

TrustForge: A Cryptographically Secure Enclave

Todd Austin, CEO
Valeria Bertacco, Chief Scientist
Alex Kisil, Director of Engineering

austin@agitalabs.com

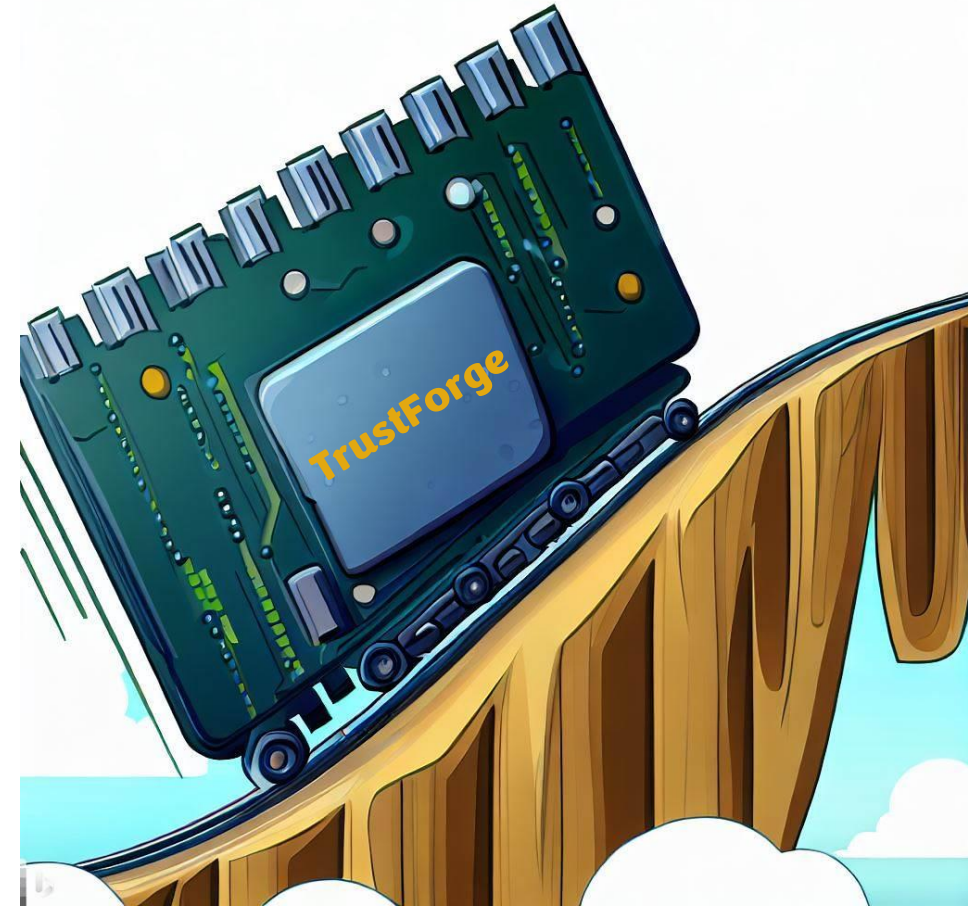


AGITA
LABS



The Little (TrustForge) Enclave that Could...

- Execute RISC-like instructions directly on encrypted data
- Transform a hackable CPU into next-gen privacy tech that surpasses FHE and ZKP
- Nullify software and hardware hacking
- Pass commercial red teaming and full formal verification with zero vulnerabilities
- Stop data breaches once and for all!
- Advance data privacy by denying software access to sensitive data it is processing
- Do this with only a 190k-gate functional unit!



*Drop by my poster
to learn more!*

An Open-Source 130-nm Fusion-Enabled Deconvolution Kernel Generator IC For Real-Time mmWave Radar Platform Motion Compensation

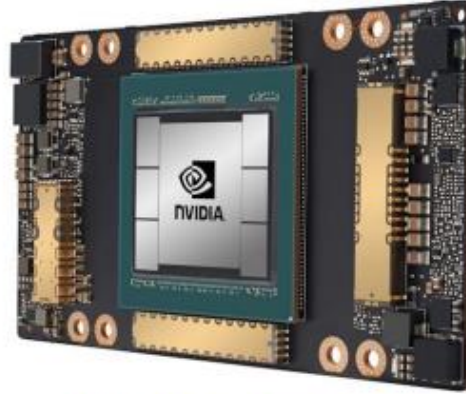
NIKHIL POOLE, PRIYANKA RAINA, AND AMIN ARBABIAN

AUGUST 29, 2023

Problem/Motivation



Autonomous Vehicles



High-Performance GPUs/TPUs



Smart Infrastructures

Modern sensor networks are powerful, consuming 100s of W of power to achieve high resolution/performance.

But how can we process the same quantity of data on resource-limited edge devices with minimal latency?

Micro-Drones



AR Glasses



Wearable/Portable Aids



Goal:
Real-time edge sensing providing high-resolution with minimal power consumption.

mmWave Radar: An Optimal Sensing Modality

cm-Scale Range Resolution



High Velocity Resolution



Mm-Scale Accuracy/Sensitivity



High Maximum Range + Velocity



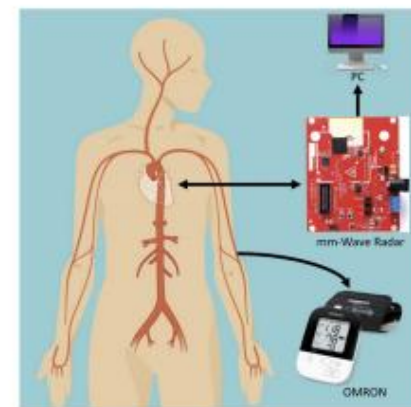
Small, Easily Integrable Form Factor



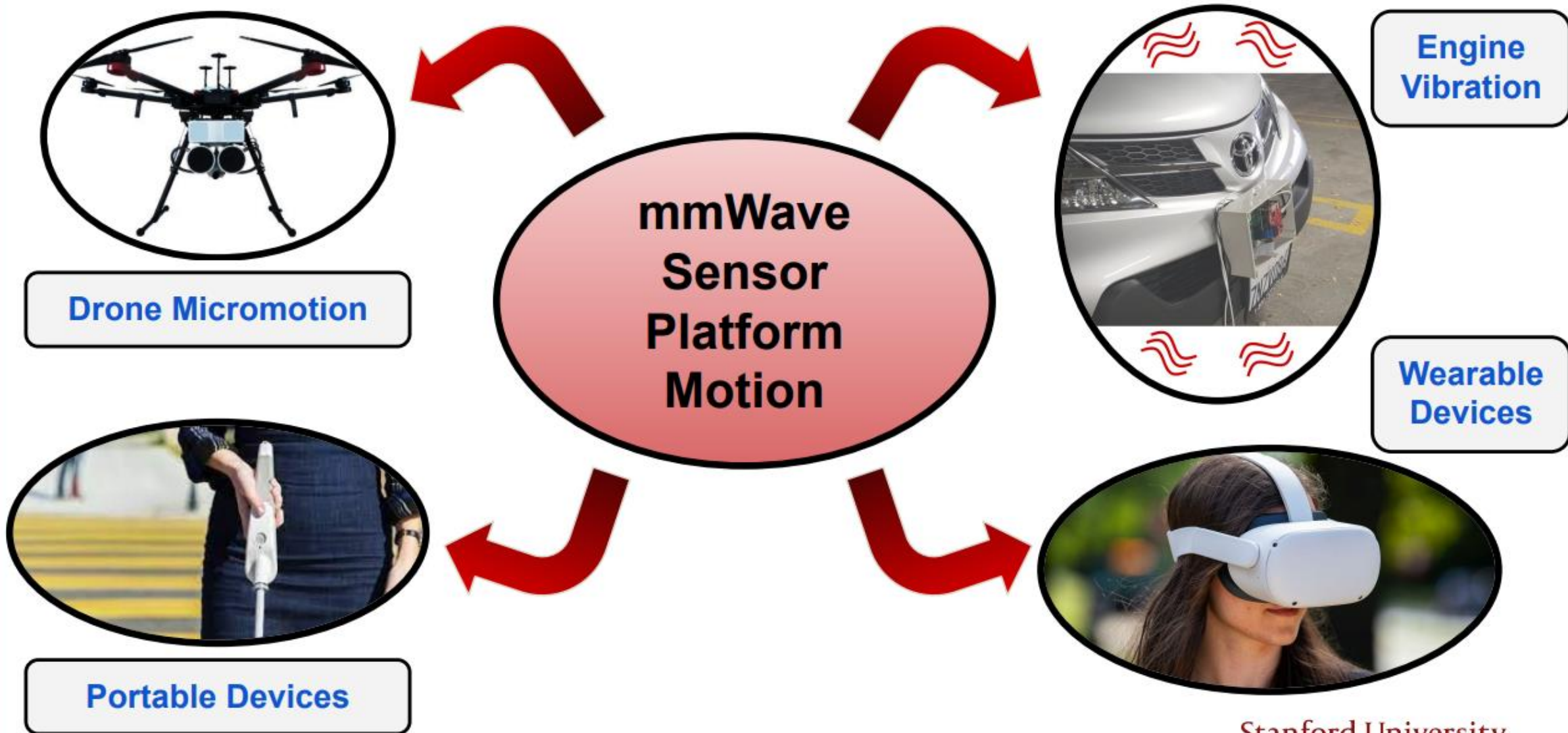
Environmentally Robust



Widely Applicable

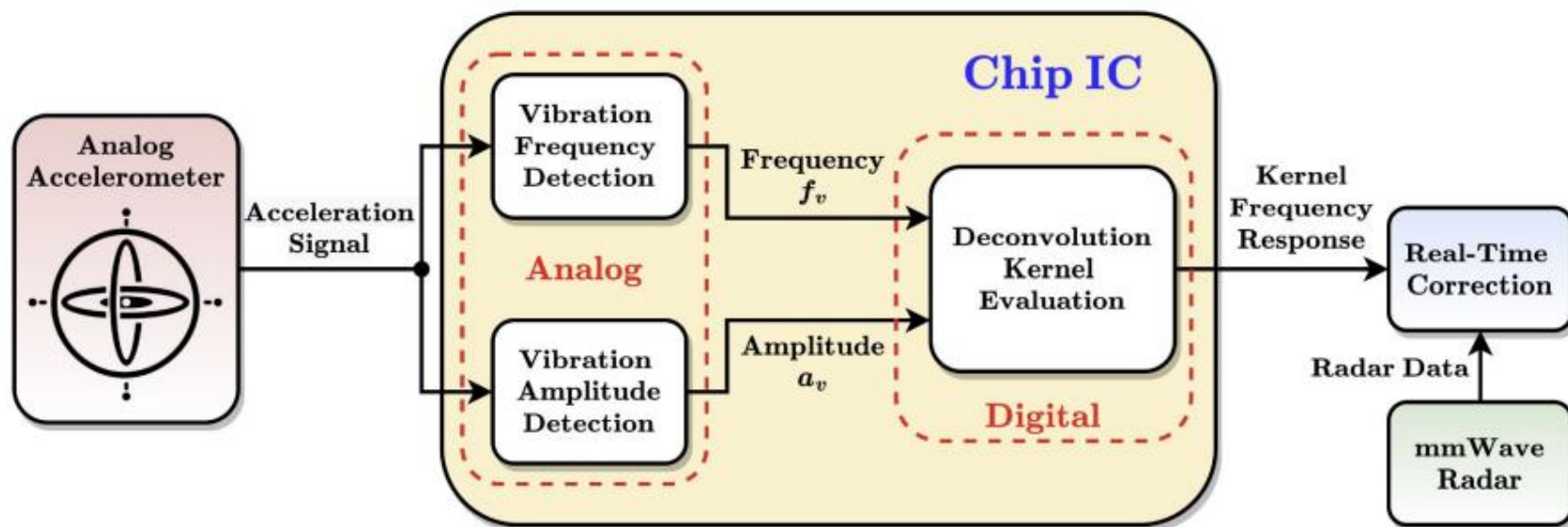


Primary Challenge With High-Resolution Edge Sensing



Solution

First custom IC enabling real-time, resource-efficient FMCW mmWave radar platform vibratory motion compensation, via early sensor fusion.



**26-35 Vibration
Suppression**

**FULLY
OPEN-SOURCE**

< 95-ms Latency

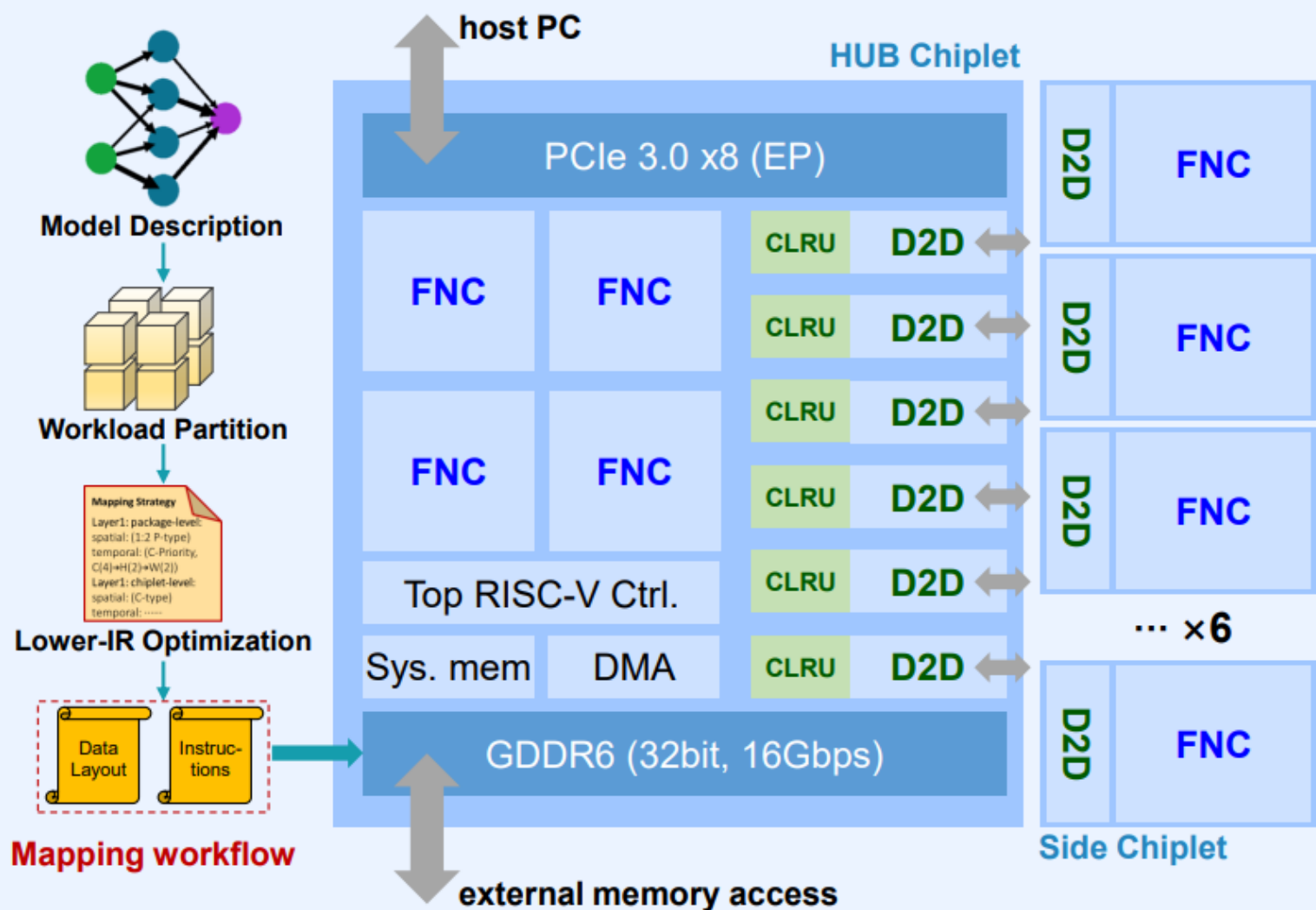
A Scalable Multi-Chiplet Deep Learning Accelerator with Hub-Side 2.5D Heterogeneous Integration

Zhanhong Tan¹, Yifu Wu², Yannian Zhang²,
Haobing Shi², Wuke Zhang², Kaisheng Ma¹

¹Tsinghua University,
²Polar Bear Tech



Overall Architecture



Flexible Neural Core (FNC)

