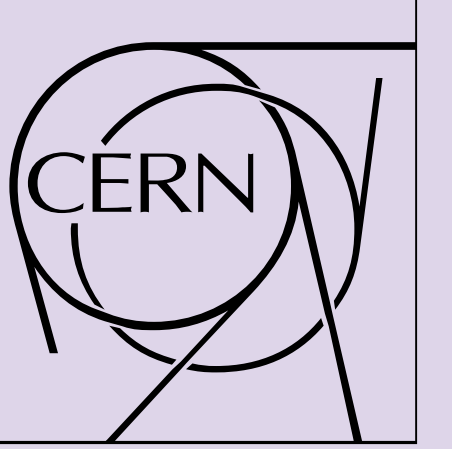


# Evaluating a File-based Event Builder to enhance the Data Acquisition in the CMS Experiment



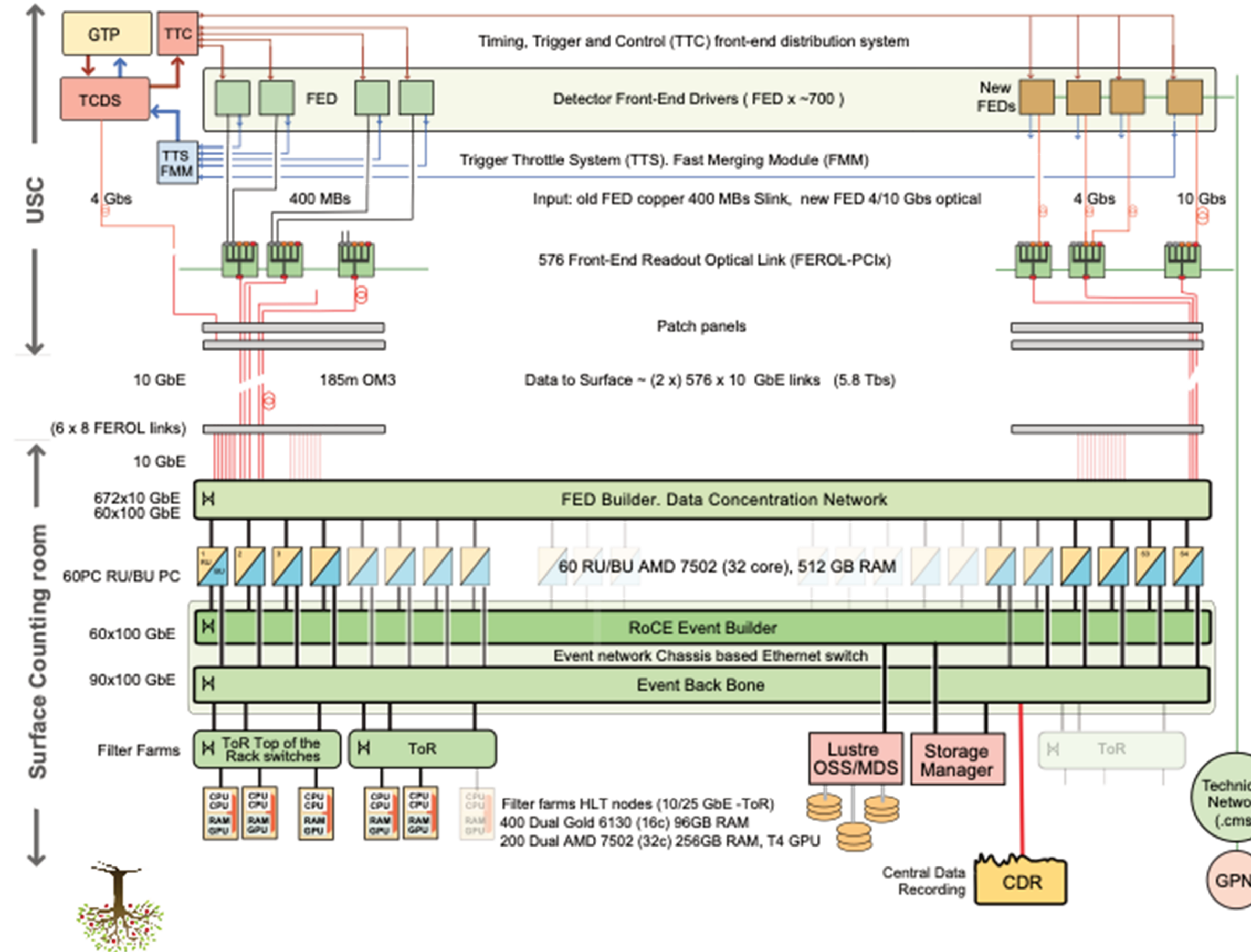
J. Alawieh<sup>1</sup>, K. Arutjunjan<sup>1</sup>, M. Bacharov Durasov<sup>1</sup>, U. Behrens<sup>2</sup>, A. Bocci<sup>1</sup>, J. Branson<sup>3</sup>, P. M. Brummer<sup>1</sup>, J. A. Bugajski<sup>1</sup>, E. Cano<sup>1</sup>, S. Cittolin<sup>3</sup>, A. Corominas I Mariscot<sup>1</sup>, G. L. Darlea<sup>4</sup>, C. Deldicque<sup>1</sup>, M. Dobson<sup>1</sup>, A. Dvorak<sup>1</sup>, C. Emmanouil<sup>1</sup>, A. Gaille<sup>1</sup>, D. Gigi<sup>1</sup>, F. Glege<sup>1</sup>, G. Gomez-Ceballos<sup>4</sup>, P. Gorniak<sup>1</sup>, J. Hegeman<sup>1</sup>, G. Izquierdo Moreno<sup>1</sup>, T. O. James<sup>1</sup>, T. Jayakumar<sup>1</sup>, W. Karimeh<sup>1</sup>, R. Krawczyk<sup>2</sup>, W. Li<sup>2</sup>, K. Long<sup>4</sup>, F. Meijers<sup>1</sup>, E. Meschi<sup>1</sup>, M. Migliorini<sup>1,7</sup>, S. Morović<sup>3</sup>, B. J. Odetayo<sup>1,5</sup>, L. Orsini<sup>1</sup>, C. Paus<sup>4</sup>, A. Petrucci<sup>3</sup>, M. Pieri<sup>3</sup>, D. S. Rabaday<sup>1</sup>, A. Racz<sup>1</sup>, T. Rizopoulos<sup>1</sup>, H. Sakulin<sup>1</sup>, C. Schwick<sup>1</sup>, D. Šimelevičius<sup>1,6</sup>, P. Tzanis<sup>1</sup>, C. Vazquez Velez<sup>1</sup>, P. Žejdl<sup>1</sup>

