barré, Christoph Paasch, and Olivier Bonaventure. ``Multipath tcp: from theory to practice''. In NETWORKING 2011, pages 444-457. Springer, 2011.
8. Amanpreet Singh, Carmelita Goerg, Andreas Timm-Giel, Michael Scharf, and T-R Banniza. ``Performance comparison of scheduling algorithms for multipath transfer''. In Global Communications Conference (GLOBECOM), 2012 IEEE, pages 2653-2658. IEEE, 2012.

Questions??

Thank you!!!

Backup Slides

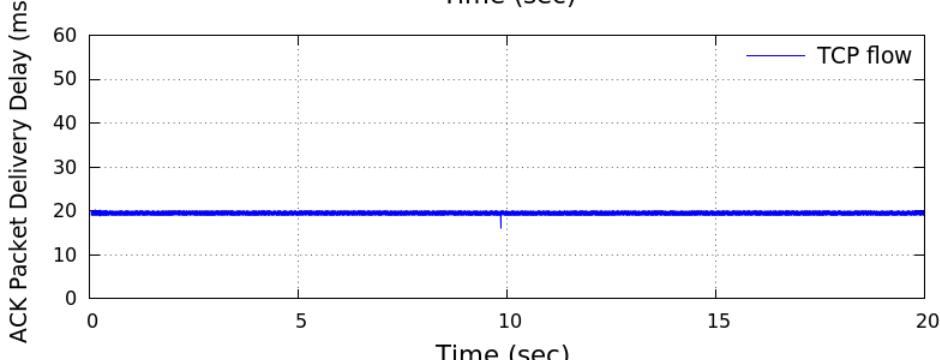
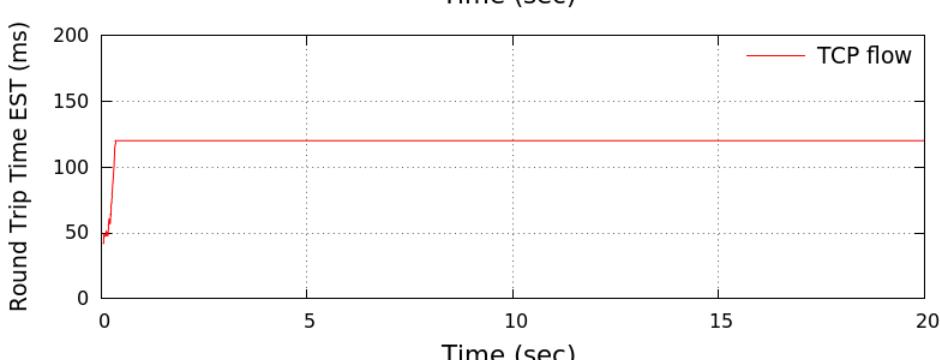
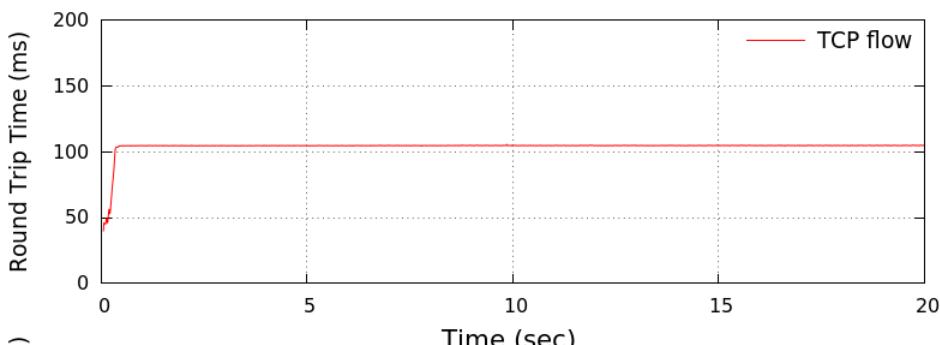
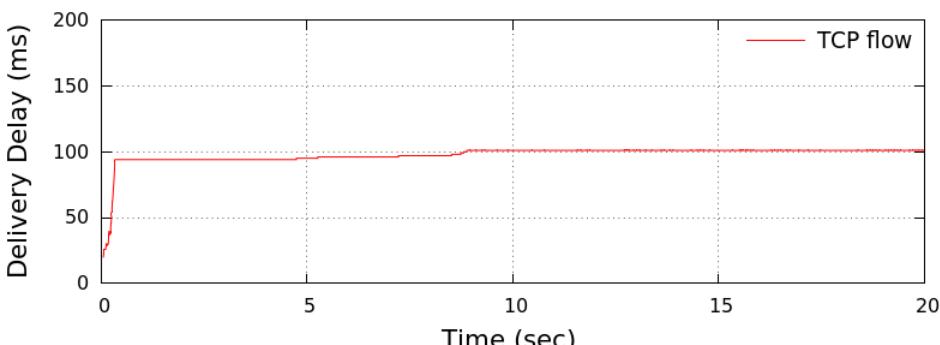
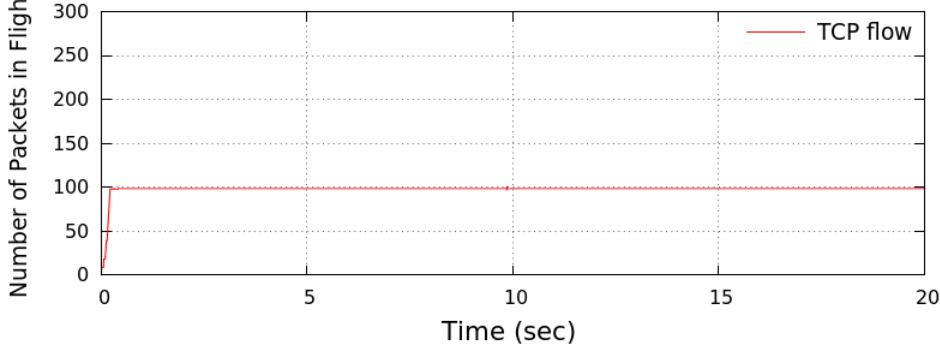
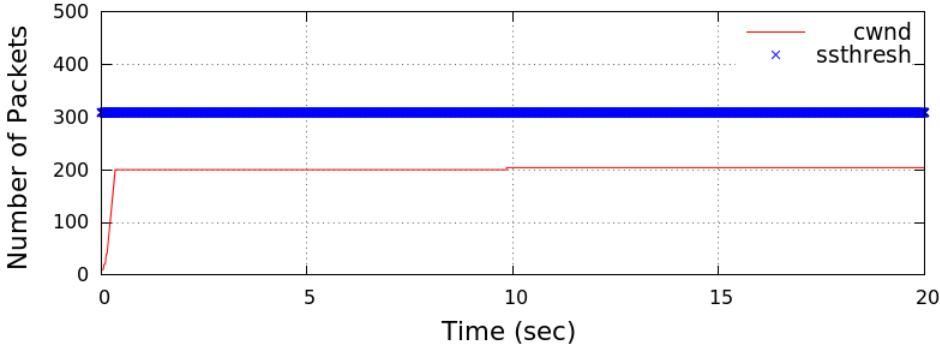
Delay - 40 ms (20 ms + 20 ms)



TCP with Router - (Mbit/sec)

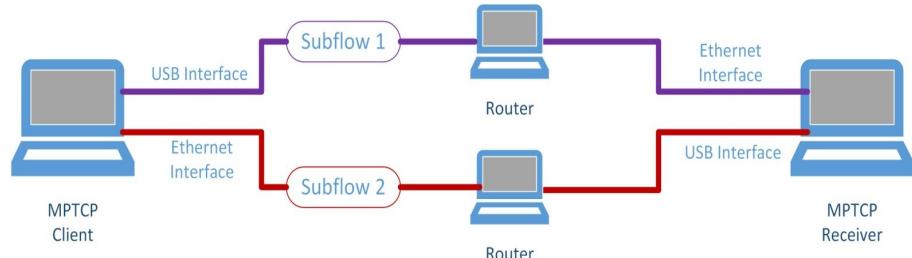
Link Capacity 10

Delay - 40 ms (20 ms ± 20 ms)



Performance Evaluation - Two Node Testbed with Router

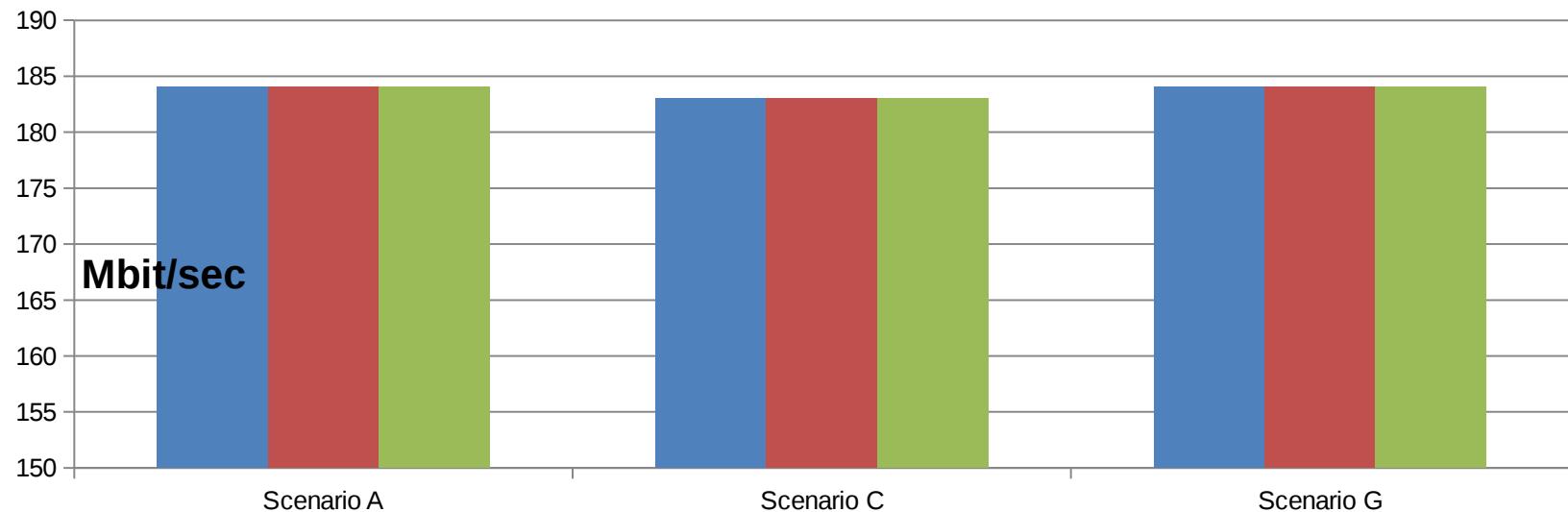
- Configuration :
 - Congestion Control = Reno
 - tc qdisc "add" dev "ethX" root netem delay "Yms" – at router



Capacity (Mbit/sec)	Name of Scenario	Subflow 1			Subflow 2		
		FPD (ms)	BPD (ms)	RTT (ms)	FPD (ms)	BPD (ms)	RTT (ms)
100/100 & 100/10 & 10/10	A	10	20	30	20	10	30
	C	10	40	50	20	20	40
	G	10	10	20	10	10	20

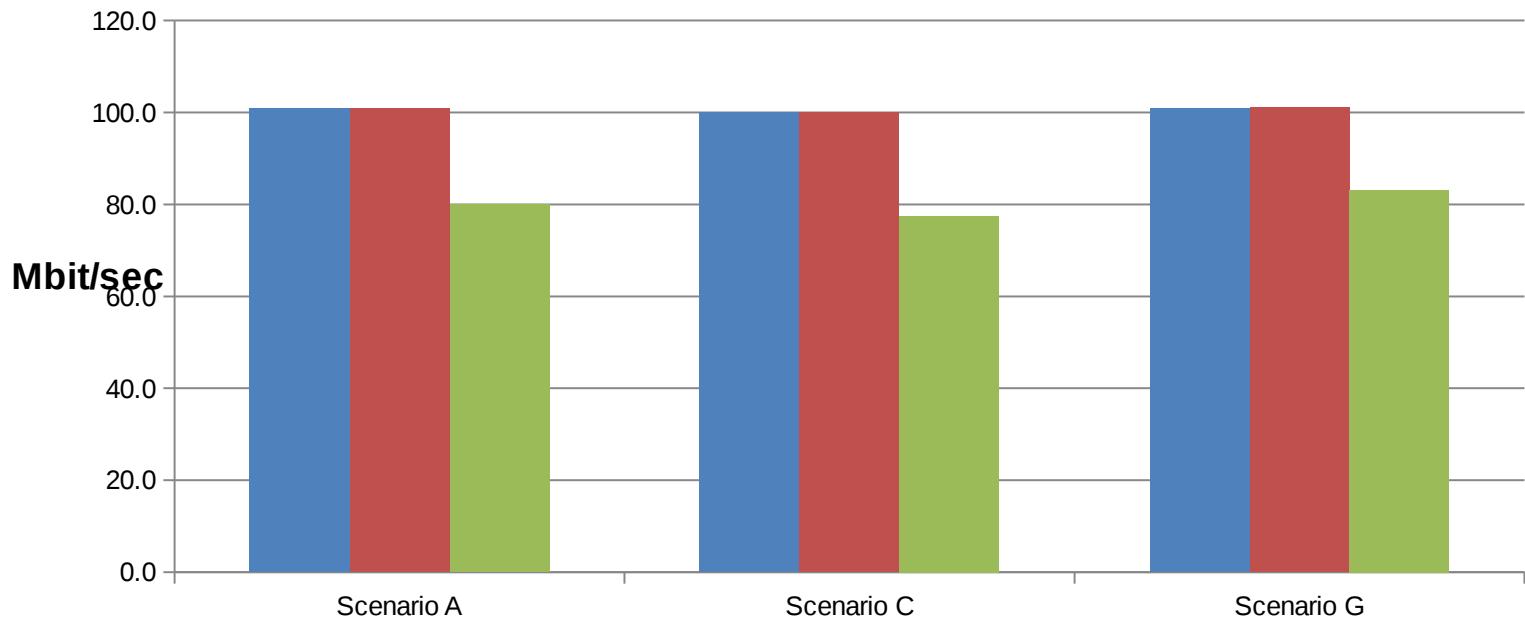
Link Capacity 100/100 (Mbit/sec)

Throughput Comparison Table



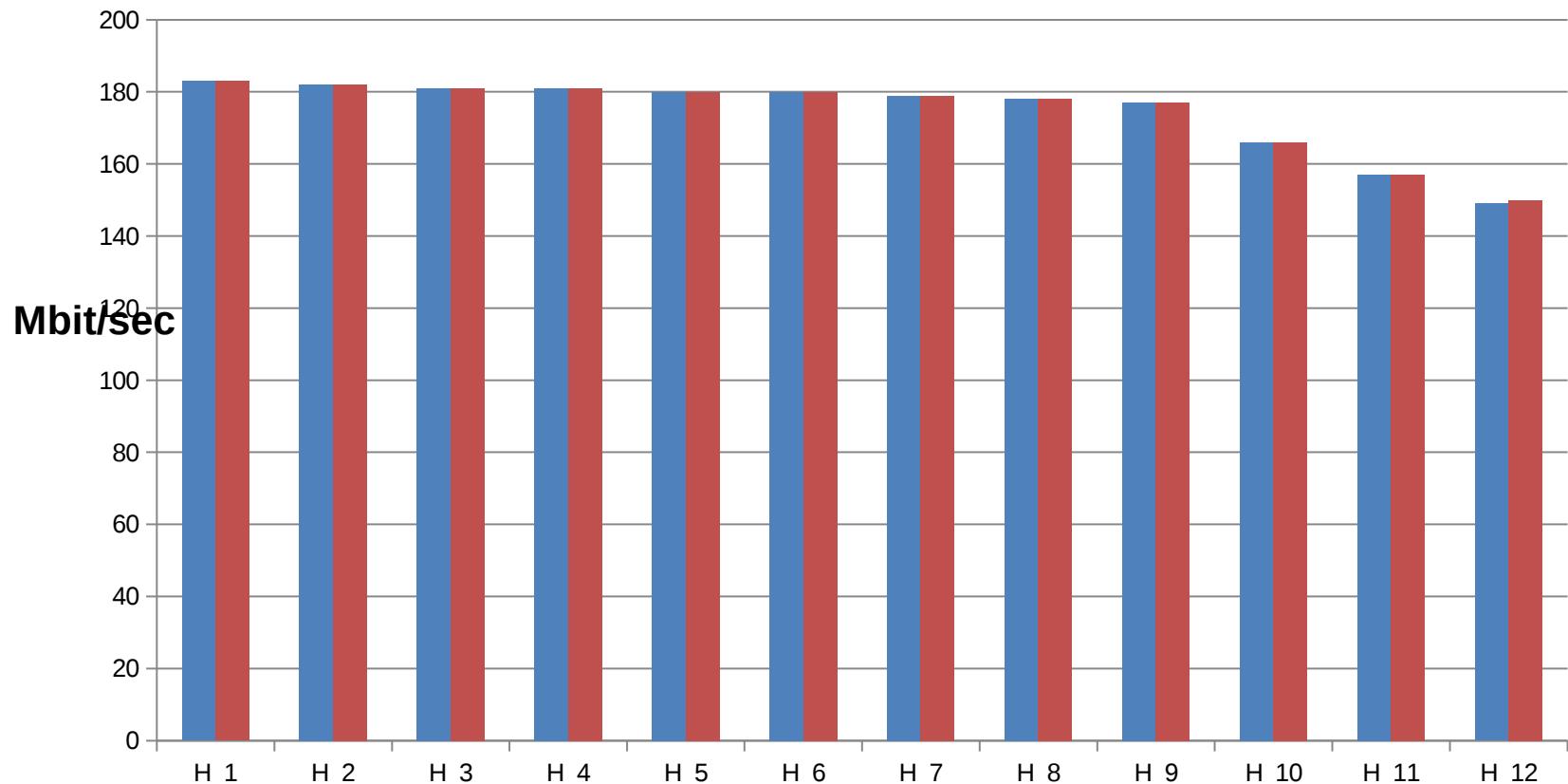
Link Capacity 100/10 (Mbit/sec)

Throughput Comparison Table



Link Capacity 100/100 (Mbit/sec)

