

Développement d'un Event-Builder à 10Gbps pour une caméra de télescope Cherenkov

Julien HOULES Dirk HOFFMANN

Centre de Physique des Particules de Marseille
IN2P3/CNRS



Context

Le Cherenkov Telescope Array

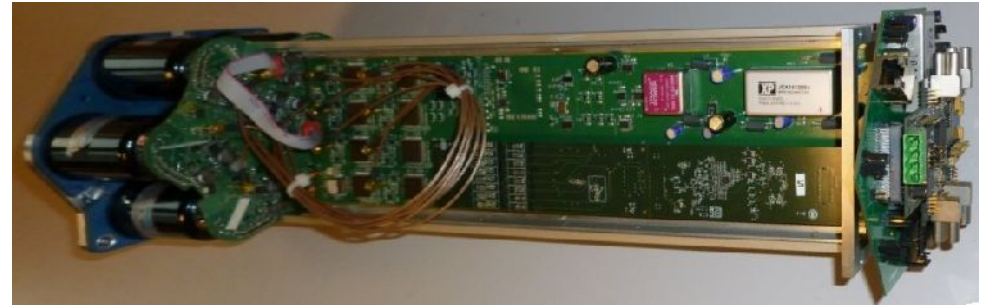
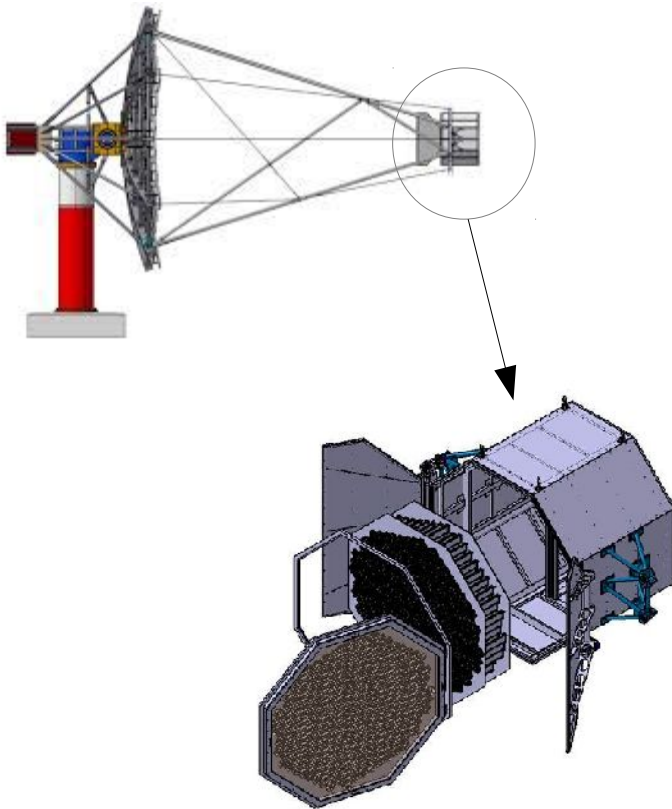


2 réseaux différents envisagés : nord et sud

~ 100 Télescopes gamma

3 tailles différentes : LST (10 GeV-100 GeV) 24 mètres,
MST (100 GeV-1 TeV) 12 mètres,
SST (> qqes TeV) 6 mètres.

Sur le MST, 2 caméras FlashCam et NectarCam

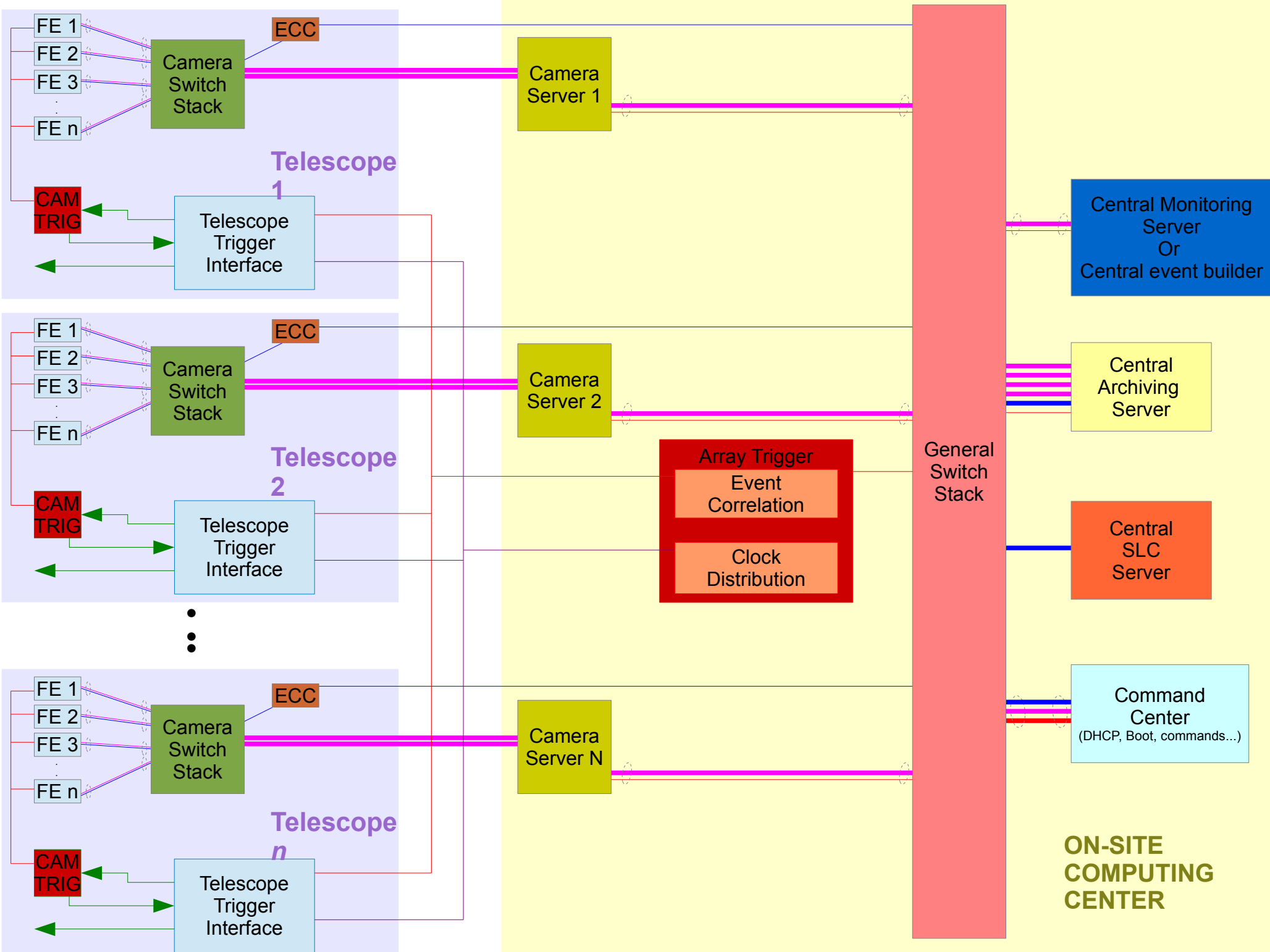


~ 1800 photomultiplicateurs (PM)

7 PM par module

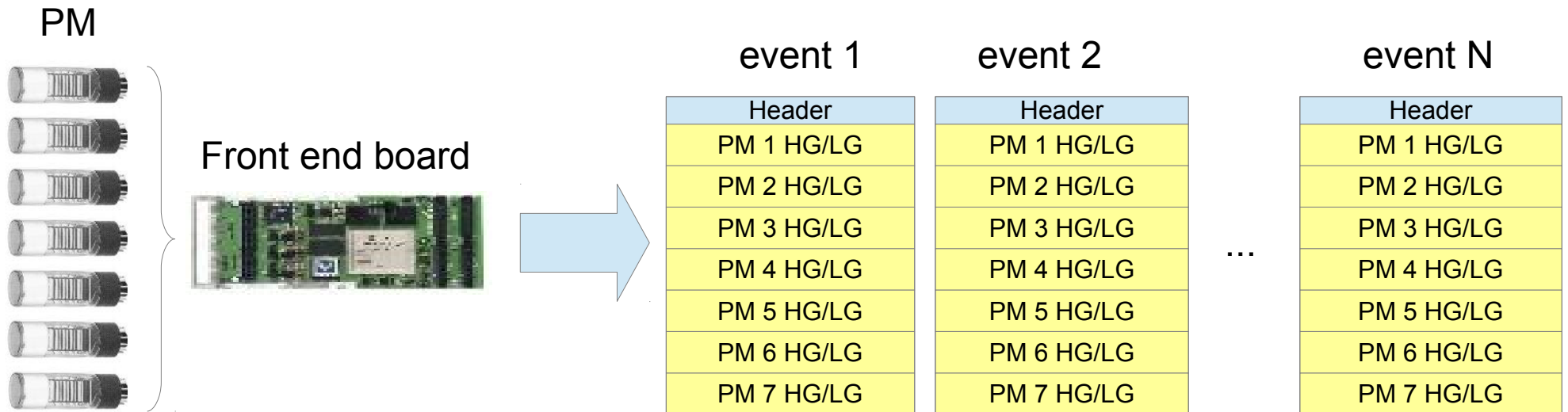
→ ~250 modules

Un lien Gigabit Ethernet par module



Front end data flow

Hypotheses



Whole Camera ~ 2016 PM -> 288 front end boards

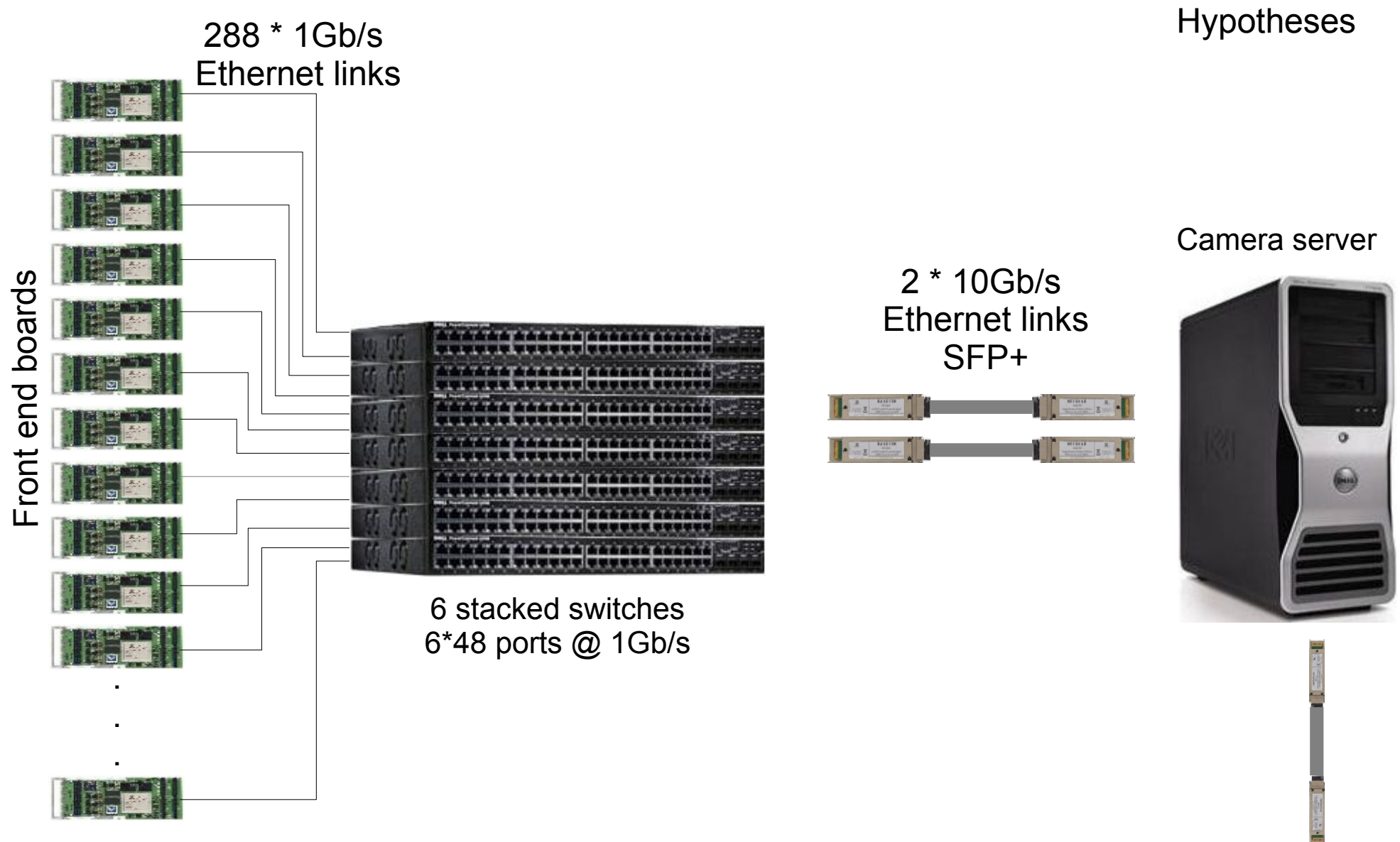
Camera DAQ Hypotheses

Hypotheses

- 2016 pixels camera
- Average camera trigger rate : 8 kHz
- Maximum size of waveforms for 1 PM :
144 bytes (16 bit * 2 gains * 36 samples)
 - Max theoretical bandwidth = $8000 * 2016 * 144 = 2.32 \text{ GB/s}$
 $= 18.6 \text{ Gb/s}$
- 7 detectors for each front end board : 288 boards/camera
 - Each board generates a flow of $2320/288 = 8 \text{ MB/s}$
 $= 64 \text{ Mb/s}$

With 1024 bytes packets (network headers neglected)

Camera DAQ Architecture



Hypotheses

Camera server

Critical points

Critical points

Poissonian distribution of the triggers :

The average data rate for a trigger rate of 8 kHz can be forwarded on the two uplinks but when instantaneous rate becomes much higher, packets must be stored.

Critical points

Poissonian distribution of the triggers :

The average data rate for a trigger rate of 8 kHz can be forwarded on the two uplinks but when instantaneous rate becomes much higher, packets must be stored.

Synchronicity of the packets sending :

Will the switches be able to process packets arriving exactly at the same time on all the ports (1 ns) ?

Critical points

Poissonian distribution of the triggers :

The average data rate for a trigger rate of 8 kHz can be forwarded on the two uplinks but when instantaneous rate becomes much higher, packets must be stored.

Synchronicity of the packets sending :

Will the switches be able to process packets arriving exactly at the same time on all the ports (1 ns) ?

Packets reception and processing at link speed :

Medium-sized packets reception with regular Linux sockets at link speed is not possible

Critical points

Poissonian distribution of the triggers :

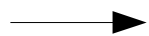
The average data rate for a trigger rate of 8 kHz can be forwarded on the two uplinks but when instantaneous rate becomes much higher, packets must be stored.

Synchronicity of the packets sending :

Will the switches be able to process packets arriving exactly at the same time on all the ports (1 ns) ?

Packets reception and processing at link speed :

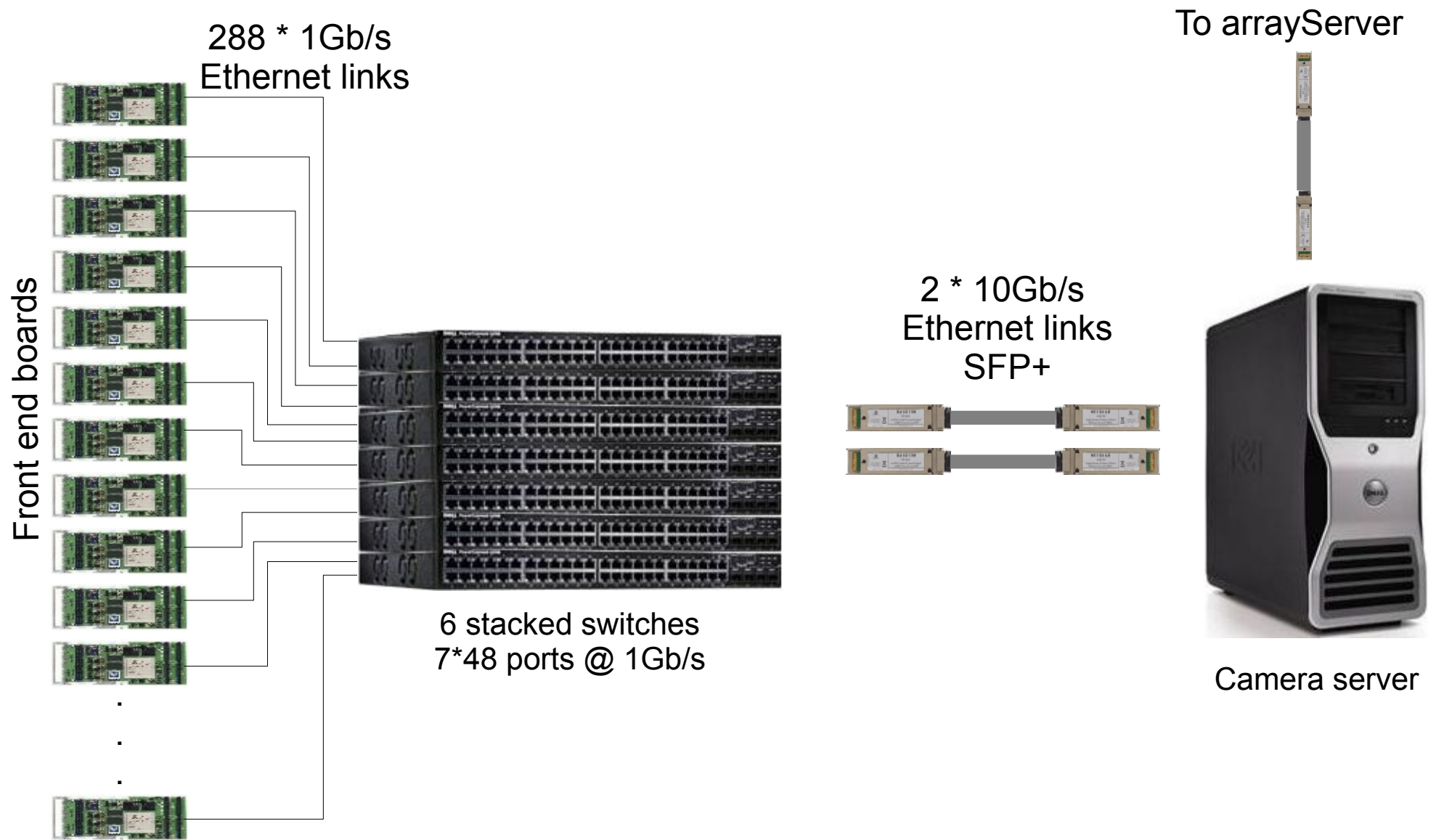
Medium-sized packets reception with regular Linux sockets at link speed is not possible



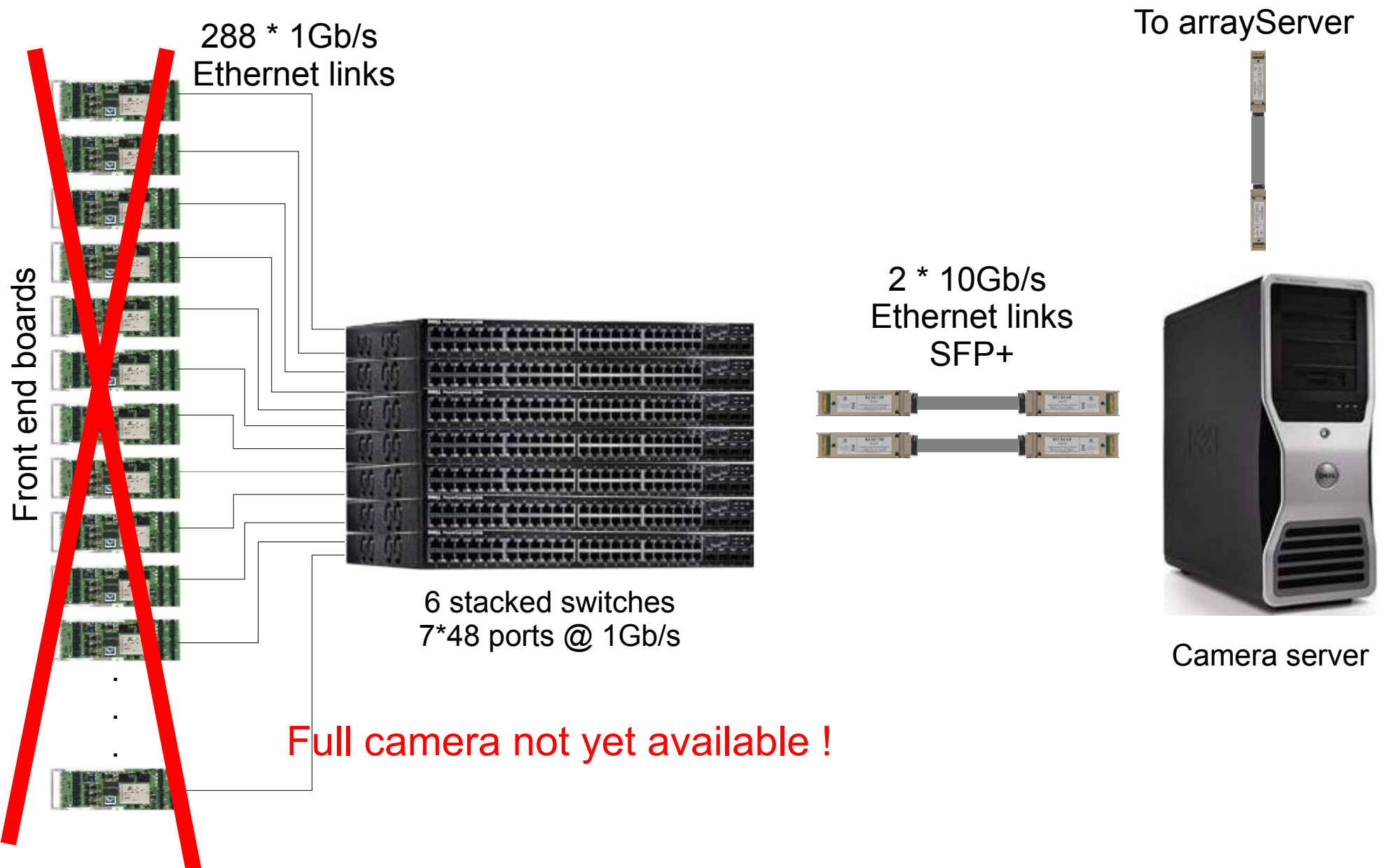
Need to build a stimulator able to validate the acquisition chain

Stimulator

Camera DAQ Architecture



Camera DAQ Architecture



Full camera not yet available !

