# ONSEN
## (Online Selector Nodes)

Dennis Getzkow[2], Thomas Geßler[2], Wolfgang Kühn[2],
Jens Sören Lange[2], Klemens Lautenbach[2], Zhen-An Liu[1],
Björn Spruck[3], Jingzhou Zhao[1], (Leonard Koch[2], David
Münchow[2]), [1]IHEP Beijing, [2]Univ. Giessen, [3]Univ. Mainz

# Outline

- Overview of PXD DAQ
- ONSEN
    - Hardware status
    - Full system test at Giessen, results
    - Processing basf2 events in ONSEN
    - Answer to questions, raised in BPAC report 10/2016
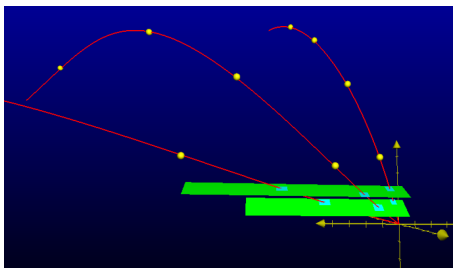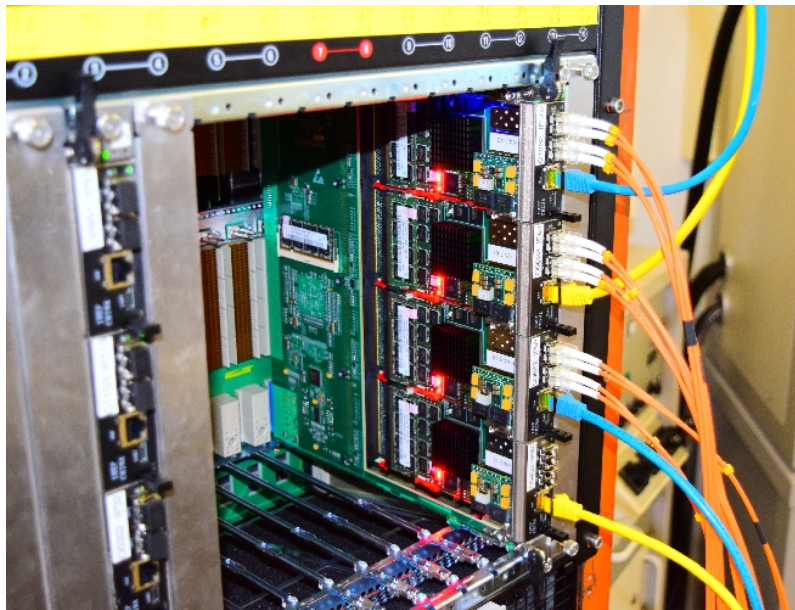
# PXD DAQ Overview

# PXD DAQ parameters

| |
|---|
| Trigger 30 kHz<br>(1/3 accept, 2/3 reject) |
| $\leq$3% PXD occupancy<br>data input $\leq$21.6 GB/s |
| ROI selection<br>(region of interest)<br>HLT (SVD+CDC), PC farm<br>DATCON (SVD only), FPGA<br>logical OR (on ONSEN)<br>data reduction factor $\geq$10 |

# ONSEN 1/8 system

# Status of ONSEN hardware



| ONSEN AMC card |
| --- |
| v4.0 (final) |
| Virtex-5 FX70T |
| 2 optical links (6.25 Gbps) |
| GbE |

| DATCON AMC card |
| --- |
| Virtex-5 LX50T |
| 4 optical links (3.125 Gbps) |

slow control / monitoring:
IPMI add-on boards (Mainz)

# Status of ONSEN hardware
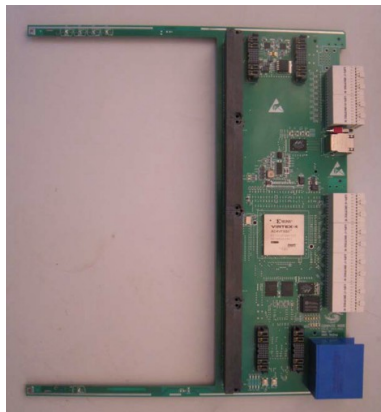
| |
|---|
| ONSEN xTCA carrier card |
| v3.3 (final) |
| Virtex-4 FX60 (switcher to ATCA backplane) |
| GbE |

| |
|---|
| add-on: RTM board power supply board |

# AMC card mass production

# ONSEN hardware status

| AMC v4.0 | |
|----|----|
| 10 | KEK |
| 8 | DESY |
| 4 | IHEP (repair) |
| 21 | Giessen |
| 43 | (total) |

