

# Ultra Low Latency AI support to drive detectors

**Peter Rakyta, Gregory Morse**

Eötvös Loránd University

HUN-REN Wigner Research

Centre for Physics

**Benjamin Mencer, Ryan Coffee**

SLAC National Accelerator Laboratory

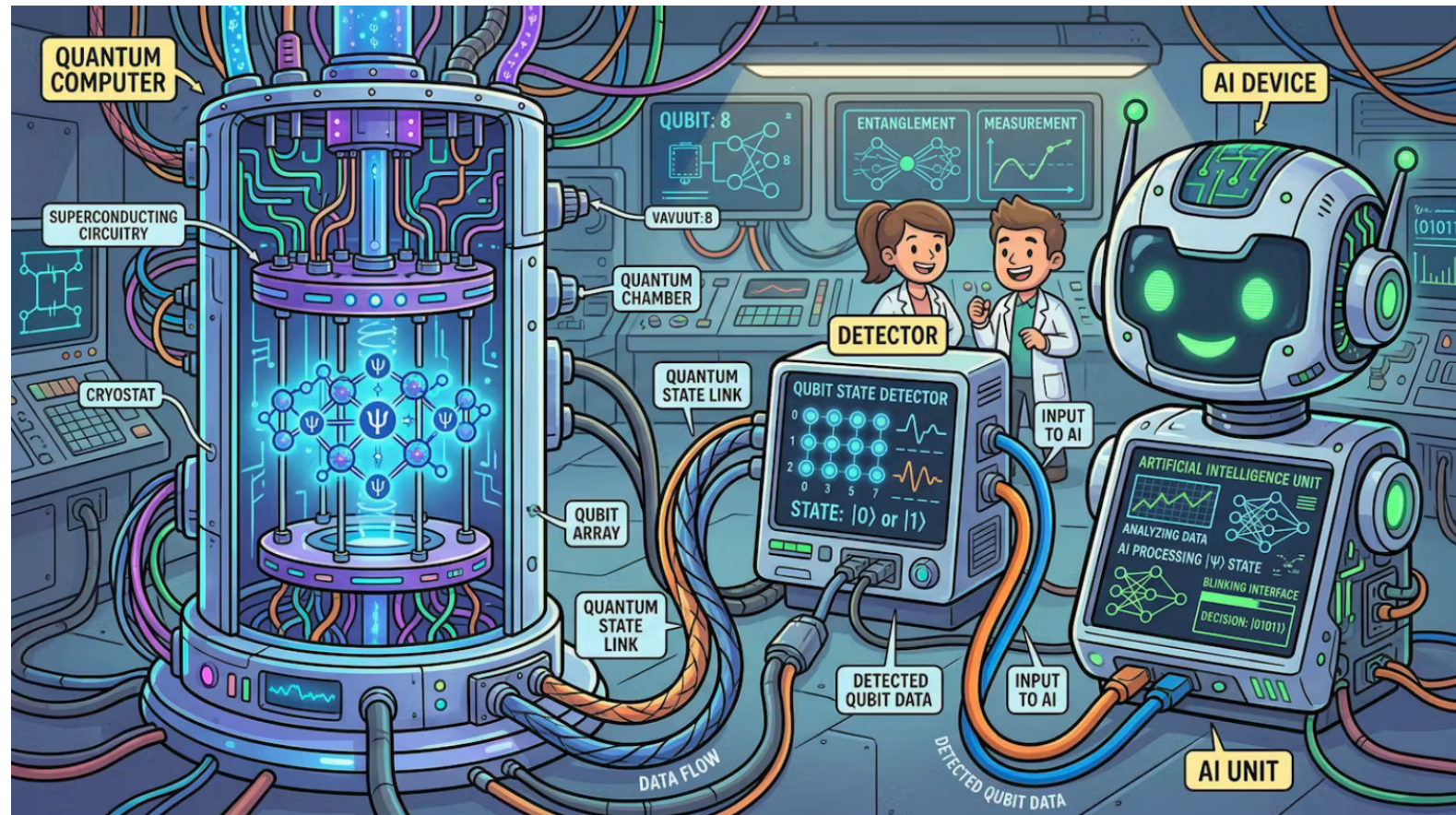


**ELTE**  
EÖTVÖS LORÁND  
UNIVERSITY



# Ultra Low Latency AI support to drive (QC) detectors

- Error correction
- Error mitigation
- Pulse generation
  
- Requires AI feedback on a  $\mu\text{s}$  time scale



ELTE  
EÖTVÖS LORÁND  
UNIVERSITY



QNL Quantum Information  
National Laboratory  
HUNGARY

# The Genesis Mission

A national initiative to build the world's most powerful scientific platform.



**The Genesis Mission unites DOE National Labs, industry, academia, and more to harness AI for breakthroughs in energy dominance, discovery science, and national security.**



## The Genesis Mission: Transforming Science and Energy with AI

**Agency:** Office of Science

[Assistance Listings](#): 81.049--Office of Science Financial Assistance Program

Last Updated: April 10, 2026

[View version history on Grants.gov](#) [🔗](#)



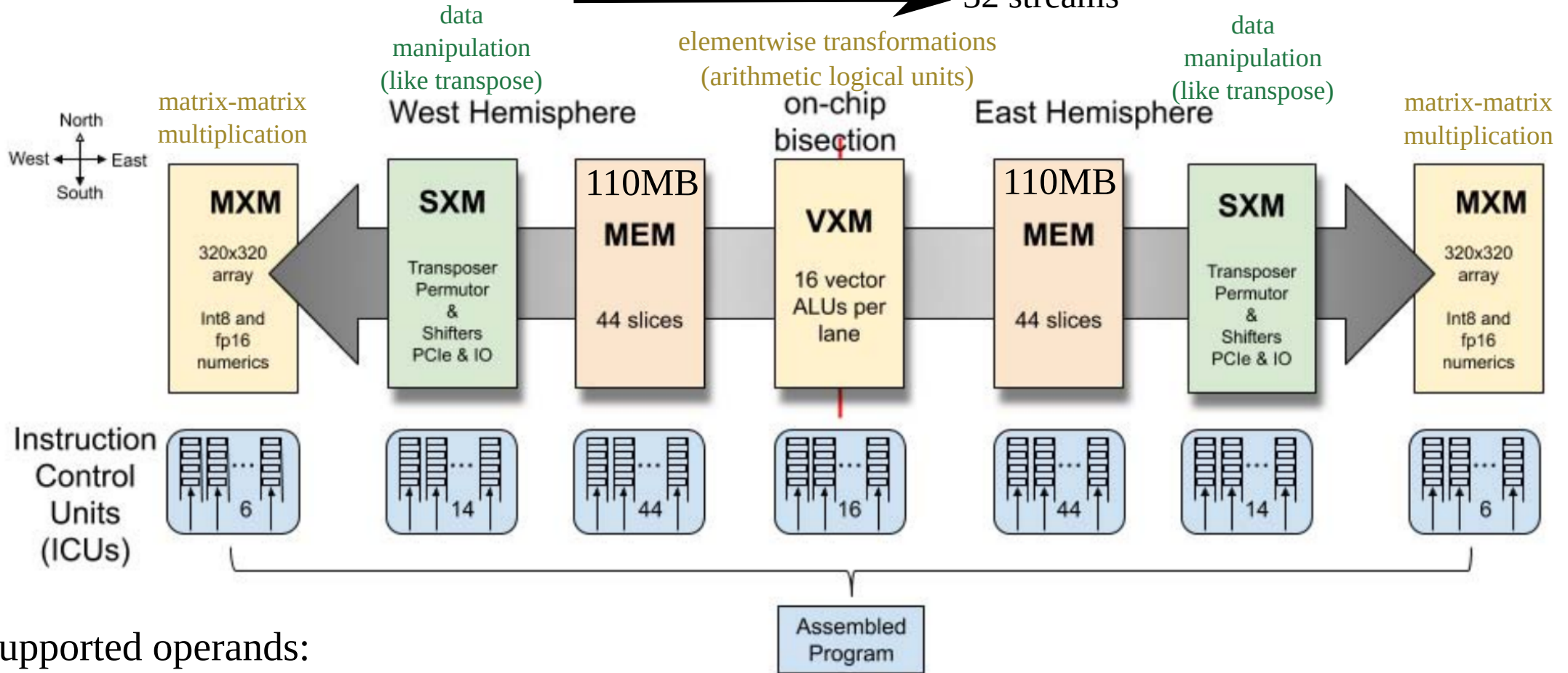
U.S. DEPARTMENT  
*of* **ENERGY**

# Think Fast: A Tensor Streaming Processor (TSP) for Accelerating Deep Learning Workloads

2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)