<sup>1</sup> CERN, Geneva, Switzerland; <sup>2</sup> Rice University, Houston, Texas, USA; <sup>3</sup> UCSD, San Diego, California, USA; <sup>4</sup> MIT, Cambridge, Massachusetts, USA; <sup>5</sup> University of Benin, Benin City, Nigeria; <sup>6</sup> Vilnius University, Vilnius, Lithuania; <sup>7</sup> University of Padova, Padova, Italy



## CMS Data Acquisition System

### Architecture for the LHC RUN 3 (2022-2026)



## Alternative File-Based Event Builder

In the LHC Phase-2, filter farm processes need to assemble events from the given orbit for analysis. Why not build the event or orbit directly within the filter farm? By doing so, the step of moving data between super-fragments and the event or orbit building process can be eliminated.

**Benefits** of building events or orbits in the filter farm process:

- No need for the Event Backbone network** resulting in **cost savings** on network line cards and NICs;
- Simplifying the event builder code** reduces its task to only building super-fragments;
- Building events or orbits** entirely in **RAM** disk on filter farm PCs **reduces** the overall number of **memory copies** in the data flow, which would not be the case if the BU applications were moved to the filter farm PCs.

### File-based Event Builder

#### Super-fragment Builder (SFB)

The SFB constructs multiple super-fragments based on the grouping of FEDs and stores them in local RAM disks.

#### Builder Filter Farm File Based (B3F)

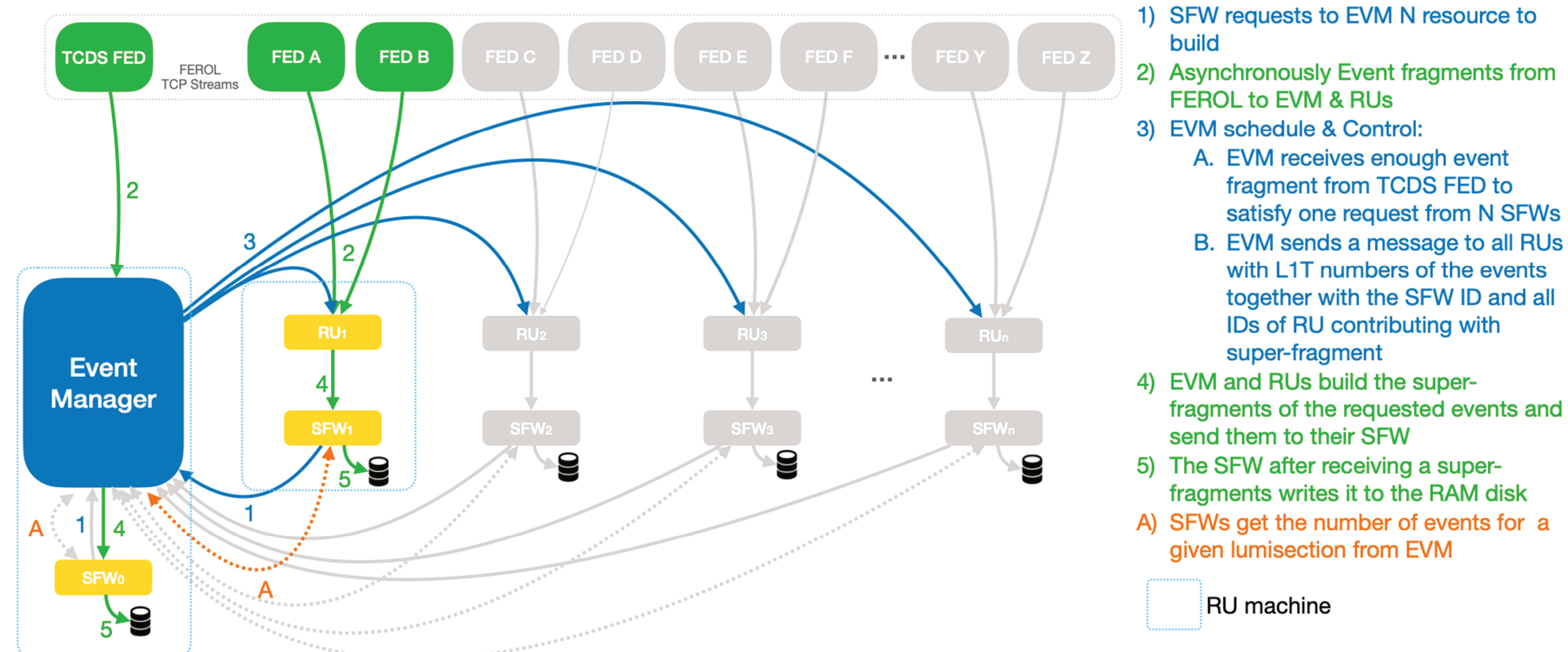
The BF3 accesses super-fragments from all RU machines via the NFS over Ethernet and builds complete events within the HLT farm.

