



# ATLAS TDAQ/DCS DataCollection

---

## An ATLAS TDAQ Candidate architecture

Document ID: DC-049  
Document Version: 3.02  
Document Date: 31st July 2002  
Document Status: Final

---

### Abstract

In this note we propose an architecture for the full scale ATLAS TDAQ system. The architecture assumes an individual access to simple ROBs, two central switches and a number of concentrating switches grouping traffic from LVL2 and EF processing nodes. Based on numbers from the ATLAS TDAQ Paper Model we propose grouping of ROBs and processing nodes not exceeding 70% of the Gigabit Ethernet bandwidth. The key components of the architecture (Gigabit Ethernet NICs, links and switches) are the commodity equipment and have been already available on the market.

### Institutes and Authors:

CERN: Robert W. Dobinson [Bob.Dobinson@cern.ch]  
INP Cracow: Krzysztof Korcyl [K.Korcyl@cern.ch]  
BNL: Micheal LeVine [M.Levine@cern.ch]

**Table 1** Document Change Record

<b>Title:</b> An ATLAS TDAQ Candidate architecture.			
<b>ID:</b> ATLAS-TDAQ-2002-049			
<b>Version</b>	<b>Issue</b>	<b>Date</b>	<b>Comment</b>
Original	1	16 <sup>th</sup> Dec 2001	initial version
1.01	1	21 <sup>st</sup> Jan 2002	Add L2SVs, SFIs, SFOs
1.02	1	1 <sup>st</sup> Feb 2002	Public release
1.02	2	4 <sup>th</sup> Feb 2002	Corrections to Table 2, conclusions
1.03	1	18 <sup>th</sup> Feb 2002	More numerical corrections, Tables 1-5
2.00	1	5 <sup>th</sup> April 2002	Corrections of ID rates along new numbers from the Paper Model.  Modifications to grouping factor of the LVL2 farms  Introduction of grouping layer of switches for the EB
3.00	1	21 <sup>st</sup> May 2002	New paper Model numbers from 10 <sup>th</sup> April
3.01	1	6 <sup>th</sup> June 2002	Informal draft note converted into ATLAS internal TDAQ note
3.02	1	31 <sup>st</sup> July 2002	Final version of the note. Changed original naming "strawman" to the Candidate architecture

## 1 Introduction

A computer network is the central part of the Data Collection system providing communication and data transfer between various components of the higher levels of the ATLAS TDAQ system. In this note we propose an architecture for the full scale ATLAS TDAQ system.

### 1.1 Purpose of the document.

We aim to present a conceptually simple, nevertheless complete, approach making maximum use of the commodity equipment. However, our proposal should be taken as representative rather than definitive. We expect, that in process of verifying the basic concepts in the bat 513 test bed, variations and optimisations will take place.

## 1.2 Glossary, acronyms and abbreviations

### 1.2.1 Glossary

#### **Read-Out Buffer [ROB]**

Standard module which receives data from the ROD, stores them and makes them available to the LVL2 trigger processors and, for LVL2-selected events, to the EF trigger processors.

#### **Data Collection [DC]**

It is a sub-system of the ATLAS TDAQ responsible for movement of event data from the ROB (or ROS) to the ATLAS High Level Triggers.

#### **Data Flow Manager [DFM]**

Part of the Data Collection sub-system, it provides the synchronization necessary to move event data from the ROB (or ROS) to the Event Filter. It distributes load between a number of SFIs

#### **Event Builder [EB]**

Part of Data Collection sub-system, it merges all the fragments belonging to a unique event number into a full event at a single destination (SFI)

#### **Event Filter [EF]**

It provides the third level of event rate reduction in ATLAS. It is based on the full event read-out and detailed event analysis

#### **Event Filter farm**

The farm of processors in which the Event Filter runs. It is managed by the SFI, which performs the EB and distributes full events among processors from the farm.

#### **Level-2 Processing Unit [LVL2 PU]**

An element of a Level-2 sub-system incorporating the calculations from which the Level-2 trigger decision is derived.

#### **Level-2 Processing Unit cluster**

A collection of LVL2 PUs attached to the same concentrating switch.

#### **Sub Farm Interface [SFI]**

Part of the Data Collection sub-system. The location where full events are built by the Event Builder. Each SFI controls EF farm.

### 1.2.2 Acronyms and Abbreviations

<b>DFM</b>	Data Flow Manager
<b>ROB</b>	Read-Out Buffer
<b>ROD</b>	Read-Out Driver
<b>LVL2 PU</b>	Level-2 Processing Unit
<b>EB</b>	Event Builder
<b>SFI</b>	Sub Farm Interface
<b>ROS</b>	Read-Out System