Stimulator : Architecture

Sync



SBC Jetway

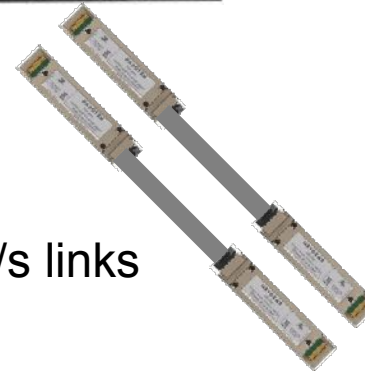


⋮



Journées informatique 2014 1Gb/s links

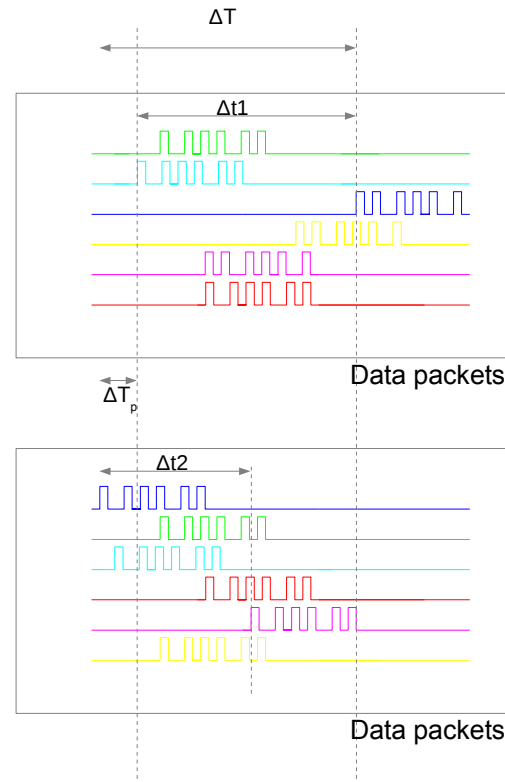
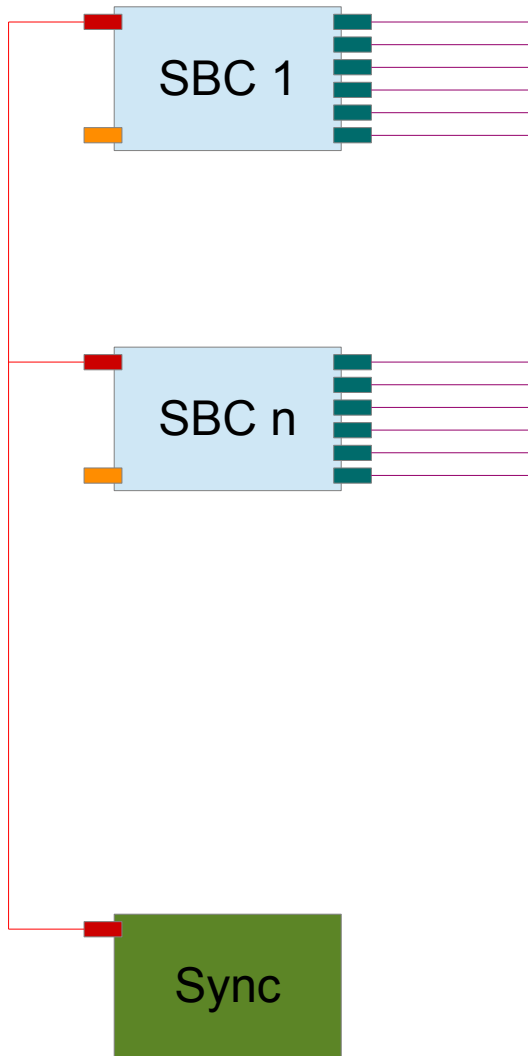
Switches stack



2 * 10Gb/s links

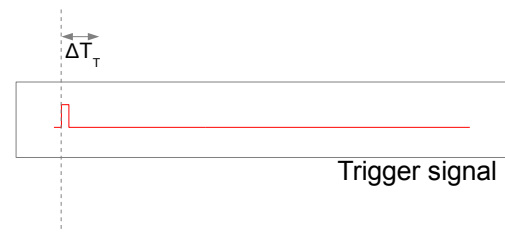


Main obstacles



$\Delta t_1, \Delta t_2$: time difference between the first packet sent and the last one on one SBC

ΔT_p : maximum time difference between the first packets sent by each SBC



ΔT_T : time difference between the « send packet signal » and the first packet sent among all the SBCs

Main obstacles

First tests with packets sending on an interrupt were very disappointing :

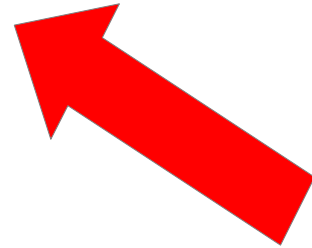
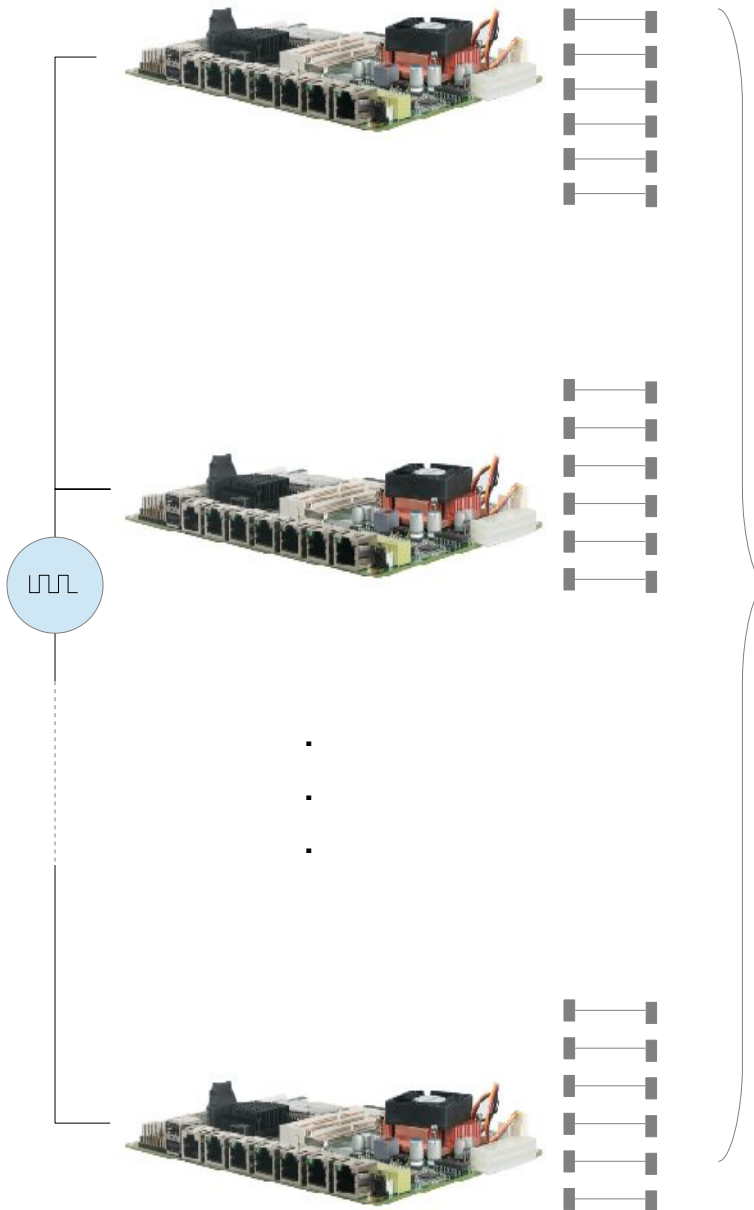
About 30 μ s with a great jitter between the first and last paquet sent by one board and even more with two boards

Results slightly improved with Xenomai (real time patch for Linux)

Need to avoid interrupt latency, serialize packets sending and bypass Linux communication layers to talk directly to hardware.

Stimulator : Architecture

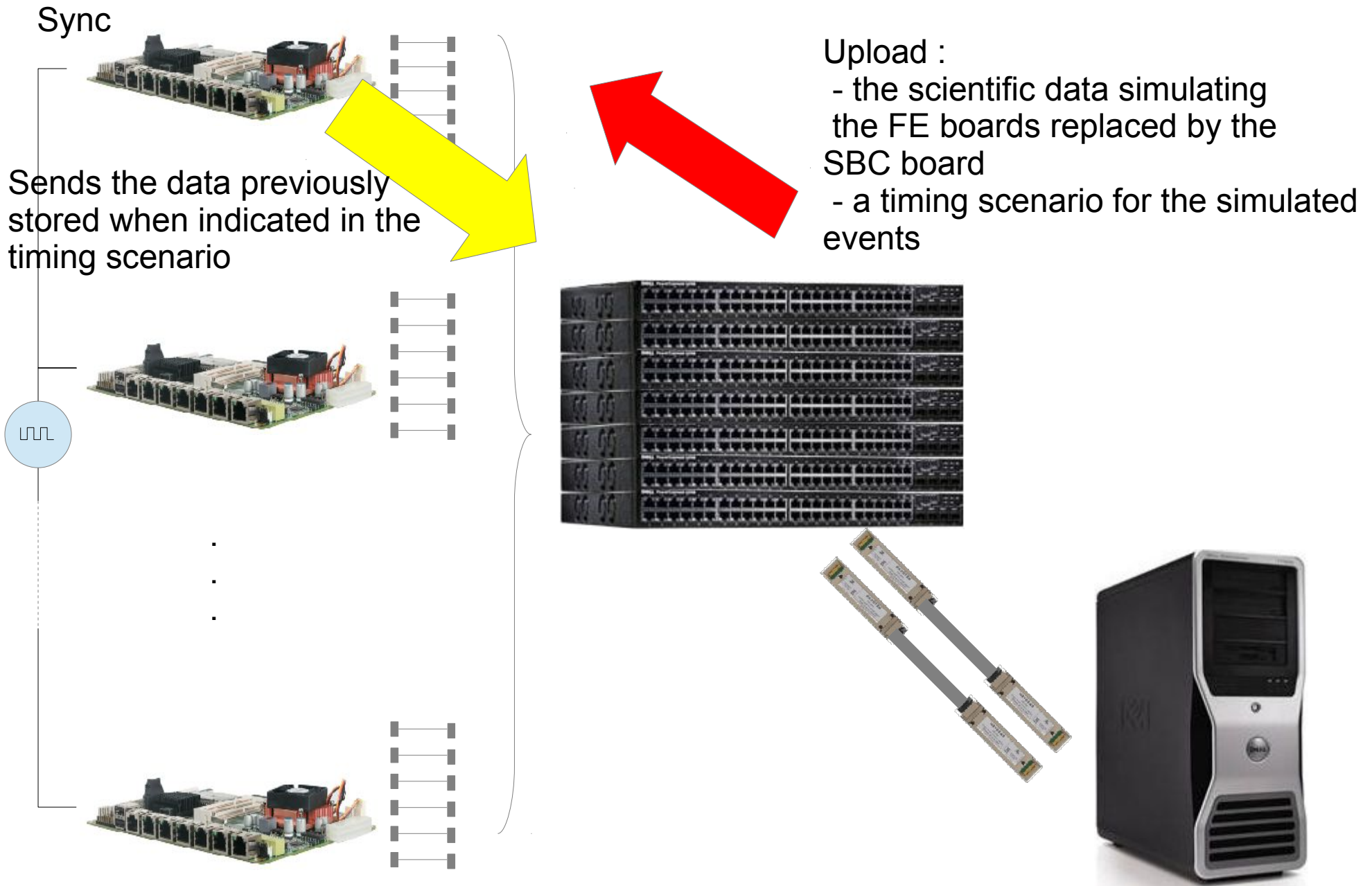
Sync



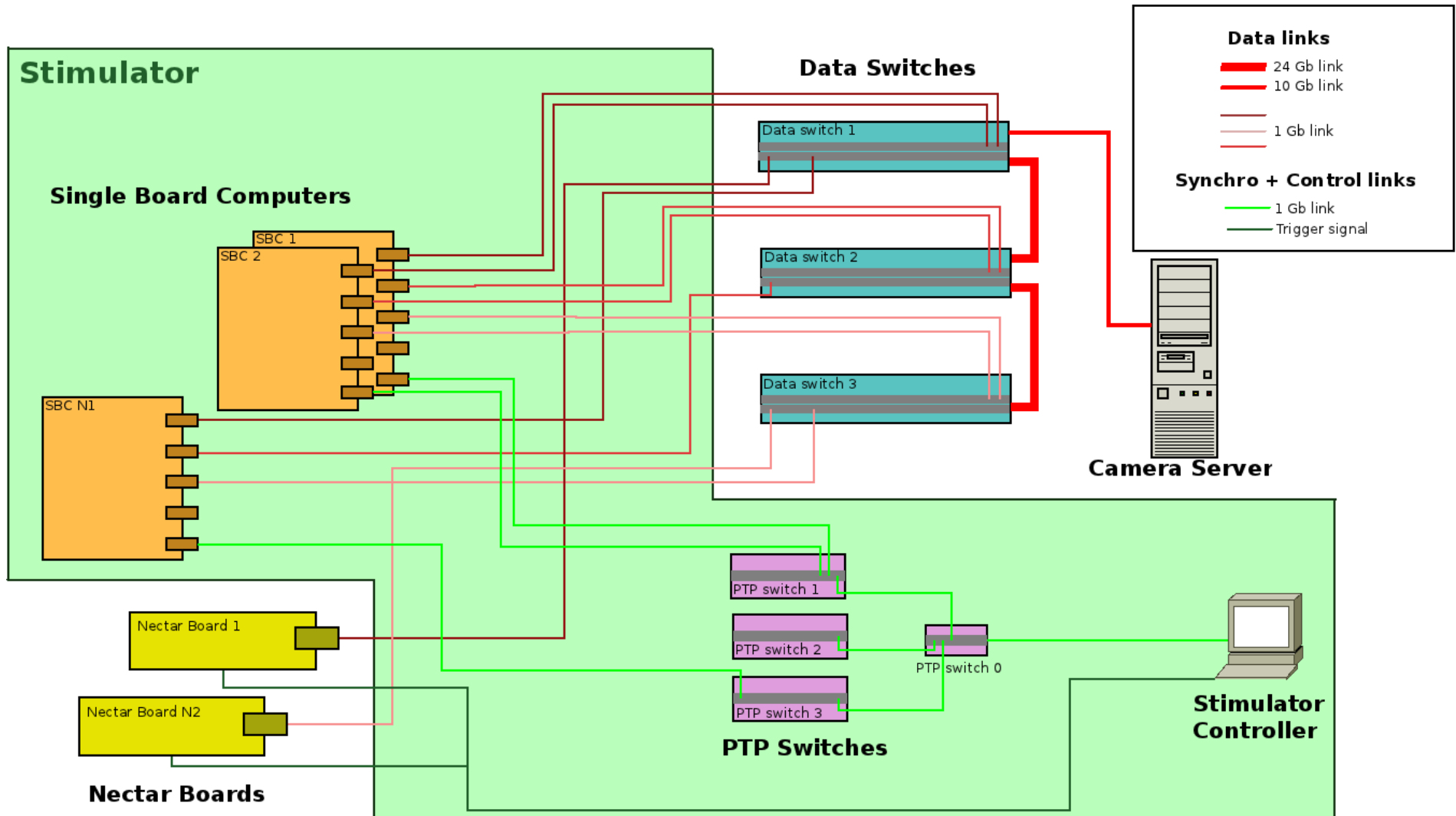
Upload :

- the scientific data simulating the FE boards replaced by the SBC board
- a timing scenario for the simulated events

Stimulator : Architecture



Lab setup



Solutions applied

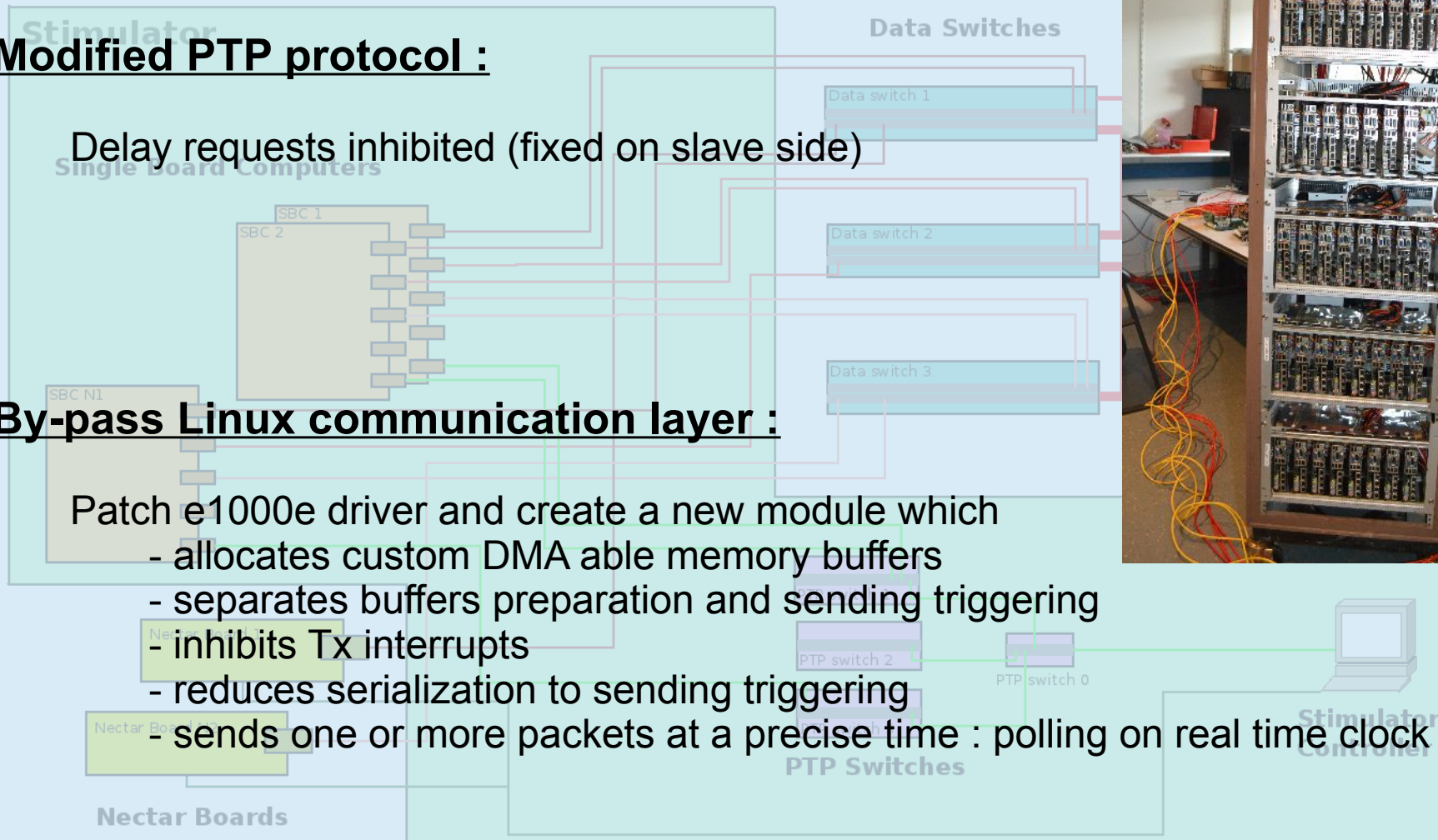
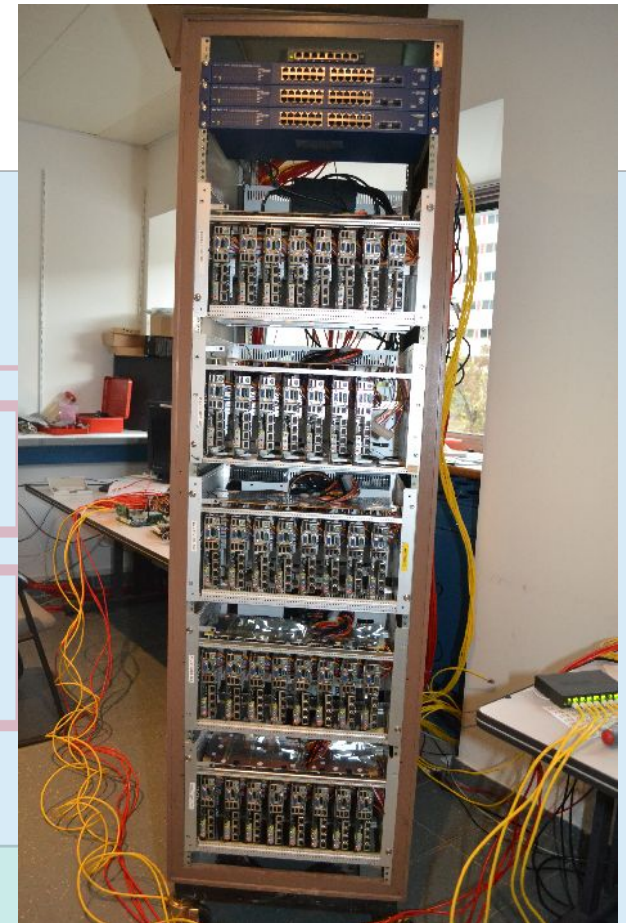
Modified PTP protocol :

Delay requests inhibited (fixed on slave side)

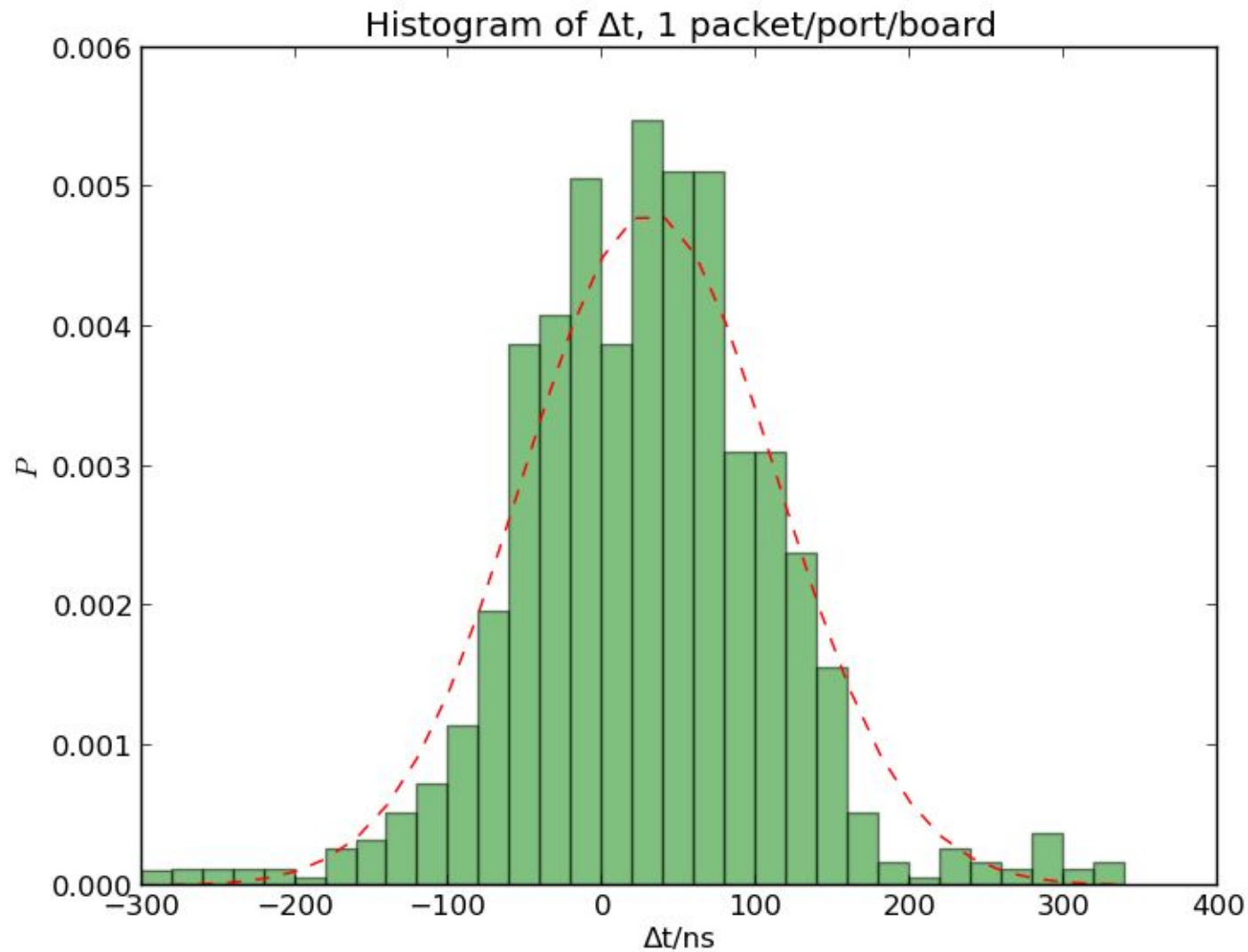
By-pass Linux communication layer :

