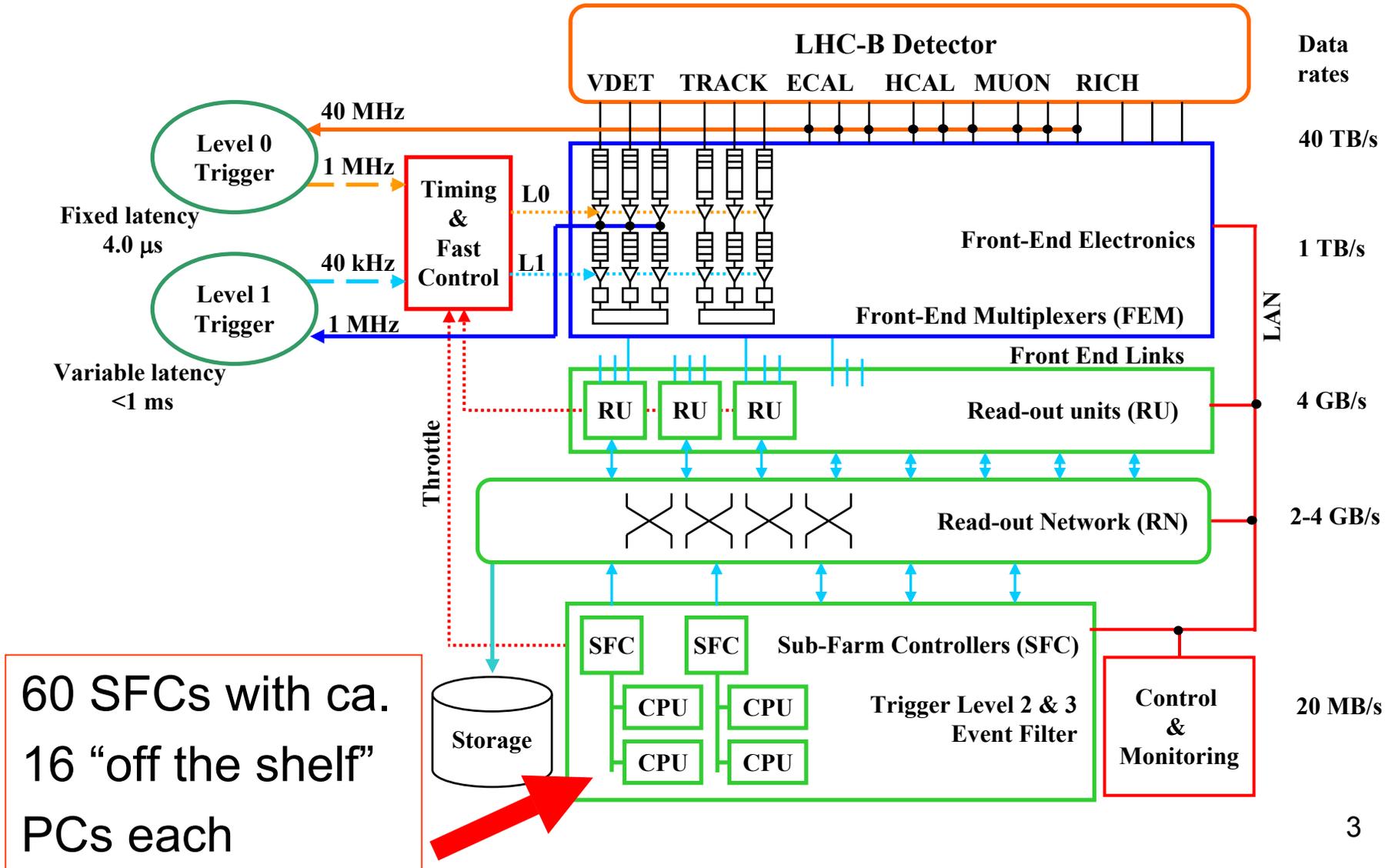# Network Performance Optimisation and Load Balancing
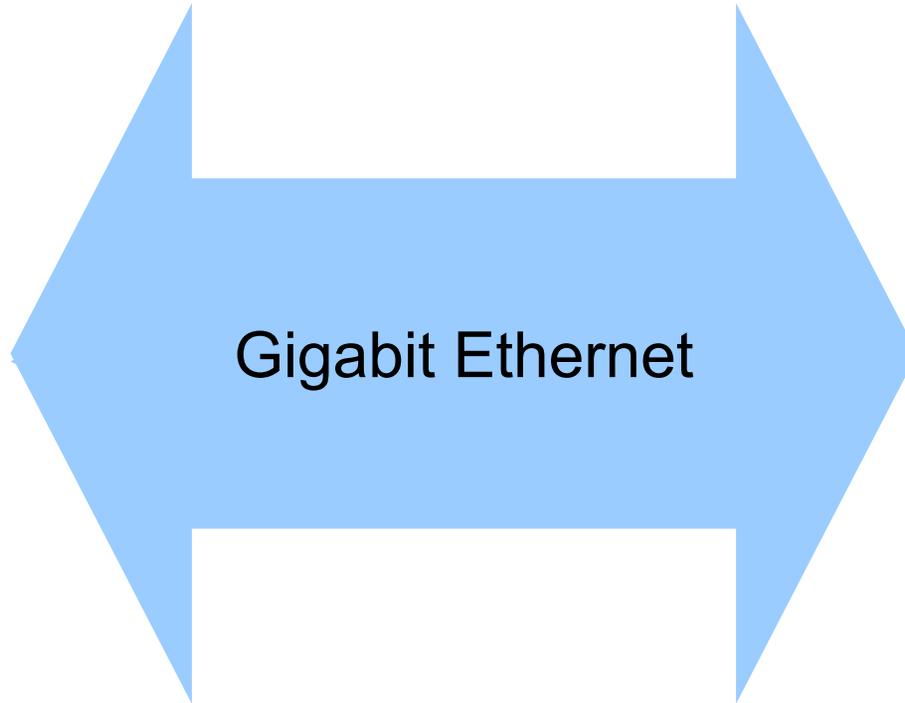
Wulf Thannhaeuser