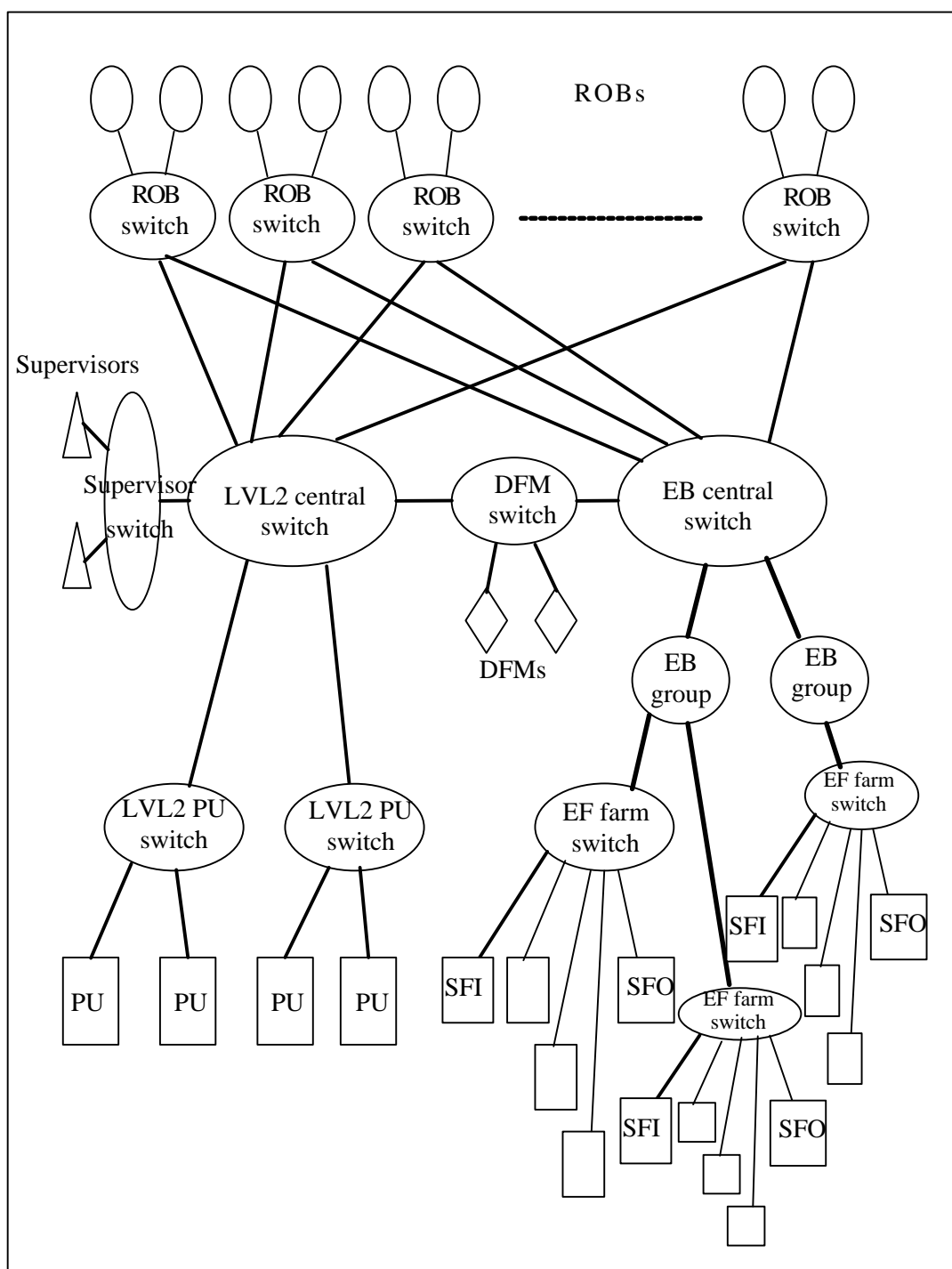
**EF** Event Filter  
**DFM** Data Flow Manager  
**VLAN** Virtual Local Area Network  
**PCI** Peripheral Card Interface

## 1.3 References

- 1 M. Abolins, J.C. Vermeulen: ‘Paper Model results for the ATLAS Trigger/DAQ system’; Draft 0.9; 2<sup>nd</sup> July 2001
- 2 M. Abolins, J.C. Vermeulen: ‘Overview of parameters used for 2001 paper model and of modifications in 2002 model’; 19<sup>th</sup> April 2002.

## 2 The TDAQ Candidate architecture

In Figure 1 we present the key elements of the architecture. In the following paragraphs we map paper model data rates and throughputs onto the model.



**Figure 1. Proposed architecture for the ATLAS TDAQ system. Bold lines represent GE links, remainder are FE links.**

## 2.1 Basic concepts

### 2.1.1 Simple ROBs with single network connection to the ROB concentrating switches

The architecture assumes that each ROB receiving data from a ROD (or a number of RODs) has a direct network connection to the central system. Each ROB can be accessed individually either by the LVL2 PU collecting ROIs or an SFI performing event building. The ROB will receive requests for data originated at the LVL2 PU or SFI and will respond with data fragment corresponding to the requested event number. The ROBs are grouped on ROB concentrating switches to make better utilization of the gigabit ports on the central switches. The ROB grouping factor depends on the load produced by a single ROB (request rate and size of replies) and varies with detector type and the accelerator luminosity. To run the network smoothly we assume that the aggregate load produced by all ROBs connected to a single concentrating switch should not exceed 70% of the throughput of the uplinks connecting concentrating switches with the central ones. The organization of ROBs is similar to the concept of the switch-based ROS, however the ROS controller, a node which relays requests and data traffic between ROBs connected to the ROB switch and the central network, is not assumed to be present in this model.

The connection of ROB to the ROB concentrating switch is done with FastEthernet as the amount of data produced for the LVL2 and the EB does not exceed 40% of nominal throughput (the highest load is generated by the EMCal ROBs for high luminosity: 4.66 MB/s corresponding to 37% of Fast Ethernet throughput – see Table 3 in the next paragraph). The latency degradation because of not going for faster links (Gigabit Ethernet) does not seem to be critical, as the main contribution to the transfer latency will be waiting time in the concentrating switch output queue. However, as the prices of Gigabit Ethernet will drop to the FastEthernet level, swapping Fast into Gigabit Ethernet can be made without any significant impact on performance of the rest of the system.

