



Filip Leonarski :: Beamline Data Scientist :: MX Data Group

Jungfraujoch: a Data Acquisition and **On-the-fly Analysis System for HDR MX**

14th August 2021



- FPGA based detector data acquisition
- Compression for diffraction images





FPGA based detector data acquisition





2019	Dectris EIGER 2 XE	16 Mpixel	400 Hz	13.5 GB/s
2020	PSI JUNGFRAU	4 Mpixel	2200 Hz	18.4 GB/s
2022	PSI JUNGFRAU	10 Mpixel	2200 Hz	46.1 GB/s



Task: JUNGFRAU 10 Mpixel in 2022

- Beamline PXIII was selected as a pilot project for SLS 2.0 (Phase 0) – it will get major refurbishment in 2022
- One of developments is tiled JUNGFRAU 10 Mpixel detector for native-SAD applications – producing up to 46 GB/s
- Need smooth integration of the detector into beamline
- 4 Mpixel system planned for installation this year





Task: JUNGFRAU integration

- What is needed to operate JUNGFRAU detector?
 - UDP receiver
 - Conversion from raw to energy/photons
 - Peak finder for fast experimental feedback
 - Bitshuffle prefilter
 - Write to memory
- JUNGFRAU at a synchrotron must be internally operated at ≥ 1 kHz frame rates for continuous measurement at synchrotron
- Each of these task is difficult at 46 GB/s
- All tasks have low computational complexity and simple logic, but data throughput is high



VMXi DLS JUNGFRAU 1M tested at multilayer monochromator high-flux beamline



JUNGFRAU 4M tested at X06DA Swiss Light Source (CH) – special box to allow to cool detector to – 15°C



Conventional hardware results

- Raw to photon conversion example
- The most powerful single server available from PSI vendor
 - 4 CPU socket, 1.5 TB RAM, NVMe SSD drives and Mellanox fiber ethernet network cards
- Profiling and performance tuning
- Outcome: the server can handle JUNGFRAU 4Mpixel at 550 Hz, but not more
- Bottleneck: memory bandwidth of CPU is too small for conversion procedure



See: "JUNGFRAU detector for brighter xray sources: Solutions for IT and data science challenges in macromolecular crystallography" Leonarski et al. Structural Dynamics (2019) https://doi.org/10.1063/1.5143480



Field programmable gate arrays (FPGA)

- FPGAs
 - Real-time device:

design guarantees throughput and latency

– Fast memory bandwidth:

up to 460 GB/s per chip with HBM2

- Vast I/O options:

100G ethernet, PCIe, ...

- Also considered GPUs, but lacked necessary functionality
- Very powerful, but significant effort in development, due to need of specific hardware design skills and hardware description languages





Alpha-Data 9H3 board





POWER9OpenCAPIFPGACPUcableboard



What difference brings OpenCAPI?

- Similar to virtual mode starting from 80286/80386
 CPUs (as opposed to real mode), which hides complexity of system memory
- Each process has its own virtual address space
 - Simplicity: virtual address space is not dependent on system physical memory
 - Safety: there is no (easy) way to access memory of another process
- Translation virtual <-> physical memory is done by CPU electronics



Source: Wikipedia



What difference brings OpenCAPI?

- Standard HW interfaces, like PCIe, are not aware of virtual address space – developer need to address issue, increasing complexity
- OpenCAPI allows external accelerators (like FPGA) to attach directly to a virtual address space of the process
- Drastically reduced complexity of interconnect

250	G (x24)	DDR4 (2 chan	PHY nels)	DD (2 ch	R4 PHY annels)	. 6
PCle (x16)	CPU CPU CORECORE	CPU CPU CORECORE	Memory Controller	CPU CPU CORECORE	CPU CPU CORECORE	
PCle (x16)	L2s L3\$	L2s L3\$		L2\$ L3\$	L2\$ L3\$	ink
999 999 999	L3\$	L3\$ L2\$		L3\$	L3\$	SMP L
an an Meist	CPU CPU CORECORE	CPU CPU CORECORE		CPU CPU CORECORE	CPU CPU CORECORE	4B
PC	le Ctrl		nterconne SMP Ctrl	ct A	celerators	nk
	L3\$	L3\$		L3\$	L3\$	MP Li
PCle (x16)	CPU CPU CORECORE	CPU CPU CORECORE	Memory Controller	CPU CPU CORECORE	CPU CPU CORECORE	4B S
250	G (x24)	DDR4 (2 chan	PHY nels)	DDI (2 ch	R4 PHY annels)	

OpenCAPI 25G

WikiChip

POWER9



FPGA as software development project

- Aim: (relatively) fast transformation of SW based algorithms to FPGA
- High-level synthesis -> compile C++ functions into hardware description language
 - Not as efficient as pure HDL
 - Easier to write
 - Easier to test and plug into CI pipeline