# Network Performance Optimisation

# Network Optimisation: Where?



**LHC-B Detector**

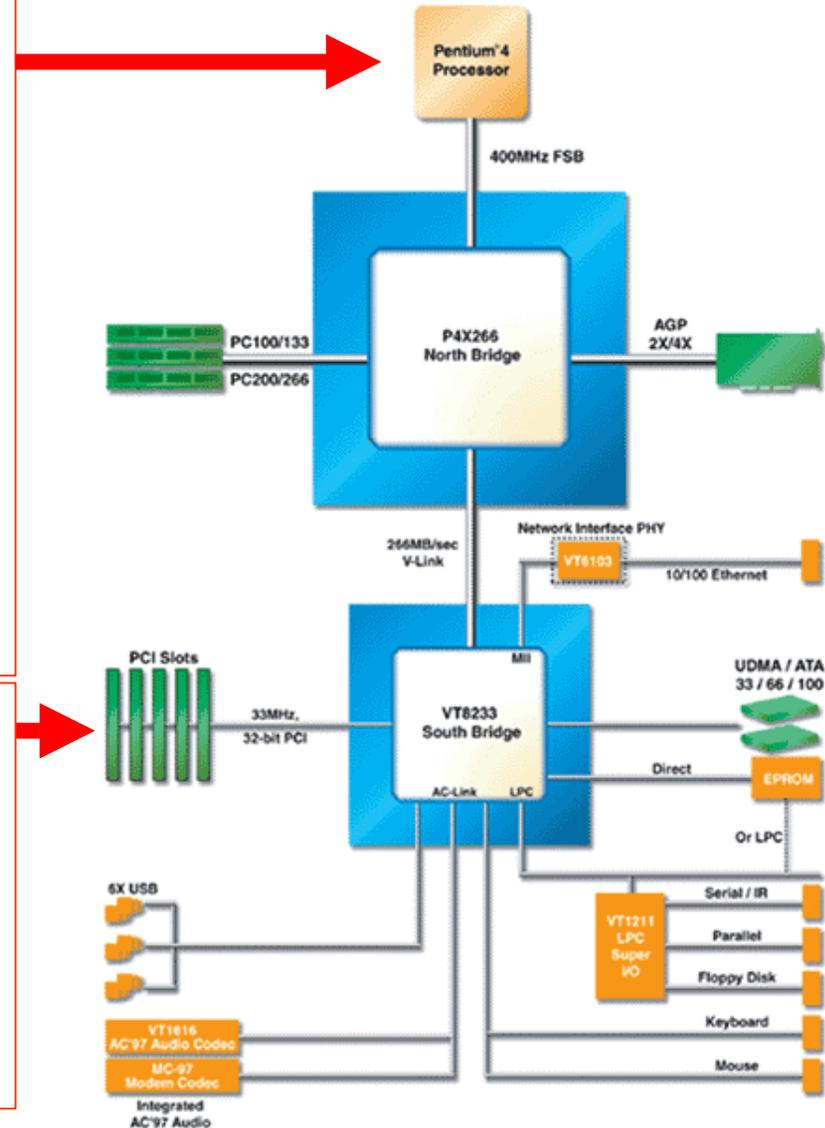VDET    TRACK    ECAL    HCAL    MUON    RICH