Throughput Comparison Table



Graphical Analysis

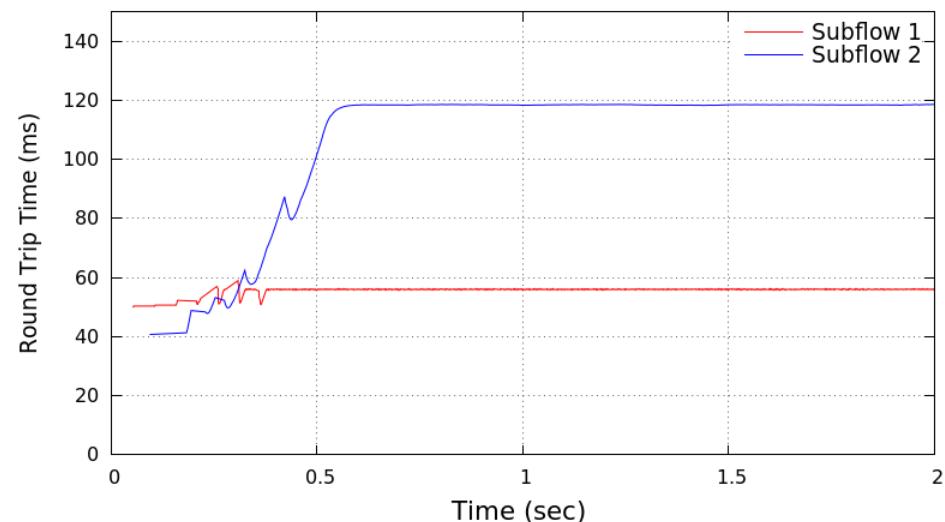
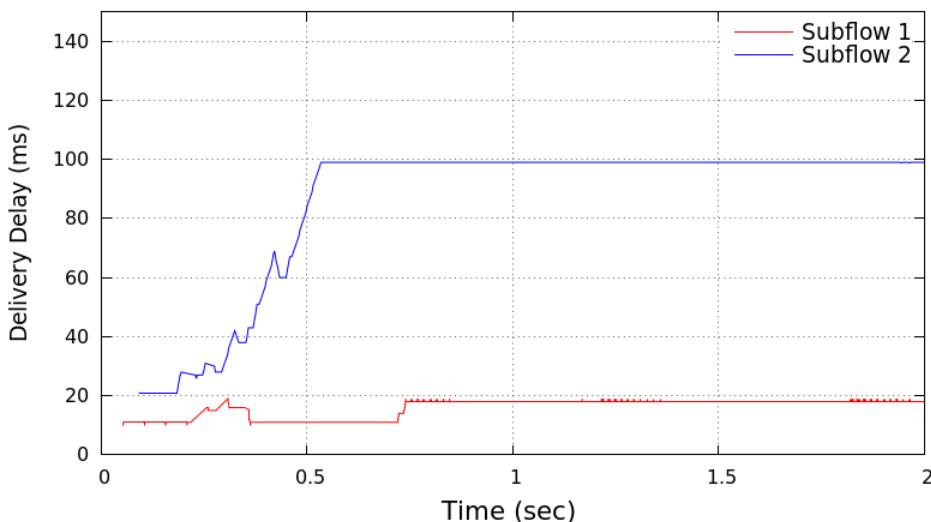
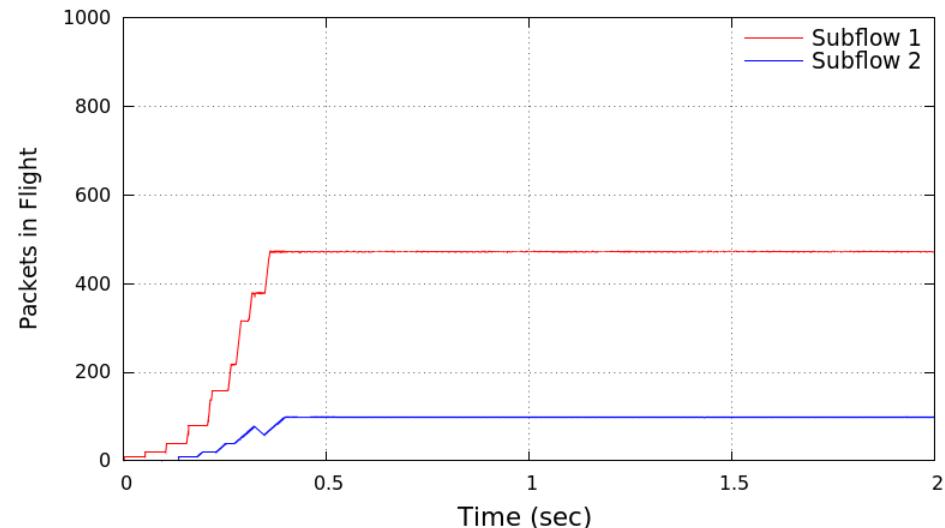
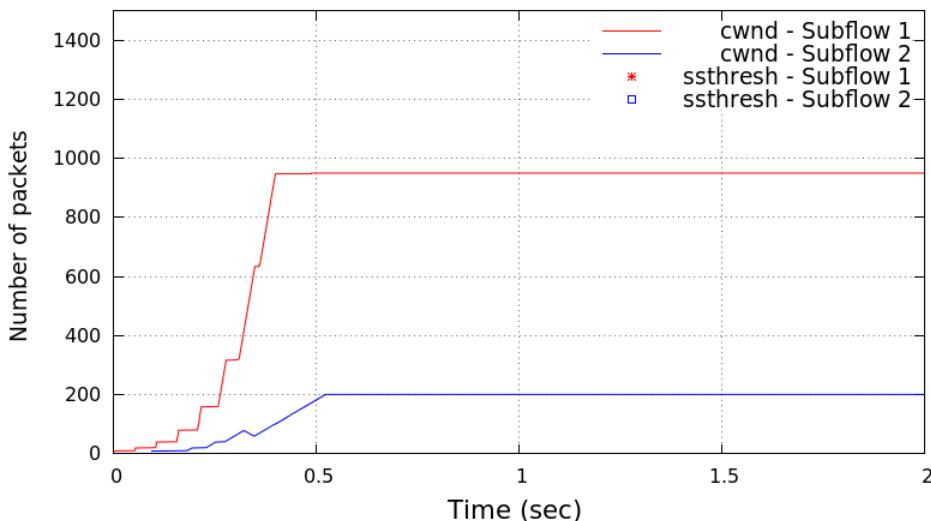
Link Capacity 100/10 (Mbit/sec)

Scenario C

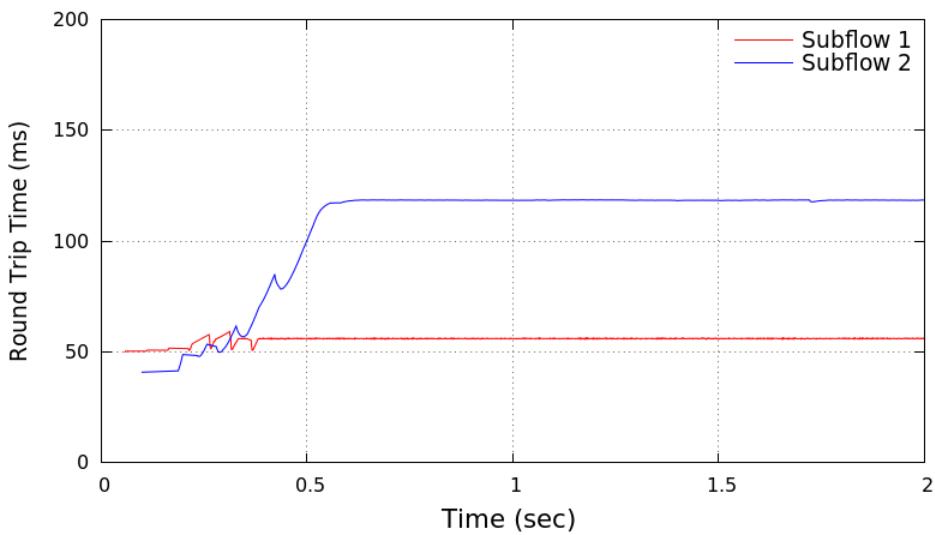
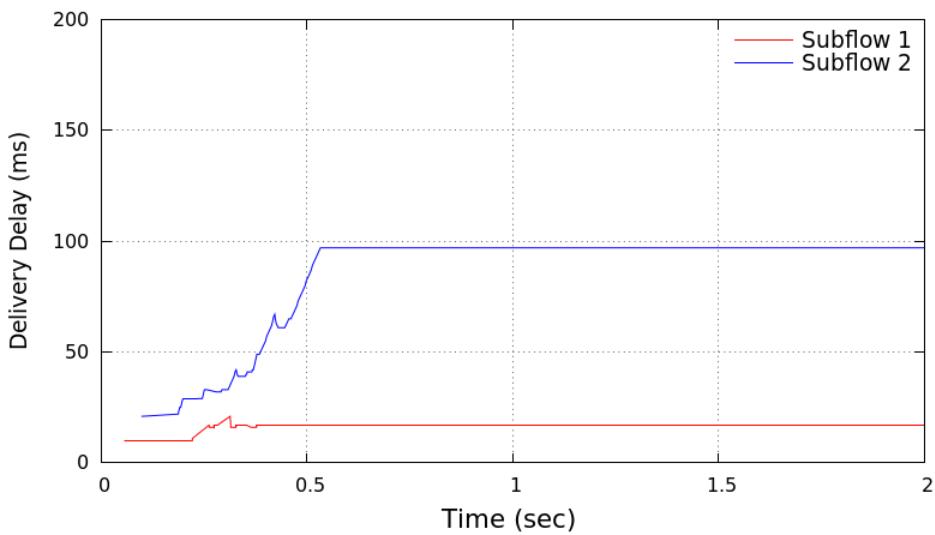
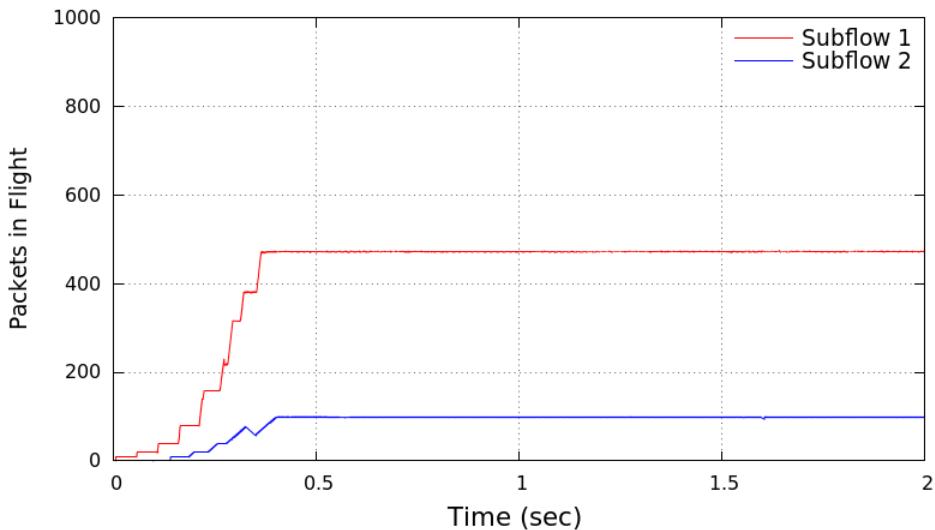
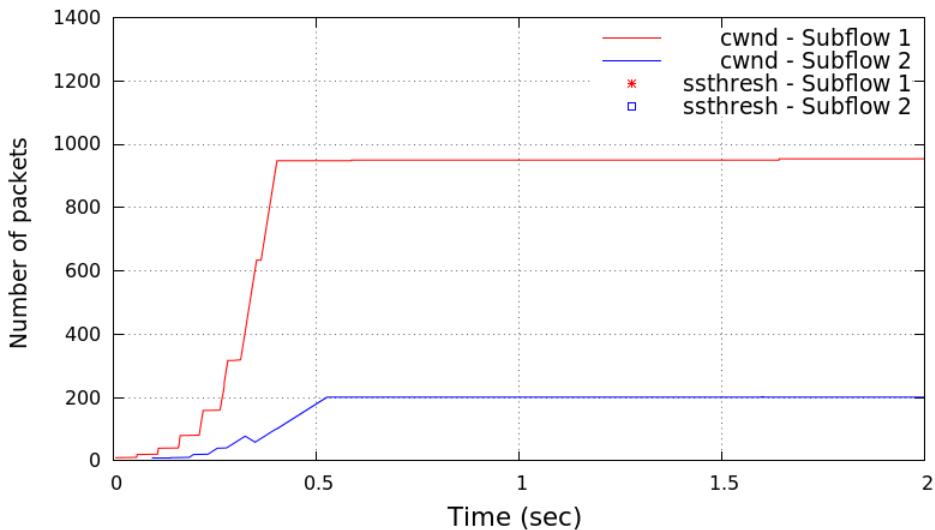
Subflow 1 RTT - 50 ms (10 ms + 40 ms)

Subflow 2 RTT - 40 ms (20 ms + 20 ms)

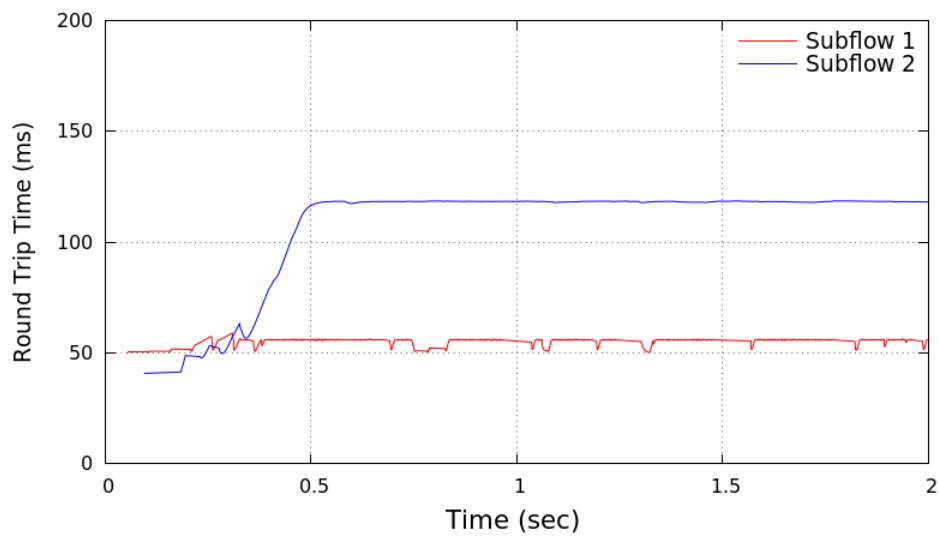
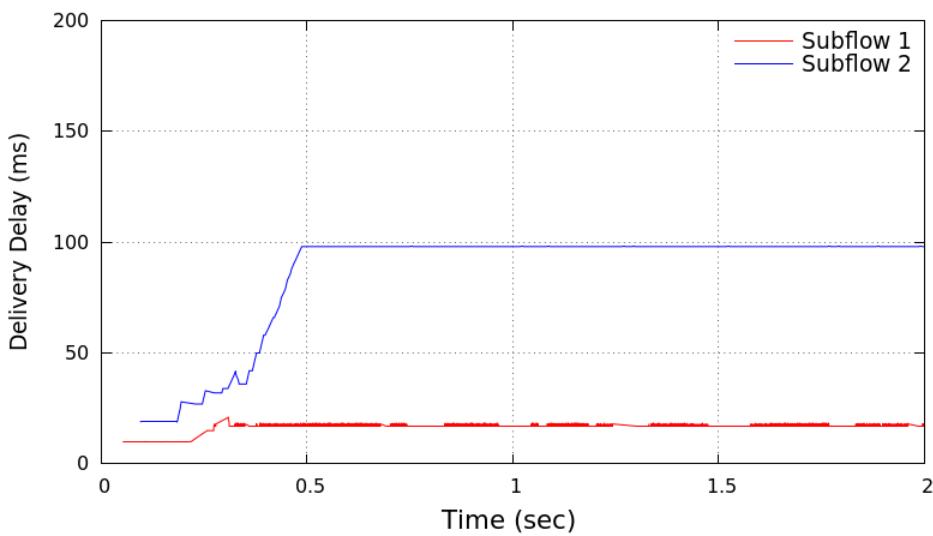
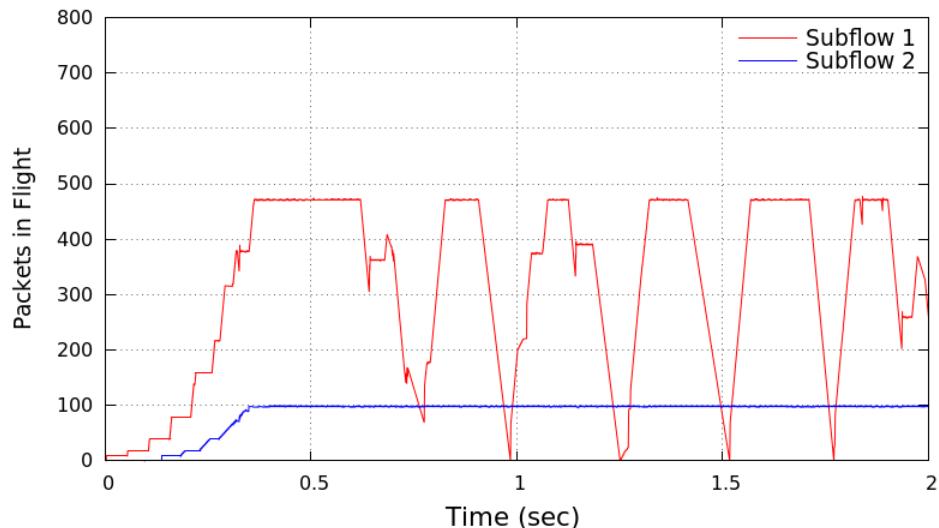
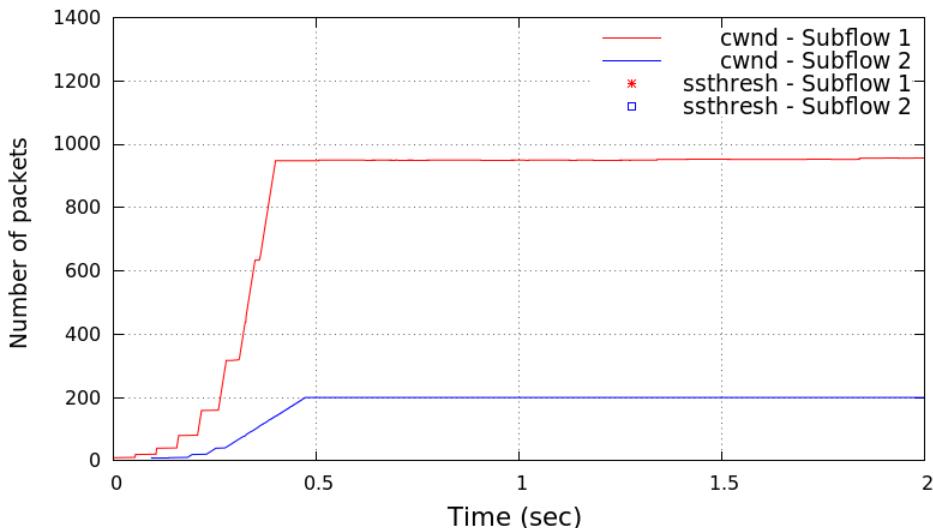
Delivery Delay Scheduler



