

# Scalable and Real-time Baseband Processing using Heterogeneous Compute Resources

**Tingjun Chen**

Nortel Networks Assistant Professor

Departments of Electrical & Computer Engineering and Computer Science

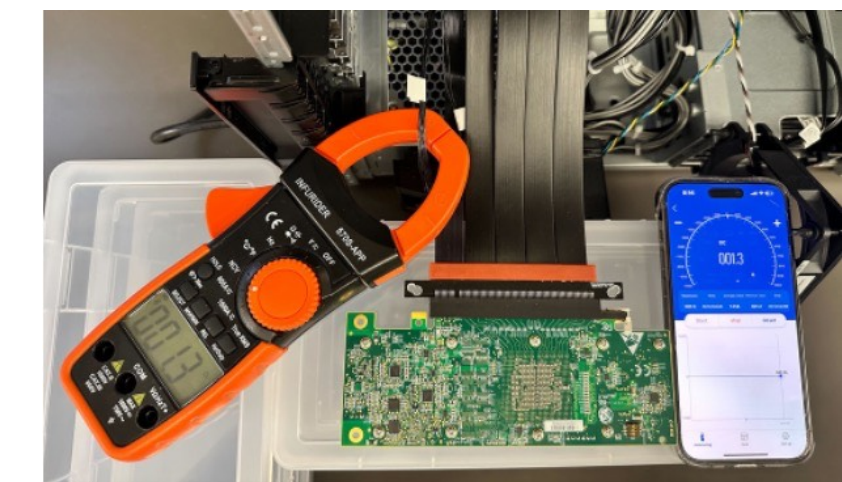
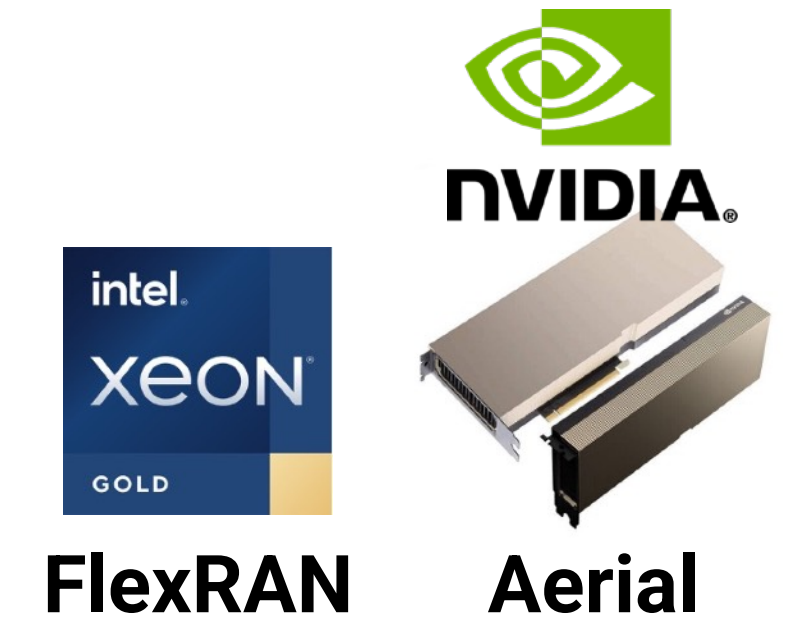
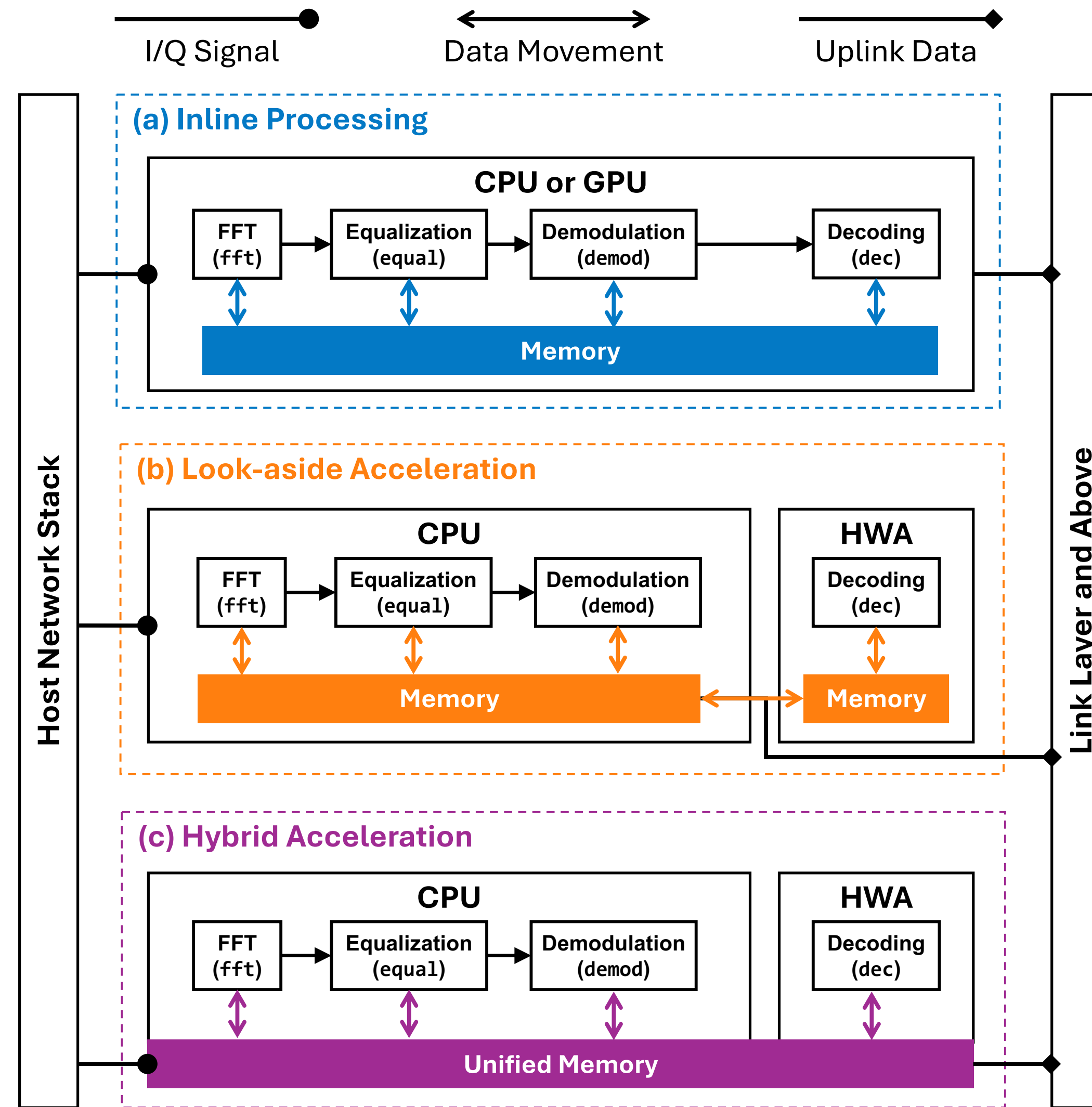
Duke University

First OAI Foundation U.S. Hands-on Workshop

November 17, 2025

# L1 (Baseband) Processing Models for vRANs

- Leveraging the powerful **hardware accelerators (HWAs)** for L1 processing in virtualized radio access networks (vRANs)
- **Cloud-native scaling** favors flexibility, software-centric upgrades, and multi-tenancy supports
- **Power-latency tradeoffs** determine the overall system efficiency, but also remains difficult to model
- Requires both significant R&D and engineering efforts!