**Data rates**

40 MHz

**Level 0 Trigger**

1 MHz

Fixed latency 4.0 μs

**Timing & Fast Control**

L0

L1

40 kHz

**Level 1 Trigger**

1 MHz

Variable latency <1 ms

Throttle

**Front-End Electronics**

**Front-End Multiplexers (FEM)**

Front End Links

RU    RU    RU    **Read-out units (RU)**

**Read-out Network (RN)**

SFC    SFC    **Sub-Farm Controllers (SFC)**

CPU    CPU

CPU    CPU

**Trigger Level 2 & 3 Event Filter**

Storage

**Control & Monitoring**

LAN

40 TB/s

1 TB/s

4 GB/s

2-4 GB/s

20 MB/s

60 SFCs with ca. 16 "off the shelf" PCs each

3

# Network Optimisation: Why?

Gigabit Ethernet

| | | |
|---|---|---|
| Ethernet Speed: | 10 Mb/s | |
| Fast Ethernet Speed: | 100 Mb/s | |
| Gigabit Ethernet Speed: | 1000 Mb/s | |
| (considering full-duplex: | 2000 Mb/s) | |

# Network Optimisation: Why?

An "average" CPU might not be able to process such a huge amount of data packets per second:
-TCP/IP Overhead
-Context Switching
-Packet Checksums

An "average" PCI Bus is 33 MHz, 32-bit wide.
Theory: 1056 Mbit/s
Actually: ca. 850 Mbit/s
(PCI overhead, burstsize)

# Network Optimisation: How?

An "average" CPU might not be able to process such a huge amount of data packets per second:
-TCP/IP Overhead
-Context Switching
-Packet Checksums

An "average" PCI Bus is 33 MHz, 32-bit wide.
Theory: 1056 Mbit/s
Actually: ca. 850 Mbit/s
(PCI overhead, burstsize)

Reduce per packet Overhead:
*Replace TCP with UDP*

# TCP / UDP Comparison

- <u>TCP (Transfer Control Protocol):</u>
  - connection-oriented protocol
  - full-duplex
  - messages received in order, no loss or duplication

  $\Rightarrow$ *reliable but with overheads*

- <u>UDP (User Datagram Protocol):</u>
  - messages called "datagrams"
  - messages may be lost or duplicated
  - messages may be received out of order

  $\Rightarrow$ *unreliable but potentially faster*

# Network Optimisation: How?

An "average" CPU might not be able to process such a huge amount of data packets per second:
- TCP/IP Overhead
- Context Switching
- Packet Checksums

An "average" PCI Bus is 33 MHz, 32-bit wide.
Theory: 1056 Mbit/s
Actually: ca. 850 Mbit/s
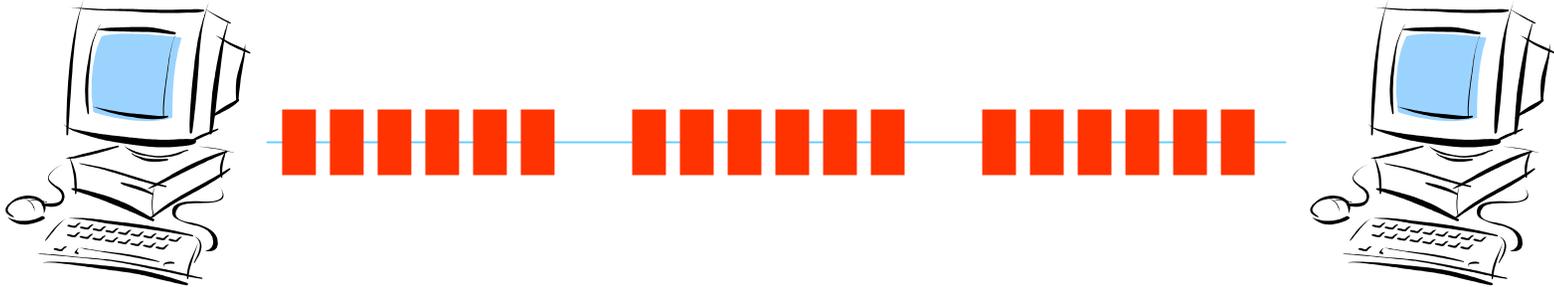(PCI overhead, burstsize)