### 2.1.2 Two (or more) central switches

The central part of the architecture is based on two central switches - one for the LVL2 data flow and the other for the EB data flow. The concept of using more than one central switch helps to cut down the size of the required central switch. The switch market observation proves that the smaller the switch the cheaper is the cost per port. Having two distinct switches for two main data flows also helps to separate the two data flows at an early stage. The concept of having two network interfaces for a ROB/ROS has been dropped as the aggregated amount of traffic it produces for the two subsystems is smaller than the bandwidth provided by the Fast Ethernet standard, hence having two distinct network interfaces seems to be an unnecessary complication.

The two central switches make the whole network more fault tolerant. In the setup with a single central switch and in case of a switch crash, the operation of the Data Collection system will stop. The two switches in the system guarantee that in case of a crash of a single central switch, the other switch can be used to continue operation of the trigger system (with decreased performance). The proposed architecture offers a possibility to rearrange resources between the LVL2 and the EB subsystems to continue operation of the two subsystems through a single switch.

The number of central switches corresponds to the number of uplinks at the ROB concentrating switches. With only two uplinks in the concentrating switches, there is no point to increase the number of central switches as the traffic to the additional switches will have to flow through these uplinks and existing switches.

Using the ROB concentrating switches with a bigger number of uplinks can encourage installation of more than two central switches. However, installing an additional switch in the central part of the architecture would not cut the central switch size by two as each installed switch would need the same number of ports to connect to the ROB concentrating switches. The possible gain in number of ports would be in connections to the LVL2 PU farms or the EF farms. However, these connections consume only half of the switch's ports at the most. To gain in the number of ports to the ROB concentrating switches one would need concentrating switches with bigger Fast Ethernet valences. In the proposed architecture we plan to use concentrating switches with two Gigabit uplinks and 48 Fast Ethernet ports. In almost all cases the amount of traffic produced by 48 ROB's does not exhaust Gigabit uplink capacity (exceptions being SCT and TRT for intermediate and high luminosity – see Table 2 and 3), therefore adding another Gigabit uplink is not necessary. Switches with more than two Gigabit uplinks and larger number of Fast Ethernet ports will be more expensive.

The proposed architecture provides means to handle more than one central switch. All the issues presented here are valid independently of the number of central switches. At the time of writing this document two central switches are the best option. However, as the Ethernet technology evolves, installation of larger numbers of central switches may become more attractive.

### **2.1.3 LVL2 PU clusters and EB grouping switches**

The idea of grouping ROB's in concentrating switches for better throughput exploitation of the central switches' ports has been mirrored onto the LVL2 PU clusters and the EF farms via the EB grouping switches. The same criterion of not exceeding 70% on Gigabit uplink nominal throughput has been used here.

### **2.1.4 DFM**

The DFM receives information on events processed by the LVL2 sub-system and controls operation of the EB sub-system. It has to connect to the LVL2 sub-system and to the EB-subsystem. In the proposed architecture, the DFM is connected to the DFM switch that connects the two central switches. The DFM switch allows the DFM to operate through a single network interface, whereas an alternative scenario with direct connection to the LVL2 sub-system and the EB

sub-system would require two interfaces. Having the DFM connected to the network via the switch allows to use more than one node running the DFM functionality without adding extra ports to the central switches (currently it is not known whether a single DFM node will not become a bottleneck).

### 2.1.5 EF farms

We propose to connect directly the EF farm switch to the EB grouping switch and further to the EB central switch (as an alternative the EF farm switch will connect via the SFI to the EB grouping switch and further to the EB central switch). In the proposed organisation the SFI will operate through a single network interface (in the alternative the SFI would require two interfaces and would act as a “router”). Placing the SFI “behind” the EF farm switch makes the EF farm structure resemble the LVL2 PU cluster structure and thus opens a possibility to use some part of the EF farms as the LVL2 PU clusters. This may be useful either in case of one central switch crash or in case some processing resources from the EF can be used to run more complicated and more time-consuming LVL2 PU filtering algorithms.

### 2.1.6 VLANs

The Ethernet standard does not allow creating loops in the network in order to prevent broadcast messages from circulating infinitely, consuming all deployed bandwidth. The two Gigabit uplinks in the ROB concentrating switches and the DFM switch create dangerous loops. The Spanning Tree algorithm, omnipresent in contemporary switches, would switch off the redundant links to break the loops what would deteriorate performance.

To avoid loops we propose to utilise VLANs in the network. Implementing VLANs will require that all network components (nodes and switches) have to recognise extended Ethernet (VLANs and priorities) and all nodes have to act as routers – before placing packets on the network, they have to decide on which network (which VLAN) they wish to put the packet. The VLANs can also be used to partition the network between sub-detectors wishing to operate independently from other sub-detectors or the rest of the system.

### 2.1.7 Traffic patterns

Special thought has been given to provide solutions minimising packet losses. Packets can be lost due to broken links, corruption, ageing in switches and the most common: buffer overflows in switches. Buffer overflows can be avoided/minimised by supervising traffic patterns flowing through the architecture. Whereas broadcast and multicast handling is switch vendor-specific and we have little or no control over it, operation with unicasts gives full control over traffic patterns flowing through the architecture.

