

GPU based Level-3 Real Time Trigger for the Belle II High Level Trigger System

Ryosuke Itoh (KEK)
on behalf of
Belle II DAQ Group

Real Time 2026, La Biodola, Elba, Italy
May 25, 2026



1. The Belle II experiment

- The Belle II experiment is a new generation B-factory experiment at KEK in Japan to search for New Physics(NP).
- A huge number of B mesons are produced in e^+e^- collisions in the SuperKEKB accelerator and NP is searched for in B meson decays by the Belle II detector complex.

KEKB to SuperKEKB

- ◆ Nano-Beam scheme
extremely small β_y^*
low emittance
- ◆ Beam current double

$$L = \frac{y_z}{2\epsilon_r} \left(1 + \frac{\sigma_x^2}{\sigma_x^2} \frac{L_x \bar{\epsilon}_{x,y}}{\beta_y^*} \frac{R_L}{R_y} \right)$$

40 times higher luminosity
 $2.1 \times 10^{34} \rightarrow 8 \times 10^{35} \text{ cm}^{-2}\text{s}^{-1}$

Injector Linac upgrade

DR tunnel

New e+ Damping Ring

Reinforce RF systems for higher beam currents 2

Redesign the lattice to squeeze the emittance (replace short dipoles with longer ones, increase wiggler cycles)

Colliding bunches

New superconducting final focusing magnets near the IP

New HER wiggler section

Replace beam pipes with TiN-coated beam pipes with antechambers

KL and muon detector:
Resistive Plate Counter (barrel)
Scintillator + WLSF + MPPC (end-caps)

Particle Identification
Time-of-Propagation counter (barrel)
Prox. focusing Aerogel RICH (fwd)

EM Calorimeter:
CsI(Tl), waveform sampling (barrel)
Pure CsI + waveform sampling (end-caps)

Vertex Detector PXD+SVD
2 layers DEPFET + 4 layers DSSD

Central Drift Chamber
He(50%):C₂H₆(50%), Small cells, long lever arm, fast electronics

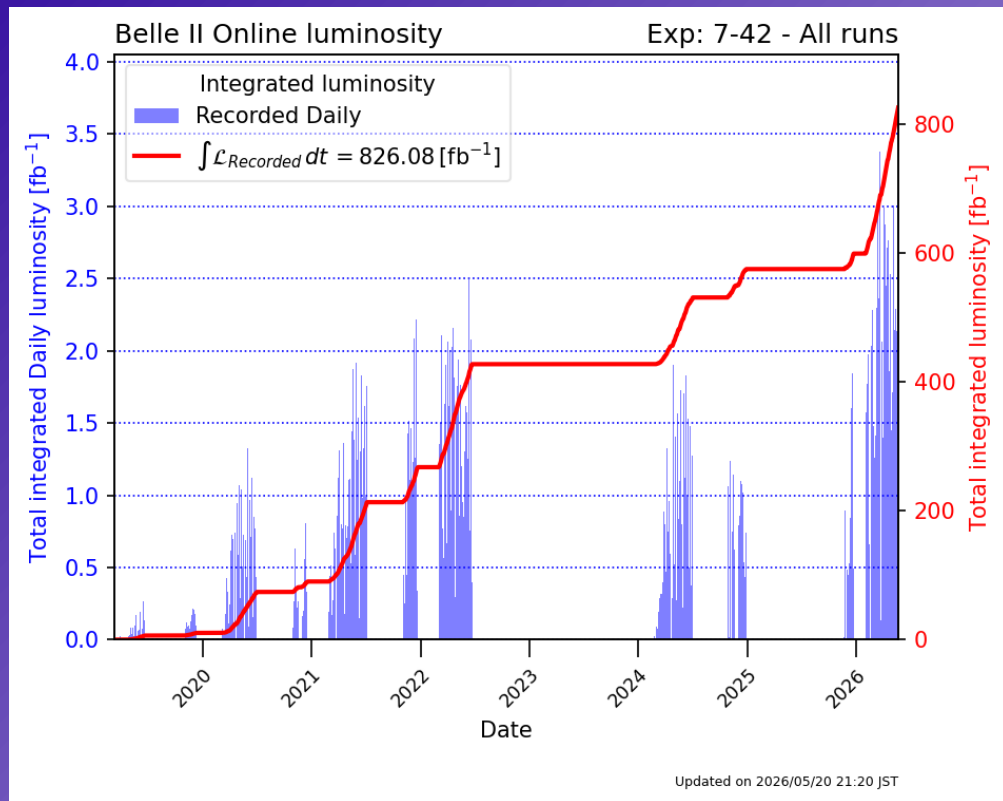
electron (7GeV)

positron (4GeV)

Beryllium beam pipe
2cm diameter

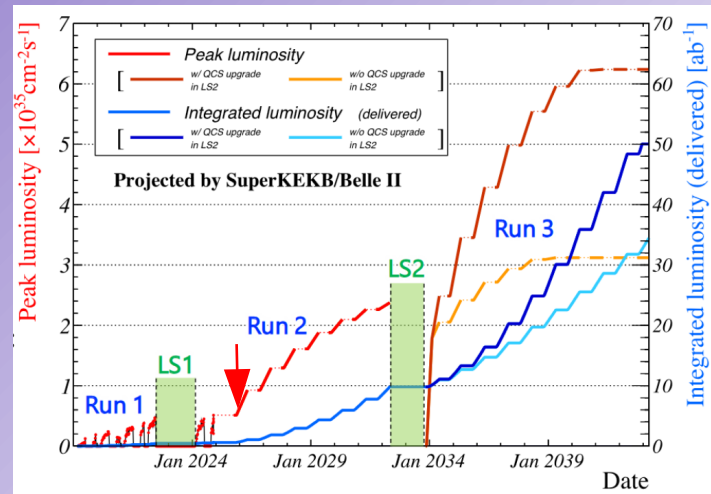
CDC

Belle II operation history

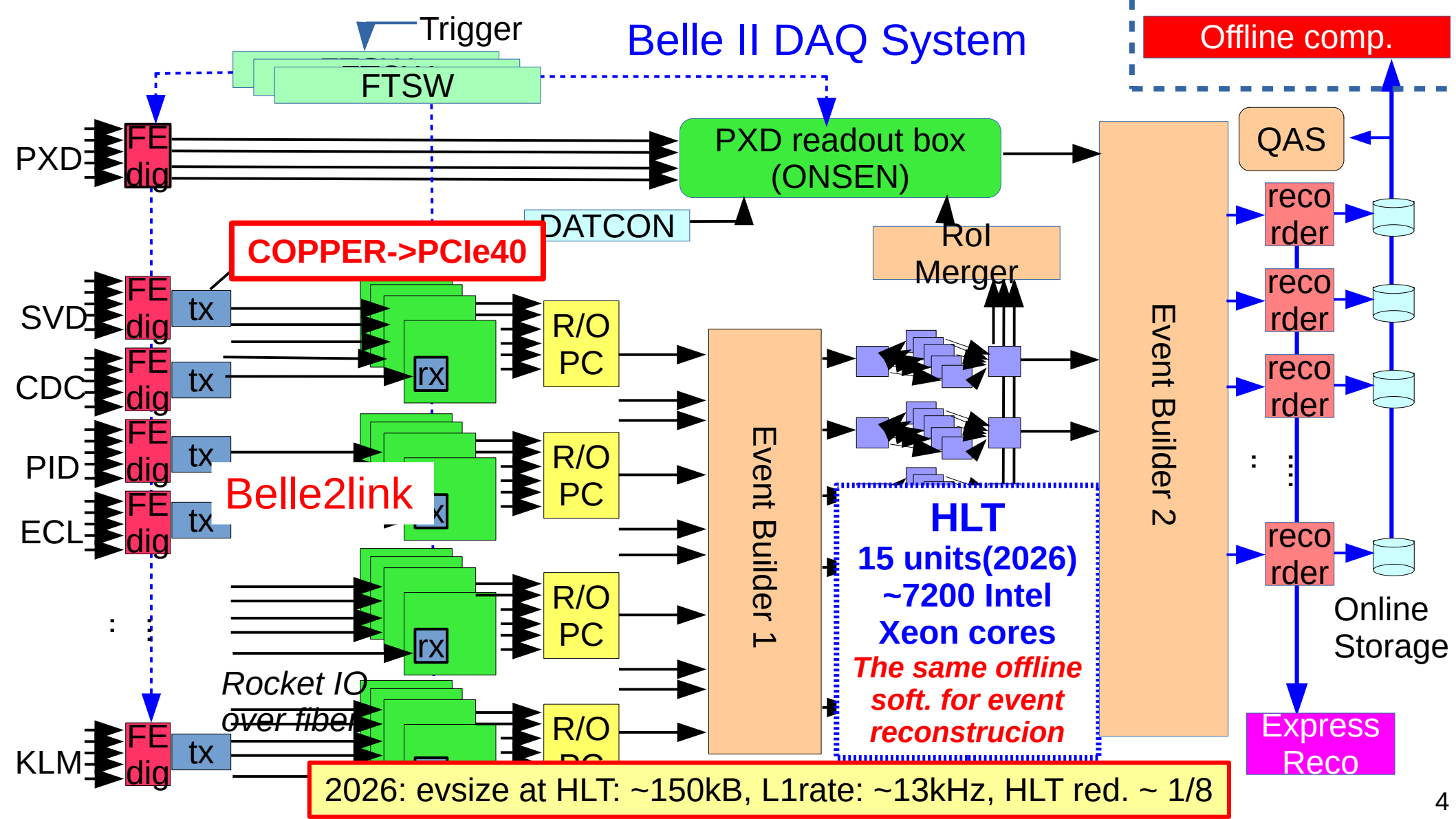


* Run 2 (Jan, 2024 -)
 $L_{\text{peak}} = 5.24 \times 10^{34} \text{ cm}^{-2}\text{sec}^{-1}$
 Integrated L $\sim 830 \text{ fb}^{-1}$
 (Target luminosity $6 \times 10^{35} \text{ cm}^{-2}\text{sec}^{-1}$)

* Many difficulties in accelerator operation.
 -> Pixel Detector (PXD) is temporarily turned off.



