

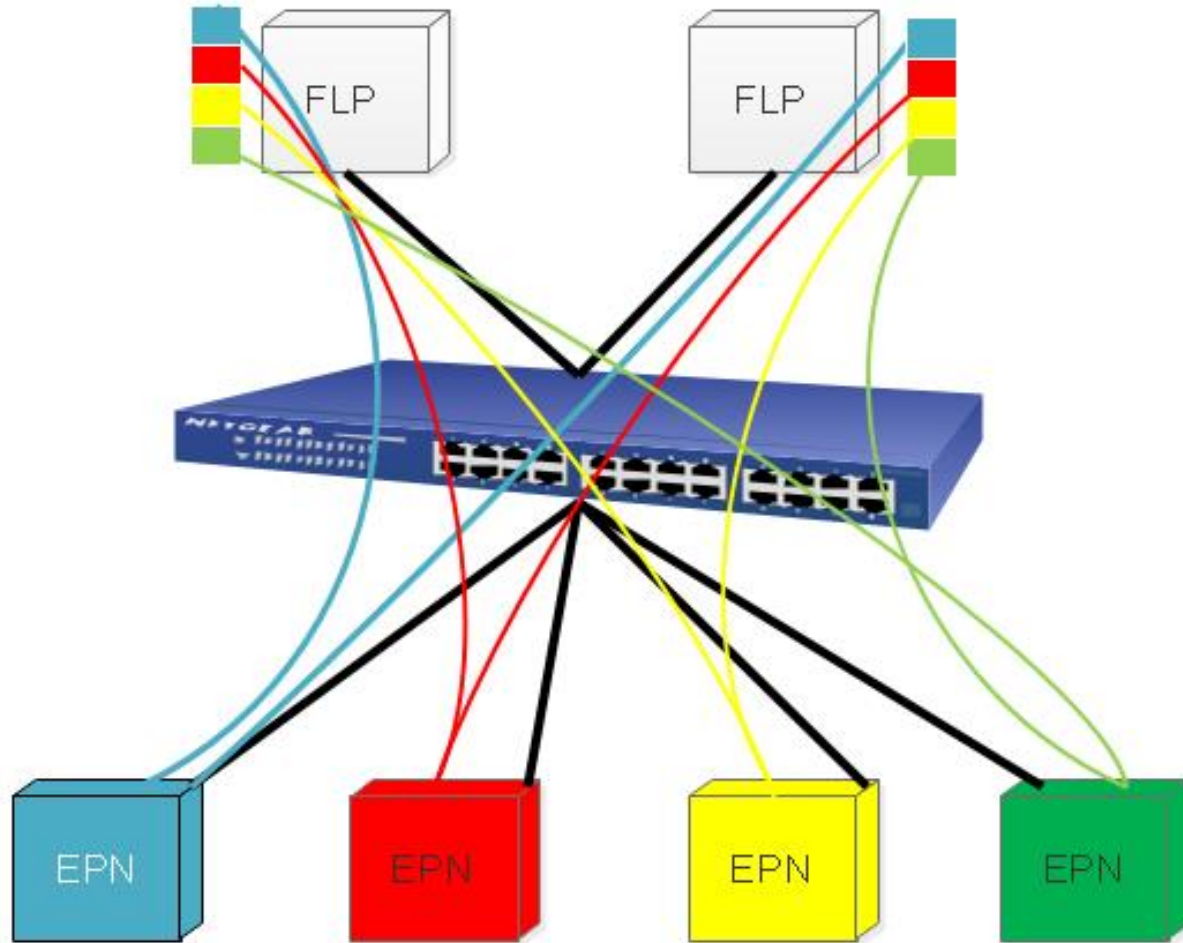
# ALICE O2 WP2 Group

Professor Syed Asad Hussain  
Dean Faculty of Information Sciences and Technology  
COMSATS University Islamabad, Pakistan

# WP2-Data Flow and System Simulation

- Headed by Dr. Iosif Legrand at CERN
- WP2 consists of three deliverables
  - a) Develop a tool for detailed simulation of network topologies for O2 facility for a quick evaluation of multiple scenarios
  - b) Develop a tool to allow for a quick estimate of computing resource needs
  - c) Estimate overall long term computing needs for future runs and Grid resources

# Simulation Scenario



# WP2-Data Flow and System Simulation

## Introduction

- Constitute a system that handles data transfer at high speed of at least 4Tbps
- Basic scenario constitute of 250 FLPs (low processing nodes) and 1500 EPNs (high processing nodes)
- Data Transmission consists of creating full time frame building tasks where the task of building the full time frame is to send the time-frame data chunks from all FLPs to one EPN unit.
- Each FLP is supposed to generate 40MB @ 50Hz (every 0.02 seconds)
- Total input rate at a FLP is ~500Gb/s and the average bisection traffic through the switching fabric is ~4Tb/s
- EPN Input speed required 400 logical connections to send all 4Tb to one EPN

# Introduction (Cont..)

- The link rate between FLP and switch is 40Gbps and from switch to EPN is 10Gbps
- TCP is being used as the primary transport protocol
- Simulation at COMSATS Lahore lab is done using NS-3 due to its flexibility and it being an open source

# Some Challenges

The main task which we are handling is (The Full-Time Frame Building), figure # 1 shows the schematic views of the data flow to create complete time frames on the EPN nodes. Following are the main challenges:

- To be able to keep up with the high input rate.
- ~ 250 FLP processing units that collect sub-detector time frames and using a control algorithm send the data to assemble the full-time frame in one EPN node. (we need to avoid congestion on the switch) (TCP throttling).
- Scalability issue (need to design a system that is scalable).