Reduce per packet Overhead: *Replace TCP with UDP*

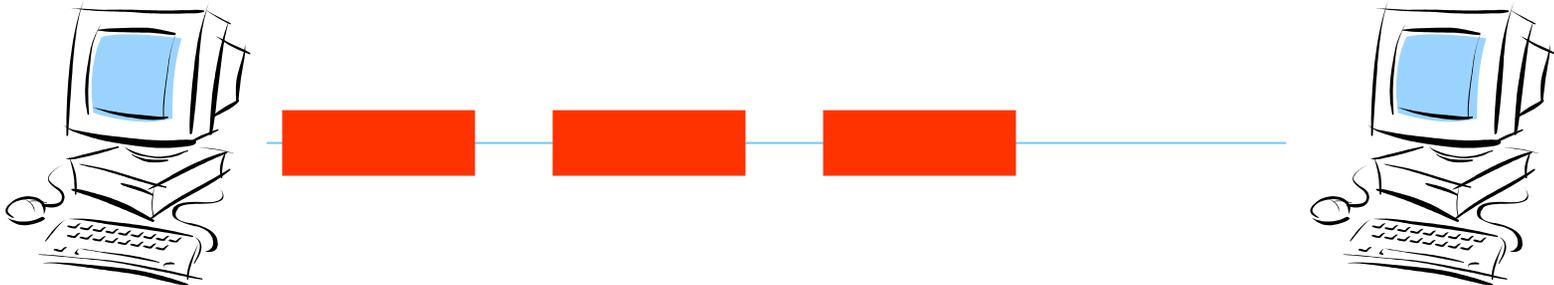Reduce number of packets: *Jumbo Frames*

# Jumbo Frames

Normal Ethernet

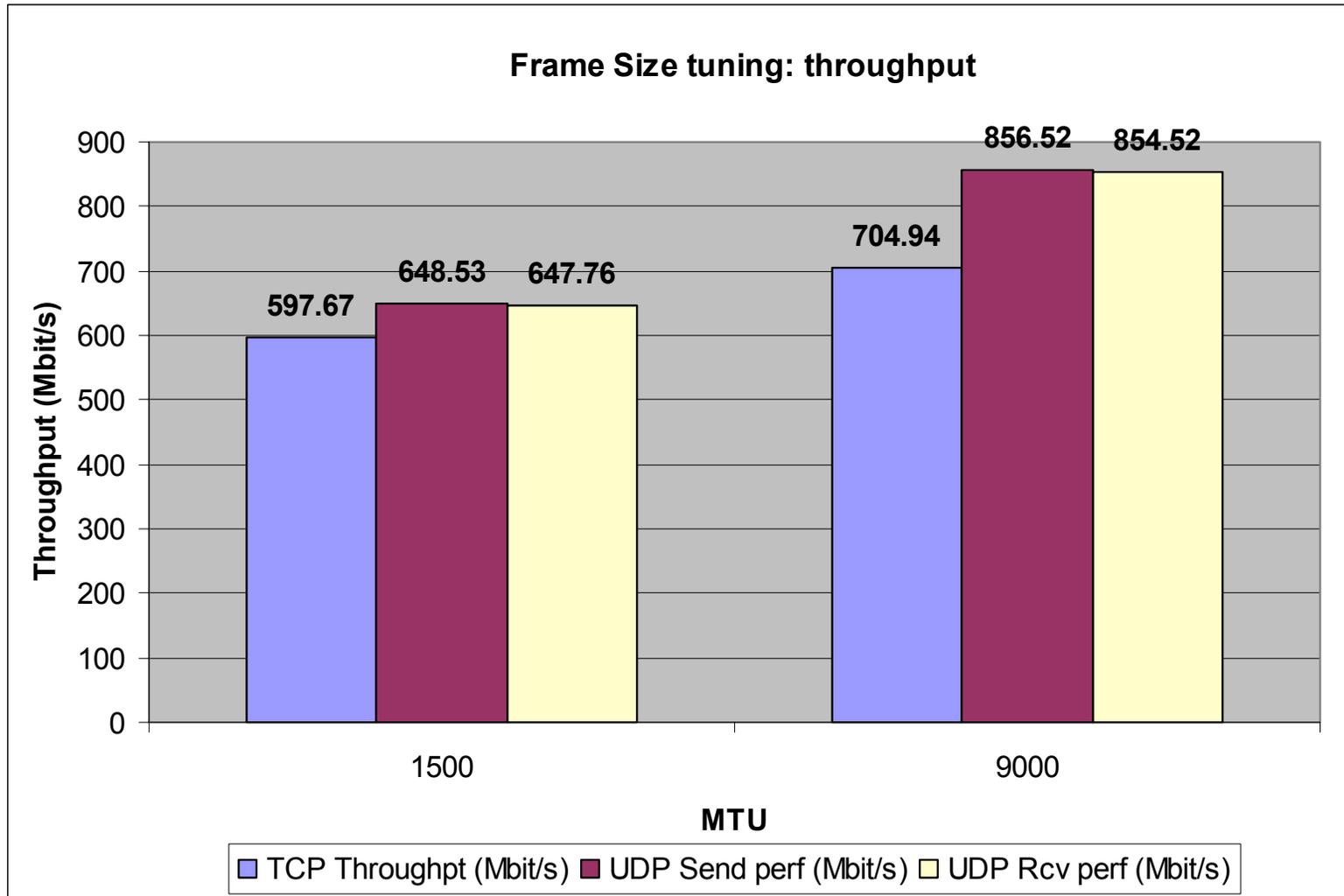Maximum Transmission Unit (MTU): 1500 bytes

Ethernet with Jumbo Frames MTU: 9000 bytes

# Test set-up

- <u>Netperf</u> is a benchmark for measuring network performance

- <u>The systems tested</u> were 800 and 1800 MHz Pentium PCs using (optical as well as copper) Gbit Ethernet NICs.

- <u>The network set-up</u> was always a simple point-to-point connection with a crossed twisted pair or optical cable.

- <u>Results were not always symmetric:</u>

  With two PCs of different performance, the benchmark results were usually better if data was sent from the slow PC to the fast PC, i.e. the receiving process is more expensive.