Patch e1000e driver and create a new module which

- allocates custom DMA able memory buffers
- separates buffers preparation and sending triggering
- inhibits Tx interrupts
- reduces serialization to sending triggering
- sends one or more packets at a precise time : polling on real time



Results



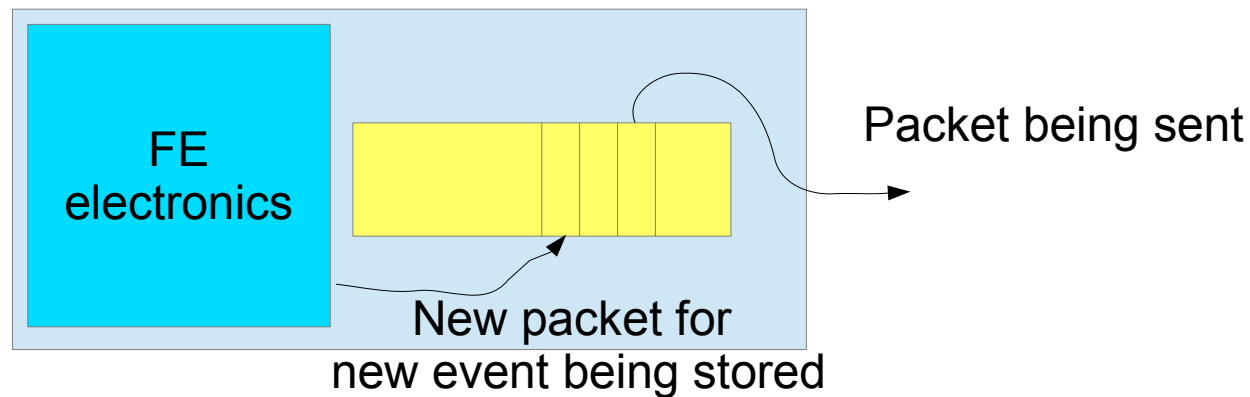
Mean = 30 ns

Sigma = 83 ns

Data bursts management

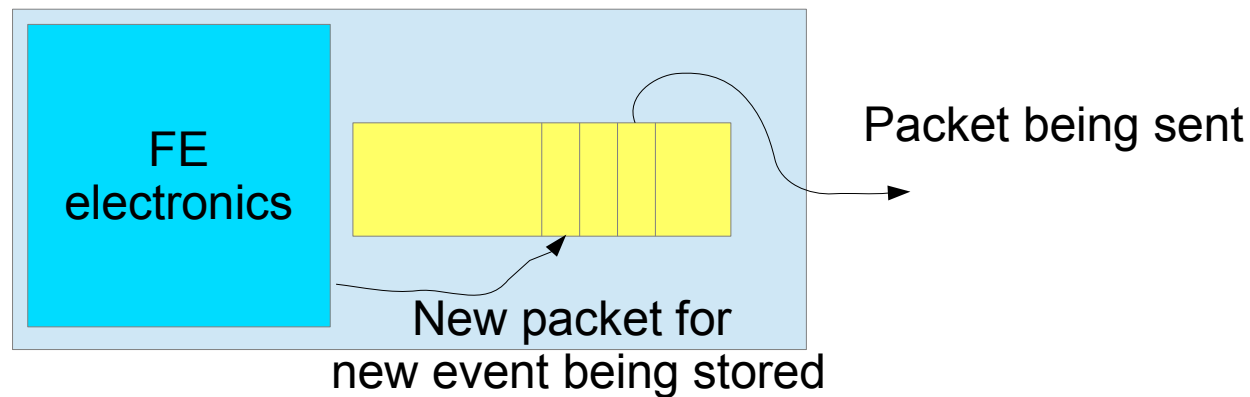
Instantaneous flow : critical point

- Average camera trigger rate : 8 kHz
- Time to send a 1024 bytes packet on a 1 gigabit link : 8 μ s
- If *FE dead time* < Δt between triggers < 8 μ s, the packet is stored in the FE packets buffer



Instantaneous flow : critical point

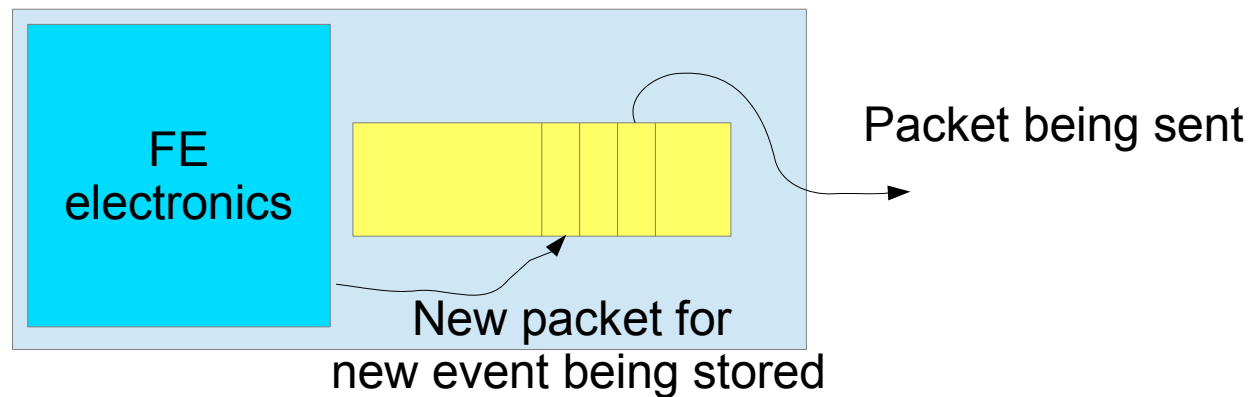
- Average camera trigger rate : 8 kHz
- Time to send a 1024 bytes packet on a 1 gigabit link : 8 μ s
- If *FE dead time* < Δt between triggers < 8 μ s, the packet is stored in the FE packets buffer



In this case, the front end board generates a **dataflow at link speed (1Gb/s)** until its packet buffer is empty.

Instantaneous flow : critical point

- Average camera trigger rate : 8 kHz
- Time to send a 1024 bytes packet on a 1 gigabit link : 8 μ s
- If *FE dead time* < Δt between triggers < 8 μ s, the packet is stored in the FE packets buffer

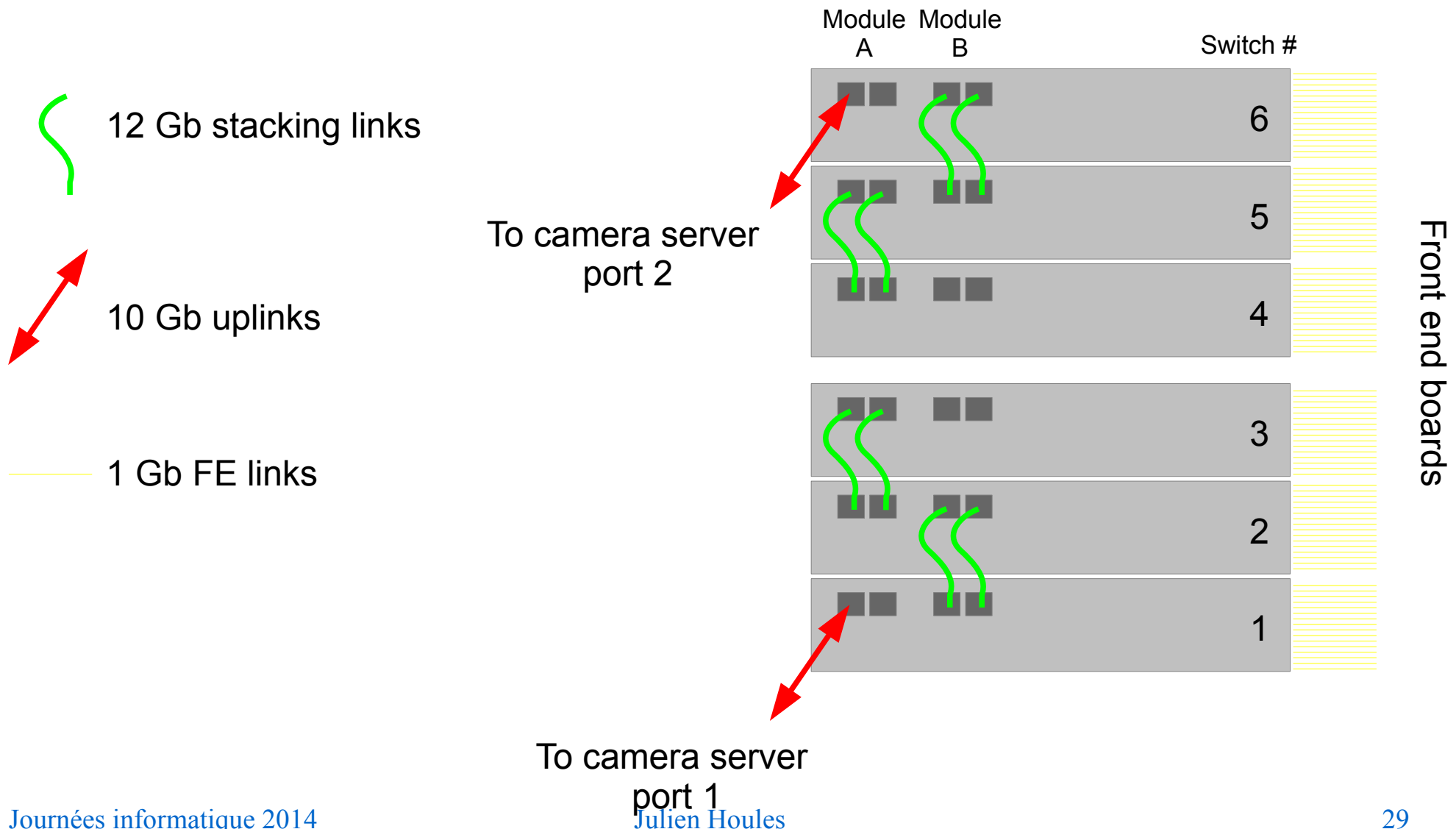


In this case, the front end board generates a **dataflow at link speed (1Gb/s)** until its packet buffer is empty.

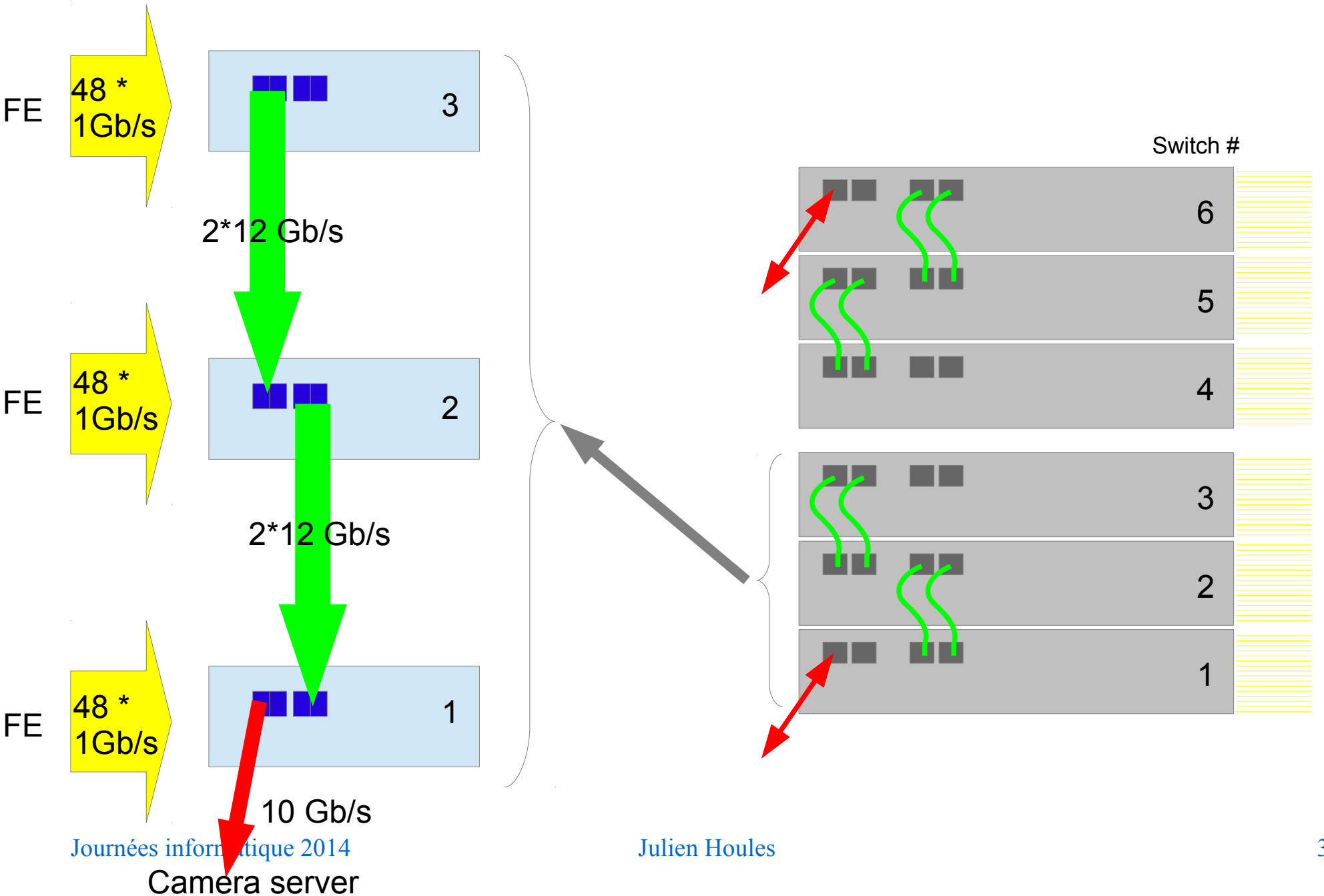
Will the switches be able to manage such an instantaneous flow **without losing packets** ?

Switches stack links

Switches stack rear panel :

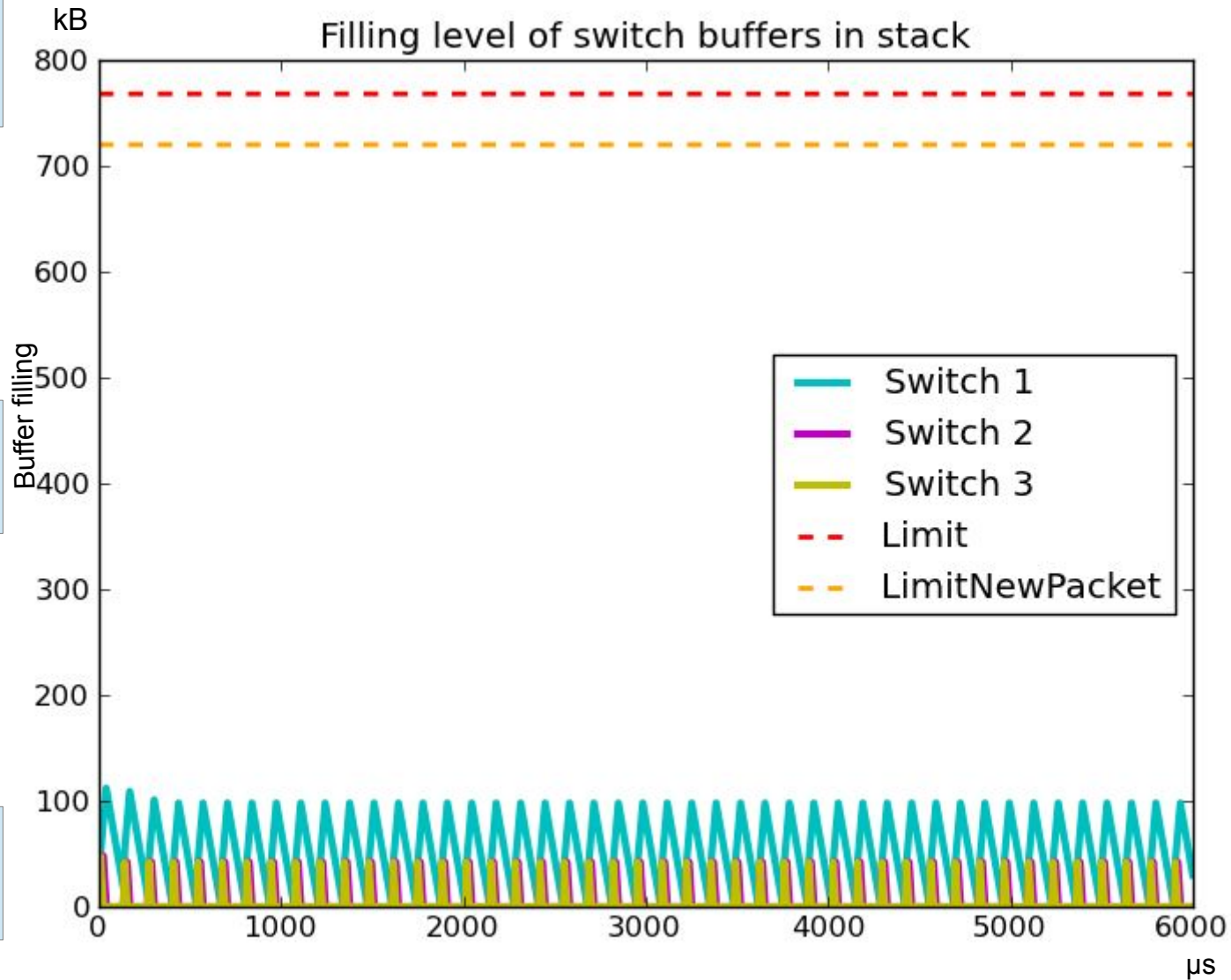
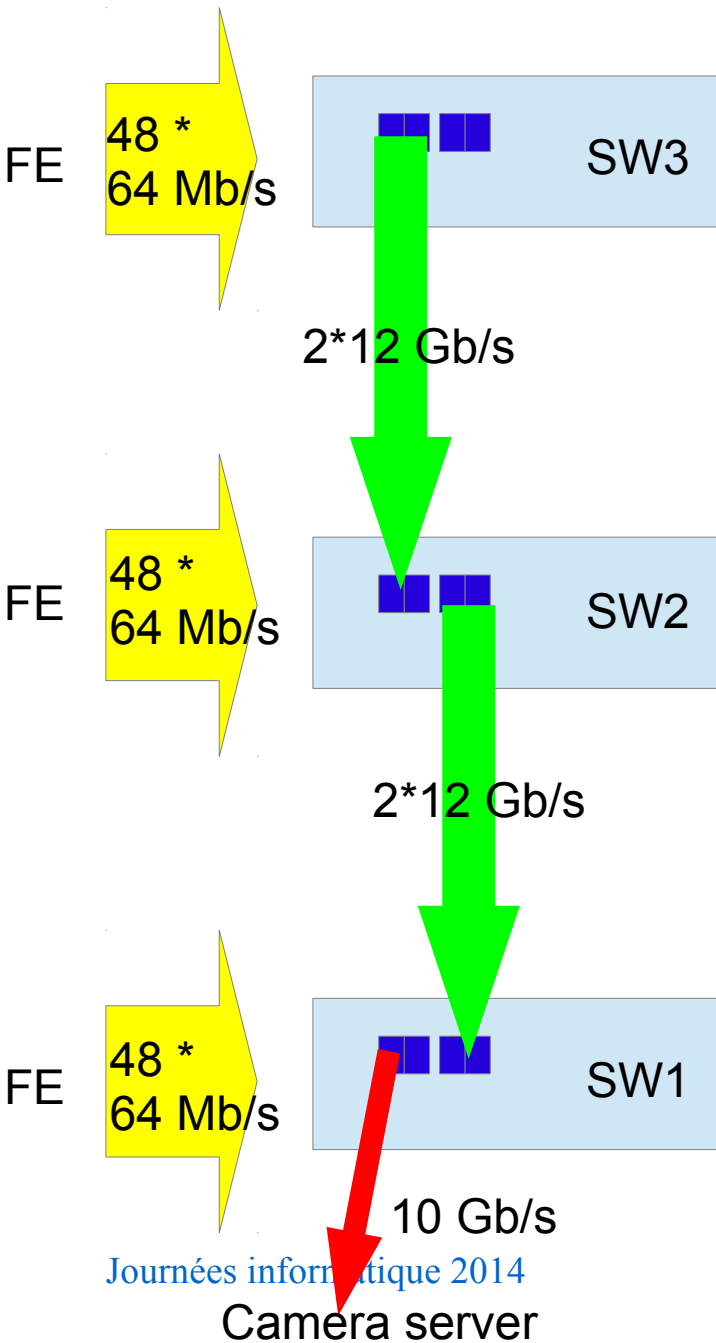