The calculations presented in the following paragraphs are based on using unicasts only.

An attempt to base the message exchange mechanism on unicasts only has the biggest impact on performance of the EB sub-system. The operation of the EB is based on request-reply mechanism for assigning events to the SFIs. Any SFI willing to process an event sends a message to the DFM. The DFM replies with the event number and the SFI generates 1500 unicasts messages to all ROB's requesting data for selected event.



## 2.2 ROB:

We use Paper Model [1] and [2] to find average rates for detector ROB.

Our calculations and assumptions for different scenarios: low, intermediate and high luminosity are presented in the Tables 2, 3 and 4 respectively.

Detector	# of ROB	RoI avg. (max) rate [kHz]	ROB avg. (max) LVL2 rate [MB/s]	ROB EB rate [MB/s]	ROBs grouping <87.5 [MB/s]	NB Of T5C <sup>1</sup>
Muon P	192	0.74 (1.48)	0.65 (1.30)	1.76	4 * 48	4
Muon T	48	1.62 (2.24)	0.75 (1.03)	0.92	1 * 48	1
Em Cal	760	0.86 (2.26)	0.71 (1.88)	1.66	15 * 48 + 40	16
Had Cal	98	1.18 (1.91)	0.98 (1.59)	1.66	2 * 48 + 2	3
TRT	256	5.25 (5.27)	1.47 (1.48)	0.56	5 * 48 + 16	6
SCT	92	5.67 (5.83)	1.59 (1.63)	0.56	1 * 48 + 44	2
Pixels	120	5.78 (6.01)	1.16 (1.20)	0.46	2 * 48 + 24	3
Total ATLAS	1566					35

**Table 2. ROB throughput calculations leading to link load estimates and ROB grouping for low luminosity.**

Detector	# of ROB	RoI avg. (max) rate [kHz]	ROB avg. (max) LVL2 rate [MB/s]	ROB EB rate [MB/s]	ROBs grouping < 87.5 [MB/s]	NB Of T5C <sup>1</sup>
Muon P	192	0.62 (1.23)	0.55 (1.08)	1.76	4 * 48	4
Muon T	48	1.34 (1.86)	0.62 (0.86)	0.92	1 * 48	1
Em Cal	760	1.28 (3.35)	1.06 (2.78)	1.66	15 * 48 + 40	16
Had Cal	98	2.10 (3.36)	1.74 (2.79)	1.66	2 * 48 + 2	3
TRT	256	4.54 (4.56)	1.86 (1.87)	0.82	5 * 47 + 21	6
SCT	92	4.96 (5.13)	2.03 (2.10)	0.82	2 * 43 + 6	3
Pixels	120	5.07 (5.33)	1.42 (1.49)	0.56	2 * 48 + 24	3
Total ATLAS	1566					36

**Table 3. ROB throughput calculations leading to link load estimates and ROB grouping for intermediate luminosity**

Detector	# of ROB	RoI avg. (max) rate [kHz]	ROB avg. (max) LVL2 rate [MB/s]	ROB EB rate [MB/s]	ROBs grouping <87.5 [MB/s]	NB Of T5C <sup>1</sup>
Muon P	192	0.33 (0.72)	0.29 (0.63)	1.76	4 * 48	4
Muon T	48	0.70 (1.08)	0.32 (0.50)	0.92	1 * 48	1
Em Cal	760	1.35 (3.61)	1.12 (3.00)	1.66	15 * 48 + 40	16
Had Cal	98	1.18 (1.95)	0.98 (1.62)	1.66	2 * 48 + 2	3
TRT	256	0.02 (0.03)	0.03 (0.04)	2.56	7 * 34 + 18	8
SCT	92	0.74 (1.01)	0.95 (1.29)	2.56	2 * 34 + 24	3
Pixels	120	0.94 (1.36)	0.55 (0.79)	1.16	2 * 48 + 24	3
Total ATLAS	1566					38

**Table 4. ROB throughput calculations leading to link load estimates and ROB grouping for high luminosity**

Notes:

- ROI average and max rates have been calculated assuming 80 kHz LVL1 rate (rates from the 40 kHz LVL1 menu have been multiplied by 2, but the ID full scan which has been retained at 5 kHz)
- 2 kHz LVL2 accept rate (EB rate) has been assumed.
- ROB data fragments produced for the LVL2 requests (RoI replies) and for the EB (event fragment replies) have the same size for a given luminosity (we assume here that any data compression/formatting is performed at the ROD level).
- Data sizes produced by ROB include:
  - 32 Bytes overhead of the message identification control information (ROB identification, event ID, message ID etc)
  - 26 Bytes of Ethernet framing (source and dest addresses, control bytes, CRC etc)
  - 20 bytes of IP header (we assume IP RAW\_SOCKETS communication protocol)

5. We assume that ROBs will be grouped using concentrating switches. The concentrating switches will have two uplinks: one for the LVL2 traffic and another one for the EB traffic. The uplinks will be Gigabit Ethernet with nominal throughput of 125 MB/s. We group ROBs into a single concentrating switch if the aggregated bandwidth produced either for the LVL2 traffic or for the EB traffic does not exceed 70% (87.5 MB/s) of nominal throughput of the uplink. Grouping in the above tables is presented in the form:  $a * b + c$ , where 'a' represents number of switches, 'b' represents number of ROBs connected to a single switch and 'c' represents the number of ROBs that have to be connected to the last switch (usually the valence of the last switch is not fully utilised).
6. The maximum total rate per ROB (max LVL2 + EB: sum of entries in column 5 in braces and column 6) does not exceed 63 % of nominal FE throughput (the highest load is for the Em Cal for high luminosity:  $1.24 + 2.46 = 3.7$  MB/s). We therefore assume that ROBs will be connected to the concentrating switches via the FE links. Switches with large FE port valences and two Gigabit uplinks (BATM *T5 Compact* like switches<sup>1</sup>) should be cost effective compared to all-Gigabit switches.