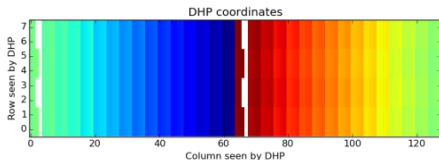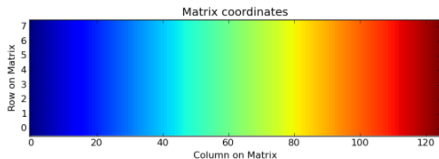| Carrier v3.3 | |
|----|----|
| 3 | KEK |
| 2 | DESY |
| 1 | IHEP (repair) |
| 6 | Giessen |
| 12 | (total) |

(status in VXD production database 12.10.2017)

- 33 AMC and 9 carrier to be sent to KEK for phase 3
  will first be sent to DESY for PXD commissioning (testpattern and cosmic), then sent from DESY to KEK
- repair: 4+2 AMC cards, problem with flash must be fixed, no automatic bitstream booting
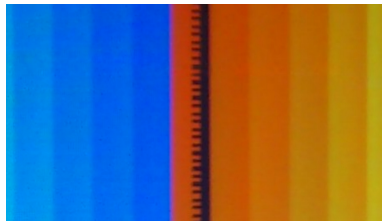- repair: 1 carrier board, 1 backplane channel not working

# ONSEN firmware: remapping

- introduced for PXD9 (1$^{st}$ time required in TB 04/2016)
- mirrored per 4 columns
- then mirrored per 64 columns
- 250 vs. 256 pixels
- different for PXD layer 1 and layer 2



Matrix coordinates



DHP coordinates
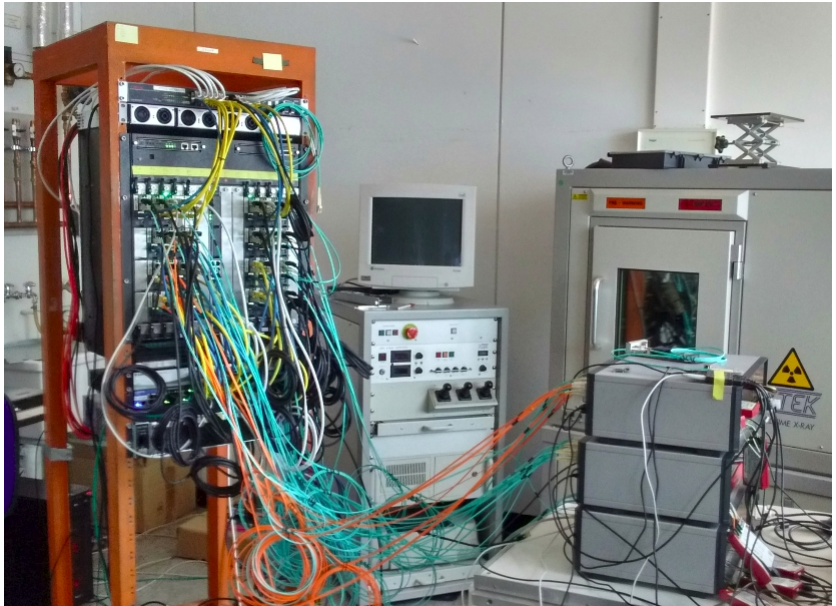
# ONSEN firmware: remapping

- implemented in basf2 unpacker (offline) in TB 04/2016
- implemented on Onsen (online) in TB 02/2017
  exact lookup tables on FPGA (no approximation)
  running stable in complete TB
- future: PXD online cluster finder will require remapping implemented on DHE (planned for phase 3)