Intel ACC100 eASIC

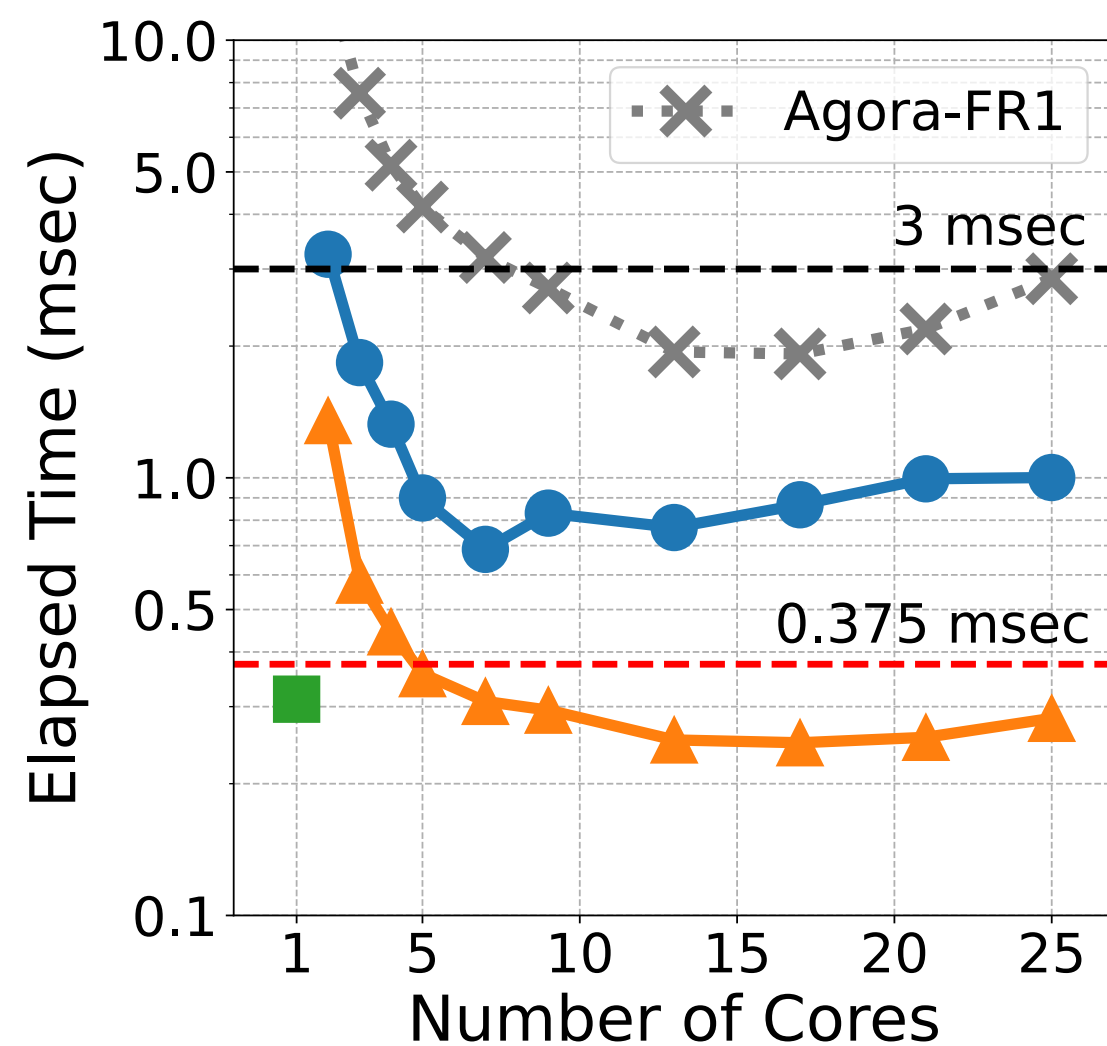


NVIDIA Jetson Orin AGX

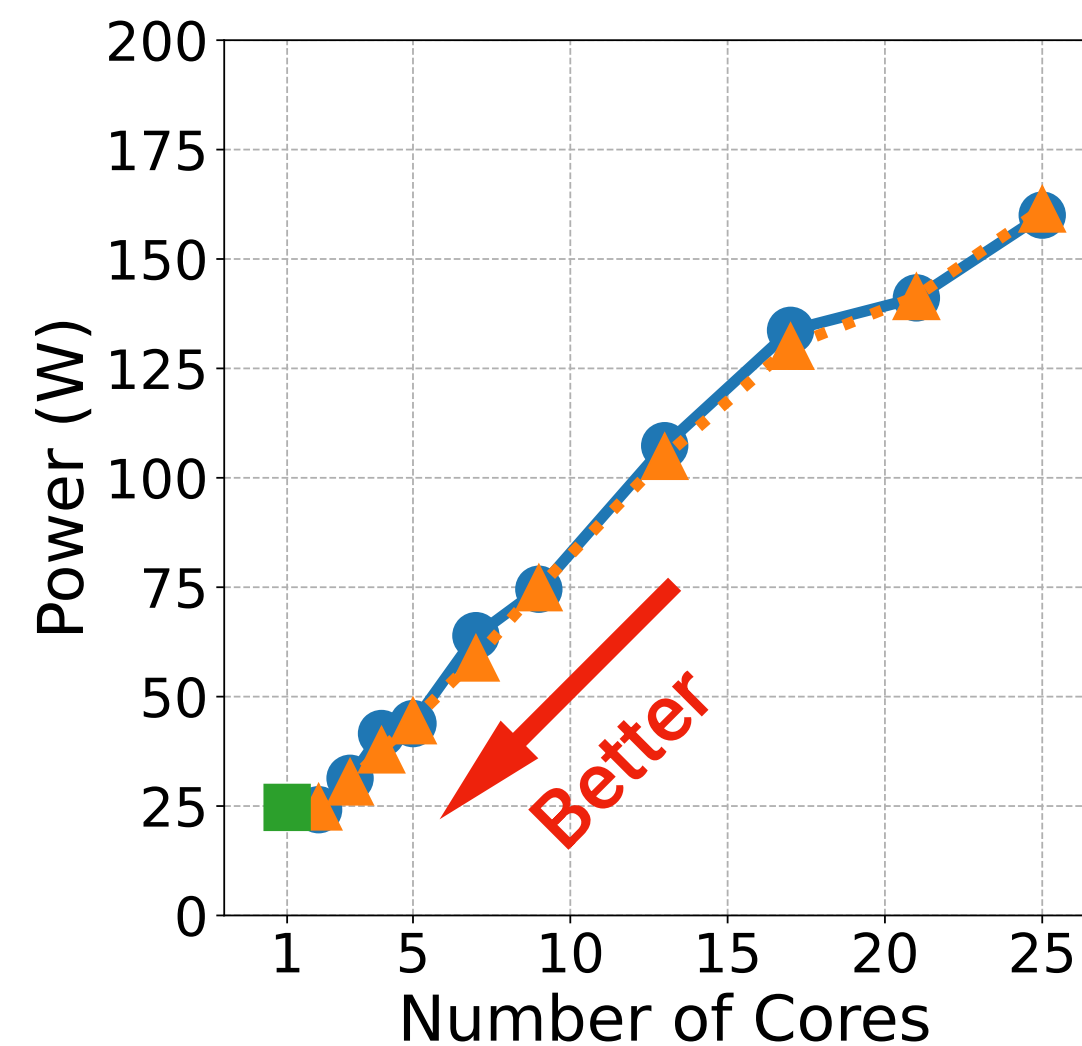
# Savannah: Real-time Efficient mmWave Baseband Processing

- Existing software-based PHY processing frameworks are designed for FR1, while FR2 has much tighter PHY processing latency requirements due to larger subcarrier spacing

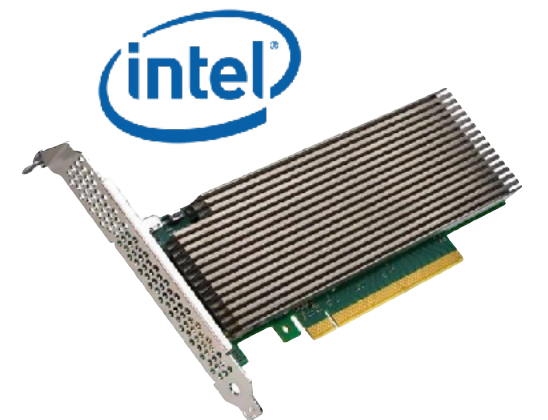
**99.9<sup>th</sup>-percentile Elapsed Time**



**Power Consumption**



Agora [ACM CoNEXT'20]	
✕	Agora-FR1 CPU only (multi-core)
●	Agora CPU only (multi-core)
▲	Savannah-mc CPU only (multi-core)
■	Savannah-sc CPU (single-core) + ASIC (ACC100)

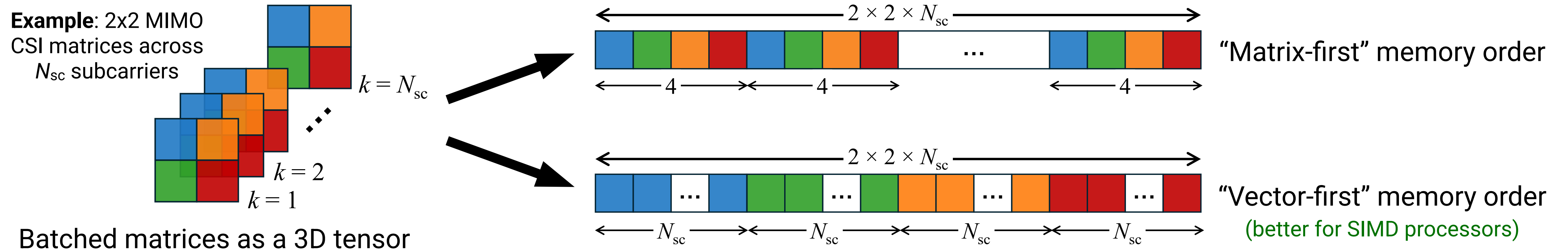


\*Same total fronthaul traffic for FR1/FR2

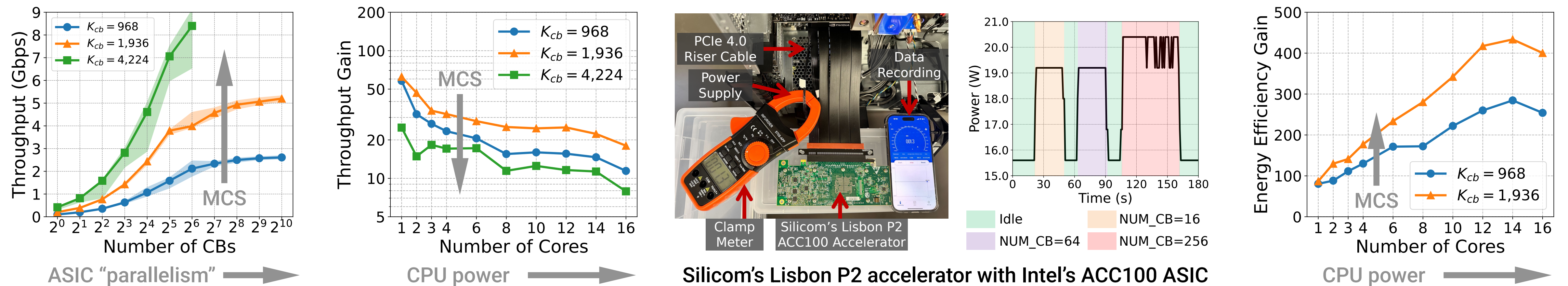
- 2x2 MIMO link with 100 MHz bandwidth, 100% uplink traffic load
- FR2 ( $\mu = 3$ ) with 0.125 ms slot, FR1 ( $\mu = 0$ ) with 1.0 ms
- PHY processing deadline set to 3-slot (frame schedule: DDDSU)