each stream carries **320 vectorized bytes**  
over 320 lanes

32 streams ← streams of data across the chip → 32 streams

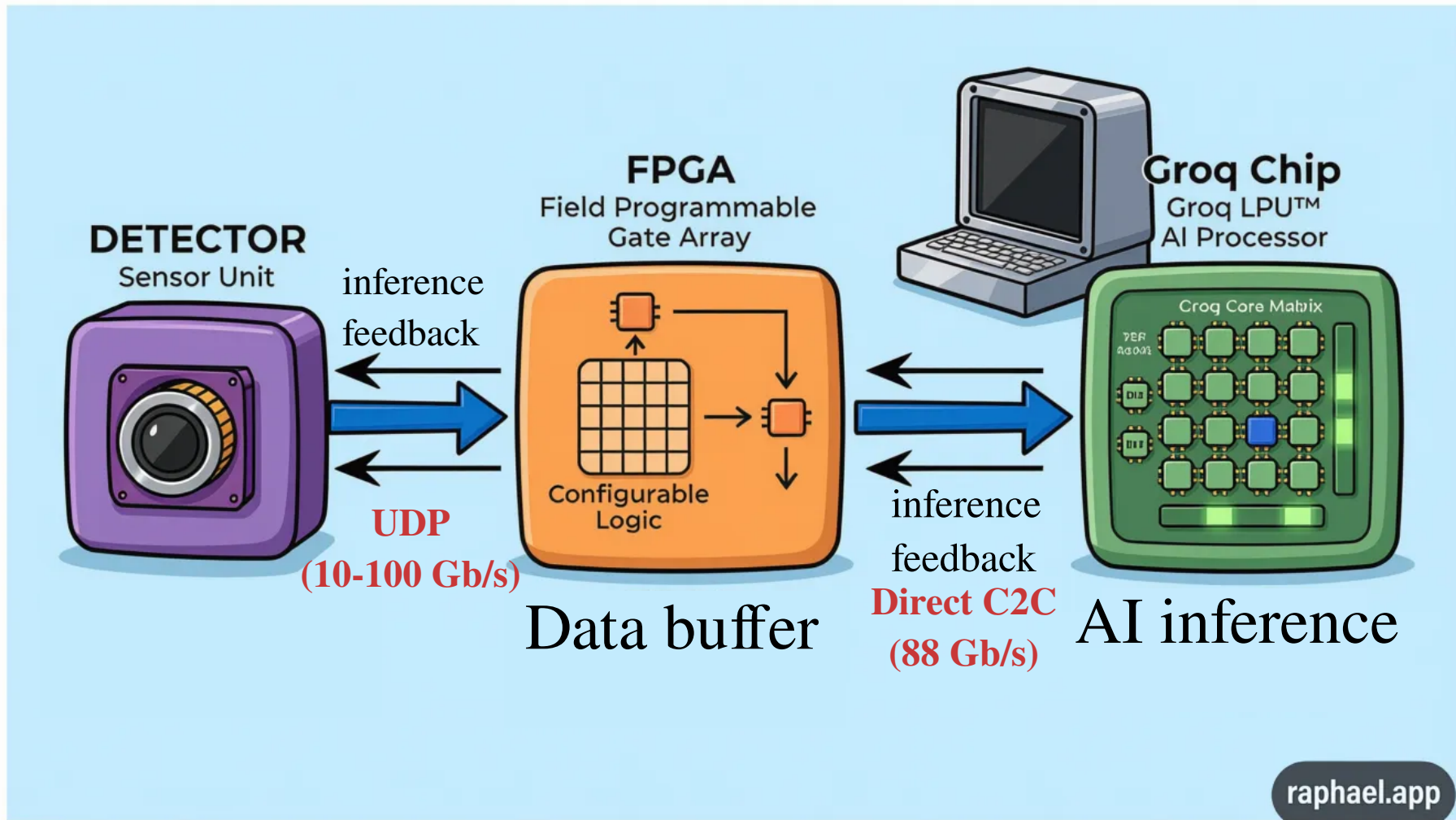


supported operands:

VXM: int8, int16, int32, uint8, uint16, uint32, float16, float32, bool8, bool16, bool32