# Results with the optimisations so far



**Frame Size tuning: throughput**

# Network Optimisation: How?

An "average" CPU might not be able to process such a huge amount of data packets per second:
- TCP/IP Overhead
- Context Switching
- Packet Checksums

An "average" PCI Bus is 33 MHz, 32-bit wide.
Theory: 1056 Mbit/s
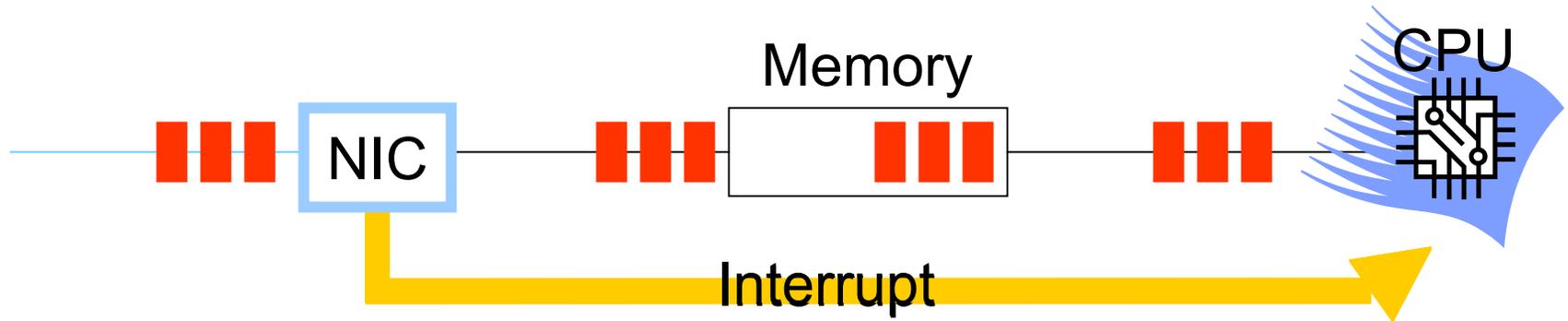Actually: ca. 850 Mbit/s
(PCI overhead, burstsize)

Reduce per packet Overhead: *Replace TCP with UDP*

Reduce number of packets: *Jumbo Frames*

Reduce context switches: *Interrupt coalescence*

# Interrupt Coalescence

Packet Processing without Interrupt coalescence:

Memory

CPU

NIC

Interrupt

Packet Processing <u>with</u> Interrupt coalescence:

Memory

CPU

NIC

Interrupt

# Interrupt Coalescence: Results



**Interrupt coalescence tuning: throughput
(Jumbo Frames deactivated)**

Throughput (Mbit/s)

- 526.09
- 587.43
- 556.76
- 633.72
- 666.43
- 593.94

32          1024

**us waited for new packets to arrive before interrupting CPU**

□ TCP Throughpt (Mbit/s)   ■ UDP Send perf (Mbit/s)   □ UDP Rcv perf (Mbit/s)

# Network Optimisation: How?

An "average" CPU might not be able to process such a huge amount of data packets per second:
-TCP /IP Overhead
-Context Switching
-Packet Checksums

An "average" PCI Bus is 33 MHz, 32-bit wide.
Theory: 1056 Mbit/s
Actually: ca. 850 Mbit/s
(PCI overhead, burstsize)

Reduce per packet Overhead:
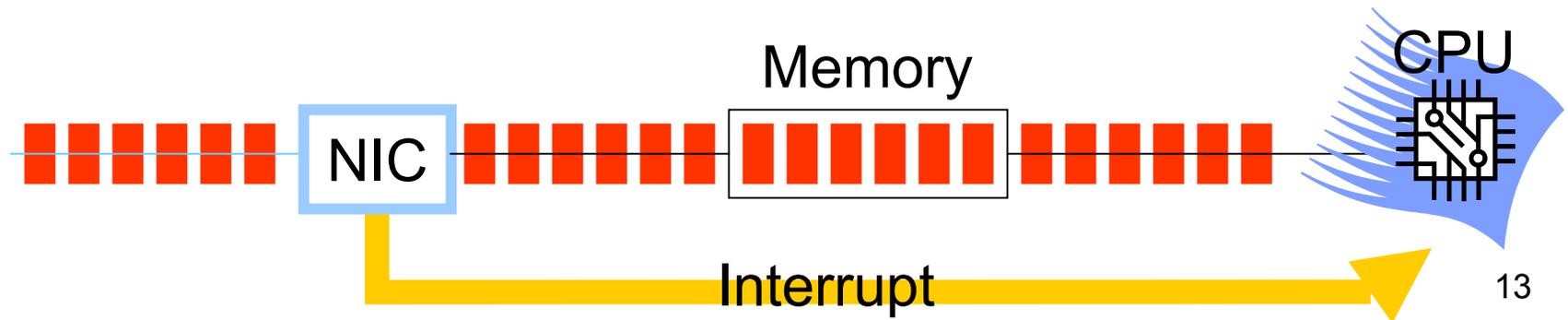*Replace TCP with UDP*