Belle II DAQ System



2026: evsize at HLT: ~150kB, L1rate: ~13kHz, HLT red. ~ 1/8

STOP

ABORT

AUTOPILOT MODE

Full control to update conf

STOP **ABORT**

Reload GUI

Exp # : 42

Run # : 45

Sub # : 0

Run control

RUNNING

TTD Status

RUNNING

of Opened Files

360

Detector states (use full ABORT before you check or uncheck a subsystem)

<input type="checkbox"/> PXD	OFF	<input checked="" type="checkbox"/> TOP	RUNNING	<input checked="" type="checkbox"/> KLM	RUNNING
<input checked="" type="checkbox"/> SVD	RUNNING	<input checked="" type="checkbox"/> ARICH	RUNNING	<input checked="" type="checkbox"/> TRG	RUNNING
<input checked="" type="checkbox"/> CDC	RUNNING	<input checked="" type="checkbox"/> ECL	RUNNING	HLT	RUNNING
				ROI_EB2	OFF

Trigger / Data status

Trigger input	# events : 40485687	Rate : 13.02 kHz	Trigger output	# events : 26017485	Rate : 12.94 kHz	Run start: 2026-04-28 16:31:58
---------------	----------------------------	-------------------------	----------------	----------------------------	-------------------------	---------------------------------------

	HLT01	HLT02	HLT03	HLT04	HLT05	HLT06	HLT07	HLT08	HLT09	HLT10	HLT11	HLT12	HLT13	HLT14	HLT15
# events :	1736276	1734825	1734834	1734869	1734937	1734936	1734999	1735045	1735069	1735137	1735186	1735244	1735244	1735322	1735345
Rate :	0.8 kHz	0.8 kHz	0.9 kHz	0.8 kHz	0.8 kHz	0.9 kHz	0.9 kHz	0.9 kHz	0.9 kHz	0.8 kHz	0.9 kHz	0.9 kHz	0.8 kHz	0.9 kHz	0.8 kHz
Flow :	23.5 MB/s	24.3 MB/s	21.2 MB/s	17.0 MB/s	22.8 MB/s	21.3 MB/s	25.1 MB/s	23.0 MB/s	20.6 MB/s	23.2 MB/s	21.1 MB/s	27.7 MB/s	21.1 MB/s	16.4 MB/s	23.7 MB/s

<p>PXD Run # : 38</p> <p>NOTREADY</p> <p>LOAD</p> <p>ABORT</p> <p>BOOT</p>	<p>ROI_EB2 Run # : 38</p> <p>NOTREADY</p> <p>LOAD</p> <p>ABORT</p>	<p>(Managed by DAQ group)</p> <p>EB2TX5 NOTREADY - EB2TX10 NOTREADY - EB2TX15 NOTREADY - EB2TX20 NOTREADY - EB2TX25 NOTREADY -</p> <p>EB2TX1 NOTREADY - EB2TX11 NOTREADY - EB2TX16 NOTREADY - EB2TX21 NOTREADY - EB2TX26 NOTREADY -</p> <p>EB2TX2 NOTREADY - EB2TX7 NOTREADY - EB2TX12 NOTREADY - EB2TX17 NOTREADY - EB2TX22 NOTREADY - EB2TX27 NOTREADY -</p> <p>EB2TX3 NOTREADY - EB2TX8 NOTREADY - EB2TX13 NOTREADY - EB2TX18 NOTREADY - EB2TX23 NOTREADY - EB2TX28 NOTREADY -</p> <p>EB2TX4 NOTREADY - EB2TX9 NOTREADY - EB2TX14 NOTREADY - EB2TX19 NOTREADY - EB2TX24 NOTREADY - EB2TX29 NOTREADY -</p>	<p>SVD Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>
---	---	--	--

<p>CDC Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>	<p>TOP Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>	<p>ARICH Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>	<p>ECL Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>	<p>KLM Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>	<p>TRG Run # : 45</p> <p>RUNNING</p> <p>STOP</p> <p>ABORT</p>
--	--	--	--	--	--

HLT Run # : 45

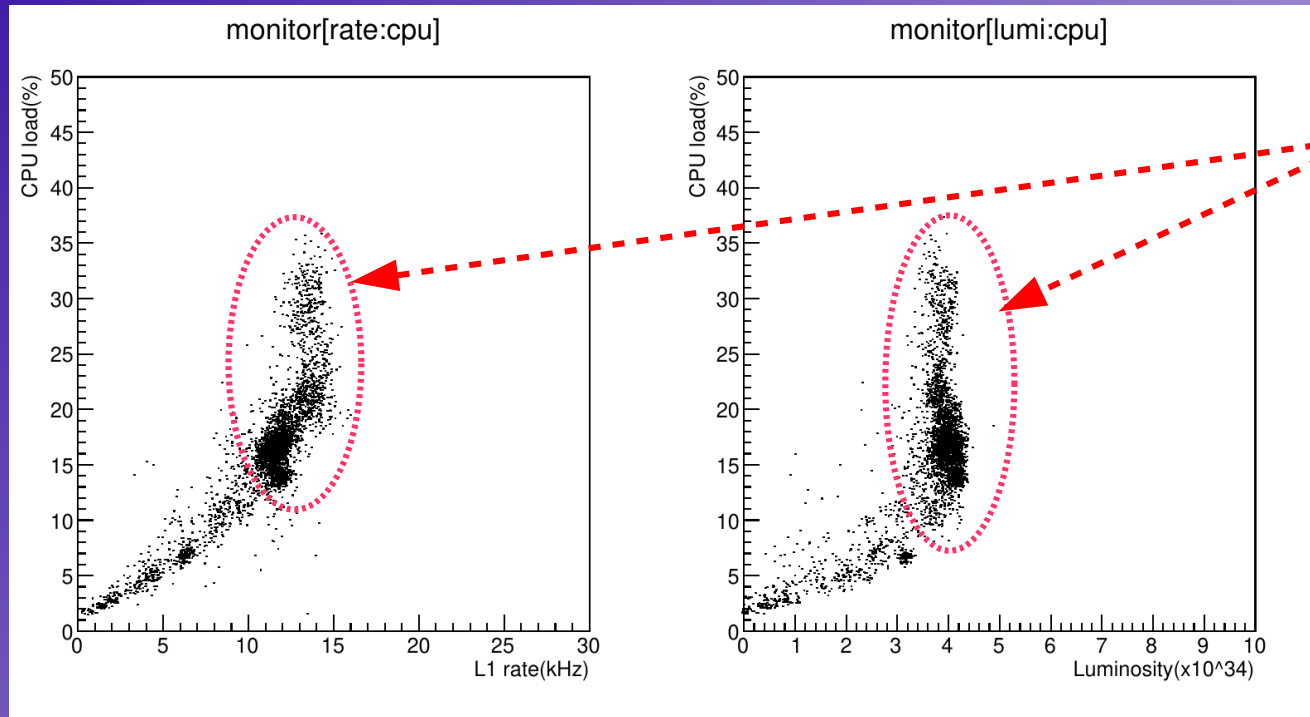
<input checked="" type="checkbox"/> RC_HLT01 RUNNING - RC_STORE01 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT06 RUNNING - RC_STORE06 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT11 RUNNING - RC_STORE11 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_ERECO01 RUNNING - <input checked="" type="checkbox"/> DQMHLT RUNNING - <input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> RC_HLT02 RUNNING - RC_STORE02 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT07 RUNNING - RC_STORE07 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT12 RUNNING - RC_STORE12 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_ERECO02 RUNNING - <input checked="" type="checkbox"/> DQMereco RUNNING - <input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> RC_HLT03 RUNNING - RC_STORE03 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT08 RUNNING - RC_STORE08 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT13 RUNNING - RC_STORE13 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_ERECOPHY RUNNING - <input checked="" type="checkbox"/> DQMPHYSICS RUNNING - <input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> RC_HLT04 RUNNING - RC_STORE04 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT09 RUNNING - RC_STORE09 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT14 RUNNING - RC_STORE14 RUNNING - <input checked="" type="checkbox"/>	
<input checked="" type="checkbox"/> RC_HLT05 RUNNING - RC_STORE05 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT10 RUNNING - RC_STORE10 RUNNING - <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RC_HLT15 RUNNING - RC_STORE15 RUNNING - <input checked="" type="checkbox"/>	

PXD is turned off

2. Motivation

- Current SuperKEKB status
 - * Although the luminosity is still low, **the background is unexpectedly higher.**
 - * This situation may persist as the luminosity continues to rise.
- The HLT CPU load is higher even though the physics event rate is far below the design value.
 - * **Current HLT (with 15 units) already has a sufficient processing power to manage “physics rate” at the final design luminosity.**
 - * The main issue is **the higher CPU consumption caused by the “unexpectedly high background”.**

Current HLT performance



Higher
CPU
consumption
than
expected

Data taken in
early 2026 runs

- Heavy noise hits caused by the background are contaminated in the taken events.
 - > Source of the unexpected higher CPU consumption.