MXM: int8 x int8 → int32, float16 x float16 → float32

# Inference loopback via data-flow deterministic architecture

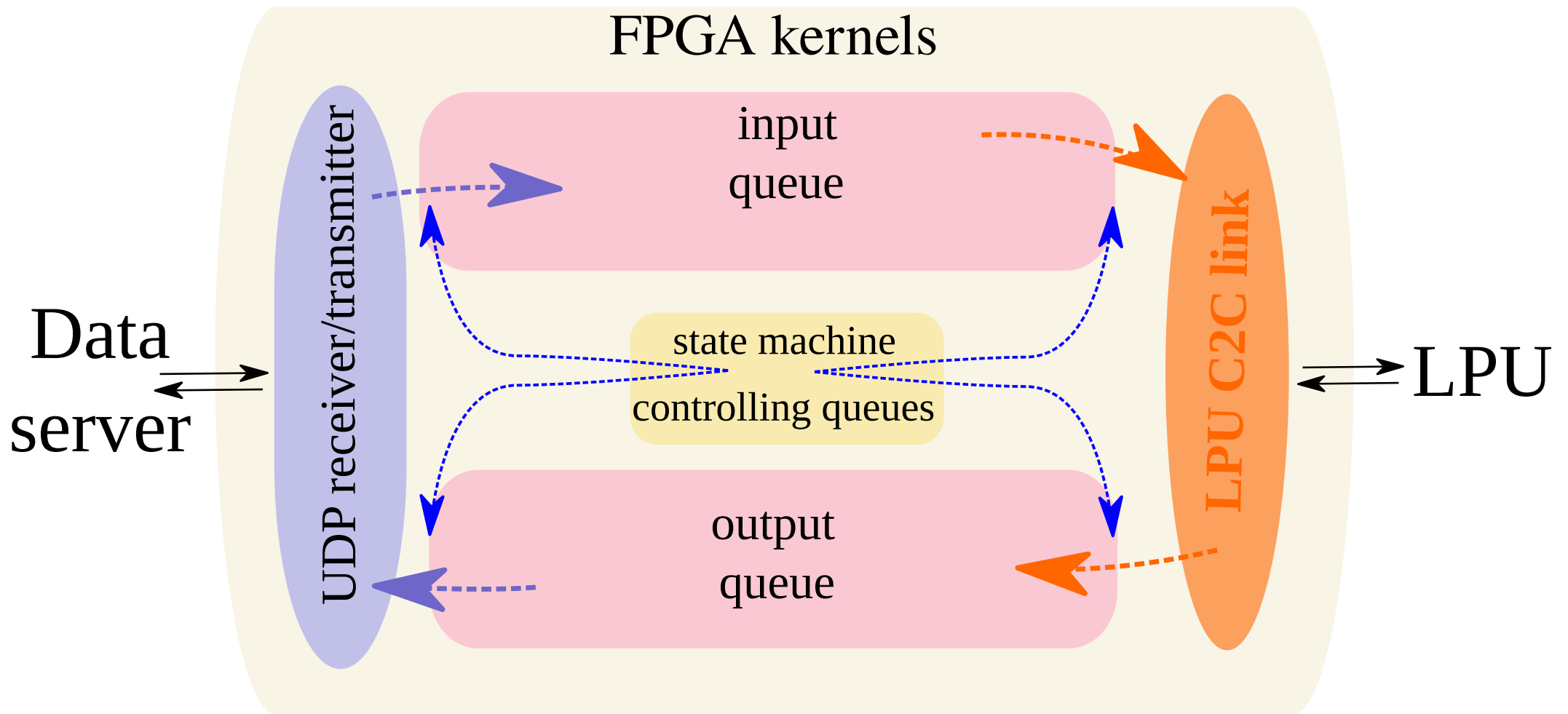


**Inference latency:**

X-ray imaging  $\sim 45 \mu s$  (22 kHz)

fusion plasma characterization  $\sim 10 \mu s$  (100 kHz)

# FPGA gate keeper



# Inference engine: workflow

- Hardware bring-up: Groq sends "**notify**" to the FPGA  
*How much input data is required for the inference?*
- FPGA collects input data incoming through the UDP channel.
- FPGA triggers the Groq chip to run the inference program after the FPGA collected the input data
- The Groq chip sends "**read**" command to the FPGA as part of the inference program
- The FPGA starts streaming the inference input to Groq (92 Gbs)
- The Groq chip does the calculations
- The Groq chip sends "**write**" command and the data to the FPGA
- The FPGA starts streaming the received inference return via the UDP channel
- Groq sends "**notify**" to the FPGA

# I/O Data compression

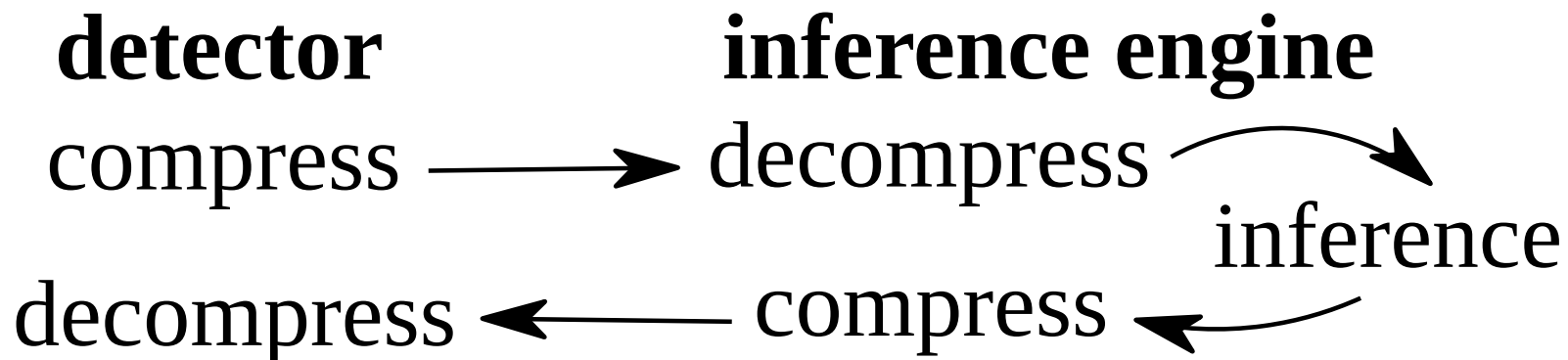
- Transmitting data between the detector and the inference engine is a bottleneck



Feature	FP4 (E2M1)	MXFP4	NVFP4
<b>Format Structure</b>	4 bits (1 sign, 2 exponent, 1 mantissa) plus software scaling factor	4 bits (1 sign, 2 exponent, 1 mantissa) plus 1 shared power-of-two scale per 32 value block	4 bits (1 sign, 2 exponent, 1 mantissa) plus 1 shared FP8 scale per 16 value block
<b>Accelerated Hardware Scaling</b>	No	Yes	Yes
<b>Memory</b>	Up to 4x less memory than FP16		
<b>Accuracy</b>	Risk of noticeable accuracy drop compared to FP8	Risk of noticeable accuracy drop compared to FP8	Lower risk of noticeable accuracy drop particularly for larger models



- We apply the following compression approach



OCP Microscaling Formats (MX) Specification  
 concept of: FP8, FP4  
 definition of mantissa  
 and exponent bits

# I/O Data compression

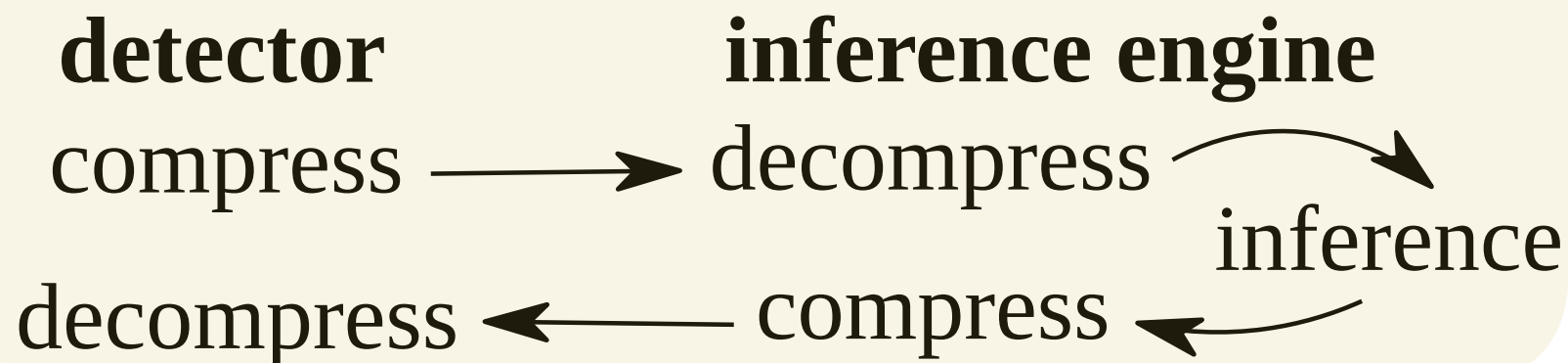


- Transmitting data between the detector and the inference engine is a bottleneck
- Our proof-of-concept solution is to use 0 exponent bits: **fixed-point number representation**

$$\text{FP32/FP16 data} = \text{scale}_{\text{FP32}} \times \text{uint8/uint4 data} + \text{offset}_{\text{FP32}}$$