# Savannah: Real-time Efficient mmWave Baseband Processing

- **Key insight #1:** Vectorizing small matrix operations with optimized memory layout



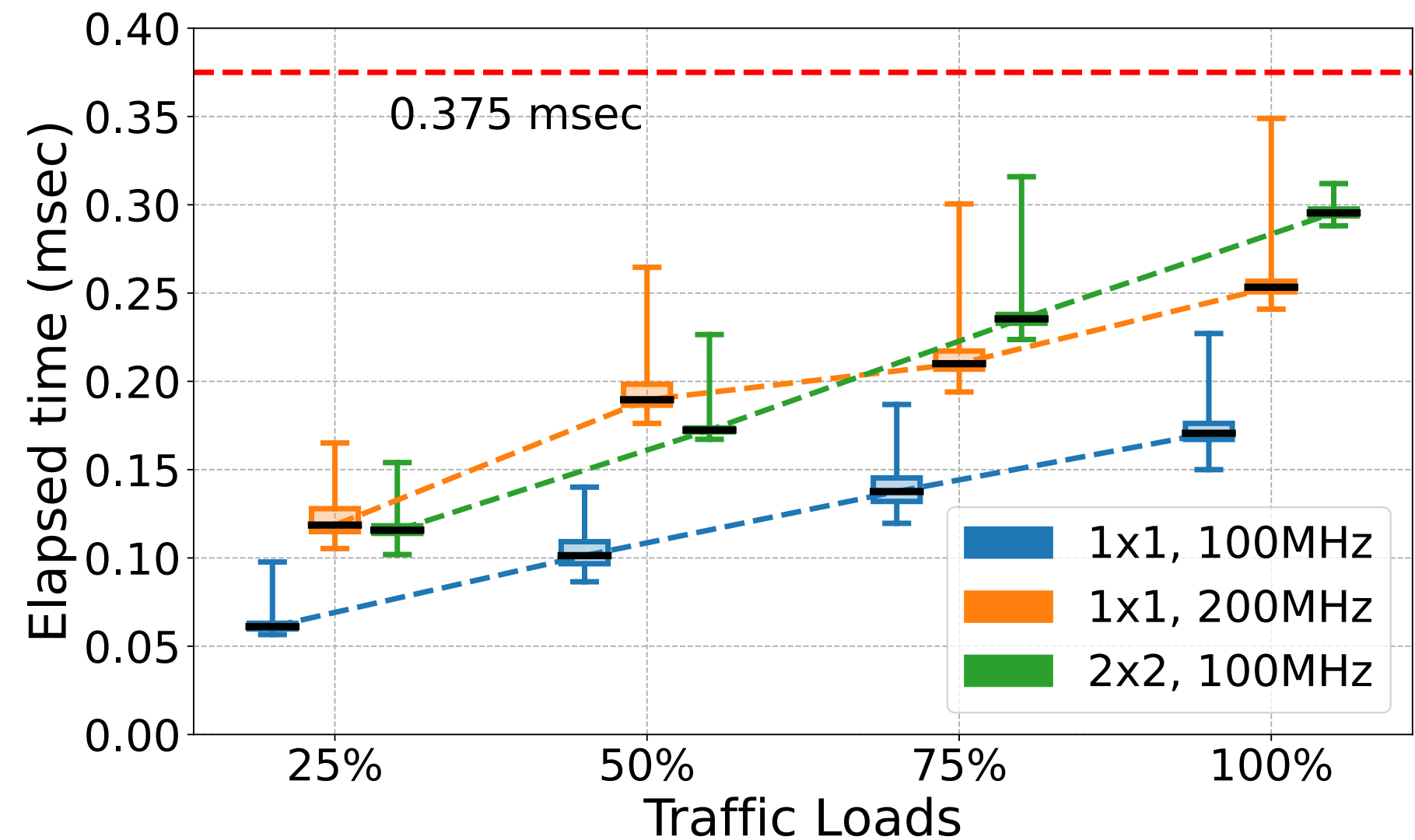
- **Key insight #2:** Leveraging heterogeneous compute resources (e.g., CPU + ASIC)



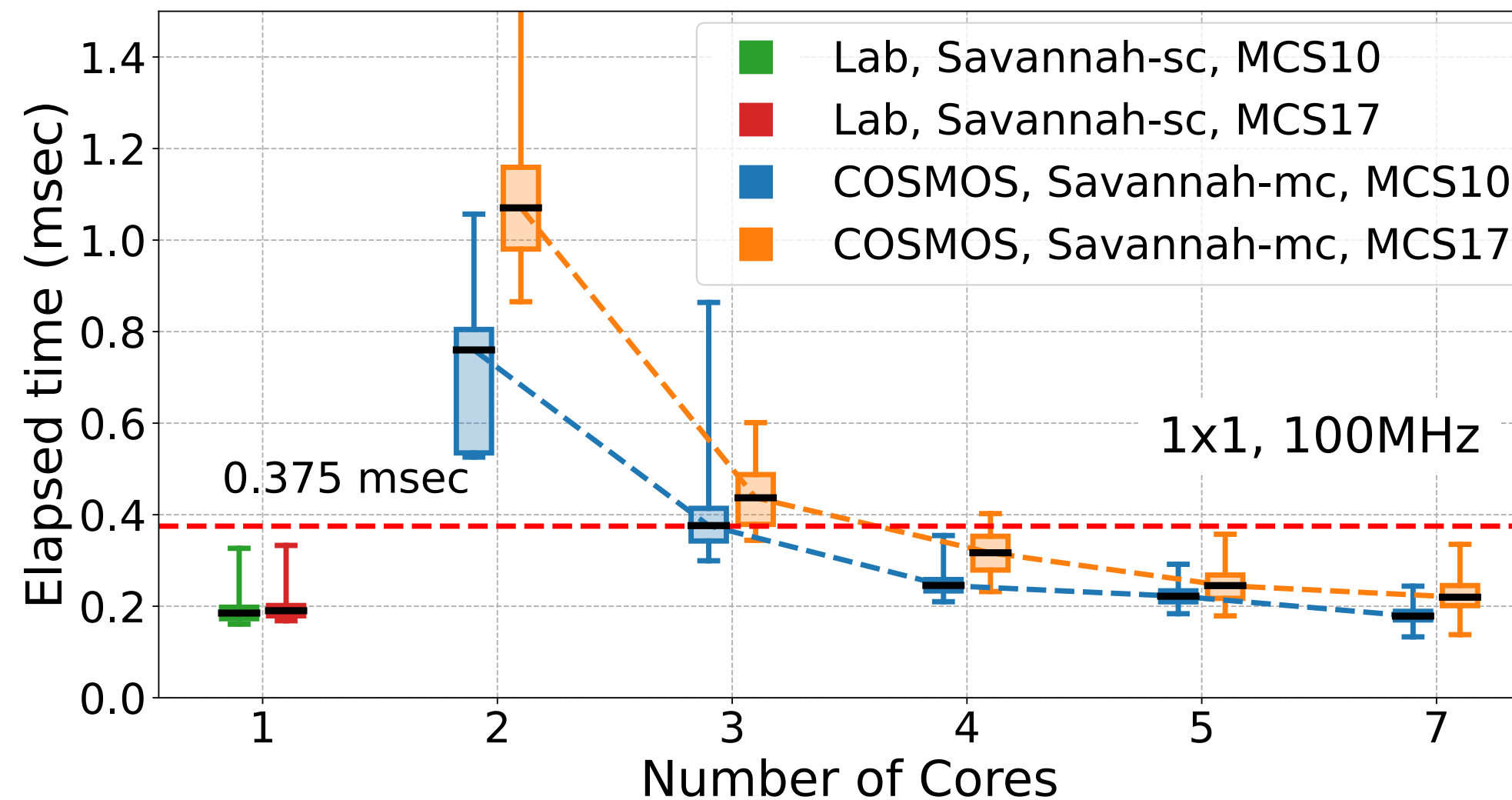
# Savannah: Real-time Efficient mmWave Baseband Processing