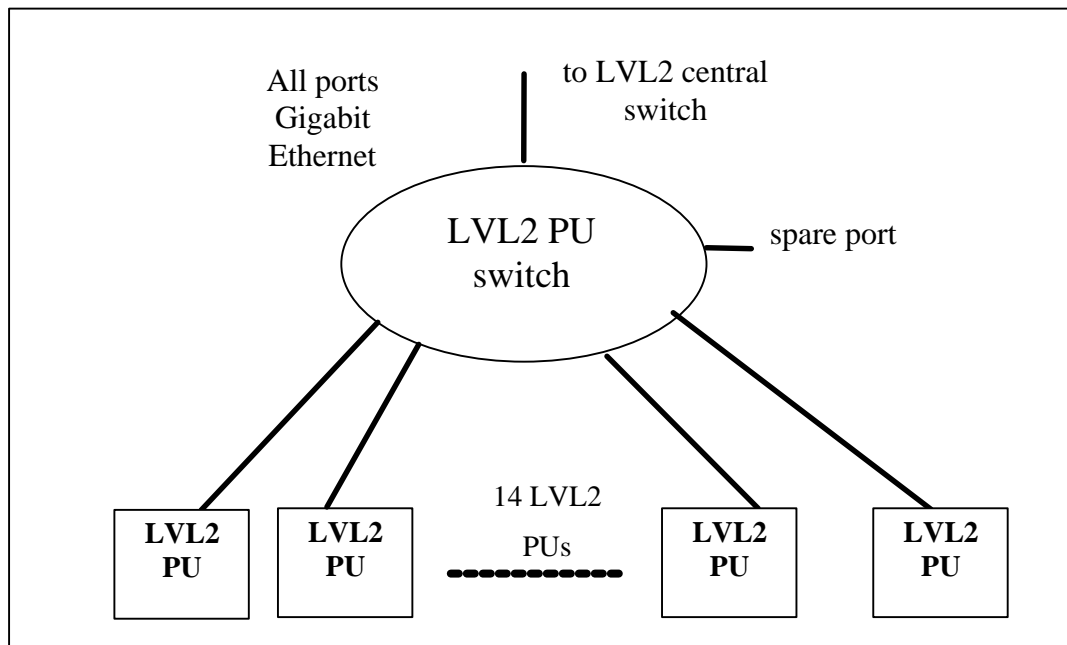
As the ROB load conditions are the heaviest for the high luminosity scenario, we will use grouping presented in Table 3. The required number of ports in the two central switches to connect ROB concentrating switches will be thus 38.

---

<sup>1</sup> T5C or BATM *T5 compact* switch has 2 Gigabit ports and 48 FastEthernet ports

## 2.3 LVL2 clusters

The current estimate for processing power required for the LVL2 tasks is determined by the power required for system operation in the intermediate luminosity mode. Table 44 in the Paper Model results note (intermediate luminosity, TRT scan after MU8 and processing time equal to mean of the processing times distribution) cites 750 4 GHz processors. We assume that 375 dual-processor boxes each with a single NIC card will be used. Taking into account 80 kHz LVL1 rate, this will result in 213 Hz event rate for a single box. The total ATLAS throughput for the LVL2 divided by 80 kHz LVL1 rate gives an average 24.3 kB data for a single event. With 320 Hz rate running each processing box this will require 5.18 MB/s throughput per box. We propose to group processing boxes using small all-Gigabit switches. As the throughput to each processing box (5.18 MB/s) is ~24 times lower than the nominal Gigabit Ethernet bandwidth we propose a grouping ratio 14:1. The throughput of 14 processing boxes remains below 60 % of the throughput of the Gigabit uplink connecting the cluster of processing units to the central switch. We propose to group 14 processing boxes on a single 16-port Gigabit switch with 1 uplink to the central switch (leaving one Gigabit port unconnected). The proposed cluster organisation is presented in Fig. 2. Grouping 375 boxes into groups of 14 will result in 27 level 2 processing clusters. Each cluster will connect to the central LVL2 switch via a single GE link, requiring 27 ports in the central LVL2 switch.



**Figure 2. Organisation of LVL2 Processing Unit farm**

## 2.4 EF farms

Organisation of the EF farm is presented in Fig. 3. We assume that 150 EF (SFIs) farms will be used. The LVL2 accept rate of 2 kHz spread across 150 farms results in 14 events per second treated by a single farm. Throughput required by a single farm where events have a size of 1.35 MB, results in 18.9 MB/s, which can easily be handled by a Gigabit Ethernet link. Also the throughput required at the SFI for collecting fragments from the ROB's and sending it out to workers ( $2 * 18.9$  MB/s) seems to be well below the limit of the PC's internal PCI bus; even the slowest PCI [32bit/33 MHz] will be able to handle that traffic. Installing 10 dual-processor PCs in a farm will allow event processing times up to 1.4 second (event rate: 0.7 Hz per processor).