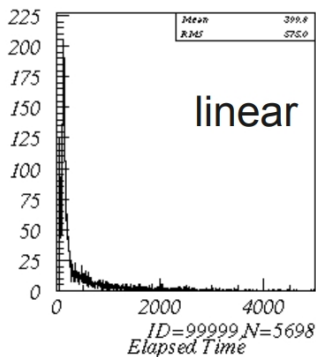There is one row alternating in DHP ID row-by-row

# Full system test, results

- 3 weeks testing (storing binary output data on SSD for crosscheck)
- 2 long runs over weekend
- Trigger rate $\leq 8$ kHz (limited by DHC aurora line rate)
  requirement 30 kHz / 4 links/DHC = 7.5 kHz
- Data rate $\sim 595$ MB/s
  540 MB/s is 3% occupancy
- Runs with HLT "send all" flag with reduced data rate of 600 Hz,
  send downscaled fraction of non-ROI processed (was problem in
  TB 2016)
  - No connection interrupts (backplane and external)
  - No buffer overflows (level $\leq 73\%$)
  - No framing errors, no data format errors
- Multiple start/stop without cold start
- Stable temperature in ATCA shelf ($\sim 60^{o}$ C at FPGA)

Simon Reiter
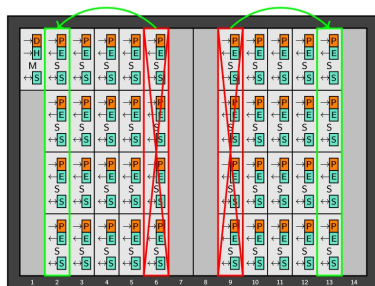
- "send ROIs" flag in HLT data (write also ROIs into the data stream for offline check) → no error

- HLT reject trigger → no error non-triggered data are removed in ONSEN, buffer is freed

- HLT trigger unordered → no error

- HLT with fixed latency ($\tau$=1 s) → no mismatch

- HLT latency according to Belle distribution, ~$10^9$ events (~8 hours, 30 kHz)
  → 7 mismatches
  → 111 "no DHC data" (but possibly HLT arrives <u>before</u> data)



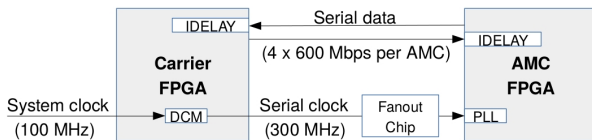$ID=99999, N=5698$
*Elapsed Time*

# Full system test, backplane link problem

- phase 3 requires scaling of
  ONSEN carrier boards
  from 2 to 9

- problem: with merger firmware
  sending to multiple boards, all
  backplane links become unstable

- → crosstalk found between
  Ethernet IO and one MGT power
  supply (on the carrier board
  FPGA, not the backplane)

- solved by avoiding that link
  → use different ATCA slots
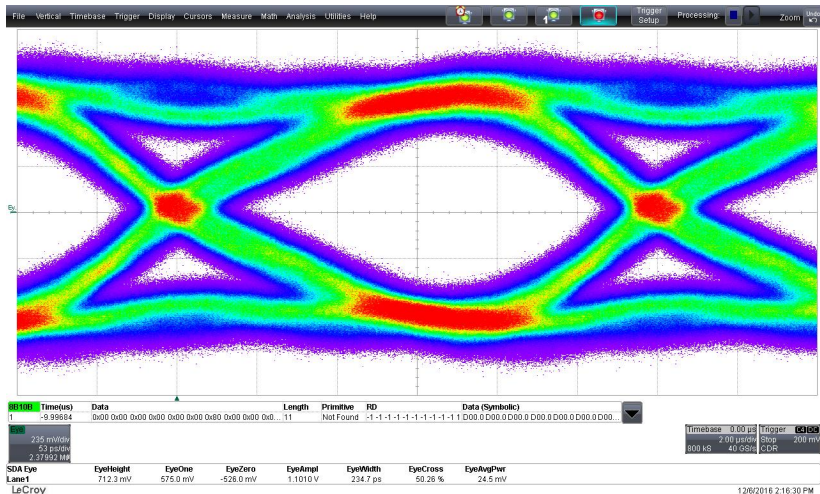  (different FPGA pins)
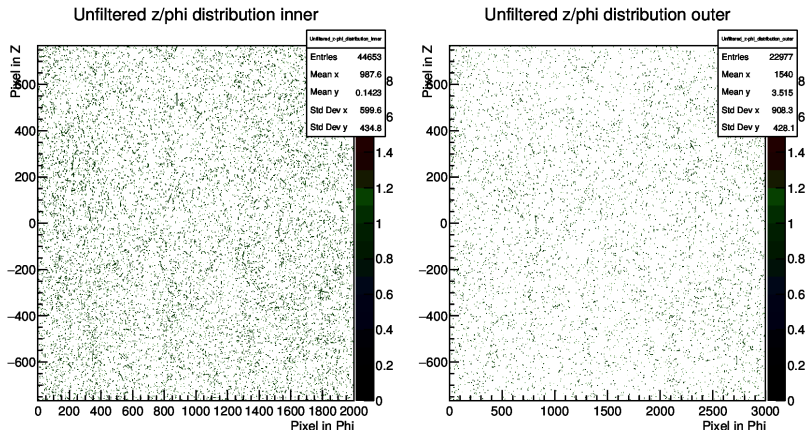
# Full system test, links Between carrier and AMCs