- **Lab:** Dell R750 server w/ Intel Xeon Gold 6348 CPU (56-core, @2.6 GHz), Silicom's Lisbon P2 ACC100 card
- **COSMOS:** Dell R740 server w/ Intel Xeon Gold 6226 CPU (48-core, @2.7 GHz)

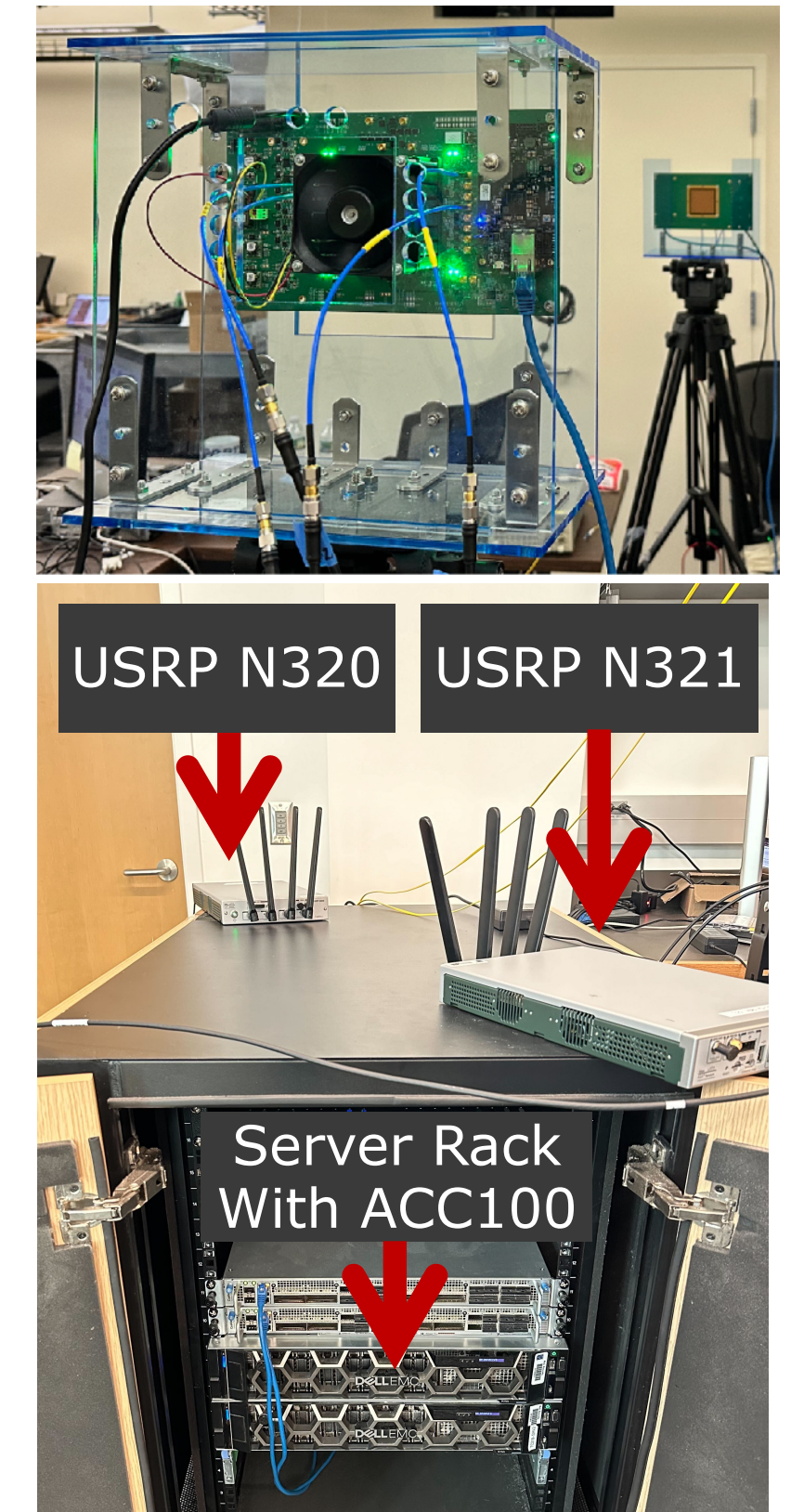
Achievable link data rate:  
 243.9 Mbps, 487.8 Mbps, 487.8 Mbps



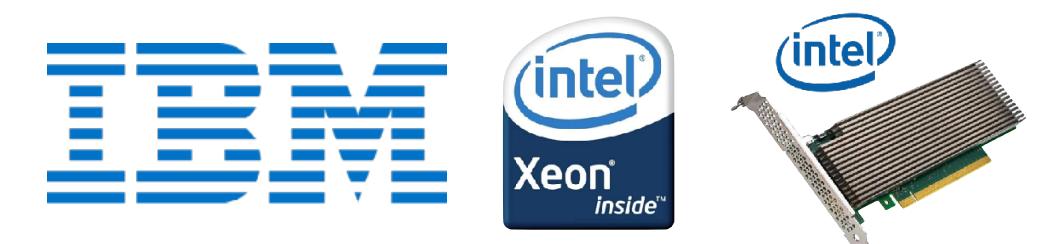
MCS10: 16QAM, LDPC code rate = 340/1,024  
 MCS17: 64QAM, LDPC code rate = 438/1,024



Savannah supports 2x2 MIMO, 100 MHz bandwidth in FR2 with a single CPU core + ASIC with 2x reduced power consumption with optimized signal processing and compute resource sharing



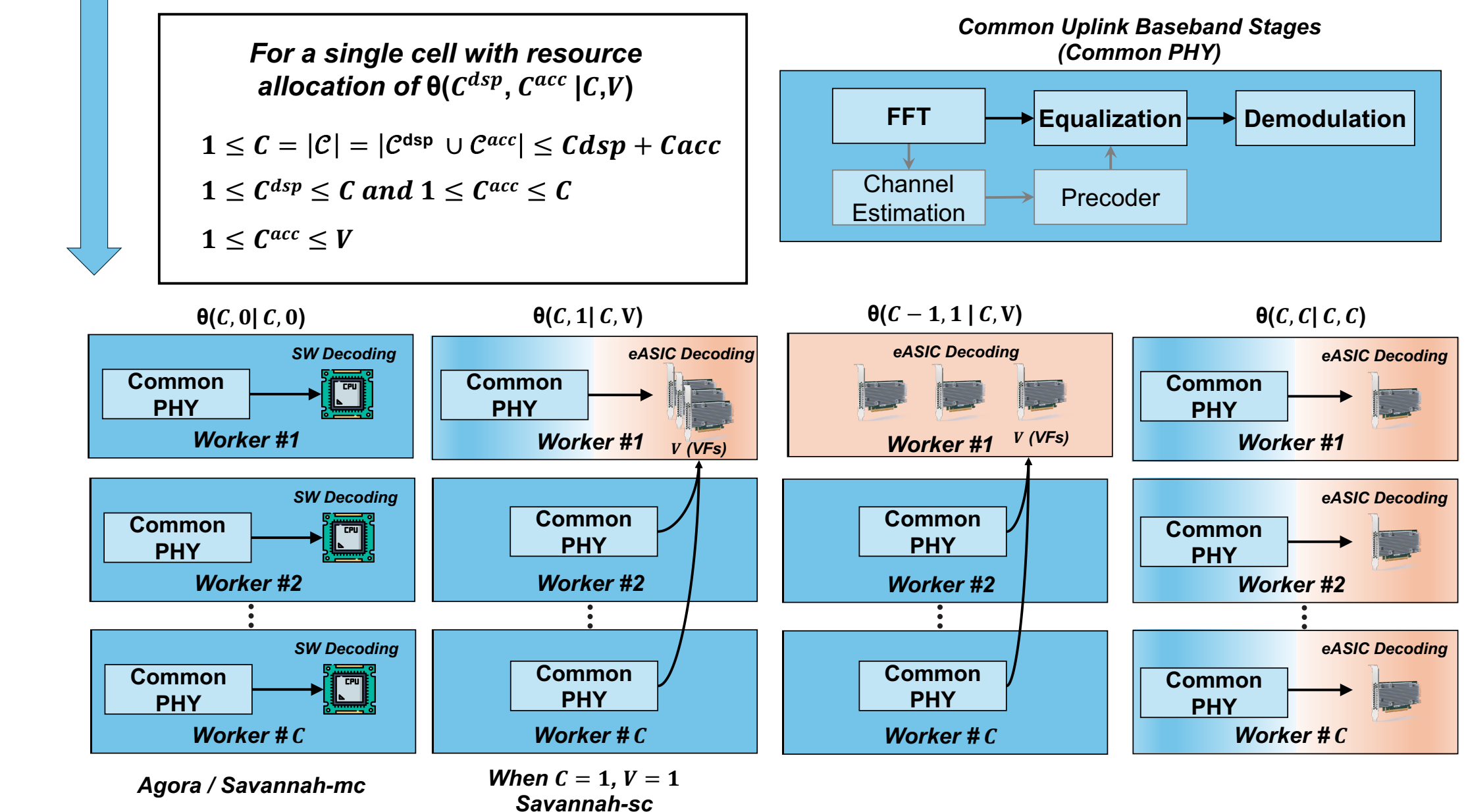
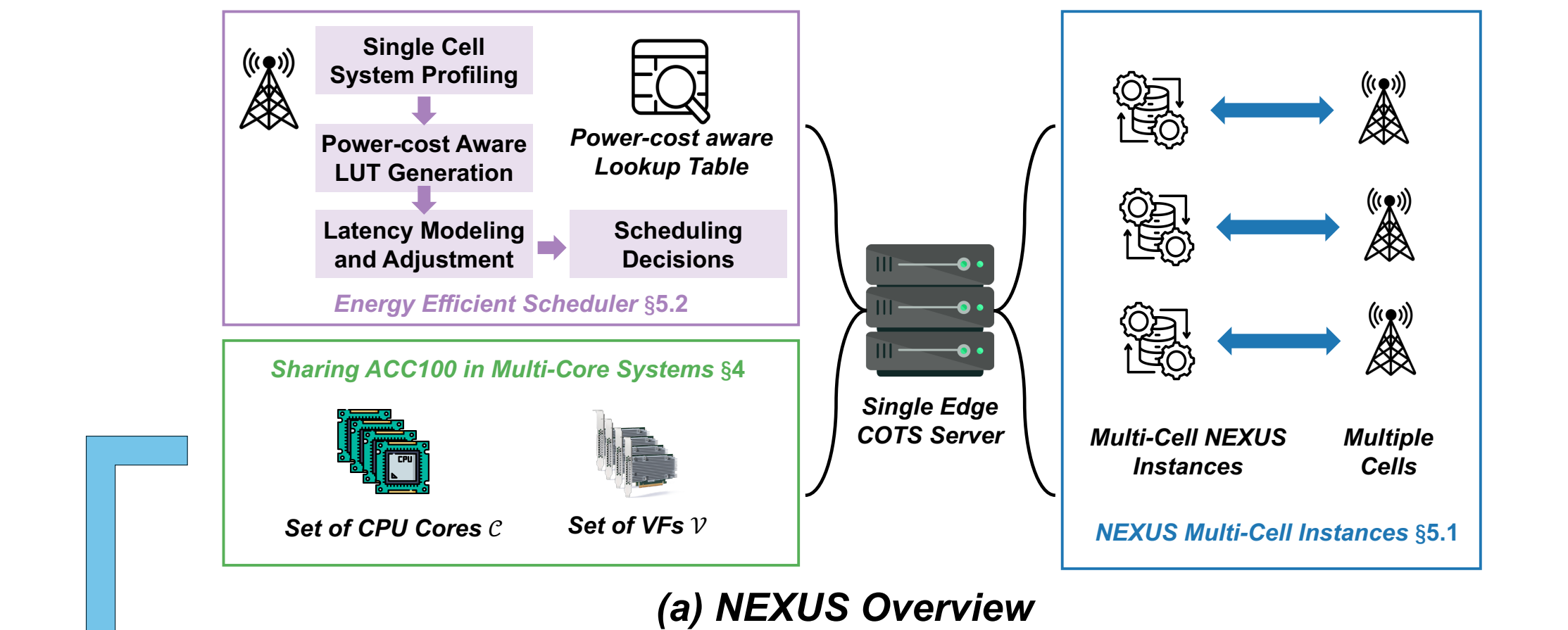
ReACT Testbed Phase 1