# Network Simulation Tools

- NS3 is being used at Lahore Campus and OMNet++ at Islamabad campus.
- A complete customized transfer helper is developed to meet the requirement.
- Maximum MTU, Jumbo Frames and Maximum Packets size is used to test the given scenario at full load.
- Helpers are created in Java, C++ and even Python to see the performance variations on these different platforms.
- High Computing Machines and Servers are dedicated for this task

# Simulation Parameters for NS3

Name	Value (TCP)
No of FLPs	1, 4, 25, 60, 70, 250
No of EPNs	6,24,150,360,420, 1500
Input Rate at FLP	40MB in 50Hz
FLP-Switch Link Rate	40Gbps
Switch-EPN Link Rate	10Gbps
Simulation Time	4 hours
Simulator	NS-3



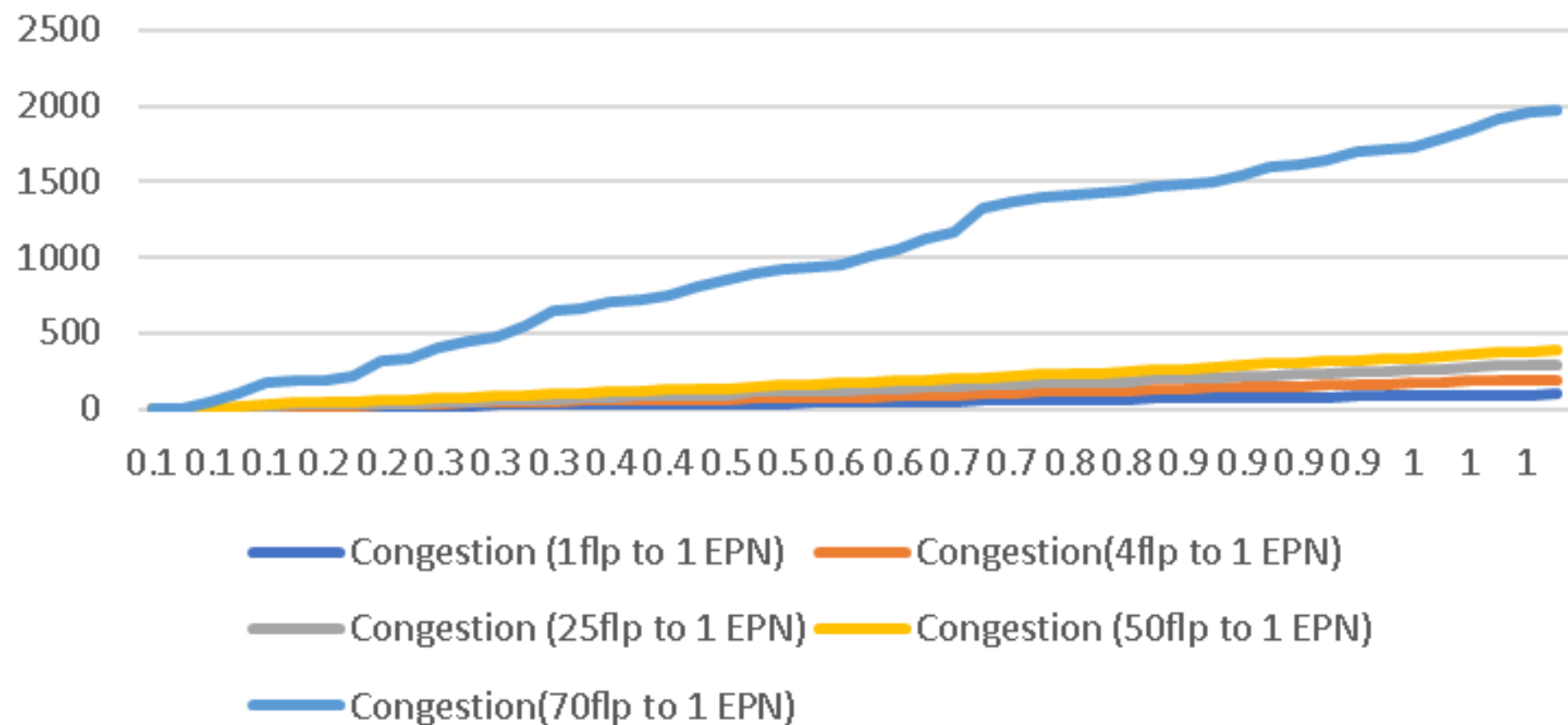
# Problems faced so far

- Each node/EPN based on its characteristics and input rate accepts only 130Gb of data without congestion and retransmission attempts.