Efficient suppression of such noise hits/events is the key to mitigate the problem.

Current approaches for the mitigation

[Hardware]

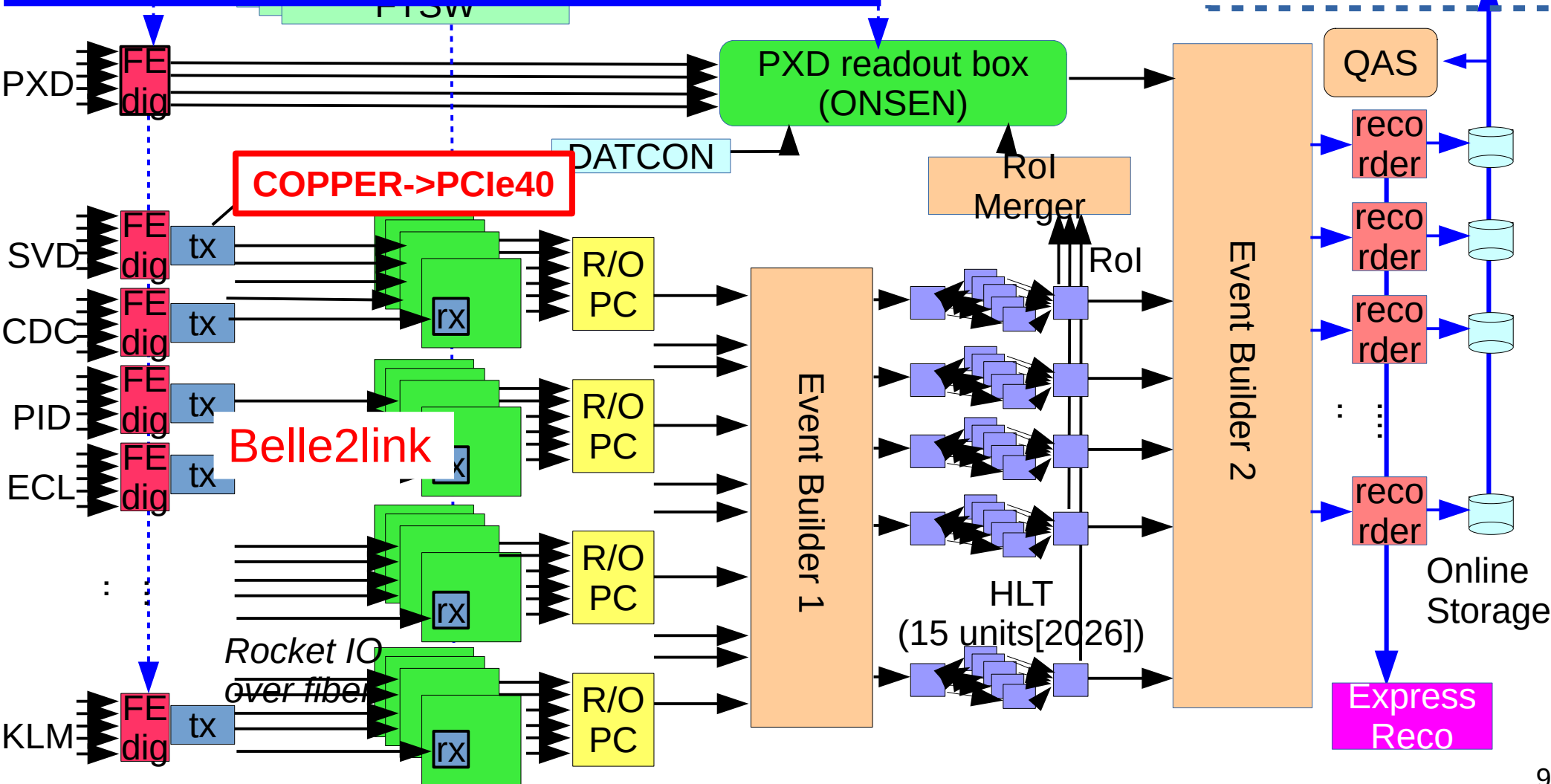
- * Continue to add more PC servers in HLT units : **baseline upgrade**
 - Difficult to prepare a large CPU power in a short period with a limited budget.
(only up to 5-10% performance increase per year)

[Software]

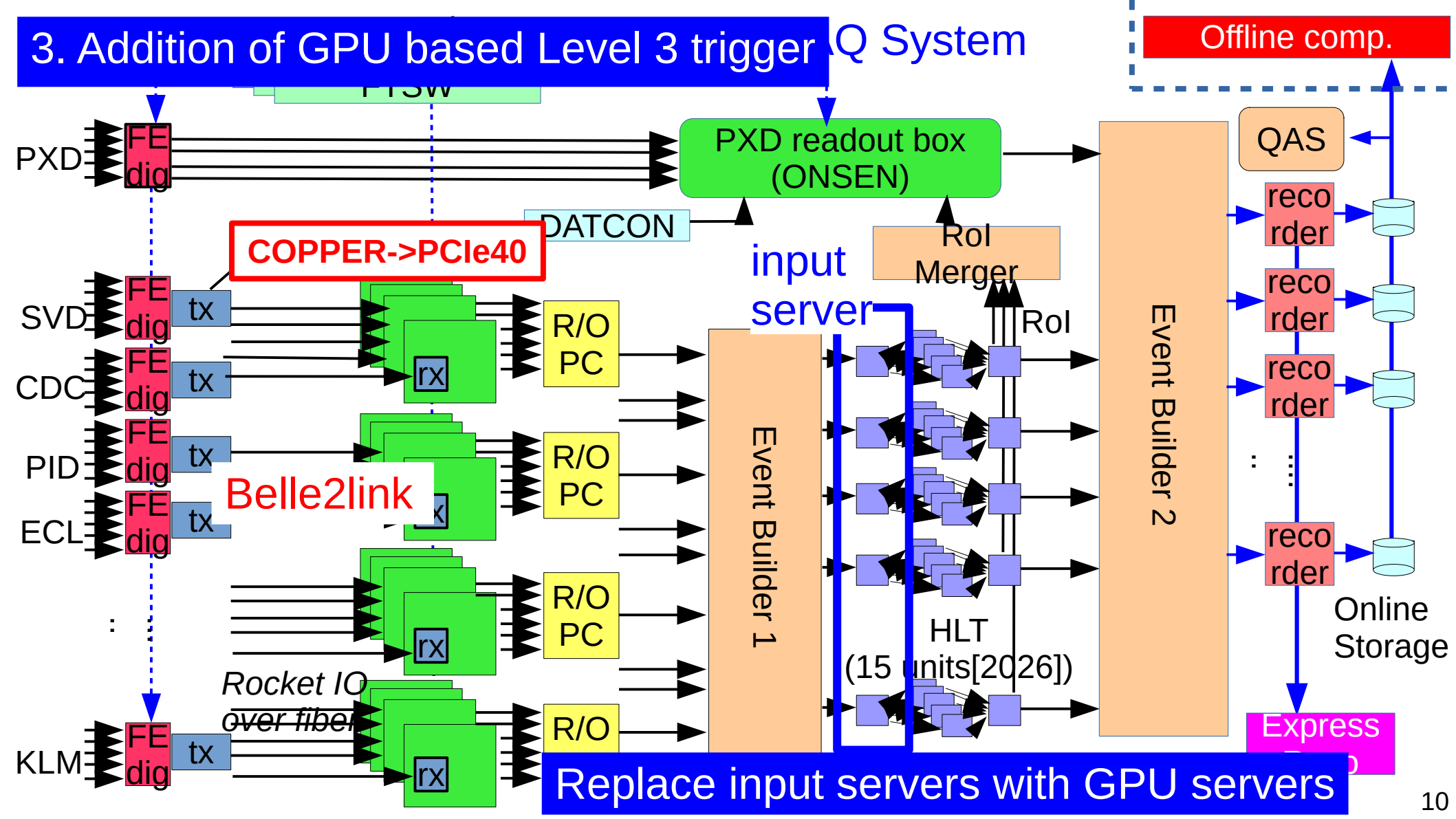
- * Signal “cleaning/filtering”
 - > already done. 10-20% reduction of CPU power.
 - * Implementation of software Level 3 trigger in HLT processing chain.
 - > being prepared
- => **But still worry for the lack of CPU power.**

3. Addition of GPU based Level 3 trigger Q System

Offline comp.