- ▶ Z. Qi\*, C.-H. Tung\*, A. Kalia, and T. Chen, "Savannah: Efficient mmWave baseband processing with minimal and heterogeneous resources," in *Proc. ACM International Conference on Mobile Computing and Networking (MobiCom'24)*, 2024.
- ▶ Code and data: <https://github.com/functions-lab/Savannah>

# NEXUS: Efficient Multi-cell mmWave Baseband Processing

- **Problem:** Use a single server with a pool of (cpu, acc100\_vf) resources supporting PHY processing of heterogeneous cells with varying (bw, mimo) configurations, traffic load, and MCSs
- **Objective:** Minimize the total power consumption while ensuring that individual cells meet the latency requirements
- **Highlights:**
  - **Co-optimized software-based DSP (CPU) and LDPC decoding offloaded to ACC100 virtual functions (VFs)**
  - **Power-aware compute resource allocation**
  - **A single server can support 16 cells with a total uplink data rate of ~5.4 Gbps using 16 CPU cores and one ACC100 HWA (16 VFs) while being scalable**



# NEXUS: Single-cell Optimization

- **Cell configuration:**  $\sigma (B, M, \rho_t, \rho_f, \chi)$
- **Resource allocation strategy:**  $\theta (C^{\text{dsp}}, C^{\text{acc}} | C, V)$ 
  - With  $C \leq C_{\text{max}}$  cores and  $V \leq V_{\text{max}}$  ACC100 VFs
  - $C^{\text{dsp}}$  cores are allocated for software DSP
  - $C^{\text{acc}}$  cores are allocated for ACC100 offloading tasks

Dimension	Notation	Range	Granularity
Channel BW	$B$	{100, 200, 400} MHz	Discrete
MIMO Size	$M$	{1, 2, 4}	Discrete
Traffic Load	$\rho_t$	[0,100] %	6.25%
Transmission BW %	$\rho_f$	[0,100] %	Multiple PRBs
MCS Index	$\chi$	{0, 1, ..., 28}	Integer

- Feasible compute resource allocation strategies:

$$\Theta(C, V) = \{ \theta (C^{\text{dsp}}, C^{\text{acc}} | C, V) \}$$

$$\text{s.t. } 1 \leq C = |C^{\text{dsp}} \cup C^{\text{acc}}| \leq C^{\text{max}},$$

$$1 \leq C^{\text{dsp}} \leq C, 0 \leq C^{\text{acc}} \leq C,$$

$$1 \leq C^{\text{acc}} \leq V.$$

Each core is dedicated to either software DSP, ACC100 offloading (LDPC decoding), or both

Each core responsible for ACC100 offloading is mapped to a unique set of VFs, and each VF cannot be managed by more than one core

# NEXUS: Single-cell Optimization

- **Power-aware compute resource optimization** for a single cell:

$$\theta^* := \arg \min_{\theta \in \Theta} : P(\theta) = \alpha_1 P_1(C^{\text{dsp}}) + \alpha_2 P_2(C^{\text{acc}}) + \alpha_3 P_3(V), \text{ s.t. } \mathcal{L}(\theta) < 3 \text{ slots.}$$

$\approx 7.0 \cdot C$  (Watt)