### Super-fragment Builder

A **super-fragment** consists of the data read by **one or more Front-Ends** and corresponding to the same **L1 accept or orbit**, and the SFB constructs multiple **super-fragments** corresponding to the number of Read-Unit (RU) machines in the DAQ system, storing them in local RAM disks. In the LHC Phase-2, the DAQ **custom-designed** hardware will **identify potential issues** during data taking. In contrast, in the LHC Run 3, this task is handled by the EVB, which needs to be integrated with the SFB.

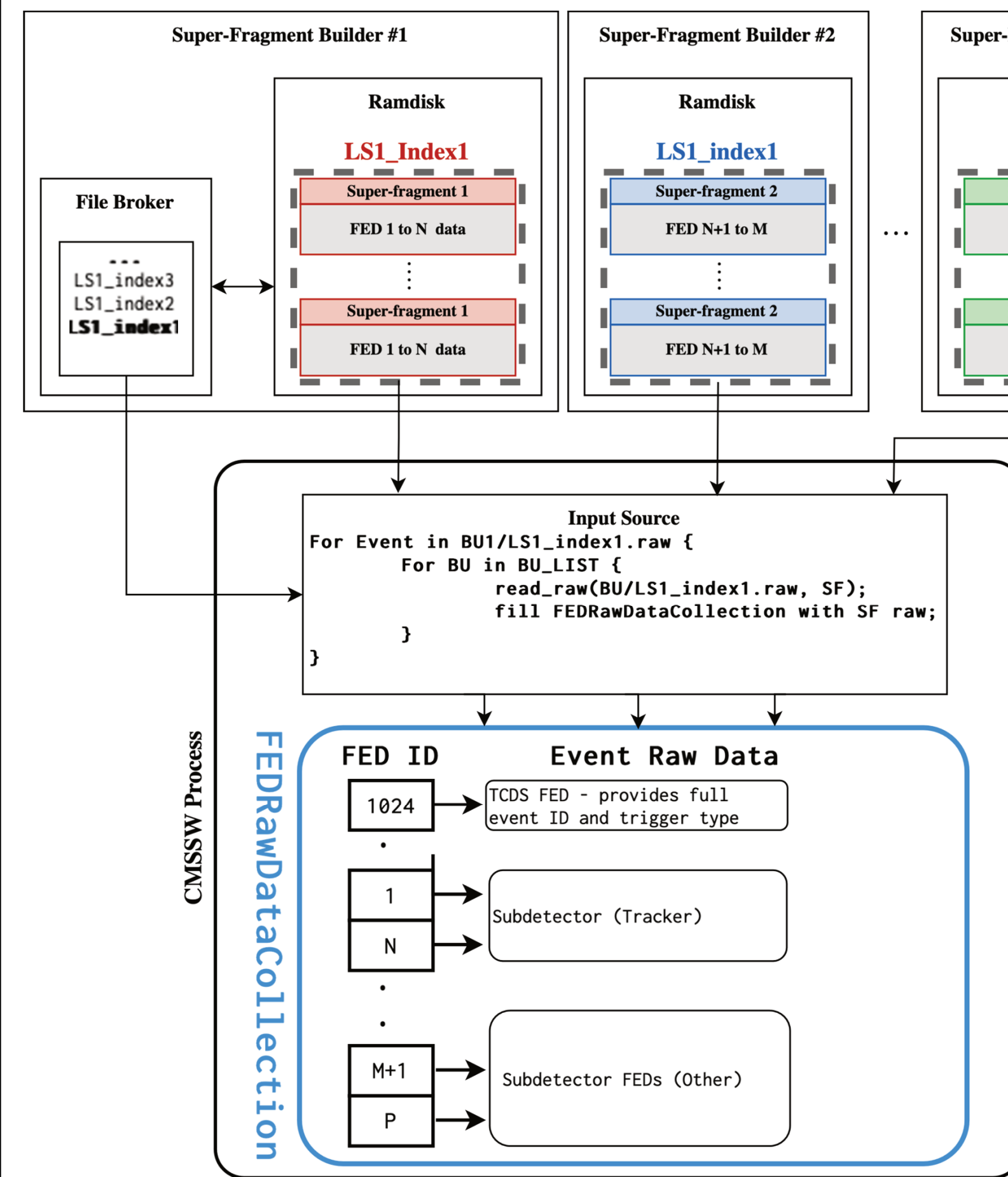
### Super-fragment Builder Protocol

Preliminary **performance tests** of the SFB using DAQ 3 hardware were conducted at the beginning of the year with the **DAQ 3 Emulator**. By **dropping data** in the **RAM disk** (without HLT), the SFB was able to reach the nominal **115 kHz L1** rate for the LHC Run 3.



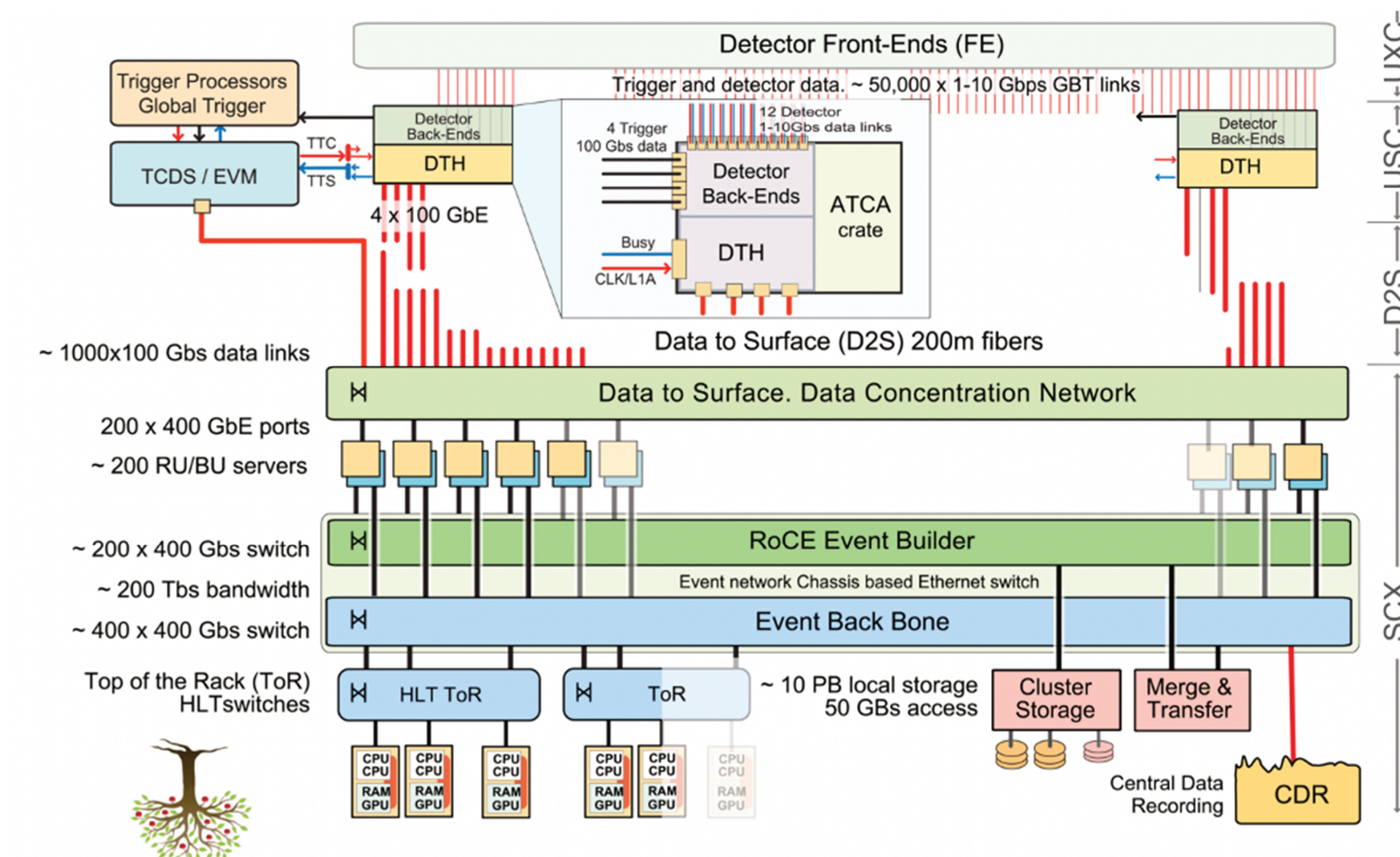
- SFW requests to EVM N resource to build
- Asynchronously Event fragments from FEROL to EVM & RUs
- EVM schedule & Control:
  - EVM receives enough event fragment from TCDS FED to satisfy one request from N SFWs
  - EVM sends a message to all RUs with L1T numbers of the events together with the SFW ID and all IDs of RU contributing with super-fragment
- EVM and RUs build the super-fragments of the requested events and send them to their SFW
- The SFW after receiving a super-fragment writes it to the RAM disk
  - SFWs get the number of events for a given lumisection from EVM