3. Addition of GPU based Level 3 trigger Q System



Choice of GPU for Belle II HLT

- * A GPU houses >10,000 CUDA cores.
- * Adding 1 GPU for each of 15 HLT input servers is enough for our purpose.
 - > in total of 150,000 CUDA cores
- * GPU with better performance in fp64 is desired but hopefully less dependent (assumption to rely mostly on fp32).
- * Vector/Tensor cores are not supposed to be used for the planned processing.

RTX 5090 : 21,760 cores, 32GB, ~100TFLOPS(fp32)

RTX 5000Ada : 12,448 cores, 32GB, 61TFLOPS(fp32)

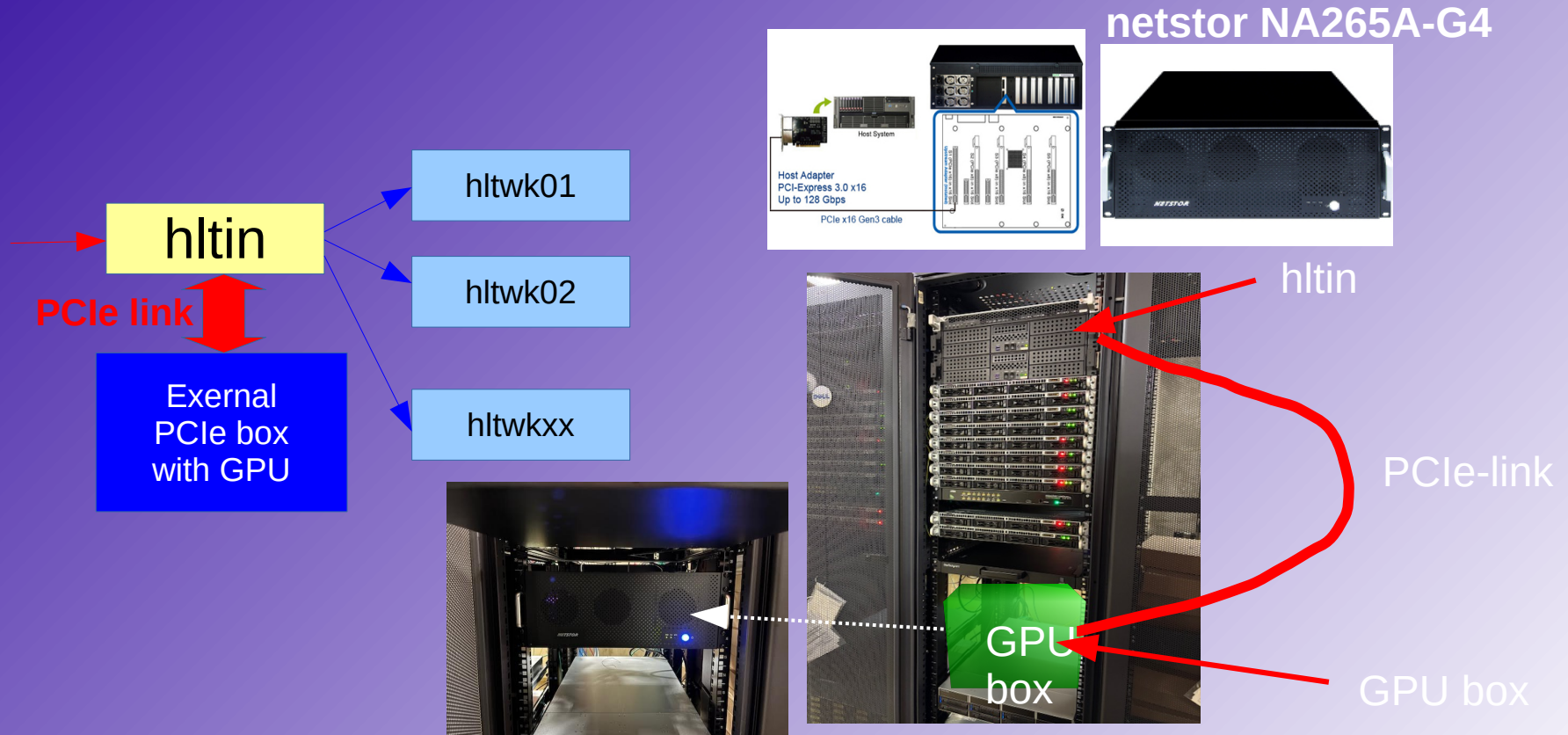
RTX 6000Ada : 18,176 cores, 48GB, 91.1TFLOPS(fp32)

H100 : 14,592 cores, 80GB, 30TFLOPS(fp32)

A2 : 1,280 cores, 16GB, 4.5TFLOPS(fp32) : (used in test bench)

GPU Implementation in existing HLT unit (unit 15)

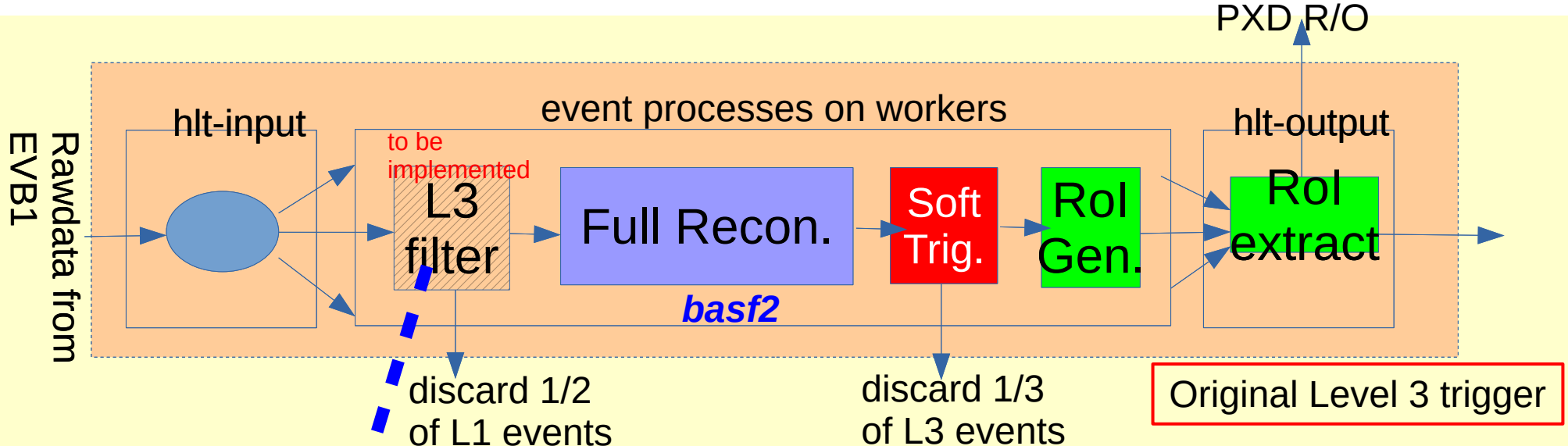
- Place GPU in an external rack-mount 4U unit (eGPU).
- Connect it to existing HLT input server via a special PCIe link.



Design Policy

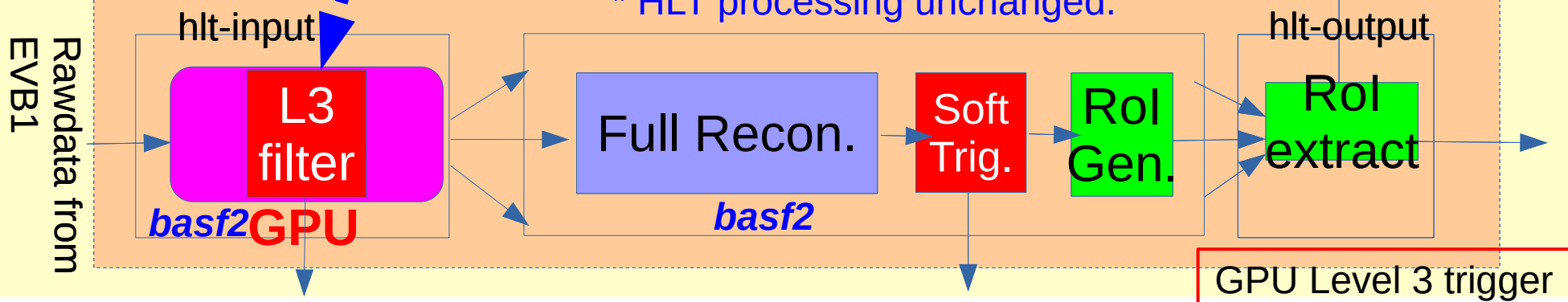
- Keep existing HLT as is.
 - * No change in the HLT hardware/software.
Current HLT software processing is not touched at all.
- Add GPU processing before the existing HLT framework as a part of the HLT system.
 - * **Level 3 trigger software is implemented in GPU.**
 - * Filter out background events before the actual HLT processing.
- Build GPU software development environment fully compatible with the Belle II software library coded mostly in C++.
“Seamless software environment between CPU and GPU”

GPU-L3 project



Level 3 prefilter: Discard background events by eliminating events from off-IP by fast event reconstruction with limited detectors (CDC+ECL).