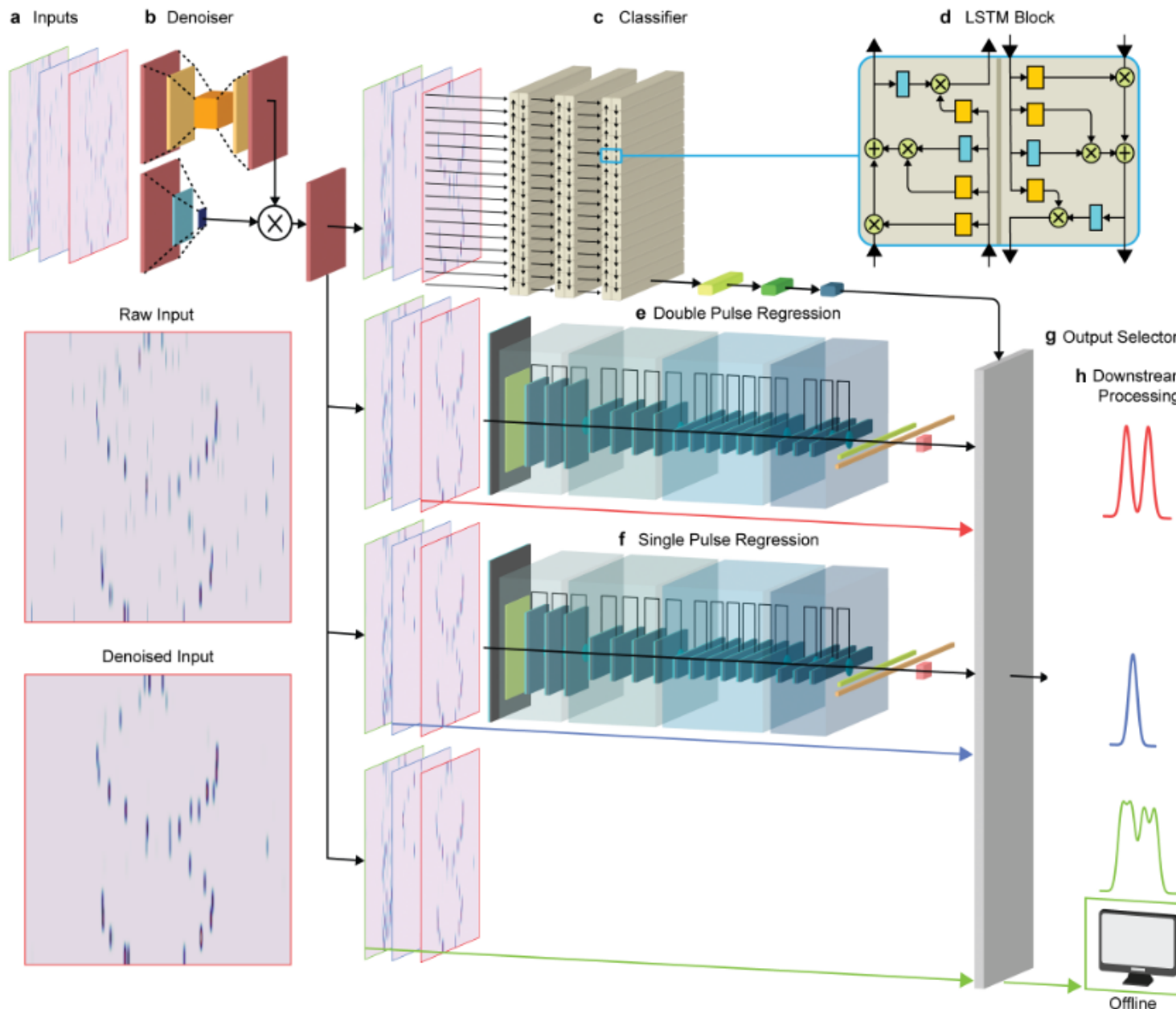
- **The compressed data are always nonzero ( $\text{offset}_{\text{FP32}}$ )**
- **Scale data within (0.0, 1.0):  $\text{scale}_{\text{FP32}}$**
- **Multiply with  $2^4$  or  $2^8$ .**
- **Rounding to nearest even integer**
- **Bundle together compressed data and scaling factors**

- We apply the following compression approach



**A Hybrid Neural Architecture: Online Attosecond X-ray Characterization**

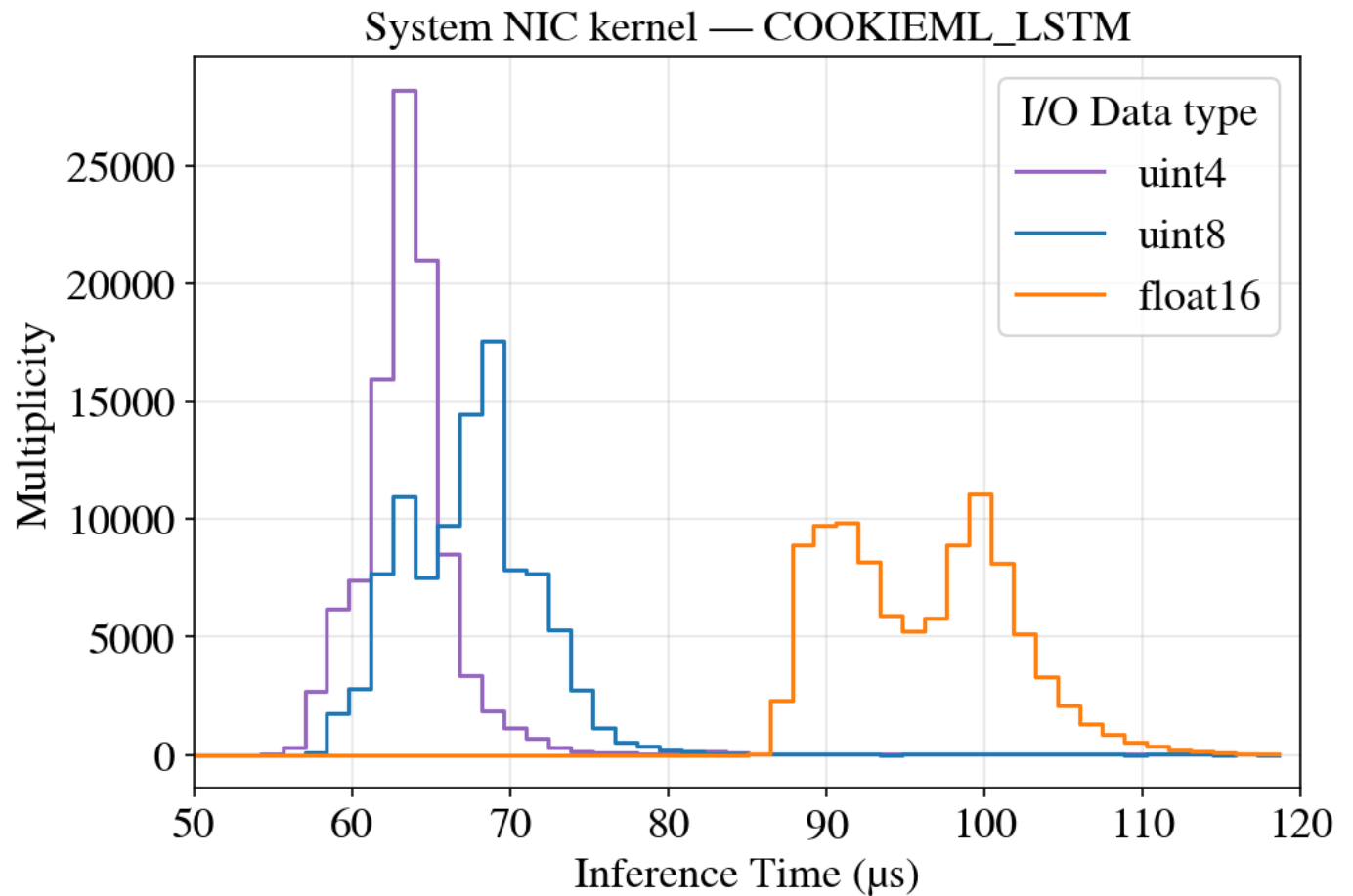
Jack Hirschman, Benjamin Mencer, Razib Obaid, Amanda Shackelford, Ryan Coffee