Switches stack model



Switches stack model

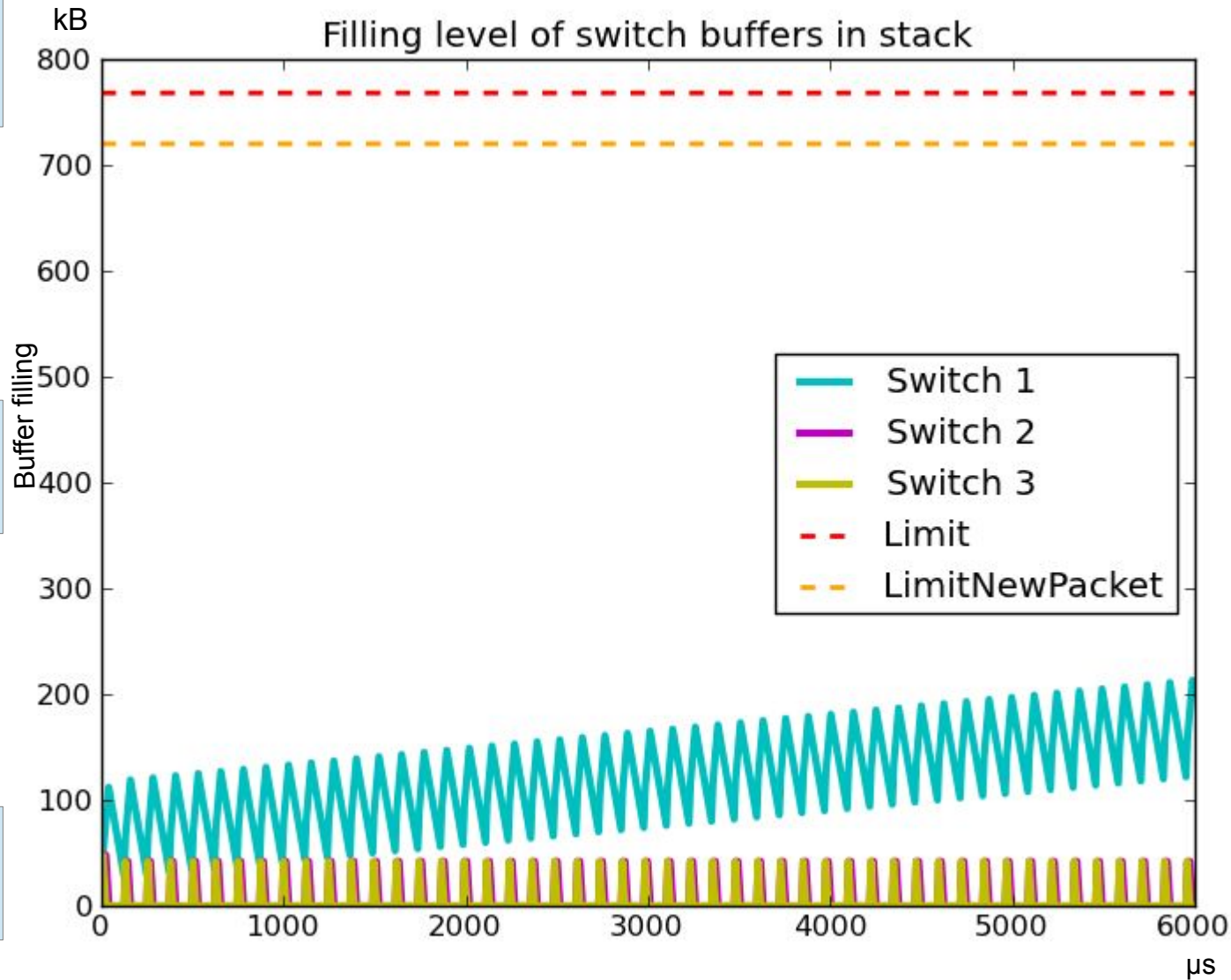
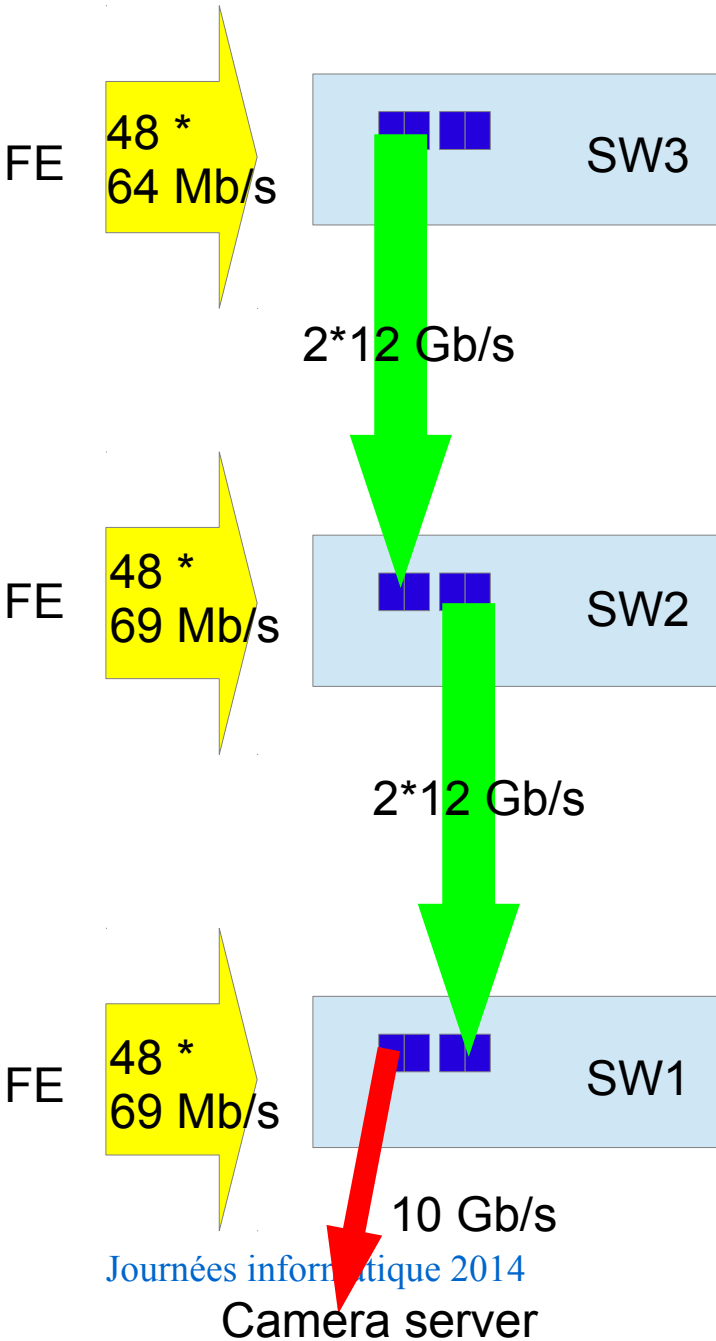
Each switch has a capacity of **768 kB**



50 consecutive incoming events within 6250 μs
1 event each 125 μs

Switches stack model

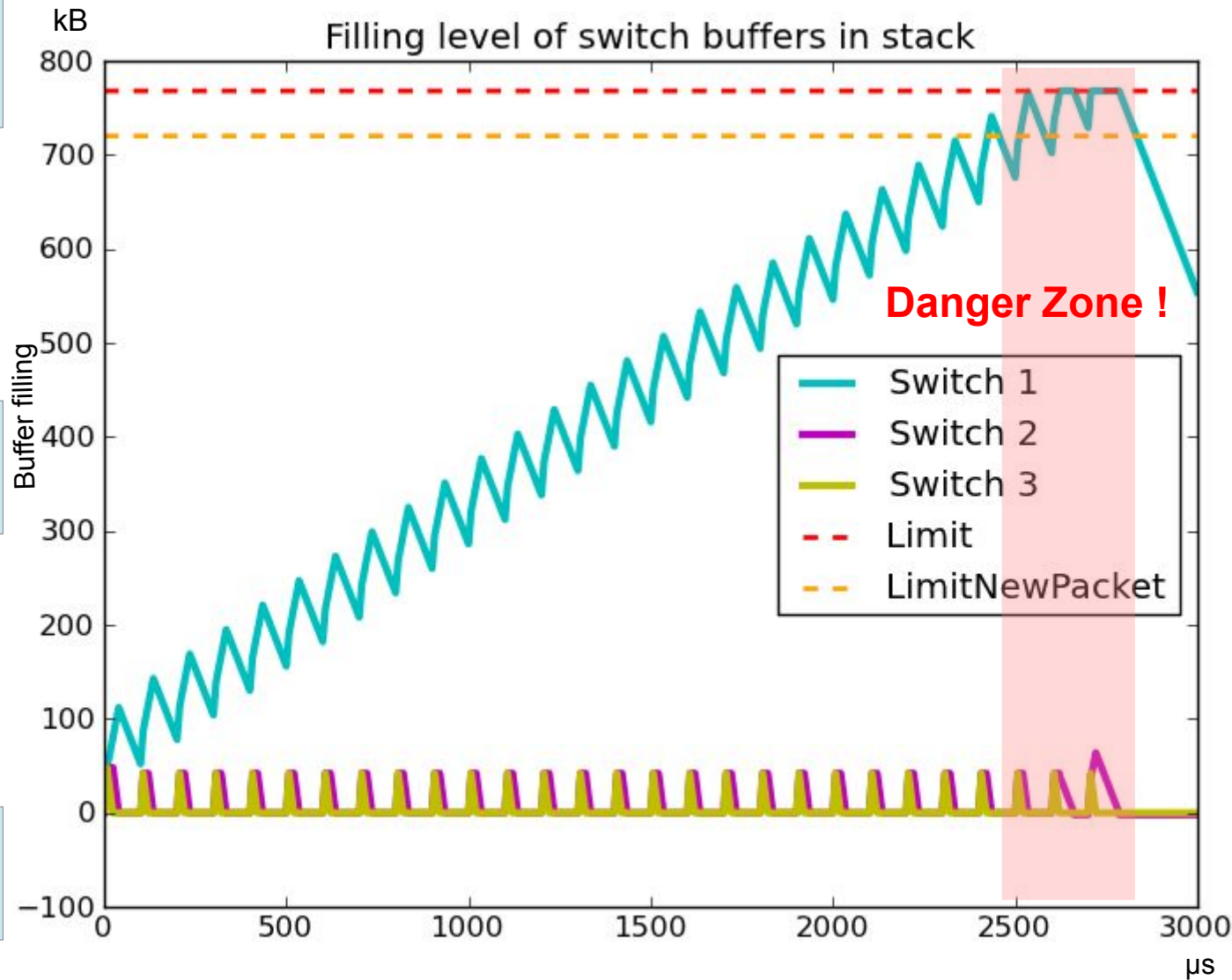
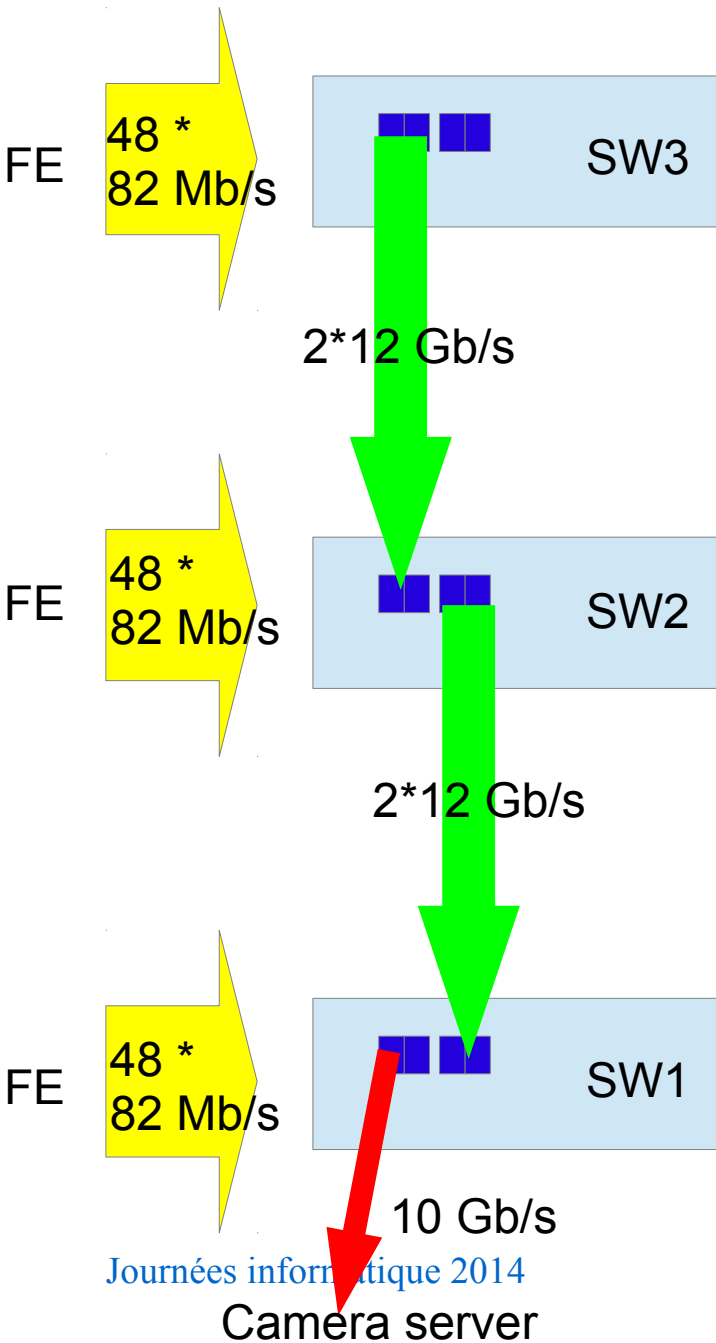
Each switch has a capacity of **768 kB**



50 consecutive incoming events within 5800 μ s
1 event each 116 μ s

Switches stack model

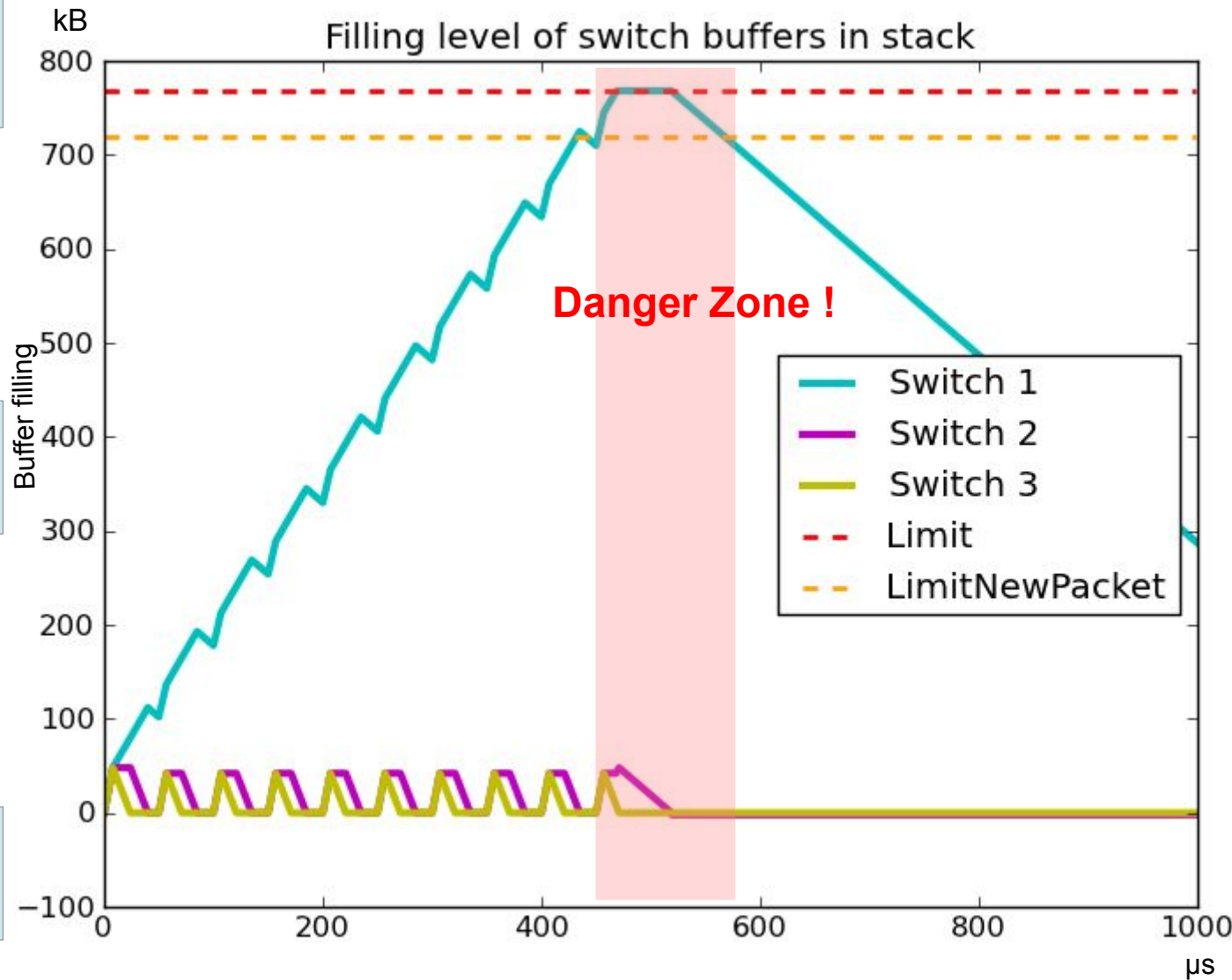
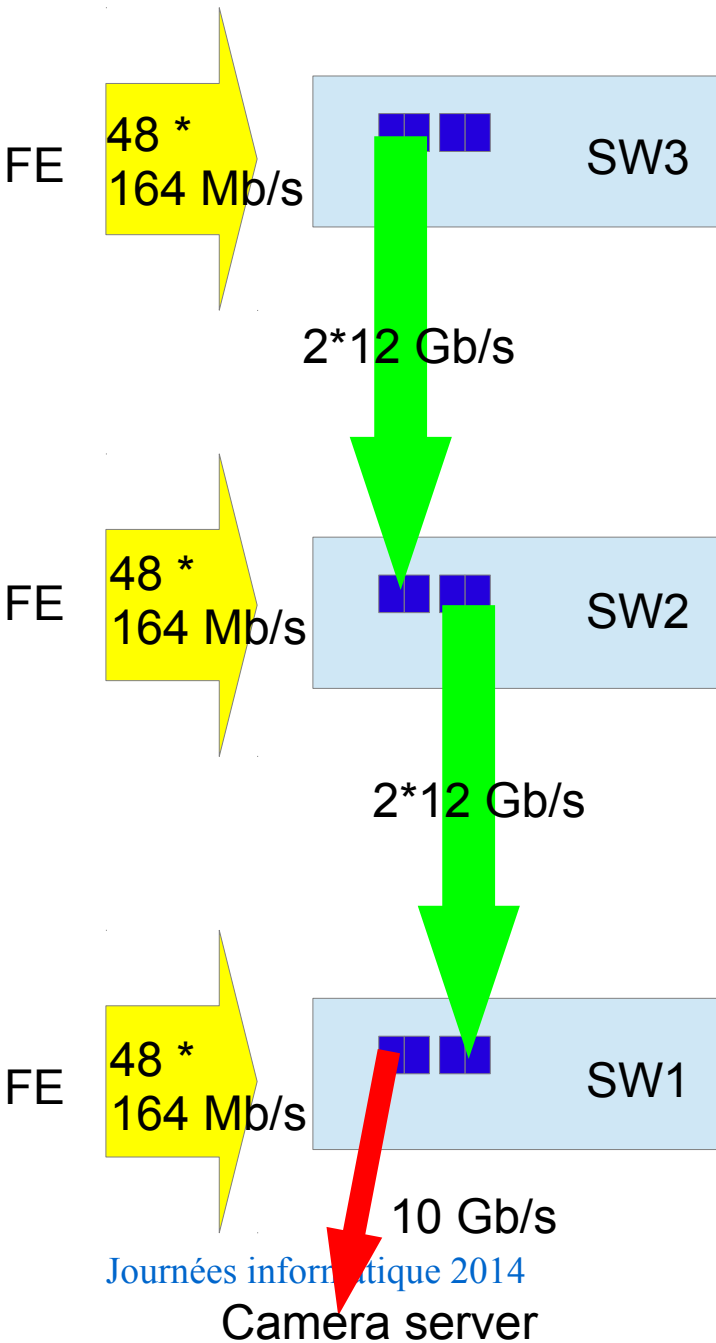
Each switch has a capacity of **768 kB**