* HLT processing unchanged.



4. Seamless Software Development Environment

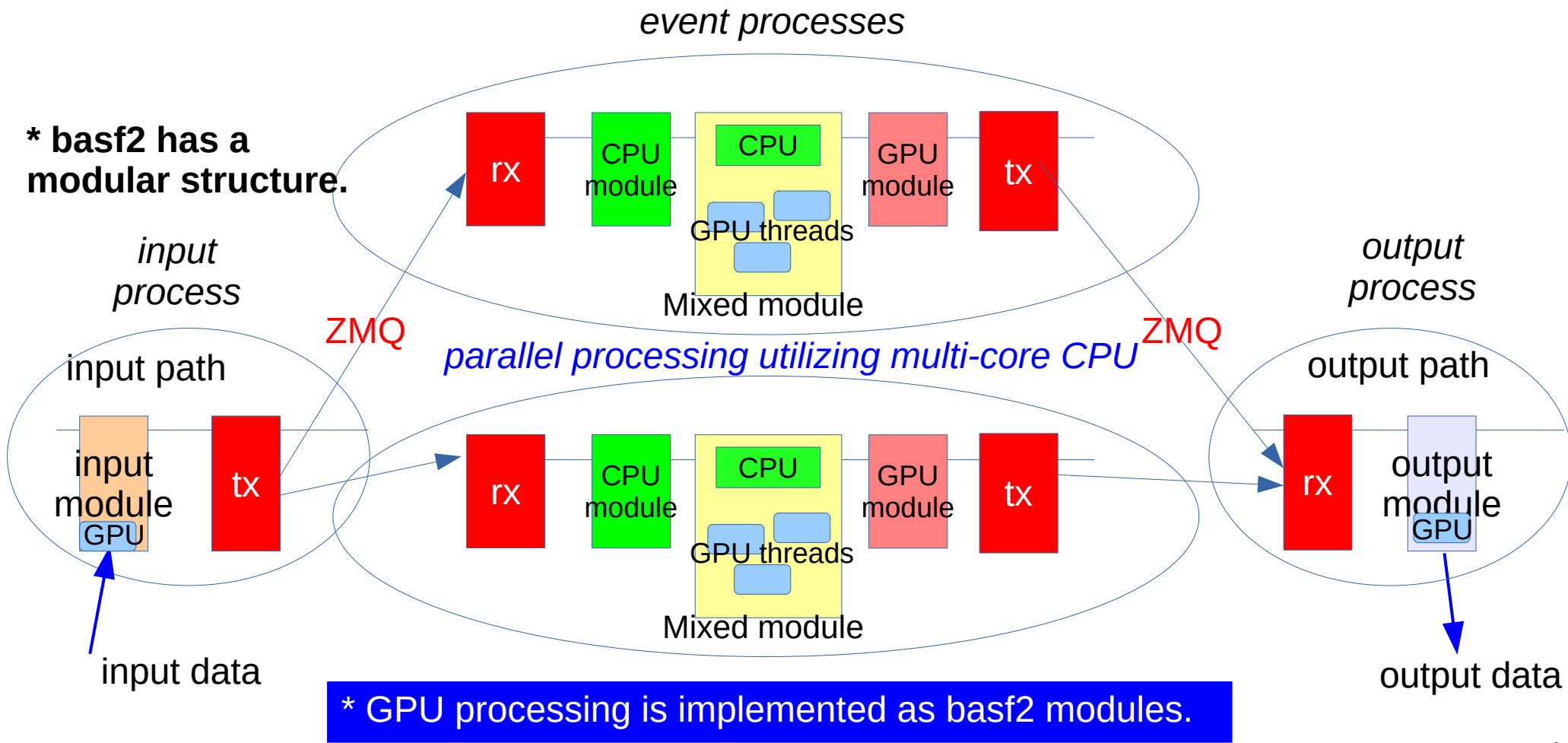
- The software on GPU is desired to be compatible with the existing Belle2 software coded mostly in C++.
- GPU is supposed to be programmed in CUDA.
- The Belle2 software is developed using “SCons” and runs on a framework called “basf2” as a set of dynamically-linked modules.



- CUDA compiler (nvcc) was added to SCons build system.
- Development environment was built so that **mixed C++ and CUDA codes can co-exist in Belle2 library and dynamically loaded in basf2.**

=> Seamless execution of CUDA-coded GPU modules and C++-coded CPU modules on the same basf2 became possible.

basf2: Belle 2 event processing framework commonly used in online and offline with the parallel processing capability on multicore CPU.



5. Data Flow for GPU

- Batch processing nature of GPU

- * GPU executes the same instruction chain for a group of different input events on a huge number of CUDA cores (threads) in parallel.
 - > “batch” processing of many events started at the same time.

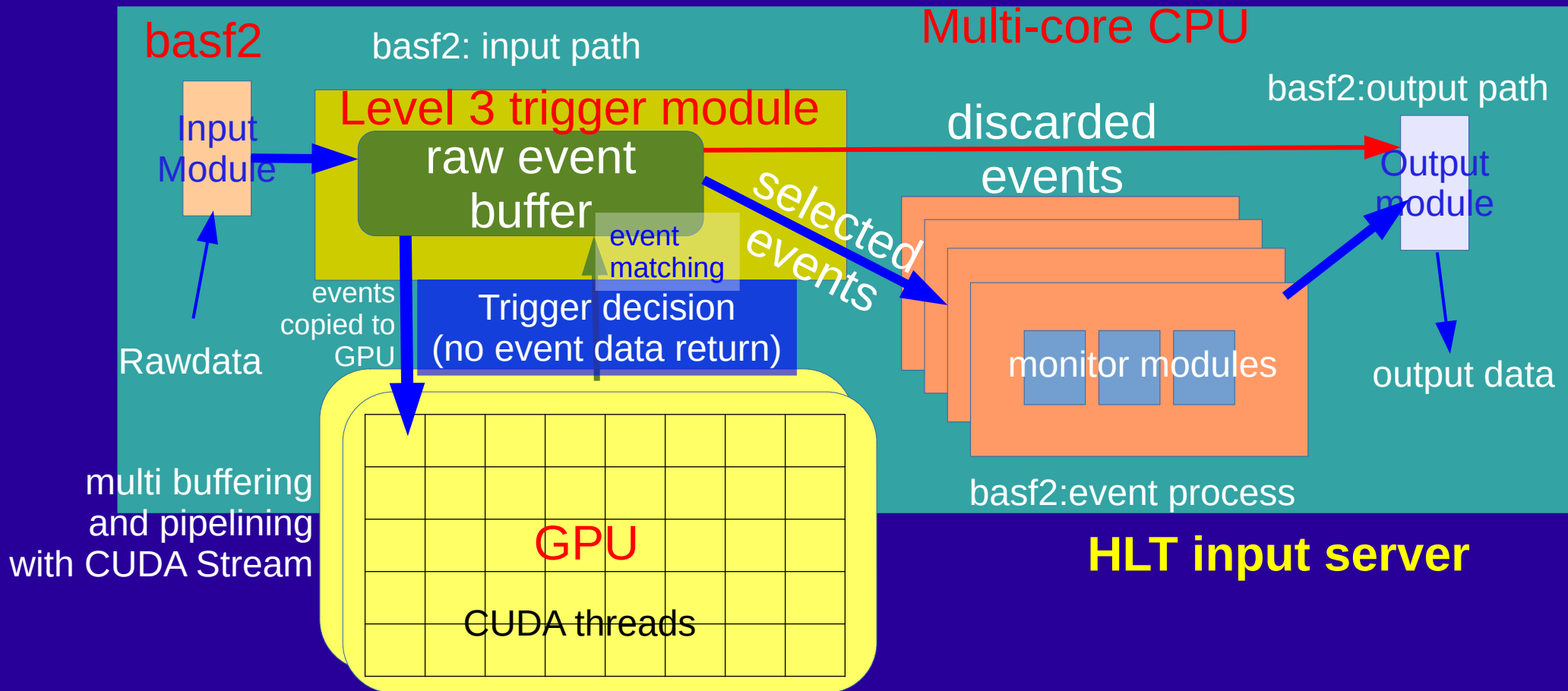
- * Processing of one “batch”

1. buffer a bunch of events in the host memory.
2. copy them to GPU memory
3. start GPU batch processing and wait for the completion
4. transfer back results and compare with events in host memory
5. send chosen events to next basf2 module

- Pipelining is necessary to average out the bursty flow rate.