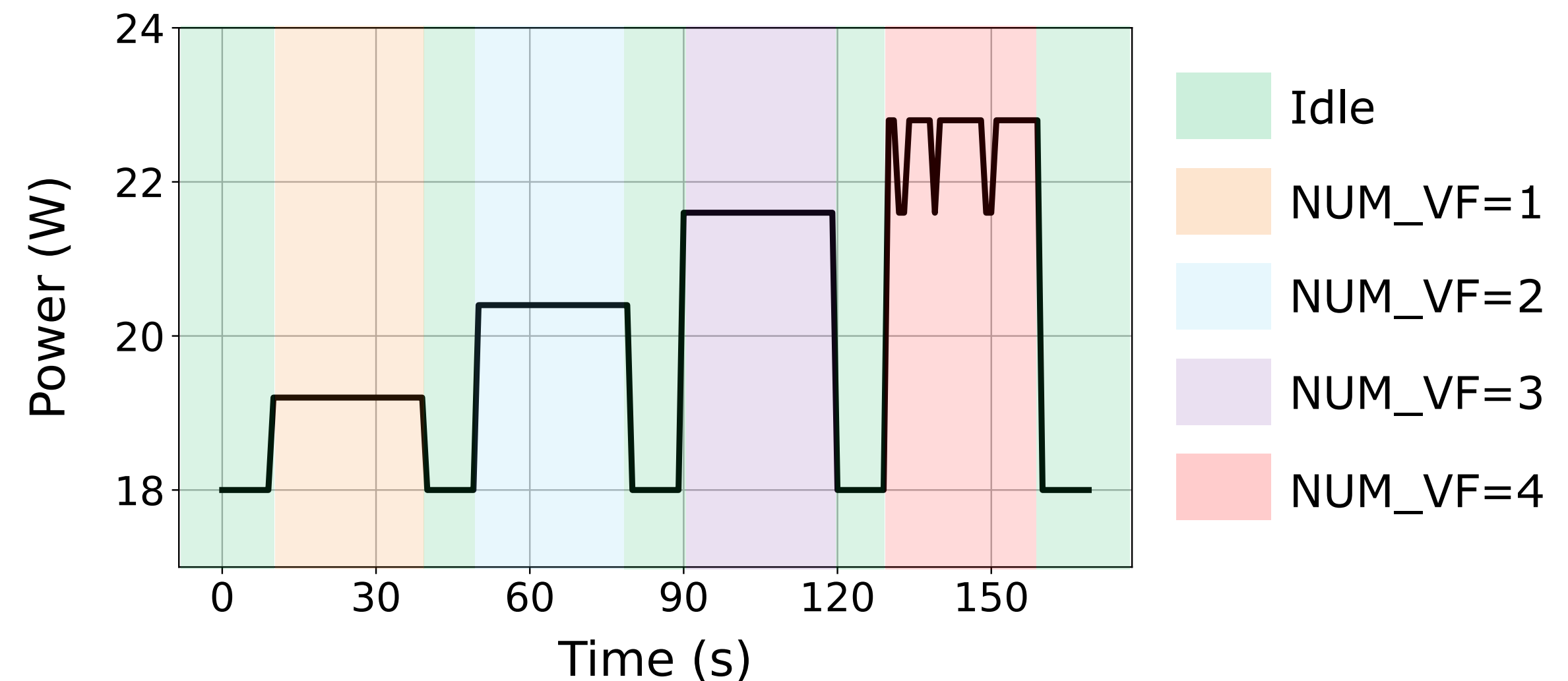
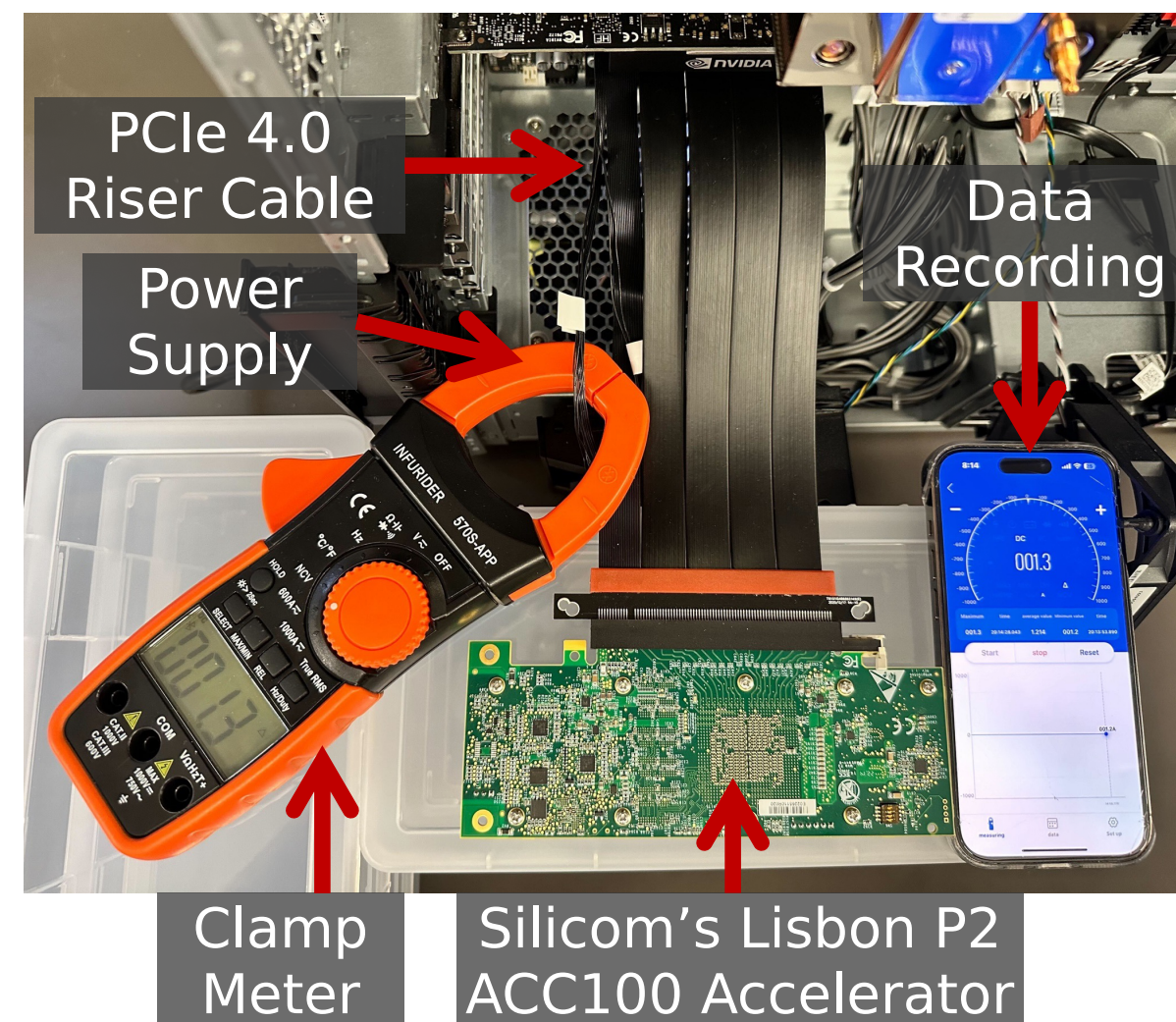
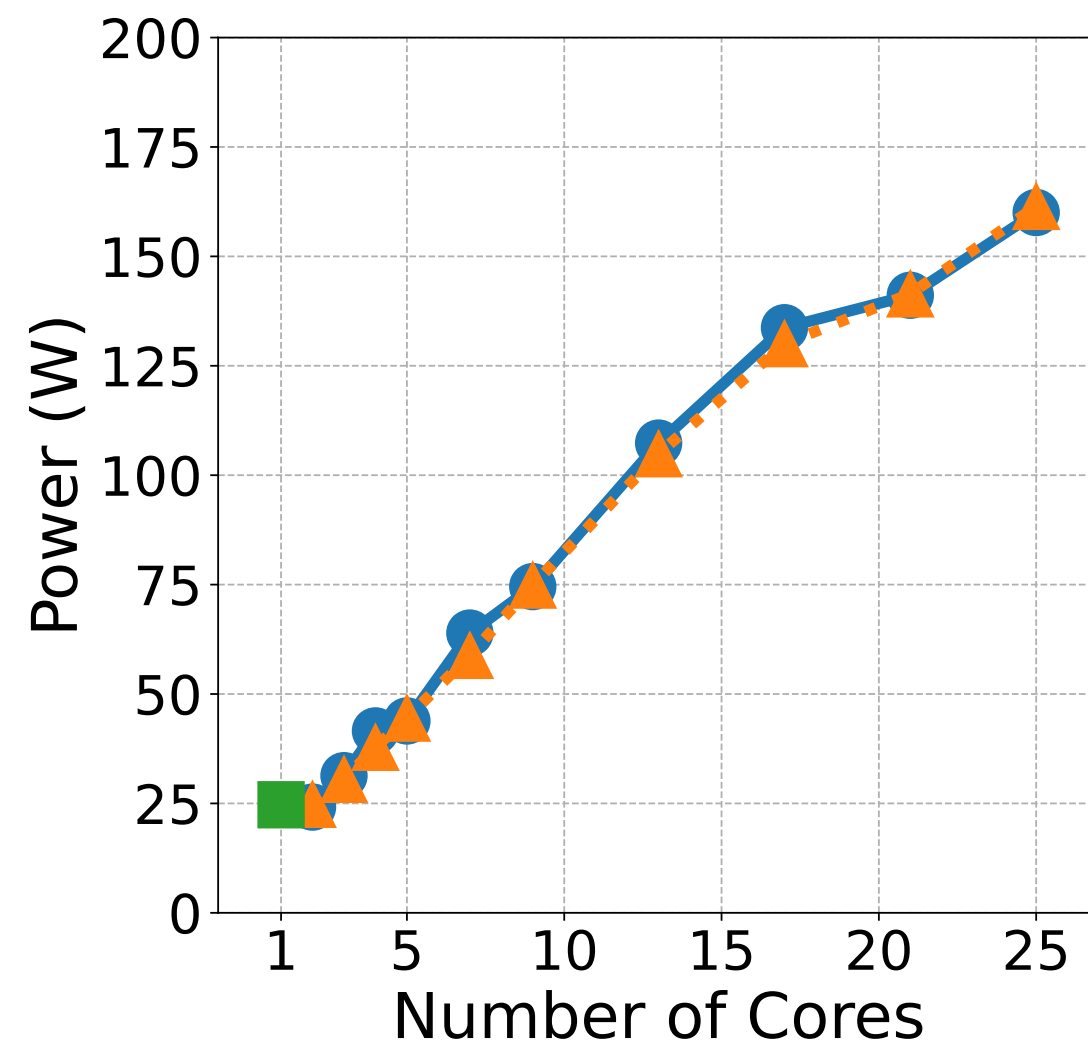
Using Intel's Performance Counter Monitor (PCM)

$\approx 1.2 \cdot V$  (Watt)

Verified method using a clamp meter



Power Consumption



# NEXUS: Single-cell Optimization

- **Power-aware compute resource optimization** for a single cell:

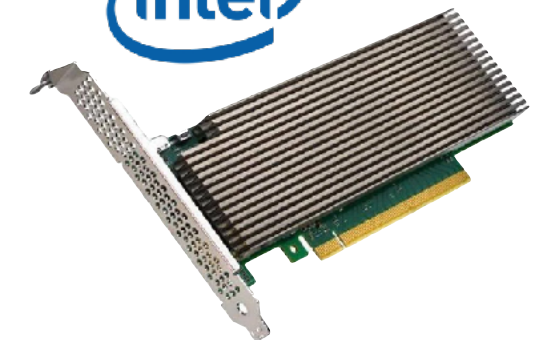
$$\theta^* := \arg \min_{\theta \in \Theta} : P(\theta) = \alpha_1 P_1(C^{\text{dsp}}) + \alpha_2 P_2(C^{\text{acc}}) + \alpha_3 P_3(V), \text{ s.t. } \mathcal{L}(\theta) < 3 \text{ slots.}$$

→ Large configuration space

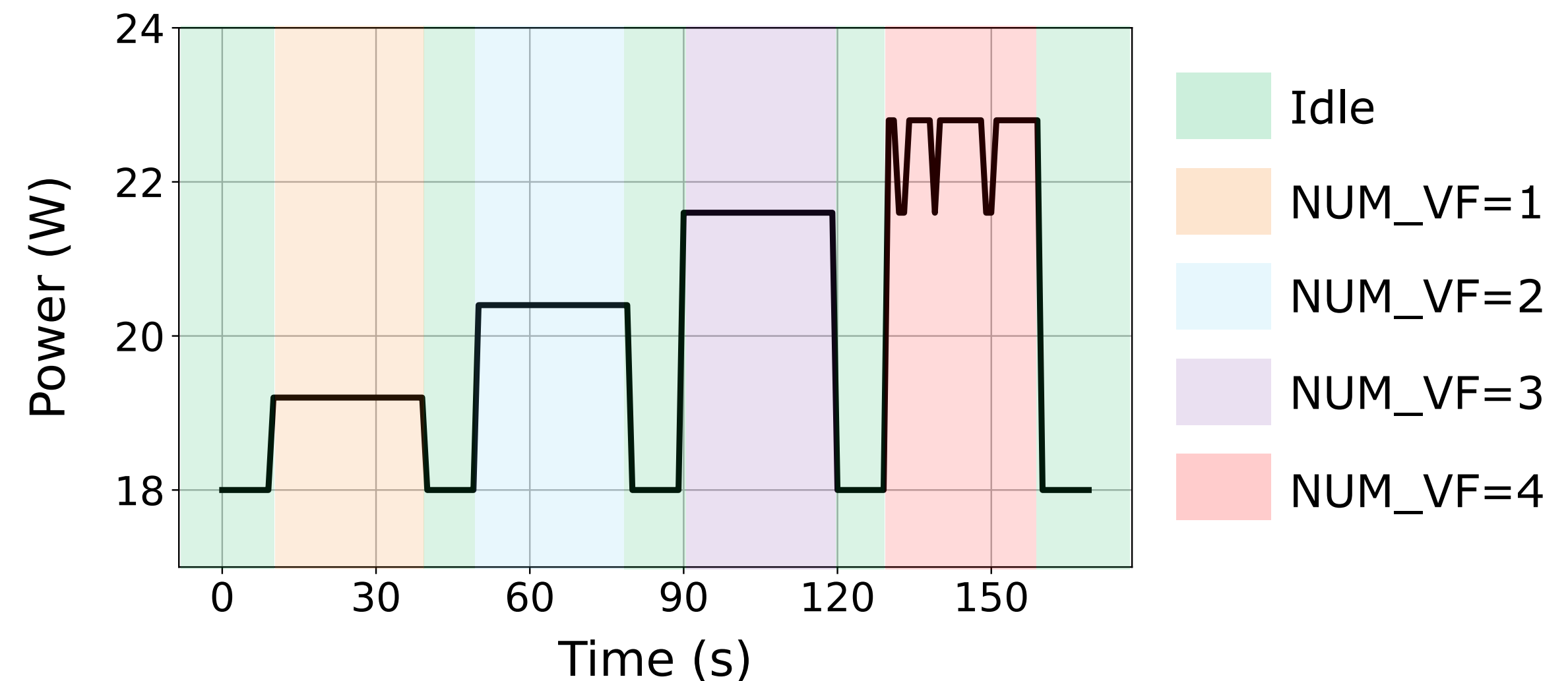
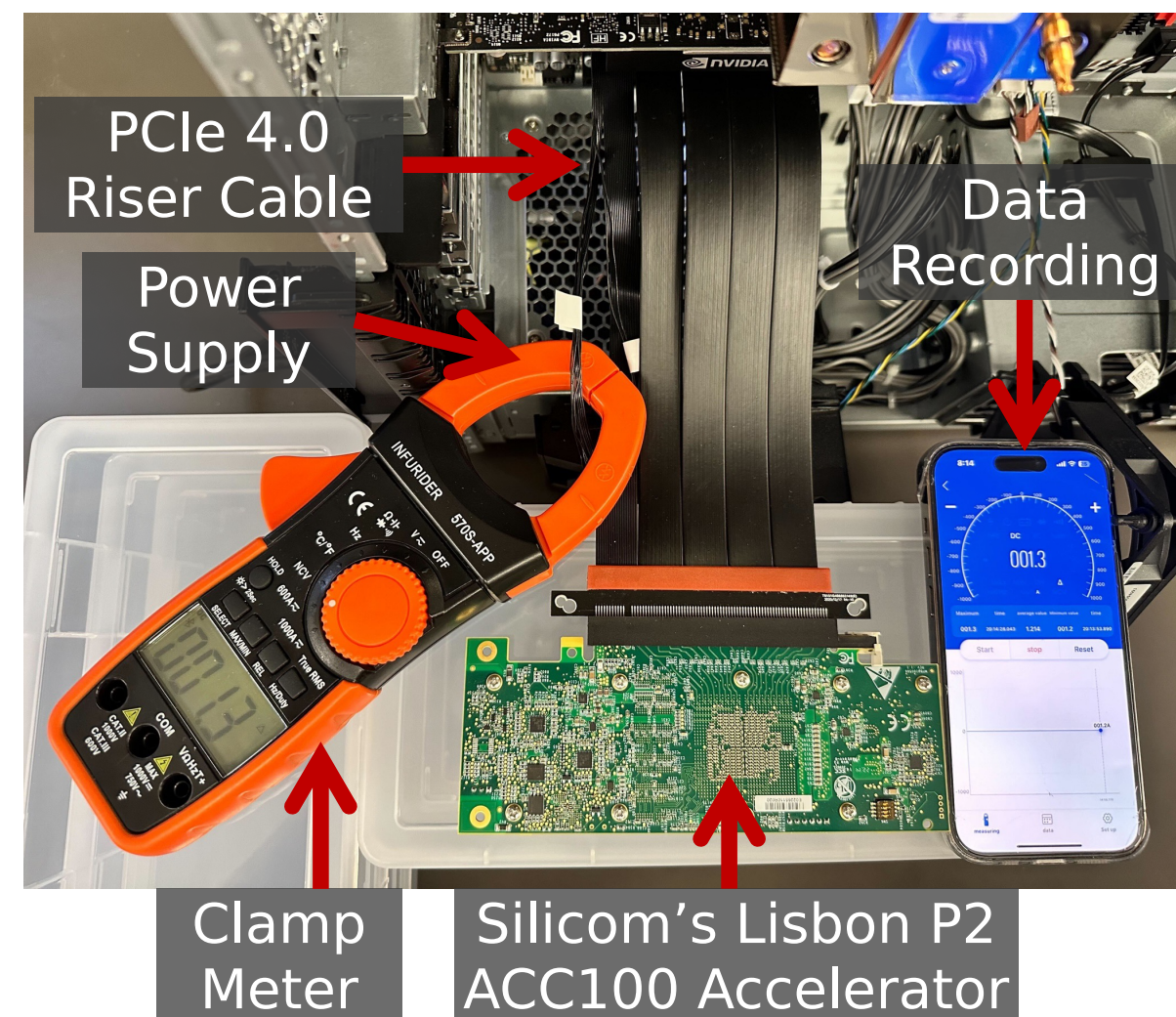
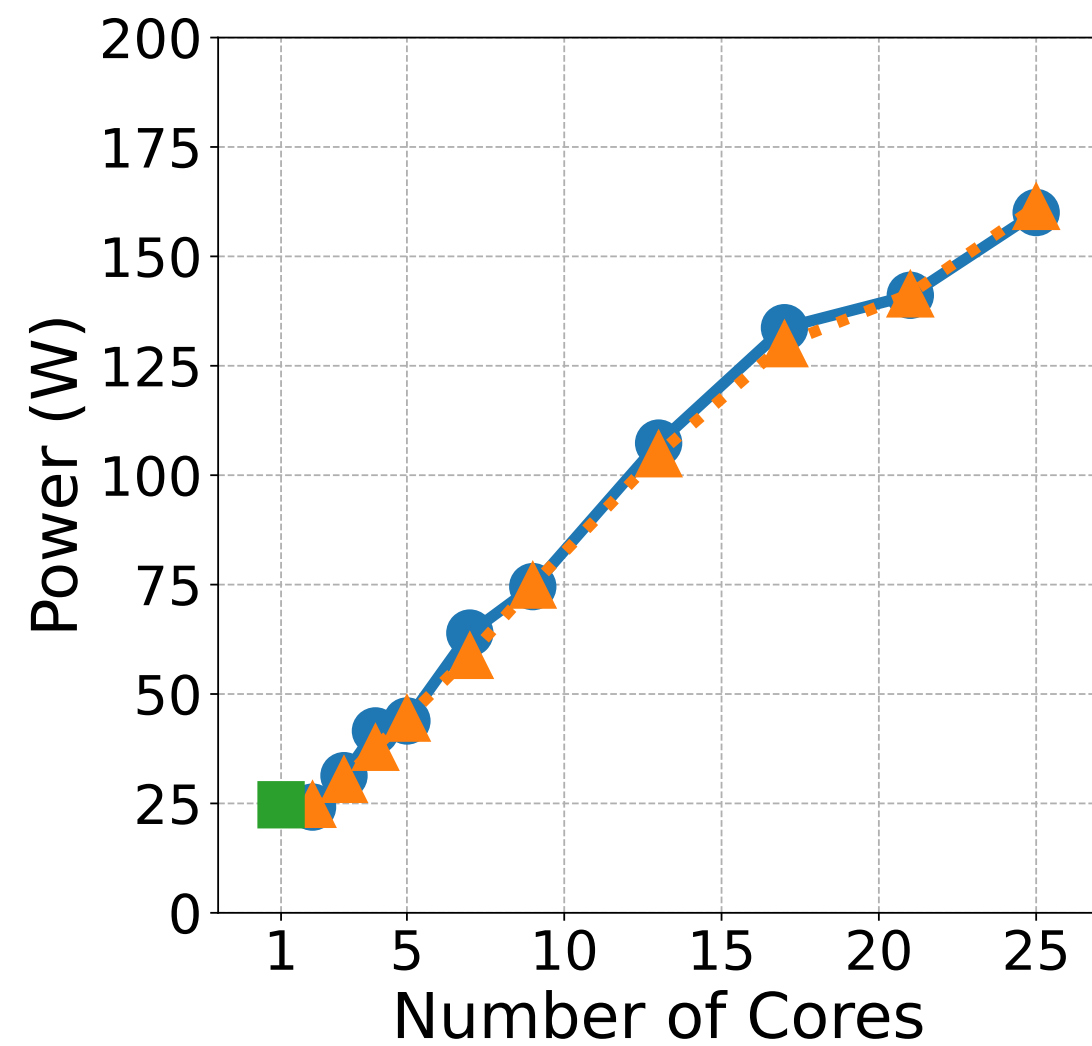
→  $\approx 7.0 \cdot C$  (Watt)  
Using Intel's Performance Counter Monitor (PCM)