Memory coherent interconnect

- (e.g. OpenCAPI, Intel CXL)
- Virtual, not physical, addressing
- CPU cache coherent
- No need for OS/kernel expertise
- Requires special hardware (at the moment)

All 589	Finished	Branches T	ags			
Filter	pipelines					c
itatus	Pipeline	Triggerer	Commit	Stages		
⊘ passe	#1981 latest	(₽blla-march2… -> 635804ad ∰ Added tools to calc.		⊘ 00:15:28 ∰ 4 days ago	Þ
⊘ passe	ed #1980	\$	P' bl1a-march2… ↔ c38c6b4a ∰ Minor modifications	 	⊘ 00:15:28 ∰ 5 days ago	Þ
⊘ passe	#1979 latest	()	Pfpga_refact ↔ 347224a2 ∰ Minor modifications		⊘ 00:15:30 5 days ago	Þ
⊘ passe	ed #1978		どfpga_refact… ↔ 7244d4a1	$\odot \odot \odot$	♂ 09:10:54	▶ - 4

Both C++ verification and HDL simulation are part of CI:

- C++ 8 minutes to cover 13 scenarios and 95% of code
- HDL 4 hours to cover single simple

scanario



Up to 50 GB/s acquisition and data analysis in a single 2U IBM POWER9 server with 1-4 FPGA boards





FPGA board with OpenCAPI interface



- Data acquisition
- Initial data analysis
 - **Pre-compression**
- (2.5 Mpixel/board for JF)



Jungfraujoch & Gold Standard



IUCrJ ISSN: 2052-2525 BIOLOGY | MEDICINE Volume 7 | Part 5 | September 2020 | Pages 784-792 https://doi.org/10.1107/S2052252520008672 OPEN O ACCESS Cited by 1

Gold Standard for macromolecular crystallography diffraction data

Herbert J. Bernstein,^{a*} Andreas Förster,^b Asmit Bhowmick,^c Aaron S. Brewster,^c Sandor Brockhauser,^{d,e,f} Luca Gelisio,^g David R. Hall,^h Filip Leonarski,ⁱ Valerio Mariani,^g Gianluca Santoni,^j Clemens Vonrhein^k and Graeme Winter^h











Commissioning in KEK (Jan – May 2021)

- Detector and data acquisition system was sent in November for an experiment in Photon Factory, KEK
- More than 2,000 datasets collected for protein targets, few real-life native-SAD structures solved
- Due to pandemic, detector support and development (including deployment of new FPGA design) was done fully remotely from Switzerland







BL-1A Photon Factory JUNGFRAU detector (up) tested in helium chamber for native-SAD measurements with 3.75 keV X-rays



Commissioning in KEK (Jan – May 2021)

 Detector and data acquisition system was sent in November for an experiment in Photon Factory, KEK



Results will be presented in another presentation:

The crystallomics pipeline, a shotgun approach on native proteomes to (re)discover the unsuspected

by Sylvain Engilberge 18.08 11:45





measurements with 3.75 keV X-rays



Possible gain from using FPGA based system



From a "state-of-the-art" conventional CPU server solution to a "FPGA boards + OpenCAPI" cutting-edge solution



18.4 GBps (4Mpixels@2.2kHz) Data acquisition + image conversion max bandwidth with 1x POWER9 IC922 2U server + 2 FPGAs solution Each FPGA acquire, convert on-the-fly and store in CPU memory 9.5GBps of images.

Comments:

PSI's published numbers reference can be found at https://doi.org/10.1063/1.5143480

4Mpixels images@2.2kHz acquisition + conversion was tested with 2 FPGAs in 2020. 10Mpixels@2.2kHz with 4 FPGAs will be tested in 2021 with high confidence. Moving the "spot finding" from post processing to the FPGA board increases the above ratios by 2 by removing actual post processing servers. Conventional CPU server solution may evolve by :

- adding network cards in parallel may add uncertainty (extra load for CPU) → performance reduction.
- New CPUs with L3 cache → performance increase.





Acquisition + conversion bandwidth increase by acquiring 4x faster (from 0.5 to 2.2kHz) 4Mpixels images

Price decrease using just 1 server + 2 FPGAs to acquire and convert 4Mpixels@2.2kHz while conventional solution would require 4 servers in parallel.