- Reconfigurable architecture for the shape diversity

Mapping dataflow

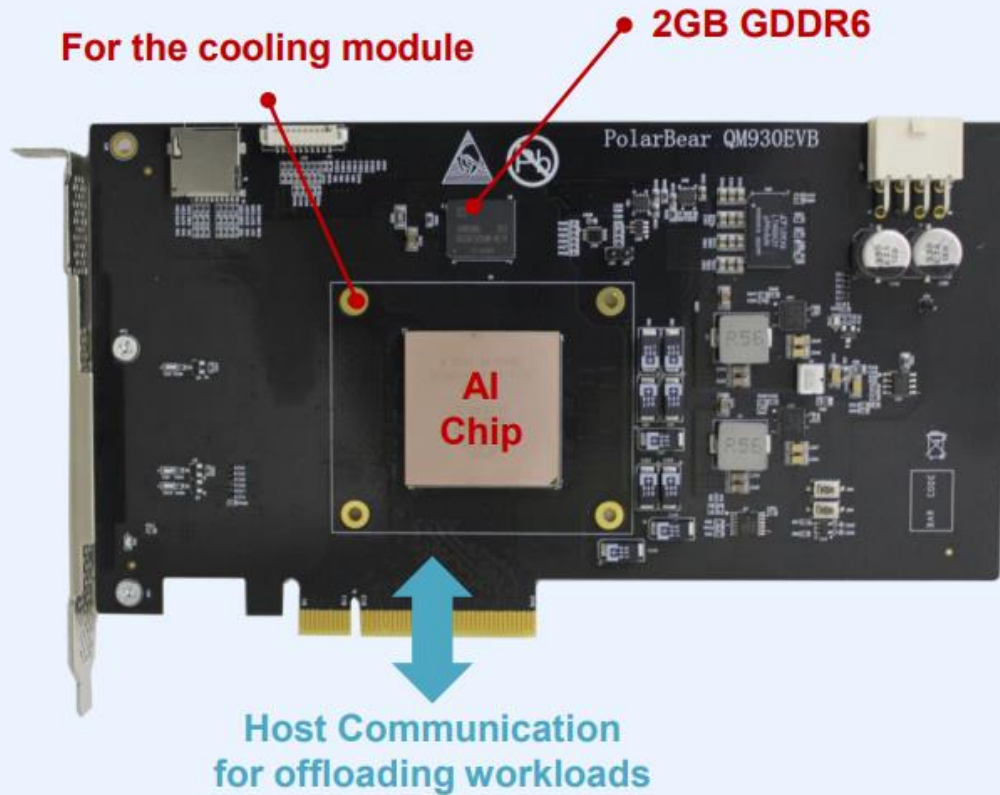
- Die-to-Die communication-aware workload generator

Interconnection

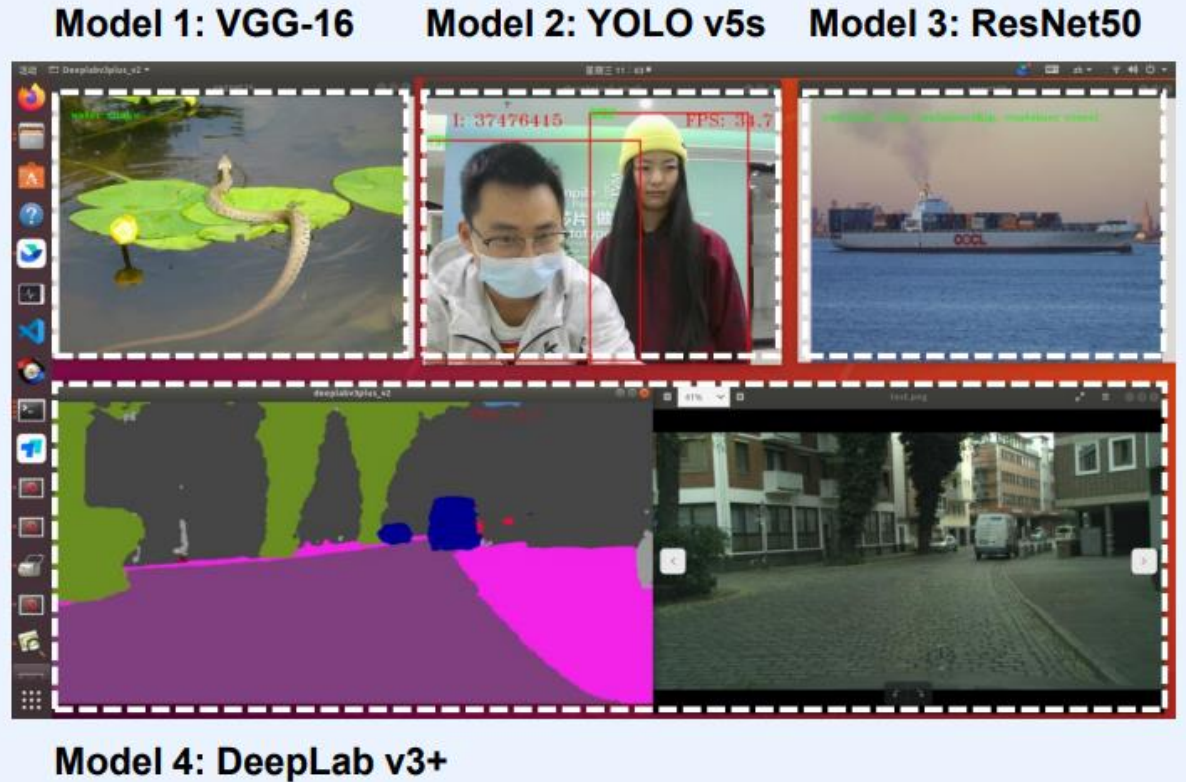
- High-bandwidth Die-to-Die based on 2.5D package
- Efficient chiplet routing unit (CLRU)

System Board and Demo

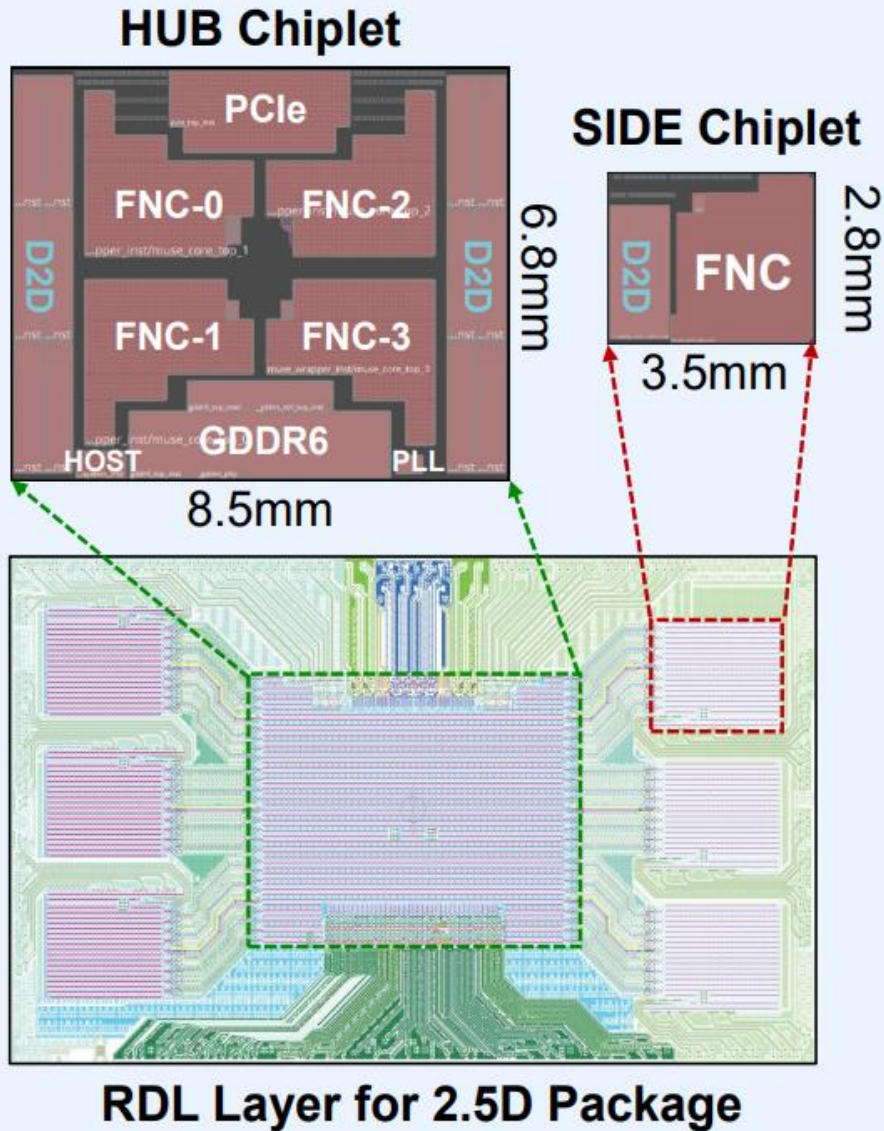
■ System PCIe-based Board



■ Demo for running concurrent 4 models



Chip Summary



Items		Specifications
Technology		CMOS 12nm
Die Area	HUB Chiplet	8.5mm × 6.8mm = 57.8mm ²
	SIDE Chiplet	3.5mm × 2.8mm = 9.8mm ²
Supply Voltage		0.8V ~ 1.2V
Frequency		100MHz – 1GHz
Peak Performance	INT4	40TOPS (8b-A × 4b-W)
	INT8	20TOPS (8b-A × 8b-W)
	INT16	10TOPS (16b-A × 8b-W)
NPU Core Efficiency		2.02TOPS/W
Power		4.5W ~ 12W
D2D Bandwidth		6×24GB/s for TX/RX
External Memory Bandwidth		64GB/s (GDDR6)
Bump Pitch for 2.5D Pkg		55μm



PHEP: Paillier Homomorphic Encryption Processors for Privacy-Preserving Applications in Cloud Computing

Guiming Shi¹, Yi Li², Xueqiang Wang², Zhanhong Tan¹, Dapeng Cao³,
Jingwei Cai¹, Yuchen Wei¹, Zehua Li³, Wuke Zhang⁴, Yifu Wu⁴,
Wei Xu^{1*}, and Kaisheng Ma^{1*}

¹Tsinghua University

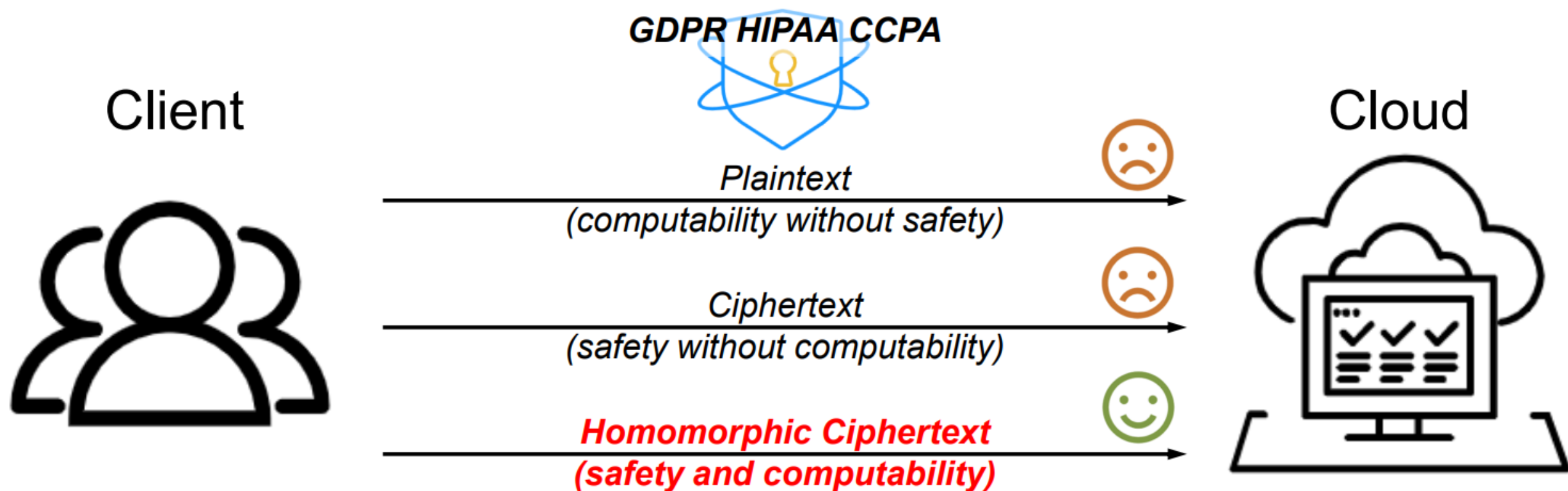
²HuaKong TsingJiao

³Xi'an JiaoTong University

⁴Polar Bear Tech

Homomorphic Encryption in Cloud Computing

- Data privacy is a critical problem in Cloud Computing.
- Paillier Homomorphic Encryption can protect the privacy of the data and enable computing on the ciphertext without decryption first.



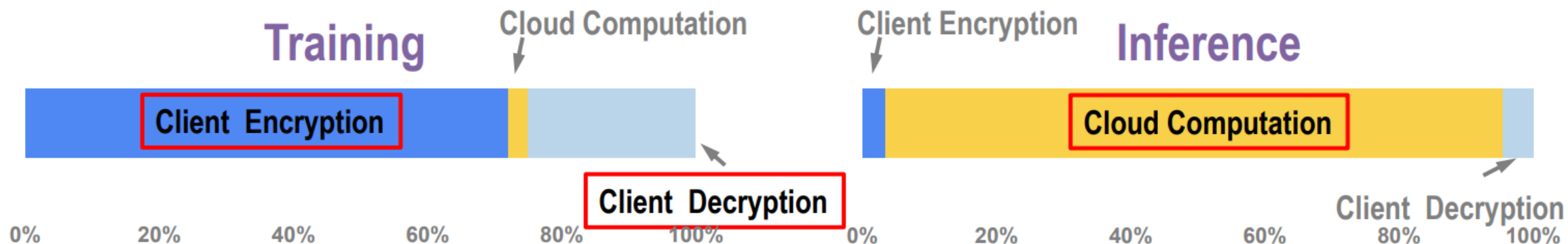
Bottleneck in Training and Inference Applications are Different

*Application-Training, Bottleneck:
Client Encryption and Decryption*

*Application-Inference, Bottleneck:
Cloud Computation*



Latency Breakdown On CPU

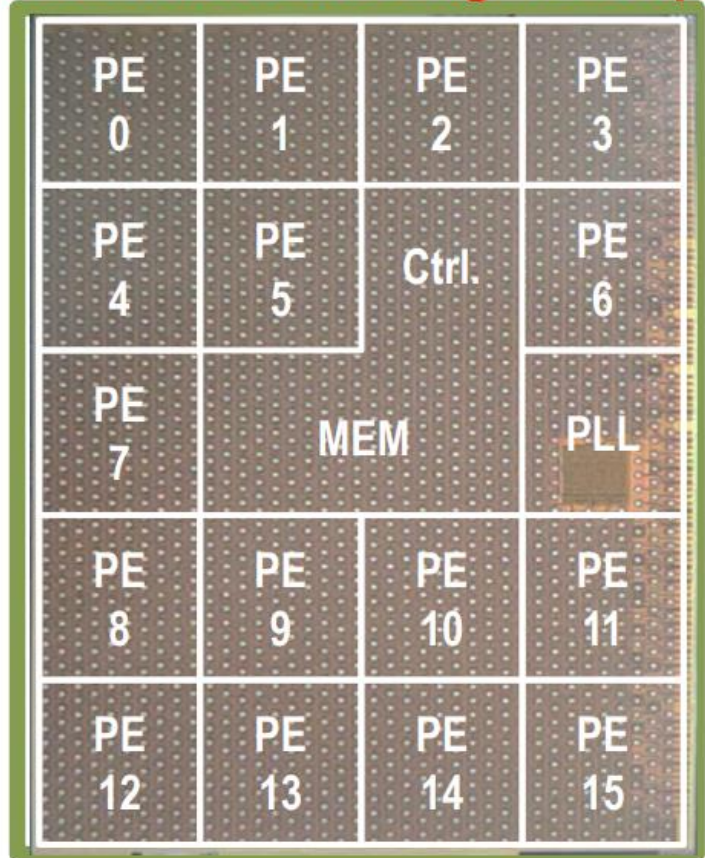


Comparison of the Two Engines: Specification

Fabricated on the Same Wafer: Significantly Reduces NRE of the Engine Chips



**Same Area: 43mm²
@ UMC 28nm HPC+**



- **Optimized for Parallelism**

Items/PE	Montgomery Unit
Algorithm	Montgomery Multiplication
Arithmetic	3*128 Bit Multiplier

- **Optimized for Performance**

Items/PE	Montgomery Unit	Stein Unit
Algorithm	Montgomery Multiplication	Stein Modular Inversion
Arithmetic	3*256 Bit Multiplier	3*4102 Bit Adder

High Performance Paillier Homomorphic Encryption Processors



PHEP Engine-1

- 480 TOPS (INT8)
- **Client Encryption: 84KOPs**
- **Cloud Computation: 402KOPs**
- **Client Decryption: 106KOPs**