28 consecutive incoming events within 2800 μs
1 event each 100 μs

Switches stack model

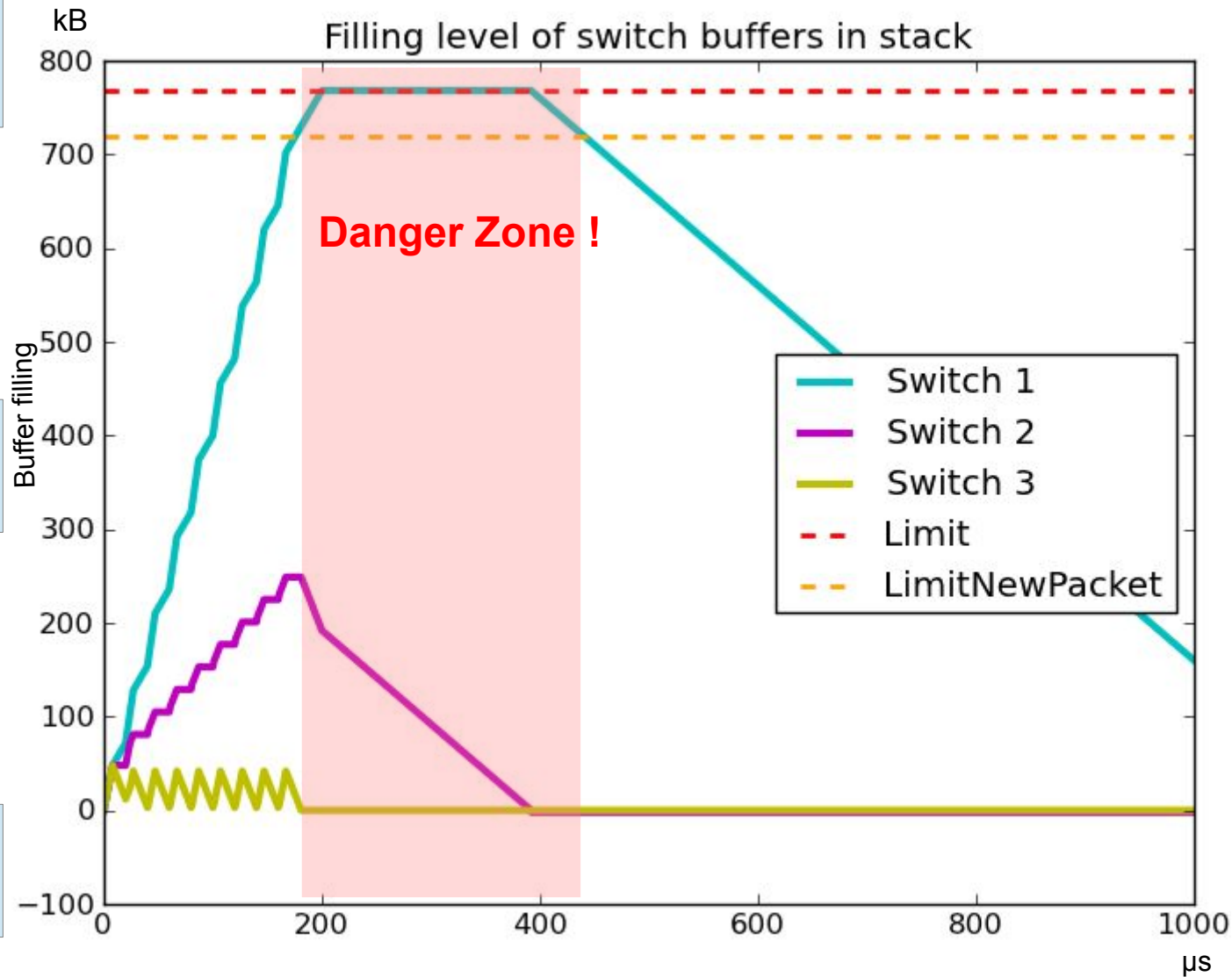
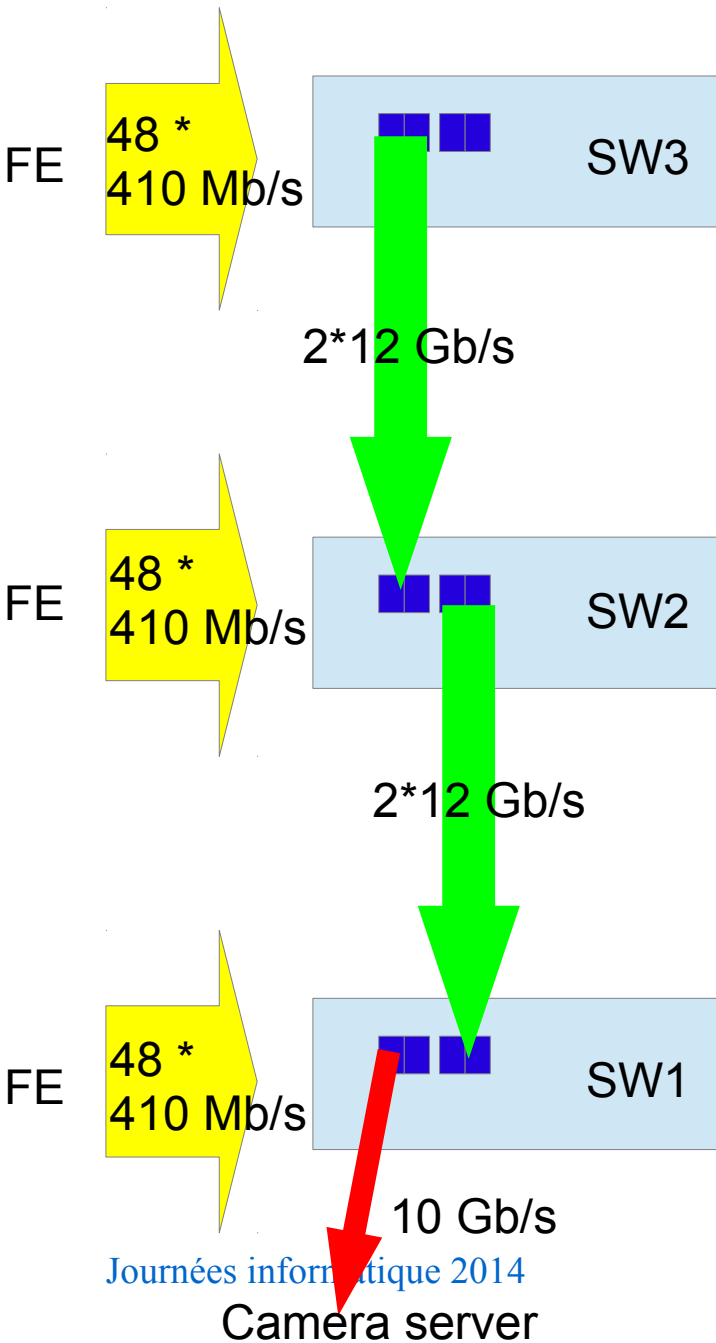
Each switch has a capacity of **768 kB**



10 consecutive incoming events within 500 µs
1 event each 50 µs

Switches stack model

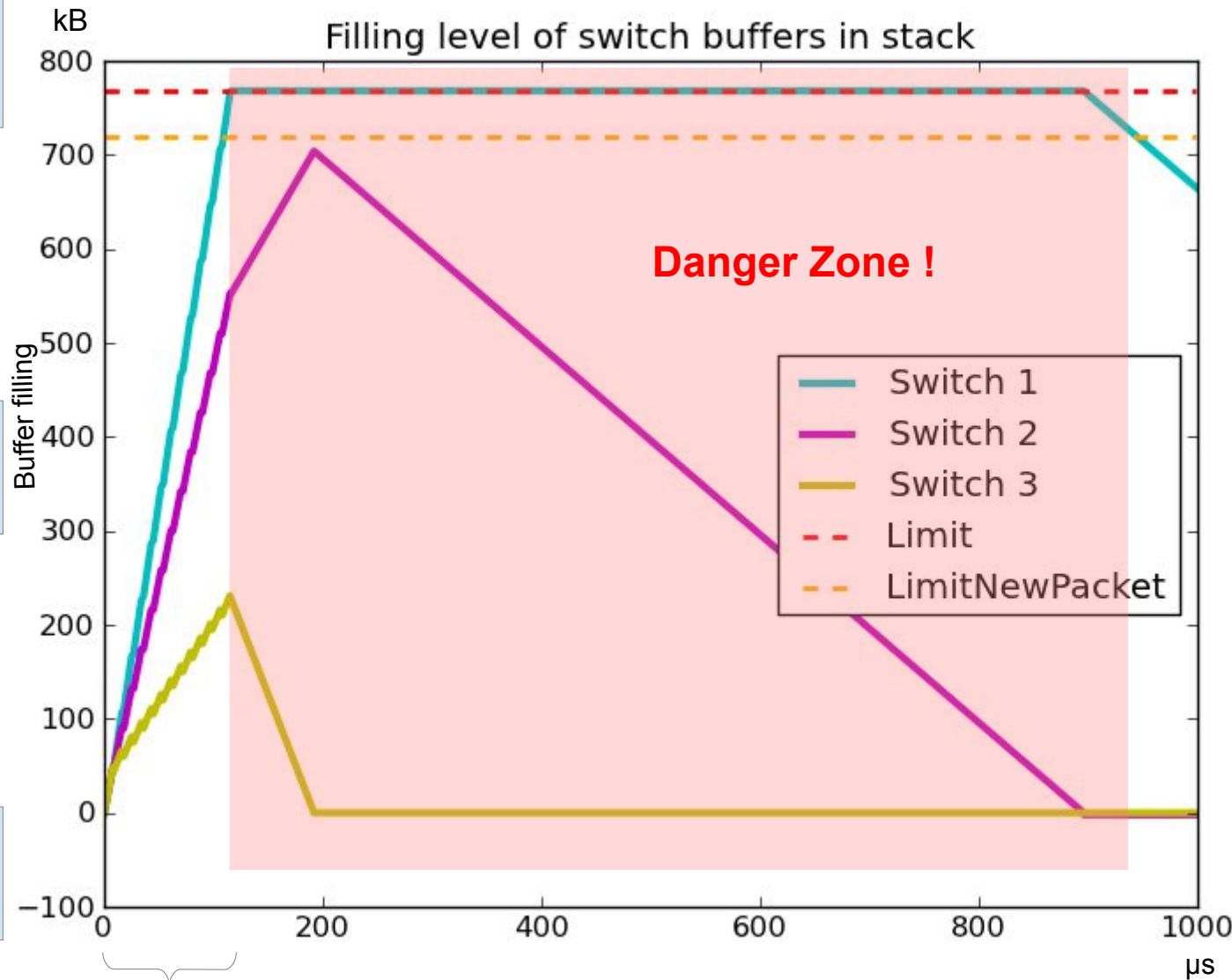
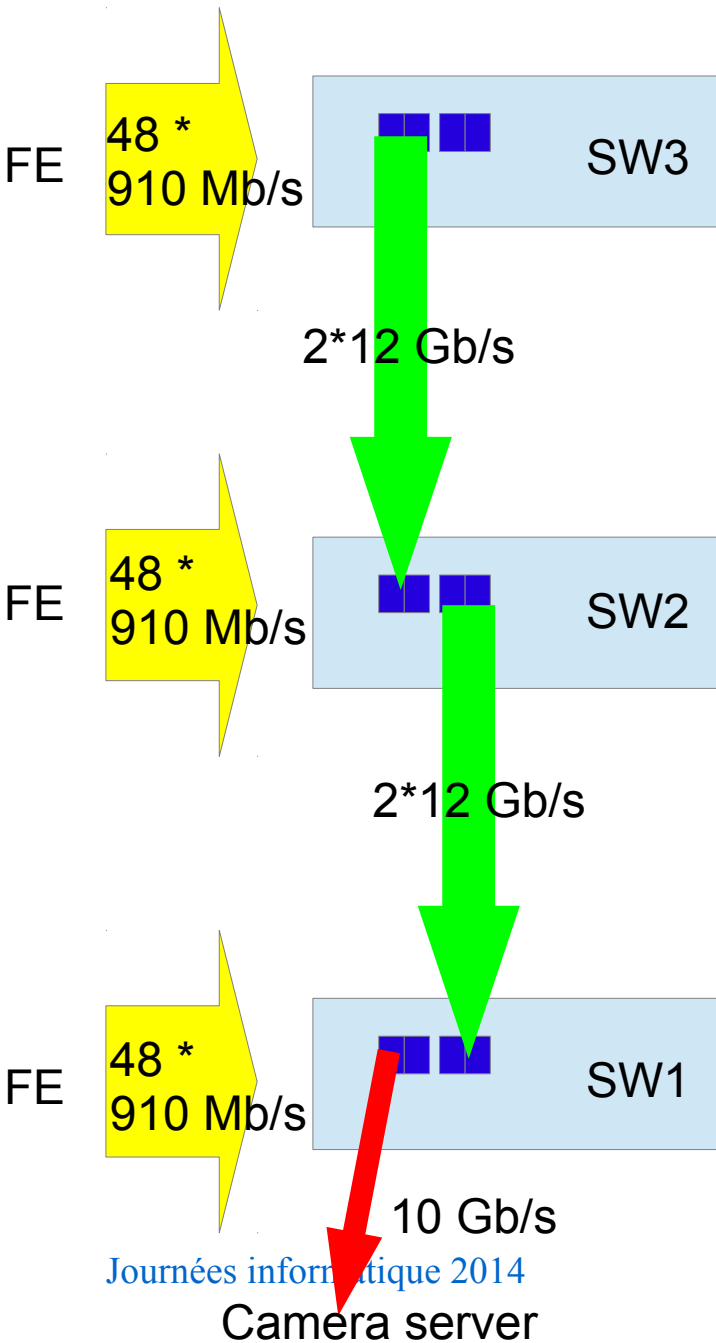
Each switch has a capacity of **768 kB**



9 consecutive incoming events within 180 µs
1 event each 20 µs

Switches stack model

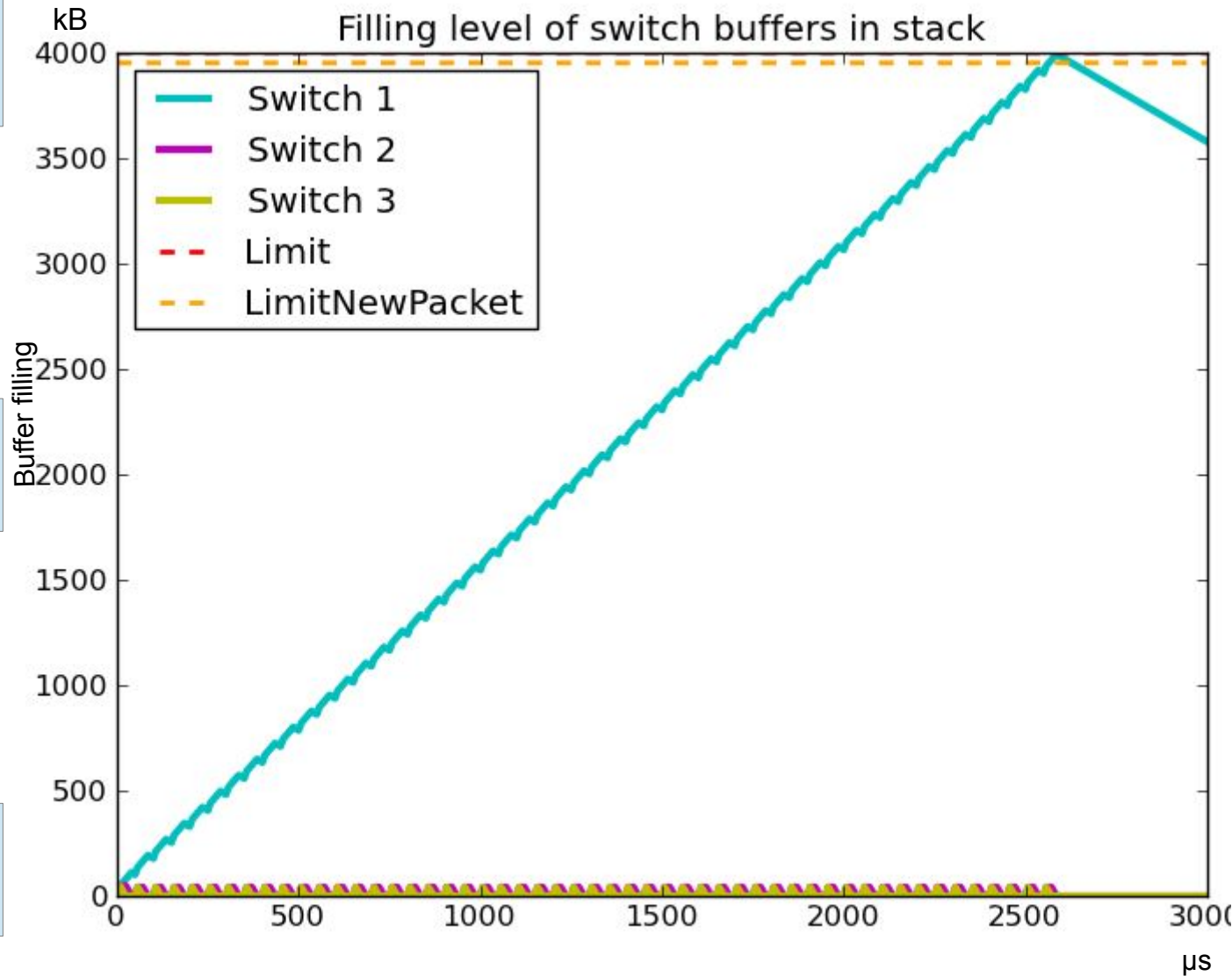
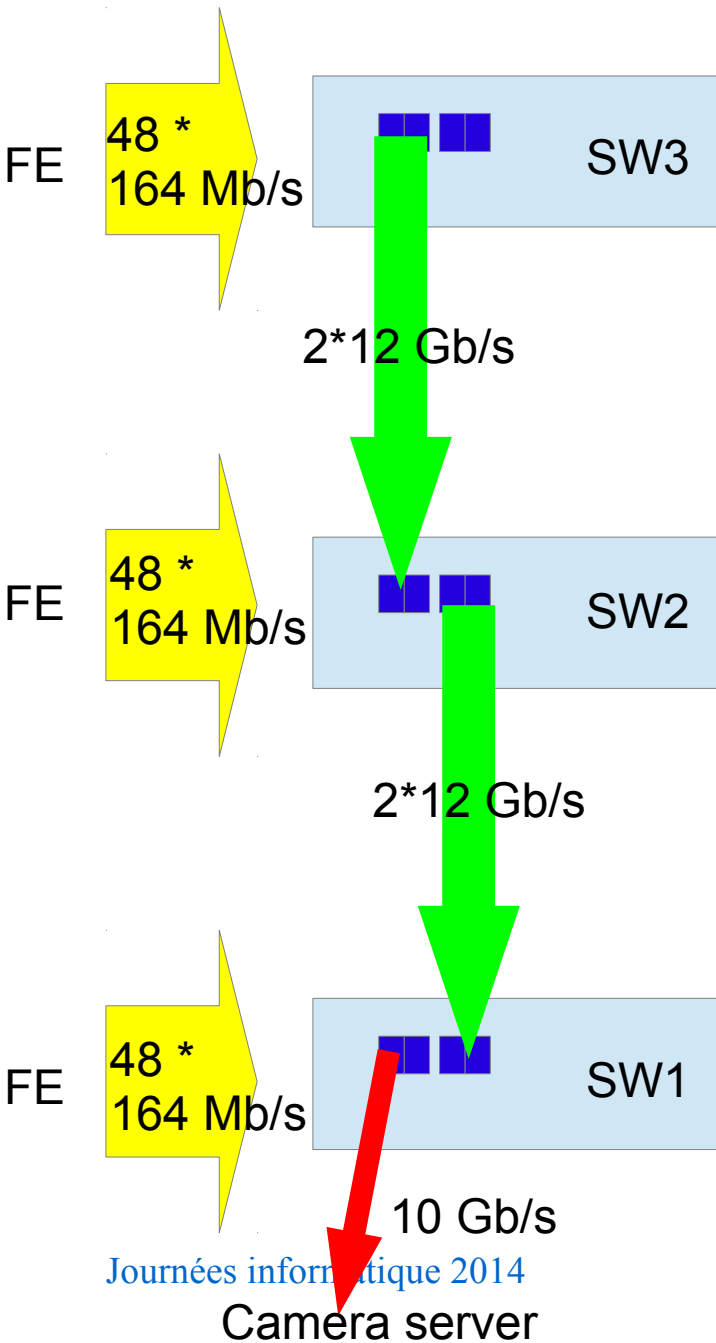
Each switch has a capacity of **768 kB**



13 consecutive incoming events within 117 μs
1 event each 9 μs

Switches stack model

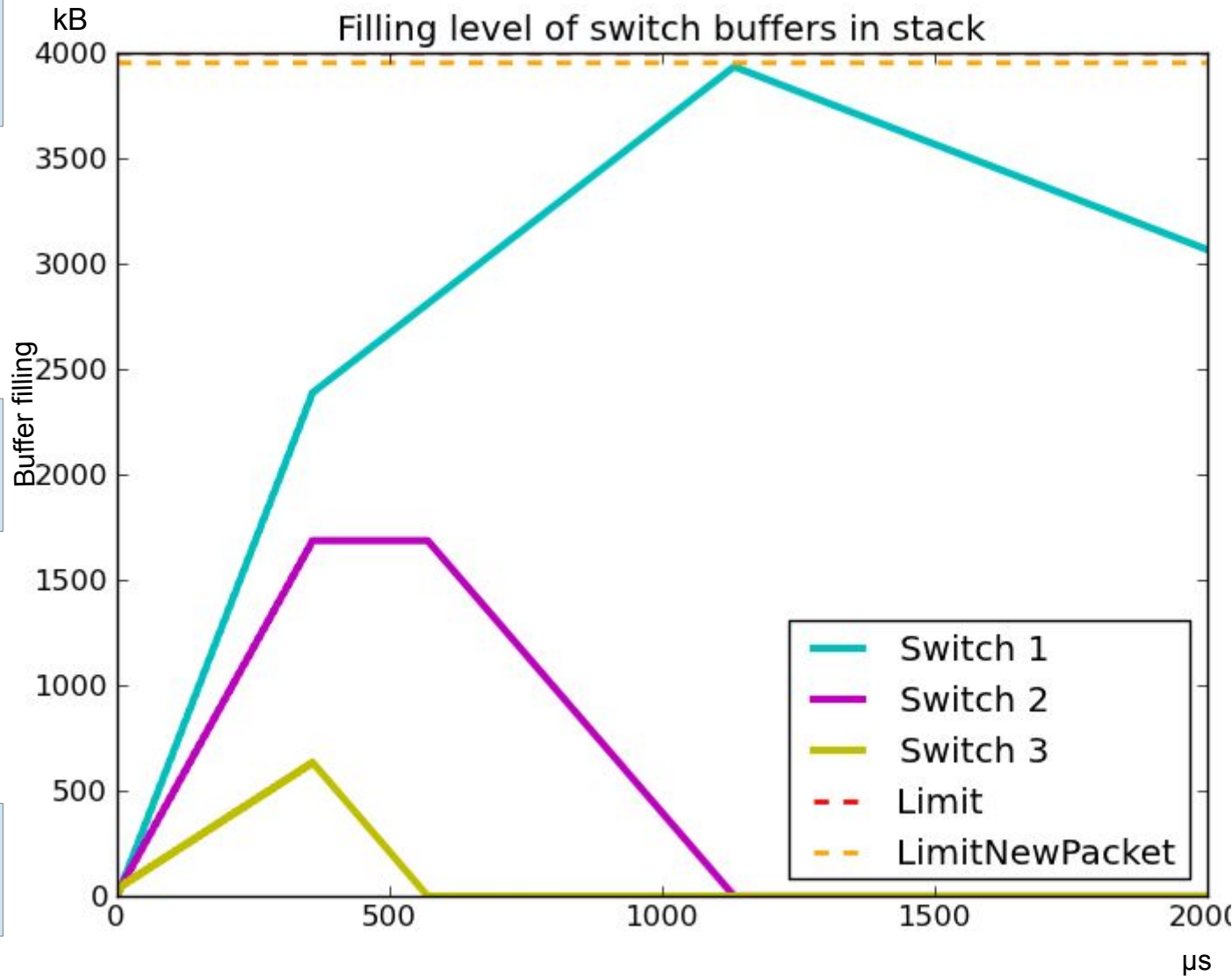
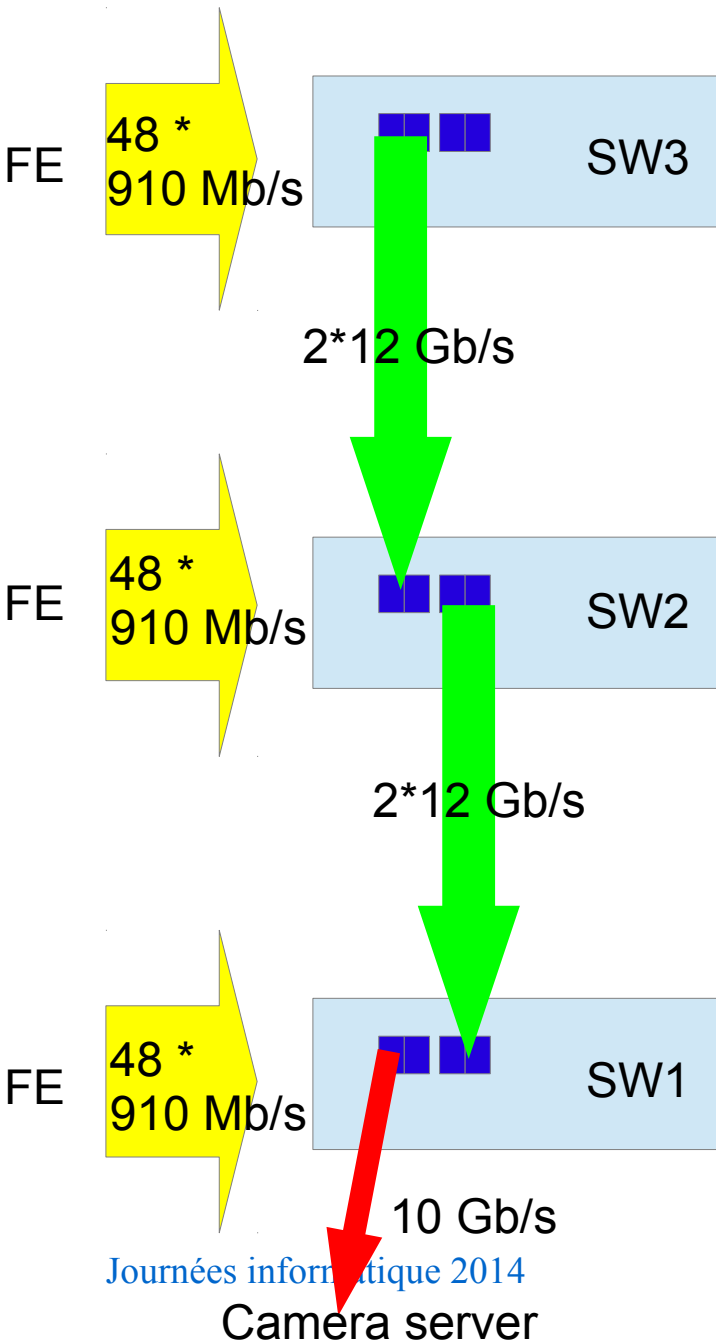
Each switch has a capacity of **4000 kB**



52 consecutive incoming events within 2600 μs
1 event each 50 μs

Switches stack model

Each switch has a capacity of **4000 kB**



40 consecutive incoming events within 360 μs
1 event each 9 μs

Switches examples

Dell PowerConnect 6248 Public price : 2600 Euros -> **15600 euros/Telescope**



BO = Buffer size = **768 kB**

Dell PowerConnect 7048 Public price : 3700 Euros -> **22200 euros/Telescope**



BO = Buffer size = **4 MB**

Problem :

Commercial switches don't seem to be well suited to our data flow model
(from big producers to single consumers)

According to the Dell engineers, the **packets storage buffer** is not common to all ports but **equally assigned to each port** (even the stacking and 10Gb ports !!??) and no flow control mechanism present on stacking links

→ High rates ports buffers very **quickly saturated** (with a few events)

Switches examples

Dell Force 10 S60

Public price : 9200 Euros -> **55200 euros/Telescope**



BO = Buffer size = **1.25 GB**

The operating system of this switch allows **custom memory allocation**

➡ With bigger buffers for the higher rates ports, **saturation should occur later**

Switches examples

Dell Force 10 S60



Public price : 9200 Euros -> **55200 euros/Telescope**

BO = Buffer size = **1.25 GB**

The operating system of this switch allows **custom memory allocation**

➡ With bigger buffers for the higher rates ports, **saturation should occur later**

We will elaborate more sophisticated models

Premiers tests

Tests de synchronicité

Yassir Moudden (CEA Saclay) a développé un générateur de paquets sur 48 ports synchronisés à la nanoseconde.
Les premiers tests montrent une perte de paquets certainement due à un récepteur pas assez fiable.