- Using different TCP variants, the congestion and retransmission has reduced the simulation time from 24 hours to 5 mins for 5 FLPs to 30 EPNs.
- A full scale system (250 FLPs and 1500 EPNs) in default configuration of NS3 took more than 5 days and we reduced this time to 24 hours by developing custom helper and than 4 hours.
- Default TCP congestion mechanism are not suitable for simulating this full scale system.

# Results (Default TCP configuration Mechanism) (initially)

No. of FLPs	Throughput at each EPN(Gbps)	Time taken for full system run	Congestion (Mbps)
1	10	2s	97
4	10	1 min	100
25	9.99	5 mins	98
60	9.99	8 mins	97
70	9.28	45 mins	1600

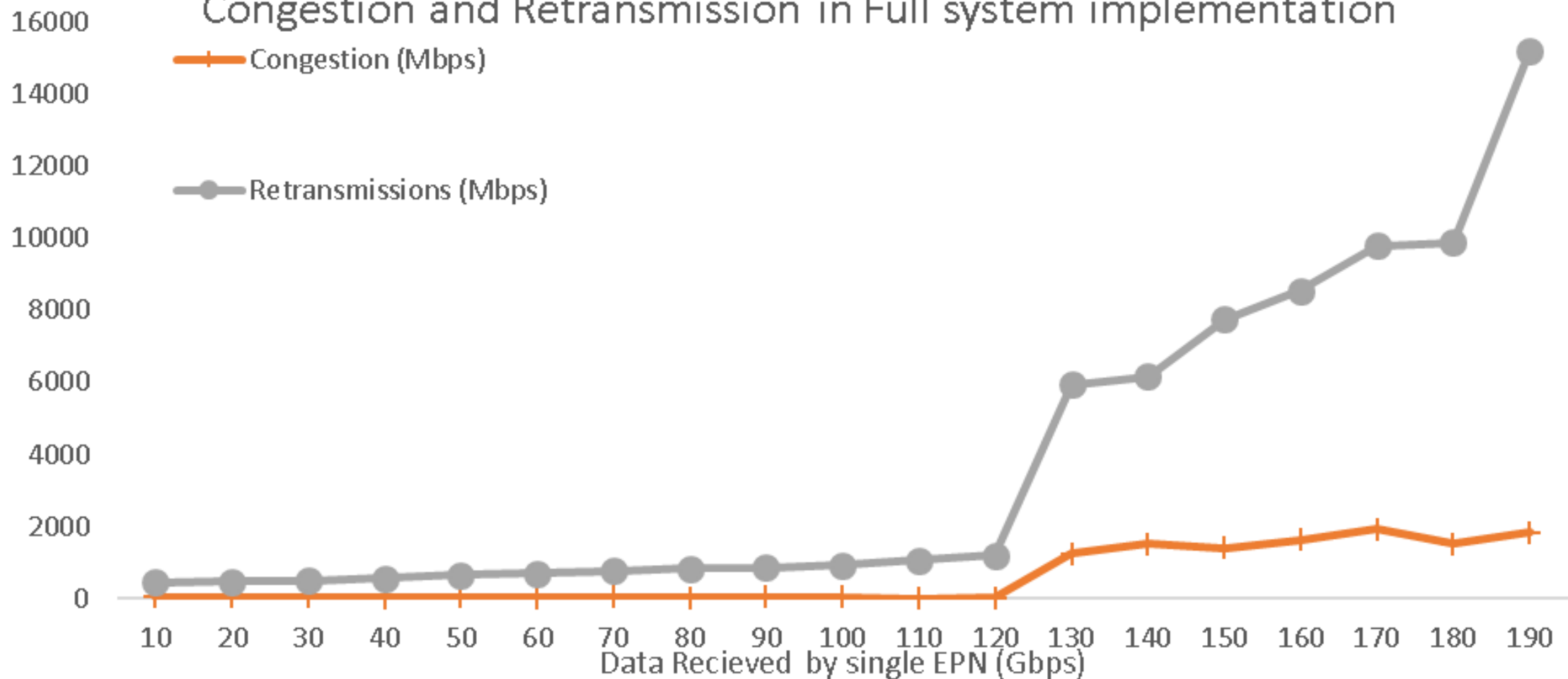
## Congestion at each EPN from number of FLPs in Mbps



### Congestion and Retransmission in Full system implementation

Congestion (Mbps)

Retransmissions (Mbps)