PHEP Engine-2

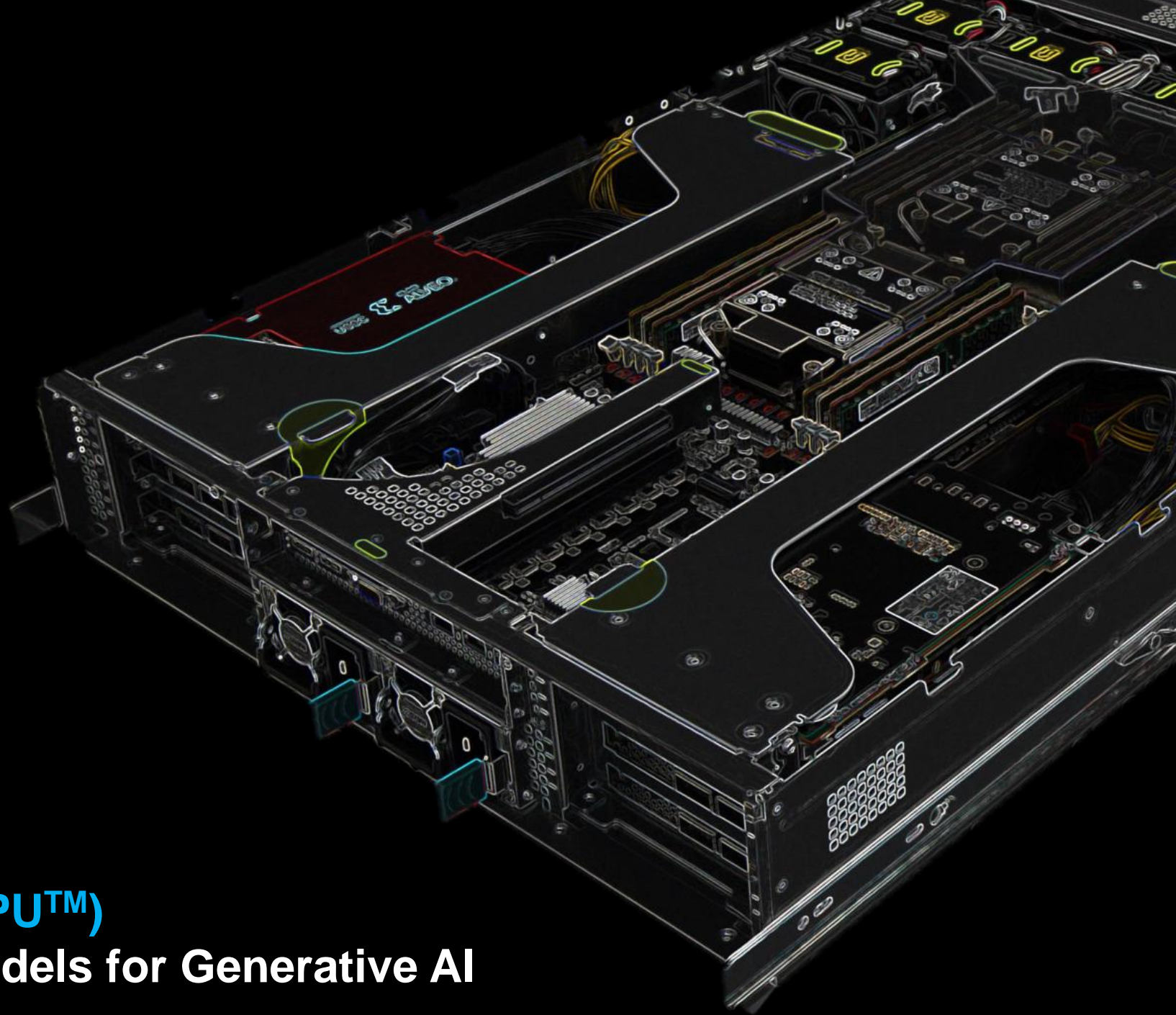
- 192 TOPS (INT8)
- **Client Encryption: 52KOPs**
- **Cloud Computation: 47MOPs**
- **Client Decryption: 48KOPs**

*Bit width of ciphertext = 4096, Bit width of plaintext = 64, Bit width of weight in Conv = 8.
Maximum Performance in Optimized Applications.*

Thank You

shigm21@mails.tsinghua.edu.cn





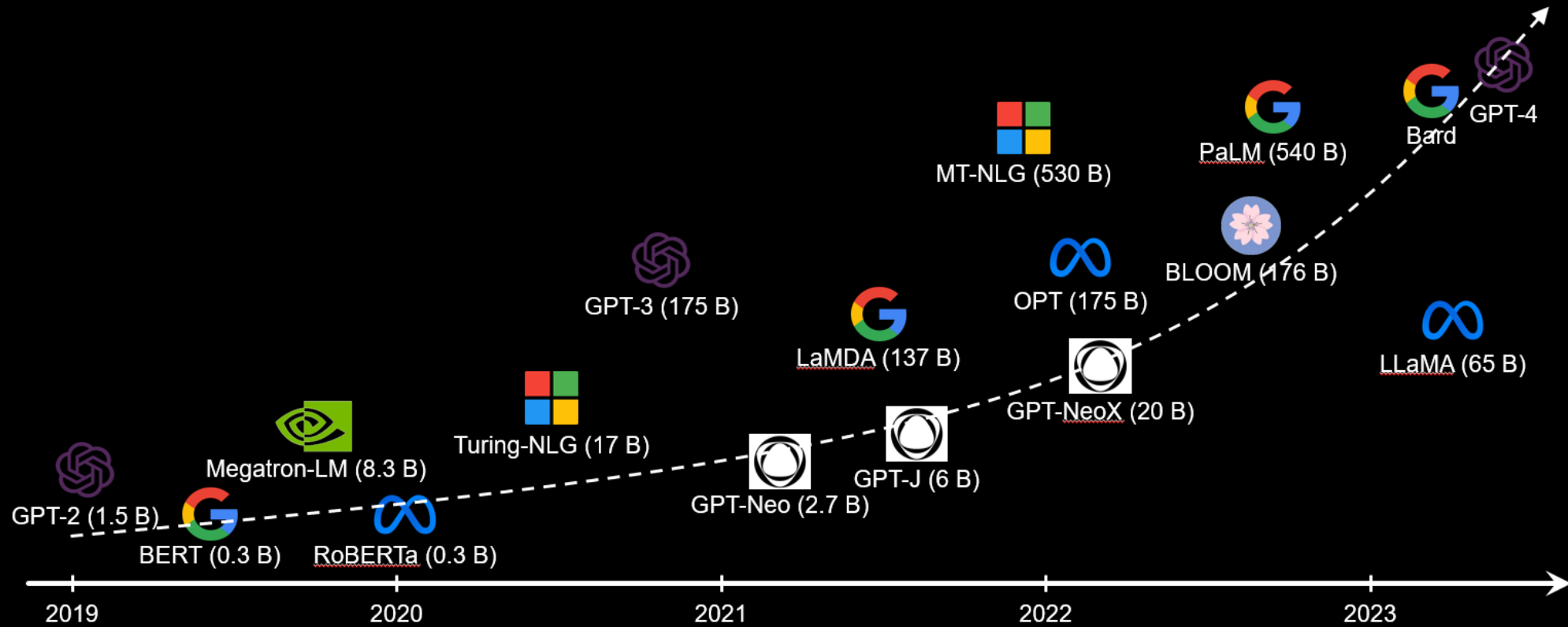
5

Seungjae Moon
Co-Founder, Research Scientist

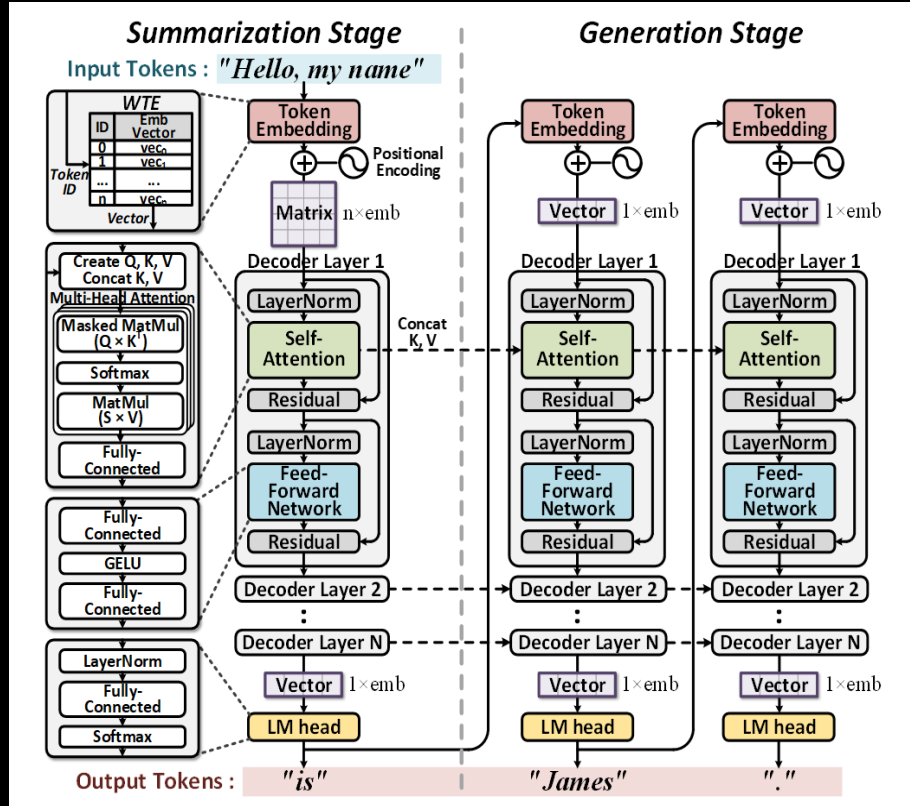
Latency Processing Unit (LPU™)
Accelerating Hyperscale Models for Generative AI

Large Language Model (LLM)

📄 Demand for highly scalable model inference hardware to support growing LLM



LLM Training & Inference Characteristics



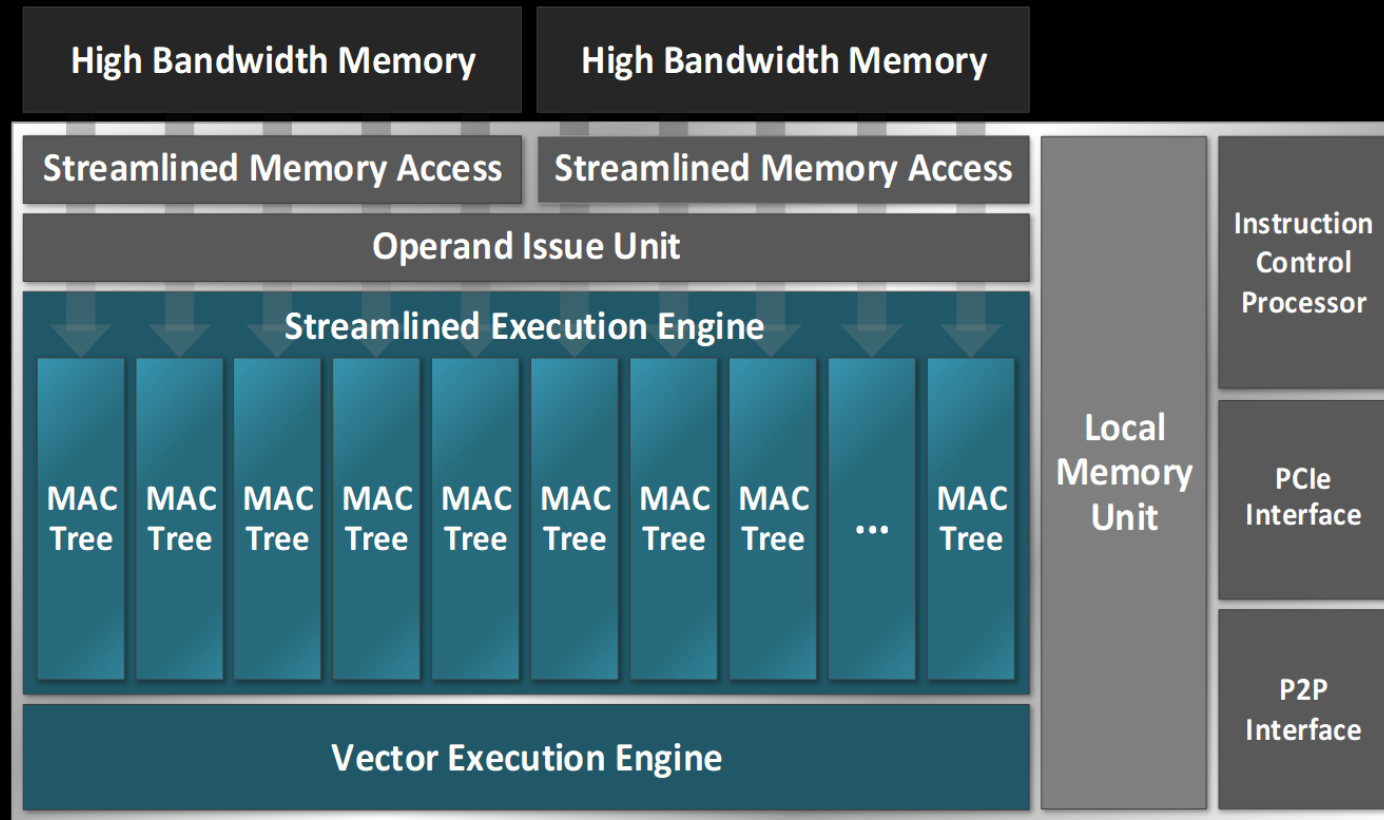
Training

- Large batch
- Compute-intensive
- Throughput-oriented hardware

Inference

- Small batch
- Memory-intensive
- Latency-oriented hardware

Latency Processing Unit (LPU)



WORLD-FIRST HARDWARE DEDICATED FOR END-TO-END INFERENCE OF LLM

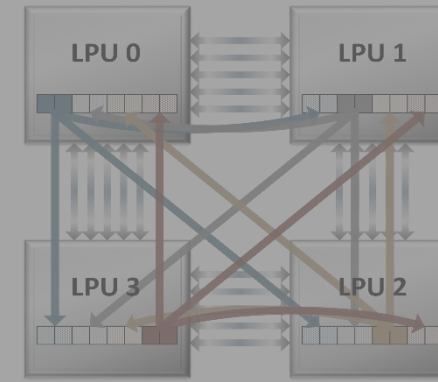
LPU Key Features

Streamlined Execution Engine

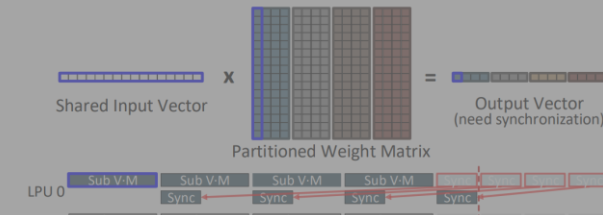
- Perfectly balances memory bandwidth and compute logic to maintain bandwidth usage of ~90%

Expandable Synchronization Link

- Lightweight full-duplex peer-to-peer communication technology with custom protocol
- Performs data synchronization with low latency and latency hiding to achieve near-perfect scalability



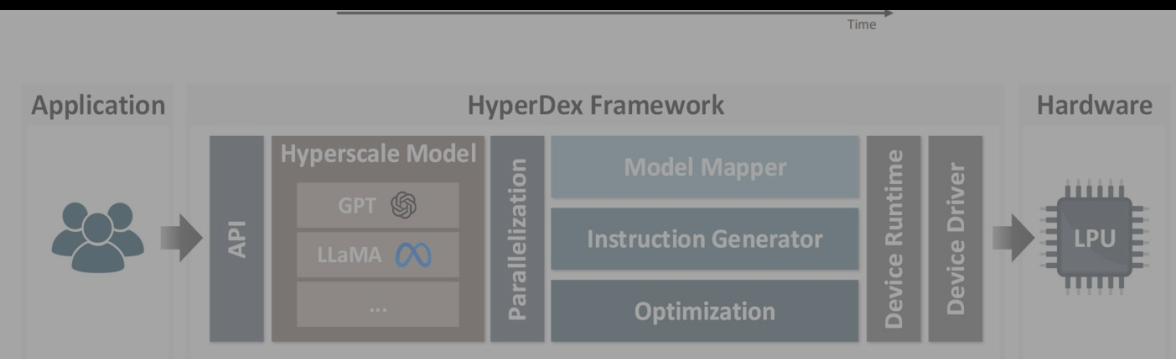
Non-Overlapping vs. Overlapping Synchronization