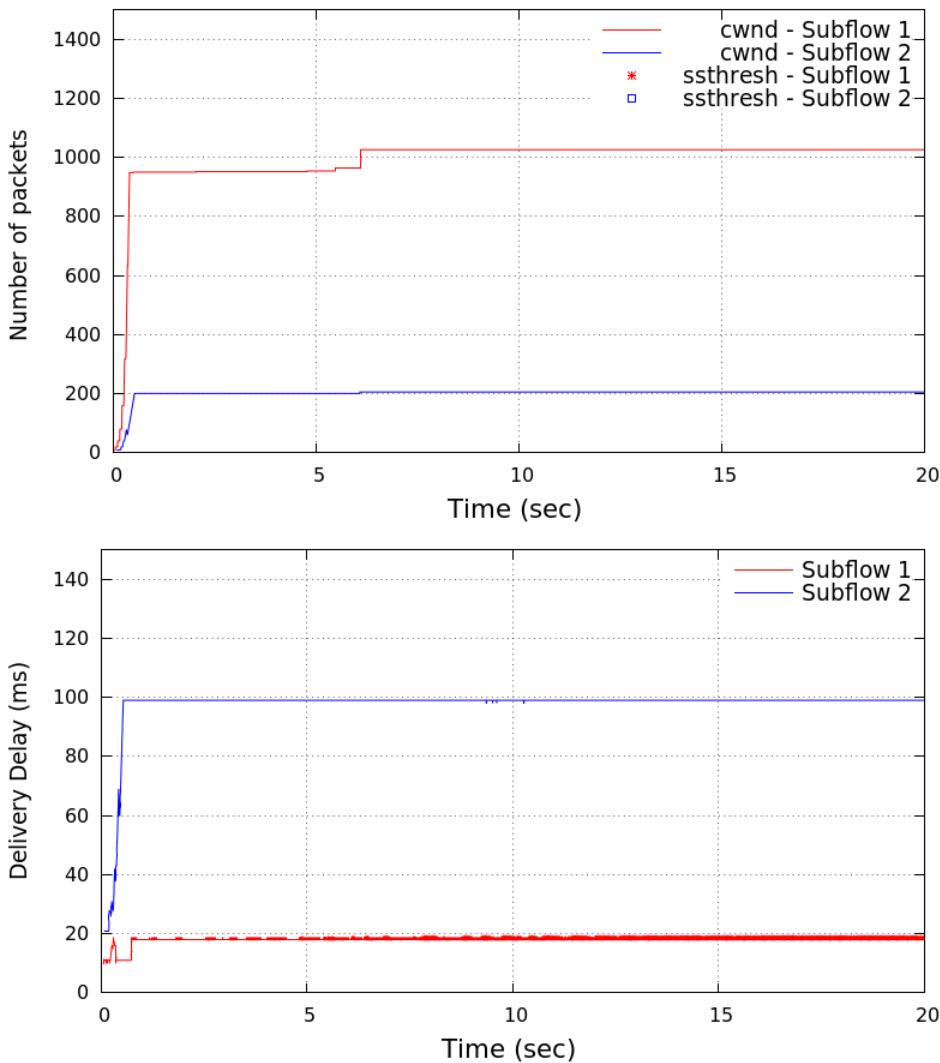
sRTT Scheduler



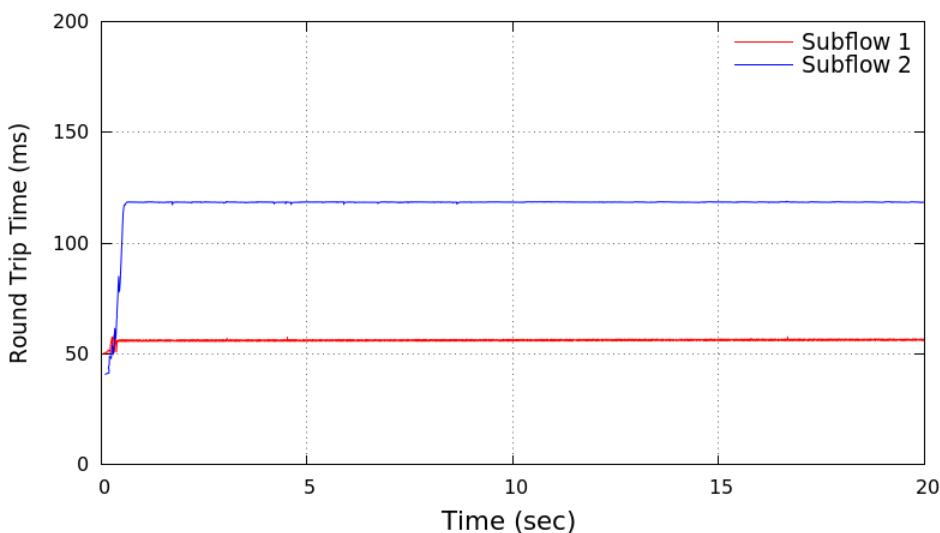
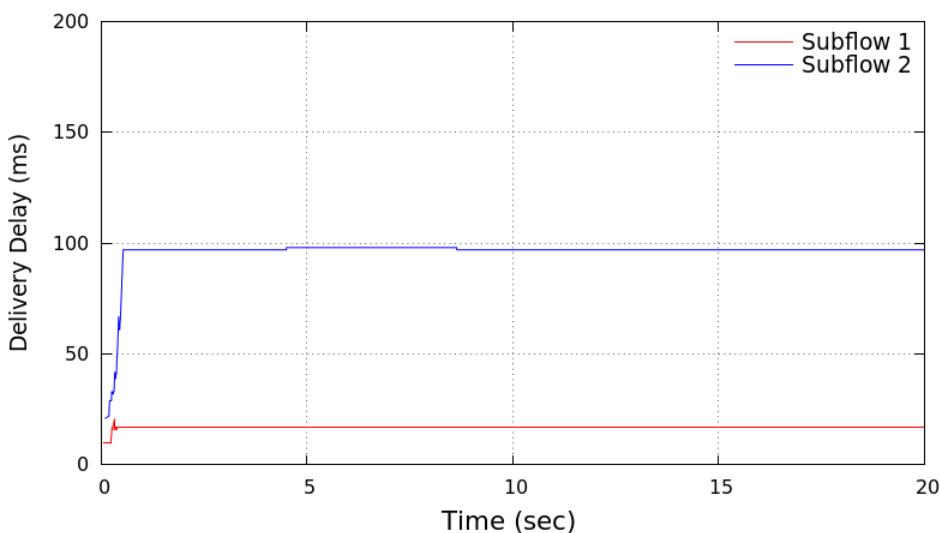
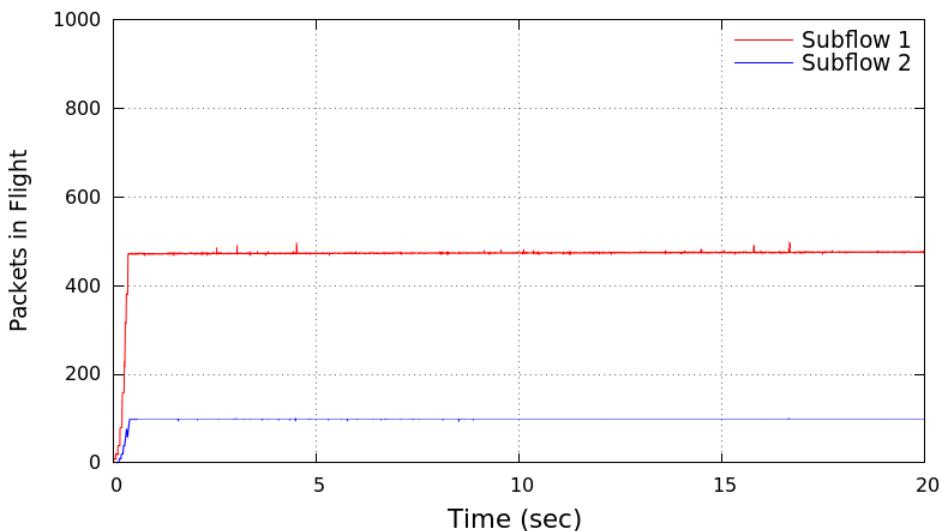
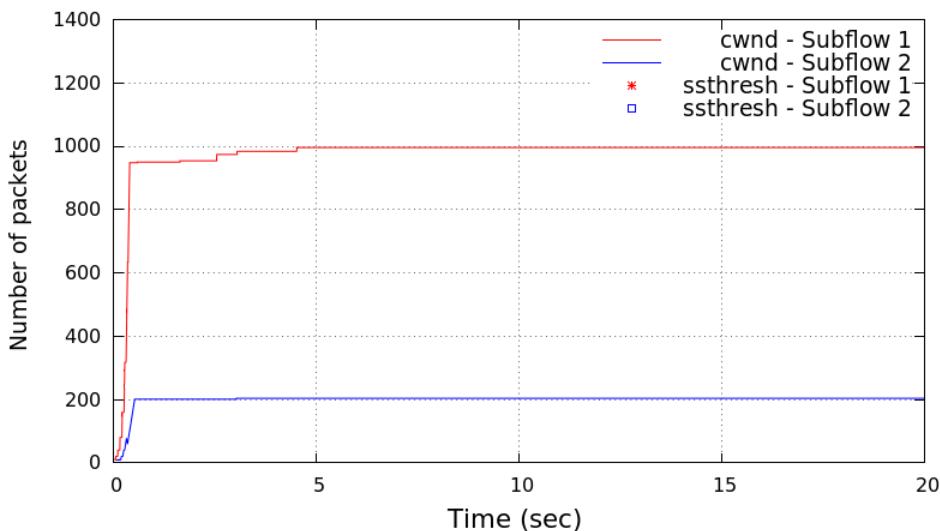
RR Scheduler



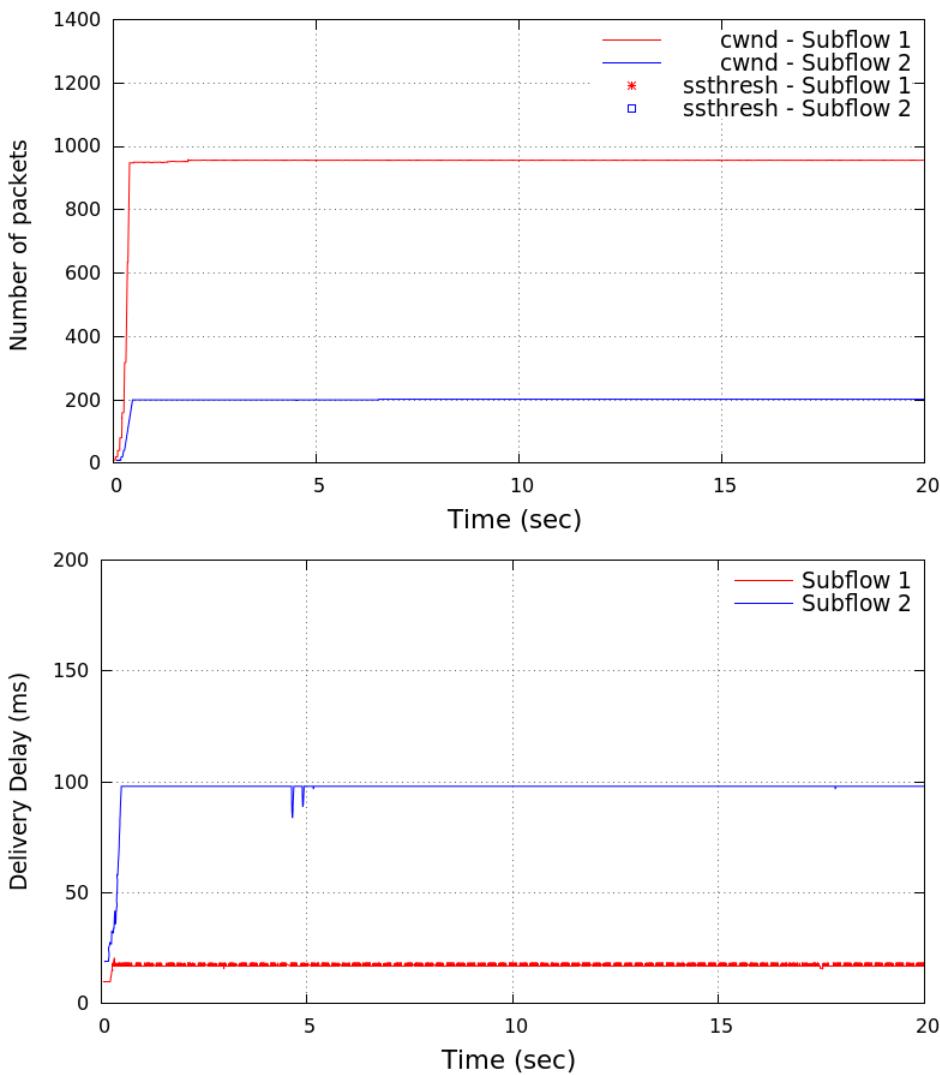
Delivery Delay Scheduler (20 sec)



Shortest Round Trip Time Scheduler (20 sec)



Round Robin Scheduler (20 sec)



Graphical Analysis

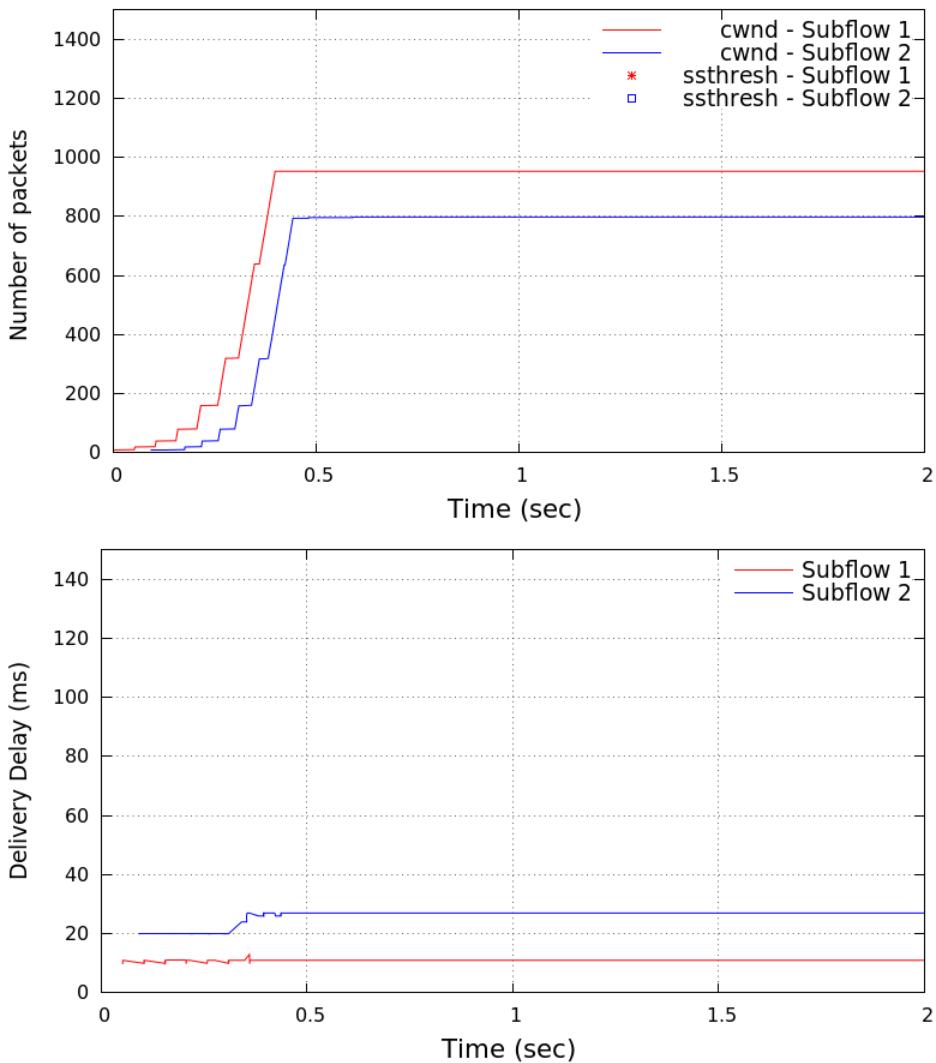
Link Capacity 100/100 (Mbit/sec)

Scenario C

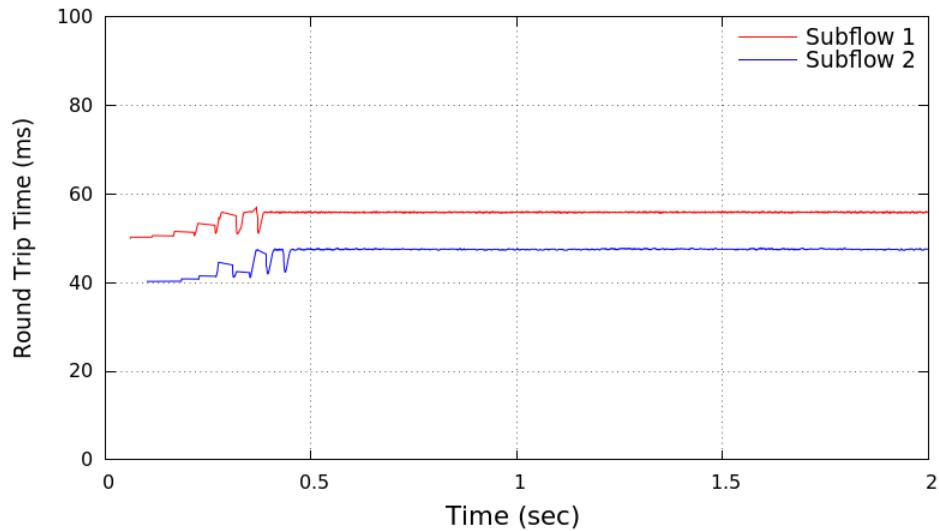
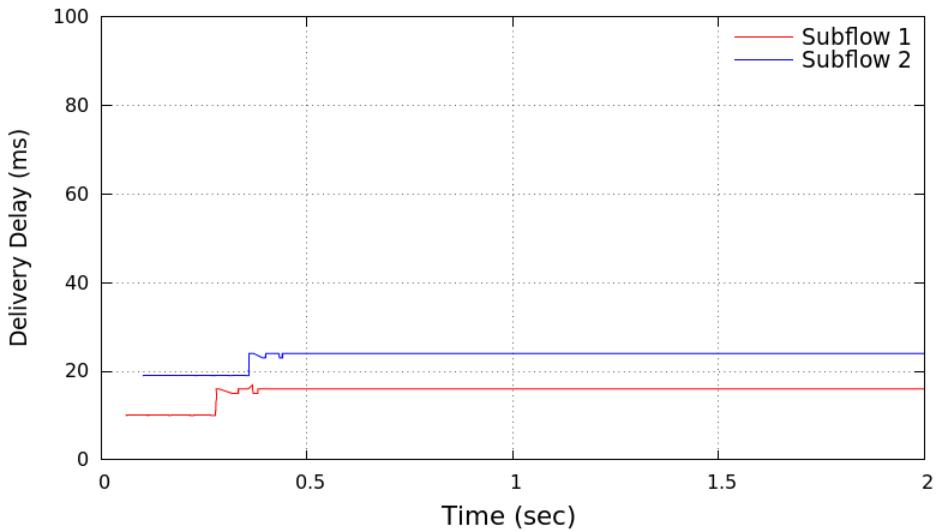
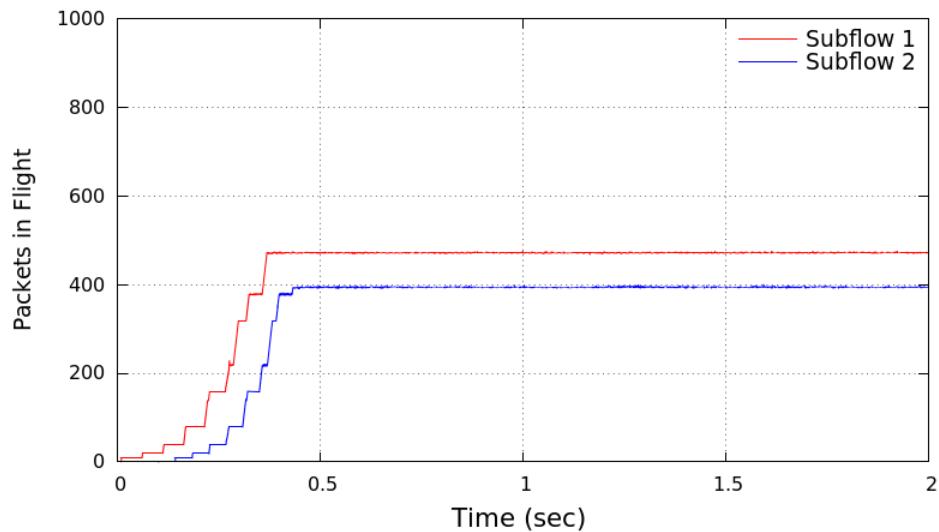
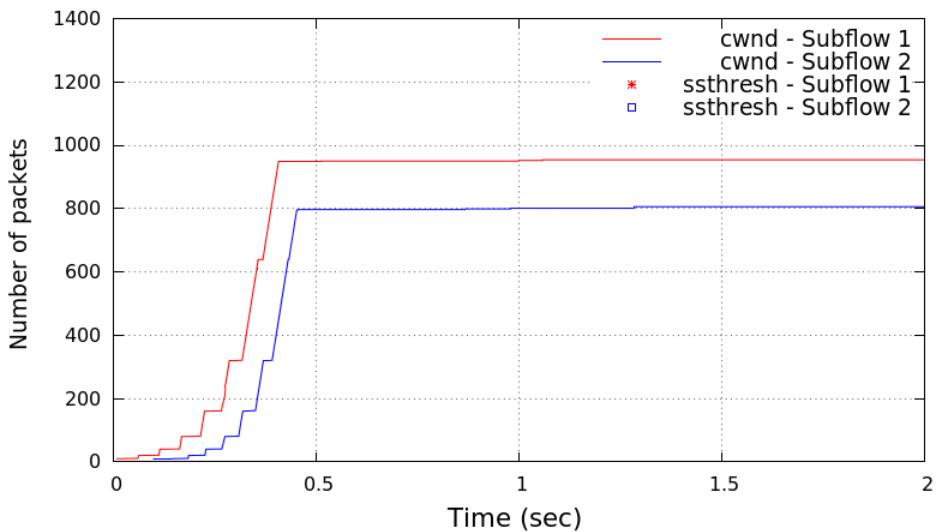
Subflow 1 RTT - 50 ms (10 ms + 40 ms)