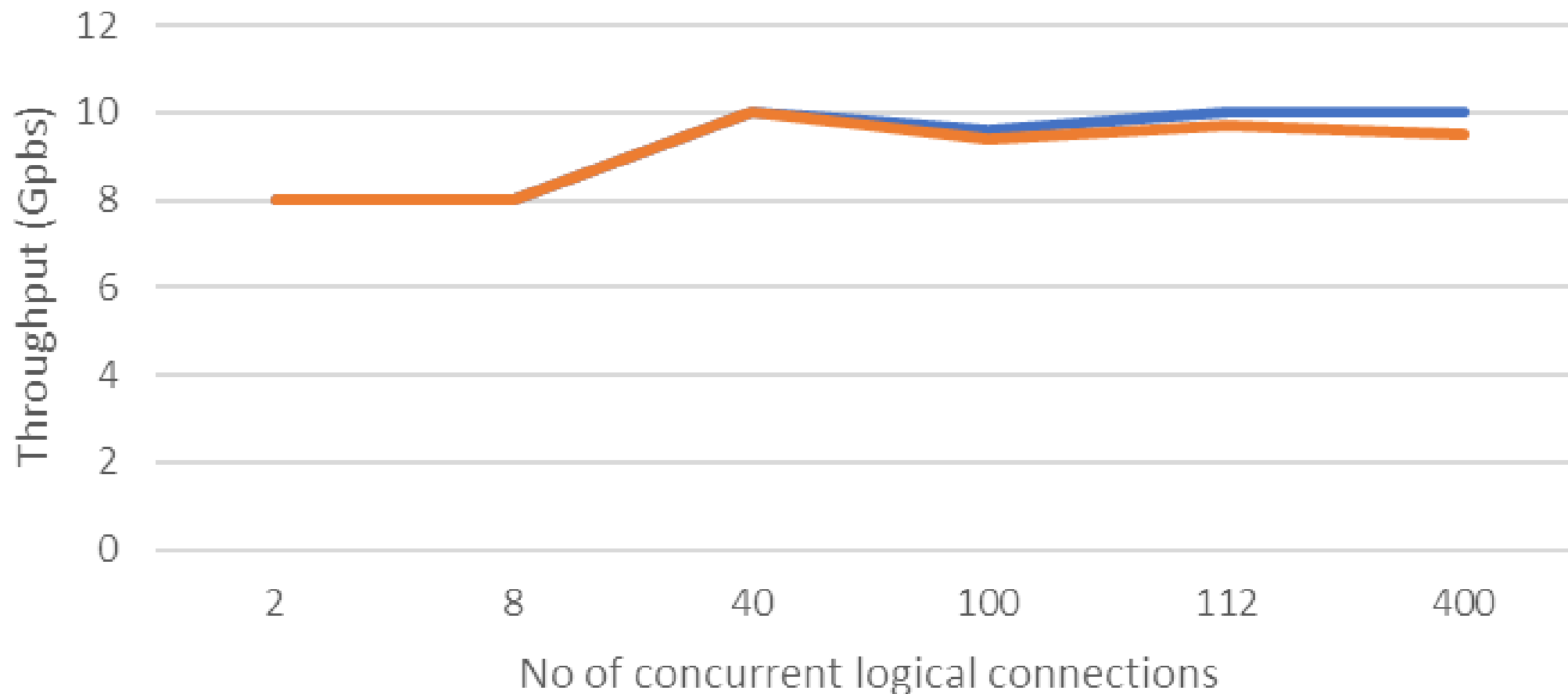
# Current Progress

- We are working on different congestion management scenarios in NS3 and looking at custom options too.
- We have successfully tested NS3 to its limits and are now in process of modifying it to full extent to implement the full system i.e. Highest MTU, Largest Packet size and Jumbo frames.
- The full scale system with only 4% retransmissions has been implemented in completion time of **4 hours** generating 4Mb congestion.
- These simulation results are 0.004% different from actual system implemented by Dr. Iosef. Legrand at CERN.

# Results

No. of FLPs	No of Logical Connections	Throughput at each EPN(Gbps)	Time taken for full system run	Congestion (kbps)
1	2	8	2s	23
4	8	8	1 min	40
25	40	10	5 mins	120
60	100	9.4	8 mins	415
70	112	9.7	45 mins	1200
250	400	9.5	4 hours	4000

# Throughput desired throughput achieved



Throughput desired (Gpbs)      Throughput achieved (Gpbs)

# Summary

- The full scale system with 250 FLPs and 1500 EPNs have successfully been implemented using NS3
- The full scale simulation took 4 hours to complete
- The congestion has been reduced from 1.6 Gb to 4Mb
- Number of retransmission attempts have been reduced.