### Builder Filter Farm File Based (B3F)



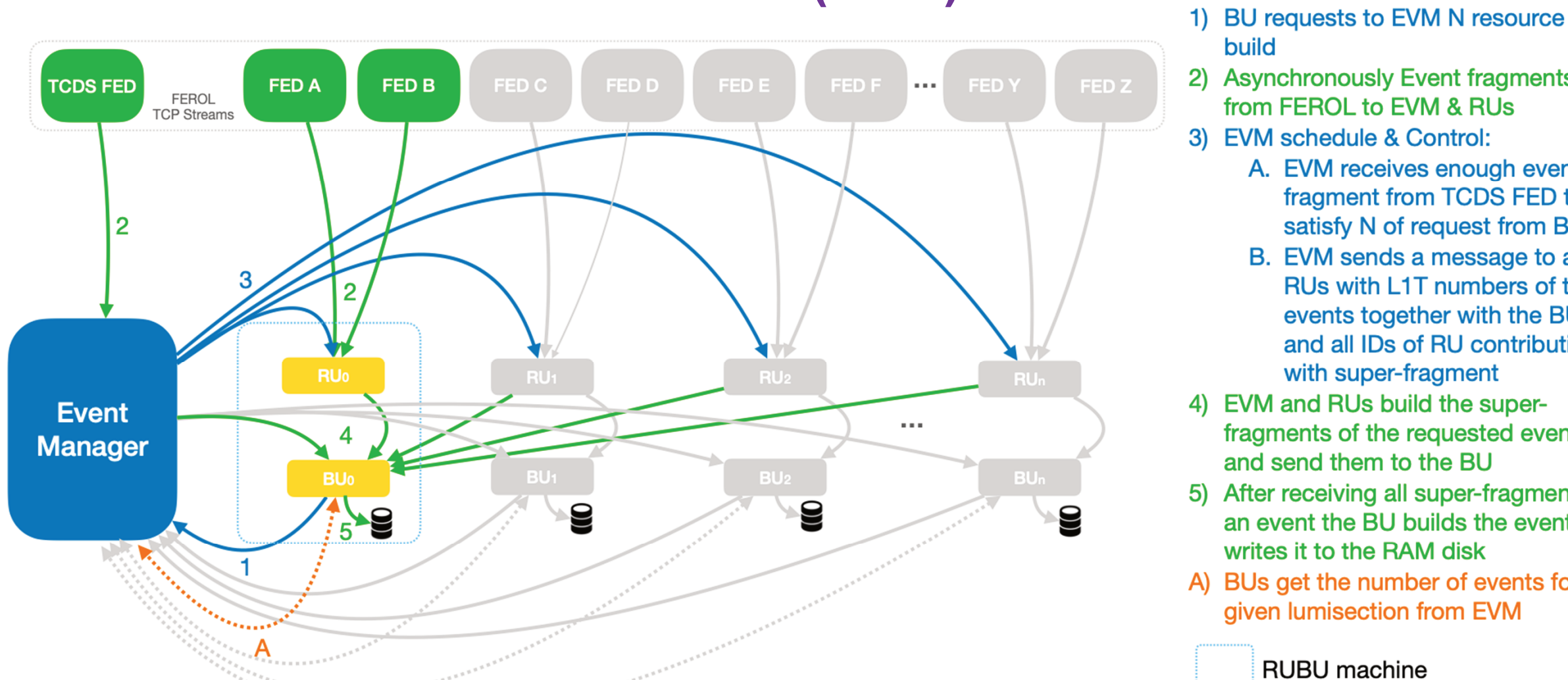
- Super-Fragment Builder:**
  - In this node, the SFB runs and creates the stream files in the RAM disk.
- File Broker:**
  - Manage the assignment of Super-Fragment stream files to Filter Farm processes
- CMSSW DAQ modules:**
  - INPUT:**
    - Read super-fragments and meta data from all RU RAM disks
    - Build events
    - Write booking metadata for HLTD
  - OUTPUT:**
    - Write output streams and meta data to local RAM disk
- High Level Trigger Daemon (HLTD):**
  - Merge output streams to Storage and Transfer System
  - Monitoring, scheduling CMSSW processes