Subflow 2 RTT - 40 ms (20 ms + 20 ms)

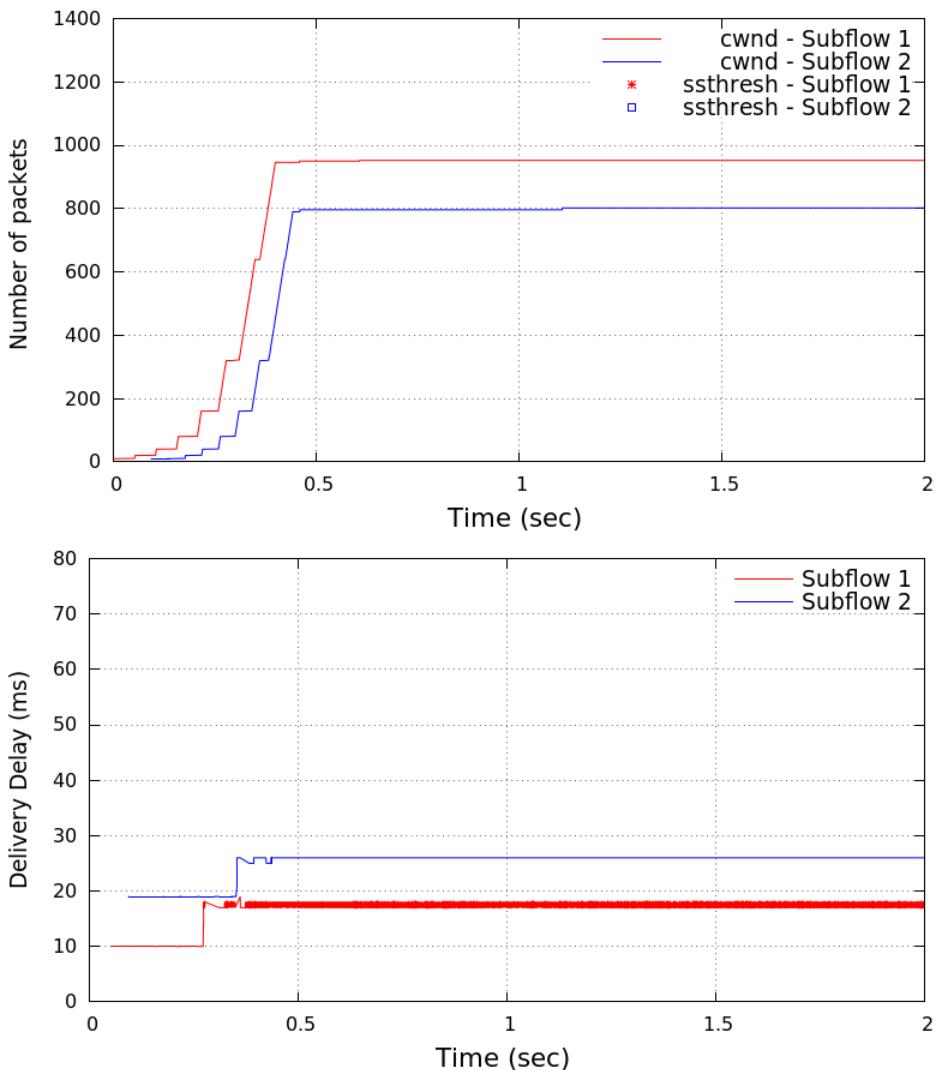
Delivery Delay Scheduler



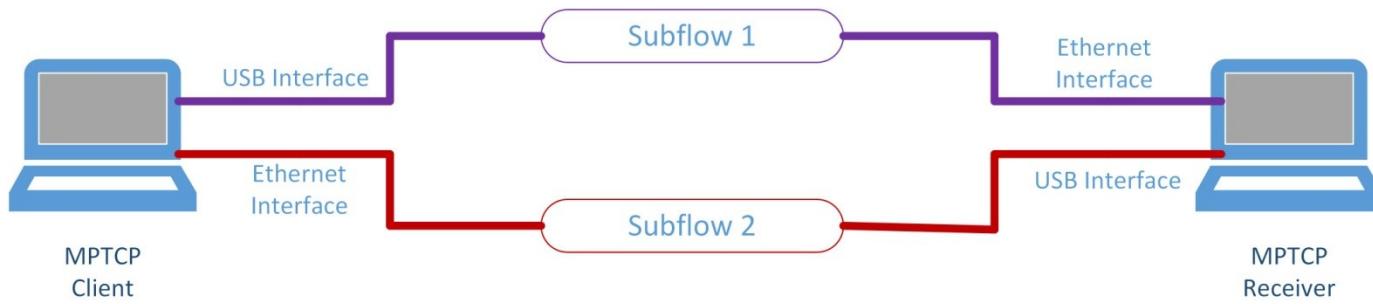
Shortest Round Trip Time Scheduler



Round Robin Scheduler



Graphical Analysis - 2 Node without Router



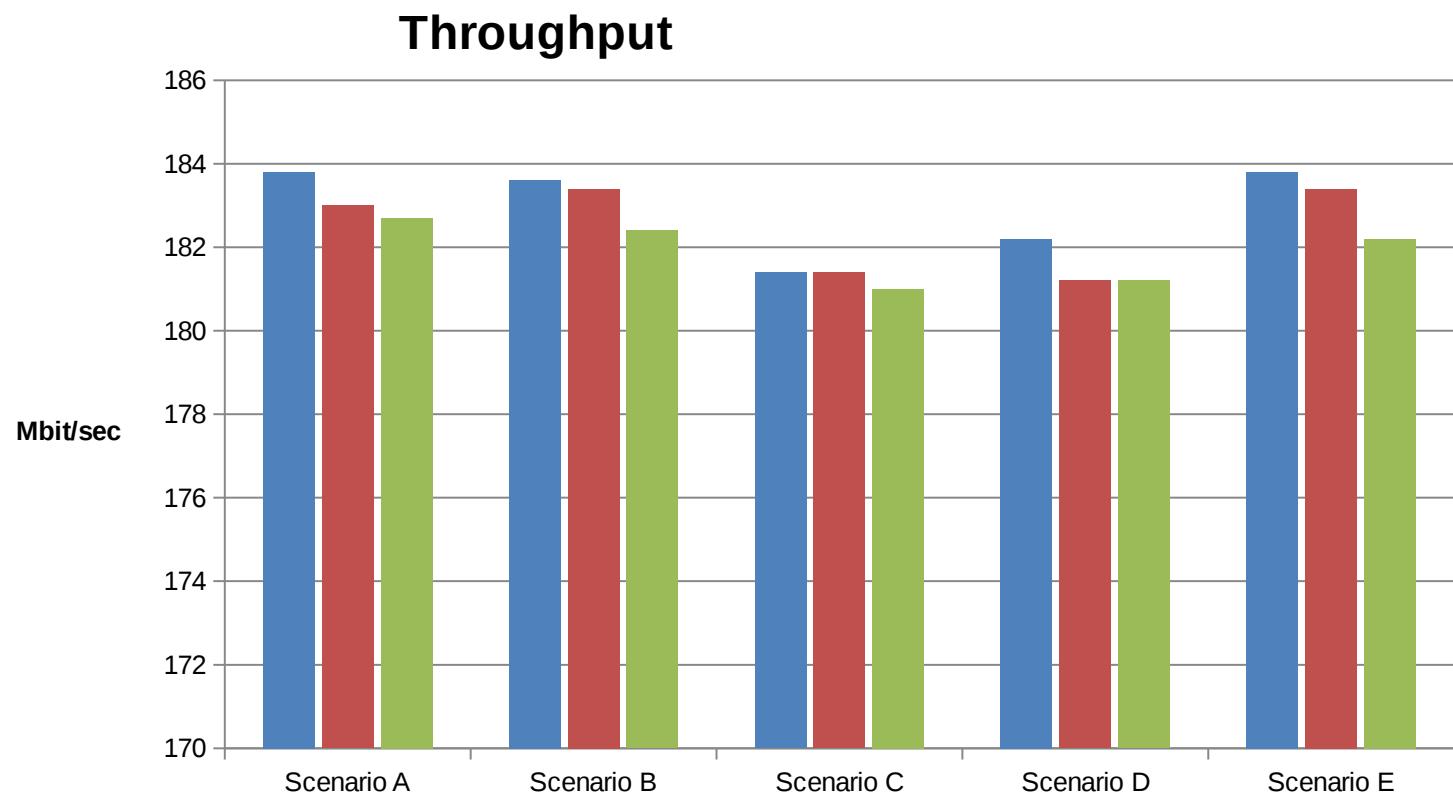
Link Capacity 100/10 (Mbit/sec)
Scenario B

Subflow 1 RTT – 20 ms (10 ms + 10 ms)
Subflow 2 RTT – 30 ms (20 ms + 10 ms)

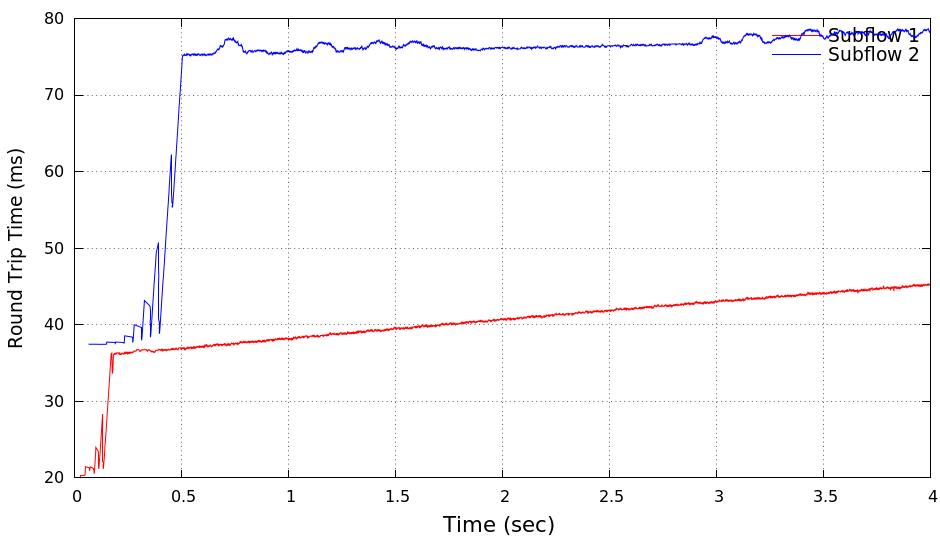
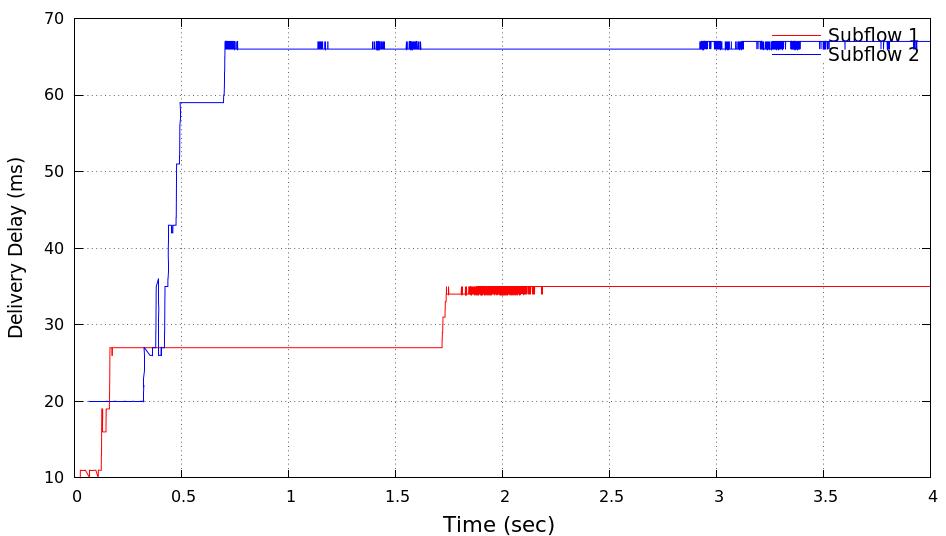
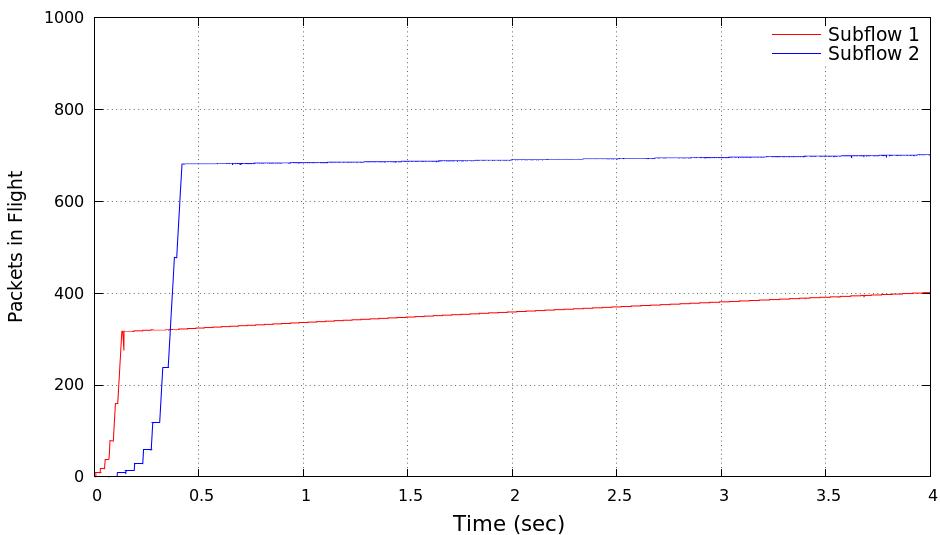
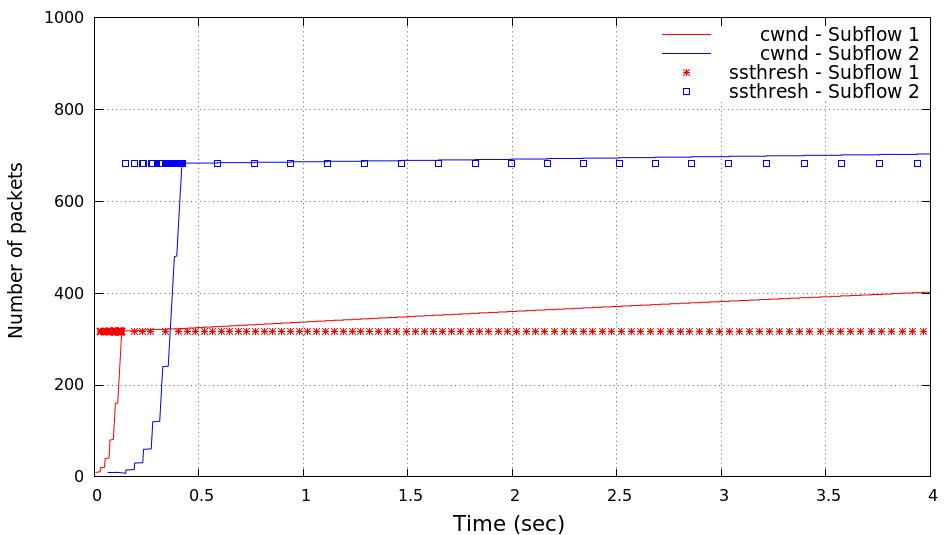
Scenarios

Capacity	Name Of scenarios	Subflow 1			Subflow 2			TCP Flow		
		FPD (ms)	BPD (ms)	RTT (ms)	FPD (ms)	BPD (ms)	RTT (ms)	FPD (ms)	BPD (ms)	RTT (ms)
100/100 & 100/10 & 10/10	A	10	20	30	20	10	30	10	20	30
	B	10	10	20	20	10	30	10	10	20
	C	10	40	50	20	20	40	10	40	50
	D	10	20	30	20	40	60	10	20	30
	E	10	10	20	20	20	40	10	10	20

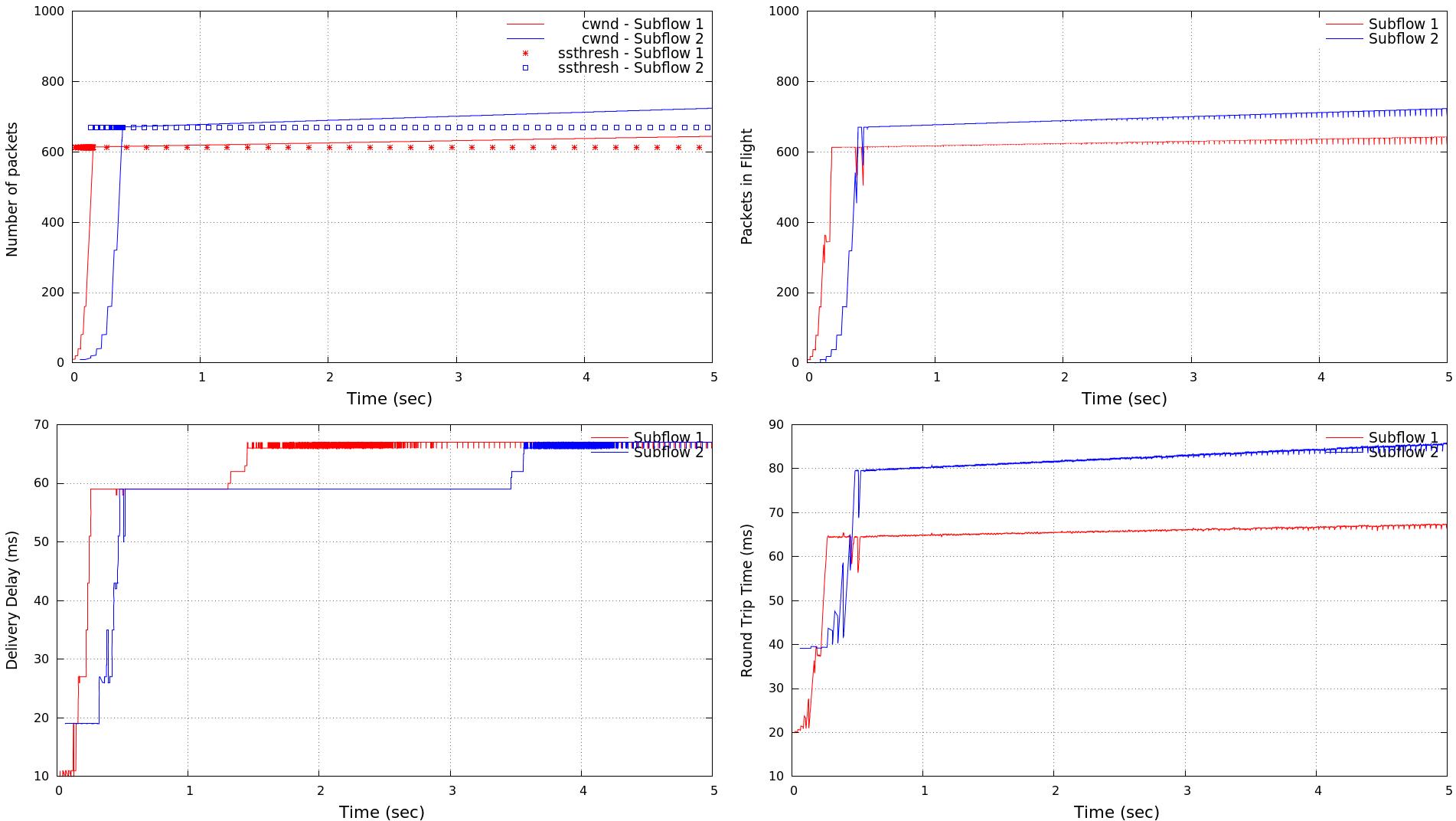
Link Capacity 100/100 (Mbit/sec)