# Simulation setup in OmNet++

*Table 1: Simulation parameters for different network scale size*

<b>Parameters</b>	<b>Architecture I</b>	<b>Architecture II</b>	<b>Architecture III</b>
<b>Number of FLPs &amp; EPNs</b>	3 FLPs & 18 EPNs	50 FLPs & 300 EPNs	250 FLPs & 1500 EPNs
<b>Number of switches</b>	Single switch		
<b>Packet size</b>	0.4 GB		
<b>Link capacity</b>	From FLP to Switch 10 Gbps From Switch to EPN 10 Gbps	From FLP to Switch 40 Gbps From Switch to EPN 10 Gbps	
<b>Traffic generation</b>	0.02 seconds		
<b>Simulation time</b>	~2.0 seconds	~1.0 seconds	~3.0 seconds

# Simulation Time Vs Real Time

*Table 2: Simulation Time Vs real time*

<b>Architecture</b>	<b>Real time taken</b>	<b>Simulation Time</b>
<b>Architecture I</b>	4 – 5 hours	~2 sec
<b>Architecture II</b>	~ 25 hours	~1 sec
<b>Architecture III</b>	+ 35 hours	~3 sec

# Output file size

*Table 3: Output (Result) file size of different Versions*

<b>Architecture</b>	<b>Size of Output file (before-Old) (Using GUI)</b>	<b>Size of Output file (after-new) (Using Cmd)</b>
<b>Version I</b>	>300 GB	-
<b>Version II</b>	>350 GB	10 MB approx.
<b>Version III</b>	>500 GB	13 MB approx.

# Simulation Results (Throughput)

Table 4: Throughput of different network scale size (*SINGLE EPN*)

<b>Architecture (FLP-EPN)</b>	<b>Architecture I (3-18)</b>	<b>Architecture II (50-300)</b>	<b>Architecture III (250-1500)</b>
<b>Throughput</b>	0.191 Gbps	3.9 Gbps	9.0 Gbps