- Connection Carrier FPGA AMC FPGA uses serial (LVDS) links
- Serial clock is distributed from Carrier to AMCs
- Clock/data phase shift is compensated by delay, determined by tuning
- Problem: strong delay difference between Carrier/AMC combinations (due to routing)
- Problem: small temperature drift of the delay
- Solution: online self-calibration mechanism vary delay, check if link is up or not
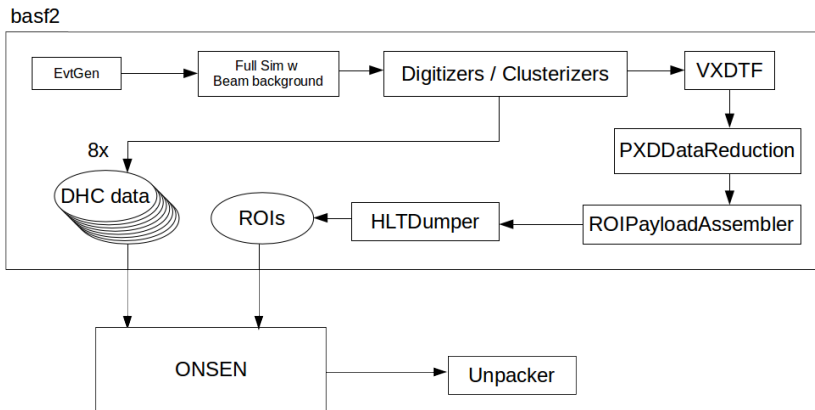
# ATCA backplane eye diagram
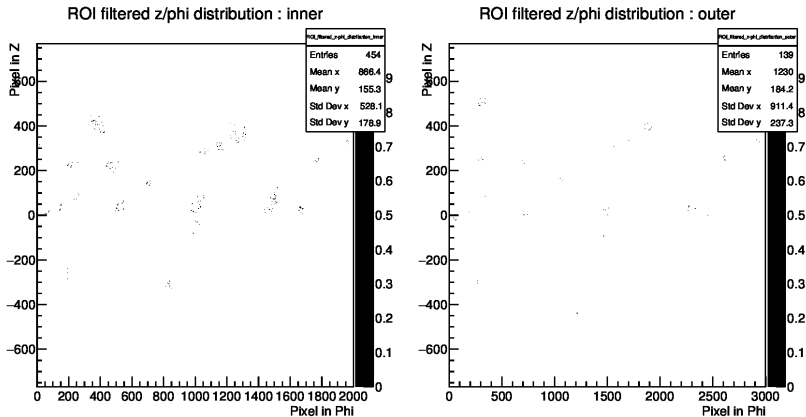
# Processing basf2 MC physics events in ONSEN



Unfiltered z/phi distribution inner

Unfiltered z/phi distribution outer

Average

occupancy 0.8% (forward), 0.4% (backward), incl. background
BonnDAQ UDP limit 128 MB/s corresponds to 0.71%
(30 kHz)

Klemens Lautenbach

# Processing basf2 MC physics events in ONSEN



Processing 5000 events (0.5 s of PXD data taking) and generate binary data required few days.
VXDTF1. Background MC8.