Reduce number of packets:
*Jumbo Frames*

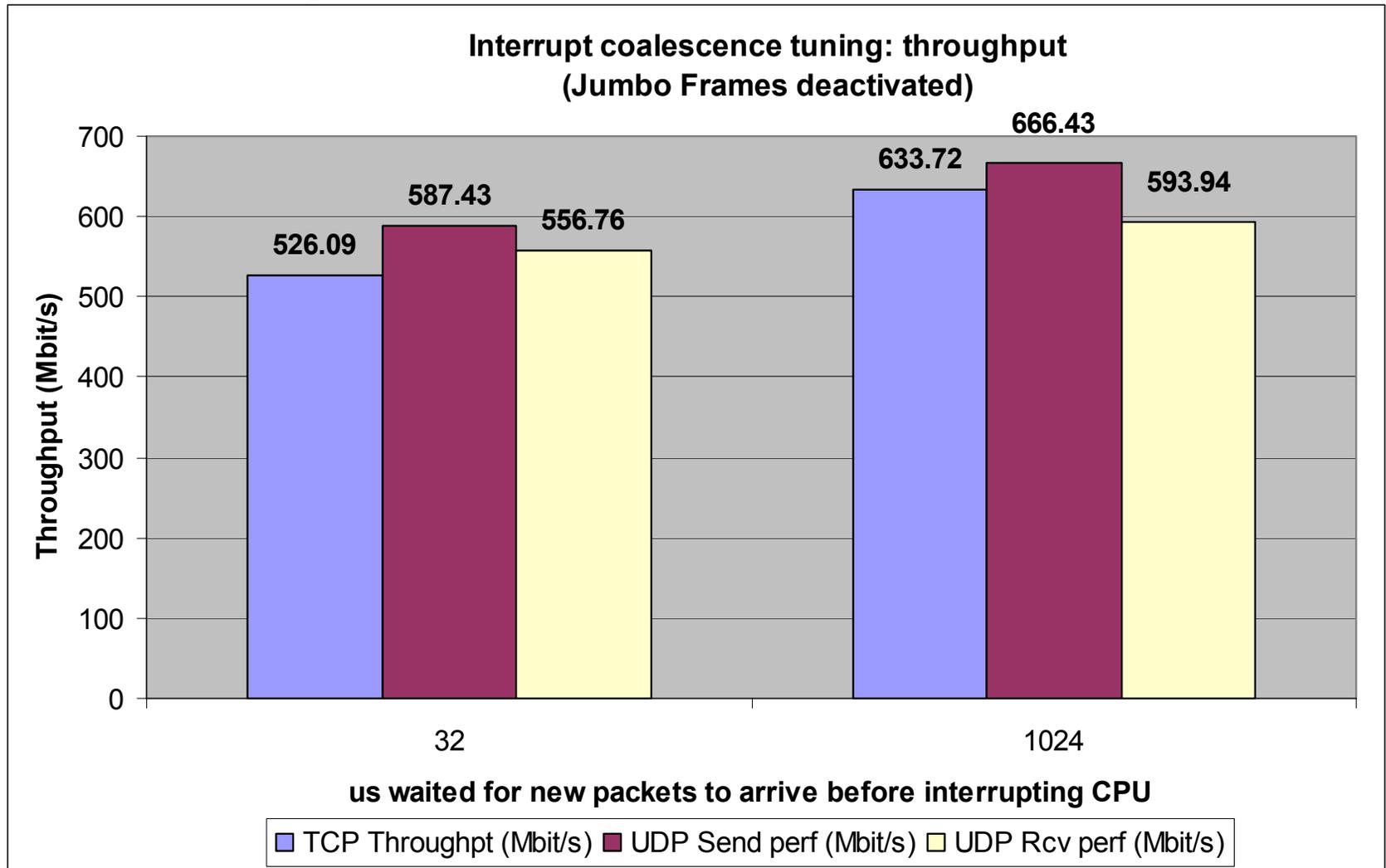Reduce context switches:
*Interrupt coalescence*

Reduce context switches:
*Checksum Offloading*

# Checksum Offloading

- A checksum is a number calculated from the data transmitted and attached to the tail of each TCP/IP packet.

- Usually the CPU has to recalculate the checksum for each received TCP/IP packet in order to compare it with the checksum in the tail of the packet to detect transmission errors.

- With checksum offloading, the NIC performs this task. Therefore <u>the CPU does not have to calculate the checksum</u> and can perform other operations in the meanwhile.