- Denoiser
- Zero pulse classifier
- LSTM based pulse num classifier
- ResNet based feature regression

# Inference latency on the LSTM pipeline

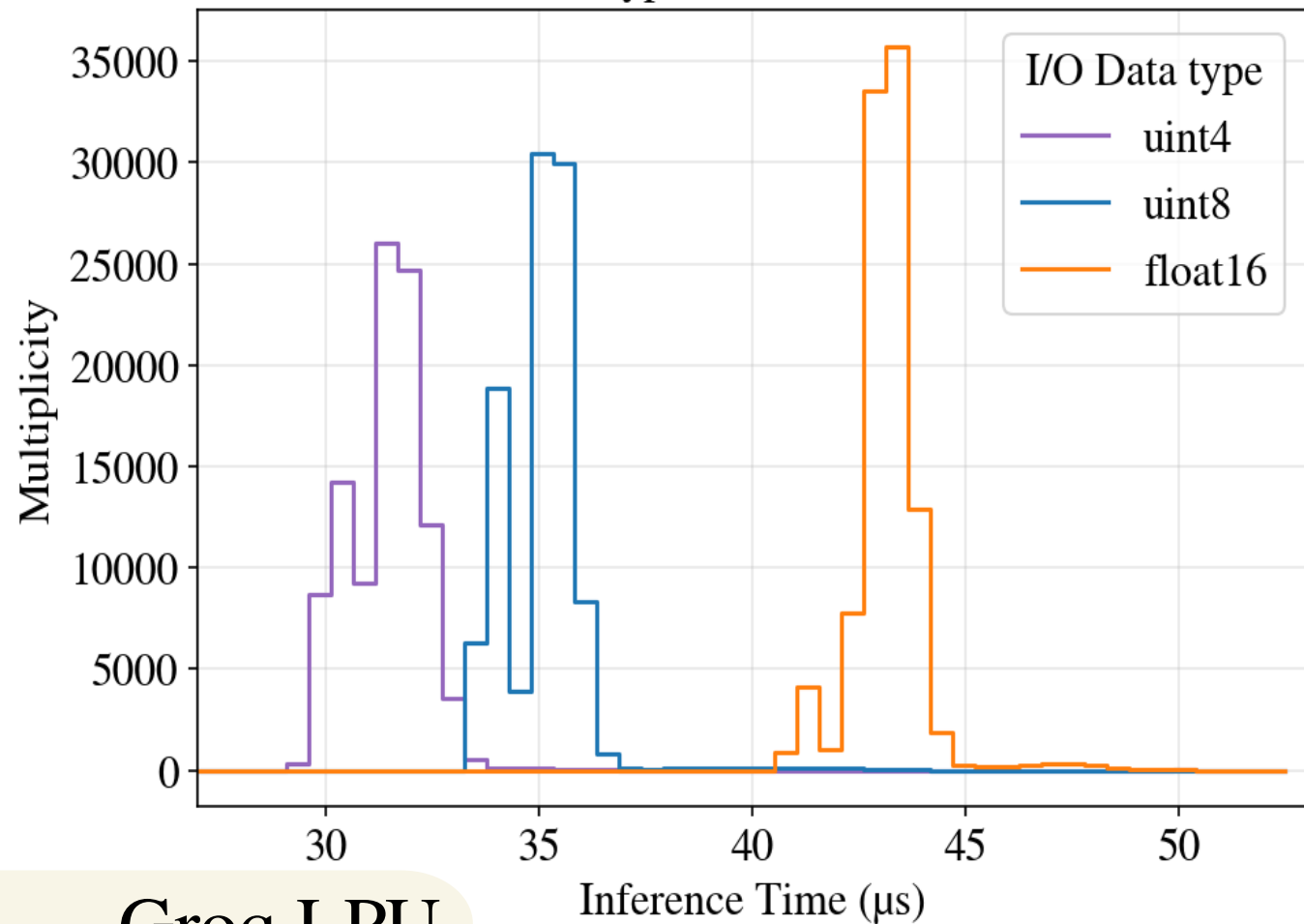
- Denoiser
- Zero pulse classifier
- LSTM based pulse num classifier
- Groq compute time:
  - int4: 19  $\mu$ s
  - int8: 20  $\mu$ s
  - float16: 20  $\mu$ s



# Inference latency on the LSTM pipeline: DPDK kernel bypass

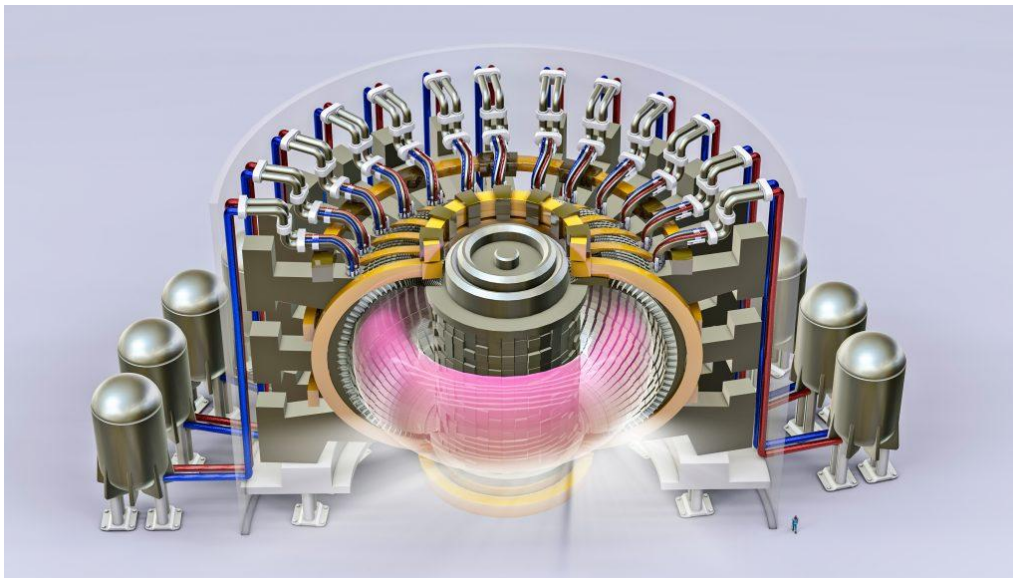
- bypasses the traditional Linux kernel networking stack to reduce latency
- DPDK is widely used in telecom, cloud networking, firewalls, routers, etc
- Groq compute time:
  - int4: 19  $\mu$ s
  - int8: 20  $\mu$ s
  - float16: 20  $\mu$ s

DPDK kernel bypass — COOKIEMML\_LSTM



# RTCAKENN

(Real-Time Convolutional Autoencoder  
Kernel-based Embedded Neural Network)



The model predicts

- pressure
- safety factor
- toroidal current density
- electron temperature
- electron density
- ion temperature
- plasma rotation profiles

- machine-learning model developed for real-time kinetic profile reconstruction in tokamak fusion plasmas
- developed by DIII-D National Fusion Facility

## Machine learning-based real-time kinetic profile reconstruction in DIII-D

JOURNAL ARTICLE · 22 December 2023 · Nuclear Fusion

DOI: <https://doi.org/10.1088/1741-4326/ad142f> · OSTI ID: 2251556

Shousha, Ricardo ; Seo, Jaemin ; Erickson, Keith; Xing, Zichuan ; Kim, SangKyeun ; Abbate, Joseph ; Kolemen, Egemen