Klemens Lautenbach

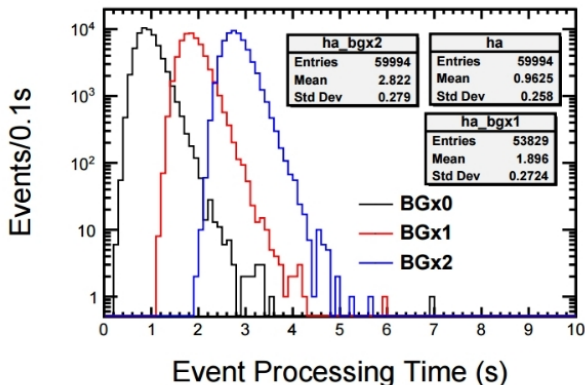# Processing basf2 MC physics events in ONSEN



factor 98.3 (inner), 121.6 (outer)
requirement $\geq 10.0 \rightarrow$ may be released

Klemens Lautenbach

- Line 363, 364, Section Event builder
  *"The ONSEN buffering capabilities should checked against the maximum estimated fluctuations."*
- HLT latency distribution from Belle ($\tau_{everage}$=1 s, $\tau_{max}$=5 sec) confirmed by Chunhua Li (Melbourne) with MC for Belle II (see next slide)
- Full system test at Giessen
  Worst case scenario: full data rate (3% occupancy), full trigger rate (30 kHz)
  $\rightarrow$ no buffer overflows (level $\leq$73%)

# Belle II, HLT latency study

New results on HLT processing time, by Chunhua Li (Melbourne)
60.000 hadronic events ($B\bar{B}$ + udsc)
For BGx2 only 4 events with t>5 s

Chunhua Li (Melbourne)

A minor concern is the connection between the ONSEN and the Event Builder 2 PCs. This was shown to be very sensitive to fluctuations in occupancy and is believed to be caused by the lack of Ethernet flow-control in the firmware used by the ONSEN system to implement this link connecting the ONSEN to the readout PCs. This firmware is called SiTCP, which started out as a KEK project, but is now closed source, and therefore cannot be easily modified

Under such conditions the only way to avoid packet-loss is the use of a very deep buffer and expensive network switches in the Event Builder 2. Even those improvements do not guarantee lossless transmission under all circumstances.

Concerning the ONSEN communication, there are other simplified TCP/IP implementations for FPGAs available, both commercial and non-commercial[1] and the team is encouraged to consider changing the TCP/IP driver in the ONSEN output FPGA to use one of those. This would greatly ease the burden on the network switches in the Event Builder 2, which also may reduce the cost[2].

- We contacted BeeBeans Technology, and very kindly received an SiTCP version (v11.0) which should recognize PAUSE frames
- This SiTCP version is installed in the present ONSEN firmware (e.g. for phase 2)
- Not tested yet, because test non-trivial
  - provoke network congestion
  - monitor, if PAUSE frames arrive
  - monitor, if SiTCP stops sending in such a case (monitor backpressure by SiTCP in chipscope ?)
  - compare old and new version of SiTCP
- Yamagata-san provided a test program to send PAUSE frames from a PC

# TB 02/2017, positive results

- final ONSEN hardware
- 2 ROI selectors parallel (2 DHCs connected)
- Onsen and DHH systems running stable for $\sim 10^9$ events
  per run up to 18 hours duration
  $\sim$1500 sroot files, 3.5 TB
  $2-3$kHz trigger rate (limited by DHC double trigger veto)
- online re-mapping (on Onsen) permanently switched on $\rightarrow$ basically
  permanently ROI selection in TB 04/2016 only 1 run ($\sim 10^5$ events)

# TB 02/2017, negative results

1. Onsen operation required cold restart for every run

   - re-upload FPGA bitstreams
   - otherwise trigger number mismatch
   - traced back to fragmented events from DHC, if ONSEN is reset, but DHC is not reset
     (DHC was not fully integrated in RC)
     not an ONSEN problem