### Architecture for the LHC Phase-2 (starting in 2030)



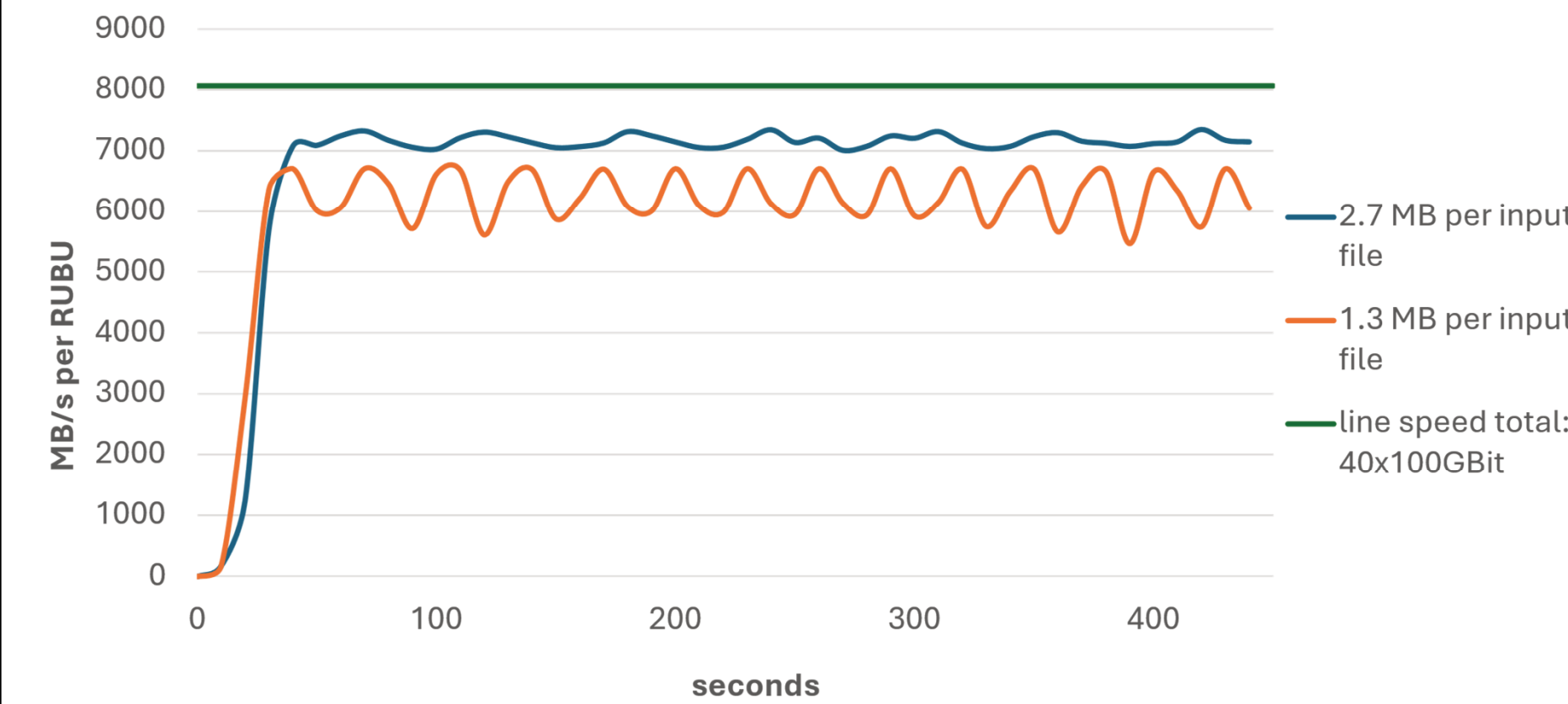
DAQ Building Information	LHC RUN 3	LHC Phase-2
L1 accept rate (maximum)	115 kHz	750 kHz
Event Size	2 MB	8.4 MB
Event network traffic	1.84 Tb/s	51 Tb/s
D2S modules	FEROL	DTH-400/DAQ-800
D2S module Ethernet ports	10 GbE	100 GbE
Number of D2S source ports	650	850
D2S network, EVB ports	100 GbE	400 GbE
EVB Network	100 GbE RoCE	400 GbE RoCE
Number of EVB nodes	~ 60	~200
Aggregation	Super-fragment	LHC Orbit (~67 Super-fragments)