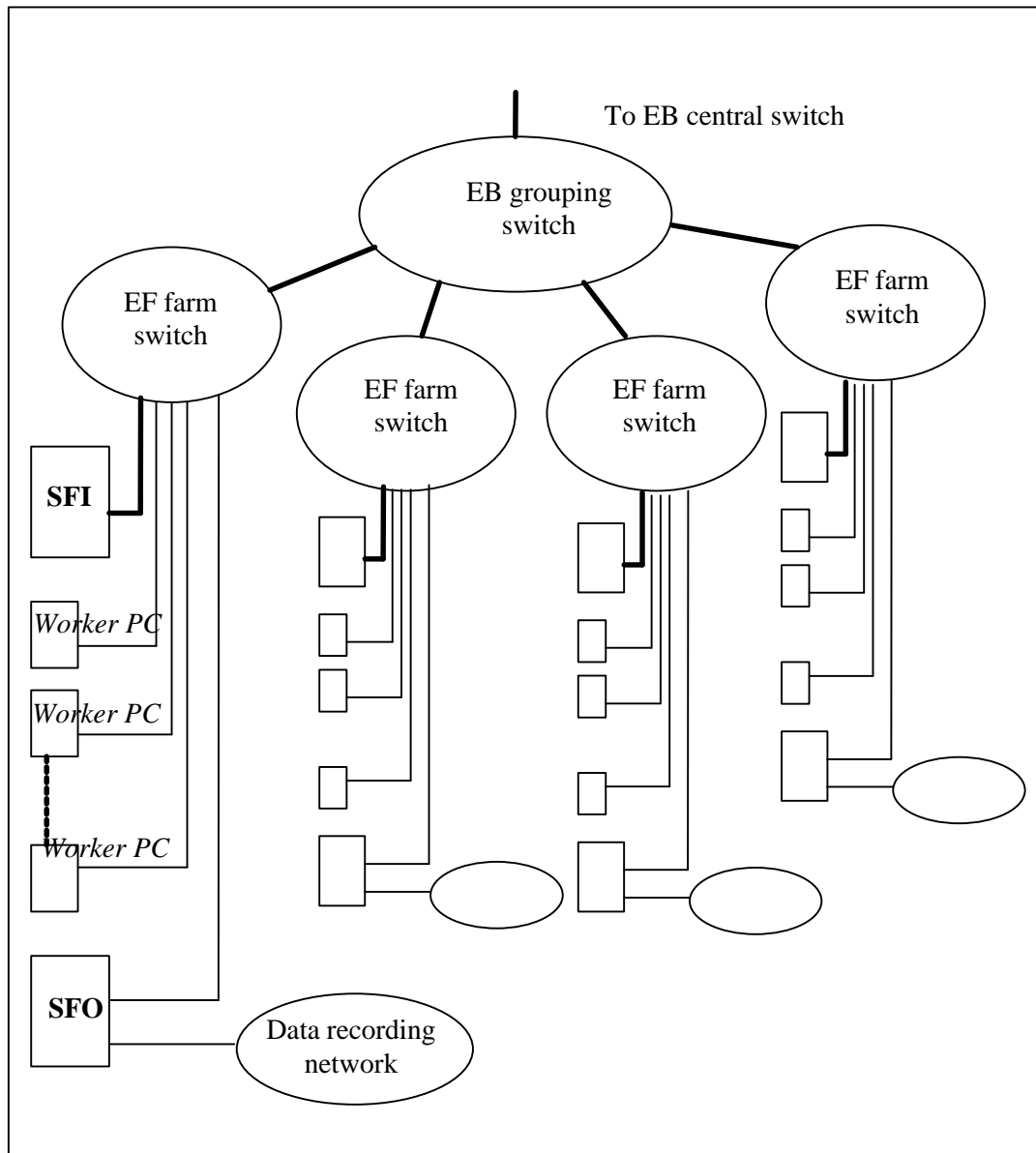
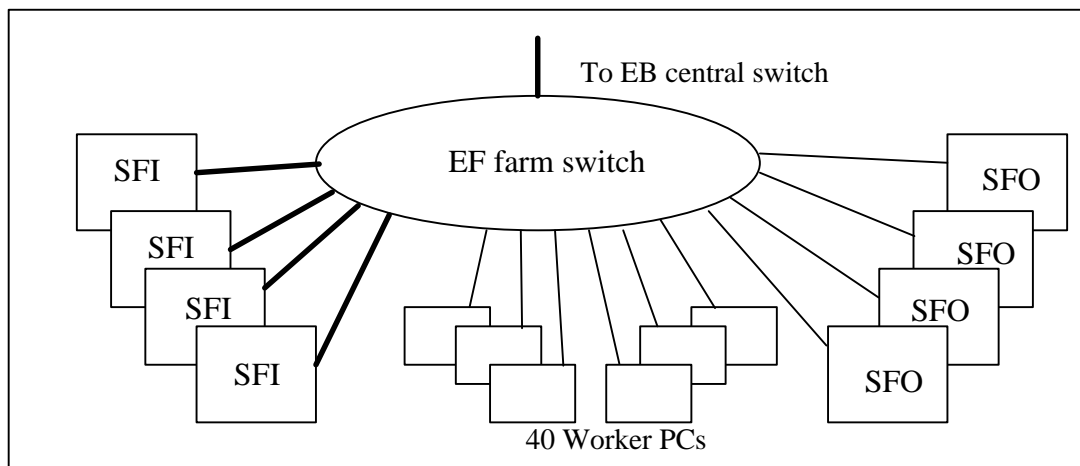