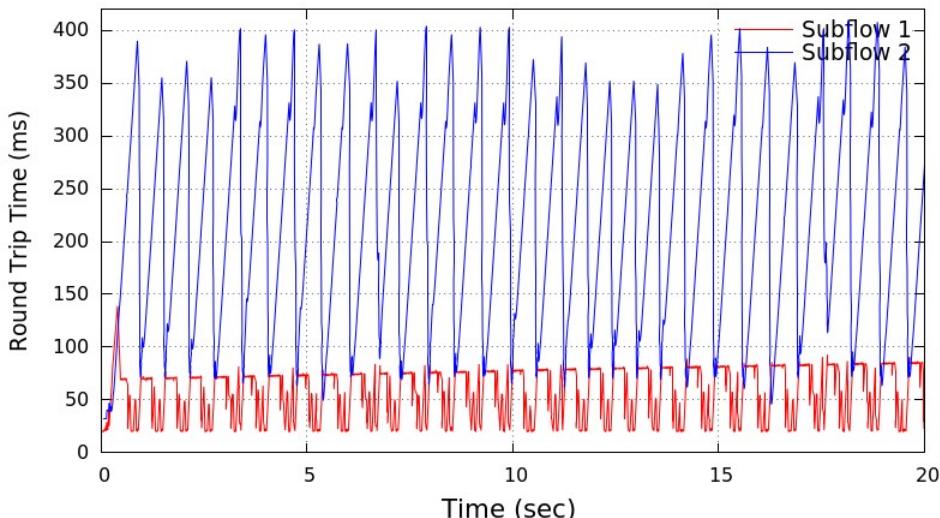
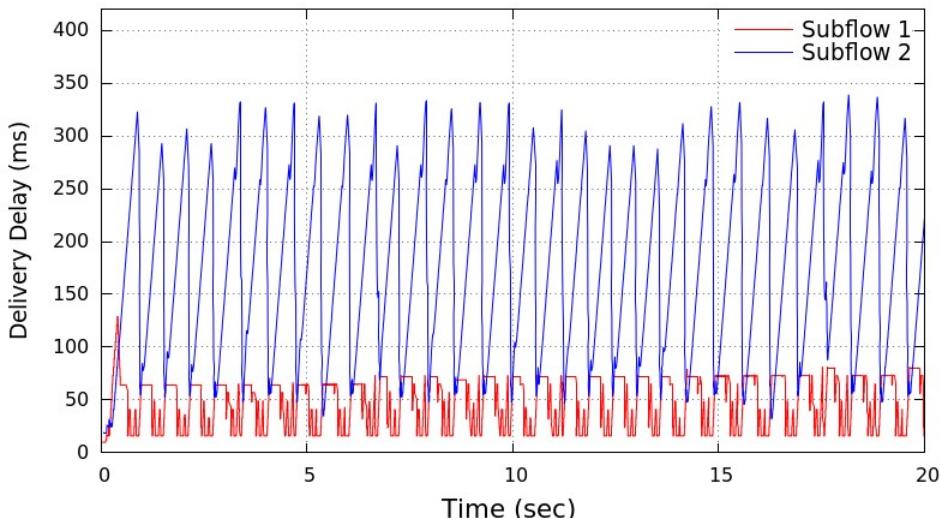
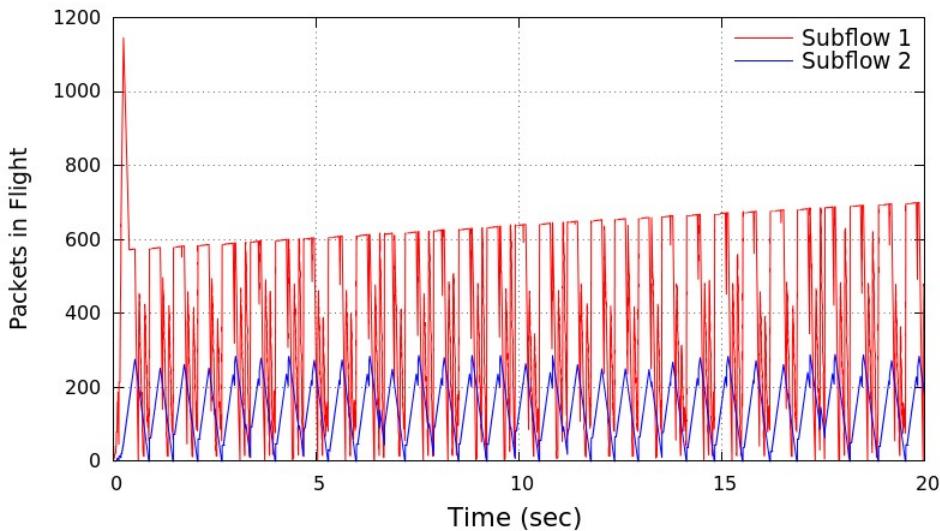
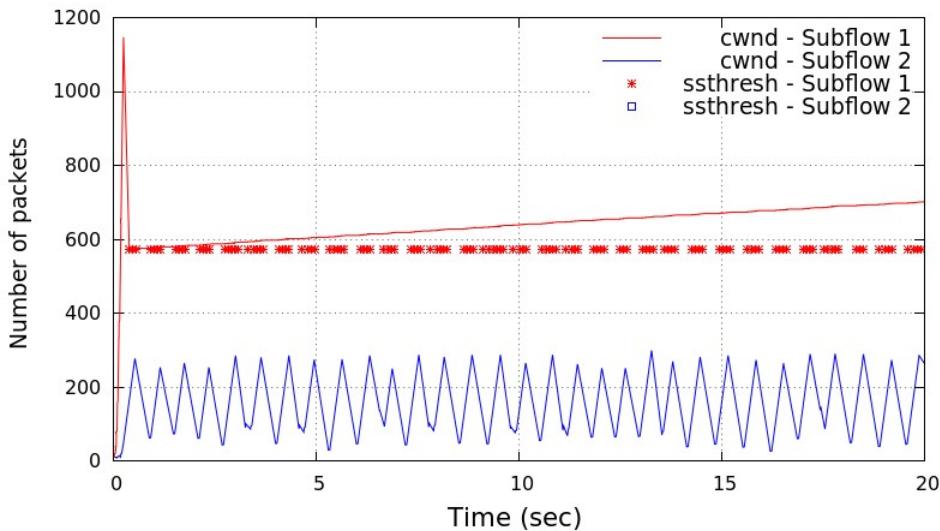
Shortest Round Trip Time Scheduler - E (100/100)



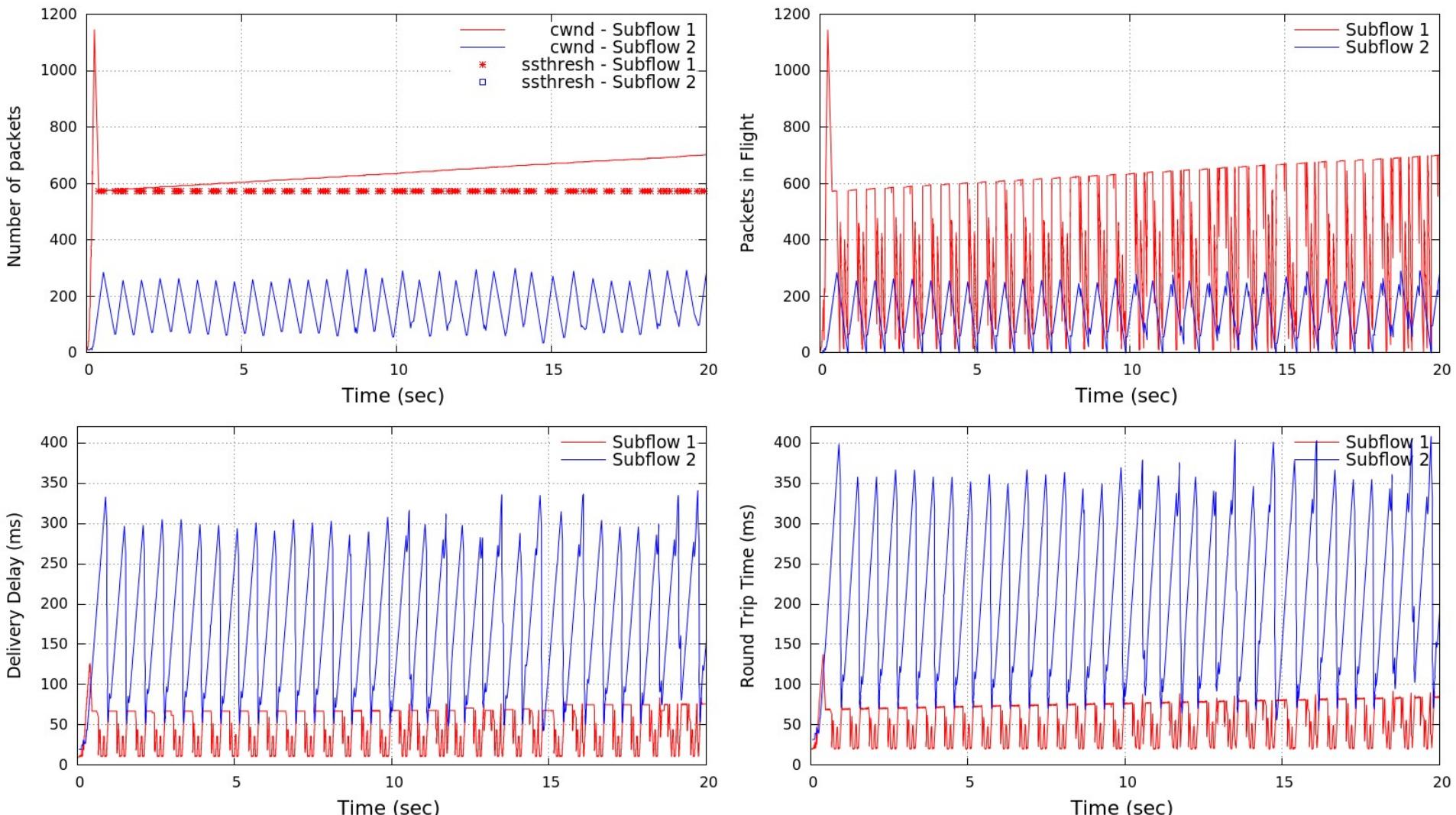
Round Robin Scheduler - E (100/100)



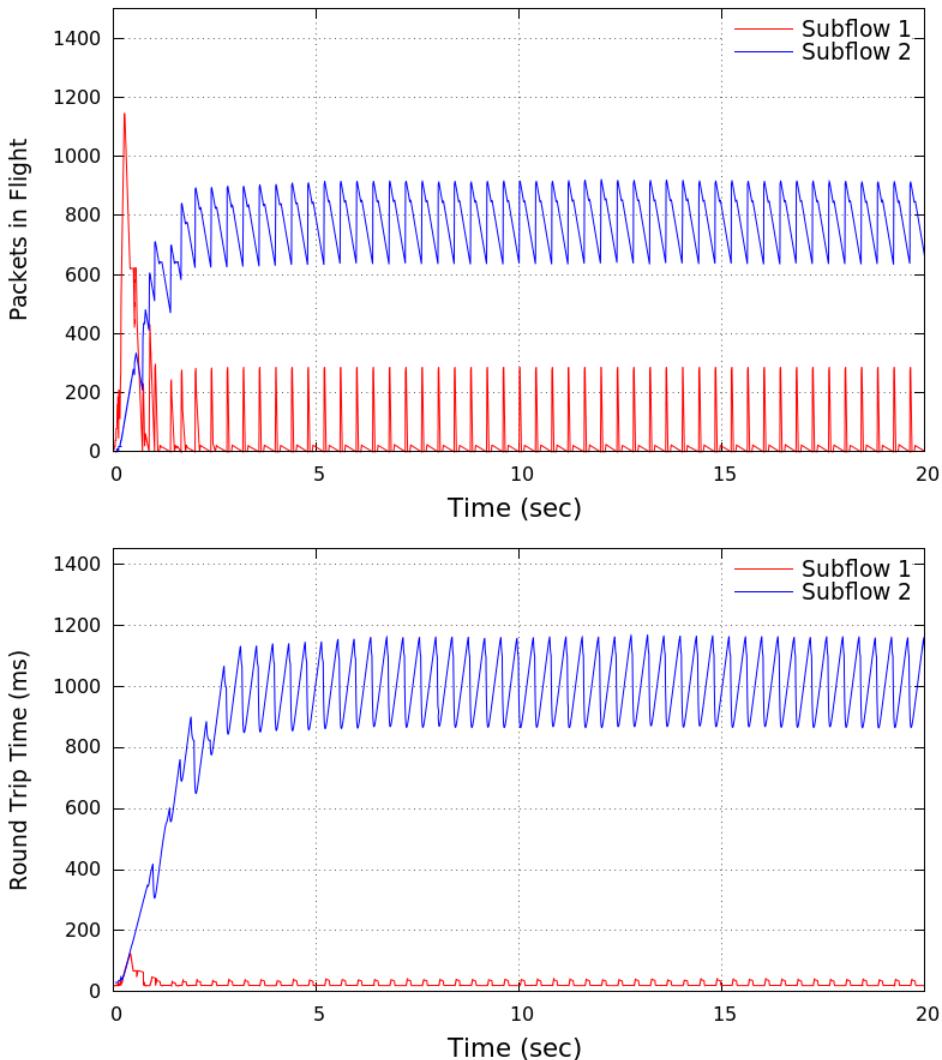
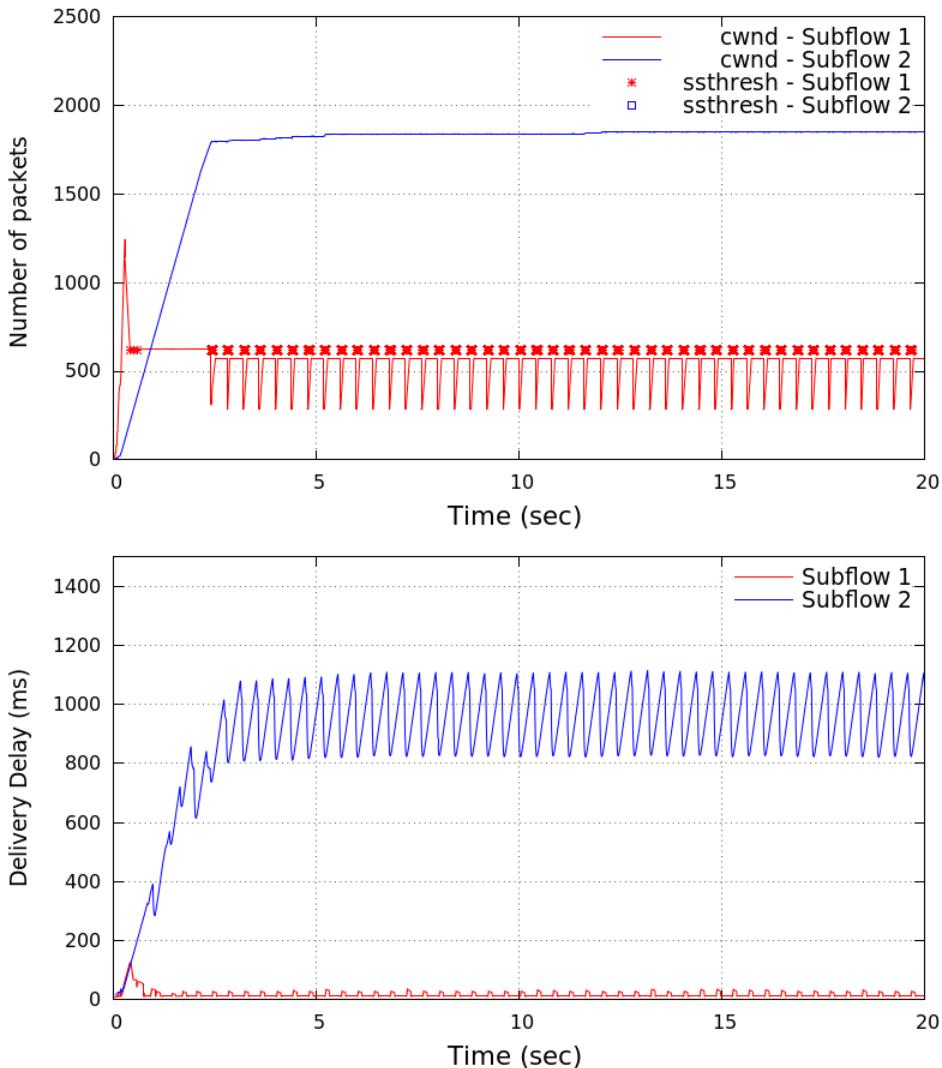
Delivery Delay Scheduler



Shortest Round Trip Scheduler



Round Robin Scheduler



Test Topology - Extended Testbed

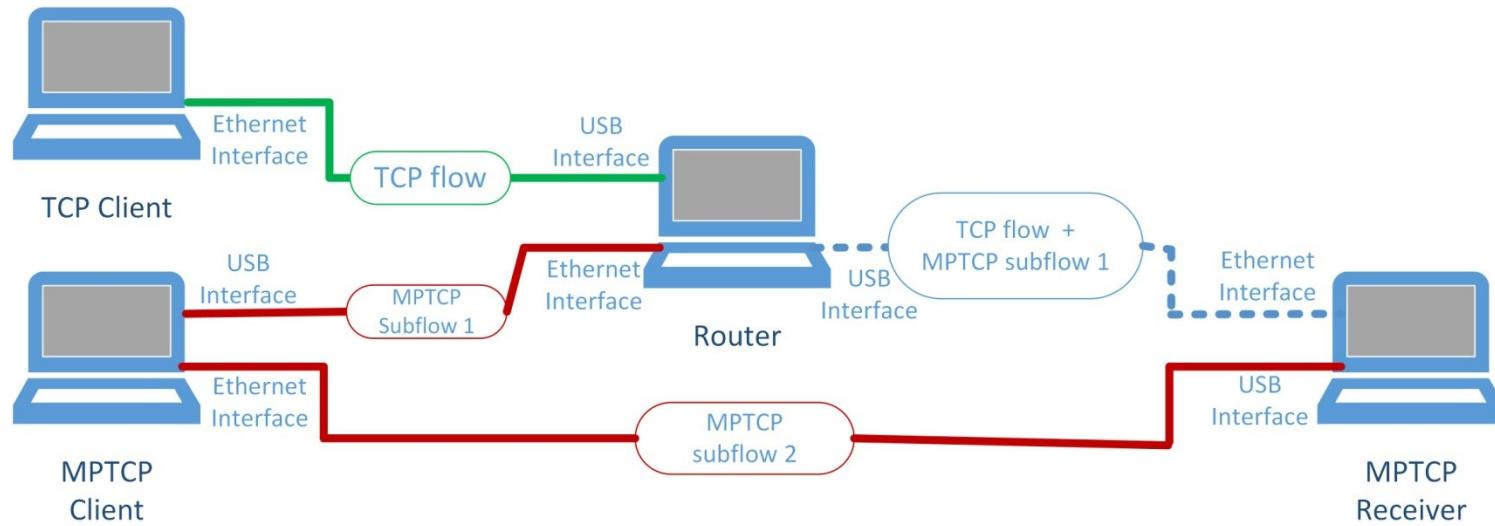
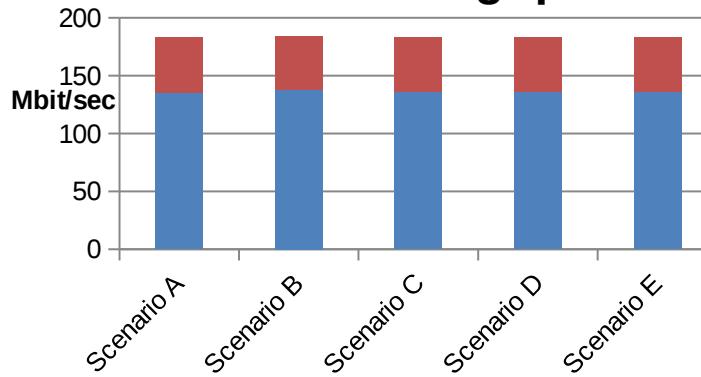