# Throughput (3 FLPs-18 EPNs)

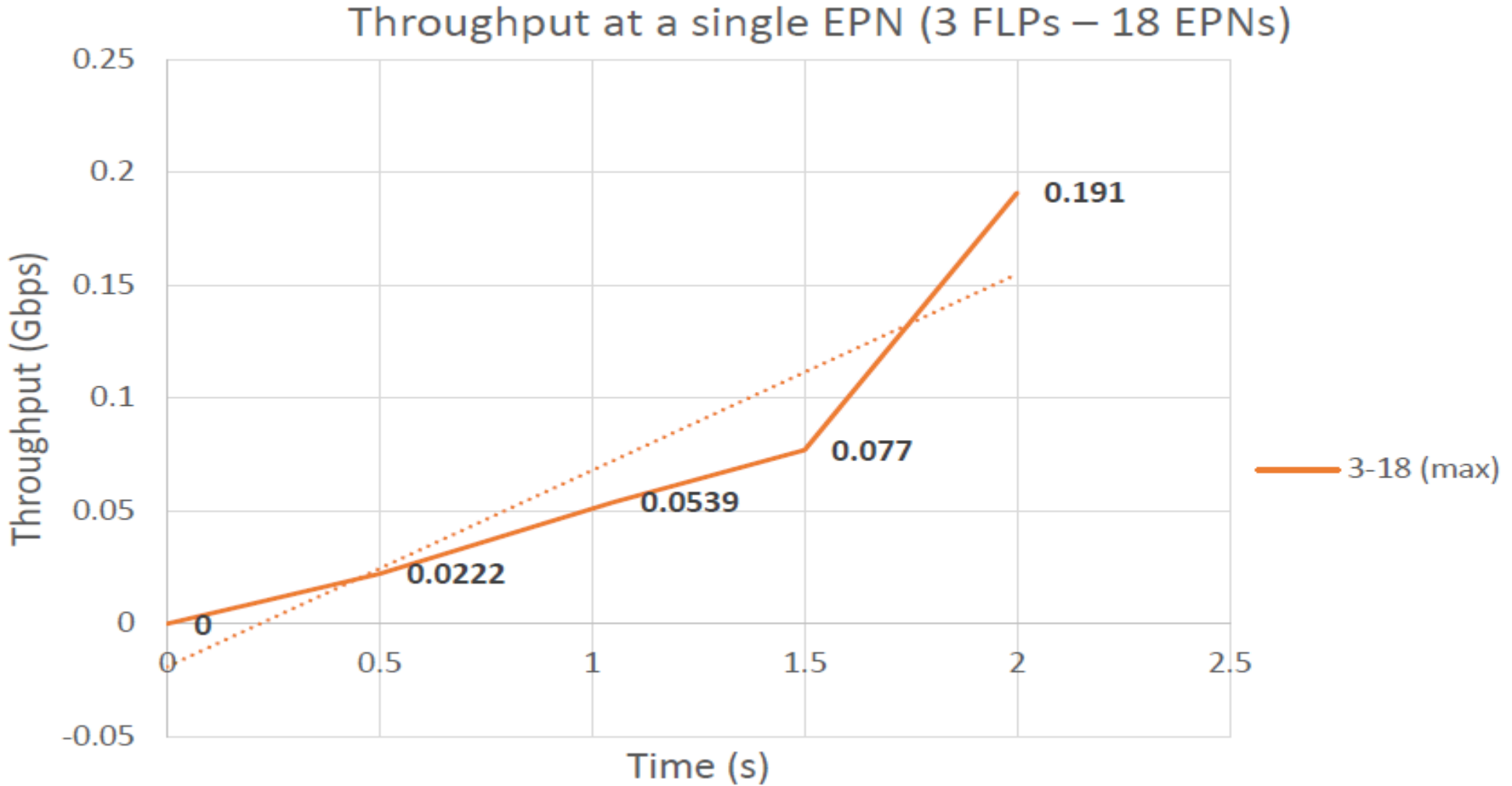


Figure 2: Throughput for network size 3Flps – 18Epn

# Throughput (50 FLPs-300 EPNs)

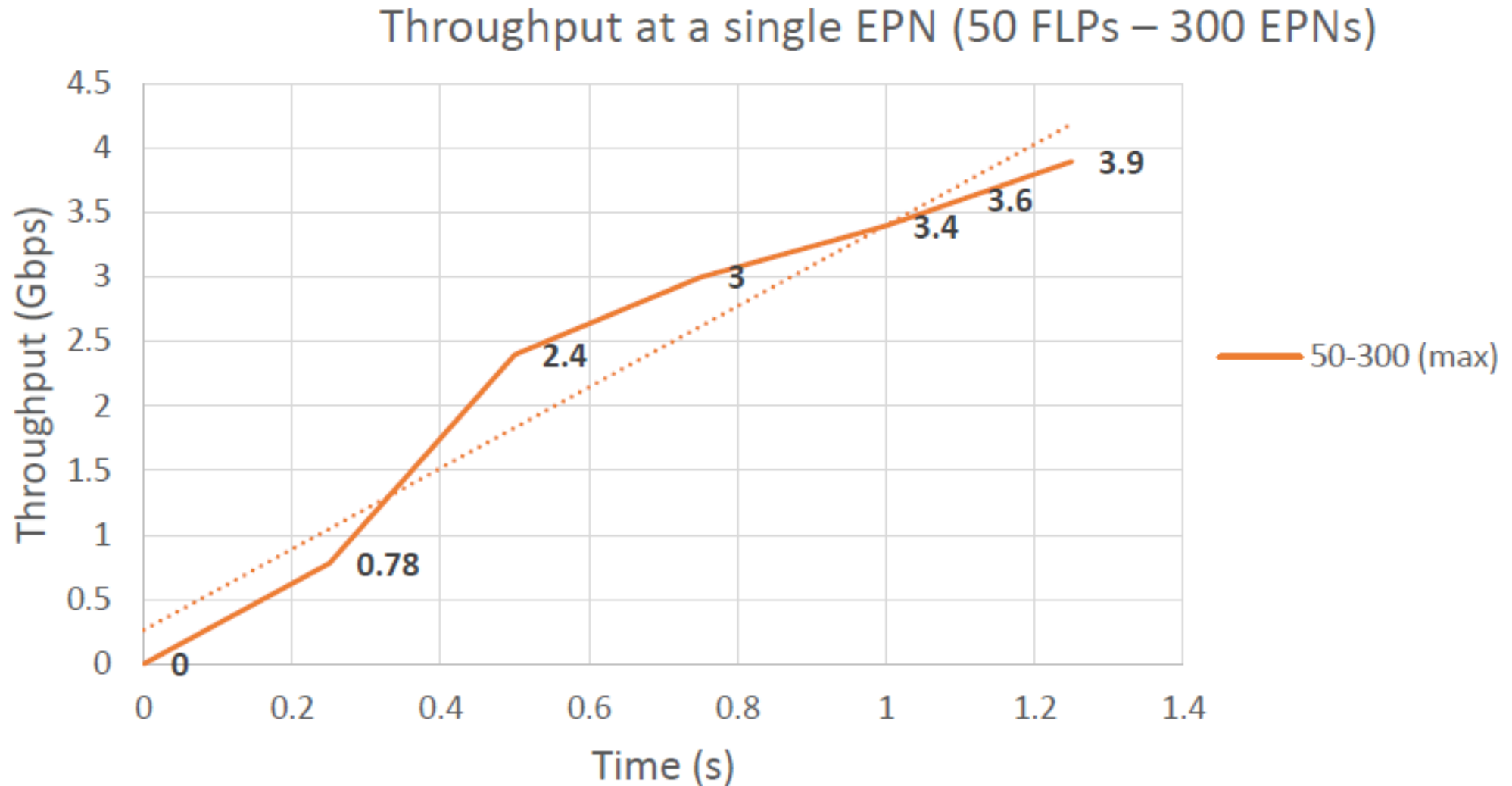


