Developments in EdgeML Inference in FPGAs, and their Applications in PET

Ryan Herbst, Ryan Coffee, Larry Ruckman SLAC National Accelerator Laboratory





SNL (SLAC Neural network Library) Framework Goals

- Provide a **set of libraries** to synthesize AI inference networks into FPGAs
- Networks of a medium size,
 - **10 20 layers**, ~100s of thousand trainable parameters
 - Total end to end latency of **1uS** to 10ms
 - Deep networks call for AI-specific ASICs, e.g. IPU or Groq, but **FPGAs still ubiquitous** at sensor

- Pipelined implementations targeting a frame rate of **100kHz 1MHz**
- **Dynamic reloading of weights** and biases that avoids re-synthesis
- Support a Keras-like API for layer definition and configuration
- Allow for a standard interface such as **HDF5** for loading weights and biases
- Allow for a modular approach with the ability to implement new and **custom layers**



3

Why Dynamic Weights & Biases?

- Our SNL implementation is targeting scientific instruments which will continuously adapt to new data and changing environments
 - High speed training to supports this goal
 - Bias and weight updates in real time
- Some AI-to-FPGA frameworks take the weights and biases as an input, pruning portions of the network structure to save resources
 - Re-synthesis is required for each new training set
 - **Risk of the FPGA implementation failing** due to increase resources usage, timing failures or massive change in internal interconnect structure
- Large FPGA designs can take hours to days for the HLS -> synthesis -> place and route cycle to complete
- This is a trade off between **robustness vs. latency** and resource usage





Scheme to Measure the X-ray Timing in LCLS

- Concept borrowed from DOE-funded LCLS (Linac Coherent Light Source) X-ray timing project ("TimeTool")
- Use frequency chirped broadband laser pulse to create a spectrum-to-time mapping
- Result is the relative timing between FEL (Free-electron Laser) X-ray pulse and the experimental optical laser



Bionta *et al.,* Opt Express. (2011) **19** (22) 21855-65. doi: 10.1364/OE.19.021855. PMID: 22109037. Bionta *et al.,* Rev Sci Instrum. (2014) **85** (8) 083116. doi: 10.1063/1.4893657. PMID: 25173255.

SLAC



Coffee *et al.*, Philo. Trans. A Math Phys. Eng. Sci. (2019) **377** (2145):20180386. doi: 10.1098/rsta.2018.0386.

Translating Method to Medical Imaging

Spectrometer

Polarizer

- The annihilation photons in PET creates a similar cascade in bulk materials as x-ray photons from LCLS, but much higher energy/photon
- Measure (with the chirped laser) the induced refractive index change directly (-45deg) rather than wait for scintillation process.
- Birefringence measured with interference for sensitivity to single photon events

BBO

crystal

YAG crystal

Mirror

 511 keV photons produce ps-scale cascades, sets optimal chirp for resolution below 0.2ps

Mirror

Focused X-ray

propagation

BBO

crystal

Polarizer

Probe

supercontinuur



Li Tao et al (2021) Phys. Med. Biol. 66 045032, DOI 10.1088/1361-6560/abd951

(b)

Birefringent

delav

Time

(a)

Linear

polarization

Time

Running the PET algorithm with LCLS Light Source

- Edge identification via conventional edge-finding or CNN algorithms
- Results for PET will be rare events (kHz MHz) individual rows in image that show rise and fall edges
- Rare event edge locations are the only relevant information
- Can we produce that information at the camera, then feature extract the framegrabber before software receives it ("lossy" compression)



Li Tao et al 2021 Phys. Med. Biol. 66 045032, DOI 10.1088/1361-6560/abd951

CoaXPress over Fiber (CXPoF) Framegrabber

- Array of optical cameras for parallel readout
- Each camera hosts ~1000 of detection elements (as signal rows)
- Streaming images are evaluated at CXPoF receiving FPGA for "live" rows
- Results from "live" rows are then used for direct PET reconstruction