2. Inconsistent states in PXD RC and global RC (READY or not-READY), in particular after Onsen cold restart

   - confusion for shift crew
   - traced back to 2 problems:
     2.1 software problem in global RC: updated state not interpreted in nsm-epics IOC
         not an ONSEN problem
     2.2 state of SiTCP connection between HLT or EB2 and ONSEN not clear ONSEN problem, but also HLT/EB2 problem

# TB 02/2017, negative results

- Solutions to problem of unknown SiTCP connection status
- FIN ACK sequence implemented and tested on ONSEN
  SiTCP terminates the TCP connection correctly, if
  - run is terminated (by run control)
  - Linux (on Onsen embedded PowerPC) is shutdown
- RBCP sideband protocol
  - enables channel status monitoring
  - implemented in SiTCP (according to documentation and
    specification), but not tested yet
  - monitoring must be done from the receiver side (HLT or
    EB2), as SiTCP connection is initiated from receiver
  - agreed with DAQ group, on TODO list

# ONSEN "sanitizer"

Protection of ONSEN against errors from other subsystems
Test system: copy of DESY setup with additional data fork
inducing errors (intentionally) from other systems

| |
|---|
| invalid CRC in HLT-, DATCON- oder DHH data |
| permanent "source ready" from DATCOM or DHC |
| $1^{st}$ DHC start-of-frame missing |
| $1^{st}$ DHC end-of-frame missing |
| every $8^{th}$ DHC start-of-frame missing |
| $1^{st}$ 4 bytes of DHC frame missing |
| DHC restart during a run (CTRL-c while sending DHC data by netcat and then restart netcat) |
| send 2 (or more) DHC start → event mismatch |
| send DHC end before DHC start |
| induced framing errors: interrupt locallink connection inside a frame (inside zero suppressed data), then send a new event to Onsen |
| different run number in HLT and DHC frames (but same trigger number, and only the 2 least significant bytes are checked) |
| $1^{st}$ DHE start double in DHC data |
| send 2x same HLT data to merger |
| send HLT header word 2x to merger (with and w/o DATCON) |
| missing HLT magic word (0xBE12DA7A + length of frame) |
| non-valid ROI from HLT |
| send DHH data with wrong DHE ID sequence |
| send many (60-70) ROI |
| send many (100 frames in 1 event) DHH data |

Dennis Getzkow

# ONSEN "sanitizer"

ONSEN firmware is now protected against 3 major external problems:

- invalid CRC in HLT frame
  $\rightarrow$ Onsen merger blocked any further incoming HLT data
- fragmented DHC data (cut in the middle of zero suppressed data block)
  $\rightarrow$ event fusion of 2 events (but no cold start required)
- double DHC start
  $\rightarrow$ event mismatch for all following events
  cold start required

Dennis Getzkow

# Test: adding a 33$^{th}$ ROI

Event 1, DHE 1.1.1 / DHE 2

Event 1, DHE 1.1.1 / DHE 2

# ONSEN emulator

- C++-Program by S. Reiter, PXD data reduction in software
- Loads test data from file (PXD/HLT/DATCON) (requires 0xBE12DA7A header)
- Similar memory management as ONSEN
- Processing time example
  1000 events, 4% PXD occupancy = 780 MB pixel data
  ONSEN: ¡ 2 seconds after sending HLT (1 Selector node)
  Emulator: (Intel i7 @ 3.4 GHz, 16 GB RAM):
  11 min, 50 s with 1 thread (factor $\leq$355)
  2 min, 40 s with 8 threads (factor $\leq$80)

Simon Reiter

# ONSEN firmware version number in bitstream

- uses 32 bits of commit-hash from the firmware git repository
- 2 files are generated: 1 bitstream, 1 linux kernel (contains epics PV definitions)
  1. hash is written into `USR_ACCESS` register ($\geq$Virtex-5). ONSEN carrier board: reading on Virtex-4 non-trivial, only by JTAG (Impact).
  2. hash is written in bitstream at a fixed adress at the end of the block-RAM. Can be read easily from PowerPC. Version is printed on console when booting and exported into epics PV. Can be logged into database: for every run it is fixed which firmware version.
  3. hash is written into version string of Linux kernel, when compiled. Kernel ELF file is also tagged with version (in addition to bitstream).
- similar mechanism for DHH:
  - store timestamp and board number in `USR_ACCESS`
  - write the same timestamp to a git tag to identify the commit