8 → with "spot finding" coded in FPGA, post processing servers can even be removed

Power consumption decrease by using just 1 server with 2 FPGAs (<500W total) rather than 4kW for 4 servers

★ 16 → with "spot finding" coded in FPGA, post processing servers can even be removed

Courtesy: B. Mesnet (IBM)



Possible gain from using FPGA based system



From a "state-of-the-art" conventional CPU server solution to a "FPGA boards + OpenCAPI" cutting-edge solution



Acquisition + conversion bandwidth increase by acquiring 4x faster (from 0.5 to 2.2kHz)

images with 2.5x more pixels (from 4 to 10

Price decrease using just 1 server + 4 FPGAs

to acquire and convert 10Mpixels@2.2kHz while conventional solution would require

★20 → with "spot finding" coded in FPGA, post

1 server with 4 FPGAs (<500W total) rather

processing servers can even be removed

Power consumption decrease by using just

processing servers can even be removed

4.5 GBps (4Mpixels@550Hz) Data acquisition + image conversion max bandwidth with a standard CPU solution (4 sockets - 1.5TB RAM) 4 Mpixels images are currently acquired at 1.1kHz rate but the too small memory bandwidth of CPU limits the conversion procedure → 4Mpixels@550Hz

46.1 GBps (10Mpixels@2.2kHz) Data acquisition + image conversion max bandwidth with 1x POWER9 IC922 2U server + 4 FPGAs solution Each FPGA acquire, convert on-the-fly and store in CPU memory 11.5GBps of images.

Comments:

PSI's published numbers reference can be found at https://doi.org/10.1063/1.5143480

4Mpixels images@2.2kHz acquisition + conversion was tested with 2 FPGAs in 2020. 10Mpixels@2.2kHz with 4 FPGAs will be tested in 2021 with high confidence. Moving the "spot finding" from post processing to the FPGA board increases the above ratios by 2 by removing actual post processing servers. Conventional CPU server solution may evolve by :

- adding network cards in parallel may add uncertainty (extra load for CPU) → performance reduction.
- New CPUs with L3 cache → performance increase.

10 servers in parallel.

than 10kW for 10 servers



Mpixels)





Compression of diffraction images



Efficient compression should be complementary to HW development

- FPGA processing can help in early stages of the pipeline, however transfer and storage remain an issue efficient compression can be a good answer
- Bitshuffle/LZ4 introduced with EIGER images is currently dominant
- There is lot of happening in the community:
 - Argonne: SZ lossy compression project for ExaFEL https://arxiv.org/pdf/2105.11730.pdf
 - Argonne: On-chip compression of X-ray data following counting statistics.
 <u>https://doi.org/10.1107/S1600577520013326</u>
 - RIKEN and FSU: TEZIP based on recurent neural networks <u>https://ieeexplore.ieee.org/document/9499386</u>
- Is HDR MX community ready to embrace these developments?
- Is there a way to tweak bitshuffle/LZ4?



How does current compression work?





Byte compression of pixels in LZ4

1	0	0	2	0	0	1	0
---	---	---	---	---	---	---	---

- Algorithm like LZ4 try to find repeatable sequences of bytes, which are replaced by some special tokens
- They focus on bytes, as texts are encoded as bytes
- The longer these sequences are the more efficient is compression
- However, it is hardly expected to find repeatable byte sequences in diffraction images



Bitshuffle filter – input sequence



Bitshuffle filter – input sequence in binary



Bitshuffle filter – shuffling operation







Bitshuffle filter - outcome

Long sequence of zero – low entropy: trivial to compress High entropy: unlikely to compress this part



Bitshuffle filter - outcome

Long sequence of zero – low entropy: trivial to compress

High entropy: unlikely to compress this part

MX diffraction images (esp. collected at high frame rates) contain mostly low counts (background) and small number of high counts (Bragg spots)



Bitshuffle implementation

- On CPU, bitshuffle is implemented as combination of some bit instructions there is no perfect CPU instruction to do it
- Bitshuffle step is generally fast (GB/s), but it will be a bottleneck, if compression is also fast
- Doing bitshuffle on FPGA is trivial just connecting wires properly between source and destination



- LZ4 is made to be extremely fast to decompress, relatively simple in construction
- Y. Collet, author of LZ4, made new algorithm called Zstandard (at Facebook)
- github.com/facebook/zstd
- Zstandard is similar to LZ4, but has more options within the algorithm
- By setting Zstdandard compression level, one can set tradeoff between performance and compression ratio
 Compression Speed vs Ratio (multi-threaded)





Zstandard for diffraction image

Lysozyme dataset collected with JF4M

	Compression ratio	Throughput
Bshuf+LZ4	7.0x	2.5 GB/s
Bshuf+Zstandard (default)	8.0x	0.6 GB/s
Bshuf+Zstandard (+10)	8.3x	0.2 GB/s
Bshuf+Zstandard (-10)	6.8x	1.9 GB/s

- Zstandard provides flexibility, not available with LZ4
- In general Zstandard compression is slower than LZ4



Zstandard for diffraction image

Lysozyme dataset collected with JF4M

	Compression ratio	Throughput
Bshuf+LZ4	7.0x	2.5 GB/s
Bshuf+Zstandard (default)	8.0x	0.6 GB/s
Bshuf+Zstandard (+10)	8.3x	0.2 GB/s
Bshuf+Zstandard (-10)	6.8x	1.9 GB/s

- Zstandard provides flexibility, not available with LZ4
- In general Zstandard compression is slower than LZ4
- Zstandard can benefit from large compression blocks. With 1 MB block compression factor of 9.0x could be reached with 1.1 GB/s throughput.
 LZ4 doesn't gain from larger block size.