Up to 937 FPS at 4096 x 2304 (8 bit) Up to 52,080 FPS at reduced resolutions

BittWare XUP-VV8

32 TX/RX Optical channels (Front Panel) On-board 64 GB DDR4 Memory PCIe GEN3 x 16 Lane

BittWare Card could support up to 32 units of S991 camera at 1 lane per Camera (1/8 rate per camera) SLAC Up to 6,510 FPS at reduced resolutions (=52,080/8) SURF's QSFP-DD CXPoF (1 of 8 lanes) PCIe GEN3 x Firmware 16 lanes AXI Up to 320 Gb/s Image **PCle** Stream Inbound Processing DMA up to 113 Gb/s (edgeML) SURF's **QSFP-DD** CXPoF (1 of 8 lanes)

Our CXPoF firmware is open source (refer to backup slide about SURF)

Our DMA firmware and kernel driver is open source (refer to backup slide about "PCIe Framework") 8

Acknowledgement

- National Institute of Health
 - Time-of-Flight PET: (PI Levin) Proj. # 5R01EB023903-04
- Department of Energy, Office of Basic Energy Science
 - CookieBox detector: FWP 100498 (PI Coffee) "Enabling long wavelength Streaking for Attosecond X-ray Science"
 - EdgeAI for CookieBox: FWP 100643 (PI Thayer) "Actionable Information from Sensor to Data Center"
 - DOE High Energy Physics Detector R&D
- Department of Energy, Office of Fusion Energy Science
 - EdgeAI for Fusion Control: FWP 100636 (PIs Koleman, Coffee, Smith, Boyer, Schneider) "Machine Learning for Real-time Fusion Plasma Behavior Prediction and Manipulation."
- SLAC-LCLS Program Development (PI Kling) "LCLS Growth"
- SLAC TID LDRD 21-007 (PI Herbst) "Edge ML for acquisition and analysis of data generated by ultra high rate detectors"

SLAC

BACKUP SLIDES

BittWare Solution: COTS FPGA PCIe DAQ Card: XUP-VV8



BittWare Card could support up to four S991 camera at 8 lane per Camera at FULL rate SLAC



Demonstration of custom CXPoF Framegrabber + S991 camera

Demonstrated 120 Hz LCLS timing fiber triggering with the S991 camera

• 4096 pixels wide x 2304 pixels high in 8-bit pixel mode (9,437,184 Byte camera image)

- Tag Release (public domain):
 - <u>https://github.com/slaclab/lcls2-coaxpress-over-fiber-apps/releases/tag/v1.0.0</u>

Rogue PyQT GUI

em	Debug Tree						
е		Mode	Туре	Value		Command	
						Read	
	enable	RW	bool	True	*		
	 CoaXPressAxiL 					Read	
	enable	RW	bool	True	*		
	RxLinkUpCnt[0]	RO	UInt12	0x0			
	RxLinkUpCnt[1]	RO	UInt12	0x0			
	RxLinkUpCnt[2]	RO	UInt12	0x0			
	RxLinkUpCnt[3]	RO	UInt12	0x0			
	TrigRate	RO	UInt32	119	Hz		
	RxLinkUp	RO	UInt4	0xf			
	TxLinkUp	RO	UInt1	0×1			
	TxLinkUpCnt	RO	UInt12	0x0			
	TxTrigCnt	RO	UInt12	0xfff			
	TxTrigDropCnt	RO	UInt12	0x0			
	RxOverflowCnt	RO	UInt12	0x0			
	RxFsmErrorCnt	RO	UInt12	0x0			
	TrigPulseWidth	RW	float	100.000	μs		
	RxNumberOfLane	RW	UInt4	4			
	TxTrigInv	RW	UInt1	0x0			
	SoftwareTrig	wo	UInt1			Exec	
	CountReset	wo	UInt1		1	Exec	
	▼ DataSteamMon				1	Read	
	enable	RW	bool	True	*		
	FrameCnt	RO	UInt64	0x1dd4	1		
	FrameRate	RO	Int32	119	Hz		
	FrameRateMax	RO	Int32	120	Hz		
	FrameRateMin	RO	Int32	119	Hz		
	Bandwidth	RO	float	8984.2 M	lbps		
	BandwidthMax	RO	float	8998.6 M	1bps		
	BandwidthMin	RO	float	8984.2 M	lbps		
	FrameSize	RO	Int32	9437184 E	Byte		
	FrameSizeMax	RO	Int32	9437184 E	Byte		
	FrameSizeMin	RO	Int32	9437184 E	Byte		
	▶ TimingRx					Read	