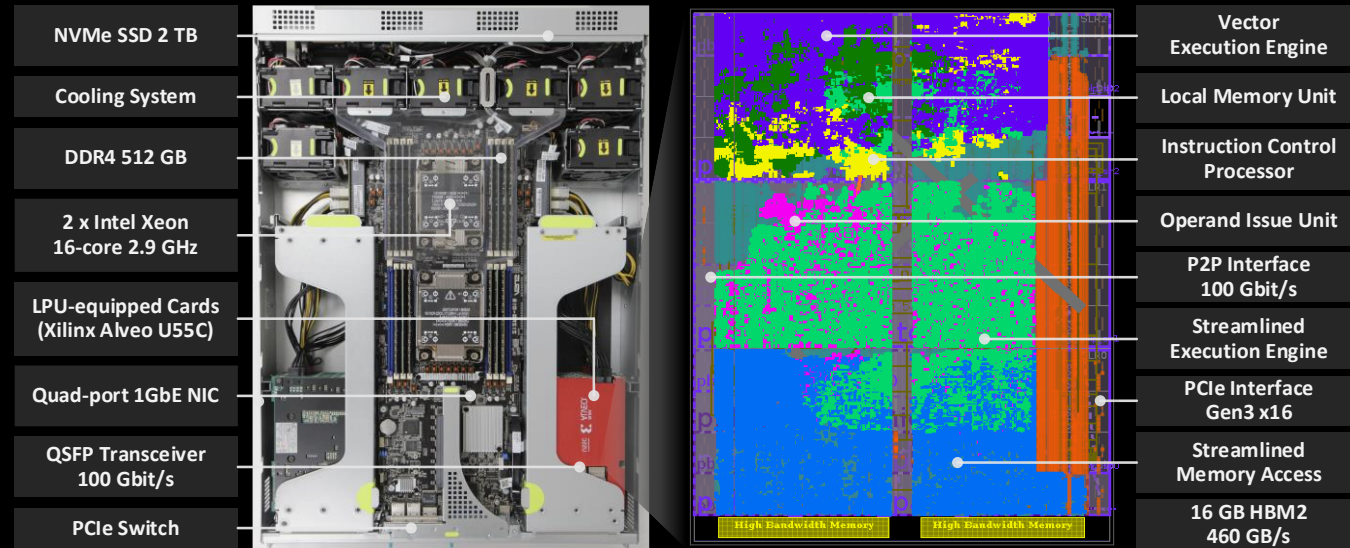
DETAILS AT OUR POSTER SESSION

HyperDex Software Stack

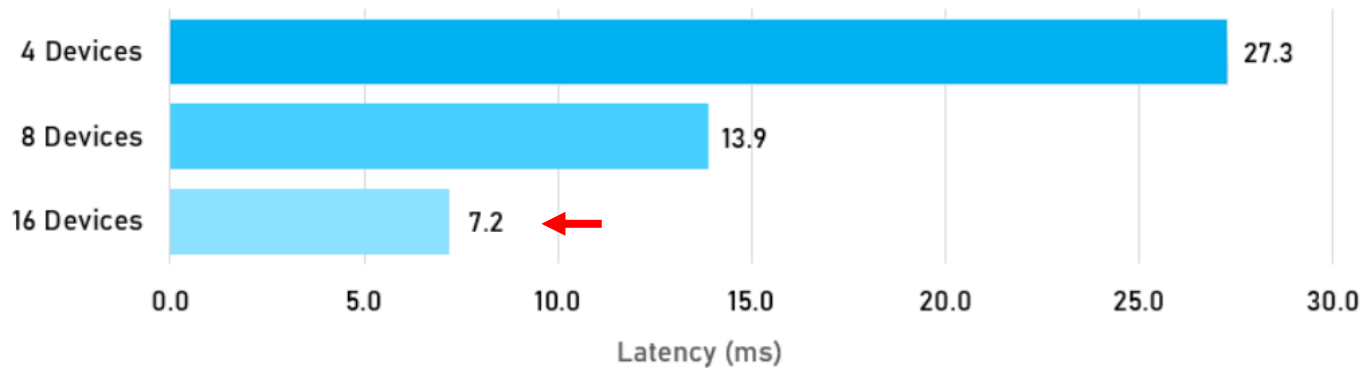
- Software framework that supports various language model structures and user requests
- Runs automated tools to create prerequisite data (e.g., parallelization, mapping, and instructions) to run the LPU at the application level



Performance



Average Latency per Output Token



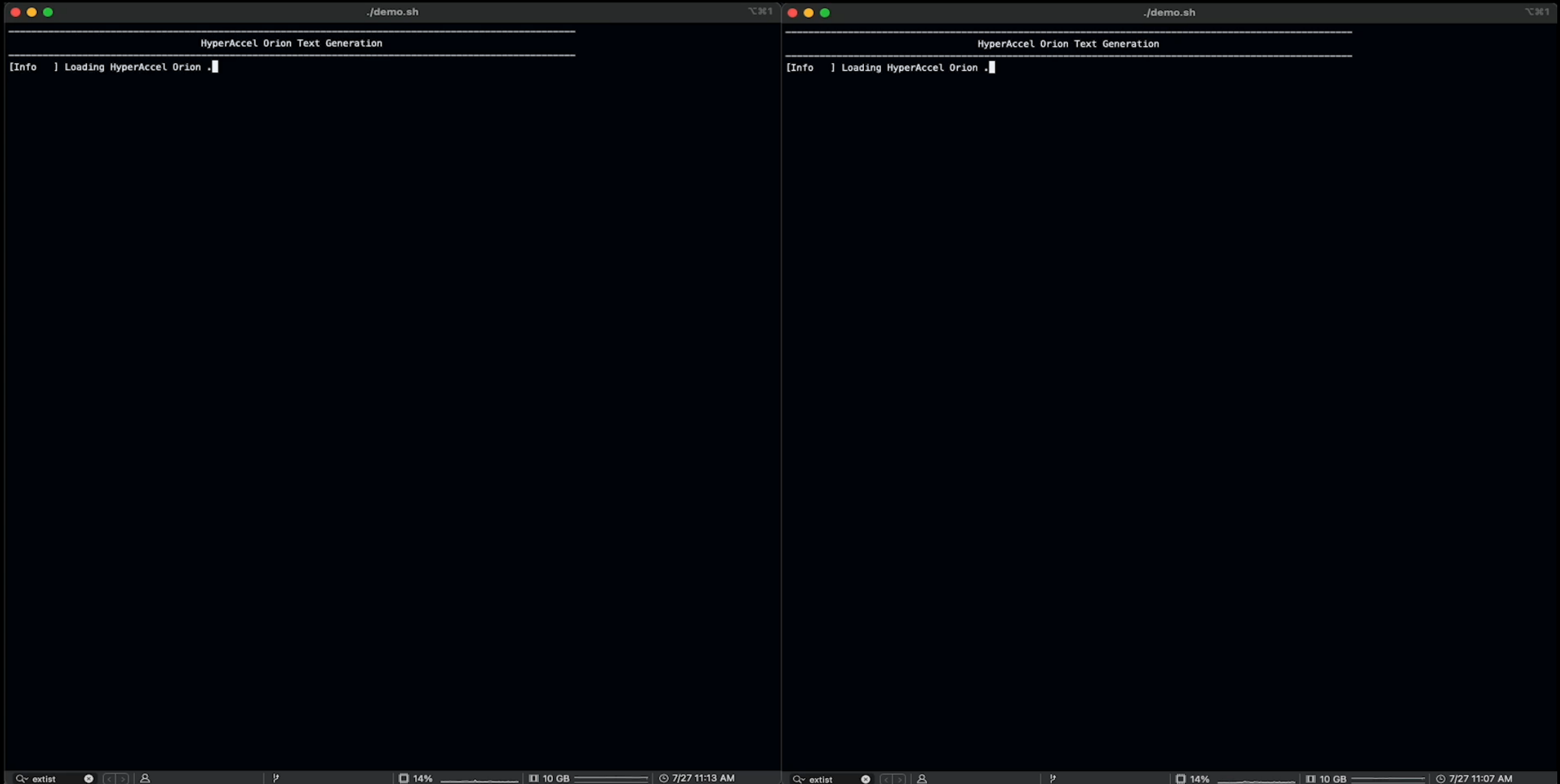
Model = OpenAI GPT3-20B, Input token size = 32, Precision = FP16

**140 TOKENS PER SECOND
FOR TEXT GENERATION WORKLOAD**

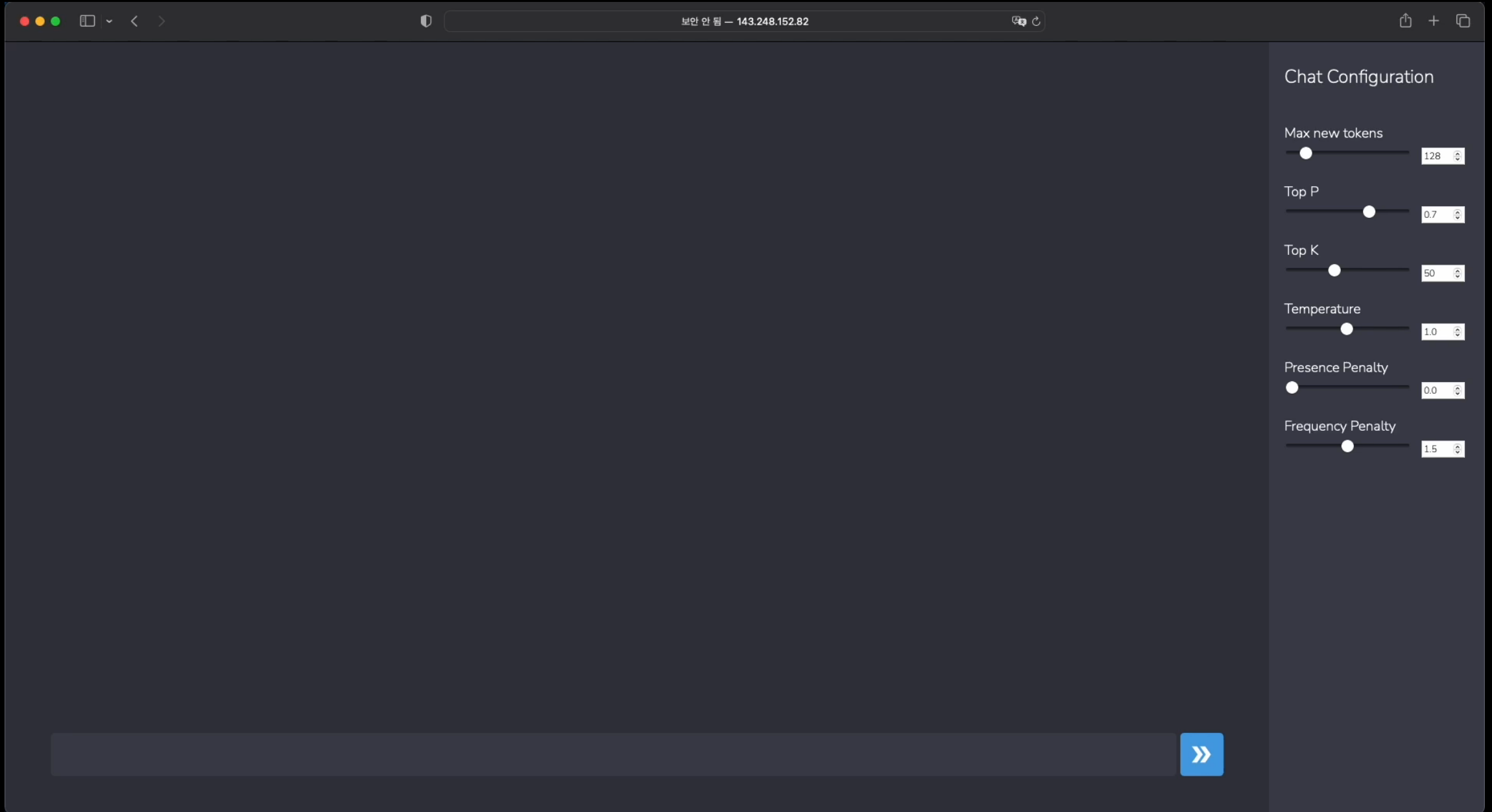
1 LPU vs. 8 LPUs on Meta AI LLaMA2 7B Model

1 LPU (49.6sec)

8 LPUs (8.7sec)



Chatbot UI





HyperAccel, a startup that provides hyper-accelerated silicon IP and solutions for emerging applications

Currently supported models (up to ~100B parameters)

- OpenAI GPT 1, 2, 3
- Meta AI LLaMA 1, 2
- Meta AI OPT
- Stanford Alpaca
- EleutherAI Polyglot
- EleutherAI GPT Neo-X

HyperAccel Latency Processing Unit (LPU™)

Accelerating Hyperscale Models for Generative AI

Seungjae Moon, Junsoo Kim, Jung-Hoon Kim, Junseo Cha, Gyubin Choi, Seongmin Hong, and Joo-Young Kim
HyperAccel, Hwasong-6, Republic of Korea

Introduction

- The fundamental goal of AI is to create human-like intelligence. Generative AI has enabled AI to do what we thought was innate to only humans: show creativity.
- Transformer-based large language models (LLM) with multi-billion parameters, such as OpenAI GPT, Meta LLaMA, can create original texts and visual contents.
- For efficient model inference, a latency-oriented and scalable hardware for small-batch memory-intensive workloads is required to meet the needs of different users.
- Latency Processing Unit, the world-first hardware accelerator dedicated for the end-to-end inference of LLM.

IP Products

- Highly flexible to reconfigure both memory types and compute resources for low-power and high-performance (baseline: GDDR with 16 lanes x 64 vector dimension MAC trees in SXE)
- Low-power: scale down memory bandwidth to that of LPDDR with fewer MAC trees in SXE
- High-performance: scale up memory bandwidth to that of HBM with more MAC trees in SXE

LPU™ Architecture

- Connects all channels of high bandwidth memory to the execution engines with datapath that exactly matches the floating bandwidth.
- Utilizes hardware-aware memory mapping and tiling to remove the need for any data reshaping and switching.
- Consists of low-latency and high-throughput custom multiply-accumulate (MAC) trees, multi-precision arithmetic function unit, and special function units.
- Out-of-order scheduling to allow simultaneous execution of independent matrix and vector operations for maximum hardware utilization.
- Achieves effective bandwidth usage of 80% during end-to-end LLM inference.

HyperDex Software Stack

- Bridges LPU platform at the application-level through standard API
- Supports various LLMs, such as GPT, OPT, LLaMA, and their variants
- Intra-layer parallelism of model parameters for parallelizable operations
- Optimal memory allocation and alignment of model parameters
- Parallel instruction chaining for minimum control overhead

Expandable Synchronization Link (ESL)

- Lightweight full-duplex peer-to-peer (P2P) communication technology that performs data synchronization with low latency and latency hiding
- Low-latency by minimal packet overhead, direct path I/O, and short dataflow
- Latency-hiding by custom protocol that enables execution and data synchronization to continuously run in tandem to hide all sync overhead except the tail-latency

Performance Results

Strong Scaling of HyperAccel Orion vs. NVIDIA DGX A100

1. This graph is an average for doubling the number of devices.

Device	Throughput (tokens/sec)	Power (W)	Efficiency (tokens/W)
HyperAccel Orion (16 LPUs)	~1.8k	~1.5k	~1.2
NVIDIA DGX A100 (8 GPUs)	~1.2k	~2.5k	~0.5

Feel free to contact us! | Email: contact@hyperaccel.ai | Website: hyperaccel.ai | LinkedIn: [linkedin.com/company/hyperaccel](https://www.linkedin.com/company/hyperaccel)

For more information, please visit our poster session!

Website: www.hyperaccel.ai

Email : sj.moon@hyperaccel.ai | contact@hyperaccel.ai



SiMa^{ai}TM

MLSoCTM - An overview

Hot Chips 35, August 28-29, 2023

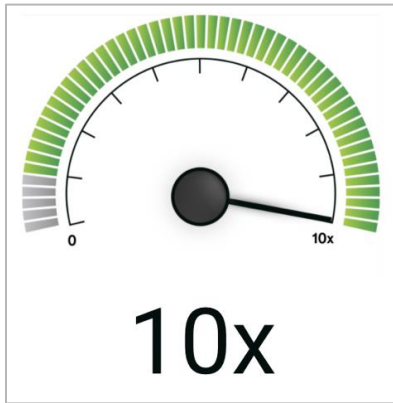
Srivi Dhruvanarayan, Victor Bittorf

Our Vision: Effortless machine learning for the embedded edge

Run **any** computer vision application, **any** network, **any** model, **any** framework, **any** sensor, **any** resolution.

Wide variety of end2end CV applications; Accuracy and common operator support

Any



Application

MLSoC: Run entire CV application efficiently

Efficient ML handling: Tile architecture.