 Do bitshuffle in FPGA and transfer bitshuffled images to host memory (currently 8 kB block)

	Compression ratio	Throughput
LZ4 (bshuf on CPU)	7.0x	2.5 GB/s
LZ4 (bshuf on FPGA)	7.0x	4.0 GB/s



Speed-up in Jungfraujoch

- Do bitshuffle in FPGA and transfer bitshuffled images to host memory (currently 8 kB block)
- Use custom Zstandard compressor only compress long sequences of zeroes (at least 5 bytes), no compression for other byte values; standard decompressor

	Compression ratio	Throughput
LZ4 (bshuf on CPU)	7.0x	2.5 GB/s
Zstandard (bshuf on CPU)	8.0x	0.6 GB/s
LZ4 (bshuf on FPGA)	7.0x	4.0 GB/s
Zstandard/RLE (bshuf on FPGA)	6.4x	7.4 GB/s



Speed-up in Jungfraujoch

- Do bitshuffle in FPGA and transfer bitshuffled images to host memory (currently 8 kB block)
- Use custom Zstandard compressor only compress long sequences of zeros (and ones), don't do any other compression; standard decompressor

Just compressing zeros of bitshuffled data accounts from most of compression and can be done extremly efficient (incl. FPGA implementation).

To get beyond this, one needs significant effort.

ompression ratio	Throughput
7.0x	2.5 GB/s
8.0x	0.6 GB/s
7.0x	4.0 GB/s
6.4x	7.4 GB/s



There is one issue - multipixels

- EIGER and JUNGFRAU modules are both 1024 col x 512 lines
- However pixels on ASIC boundaries are larger, leading to double sized pixel
- In post-processing these pixels are replaced with 2 (edges) or 4 pixels (corners)
- Resulting 1030x514 pixel image better represents geometry, but has worse properties for splitting into packets and blocks (not power-of-two)
- It is relatively complex to add these multipixels in electronics (FPGA) + they add small to image size without extra information



There is one issue - multipixels

• EIGEF 1024			
 Howe to do In po (edge Resul geom packe 	Requirement to have multipixels precludes bitshuffling on FPGA at the moment. The solution is simple – add these at decompression with special HDF5 plugin (in pipeline for development). Is there a more elegant solution?		
 It is re electe withe 			



There is one issue - multipixels

• EIGEF 1024			
• Howe			
to do	Requirement to have mul	tipixels precludes	
• In po	Ditshuming off FPGA a	t the moment.	
(edge	The solution is simple		
 Resul geom 	decompression with speci pipeline for deve		
packe	Is there a more elega	ant solution?	
 It is relected with 		How about placeho These are even big less But here GOLD ST	olders for module gaps? gger problem and even s useful. ANDARD is an answer!



Jungfraujoch is a fully integrated DAQ solution for kilohertz frame rate JUNGFRAU detectors – deployed on X06DA beamline this summer for user operation

More focus is necessary for compression development – flexibility of Zstandard can come handy for future data rates





Acknowledgements

MX Group (PSI)

- Vincent Olieric
- Takashi Tomizaki
- Chia-Ying Huang
- Sylvain Engilberg
- Justyna Wojdyła
- Meitian Wang

Detector Group (PSI)

- Aldo Mozzanica
- Martin Brückner
- Carlos Lopez-Cuenca
- Bernd Schmitt

Science IT (PSI)

• Leonardo Sala

Controls (PSI)

- Andrej Babic
- Leonardo Hax-Damiani

SLS management (PSI)

• Oliver Bunk

Photon Factory, KEK

- Naohiro Matsugaki
- Yusuke Yamada
- Masahide Hikita

MAX IV

- Jie Nan
- Zdenek Matej

Uni Konstanz

• Kay Diederichs

LBL

Aaron Brewster

DLS

• Graeme Winter

ESRF

• Jerome Kieffer

Dectris

- Stefan Brandstetter
- Andres Förster

IBM Systems (France)

- Alexandre Castellane
- Bruno Mesnet

InnoBoost SA

• Lionel Clavien