SNL Usage

- User defines layers using a **collection of C++ templates** that define each layer type and the associated activator for each layer
- Current layer types:
 - o Conv2D
 - MaxPooling
 - AveragePooling
 - Dense
 - Reservoir
- Current activators:
 - LeakyRelu
 - o Relu
- Each layer interface matches closely to the Keras model
 - Allows user to use Keras documentation as a reference
- Additional configurations are required to control the **FPGA specific configurations**





Why Streaming & pipelining?

- Certain layer types work better with streaming based upon their **data flow** model
 - Convolution (2D & 3D) and pooling layers can process **data as it arrives**
 - **Partial data** is ready as each region of the layer completes
 - Pass to next layer can begin when indices reach roughly ½ the kernel dimensions
- Other layers such as **Dense fully connected** need all of the input data to have arrived before output can begin... this suffers a **high latency penalty**



SLAC Ultimate RTL Framework (SURF) Firmware: Open source

(IPv4, ARP, DHCP, ICMP, UDP)

(ADI, Micron, SiliconLabs, TI, etc.)

(clock managers, SEM, DNA, IPROG)

(DMA, MUX, FIFO, etc)

(Crossbar, AXI4-to-AXI4-Lite bridge, etc.)

(Synchronize bits, buses, vectors, resets, etc)

(I2C, SPI, UART, line-code, JESD204B, etc)

https://github.com/slaclab/surf/tree/master/protocols/coaxpress

- <u>https://github.com/slaclab/surf</u>
- HUGE VHDL library for FPGA development
- Used in Xilinx FPGAs, Intel FPGAs and ASIC digital designs
- VHDL-based IPs for commonly implemented modules
 - Ethernet Library:
 - AXI4 Library: (Crossbar, DMA, FIFO, etc.)
 - AXI4-Lite Library:
 - AXI4 stream Library:
 - Device Library:
 - Synchronization Library:
 - Wrapped Xilinx Library:
 - Serial Protocols Library:
 - CXPoF Protocol:
- SURF is managed and maintained by TID-ID Electronics Systems Department
- New features and bug fixes on a weekly basis



-SLAC

Rogue Software (Open source)

- <u>https://github.com/slaclab/rogue</u>, <u>https://slaclab.github.io/rogue/</u>
- Software tools for both rapid prototyping and experiment deployment
- "Rogue" is not an acronym
- Operates either in a python/C++ hybrid or C++ only mode
 - Higher level python for connecting the high performance C++ modules togeth
- x86-64, ARM32 and ARM64 support
- Able to run on Linux, MAC or windows
 - Windows requires either WSL2 or Linux virtual machine
- Lots of ways to run the rogue software
 - 100% C++ only code
 - python script (non-GUI)
 - EPICS
 - Python GUI (e.g. PyQT, PyDM, etc)
- Rogue is managed and maintained by TID-ID Electronics Systems Department
- New features and bug fixes on a weekly basis



PCIe FW/SW Framework (Open source)

- <u>https://github.com/slaclab/axi-pcie-core</u> (Firmware)
- <u>https://github.com/slaclab/aes-stream-drivers</u> (Linux Kernel Driver)
- Firmware framework for BAR0 AXI-Lite interface with up to 8 DMA lanes
 - 256 TDEST per DMA lane (up to 2048 destinations total)
- Provide a "common platform" firmware frame and software kernel driver for any Xilinx PCIe card
 - Common Xilinx Dev board support
 - AC701, KC705, KCU105, KCU116, VCU128
 - Common Xilinx Data Center Card support
 - U50, U55C, U200, U250, U280, C1100
 - Much more card support can be added as needed
- Demonstrate up to 113 Gb/s for large frames (PCIe GEN3 x 16, 1MB frames)
- Demonstrate > 1MHz frame rate for small frame (<128B) without frame batching
- Managed and maintained by TID-ID Electronics Systems Department