Figure 3. Detailed view of the EB part of the system with EB grouping switches

Connections to the EF farms significantly contributes to the size of the EB central switch; we would need 150 ports in the EB central switch to connect directly each farm's switch to the central EB switch using a Gigabit connection. Since the bandwidth of these ports would be lightly utilised [ $18.9/125 = 15\%$ ], we propose to use an additional layer of concentrating switches to group traffic from a number of EF farms. To keep the usage of Gigabit Ethernet link within 70% of its nominal bandwidth the grouping factor is 4. We plan to use small all-Gig (8 port) switches as the EB grouping switches. We also propose to connect the SFI to the farm's switch via a Gigabit link.

In such a scenario, traffic for a particular SFI would flow from the central EB switch to the EB grouping switch, then further down to the farm's switch and then to the SFI. Workers can be connected to the farm's switch via FE links, as required bandwidth for a single worker is only 1.89 MB/s (1.4 events/second). We propose to use switches like the BATM *T5compact* (2 Gigabit ports + 48 FE ports) as the EF farm switches.

An advantage to the EF farm organisation as depicted in Fig. 3 is the possibility to use it as a LVL2 processor cluster. Changing functionality of the farm would require loading LVL2 PU software into SFI and Worker PCs followed by a VLAN reconfiguration for the central switches. However, as the connection to the worker PCs is via Fast Ethernet, assigning events with high data transfer rates (TRT scan) to these workers should be avoided.



**Figure 4. Detailed view of the EF part of the system with EF farm switch only**

A possible evolution of the EB part of the system is presented in Figure 4. Assuming availability of switches with a greater number (5) of GE uplinks and big valences of Fast Ethernet ports, a number of EF farms could be grouped onto single EB farm switch.

## 2.5 LVL2 Supervisors

We propose to connect a number of LVL2 supervisors to the system via an all-Gigabit switch, of the type proposed for the LVL2 PU cluster (8 ports). Using the switch will make the system more flexible regarding changing the number of Supervisors. Connecting the LVL2 Supervisor switch to the LVL2 central switch will require 1 port, leaving 7 ports to connect LVL2 Supervisors.

## 2.6 DFM

We propose to connect a number of DFM nodes to the system via an all-Gigabit switch, of the type proposed for the LVL2 PU cluster (8 ports). Using the switch will make the system more flexible regarding changing the number of DFMs. Connecting the DFM switch to the central LVL2 switch requires 1 port in the central EB switch, leaving 6 ports to connect DFMs.

## 2.7 Central network architecture

In Table 3 we showed grouping of ROB's which results in number of Gigabit uplinks transferring data from ROB's to the central switches. The total number of gigabit ports in the central switches to receive ROB data (either the LVL2 data or the EB data) is 38.

The central LVL2 switch will have to have large enough valence to connect:

- 38 ports from ROB switches
- 1 port from the LVL2 Supervisor switch
- 1 port from the DFM switch
- 27 ports from LVL2 processing units switches

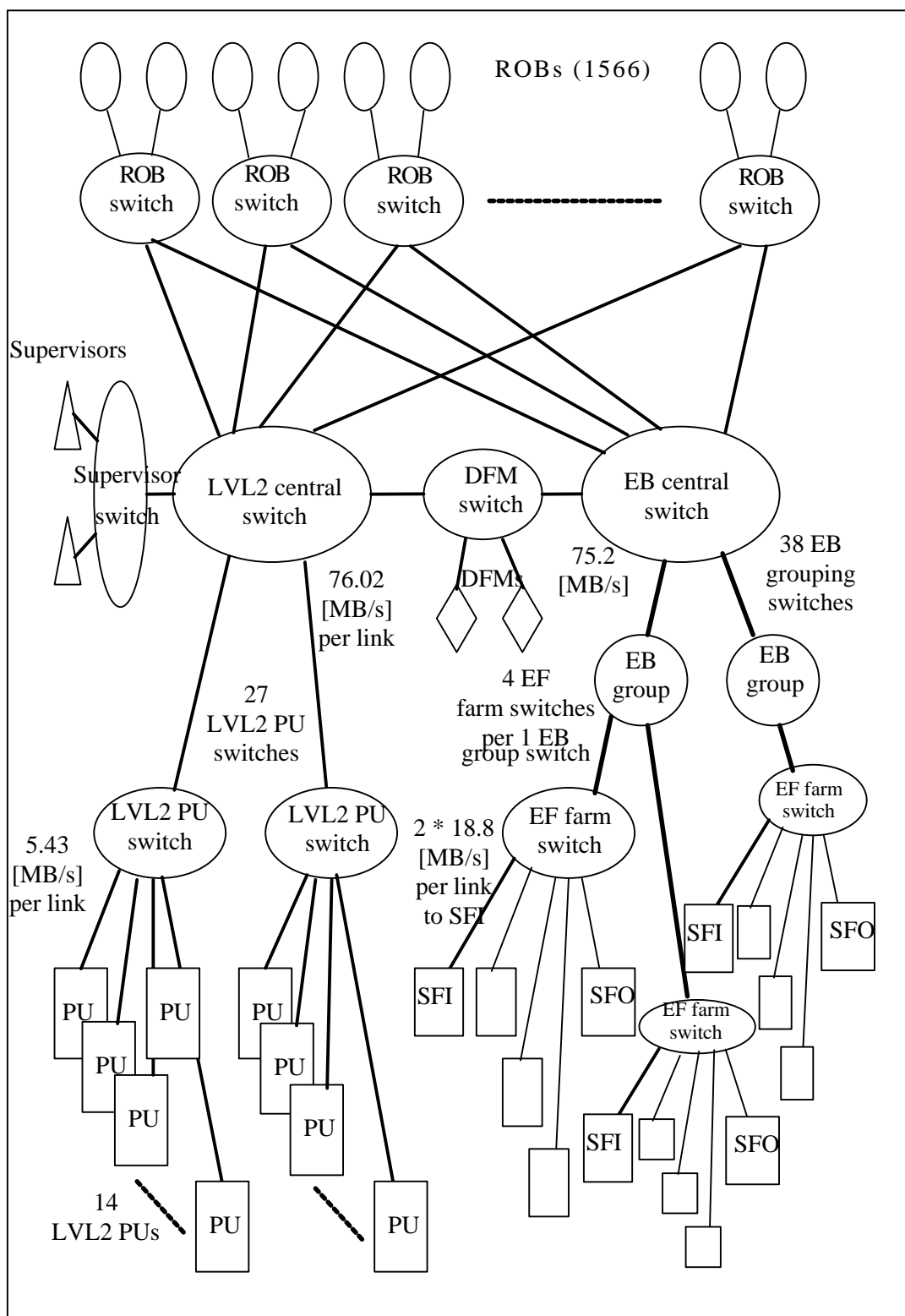
The total number of required ports is thus 67. Taking into account the modular structure of large switches (modules of 8 ports) we plan to use a 72 port switch.