Scheduling: Operations efficiently pipelined, scheduled

Performance

Data handling: On chip dec/enc, compliant with AVC, HEVC

Pre/Post processing: Dedicated CVU

Efficient static scheduling: Maximize compute while minimizing data movement;

Low power

Quantization scheme: Patented accurate low-power scheme

Fully INT8 inference: 100% of total compute dedicated

Patented cache usage: Manage memory hierarchy



Low code productization; Ready to use models, production ready platforms and dev kits

Pushbutton

SiMa.ai key innovations

ANY and 10x

Highly flexible ML accelerator

Secure, self-contained SoC

Pushbutton

Simple to develop & deploy

Fully programmable MLA



Supports full complement of CV applications at lowest power

SW controlled data movement, scheduling & synchronization



SW control of cache/memory hierarchy, data movement: minimal data movement, small cache, high compute efficiency and lowest power

Seamless heterogeneous compute



Enables ML capabilities for legacy apps; future proofs applications

Optimized end-to-end SoC pipeline



Highly optimized building blocks enable best in class end-to-end performance

Effortless customer integration



Low code ML
Customers can develop & deploy applications without having to understand details of HW

Low code customer evaluation



Vision Development Platform (VDP)
Allows customers to build apps without having to write code

Purpose built for ML edge at embedded edge

Acquire ANY Data

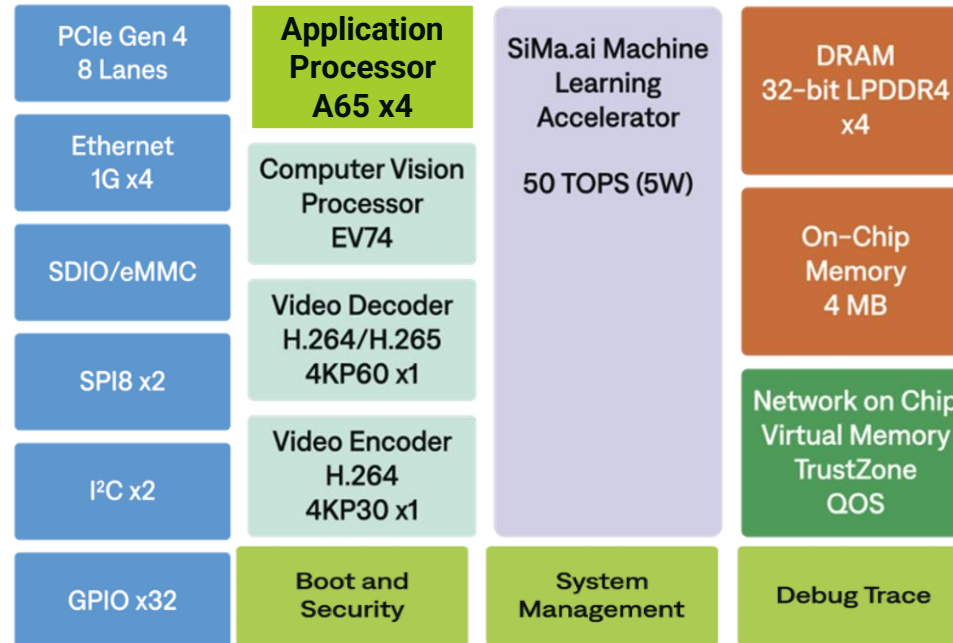
PCIe Gen4, 1GB Ethernet
Dedicate: I2C, SPI, GPIO, SDIO

10X ML Processing

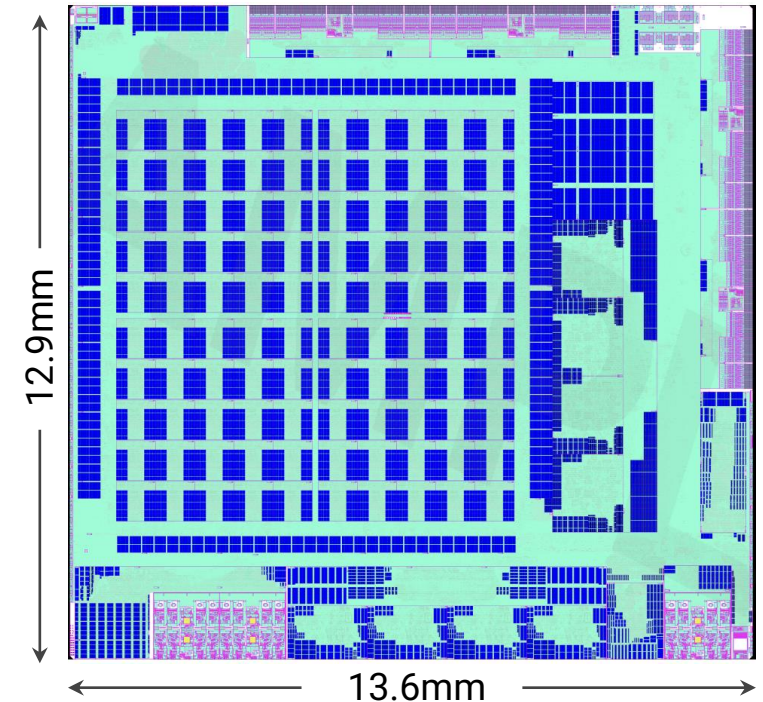
Video, CV & ML Processors
ML Accelerator - 50 TOPS INT8
Quad Vector DSP - 600 GOPS
HW video encode/decode
16GB LPDDR4

Decide, Control & Update

Quad A65E ARM8.3
Dedicated secure boot processor



TSMC - 16nm
TDP - 15-20W
Typical ML Workloads, CV Pipelines - 8-10W

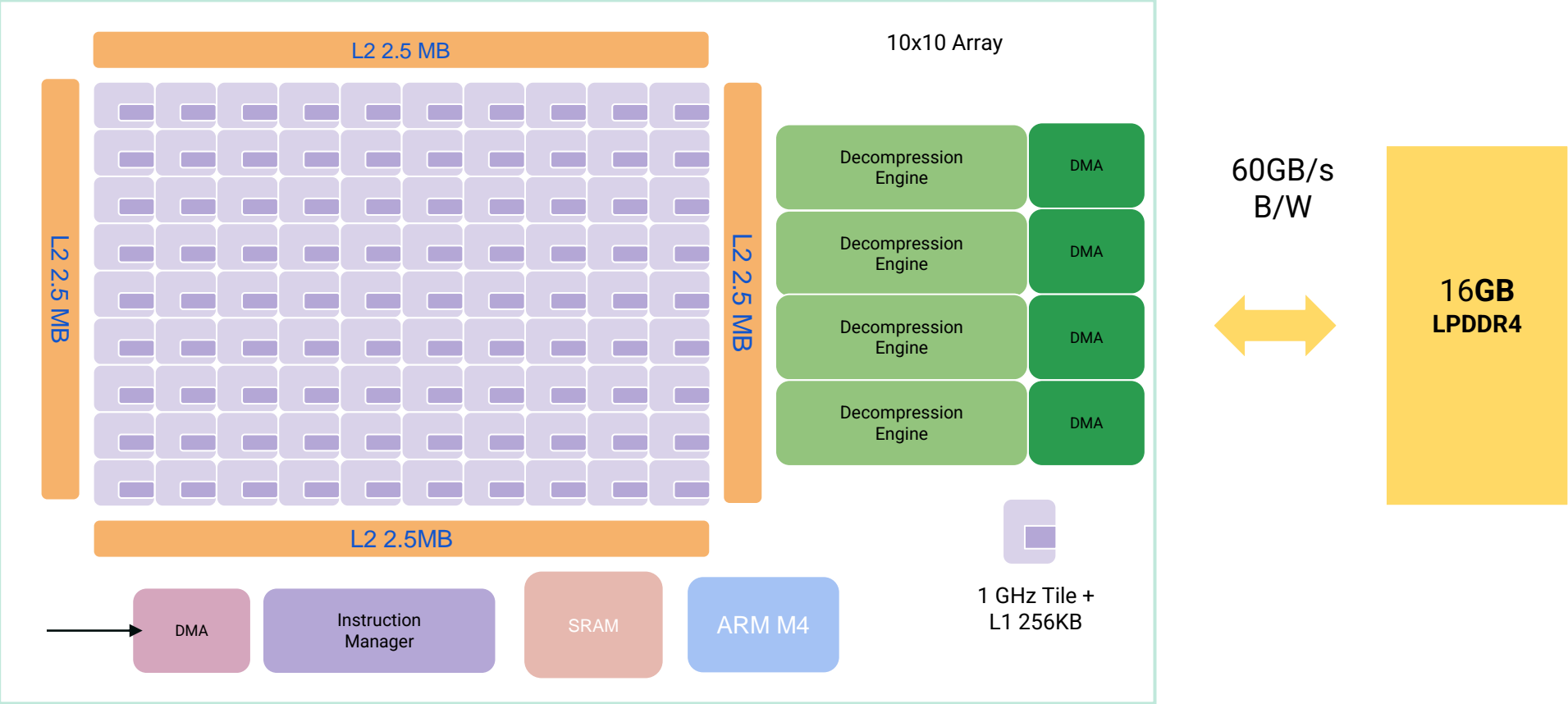


MLSoC™ Machine Learning System-on-Chip

Silicon Overview - 10x Performance for ML Processing

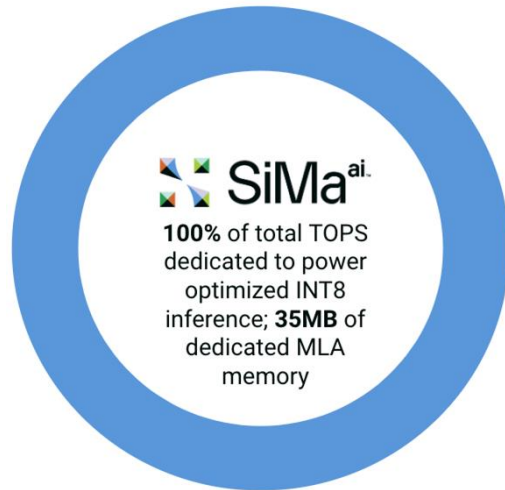
Machine Learning Accelerator

50 INT8 TOPS

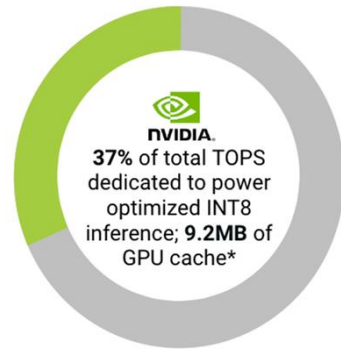


MLPerf: SiMa.ai delivers advantage over NVIDIA

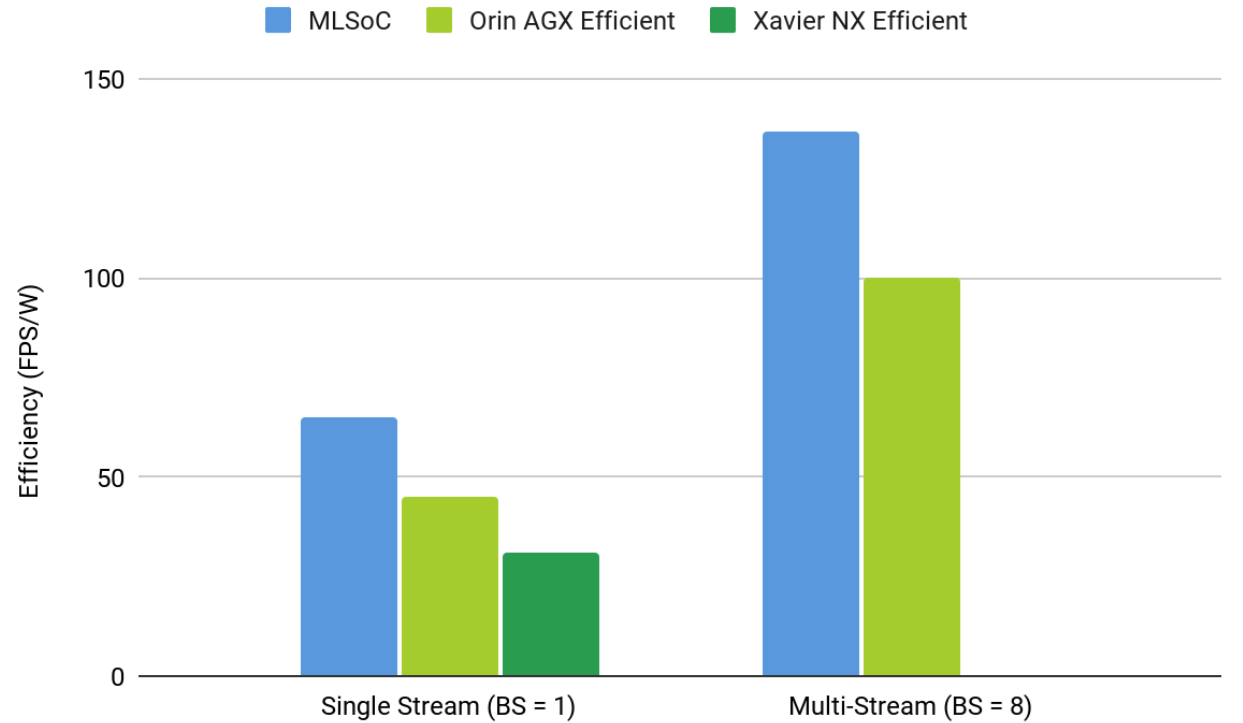
SiMa.ai MLSoC (N16) compiled results unseats Orin (8nm) **on both performance and power**



Purpose built
Low power + small area



General purpose (GPU)
High power + large area

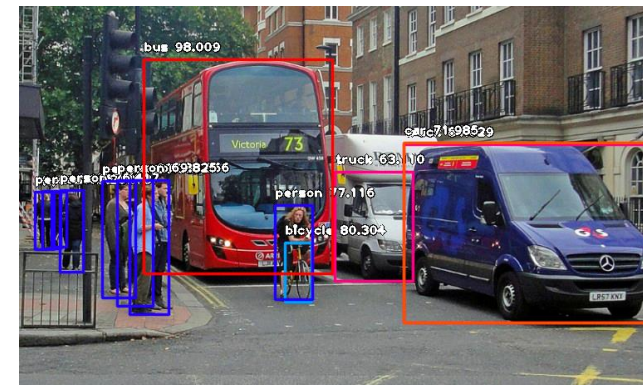
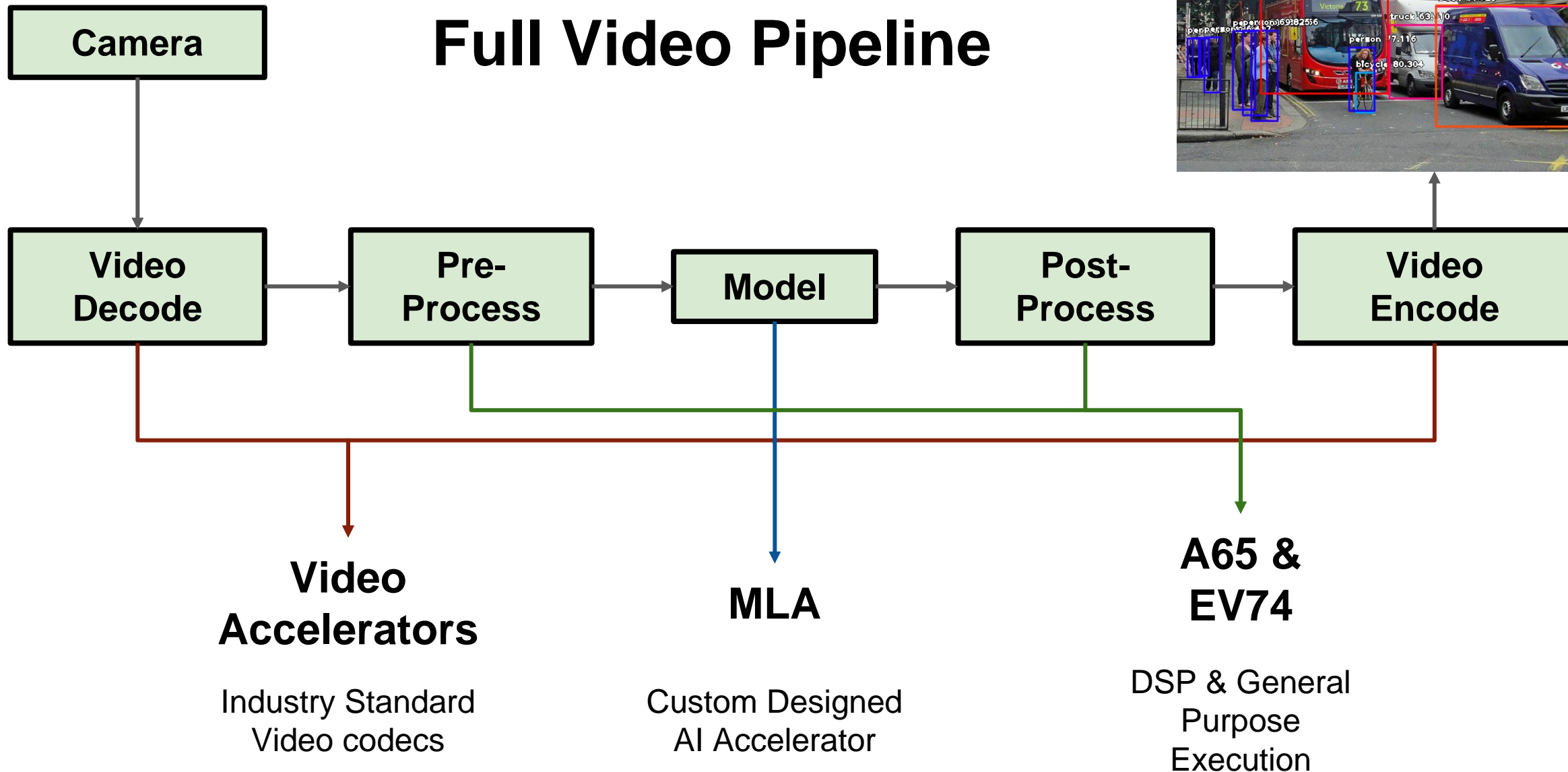


1 Camera
1.4x Orin
2.1x Xavier

8 Cameras
1.37x Orin
Xavier data not published

[Link](#) to MLPerf results

Full Video Pipeline





Shaheen: An Open, Secure, and Scalable RV64 SoC for Autonomous Nano-UAVs

University of Bologna

L. Valente, A. Veeran, M. Sinigaglia, Y. Tortorella, A. Nadalini, N. Wistoff, B. Sá, A. Garofalo, R. Psiakis, M. Tolba, A. Kulmala, N. Limaye, O. Sinanoglu, S. Pinto, D. Palossi, L. Benini, B. Mohammad, D. Rossi


luca.valente@unibo.it

7

PULP Platform

Open Source Hardware, the way it should be!