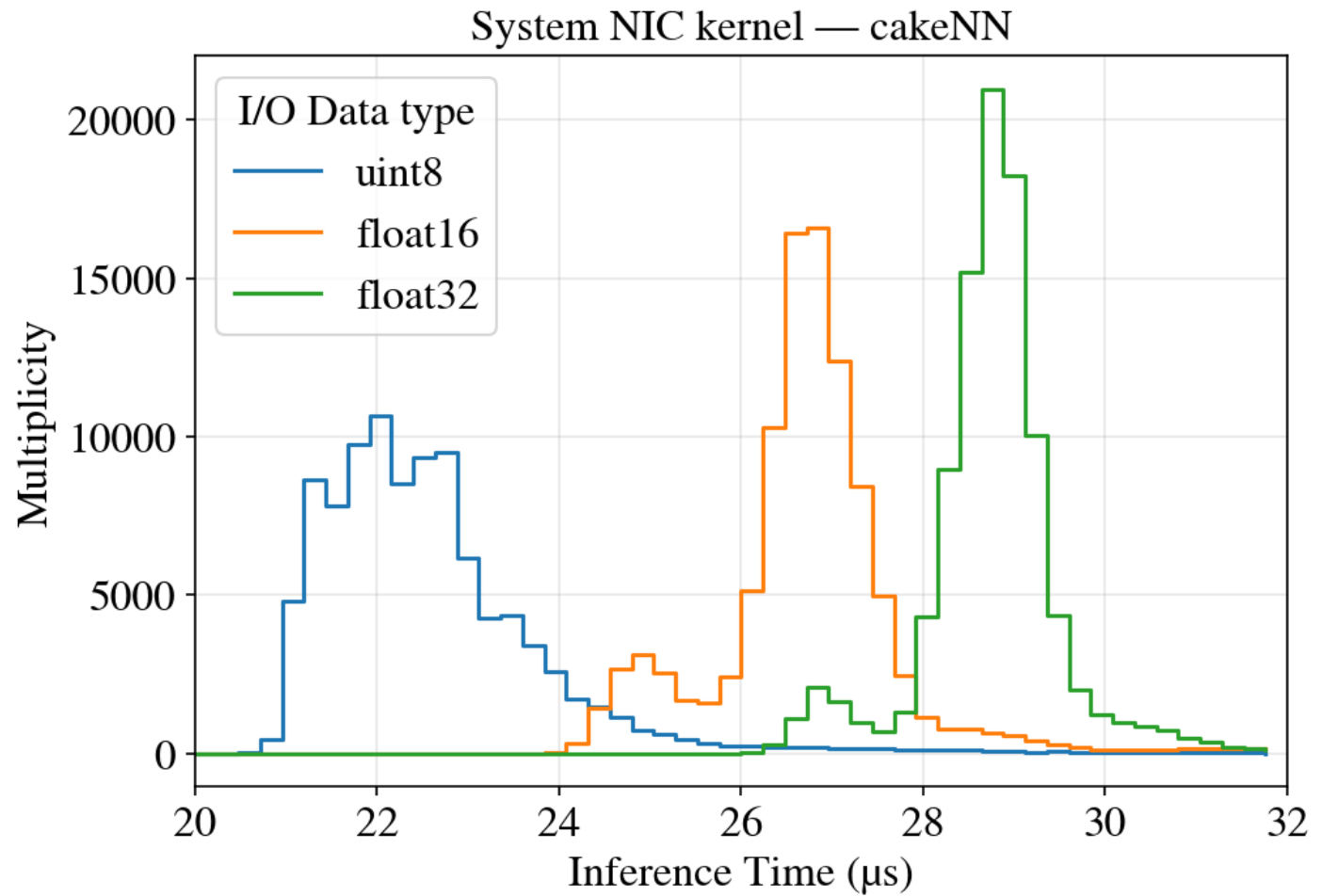
**Inference runs in under about 8 ms** (10 runs to get averaged output)

# Inference latency on the cakeNN pipeline

- pressure
- safety factor
- toroidal current density
- electron temperature
- electron density
- ion temperature
- plasma rotation profiles

## ● Groq compute time:

int8: 4  $\mu$ s  
float16: 4  $\mu$ s  
float32: 4  $\mu$ s

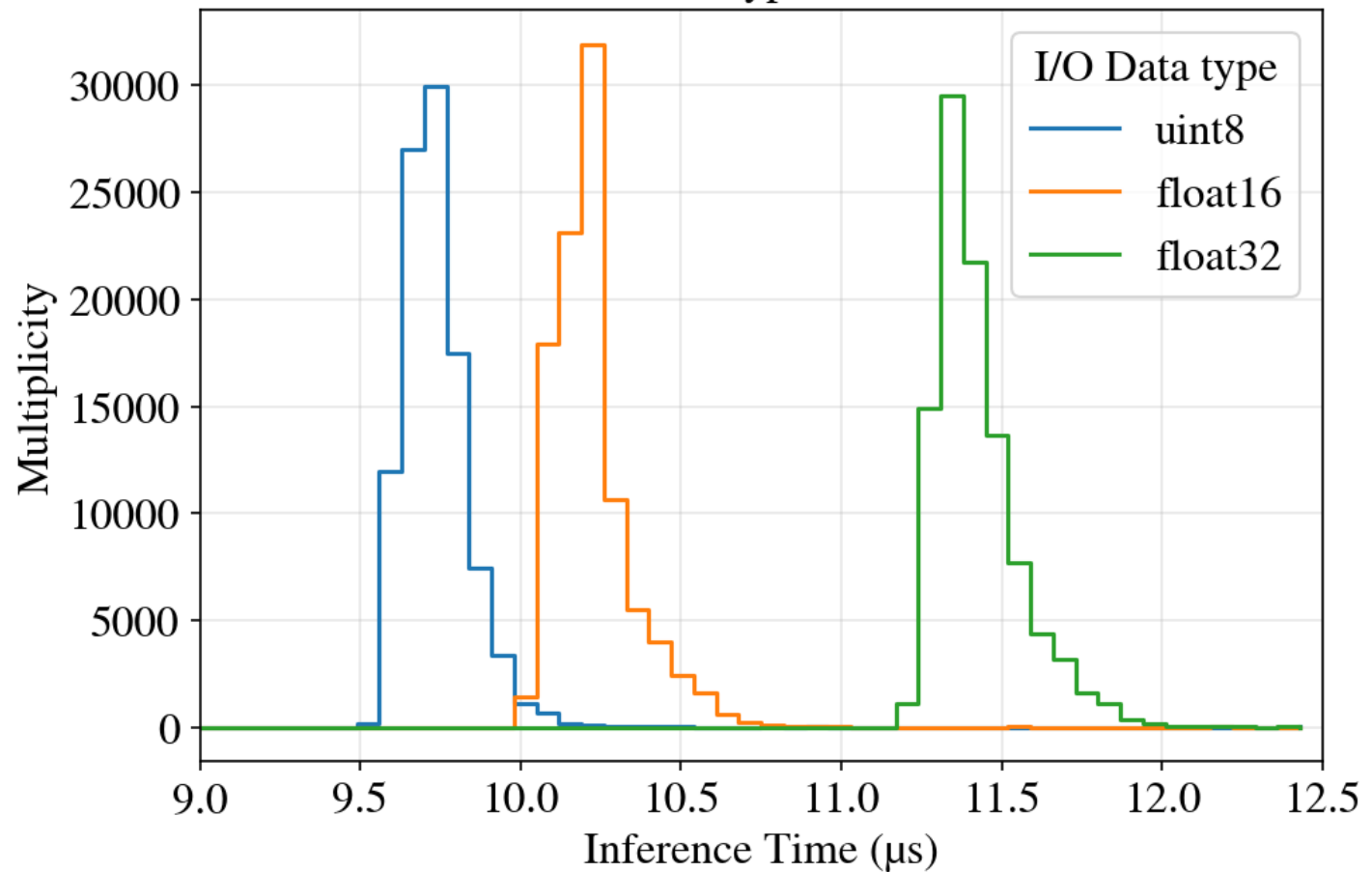


previously: ms time scale

# Inference latency on the cakeNN model: DPDK kernel bypass

- The I/O overhead is still too big
- further improvements:
  - FPGA ↔ FPGA
  - reimplement the FPGA gate keeper
- Groq compute time:
  - int8: 4 μs
  - float16: 4 μs
  - float32: 4 μs