 - > “**CUDA Stream**” is utilized to implement multi-buffering.

Actual data flow for basf2+GPU in HLT input server



“Batch” processing is started when all thread buffers are filled.

6. Implementation of Level 3 trigger code

- “Reference” Level 3 trigger software is now being implemented on GPU by recycling the old code developed in C++ for the practice of code porting to GPU and performance evaluation.
 - <- *Consisting of fast CDC tracking and ECL clustering.*

Various tips in the porting of C++ code to CUDA

- Database access from GPU is implemented by copying the parameter memory block initialized on host from the database to GPU’s global memory.
- Dynamic object creation by C++ new() is replaced with a private GPU memory management using the “placement new()” so that the objects are created in the GPU’s global memory.
- STL (such as std::vector) / ROOT functions are replaced with home-grown libraries.
- “Double precision” floating point operations are replaced with single precision ones so as to maximize the GPU performance.

Geometry/Constants DB access from GPU

Host side

The same geometry/constant class(CUDA)

GPU side

Constructor GPU CDC Geometry

Constructor GPU CDC Geometry

`__host__ set_block()`

- pointer to data blk

`__host__ load_params()`

- establish link to DB

- load parameters
in data blk.

`__device__ set_block()`

- pointer to data blk

`__host__ __device__ get_param()`

- obtain parameters



allocated
outside

Host Memory
struct
(data block)

cuda
Mem
cpy

GPU Global Memory
(shared by all threads)
struct
(data block)

allocated
outside

Current status of GPU-L3 project

Item	Status	Remark
Seamless software environment	Done and working	Still some problem (-j 1 is necessary...)
Data transport from host to GPU (and reverse)	Done and working	Pipelined data flow is implemented/tested.
Raw data unpacking (CDC/ECL/TRG)	CDC -> done ECL/TRG -> not yet	channel->wire conversion with DB
Geometry database access	CDC -> done ECL-> next	Full CDC geometry is available on GPU
CDC Track Finding/Fitting	In progress	Compilation passed Debugging in progress
ECL clustering	not yet	
TRG processing	not yet	
Level 3 decision	not yet	

7. Summary

- A significant shortage in the processing power of Belle II HLT is foreseen because of the unexpected high background.
- The addition of GPU in HLT is being planned by porting existing C++ based reconstruction software.
- A seamless development environment between CPU and GPU has been implemented in Belle2 software library.
- The pipelined data flow scheme for GPU on basf2 framework has been developed and established.
- C++ based Level3 trigger software is being ported to GPU.
- The performance will be studied in beam from this fall.

Backup Slides

Application of Level 3 reconstruction

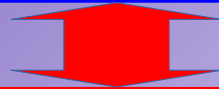
- Noise hits within the physics events also cause a high CPU consumption in the event processing.
- The results of Level 3 tracking can be used for the reduction of such noise hits.
- The tracking results are sent to main HLT processing as “Region of Interest”. The hits outside the “RoI” are discarded in the main tracking so as to reduce the CPU consumption.

Validation of GPU code in offline

- The insertion of GPU based L3 trigger may affect on the physics quality and it has to be validated in MC simulation.
- But it is difficult to prepare GPUs in offline computing facility for the MC generation.
- CUDA can generate the binary for CPU and GPU simultaneously from the same source code. (`__host__ __device__` function).
- Compare the results between GPU and CPU execution and if the difference is within a tolerable range, we can rely on the CPU code only for MC.

Machine learning approach for Level 3 trigger

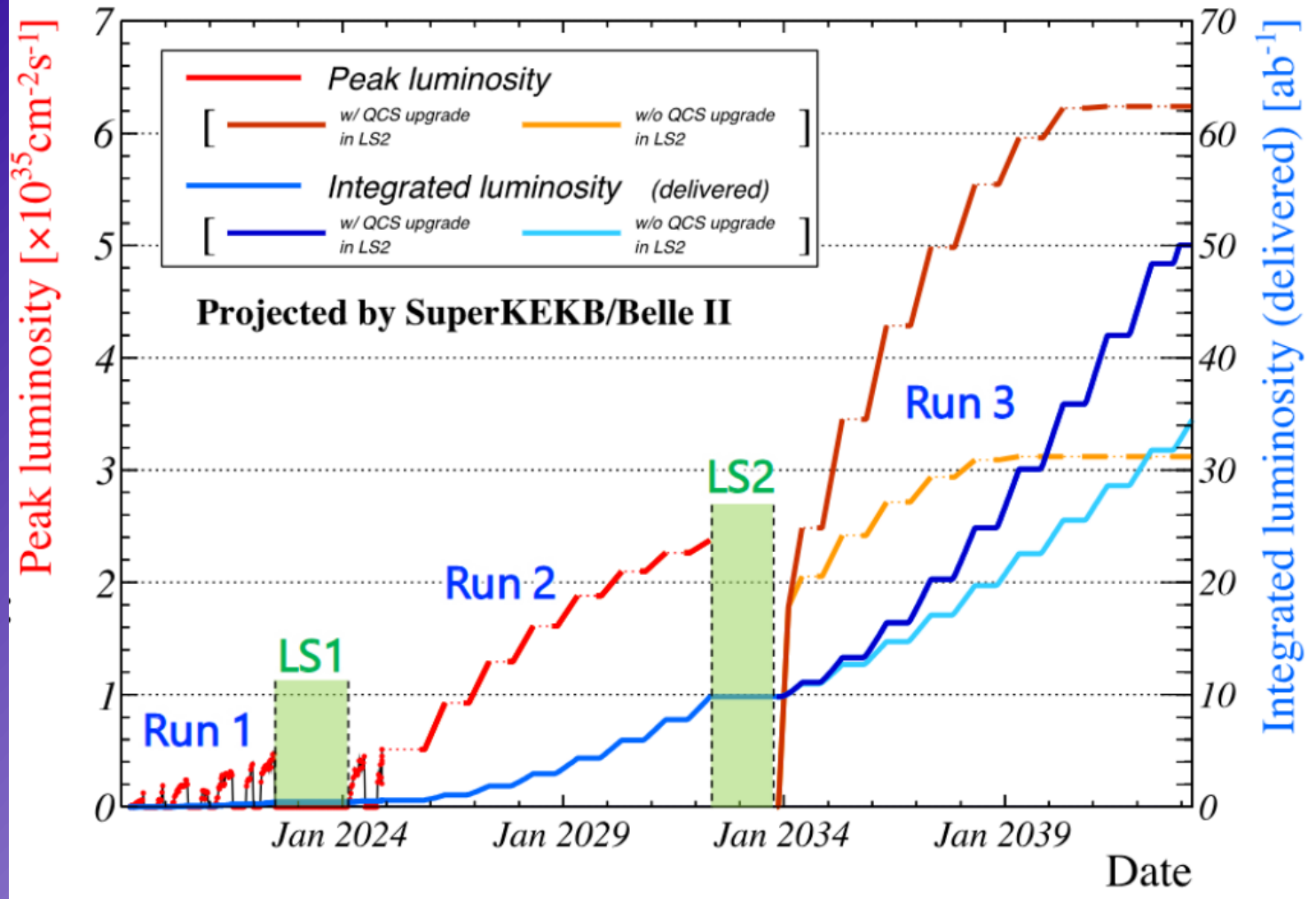
- Current design of Level 3 code is based on the conventional CDC tracking + ECL clustering approach (reference design).
- But new algorithm can be developed using the machine learning technique. Study is on-going in Belle II (Shandong U).
- The GPU-capable Belle II software environment can be utilized for the development.
- Possibility to utilize vector/tensor cores in GPU in this case.



However, the development and validation of new code will take significantly longer time...

Coding Notes:

- CUDA is mostly compatible with gcc and existing C++ codes can be compiled by “nvcc” without modifications.
 - * But compiler version should be consistent.
- Functions which run on GPU are marked with a special header key `__global__` or `__device__`.
- A basf2 module can be compiled with nvcc where C++ codes on host and CUDA codes on GPU are mixed.
- The compiled objects are linkable as a shared object which can be dynamically loaded as a basf2 module.