La reproduction au CPPM avec le stimulateur montre les mêmes limites du récepteur. Par contre, en utilisant une carte Napatech NT4E en réception sur un port Gigabit, aucune perte de paquets jusqu'à 40 paquets consécutifs (à 1 Gb/s) sur 42 ports.

Tests de profondeur de buffer

Des tests de profondeur de buffer vont être réalisés avec des scénarios aléatoires lancés sur le stimulateur.

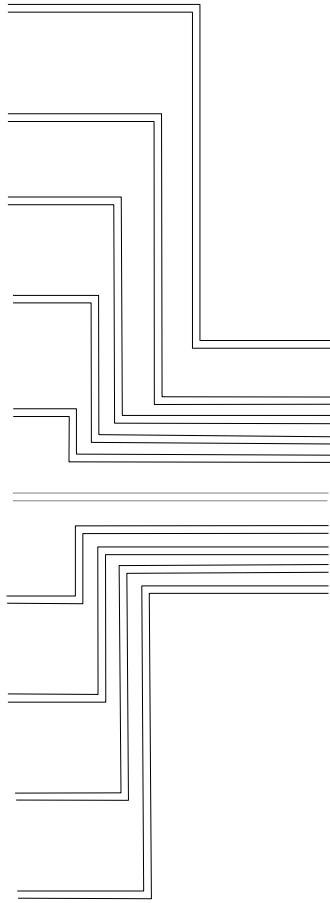
Camera server software prototype

Test configuration

10 servers



10 * 2 * 1Gb/s
Ethernet links



Powerconnect 6248

Or

Force10 S60

2 * 10Gb/s
Ethernet links
SFP+



Dell T7500 workstation

Event building measurements

Online Event building speed

Jumbo frames (8192 bytes) :

19.2 Gb/s (2.4 GB/s) with no loss

Average CPU usage :
160 % (1.6 cores/12)

~ 8000 events/s built

Regular frames (1024 bytes) :

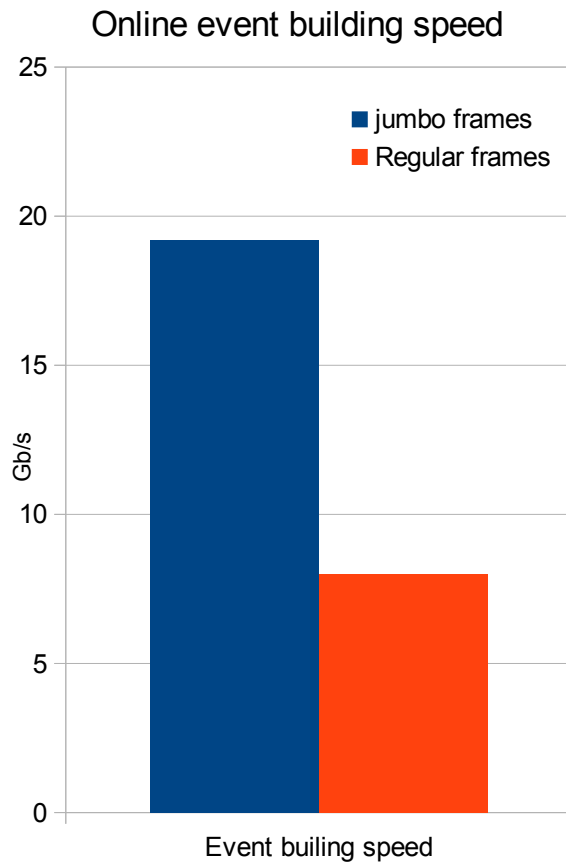
8 Gb/s (1 GB/s) with no loss

Average CPU usage :
170 % (1.7 cores/12)

~ 3300 events/s built

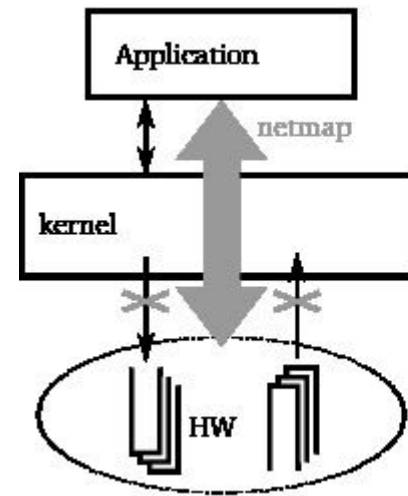
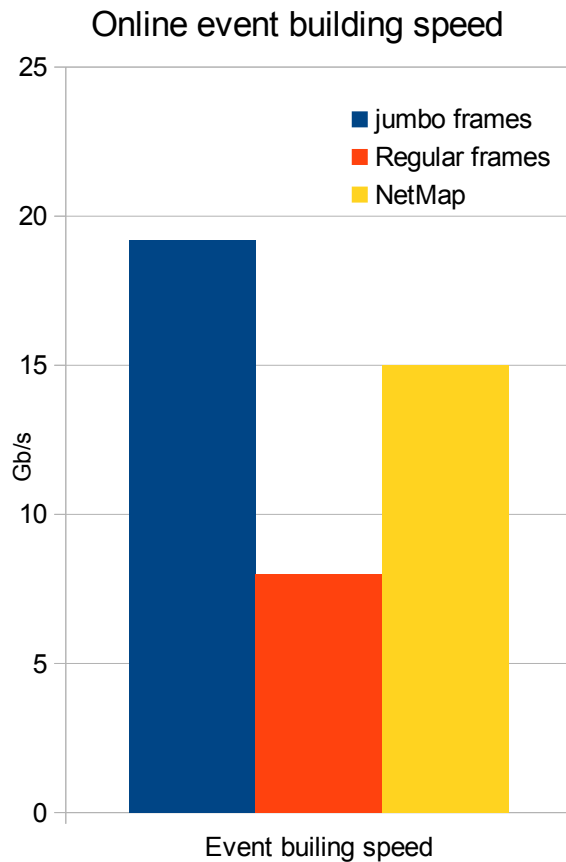
- Results obtained with standard libraries
- 300 stimulation nodes generate data
- Incoming data through two 10 Gb adapters
- Two « channels » used
- CPU load spread on 2 cores (on the same multiprocessor)

The regular frames issue



- Better results with jumbo frames
 - But must be compatible with electronics
- Must improve smaller packets reception
- Replace inefficient network software architectures by direct access to network components

The regular frames issue : Netmap



- We reach an event building at 15 Gb/s with 1024 bytes packets
- CPU usage ~ 170 %

Netmap : a novel framework for fast packet I/O

Luigi Rizzo, Università di Pisa, Italy

Proceedings of the 2012 USENIX Annual Technical Conference, June 2012

Last results for data reception only

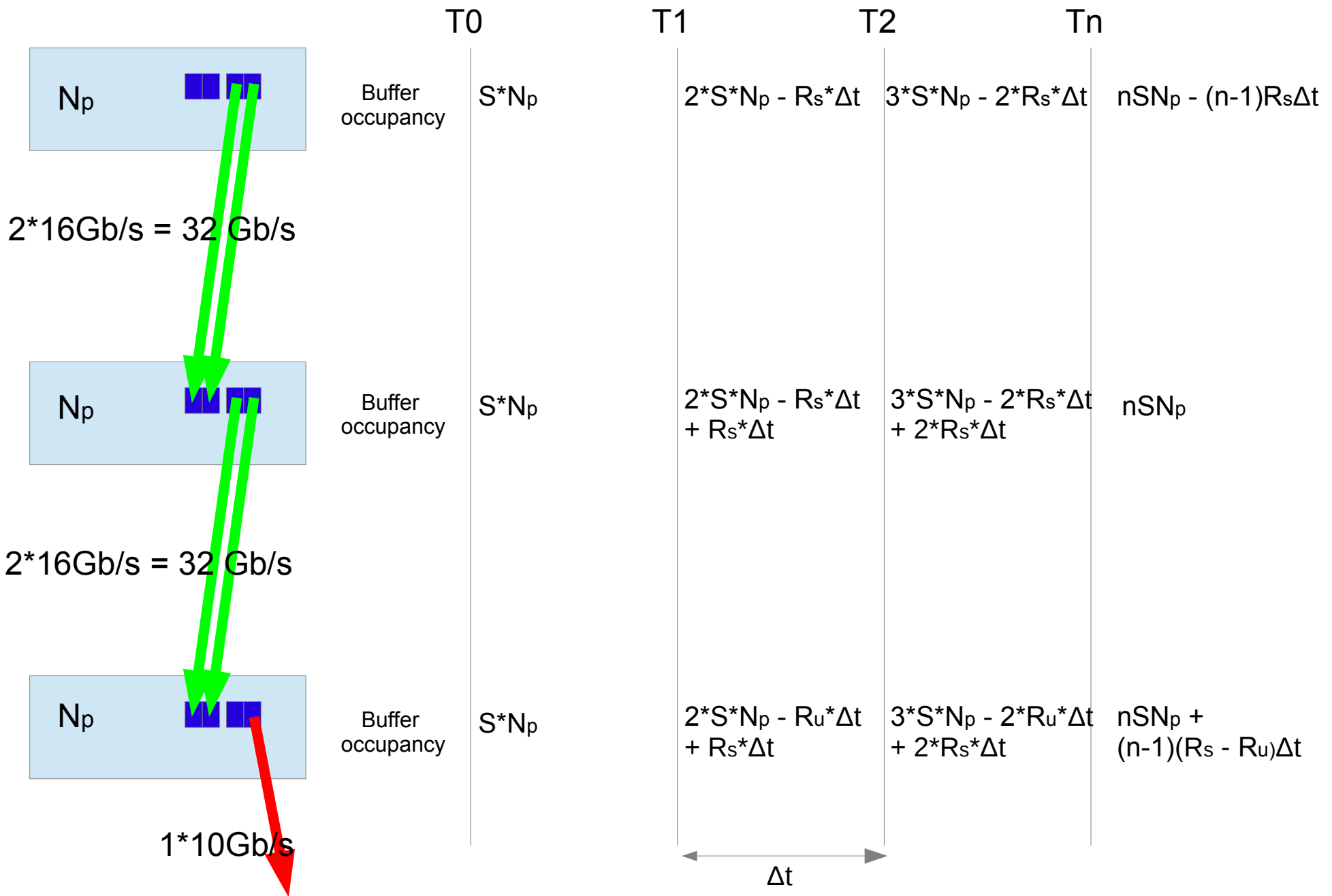
Last results with Netmap and stimulator:

Data reception at 19.2 Gb/s with small data loss (less than 0.03 %)

On two tasks (1 per 10Gb port) with isolated CPU, NUMA aware allocation...

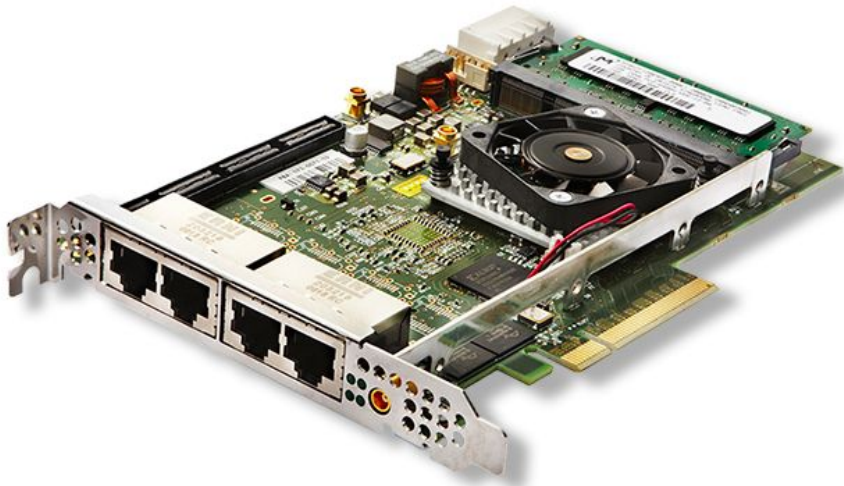
→ We should avoid data loss by distributing packets reception on several tasks (CPUs) with netmap and RSS

Backup



Stimulator : hardware time measurement

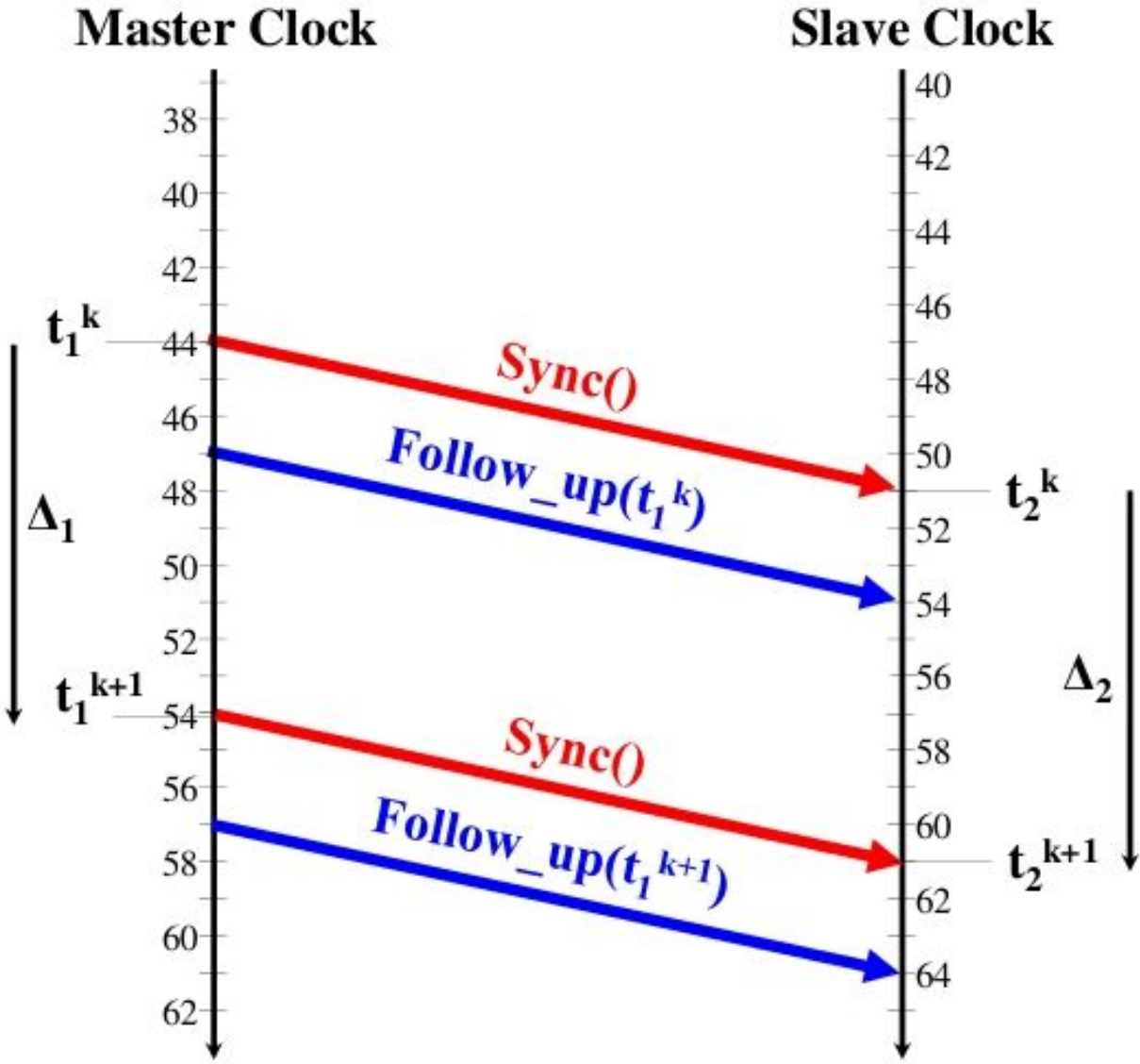
Napatech NT4E-4T



- 4 gigabit Ethernet ports
- Embedded memory for packets buffering
- Packets time stamping accuracy : 7 ns

Stimulator : Inter-boards synchronization

Clock drift determination



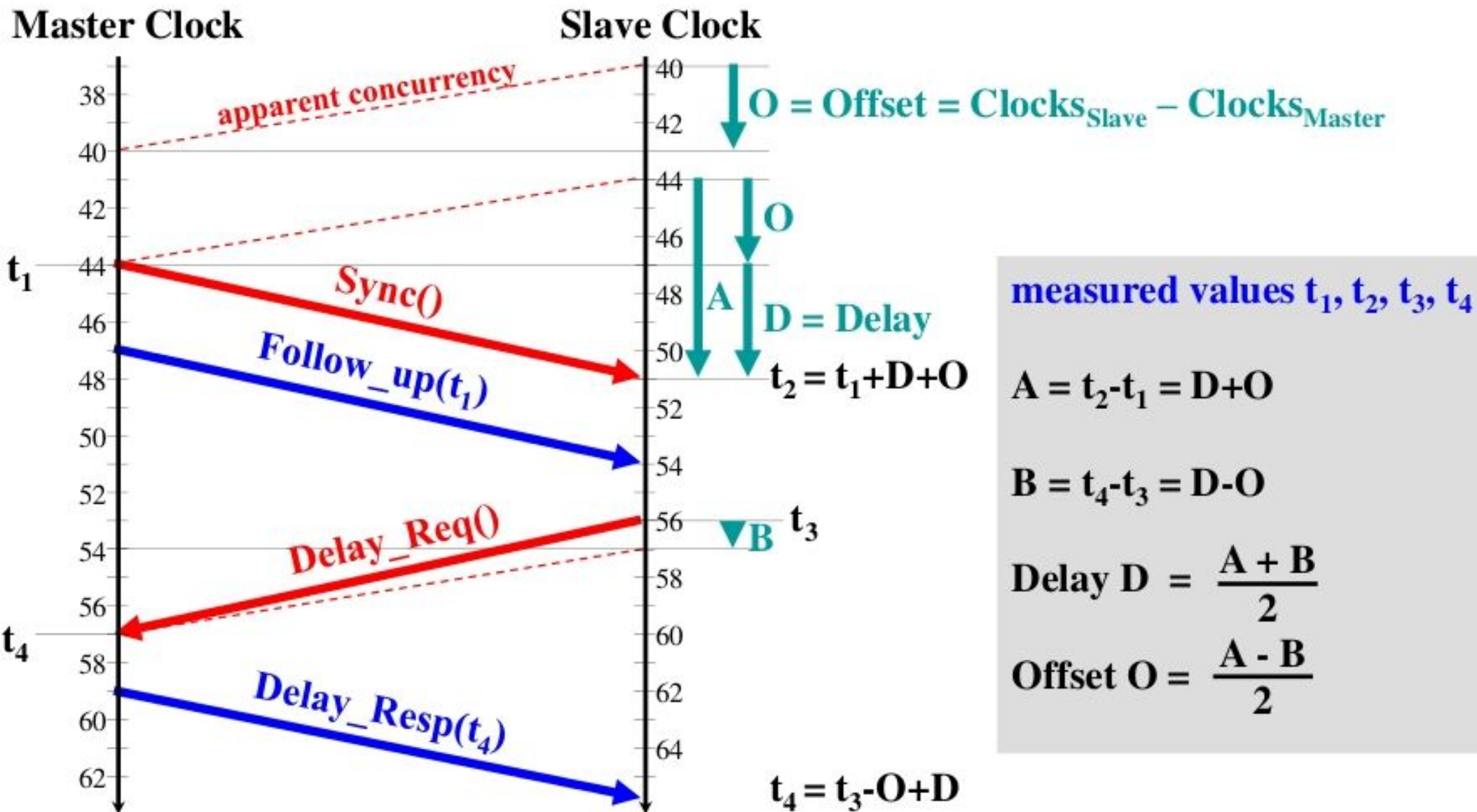
$$\Delta_1 = t_1^{k+1} - t_1^k$$

$$\Delta_2 = t_2^{k+1} - t_2^k$$

$$\text{Drift} = \frac{\Delta_2 - \Delta_1}{\Delta_2}$$

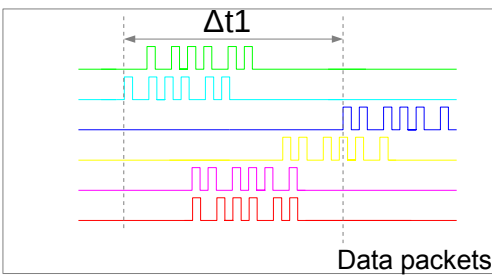
Stimulator : Inter-boards synchronization