# Phase 3 preparations

- pedestal events (full frame events, in phase 2 recorded by BonnDAQ) requires FTSW-DHC communication (switch DHPT mode to memdump)

- load balancing, $5 \rightarrow 4$

  requires RTM in DHC ATCA system
  requires ROI distribution system on ONSEN



- hit-based format $\rightarrow$ cluster-based format

  - non-trivial data format change: start-of-cluster adress requires in remapped coordinates 10 bits, but only 8 bits reserved
  - new logic in ONSEN: hit inside-cluster but outside-ROI $\rightarrow$ new cluster buffer in ROI selection
  - requires cluster finder on DHE
  - remapping must be changed from ONSEN to DHE (cluster finder needs remapped coordinates)

# ROI distribution in phase 3



Uses additional "DHH ID filter" in front of ROI selector

(master thesis D. Getzkow)

# ROI distribution in phase 2

# ONSEN future development

- why?
    - almost no spares, but Virtex-4/5 at some point not available anymore
    - FPGA resources at limit
      e.g. presently no multiport memory controller for 2nd 2GB DDR2 RAM

- when? probably 2021 (planned PXD upgrade)

- new carrier board development for $\overline{P}$ANDA
  (IHEP Beijing and Univ. Giessen)
  remain compatible with existing AMC
  $\rightarrow$ Kintex Ultrascale, next slide

- upgrade link from DHC to ONSEN
  cluster-based format will increase required bandwidth by 30-50%
  (10 bit SOC adress)

# ONSEN future development

|  | Virtex-4 FX60 (CNCB) | Virtex-5 FX70T (xFP) | Kintex UltraScale 060 (Upgrade) |
|---|---|---|---|
| **Registers** | 50k | 44k | 663k |
| **LUTs** | 50k × 4–input | 44k × 6–input | 332k × 6–input |
| **DSP Slices** | 128 | 128 | 2760 |
| **BRAM** | 4 Mb | 5 Mb | 38 Mb |
| **MGT** | 16 × 6.5 Gbps | 16 × 6.5 Gbps | 32 × 16.3 Gbps |
| **CPU** | PPC405 | PPC440 | - |

# New physics rescue system

- CLUSTER RESCUE (high dE/dx $\rightarrow$ low $p_T$, no ROI)
  multilayer preceptron
  (input cluster size. cluster shape, seed charge, etc.)
  DHC or dedicated ONSEN carrier board

# Belle II Onsen Confluence (wiki system)



here: Onsen User Guide (not completely finished)
Onsen Ph. D. and master theses are on
`https://belle2.docs.org`, googleable

# Belle II Onsen Stash (git repository)



`https://stash.desy.de/projects/B2ON/repos/onsen/browse`
automatic bitstream build (Xilinx planAhead installed on DESY servers)
before phase 2: "release" (only event filter is missing)
"super onsen" git clone $\rightarrow$ checkout everything

# Belle II Onsen JIRA (issue management system)

# BACKUP

If new SiTCP version does not recognize PAUSE frames

- Problem is non-fatal
  - Communication is lossless, as siTCP includes retransmission
  - The problem is nonfatal: worst case in case of switch overload, reminder: there is 4 GB buffer on Onsen. If Onsen buffer full, back-pressure BUSY is issued (stop triggers), but there is no abort condition or data drop

- Other solutions?
  - CMS solution is only for sending, not receiving, but we need to receive HLT data
  - Advantage of siTCP: light weight, FPGA resources 15-20%, more complex protocol would require (non-available) resources
  - Long-term solution: use TCP on a PC with PCIe cards, input 32 optical links, output 10G uplink to event builder prototype existing and tested at Giessen (ALICE C-RORC)

# Test: HLT and DATCON ROis.



Event 16, DHE 1.1.2 / DHE 3