# Network Optimisation: How?

An "average" CPU might not be able to process such a huge amount of data packets per second:
- TCP/IP Overhead
- Context Switching
- Packet Checksums

An "average" PCI Bus is 33 MHz, 32-bit wide.
Theory: 1056 Mbit/s
Actually: ca. 850 Mbit/s
(PCI overhead, burstsize)

Reduce per packet Overhead:
*Replace TCP with UDP*
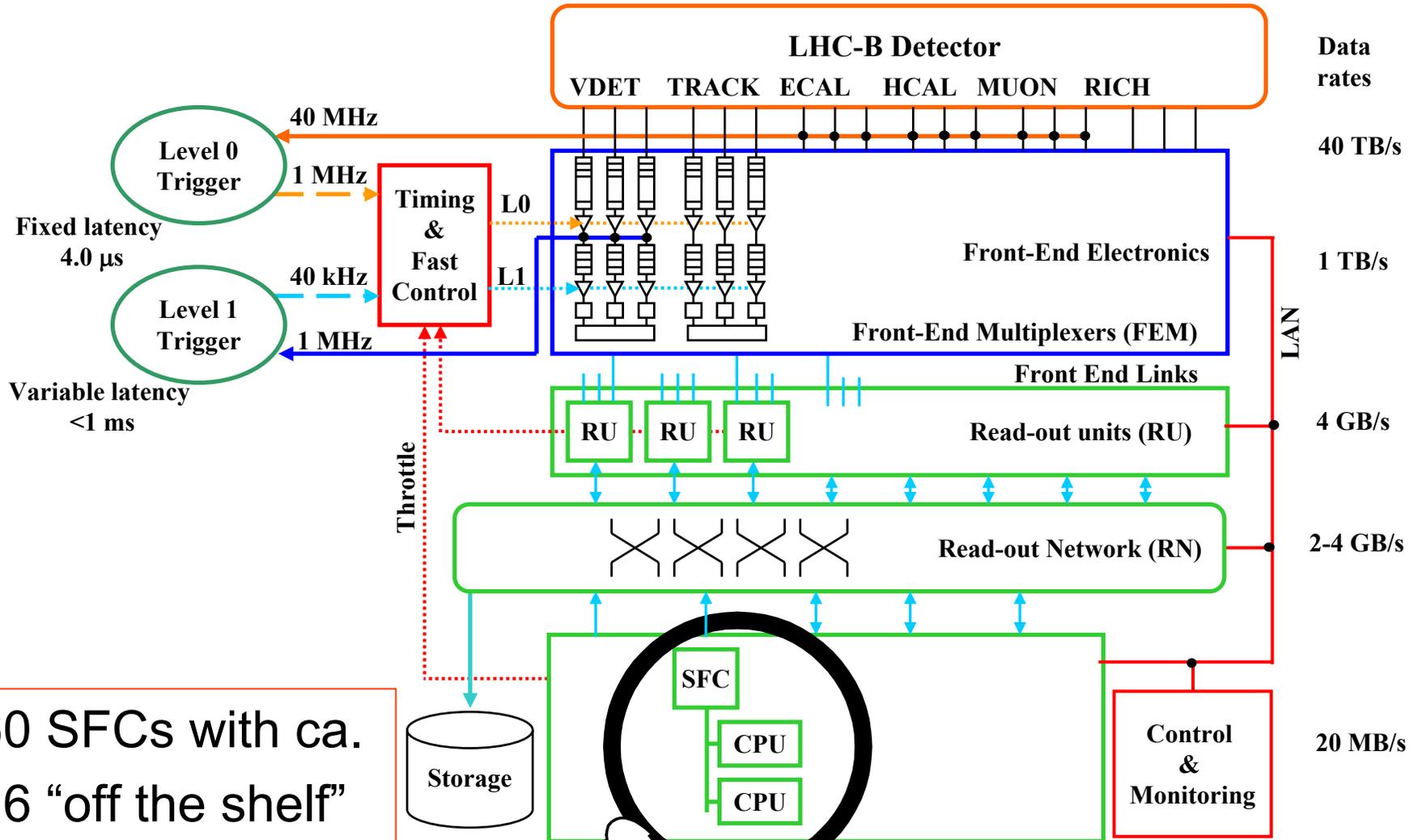
Reduce number of packets:
*Jumbo Frames*