Clock offset determination

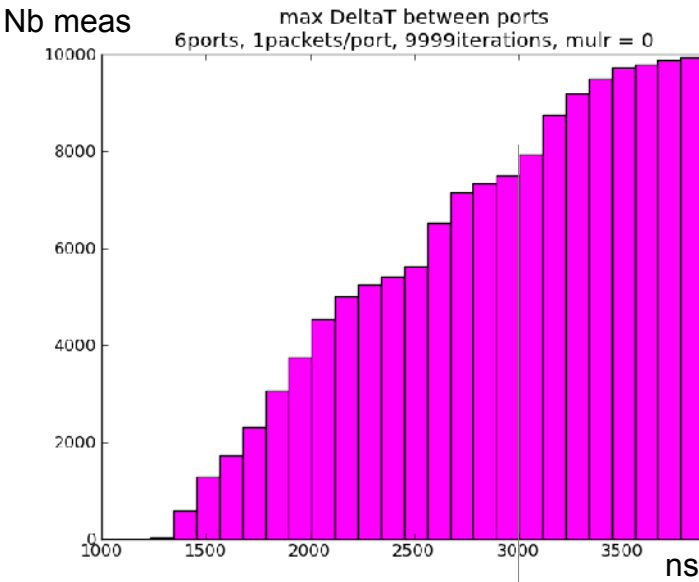


Stimulator : Board evaluation

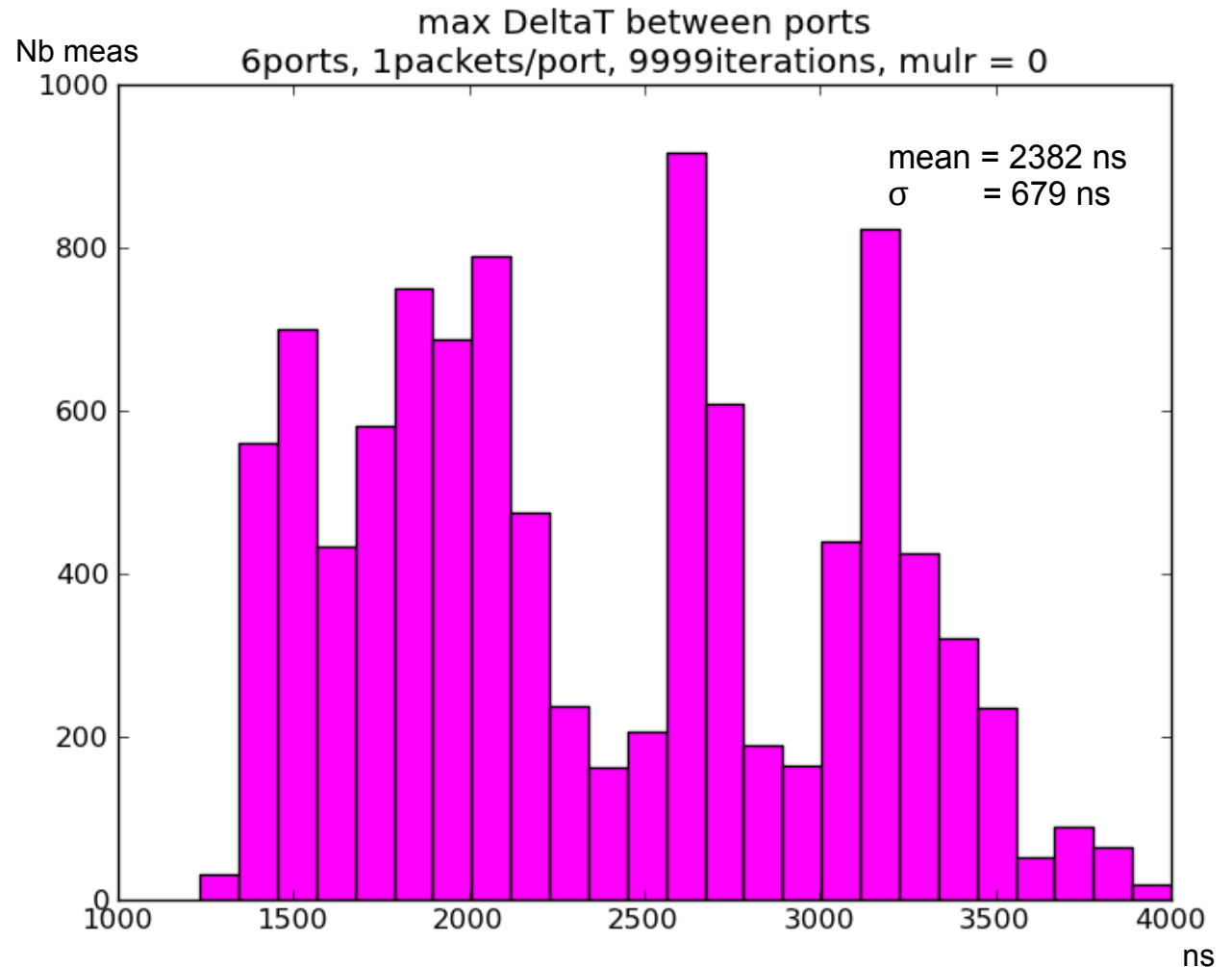
Δt_1 : time difference between the first packet sent and the last one on one SBC



10000 measurements **6 ports** 1 packet per port



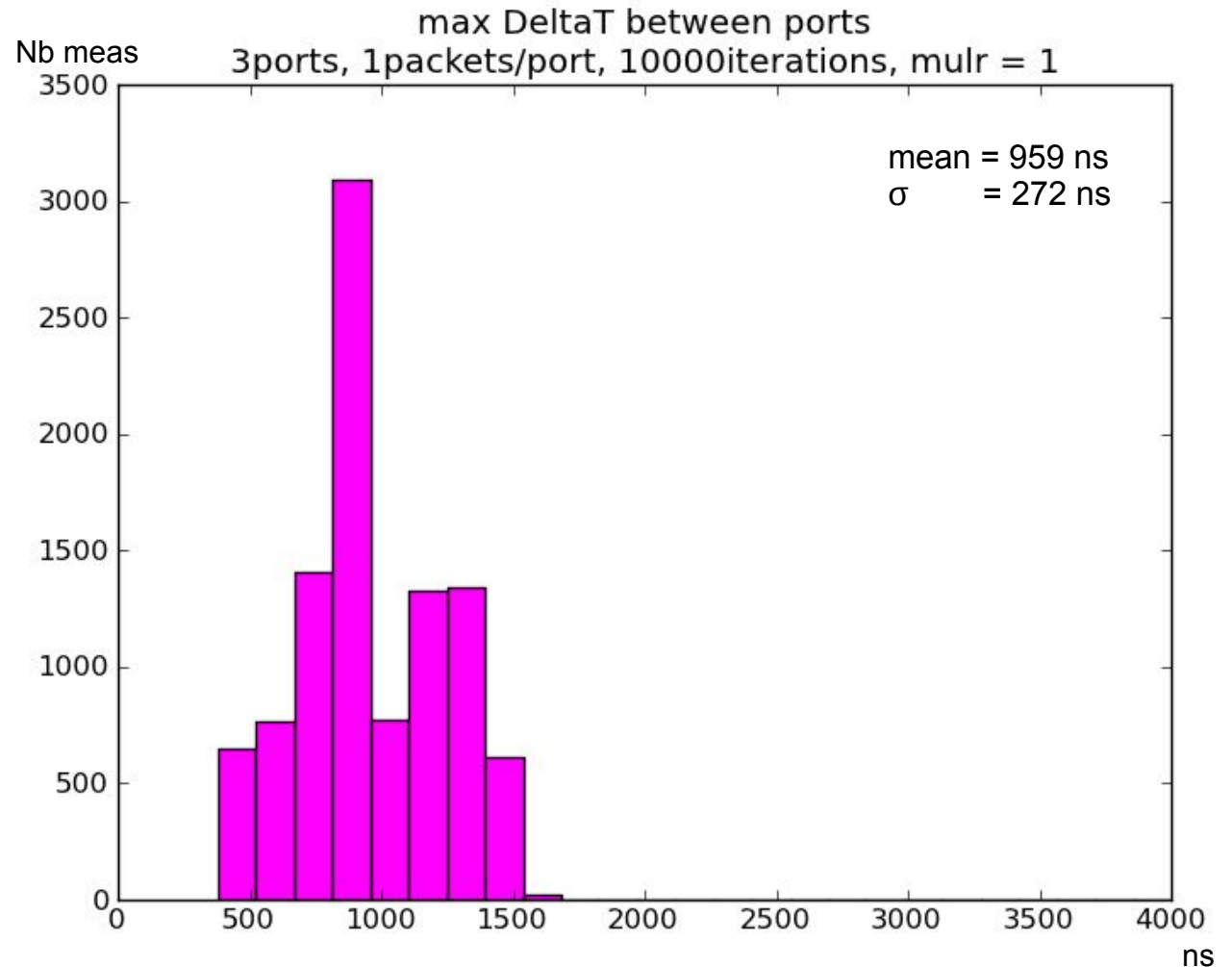
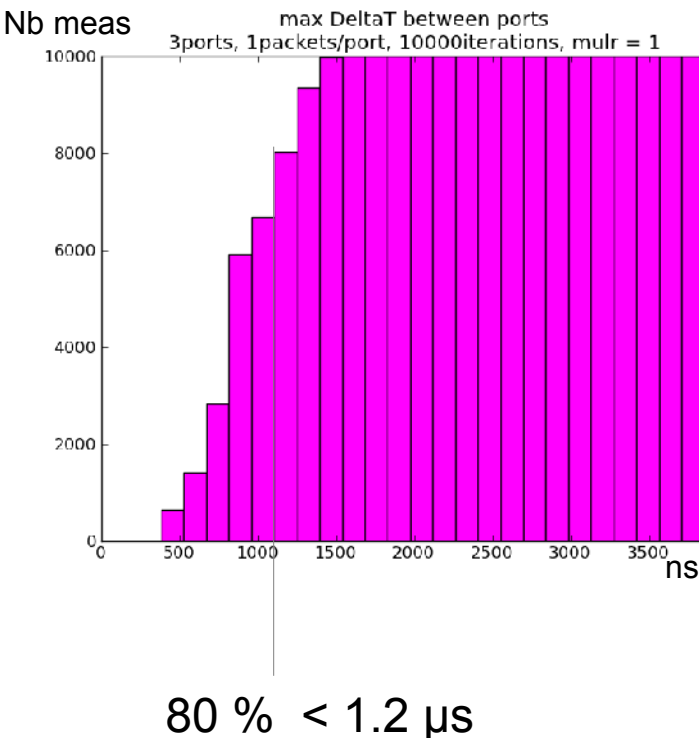
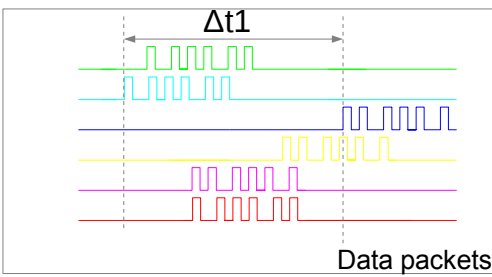
80 % < 3 μ s



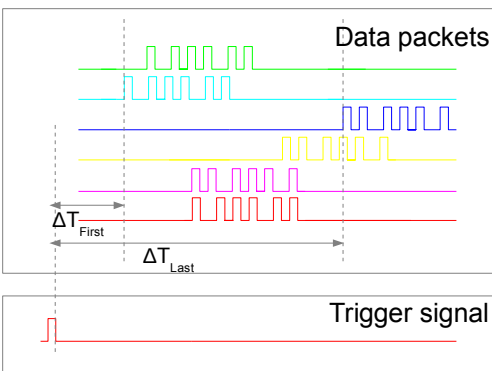
Stimulator : Board evaluation

Δt_1 : time difference between the first packet sent and the last one on one SBC

10000 measurements **3 ports** 1 packet per port



Stimulator : Board evaluation

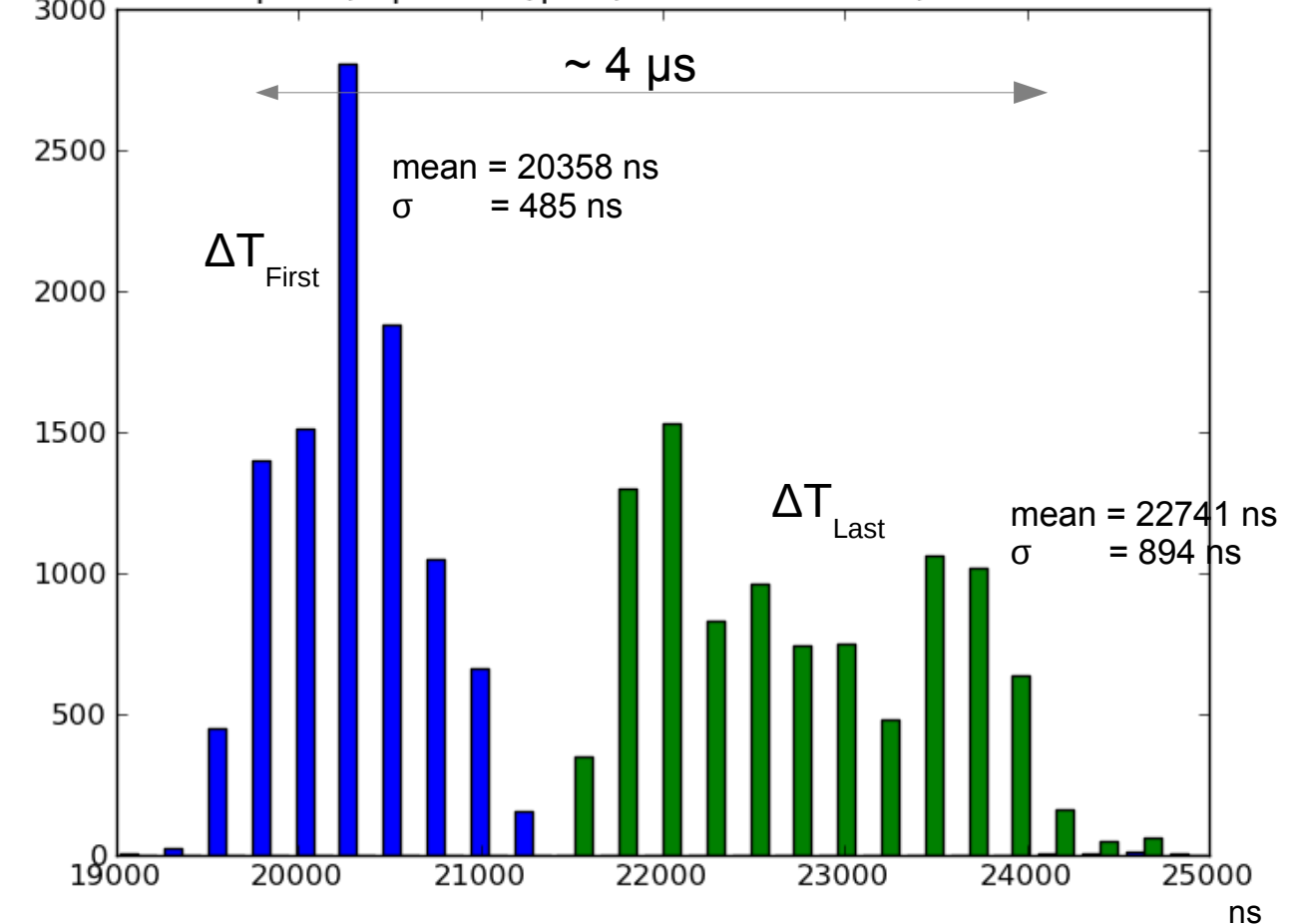


ΔT_{First} : time difference between the send command and the first packet sent

ΔT_{Last} : time difference between the send command and the last packet sent

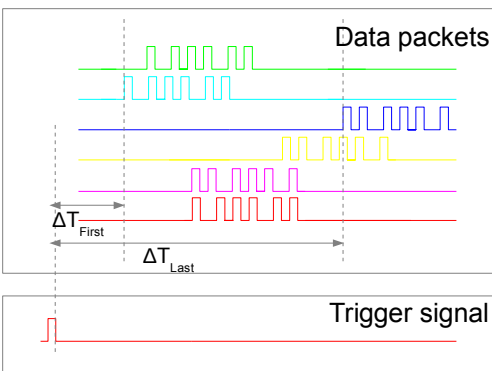
10000 measurements **6 ports** 1 packet per port

DeltaT between trigger and first packet (blue) last packet (green)
6ports, 1packets/port, 10000iterations, mulr = 0