Number of physical cores in HLT in 2025c-2026ab run

HLT01: (already renovated)

$$28 \text{ cores} * 1 + 36 * 2 + 40 * 3 + 48 \text{ cores} * 5 = 460 \text{ cores/u}$$

HLT02-05

$$20 \text{ cores} * 16 + 36 \text{ cores} * 2 + 40 \text{ cores} * 2 = 472 \text{ cores/u} \rightarrow 1888$$

HLT06-10

$$28 \text{ cores} * 12 + 36 \text{ cores} * 2 + 40 \text{ cores} * 2 = 488 \text{ cores/u} \rightarrow 2440$$

HLT11-12, 14-15

$$48 \text{ cores} * 10 = 480 \text{ cores/u} \rightarrow 1920$$

HLT13 (reserved as a test bench)

$$16 \text{ cores} * 3 + 20 \text{ cores} * 6 + 48 \text{ cores} * 8 = 452 \text{ cores/u} \rightarrow 552$$

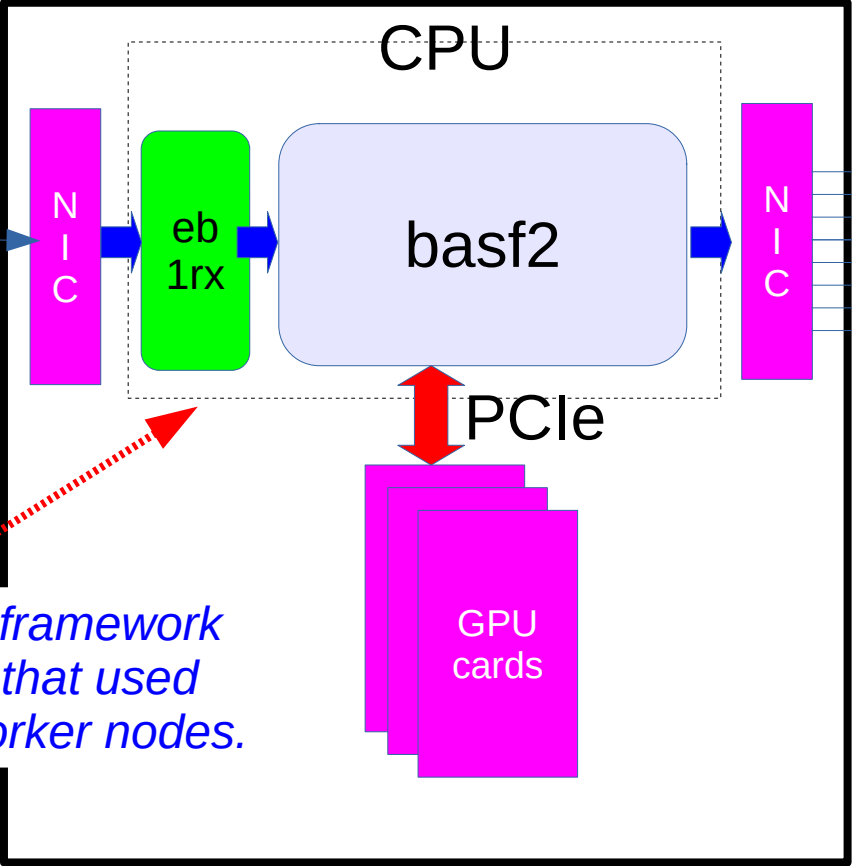
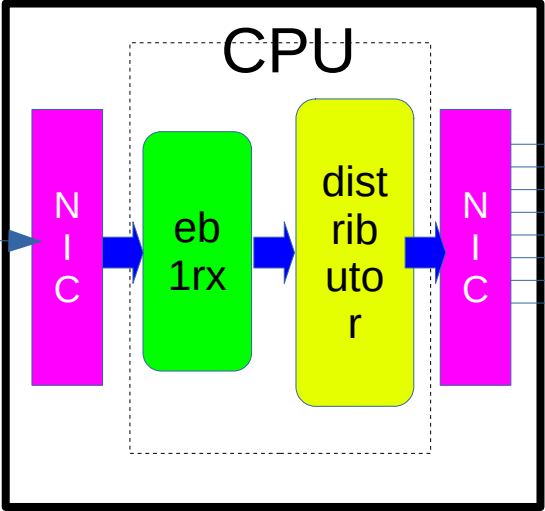
 64 x "20 core" servers were purchased in 2014-2015

-> Renovation (rep. with ≥ 48 core servers) is planned from 2025

In total of 7260 cores (c.f. design:6400 cores)

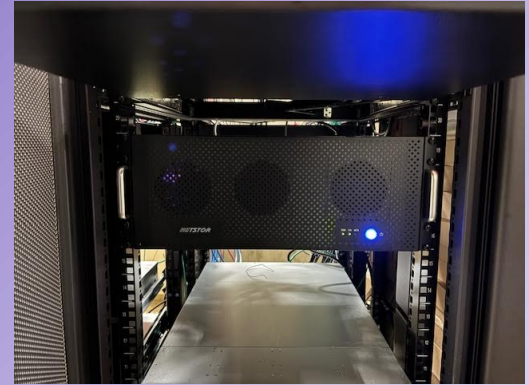
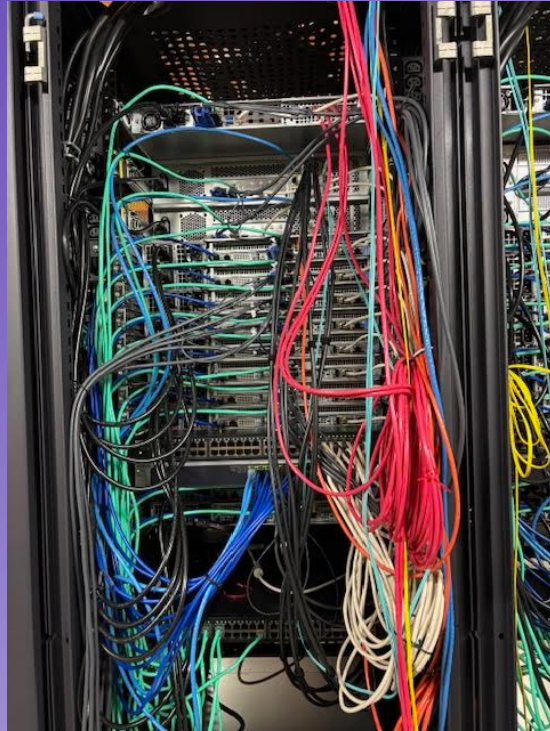
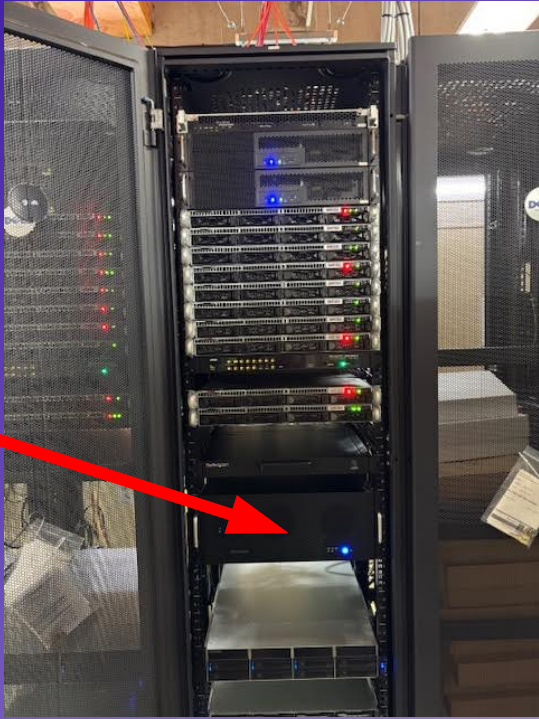
hltin with GPU

current hltin



Software framework similar to that used in HLT worker nodes.

HLT15



Test bench system

*recycled from
another purpose*



```
itohepyc% lscpu
Architecture:          x86_64
CPU op-mode(s):       32-bit, 64-bit
Address sizes:        48 bits physical, 48 bits virtual
Byte Order:           Little Endian
CPU(s):               96
On-line CPU(s) list:  0-95
Vendor ID:            AuthenticAMD
Model name:           AMD EPYC 7643 48-Core Processor
CPU family:           25
Model:                1
Thread(s) per core:  2
Core(s) per socket:  48
Socket(s):            1
Stepping:             1
Frequency boost:      enabled
CPU(s) scaling MHz:  63%
CPU max MHz:          3640.9170
CPU min MHz:          1500.0000
BogoMIPS:             4599.87
```

CPU : AMD EPYC 7643
(48 cores/96 threads)

```
itohepyc% cat /etc/redhat-release
Rocky Linux release 9.5 (Blue Onyx)
```