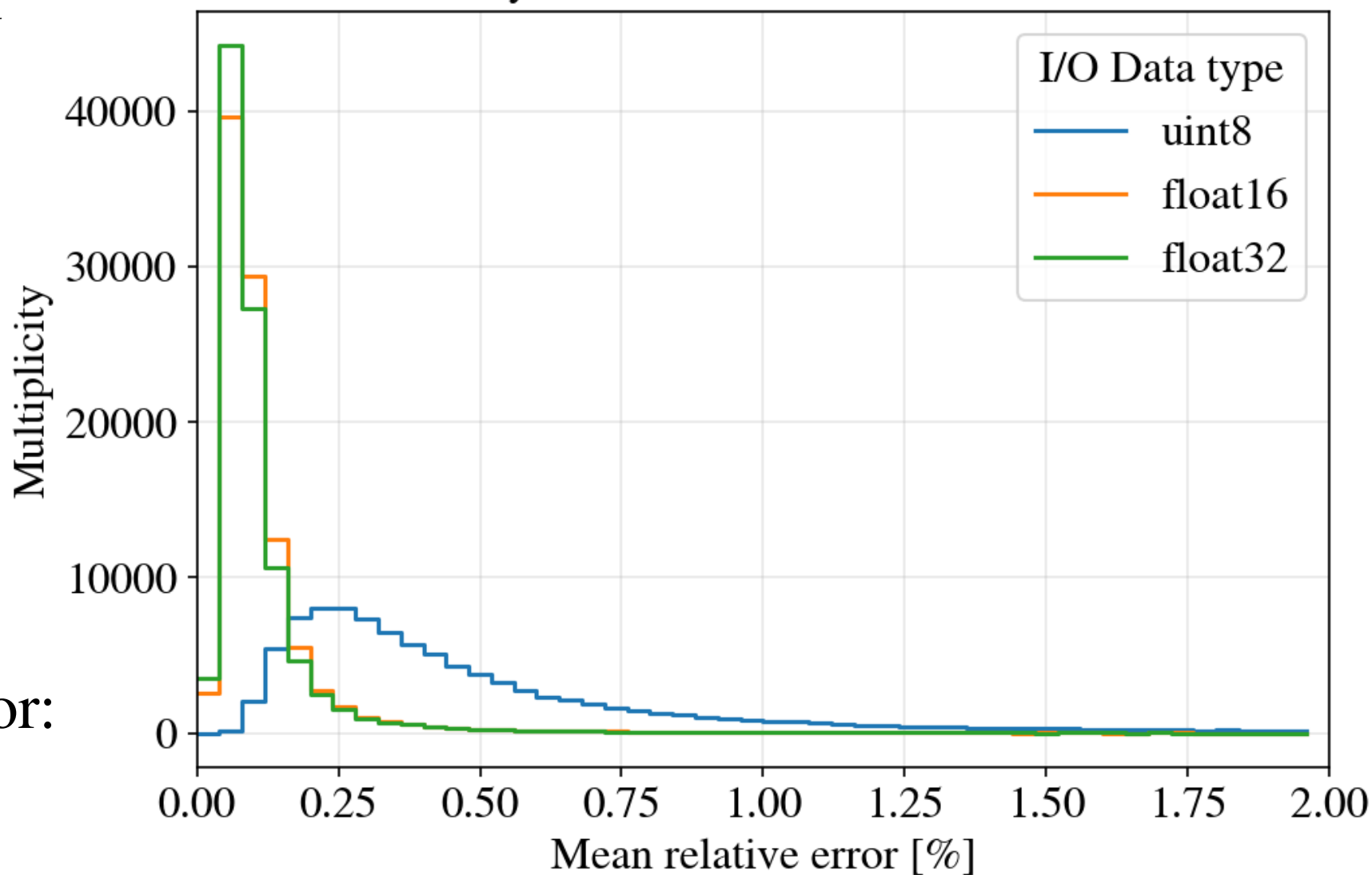
DPDK kernel bypass — cakeNN



# Numerical accuracy: cakeNN model

Element-wise comparison with  
float32 torch model  
evaluated on CPU

System NIC kernel — cakeNN

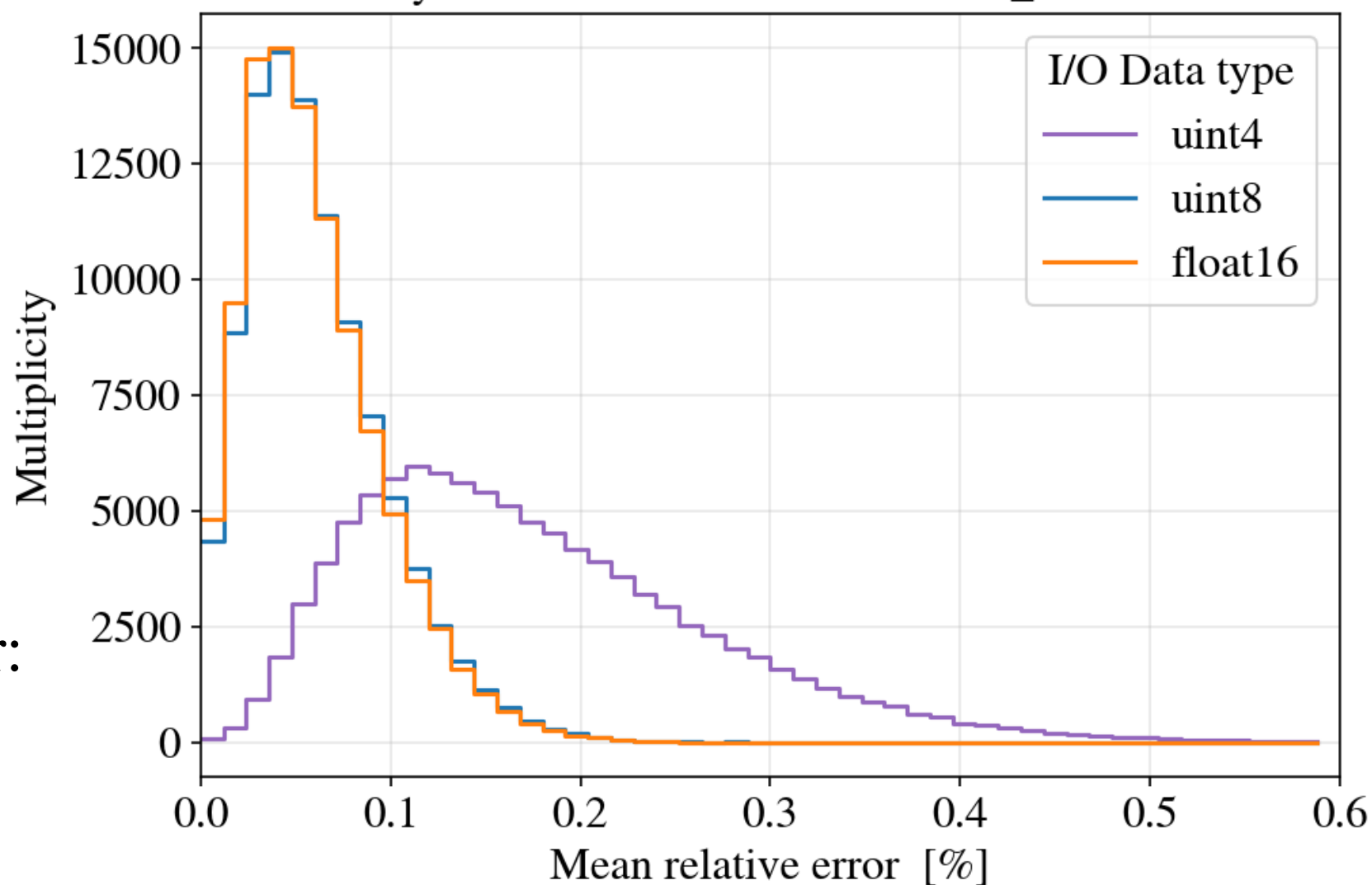


max mean error:  
less than 4-6%

# Numerical accuracy: CookieML LSTM model

Element-wise comparison with  
float16 torch model  
evaluated on CPU

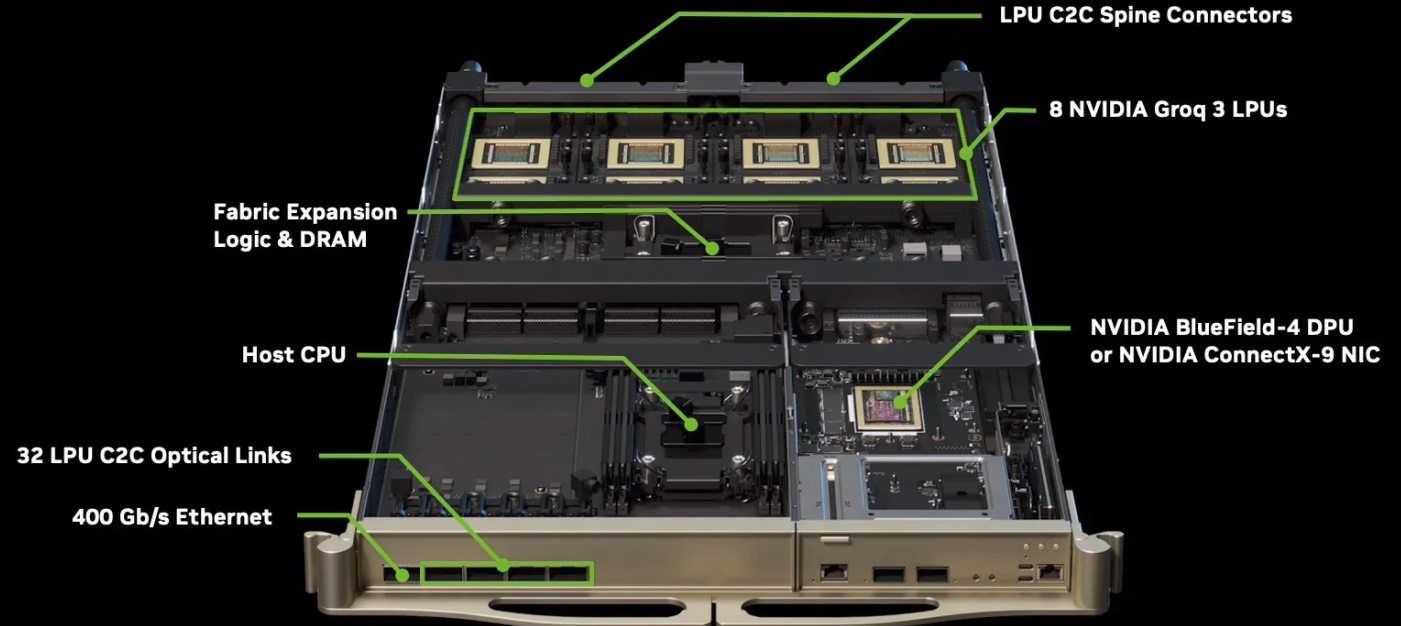
System NIC kernel - COOKIEMML\_LSTM



max mean error:  
less than 1-2%

# NVIDIA Groq 3 LPX Compute Tray

Rack-ready, liquid cooled, 1U compute engine for low-latency inference

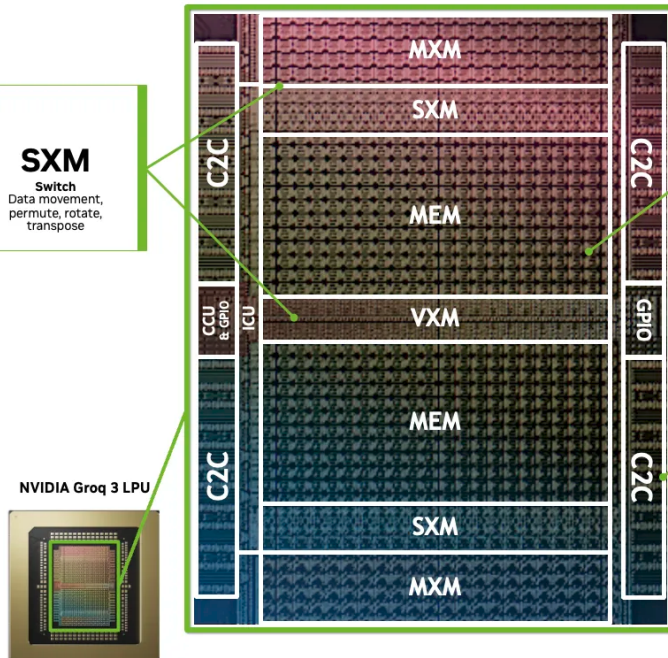


LPX compute tray specifications and configuration are preliminary and subject to change;

## NVIDIA Groq 3 LPU Chip Architecture

**Tensor-First Compute**

- MXM**  
Dense MatMuls  
1.2 PFLOPs FP8
- VXM**  
Vector Math  
Pointwise math, activations, type conversions
- SXM**  
Switch  
Data movement, permute, rotate, transpose



**Ultra High SRAM Bandwidth**

- 500 MB**  
On-Chip SRAM  
Primary weight, activation, and instruction storage
- 150 TB/s**  
Memory Bandwidth  
Enables Ultra Fast Token Generation

**Chip-to-Chip Interconnect**

- 96**  
LPU-to-LPU Connections  
6x increase in communication radix
- 112 Gbps**  
SerDes Speed  
Single link per connection

Specification	NVIDIA Groq 3 LPX
AI inference compute	315 PFLOPS
Total SRAM capacity	128 GB
On-chip SRAM bandwidth	40 PB/s
Scale-up density	256 chips
Scale-up bandwidth	640 TB/s

# Acknowledgement



This research was supported by the Ministry of Innovation and Technology and the National Research, Development and Innovation Office within the Quantum Information National Laboratory of Hungary and Grants No. 2022-2.1.1-NL-2022-00004, by the ÚNKP-24-5 New National Excellence Program of the Ministry for Culture and Innovation from the source of the National Research, Development and Innovation Fund, by the Hungarian Scientific Research Fund (OTKA) Grant No. K134437 and by the Hungarian Academy of Sciences through the Bolyai János Stipendium (BO/00571/22/11).

We acknowledge the computational resources provided by the Wigner Scientific Computational Laboratory (WSCLAB) (the former Wigner GPU Laboratory)

**contact:** Peter Rakyta, [rakyta.peter@wigner.hu](mailto:rakyta.peter@wigner.hu)