Figure 3: Throughput for network size 50Flps – 300Epnns

# Throughput (250 FLPs-1500 EPNs)

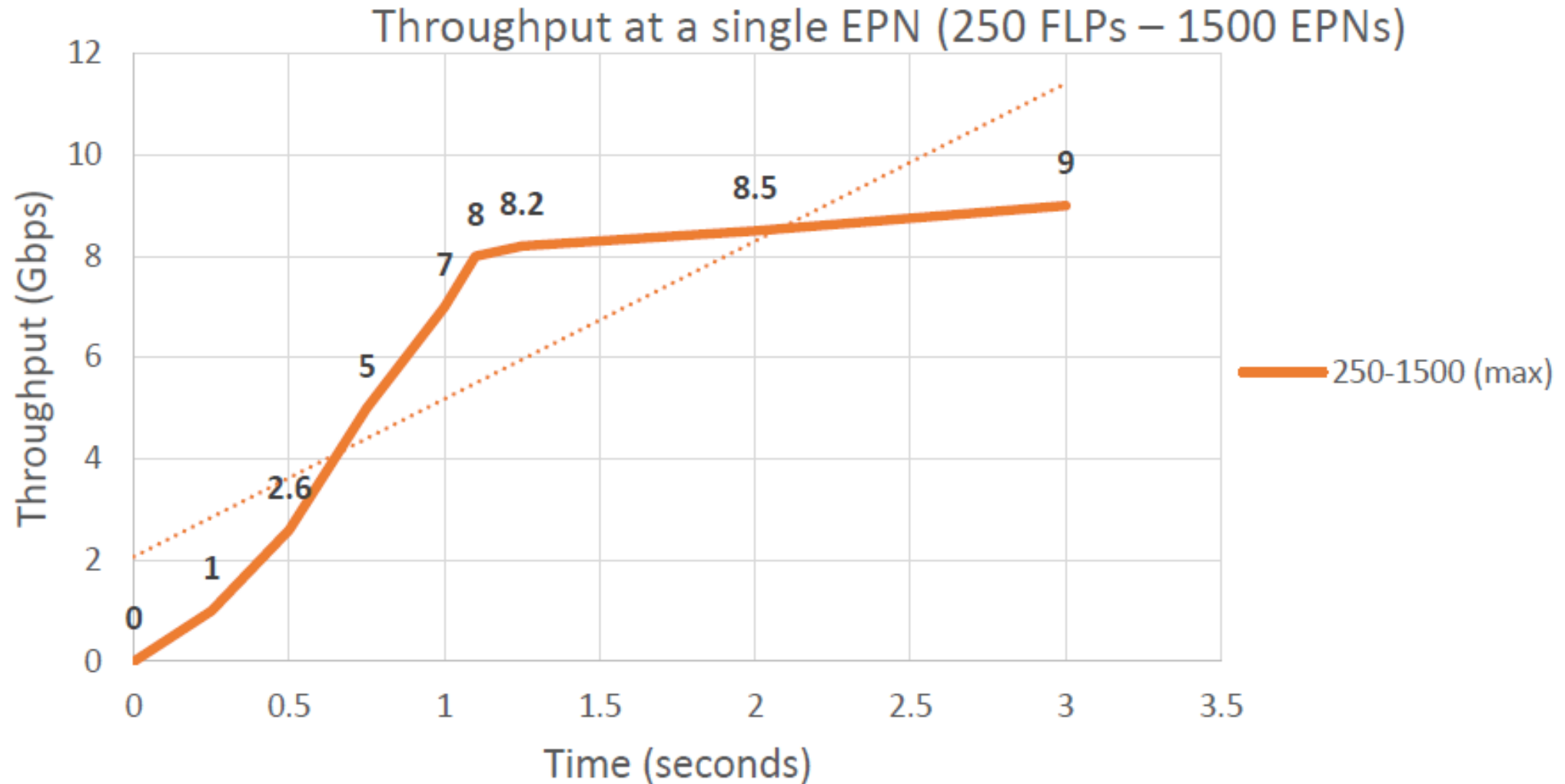


Figure 4: Throughput for network size 250Flps – 1500Epn

More Details:

**Architecture(s) proposed by COMSATS Islamabad.**





# Congestion control architecture

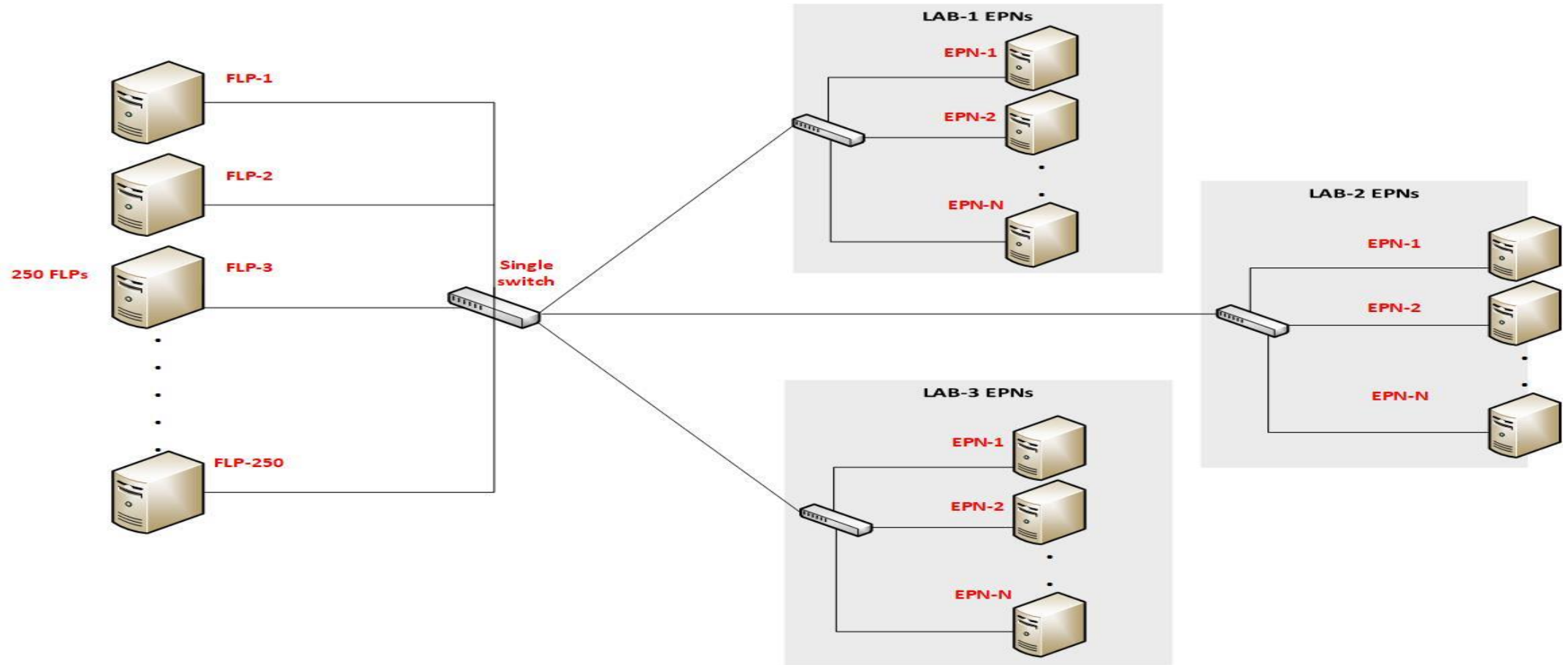


Figure 3: Congestion control mechanism (50 labs i.e. containers)

# Summary of work / meetings of Islamabad Campus

Quarters	Agenda	Description
4 <sup>th</sup> 2018	<ul style="list-style-type: none"> <li>• New Architecture discussion</li> <li>• Complete report on FLP-EPN (all simulations).</li> <li>• COMSATS Islamabad presented perfect switch FULL scale simulation</li> </ul>	<ul style="list-style-type: none"> <li>• New architecture (FLPs and EPNs are placed at different locations).</li> <li>• Few other changes like for example change of bandwidth between FLPs and EPNs and room like structure for EPNs.</li> <li>• Presented a report on complete simulation from first to last simulation results in order to show the similarity and effectiveness of results.</li> </ul>
3 <sup>rd</sup> 2018	<ul style="list-style-type: none"> <li>• COMSATS Islamabad group discussion on simulation time and storage</li> </ul>	<ul style="list-style-type: none"> <li>• Successfully reduce the time of simulation by using the command line instead of GUI simulation</li> <li>• Result file was reduced to MBs (Initially 300-400 GBs).</li> </ul>
1 <sup>st</sup> and 2 <sup>nd</sup> 2018	<ul style="list-style-type: none"> <li>• COMSATS Islamabad proposed multi-layered model for FLP-EPN simulation.</li> </ul>	<ul style="list-style-type: none"> <li>• Proposed multi-layered architecture by COMSATS Islamabad</li> <li>• Results were also discussed.</li> <li>• The building of software for O2.</li> <li>• Simulation workflow (building time frame).</li> <li>• RDMA protocol (DMA transfer b/w nodes).</li> <li>• Omnipath of Intel (discussion on its benefits).</li> <li>• Comparison of Omnipath and InfiniBand</li> </ul>

# Current Standing

<b>Task Assigned</b>	<b>Progress</b>
<b>New architecture implementation and Development (FLPs and EPNs are physically separated)</b>	<b>In Progress (50% completed)</b>
<b>Full scale architecture simulation. (Single Switch Implementation)</b>	<b>Completed</b>
<b>Multi-level hierarchical network topology along with traffic shaping.</b>	<b>Completed</b>
<b>Small-scale simulation results of the proposed model</b>	<b>Completed</b>
<b>Implement the scale up simulation using the HPC facility available at the CUI Islamabad.</b>	<b>Completed</b>
<b>Perfect switch design implementation.</b>	<b>Completed</b>
<b>R&amp;D on simulation tools available.</b>	<b>Completed</b>

Thank you