@pulp_platform 

pulp-platform.org 

youtube.com/pulp_platform 

Autonomous Nano-UAVs



- **Versatility, safety, and cost-effect:**

- small and agile
- ideal for accessing hard-to-reach areas or tight spaces (inspection/maintenance)
- relatively inexpensive to produce and operate

- **Requirements for future generation of nano-UAVs:**

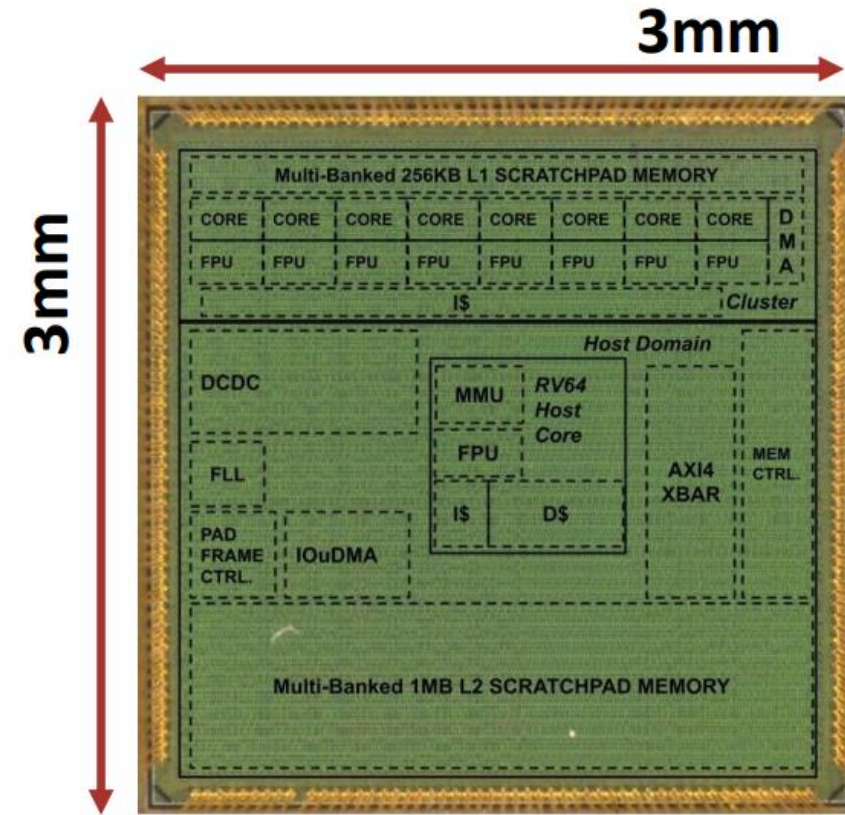
- Run increasingly complex multi-tasking workloads with large memory footprint
- Within a few hundred mW power budget
- Support for virtualization and secure operations in uncontrolled/hostile scenarios



Shaheen: an Open, Secure, and Scalable RV64 SoC for Autonomous Nano-UAVs



- **9mm² SoC in 22nm FDSOI technology with:**
 - A **RV64** Linux-capable CPU enhanced with
 - **Hypervisor** support
 - **Timing-channels mitigation**
 - An **energy efficient** programmable multi-core accelerator (PMCA) based on **8 RV32 cores with ML and DSP extensions**
 - Up to **512MB** of low-power off-chip **main memory**
 - **Logic locking** on key IPs within the architecture
 - **200mW** power envelope



RV64 and custom RV32: the best of both worlds



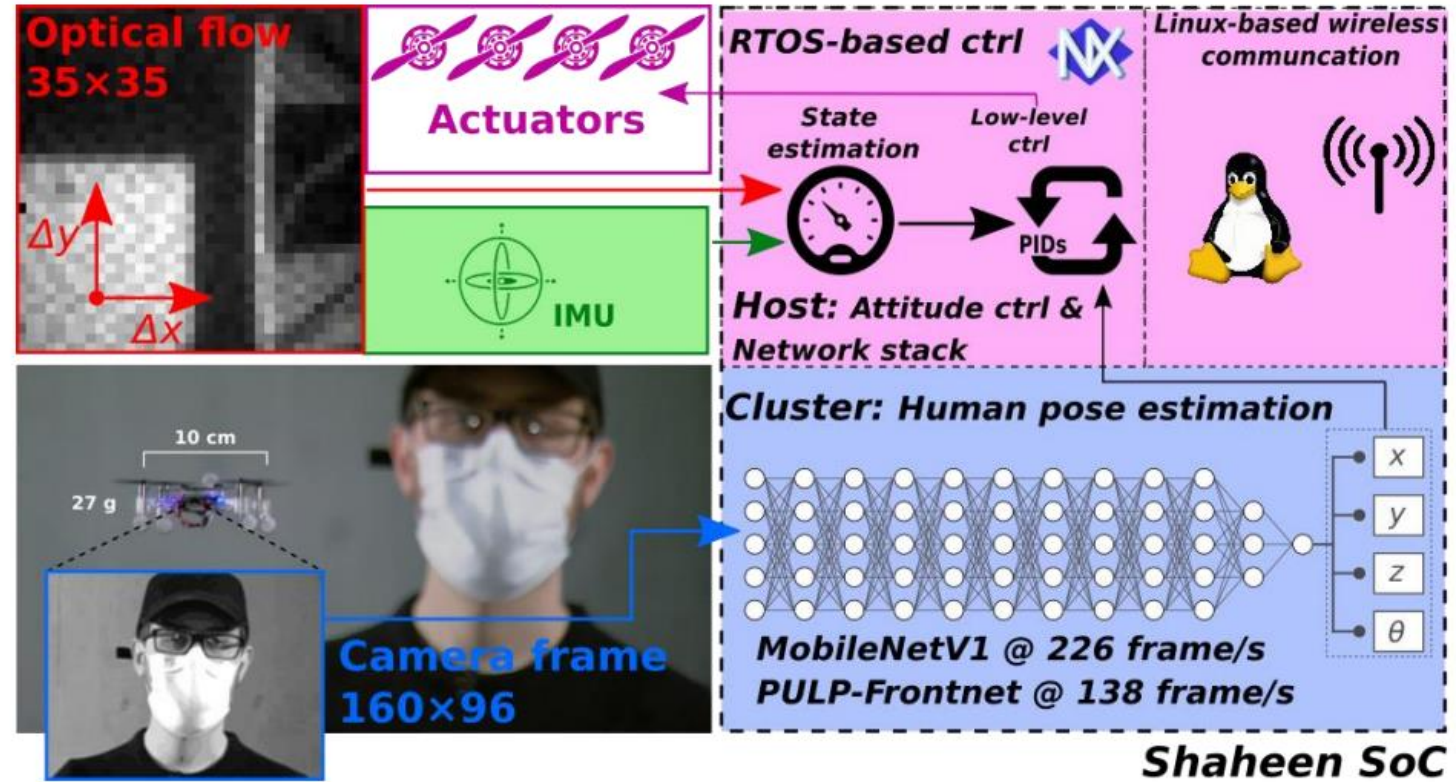
- **Host:**

- On top of the **hypervisor:**

- Attitude control (**RTOS-based**)
- **Linux-based** legacy software such as wireless network stack.

- **PMCA:**

- The PMCA runs the **CNN-based** pose estimation task fed by a low-resolution front-looking camera.



Luca Valente luca.valente@unibo.it



Institut für Integrierte Systeme – ETH Zürich

Gloriastrasse 35
Zürich, Switzerland

DEI – Università di Bologna

Viale del Risorgimento 2
Bologna, Italy

@pulp_platform



pulp-platform.org



youtube.com/pulp_platform



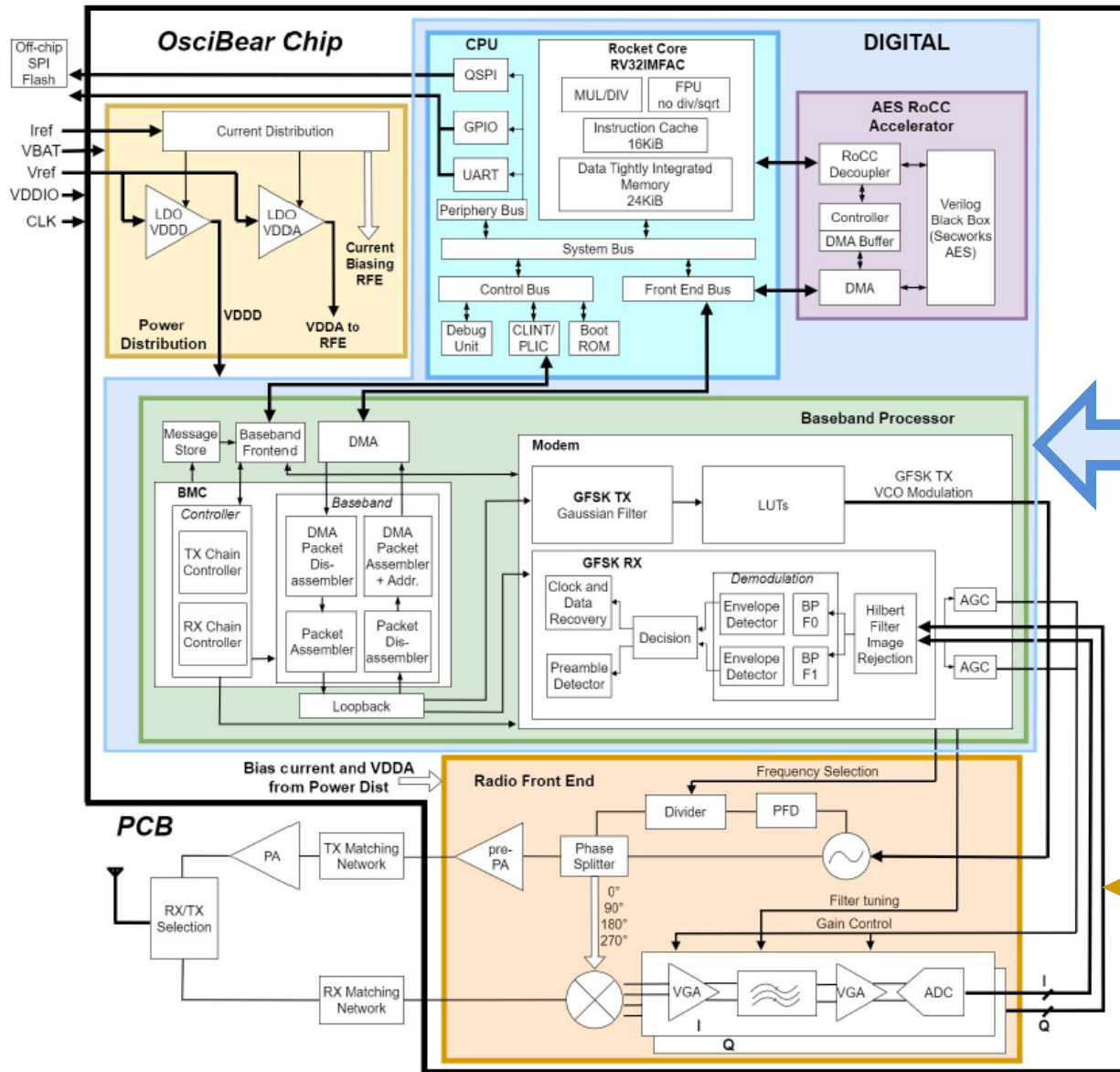


Berkeley
Wireless Research Center

A Heterogeneous SoC for Bluetooth LE in 28nm

Felicia Guo, Nayiri Krzysztofowicz, Alex Moreno, Jeffrey Ni, Daniel Lovell, Yufeng Chi, Kareem Ahmad, Sherwin Afshar, Josh Alexander, Dylan Brater, Cheng Cao, Daniel Fan, Ryan Lund, Jackson Paddock, Griffin Prechter, Troy Sheldon, Shreesha Sreedhara, Anson Tsai, Eric Wu, Kerry Yu, Daniel Fritchman, Aviral Pandey, Ali Niknejad, Kristofer Pister, and Borivoje Nikolic

Chip in a Semester

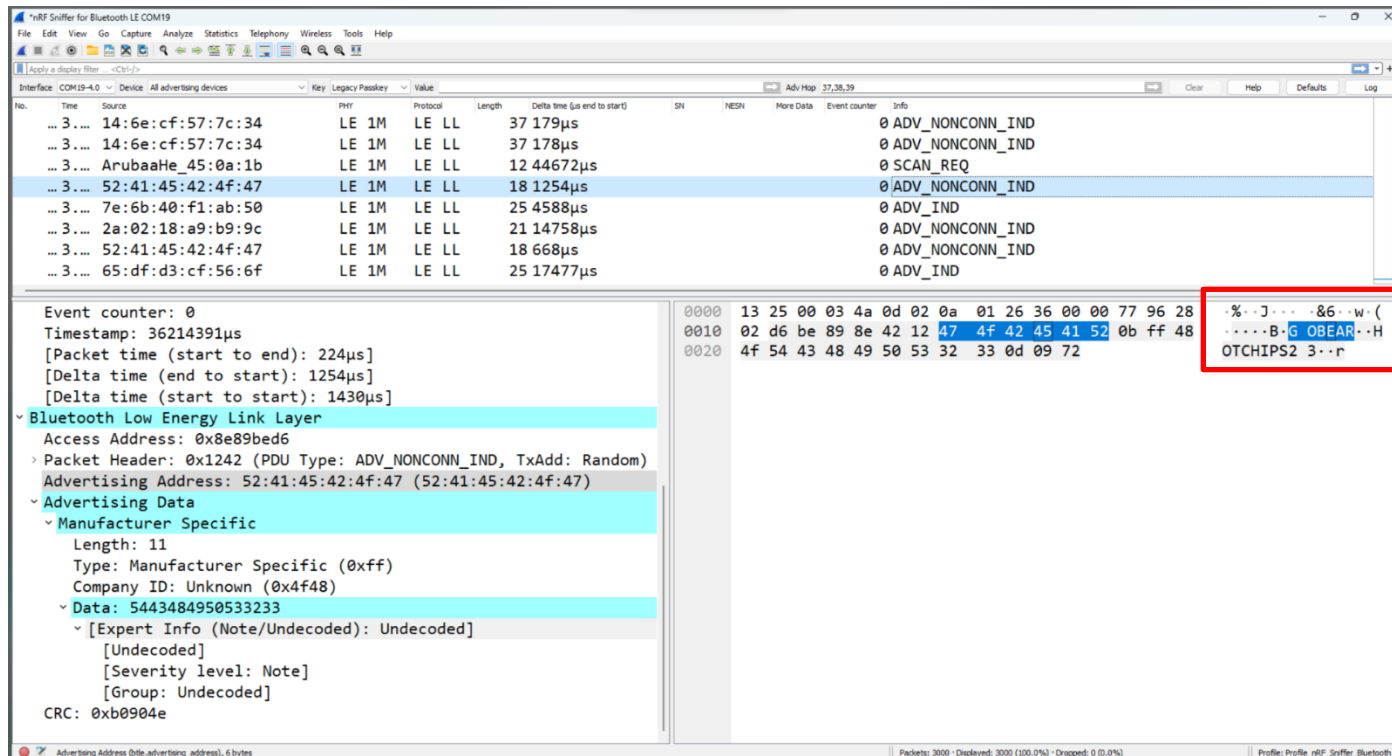


- 19 students in 14 weeks
- 1 mm^2 chip in TSMC 28nm
- Berkeley developed tools



BLE Measurement Results

- Transmit packet decoding done two ways:
 1. Software Testing
 2. Commercial (Nordic nRF52840 DK) receiver
- Receiver tested from PCB antenna port to differential output of final VGA (pre-ADC)