Almost all paquets sent
in a ~ 4 μ s interval

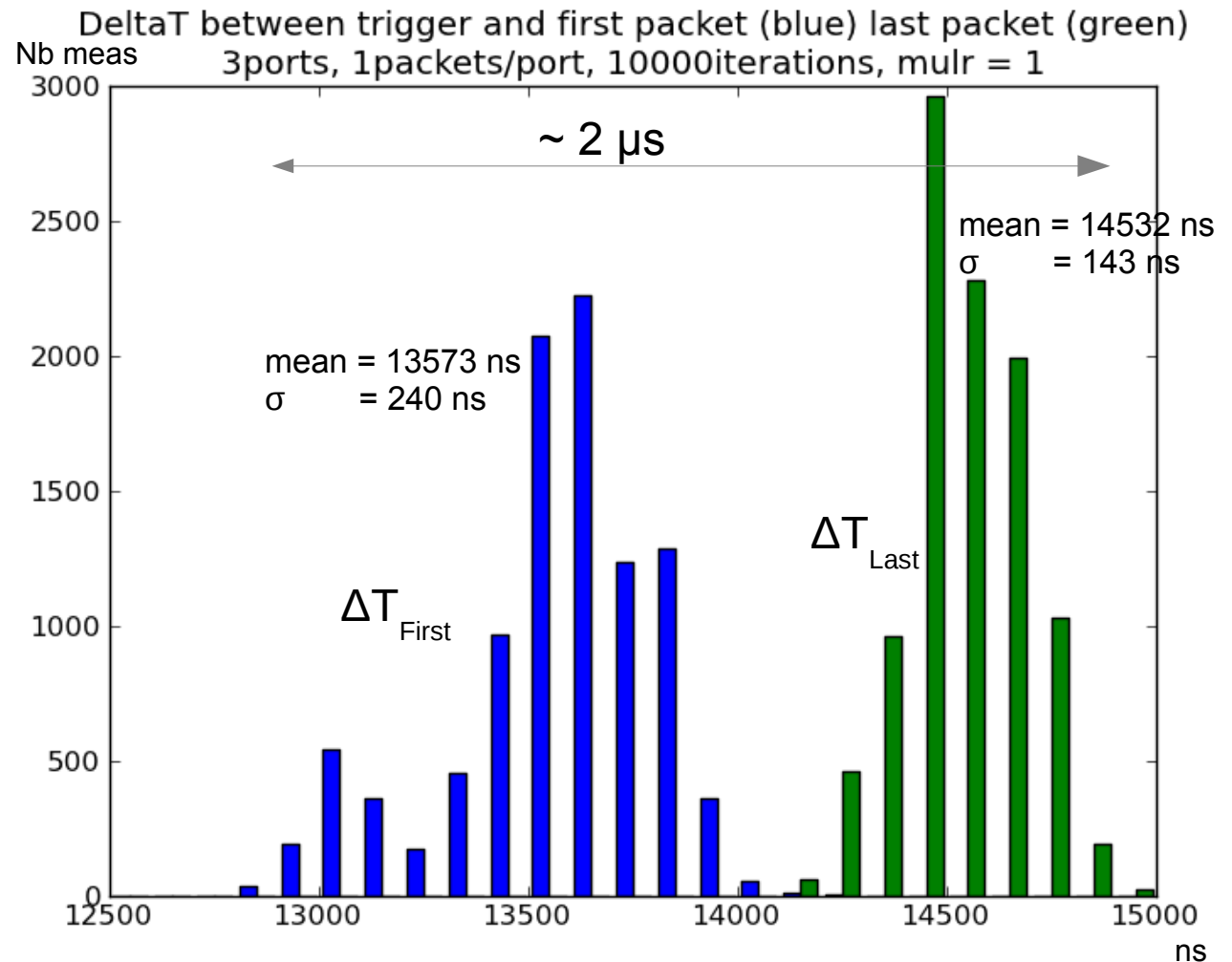
Stimulator : Board evaluation



ΔT_{First} : time difference between the send command and the first packet sent

ΔT_{Last} : time difference between the send command and the last packet sent

10000 measurements **3 ports** 1 packet per port



Almost all paquets sent
in a $\sim 2 \mu\text{s}$ interval

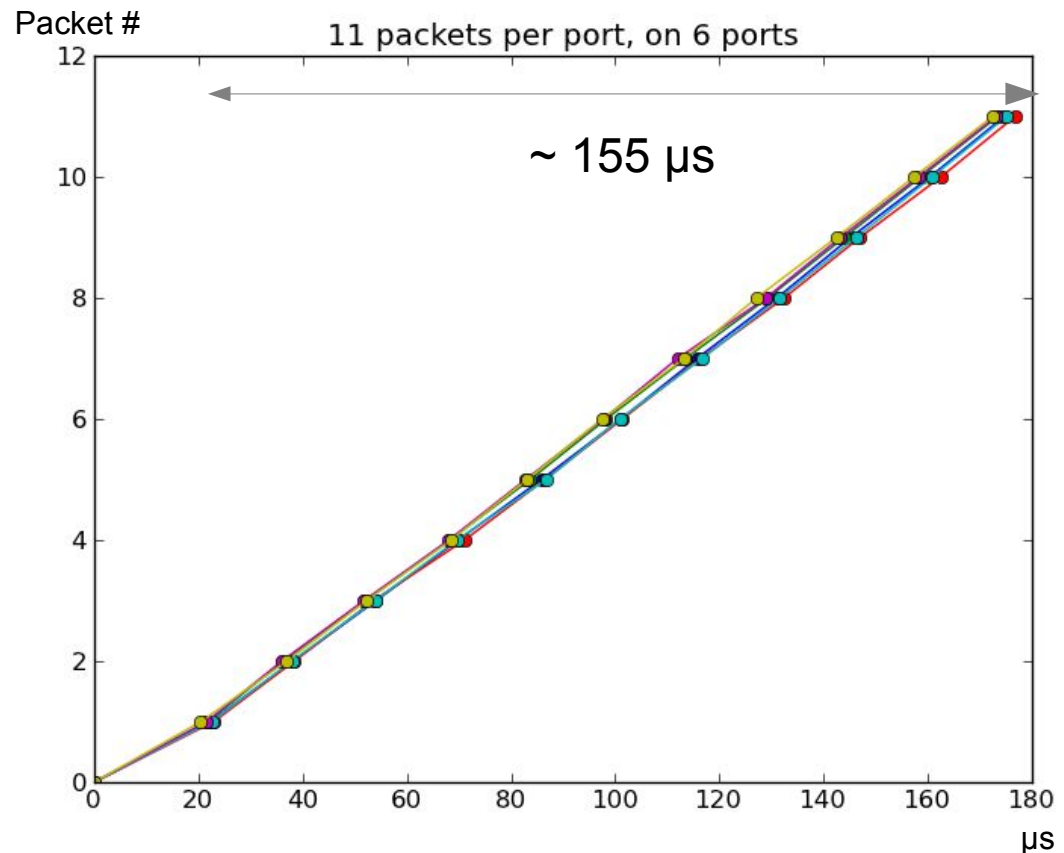
Stimulator : Board evaluation

Sending of 11 consecutive packets on each port

10000 measurements **6 ports** 11 packets per port

11 packets transferred in 155 μ s

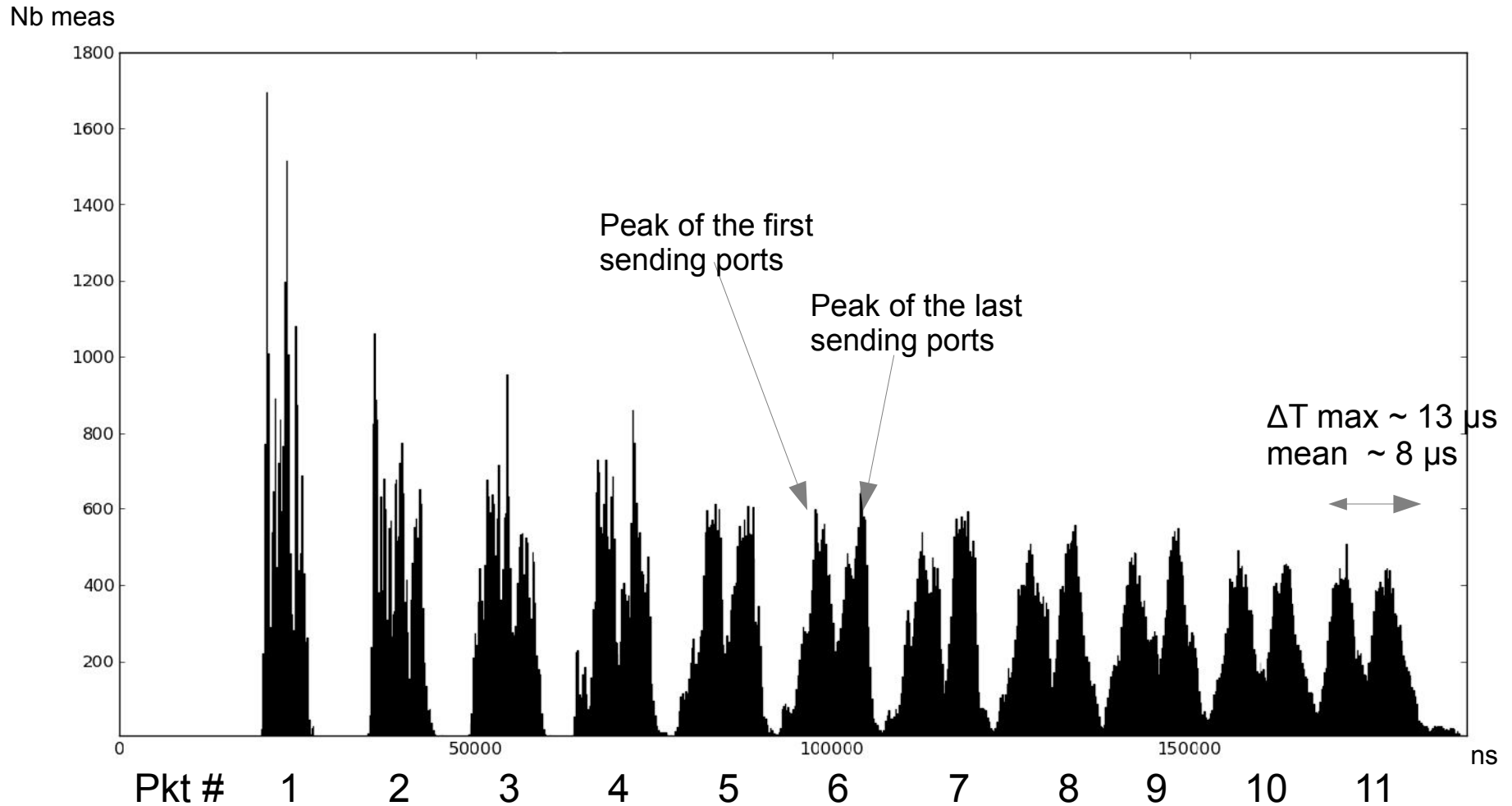
Data flow : 585 Mb/s



Stimulator : Board evaluation

Sending of 11 consecutive packets on each port :
First and last sending ports on one SBC

10000 measurements **6 ports** 11 packets per port



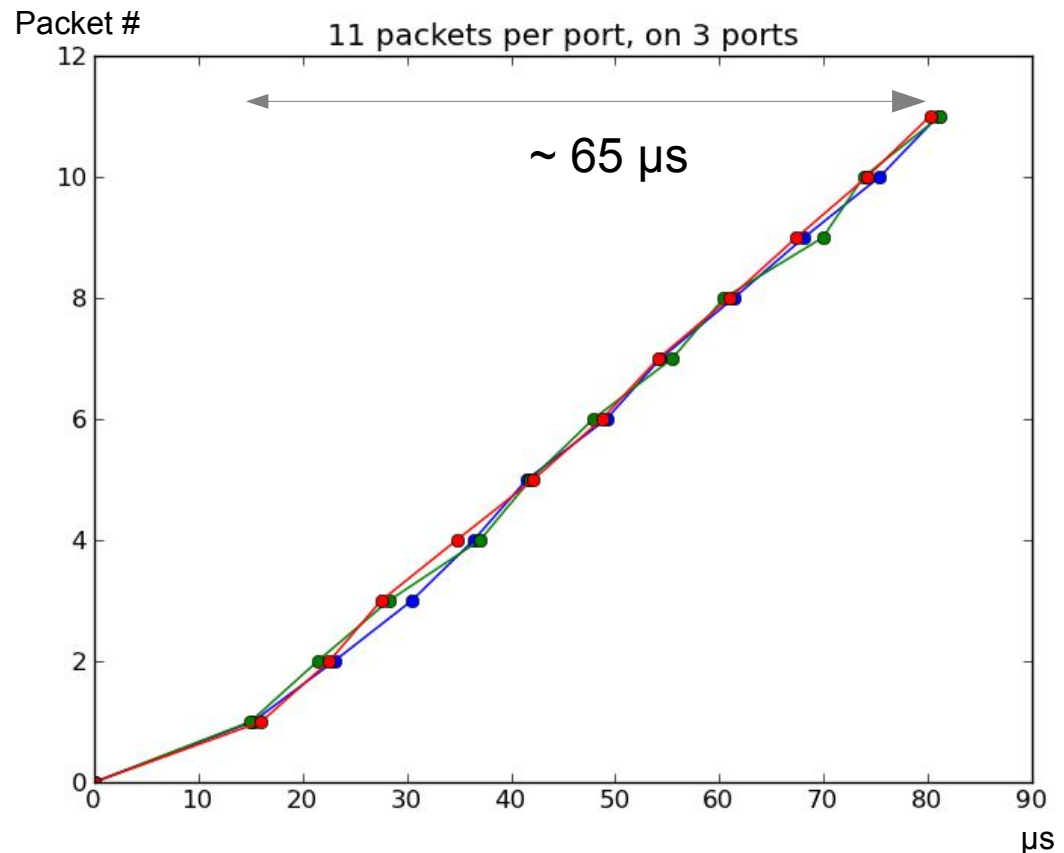
Stimulator : Board evaluation

Sending of 11 consecutive packets on each port

10000 measurements **3 ports** 11 packets per port

11 packets transferred in 65 μ s

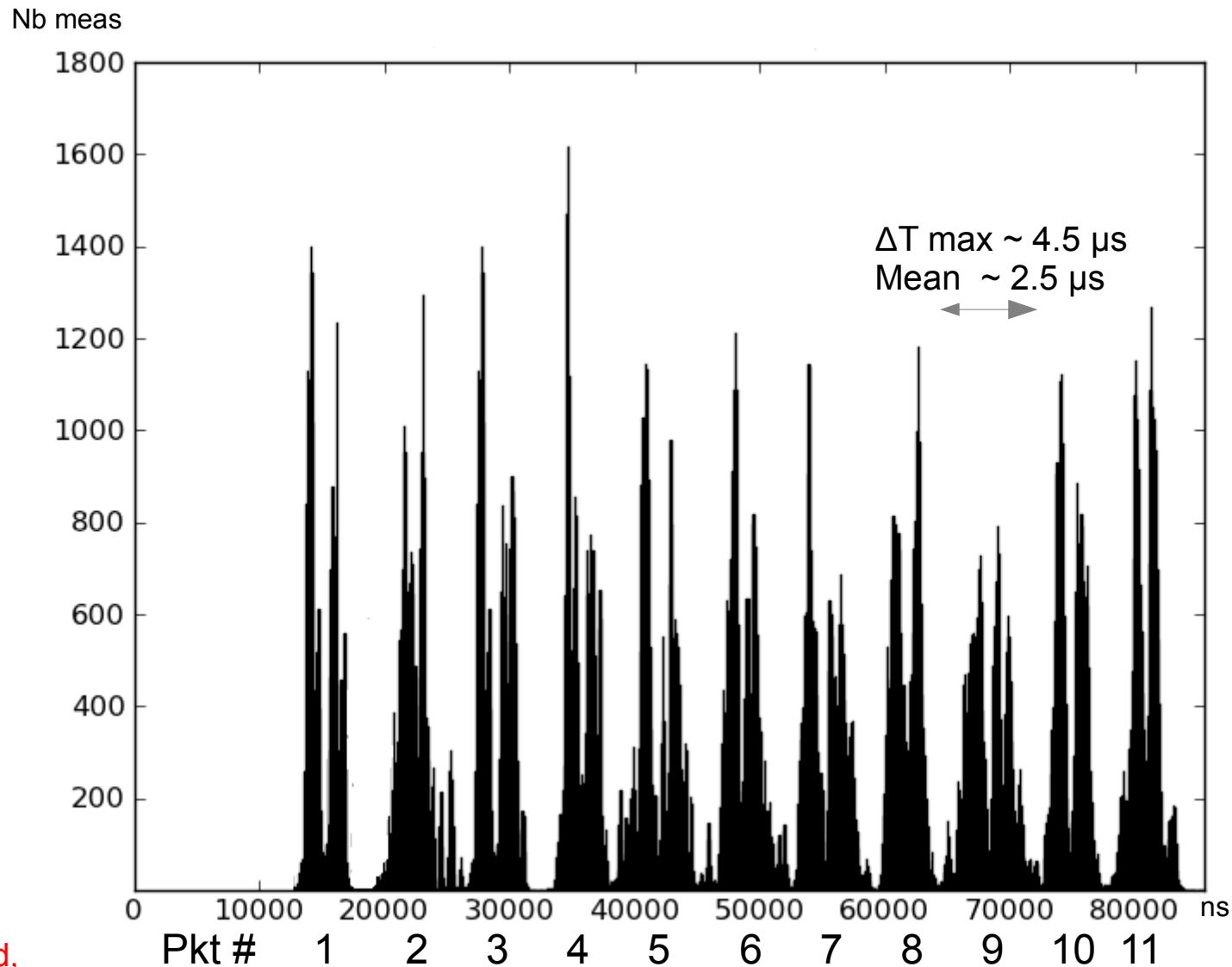
Data flow : 1.4 Gb/s



Stimulator : Board evaluation

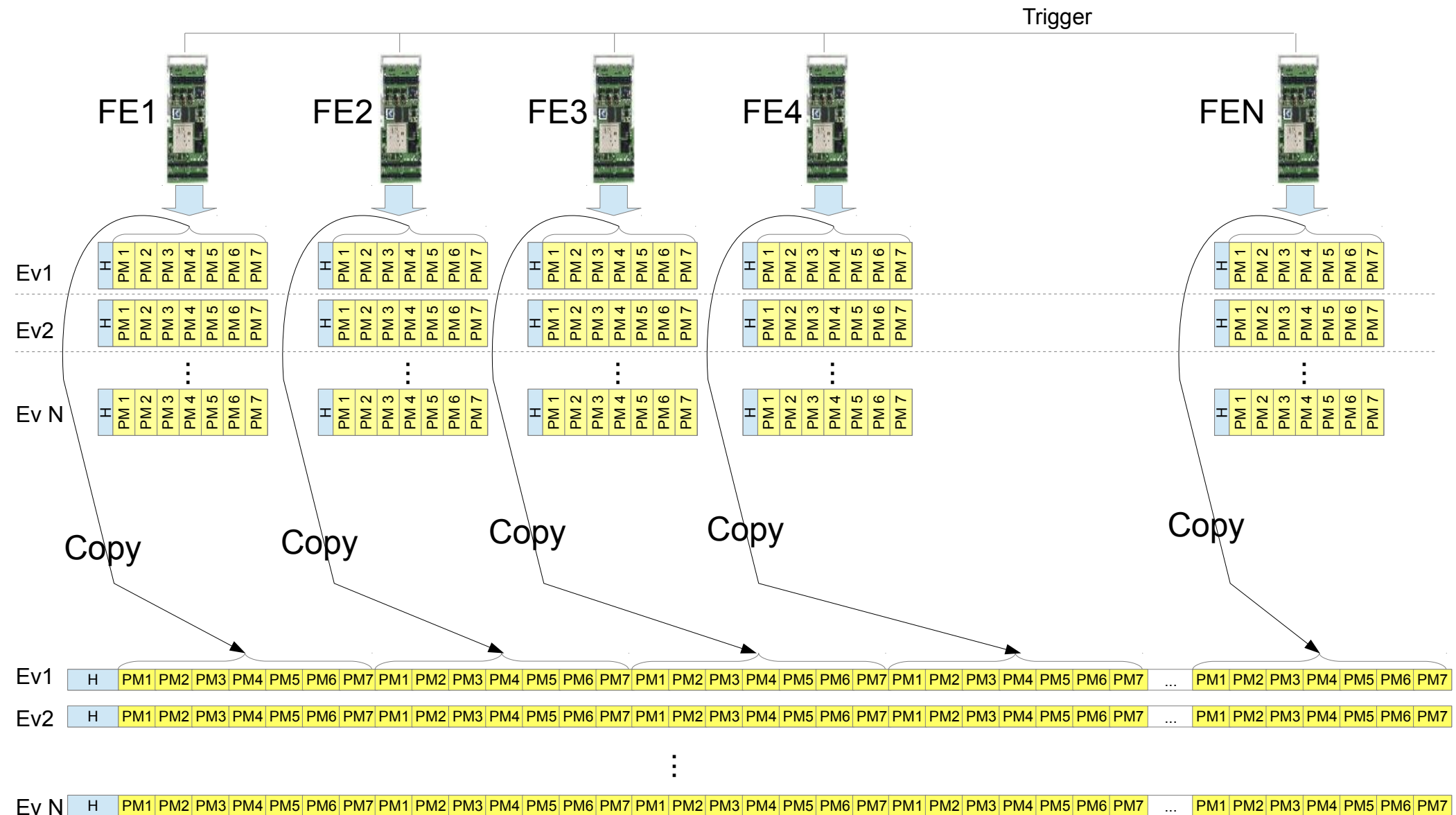
Sending of 11 consecutive packets on each port :
First and last sending ports on one SBC

10000 measurements **3 ports** 11 packets per port

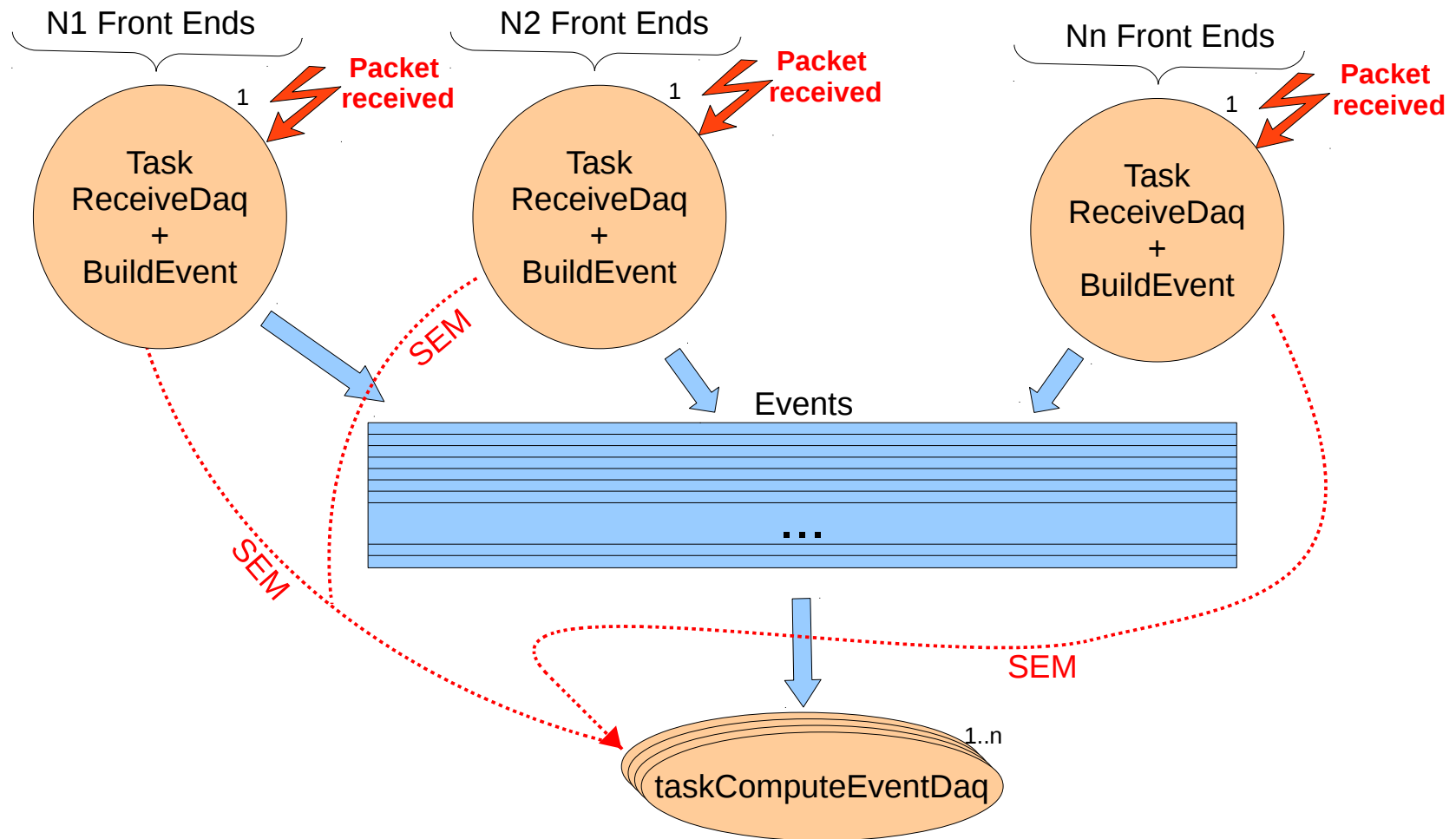


Journées in
With EVOC Board,
Results slightly better with Jetway board

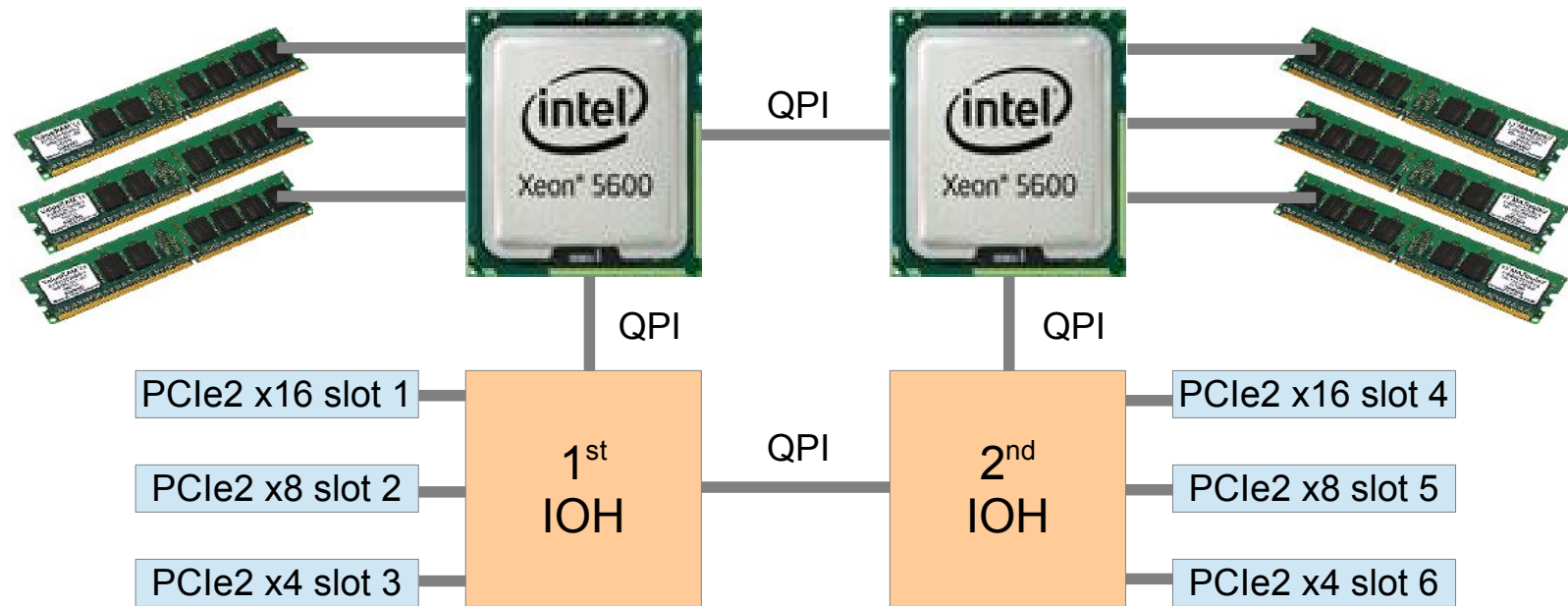
Event builder



Software overview : 2nd architecture



Non Uniform Memory Access



QPI @ 6.4 GT/s bandwidth < DDR3-1333 memory bandwidth