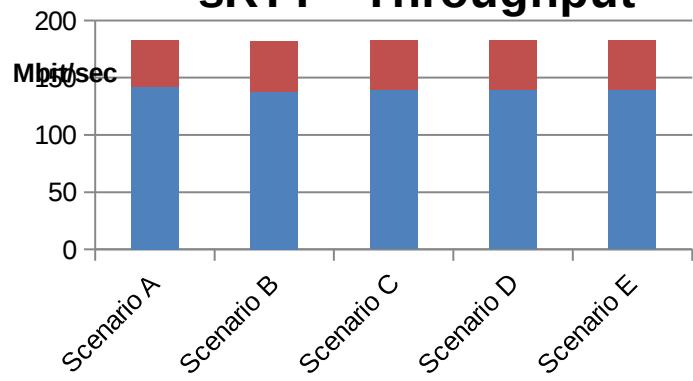
Figure : Extended Testbed

Throughput Comparison

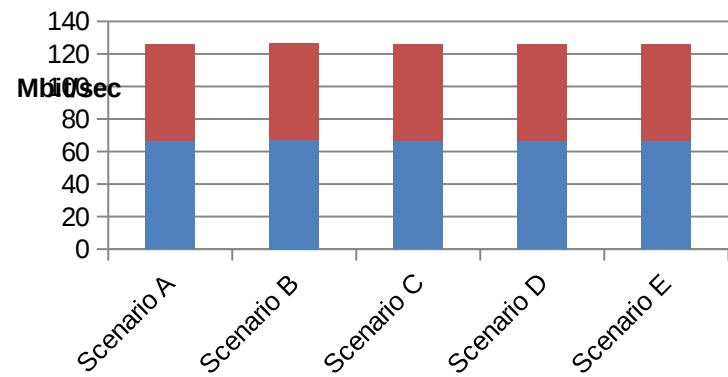
DD - Throughput



sRTT - Throughput



RR - Throughput



Graphical Analysis - Extended Testbed

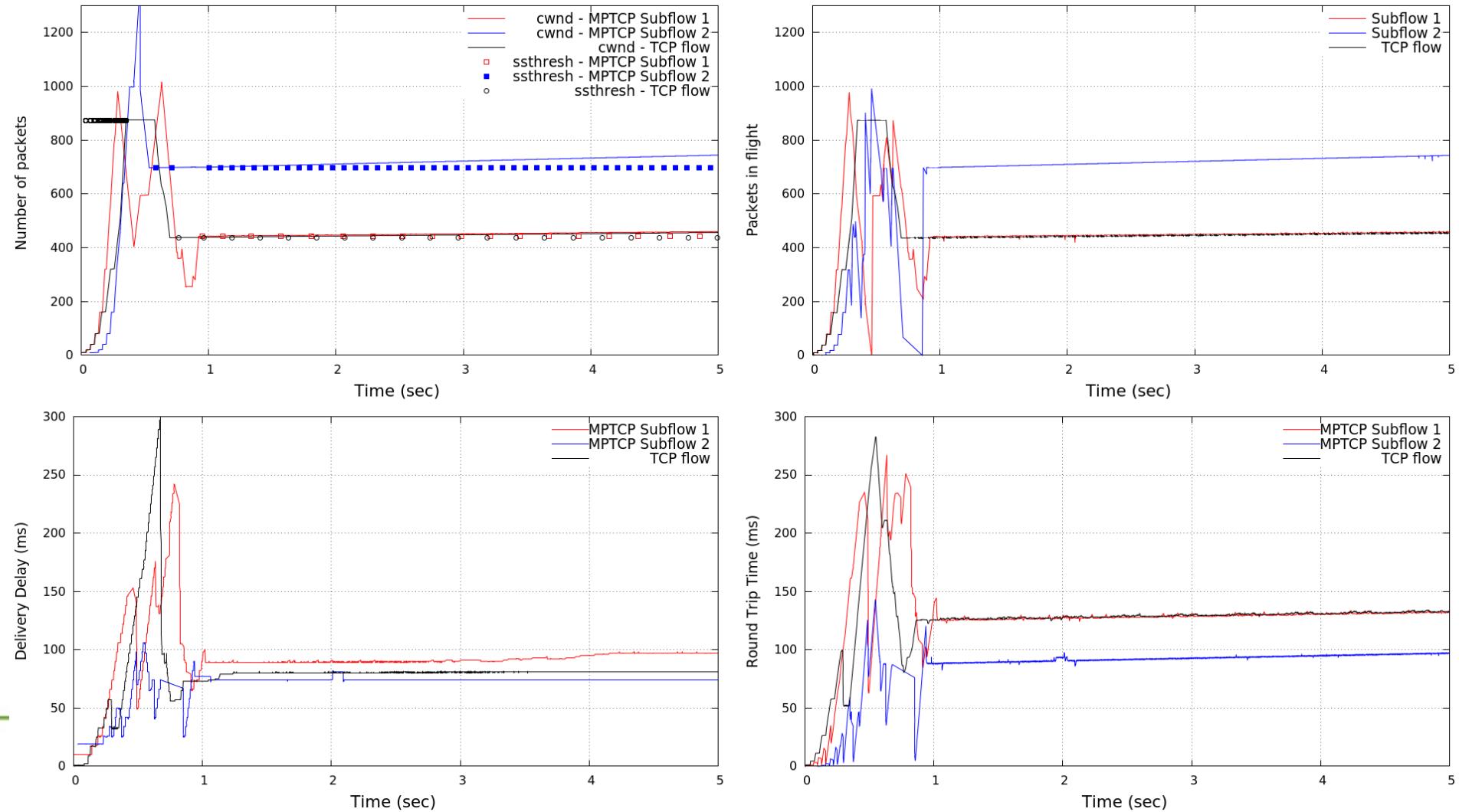
Link Capacity 100/100 (Mbit/sec)

Scenario A

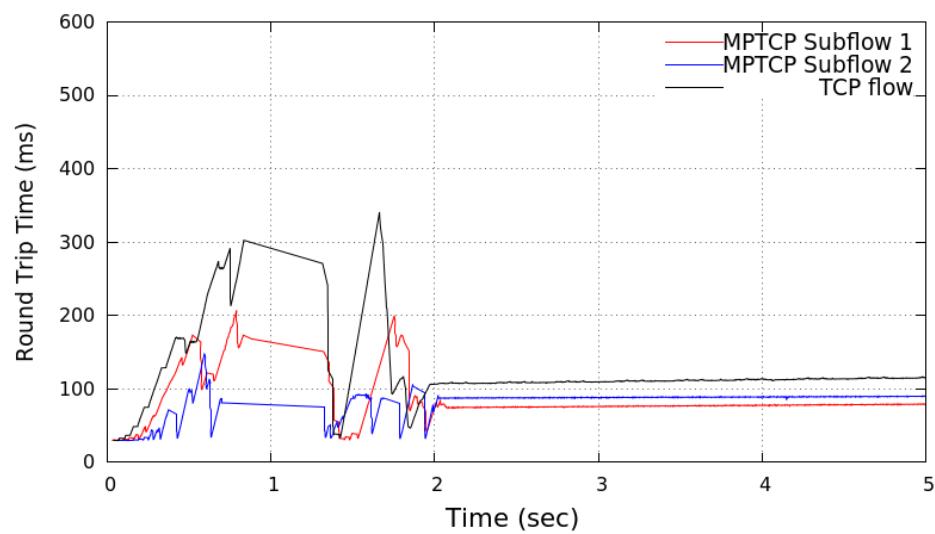
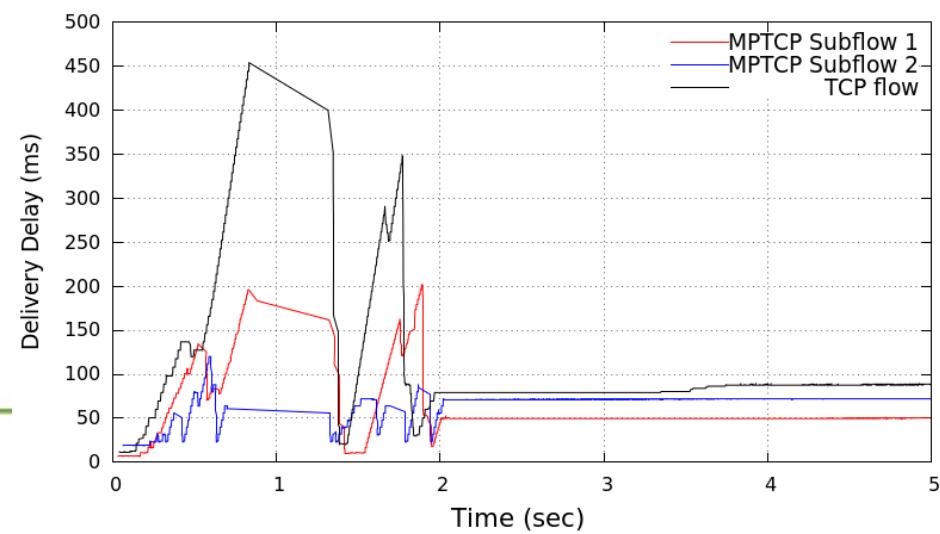
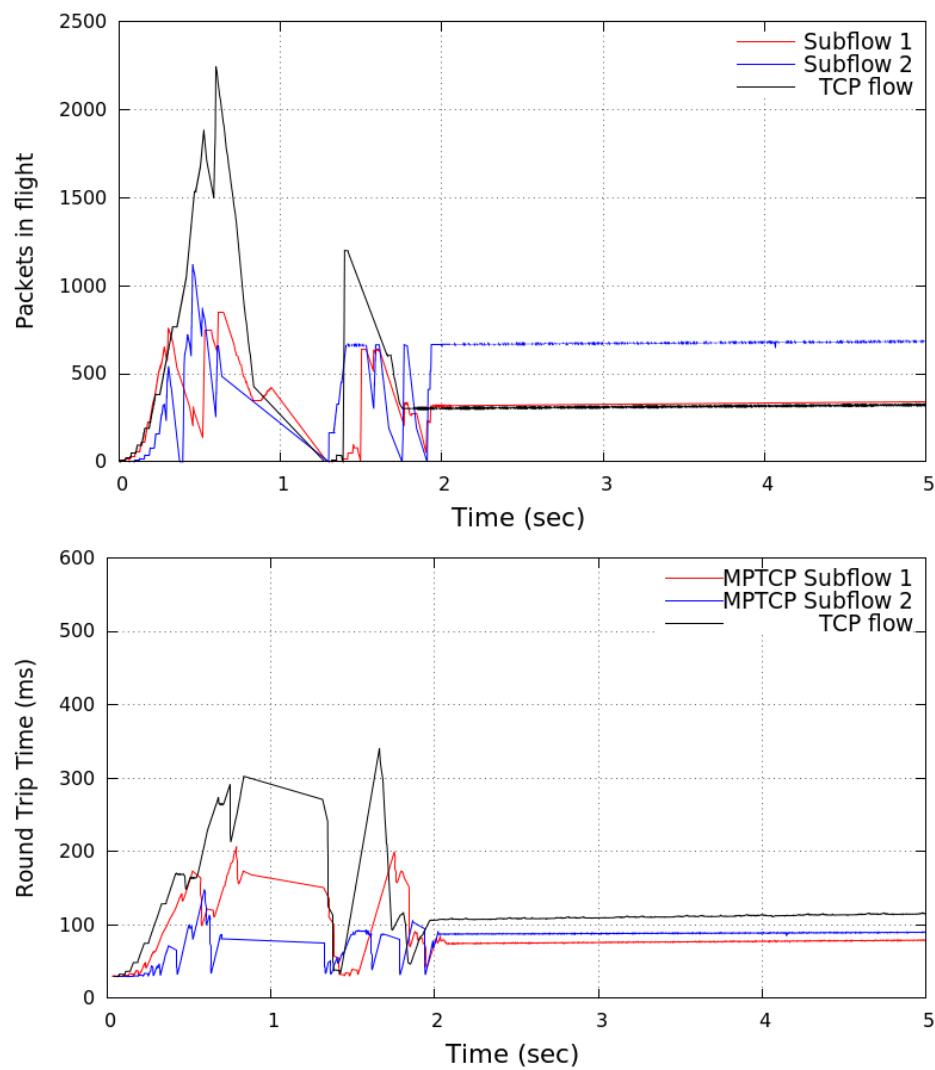
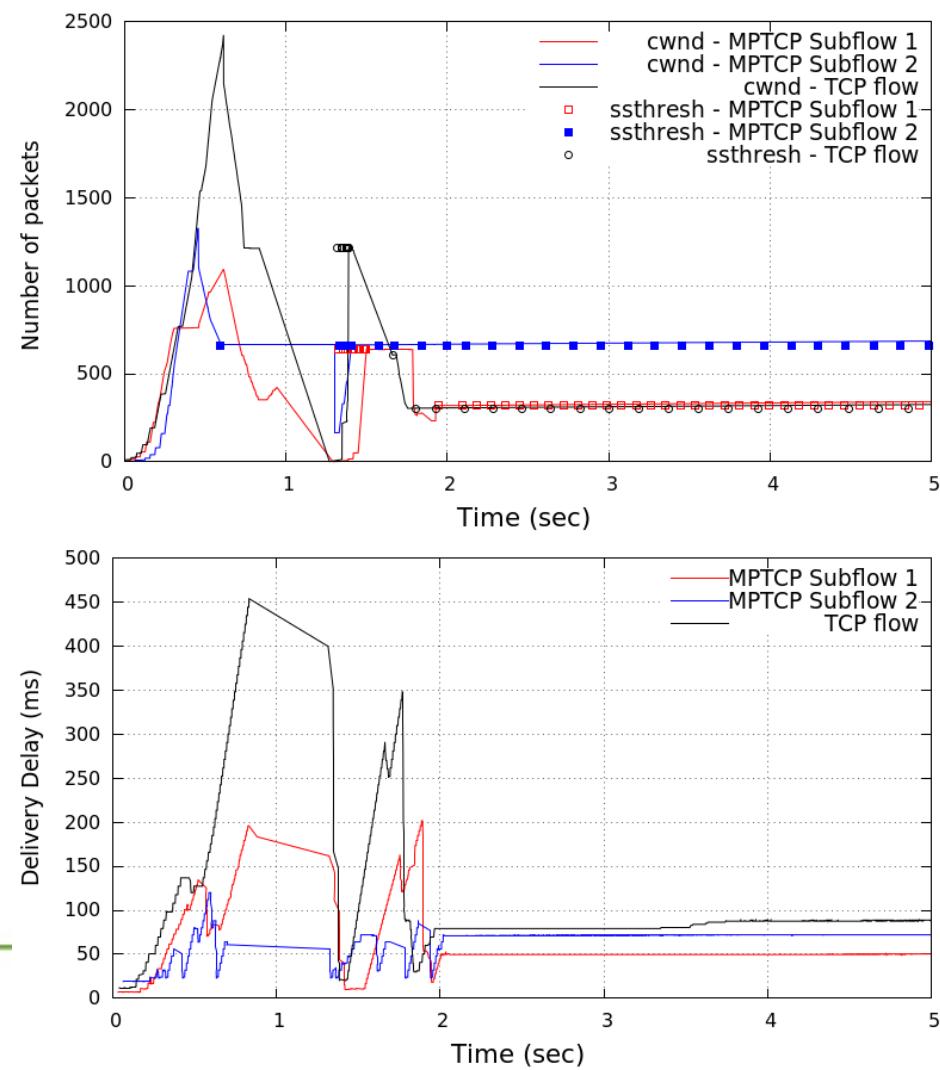
Subflow 1 RTT - 30 ms (10 ms +20 ms)

Subflow 2 RTT - 30 ms (20 ms + 10 ms)