The screenshot shows the nRF Sniffer interface with a captured BLE packet. The packet details are as follows:

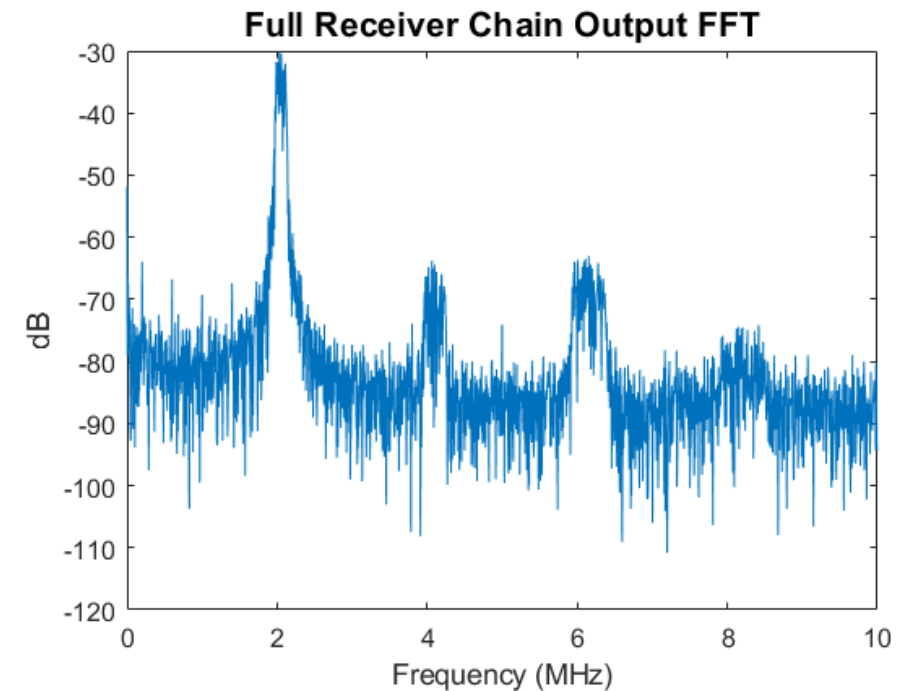
No.	Time	Source	PHY	Protocol	Length	Delta time (s to start)	SN	NESN	More Data	Event counter	Info
...	...	14:6e:cf:57:7c:34	LE 1M	LE LL	37 179μs						0 ADV_NONCONN_IND
...	...	14:6e:cf:57:7c:34	LE 1M	LE LL	37 178μs						0 ADV_NONCONN_IND
...	...	ArubaaHe_45:0a:1b	LE 1M	LE LL	12 44672μs						0 SCAN_REQ
...	...	52:41:45:42:4f:47	LE 1M	LE LL	18 1254μs						0 ADV_NONCONN_IND
...	...	7e:6b:40:f1:ab:50	LE 1M	LE LL	25 4588μs						0 ADV_IND
...	...	2a:02:18:a9:b9:9c	LE 1M	LE LL	21 14758μs						0 ADV_NONCONN_IND
...	...	52:41:45:42:4f:47	LE 1M	LE LL	18 668μs						0 ADV_NONCONN_IND
...	...	65:df:d3:cf:56:6f	LE 1M	LE LL	25 17477μs						0 ADV_IND

The decoded payload for the selected packet is shown below:

```

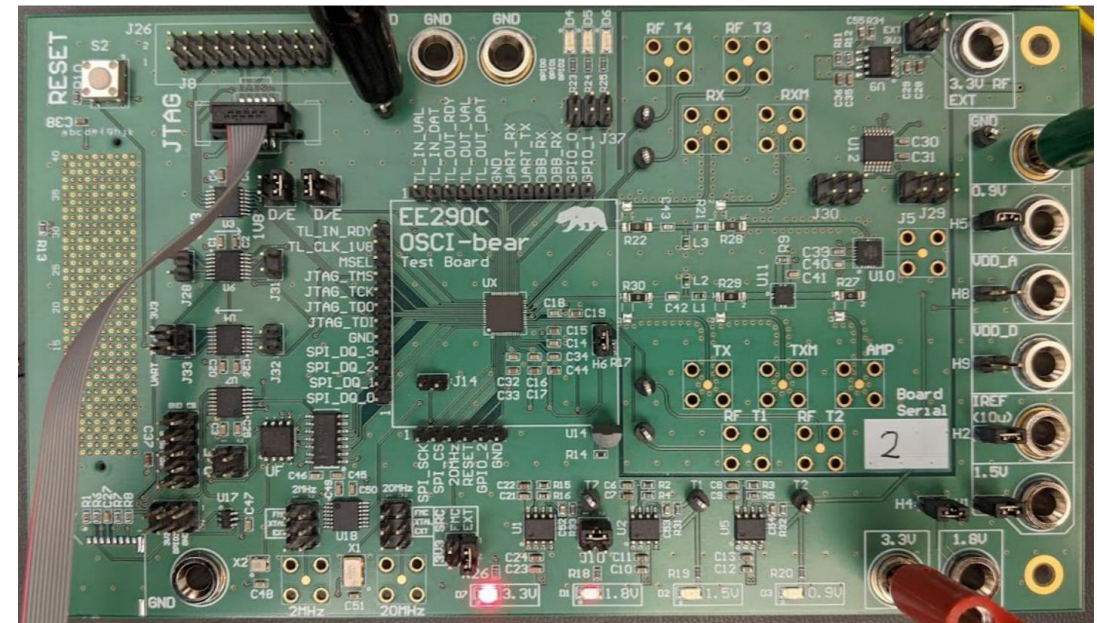
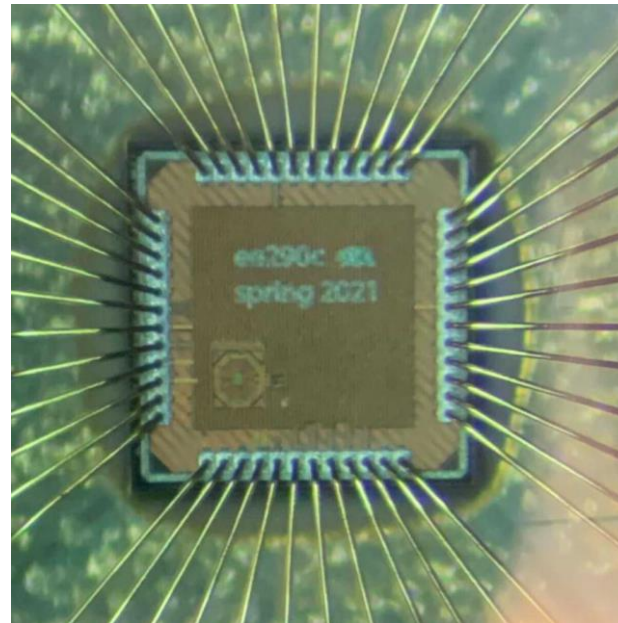
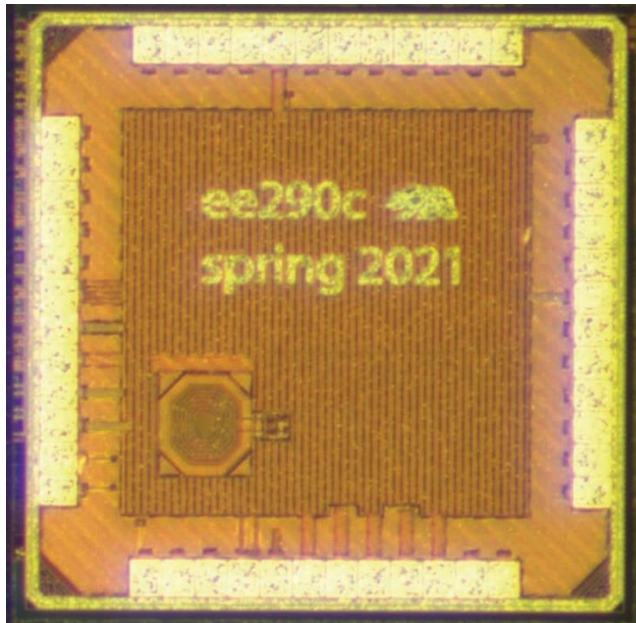
Event counter: 0
Timestamp: 36214391μs
[Packet time (start to end): 224μs]
[Delta time (end to start): 1254μs]
[Delta time (start to start): 1430μs]
Bluetooth Low Energy Link Layer
  Access Address: 0x8e89bed6
  Packet Header: 0x1242 (PDU Type: ADV_NONCONN_IND, TxAdd: Random)
  Advertising Address: 52:41:45:42:4f:47 (52:41:45:42:4f:47)
  Advertising Data
    Manufacturer Specific
      Length: 11
      Type: Manufacturer Specific (0xff)
      Company ID: Unknown (0x4f48)
      Data: 5443484950533233
        [Expert Info (Note/Undecoded): Undecoded]
          [Undecoded]
          [Severity level: Note]
          [Group: Undecoded]
    CRC: 0xb0904e
  
```

The hex payload is: 0000 13 25 00 03 4a 0d 02 0a 01 26 36 00 00 77 96 28 47 4f 42 45 41 52 0b ff 48. The ASCII representation is: .% . . J &6 . . w (. B : G OBEAR . . H OTCHIPS2 3 . . r

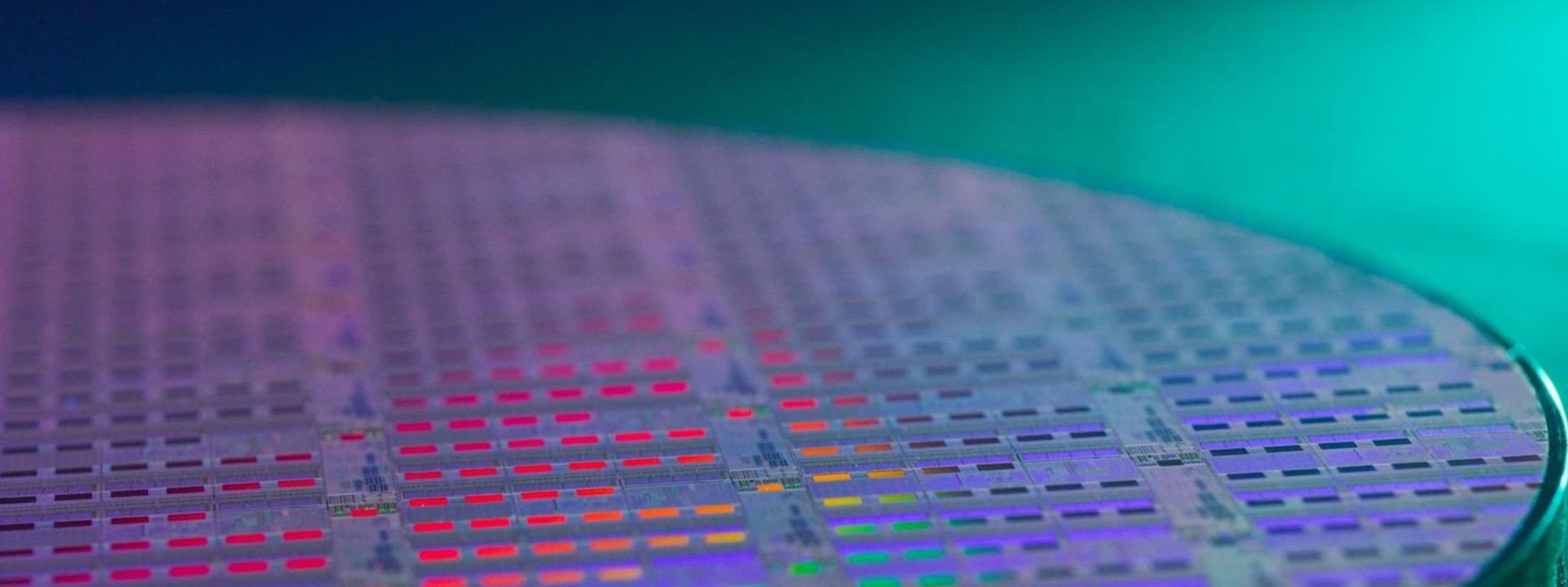


Commercial receiver decoding correctly formed packet sent with "GOBEAR HOTCHIPS23" as payload

Thank You!



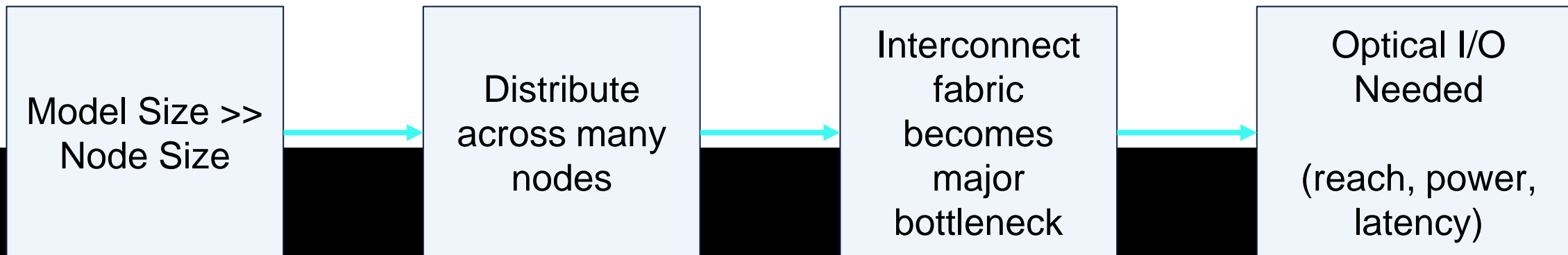
We would like to acknowledge and thank Apple for supporting integrated circuit engineering classes at UC Berkeley EECS department through Apple's New Silicon Initiative program. This work is also funded by NSF CCRI ENS Chipyard Award #2016662.