Reduce context switches:
*Interrupt coalescence*

Reduce context switches:
*Checksum Offloading*
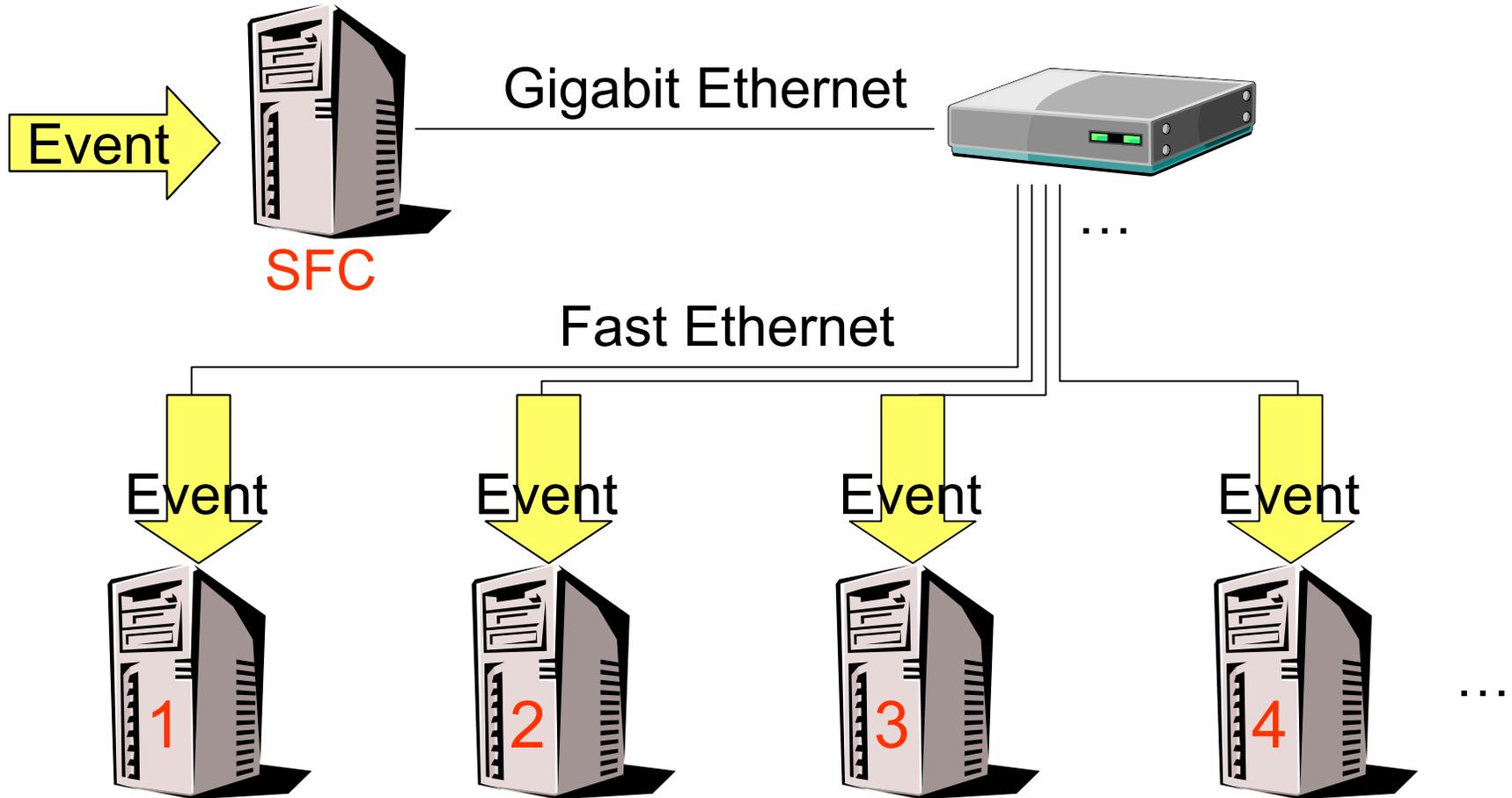
Or buy a faster PC with a better PCI bus… ☺

17

# Load Balancing

# Load Balancing: Where?



**LHC-B Detector**
VDET   TRACK   ECAL   HCAL   MUON   RICH

**Data rates**

40 MHz

**Level 0 Trigger**

1 MHz

**Timing & Fast Control**

L0

L1

40 kHz

**Level 1 Trigger**

1 MHz

Fixed latency 4.0 μs

Variable latency <1 ms

Throttle

40 TB/s

**Front-End Electronics**

1 TB/s

**Front-End Multiplexers (FEM)**

**Front End Links**

LAN

**RU   RU   RU**   **Read-out units (RU)**

4 GB/s

**Read-out Network (RN)**

2-4 GB/s

SFC

CPU

CPU

**Storage**

**Control & Monitoring**

20 MB/s

60 SFCs with ca. 16 "off the shelf" PCs each

19

# Load Balancing with round-robin

Event →

SFC

Gigabit Ethernet

…

Fast Ethernet

Event

Event

Event

Event

1

2

3

4

…

Problem: The SFC doesn't know if the node it wants to send the event to is ready to process it yet.

# Load Balancing with control-tokens



**Event**

Gigabit Ethernet

SFC

Fast Ethernet

Token
Event

Token
Event

Token
Event

1    2    3    4

...

With control tokens, nodes who are ready send a token, and every event is forwarded to the sender of a token.

# LHC Comp. Grid Testbed Structure

100 cpu servers on GE, 300 on FE, 100 disk servers on GE (~50TB), 10 tape server on GE



**64 disk server**

**1 GB lines**

**200 FE cpu server**

**Backbone Routers**

**3 GB lines**

**3 GB lines**

**8 GB lines**

**10 tape server**

**100 GE cpu server**

SFC

**36 disk server**

**100 FE cpu server**