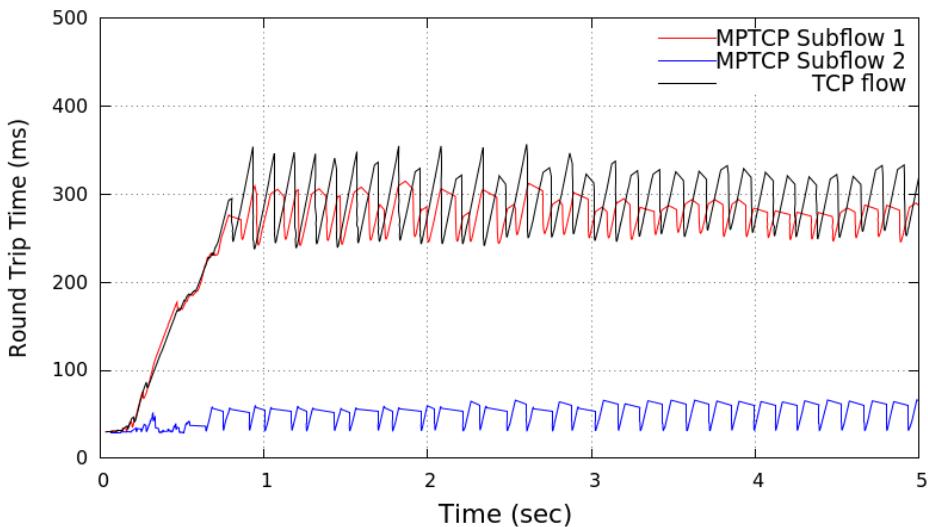
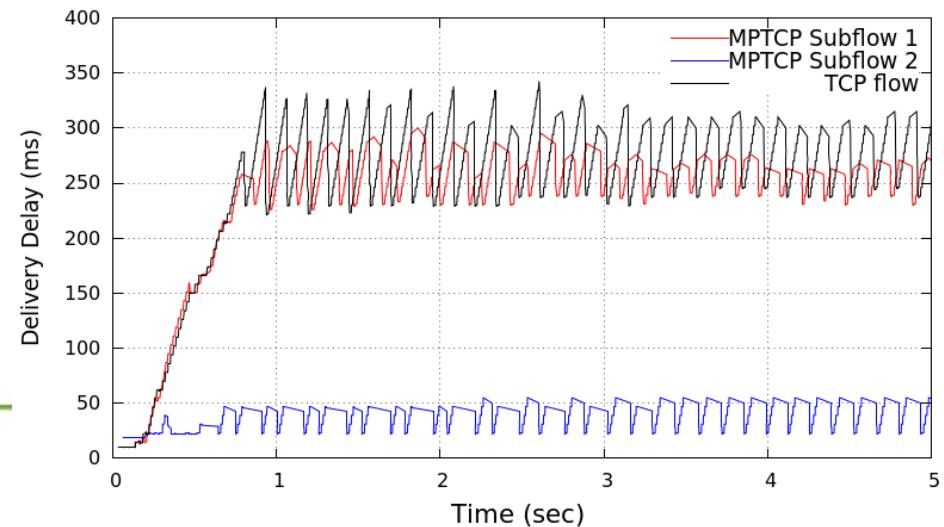
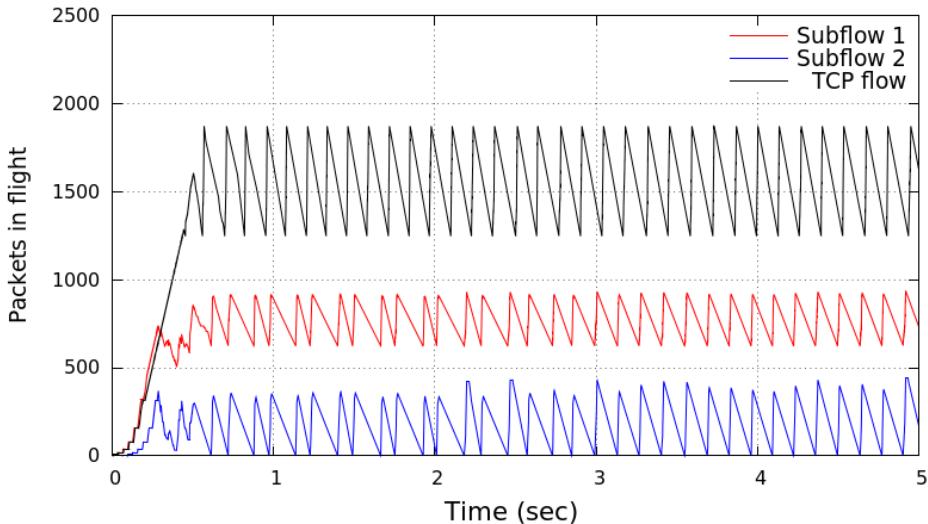
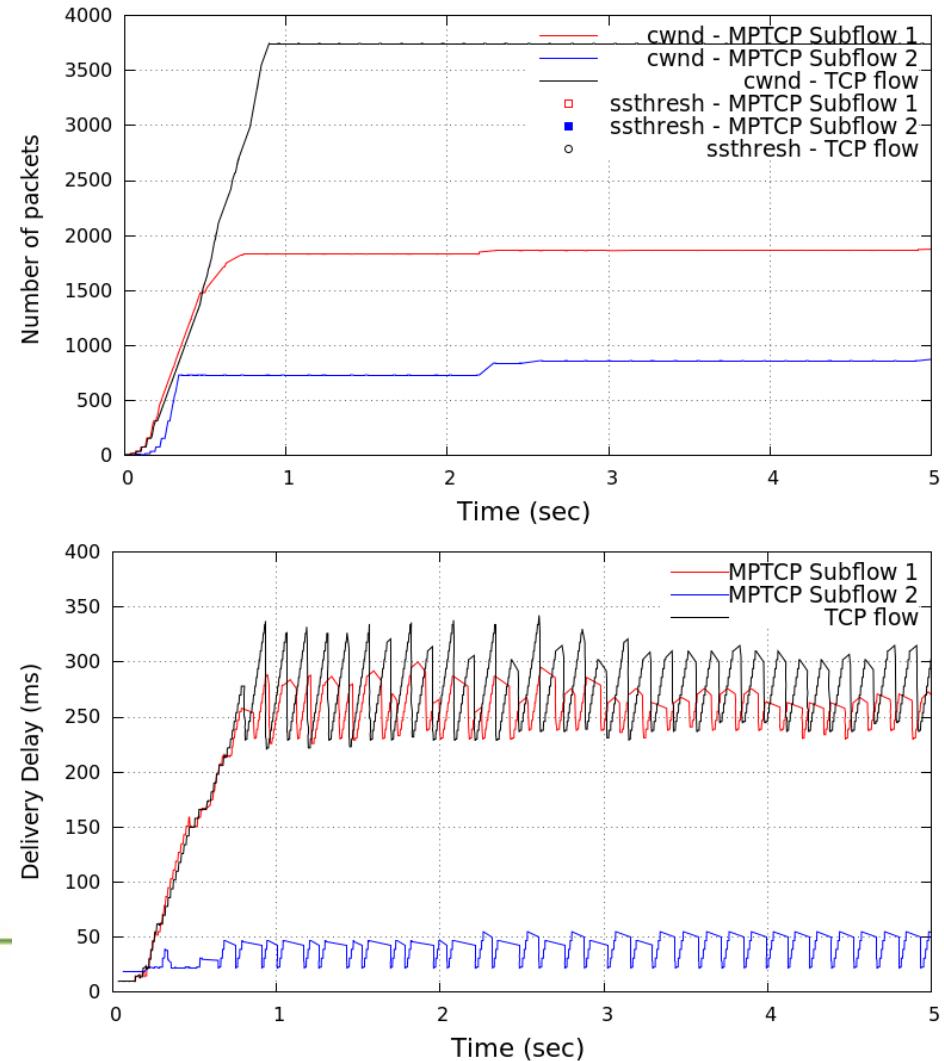
Delivery Delay Scheduler



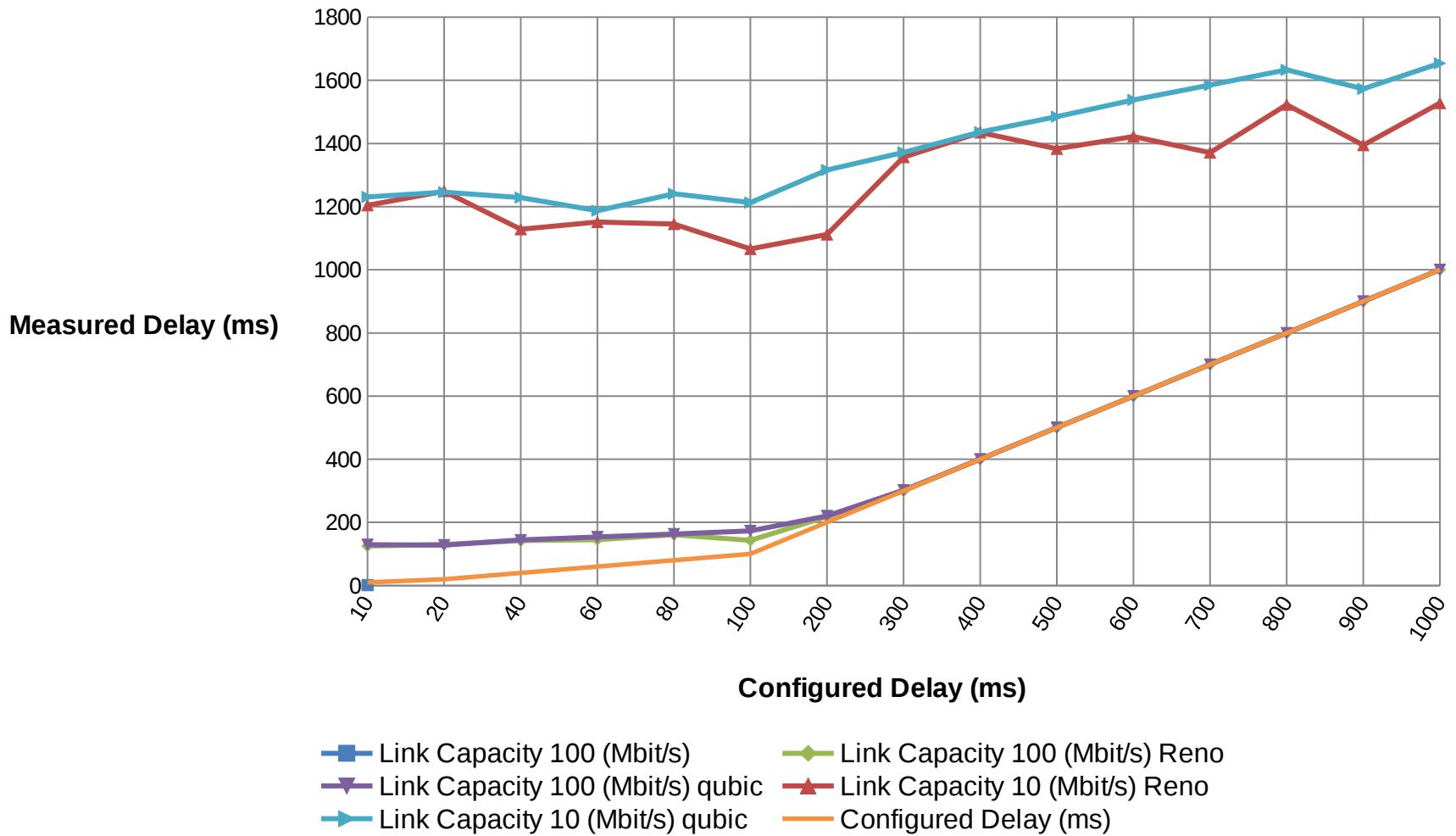
Round Trip Time Scheduler



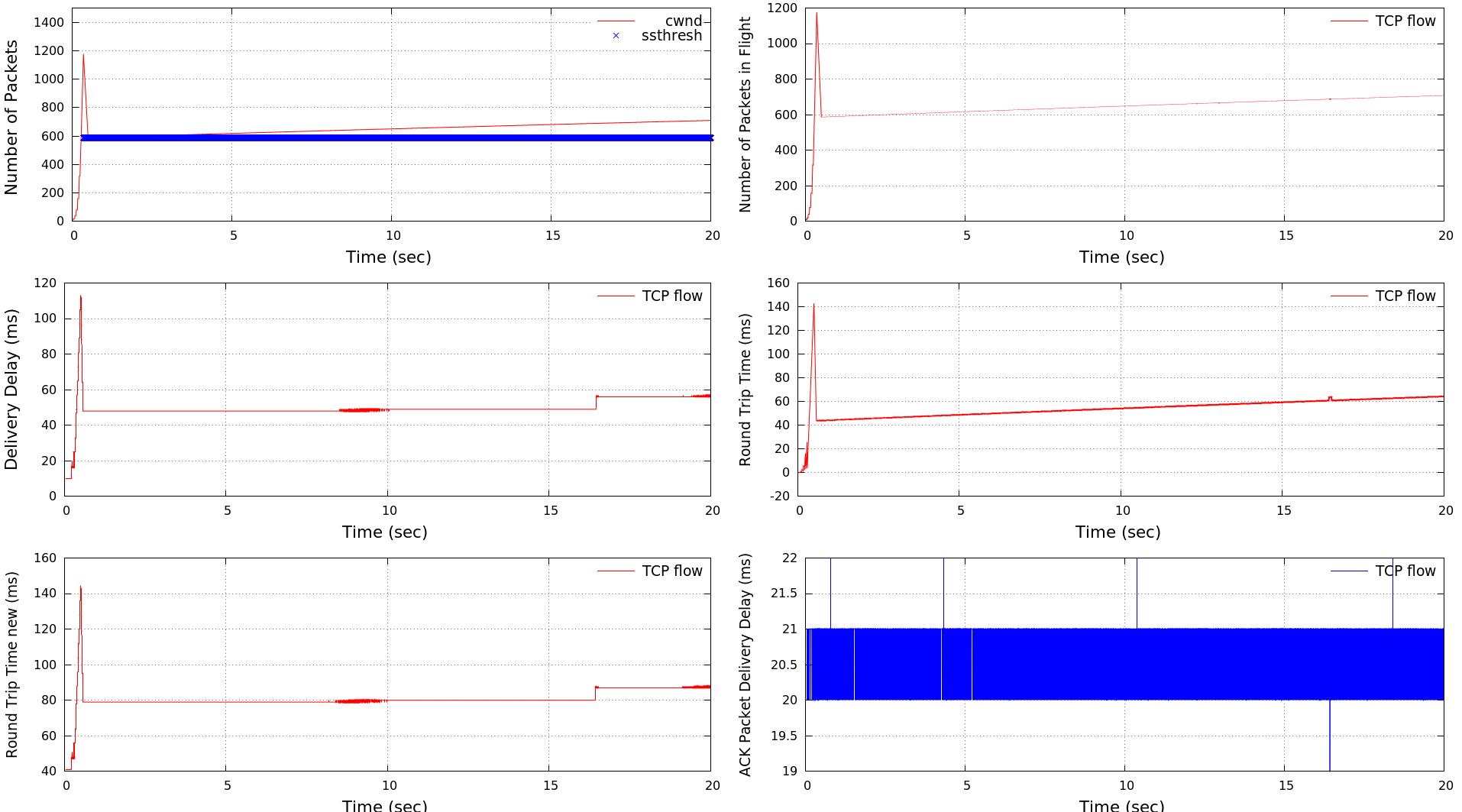
Round Robin Scheduler



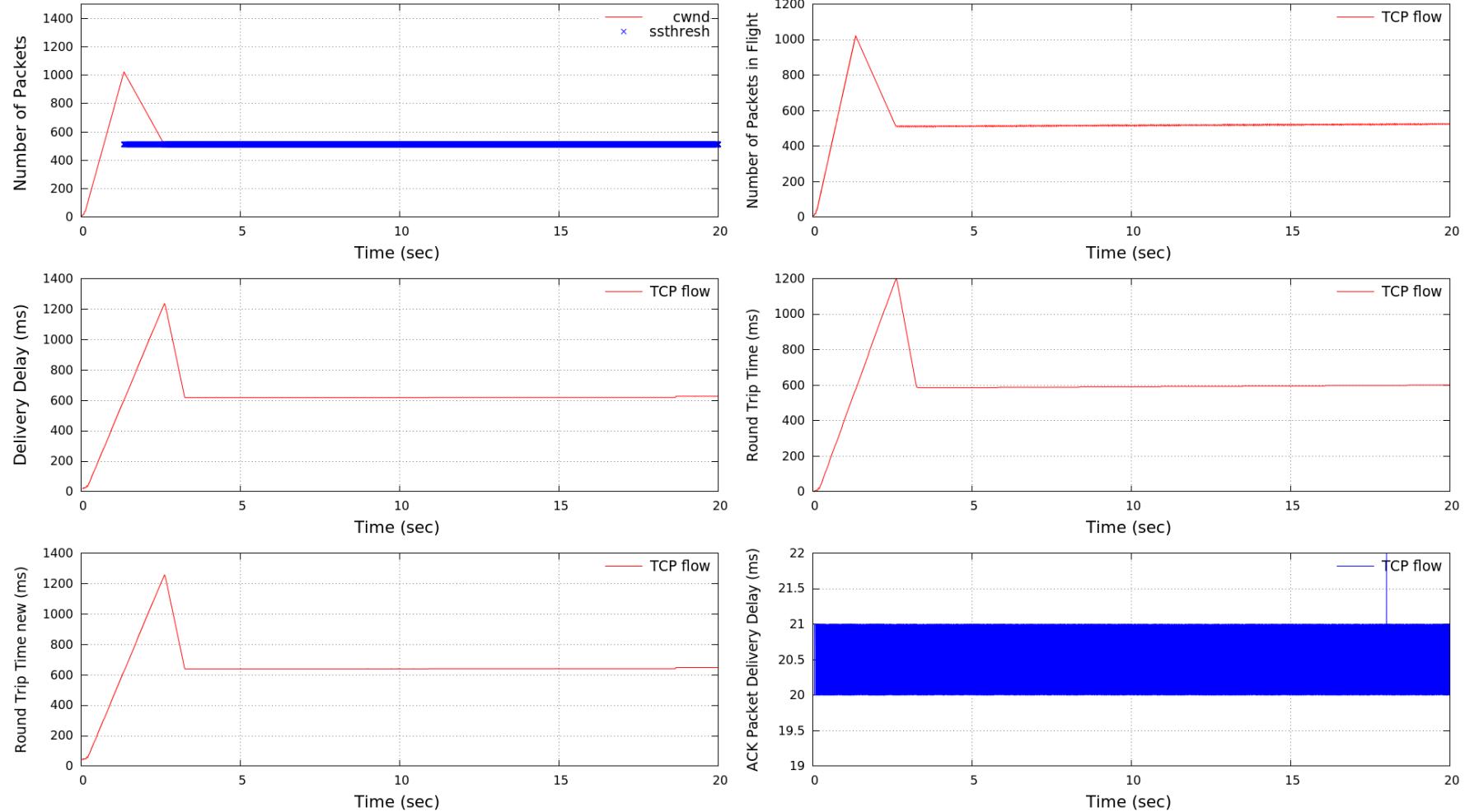
TCP - Measured Delay Vs Configured Delay



TCP – Link Capacity 100 (Mbit/sec), Delay - 40ms (20ms + 20ms)



TCP – Link Capacity 10 (Mbit/sec), Delay - 40ms (20ms + 20ms)



Simple Testbed

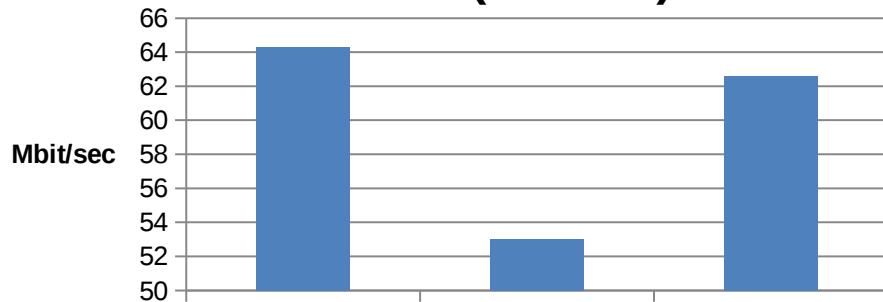
Link Capacity - 100 / 100

(Mbit/sec)

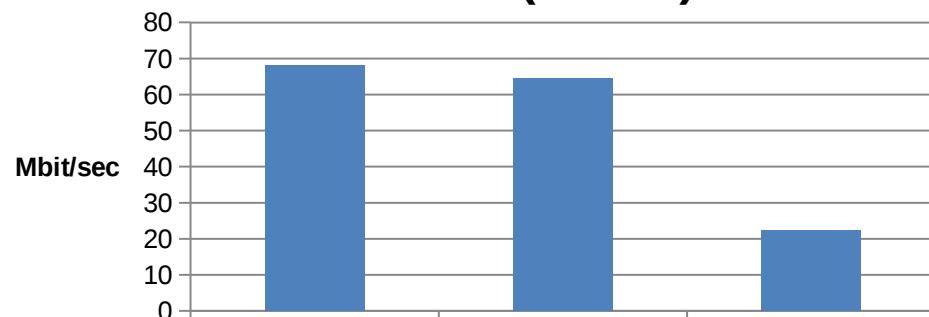
Scenario F

Hypothetical Scenario F

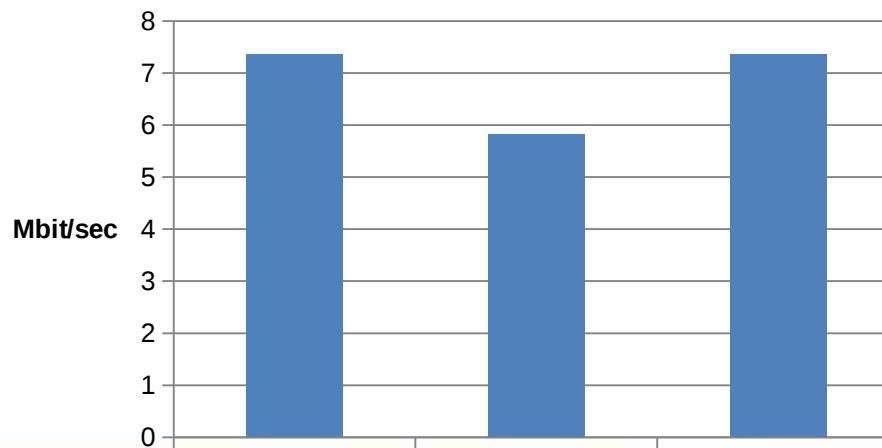
Scenario F (100/100)