### Event Builder (EVB) Protocol



- BU requests to EVM N resource to build
- Asynchronously Event fragments from FEROL to EVM & RUs
- EVM schedule & Control:
  - EVM receives enough event fragment from TCDS FED to satisfy N of request from BU
  - EVM sends a message to all RUs with L1T numbers of the events together with the BU ID and all IDs of RU contributing with super-fragment
- EVM and RUs build the super-fragments of the requested events and send them to the BU
- After receiving all super-fragments for an event the BU builds the event and writes it to the RAM disk
  - BUs get the number of events for a given lumisection from EVM

All to all HLT input over NFS test (62 RU x 196 FU) with Run 3 DAQ



This plot shows the B3F performance using DAQ 3 RUBU nodes of **all-to-all** HLT NFS using TCP/IP operation with a script running in each node to **generate RAM disk files** using real event sample:

- Total throughput Limited by top of the rack switch uplinks 40 x 100 GbE (8 links per rack);
- Better result with larger files (nearly 90% of the line speed);
- Next step to explore NFS over RDMA.

### Future Work

In the **LHC Run 3**, the main goal is to develop and commission the **File-based Event Builder** during periods when data is not being taken. The target for the **File-based Event Builder** is to go **operational** for data taking during the **final years** of the LHC Run 3.

Looking ahead to the **LHC Phase-2**, depending on the **success** of the **File-based Event Builder** during the LHC Run 3, the **DAQ** architecture will need to be **updated** and take advantage of **new technologies** as part of a **redesigned** system.



"Track 2 - Online and real-time computing" as Poster 14 in session "Poster" on Tuesday, the 22<sup>nd</sup> of October 2024

### CONTACTS



Name: Andrea Petrucci  
Organization: UCSD  
Email: andrea.petrucci@cern.ch  
Phone: +41227662563



Name: Srećko Morović  
Organization: UCSD  
Email: srecko.morovic@cern.ch  
Phone: +41754116276