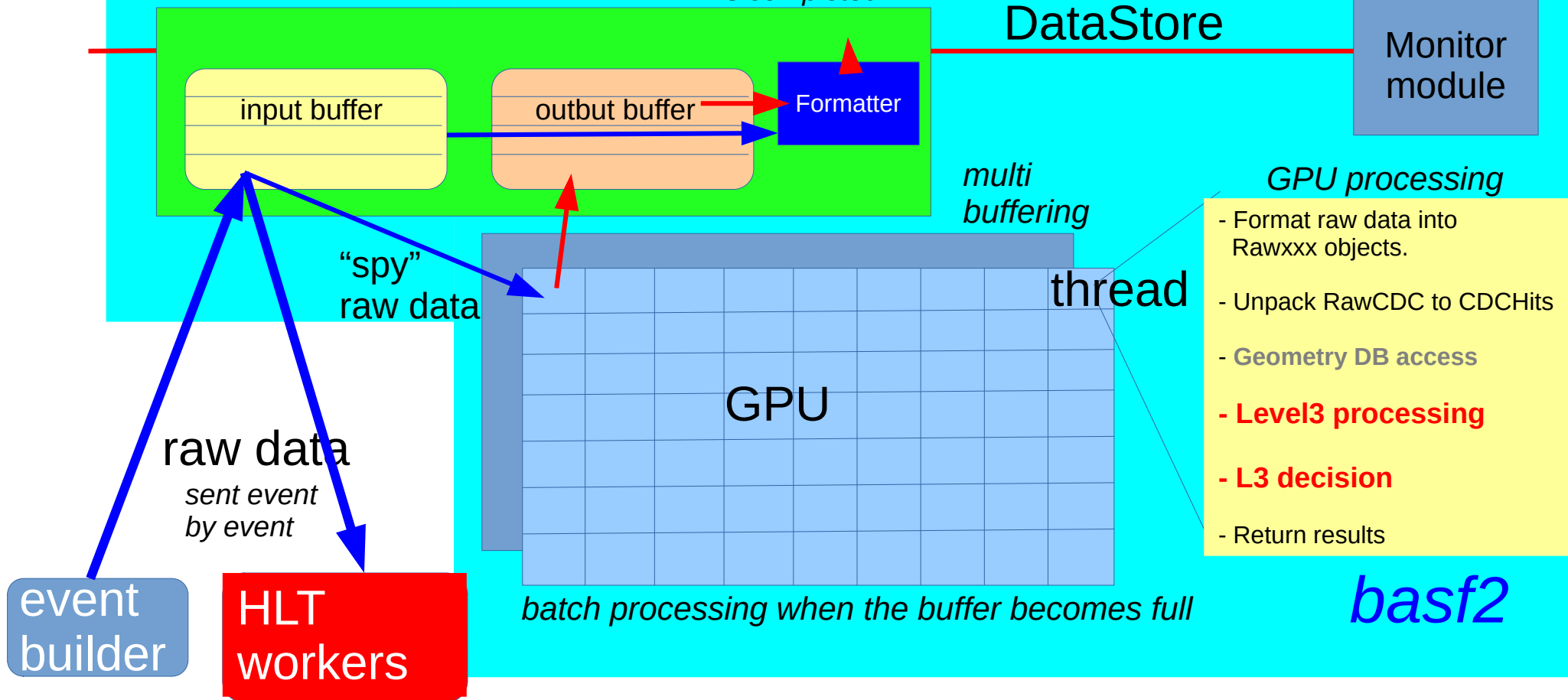
```
itohepyc% nvidia-smi
Thu Jan 30 13:59:22 2025
```

NVIDIA-SMI 565.57.01			Driver Version: 565.57.01		CUDA Version: 12.7			
GPU	Name	Perf	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr.	ECC
Fan	Temp		Pwr:Usage/Cap		Memory-Usage	GPU-Util	Compute M.	MIG M.
0	NVIDIA A2	P8	Off	00000000:81:00:0	Off	0%	Default	0
0%	33C		5W / 60W	1MiB /	15356MiB			N/A

GPU : NVIDIA A2 Tensor
(1280 CUDA cores;
10 RT cores;
40 Tensor cores)

basf2: Special Input module for performance test

HLT15:hltin



* Normal HLT data flow is not disturbed by this performance test.

Cost comparison between 1xPC server and 1xGPU

(1 PC server(Xeon Gold 6336Yx2) vs 1 GPU(RTX6000Ada))

🧠 Comparison: Xeon Gold 6336Y × 2 vs NVIDIA RTX 6000 Ada

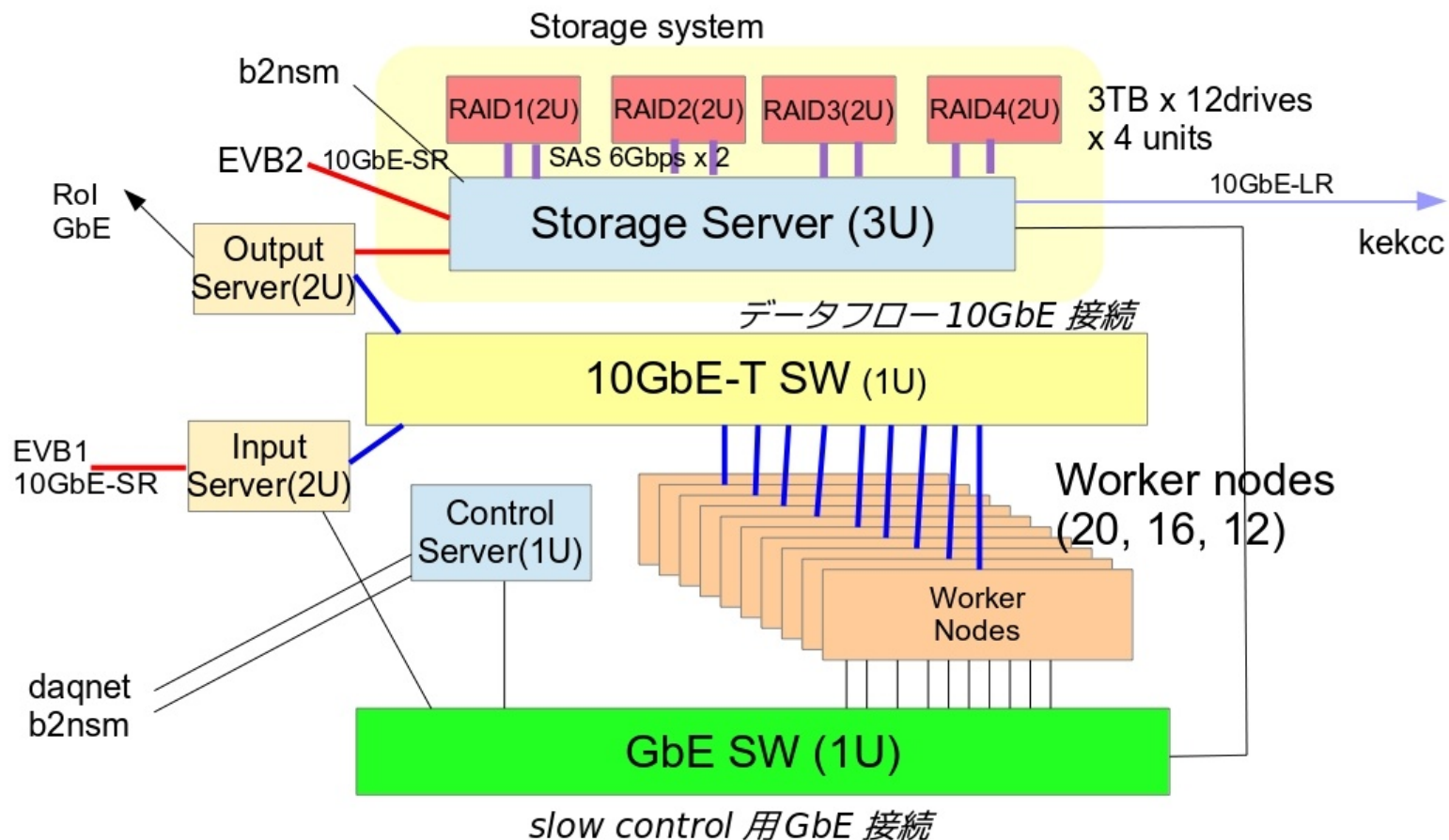
Category	Xeon Gold 6336Y × 2	NVIDIA RTX 6000 Ada
Architecture	Ice Lake-SP (10nm)	Ada Lovelace
CPU/GPU Cores	48 cores / 96 threads (24C×2)	18,176 CUDA cores
Clock Speed	2.4 GHz (Base), 3.6 GHz (Turbo)	~2.5 GHz (Boost)
FP64 (Double Precision)	~1.1–1.4 TFLOPS (combined)	~2.0 TFLOPS
FP32 (Single Precision)	~3–4 TFLOPS (combined, effective)	91.1 TFLOPS
Memory Capacity	Up to 2TB+ DDR4 ECC	48GB GDDR6 ECC
Memory Bandwidth	~200 GB/s (DDR4 x2 CPUs)	960 GB/s
Power Consumption (TDP)	~410W (total for 2 CPUs)	~300W
PCIe Interface	PCIe Gen 4.0 × 64 total (2 CPUs)	PCIe Gen 4.0 × 16
Price (JPY)	¥1,600,000–2,000,000 (CPU, RAM, motherboard, etc.)	¥850,000–900,000

- The processing power/GPU core is a bit lower than that of CPU, but number is incredibly more!
- >30 times faster performance (fp32).

* Processing power can be doubled only with a cost of 1-year HLT reinforcement!

- But fp64 performance is not so much....

Actual Implementation of HLT/Storage unit



Equipped in a single rack