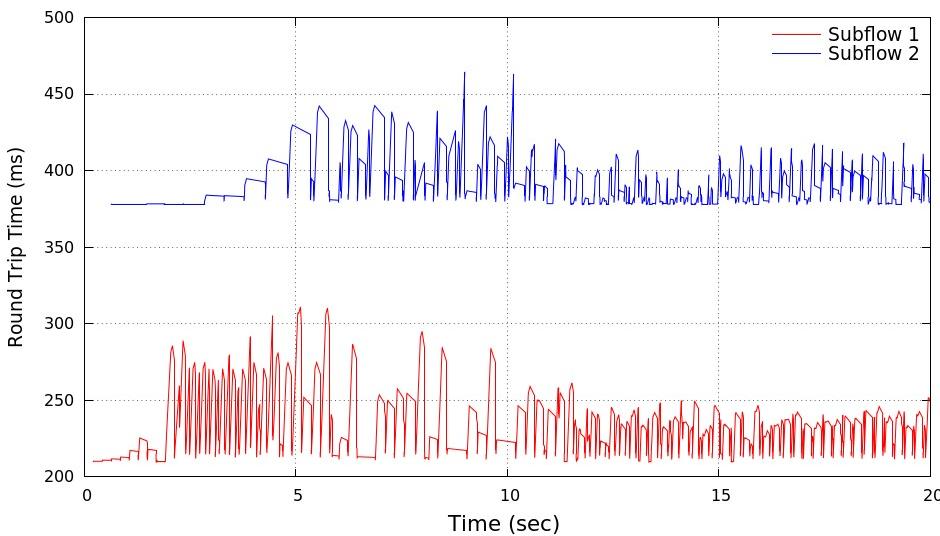
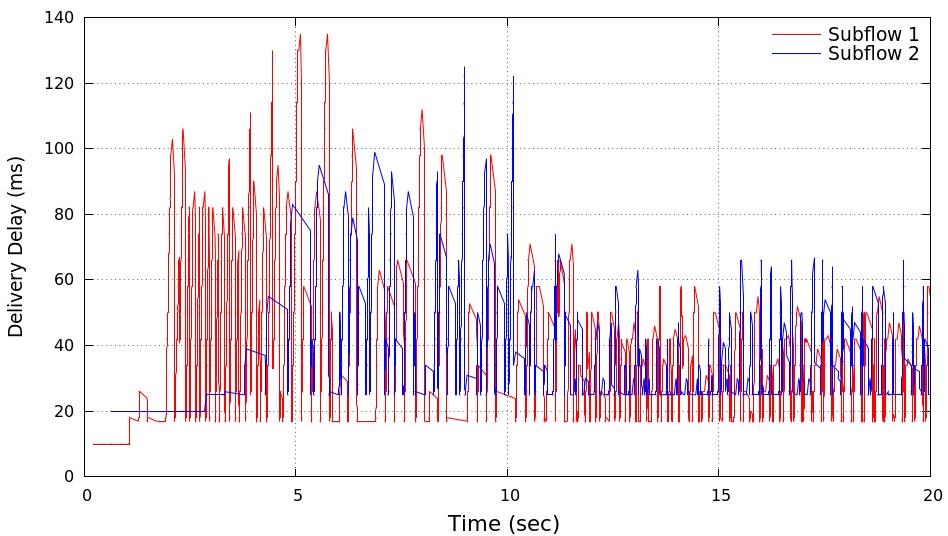
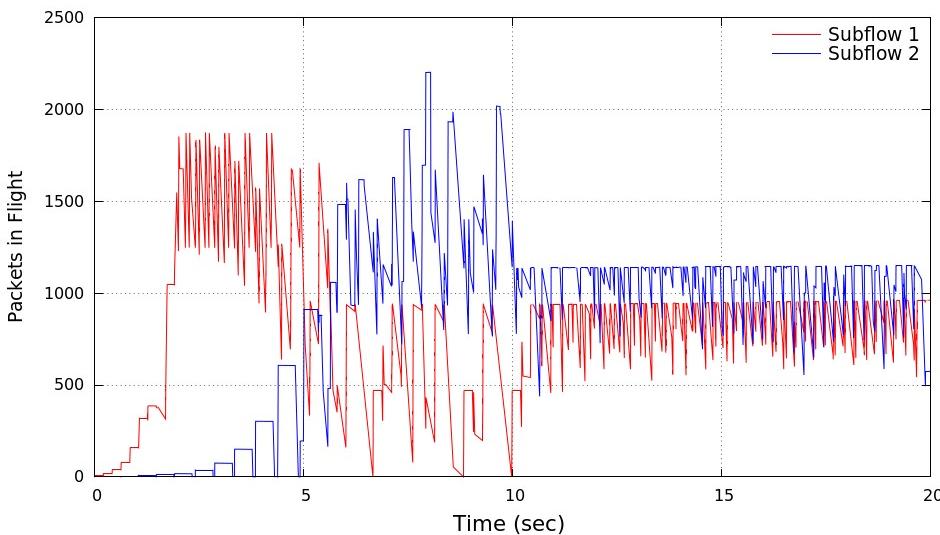
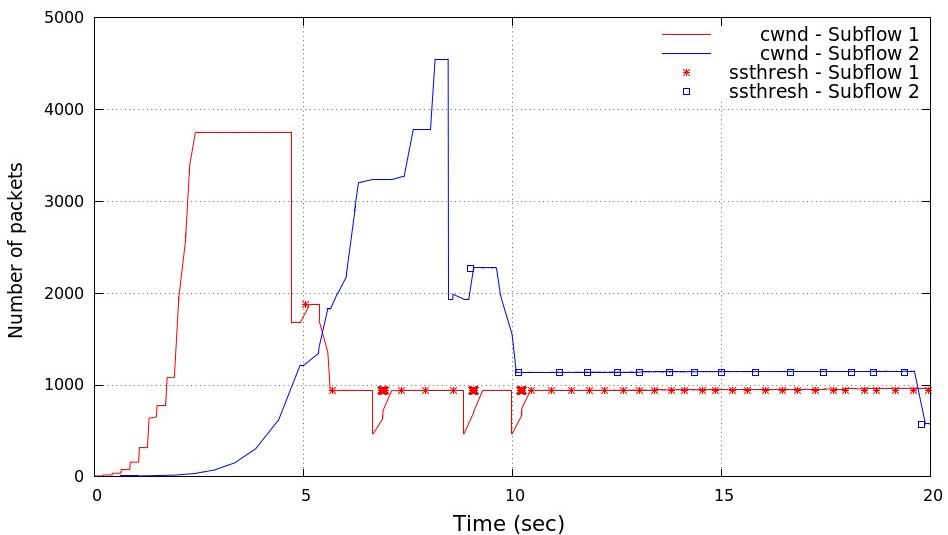


Scenario F (100/10)

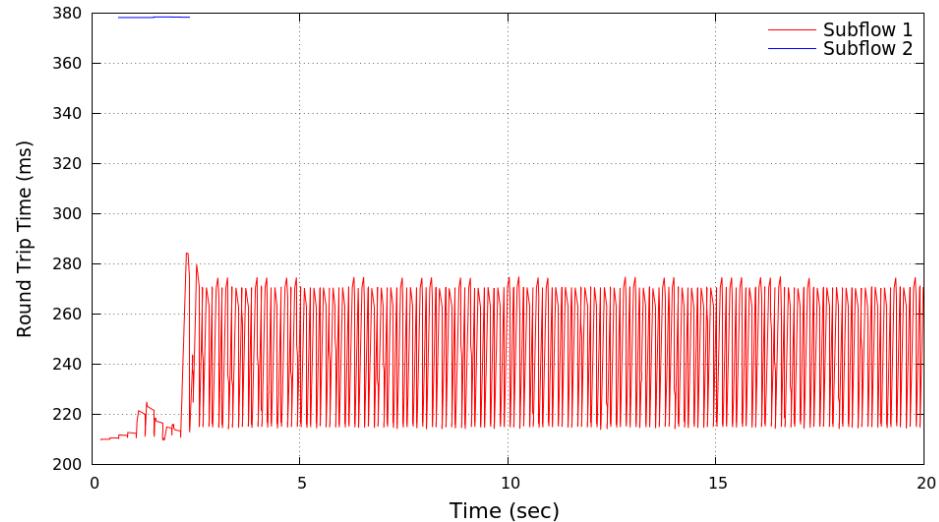
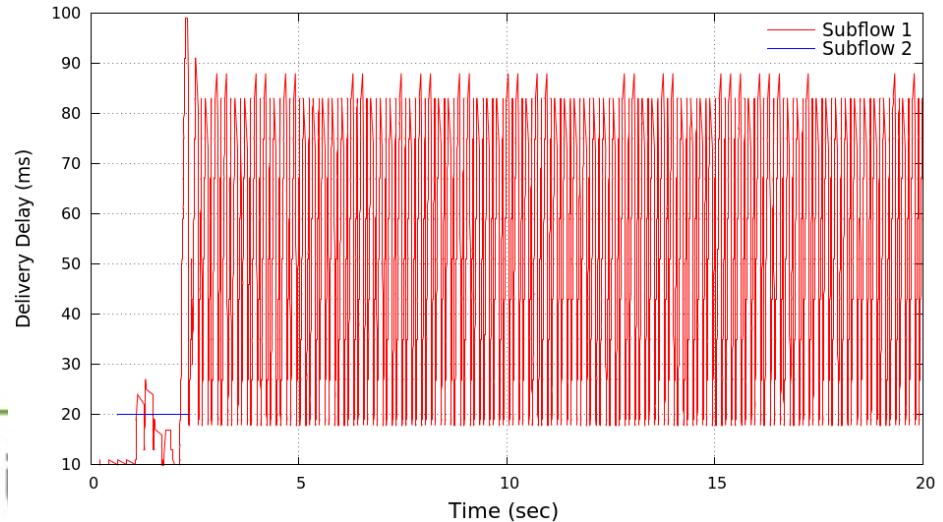
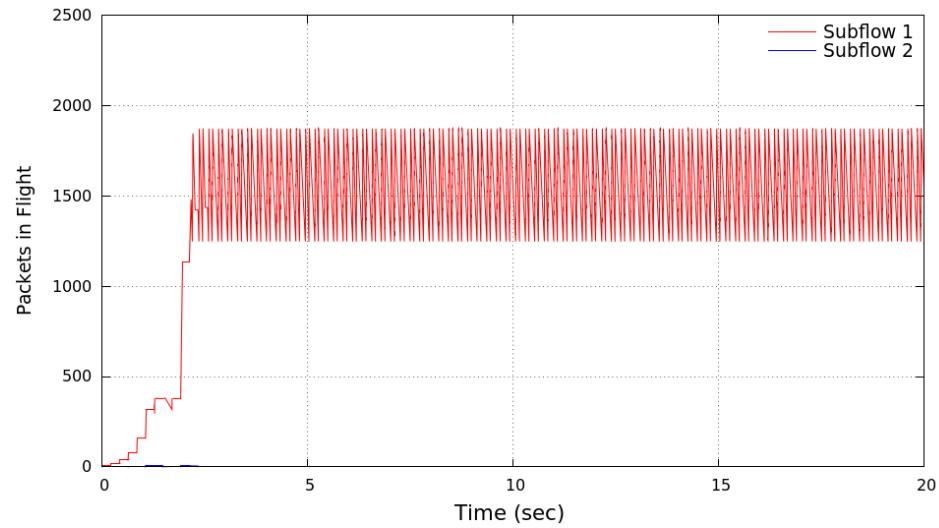
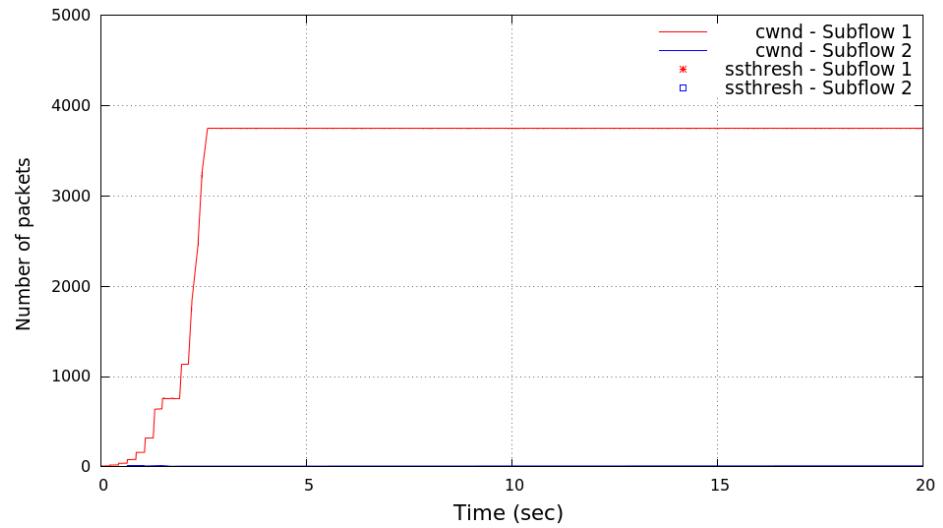


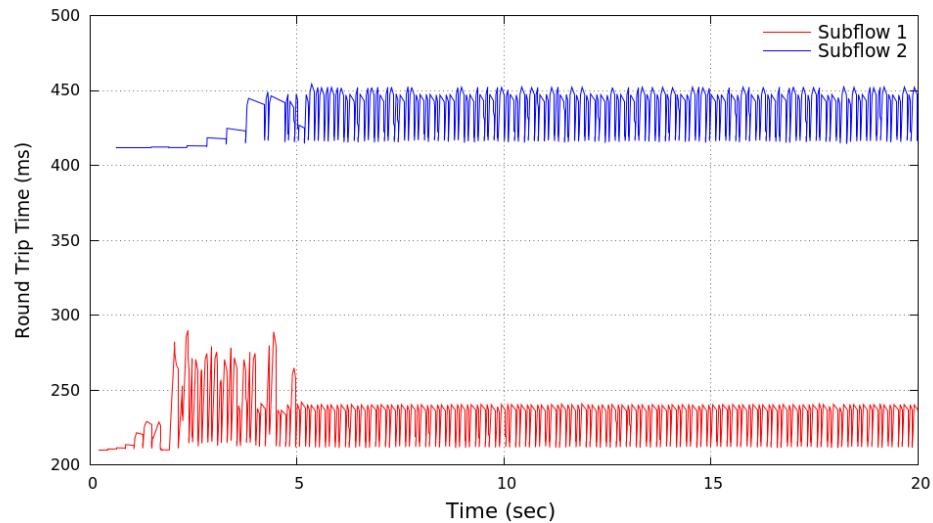
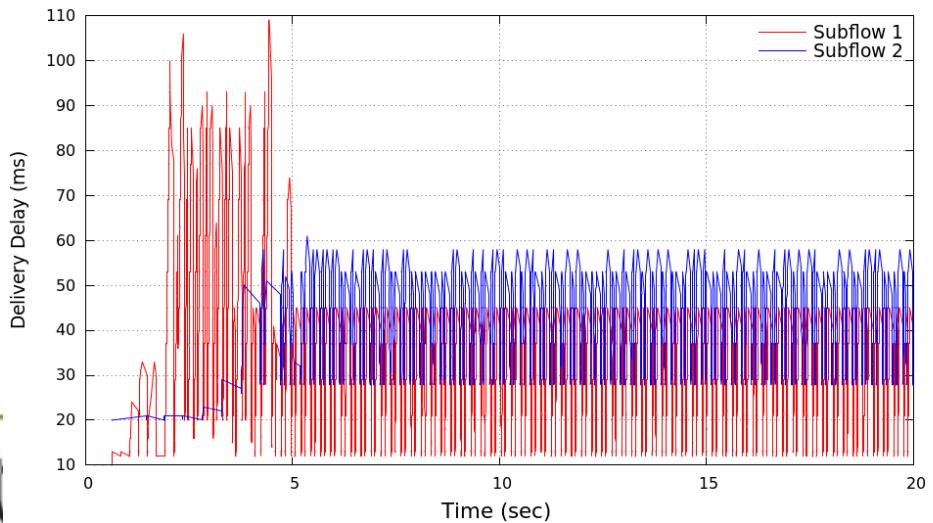
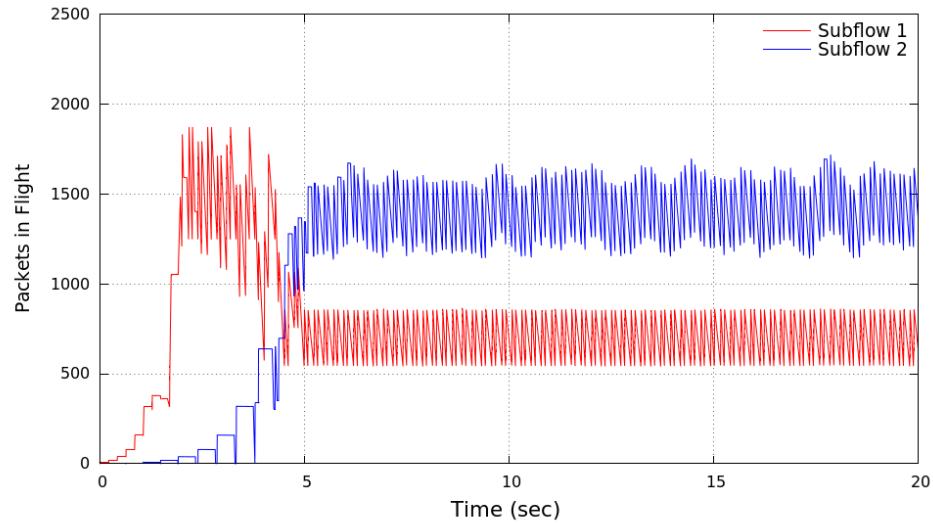
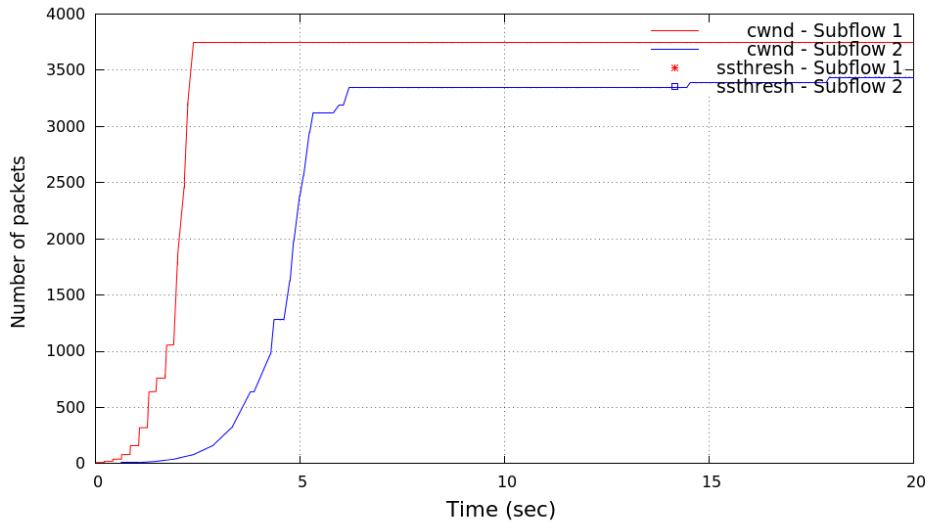
Scenario F – 10/10



DD

sRTT





Outlook (DD Scheduler)

Packet no	Reordering Delay (ms)	End-to-end Delay (ms)
1	0	50
2	0	50
3	0	50
4	0	100
5	0	100
6	0	100
7	0	100
8	0	100
9	40	50 + 40
10	30	50 + 30
11	20	50 + 20
12	0	50