The central EB switch will have to have large enough valence to connect:

- 38 ports from ROB switches
- 1 port from DFM switch
- 38 ports from EB grouping switches

The total number of required ports is thus 77. Taking into account the modular structure of large switches (modules of 8 ports) we plan to use a 80 port switch.

In Fig. 5 we present the same architecture as in Figure 1 but with quantitative information based on our calculations presented in this note.



**Figure 4. Quantitative view of the proposed Candidate architecture for the ATLAS TDAQ system. For the Gigabit uplinks loads see Table 6 below.**



### 3 Comments

We have presented the Candidate architecture using network based ROBs grouped using switches. We have assumed the main data flow is always from the ROB directly to the processing node. The presence of a local (ROS) controller is not excluded. Its functionality in our view would be limited to assuring secure broadcast and multicast, although extra functionality is not excluded (for example partial event building). However, any data flow through local controller risks introducing a bottleneck. On the other hand, allowing direct ROB addressing from the LVL2 PUs creates high I/O overhead in terms of message rates. Based on the numbers from low luminosity scenario (Table 1) we aggregated rates in Table 5.

Detector	Number of ROBs * avg. RoI request rate per ROB [kHz]
Muon Prec	53,76
Muon Trig	28,80
Em Cal	653,60
Had Cal	115,64
TRT	1344,00
SCT	520,72
Pixels	694,80
TOTAL	3411,32

**Table 5. Average ROI request rates per detector when addressing each ROB individually for low luminosity scenario. For other luminosities (Table 2 and 3) we have 3034 kHz and 1424 kHz respectively as the total ATLAS ROI request rates.**

Assuming 375 LVL2 PU boxes, the ATLAS total ROI requests rate corresponds to 13.65 kHz/box. Taking into account responses coming back from individual ROBs, we arrive at 27.30 kHz message rate per box (request + response). Introducing Local Controllers may make the job of LVL2 PU easier, but the Local Controller may then a bottleneck, as it then has to generate individual requests for each ROB it controls, and later collect responses.

We may have a potential problem of handling high message rates at hot spots in the system, Interrupt coalescence at the LVL2 PU will certainly be very helpful. But other solutions including RoI collection in the NIC, as implemented by David Botterill, may also be applicable.

Detector	Load on uplink to LVL2 [MB/s]	Load on uplink to EB [MB/s]
Muon P	12.96	79.68
Muon T	13.92	39.36
Em Cal	50.88	75.36
Had Cal	44.64	75.36
TRT	57.05	86.10
SCT	56.70	86.10
Pixels	80.16	50.88

**Table 6. Maximum load on the Gigabit Ethernet uplinks for the ROB concentrating switches resulting from the proposed grouping of ROBs for the low luminosity conditions (supports Figure 5 above).**

## 4 Verification in the testbed in bld 513 using the data collection SW

The prime aim for tests in the testbed in building 513 will be to validate the basic concept of the Candidate architecture. We plan to demonstrate the architecture can provide the required connectivity between all components. Having done that, we plan to check whether the individual components can achieve the required bandwidth and message rates as predicted for the whole system. As stated at the start of this note we would envisage variants and optimisations of the basic architecture HW and SW.

This document has been prepared using the Short Note Template provided and approved by the ATLAS TDAQ and DCS Connect Forum. For more information, go to <http://atlas-connect-forum.web.cern.ch/Atlas-connect-forum/>.

This template is based on the SDLT Single File Template that has been prepared by the IPT Group (Information, Process and Technology), IT Division, CERN (The European Laboratory for Particle Physics) and then converted to MS Word. For more information, go to <http://framemaker.cern.ch/>.