Driving Compute Scale-out Performance with Optical I/O Chiplets in Advanced System-in-Package Platforms

Mark Wade, PhD | President, CTO, Co-Founder | August 28, 2023

Hundreds/Thousands of Sockets Computing as One



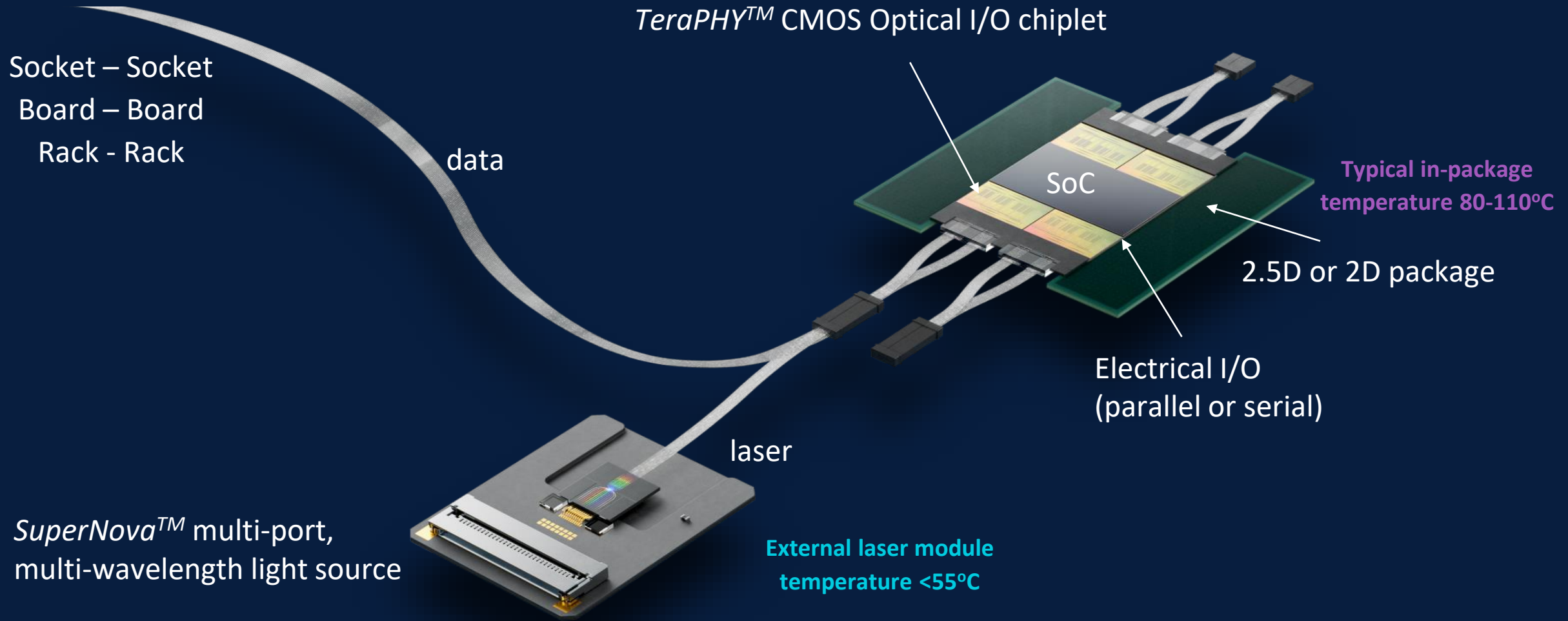
15 meters (50 feet)



DGX GH200 SuperPOD (Source: Nvidia)

Nvidia DGX GH200: 256 GPUs acting as 1 “Mega GPU”

Ayar Labs Optical I/O (OIO) Chiplet

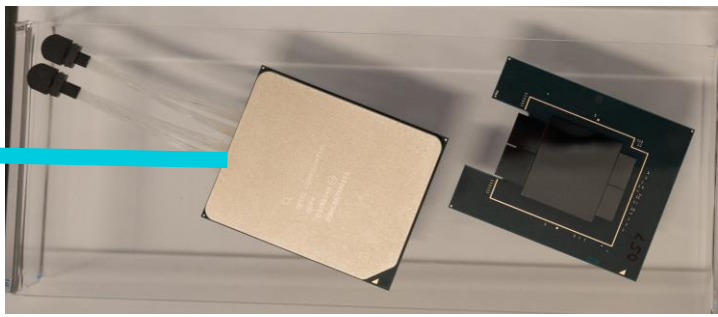
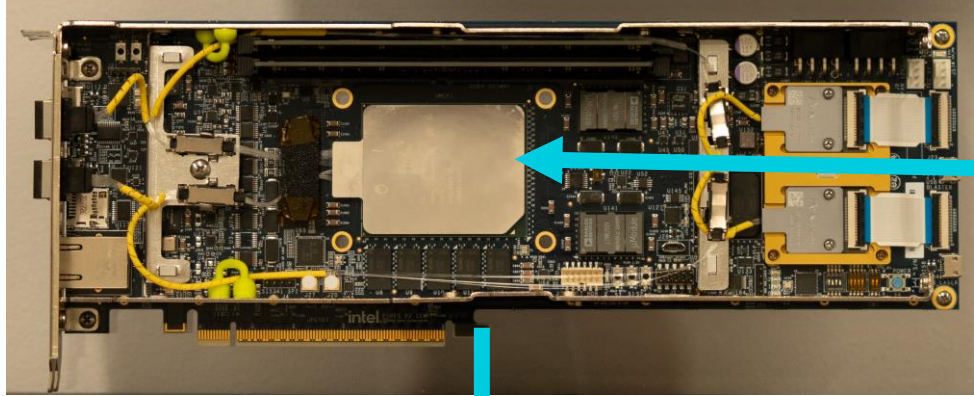


Optical FPGA PCIe Card

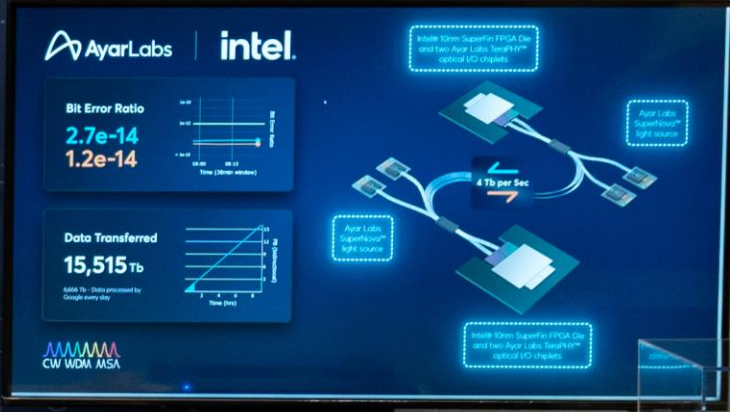
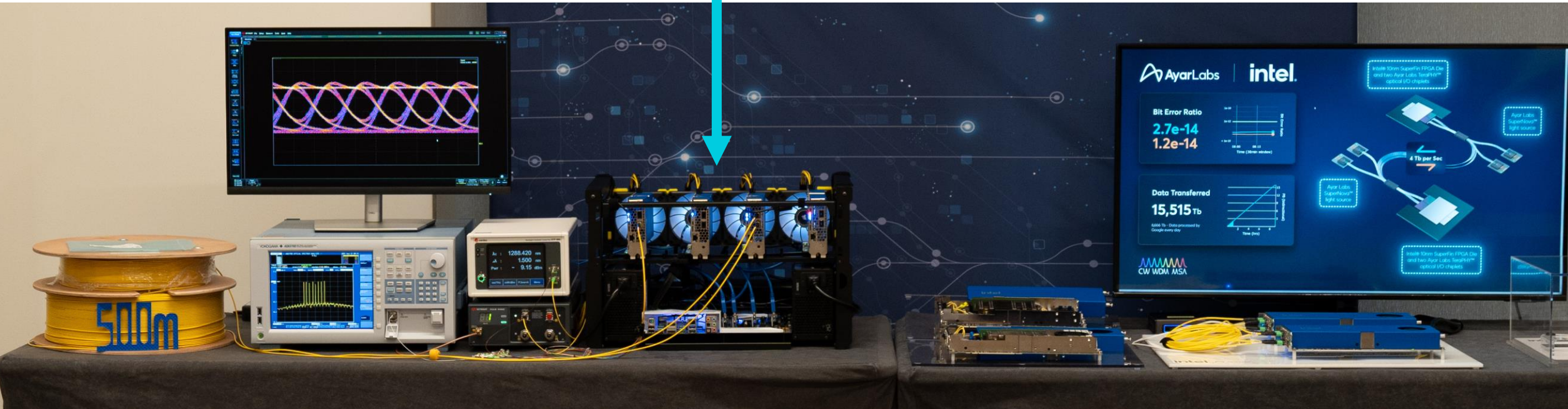
2 x 4Tbps Optical I/O Chiplets

<1e-12 BER pre-FEC

10ns + TOF latency

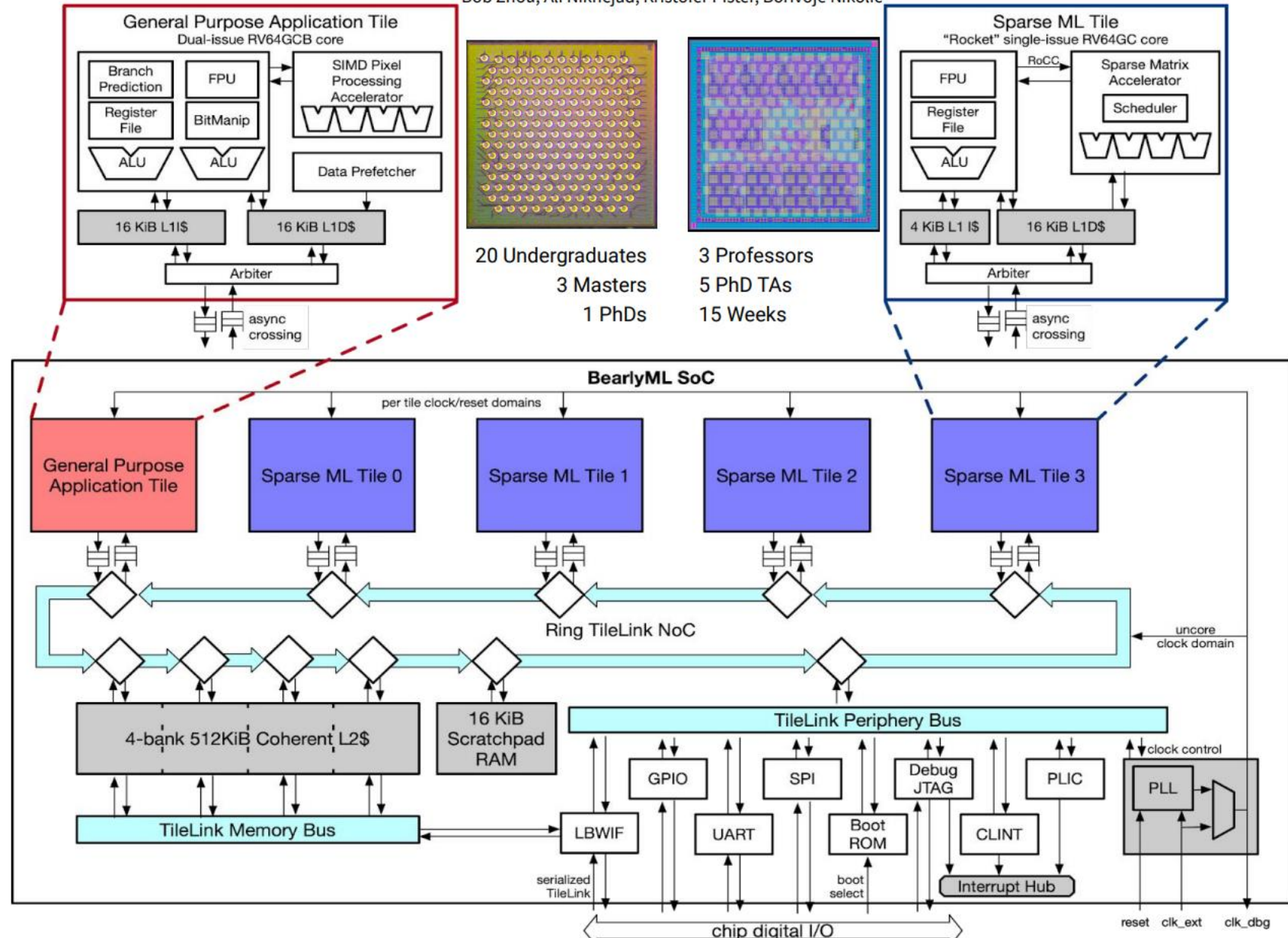


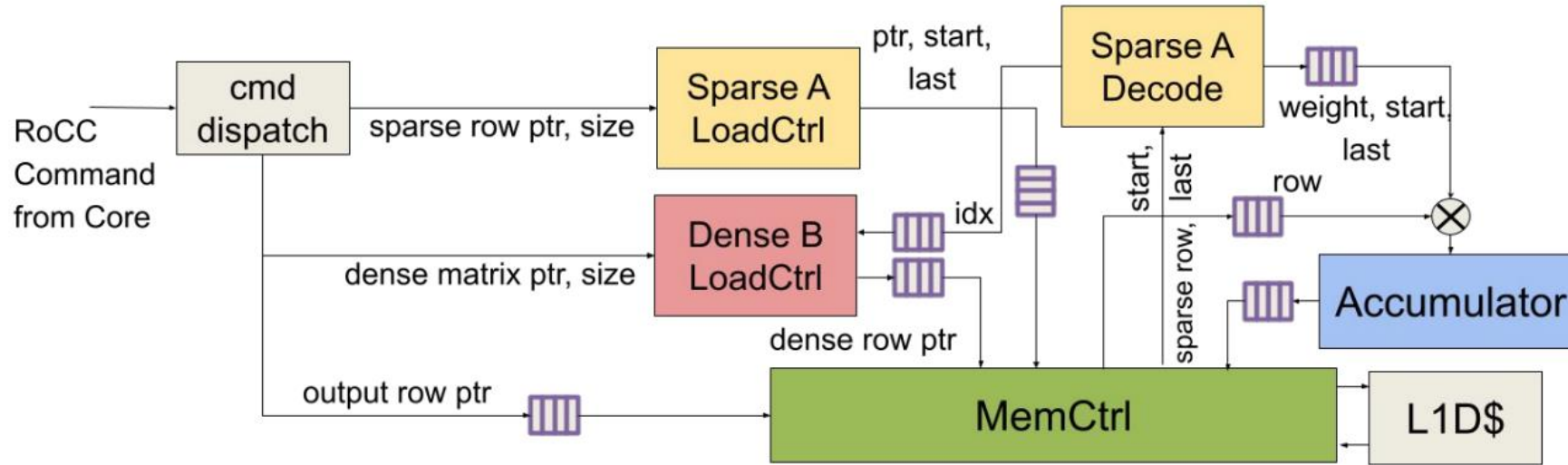
Running Live Outside



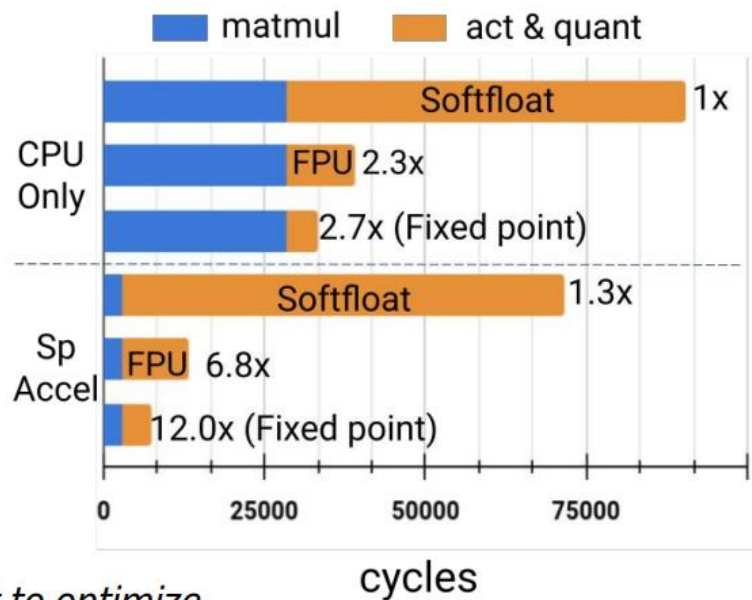
Optical I/O Enables Fully Disaggregated Distributed Compute

Yufeng Chi, Franklin Huang, Raghav Gupta, Ella Schwarz, Jennifer Zhou, Reza Sajadiany, Animesh Agrawal, Max Banister, Michelle Boulos, Jason Chandran, Jessica Dowdall, Leena Elzeiny, Claire Gantan, Anthony Han, Roger Hsiao, Chadwick Leung, Edwin Lim, Jose Rodriguez, Tushar Sondhi, Mitchell Twu, Rongyi Wang, Mike Xiao, Ruohan Yan, Paul Kwon, Zhaokai Liu, Jerry Zhao, Bob Zhou, Ali Niknejad, Kristofer Pister, Borivoje Nikolić

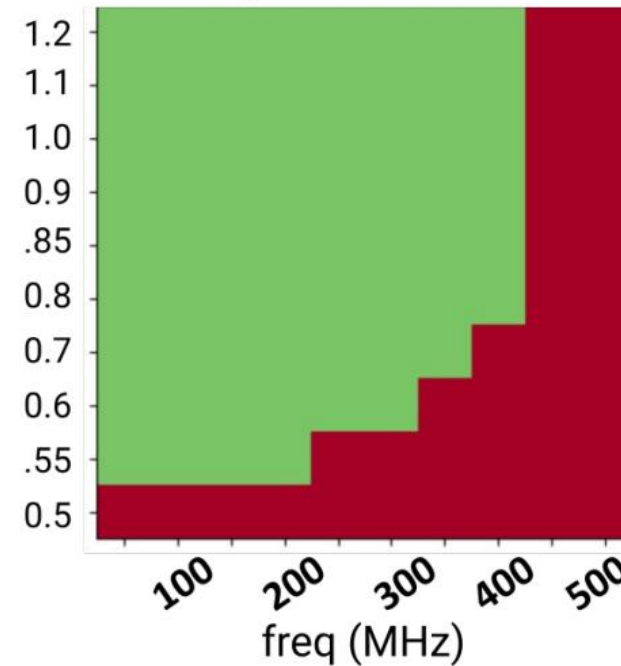




MNIST Cycle Breakdown



Volts BearlyML Shmoo Plot



don't forget to optimize non-matmul operators!



End of Poster Lightning Talks