→  $\approx 1.2 \cdot V$  (Watt)  
Verified method using a clamp meter

→ Difficult to model

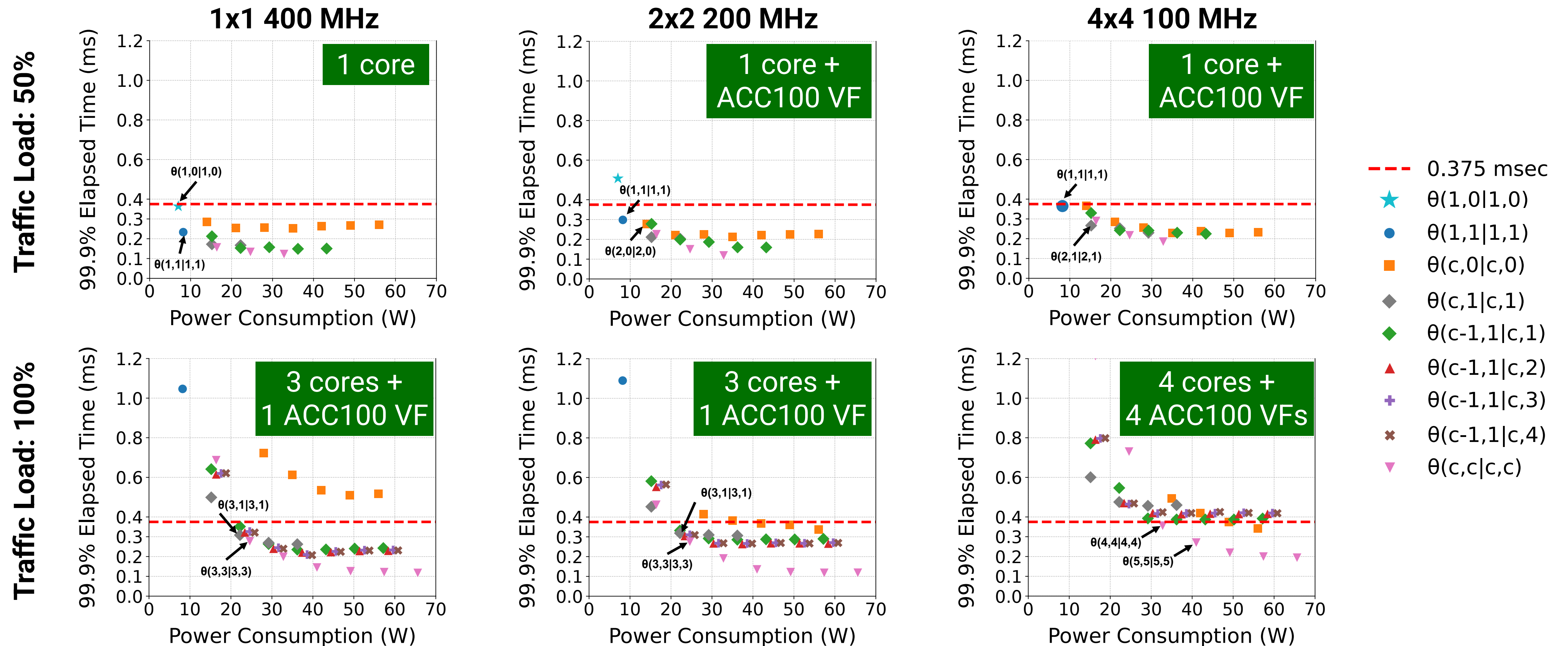


**Power Consumption**



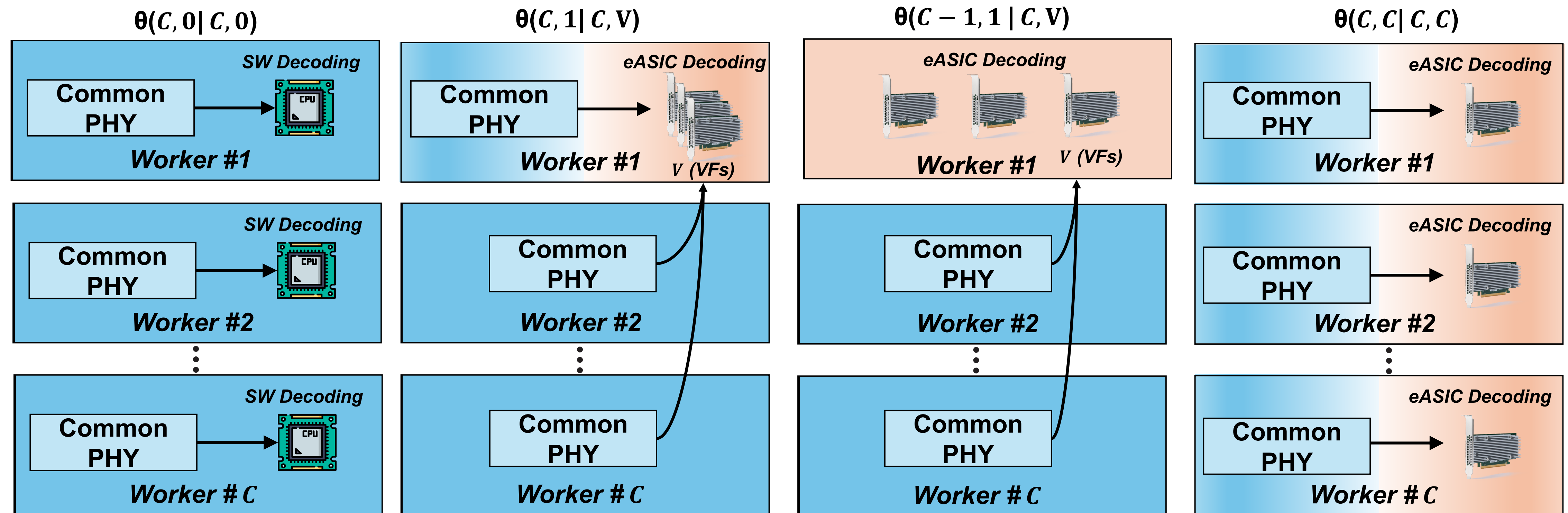
# NEXUS: Single-cell Optimization

- **Configuration space reduction** based on key observations from comprehensive system profiling
  - E.g., from 152 to 15 with  $(C, V) = (26, 16)$



# NEXUS: Single-cell Optimization

- **Configuration space reduction** based on key observations from comprehensive system profiling
  - E.g., from 152 to 15 with  $(C, V) = (26, 16)$
- We consider four resource allocation strategies:

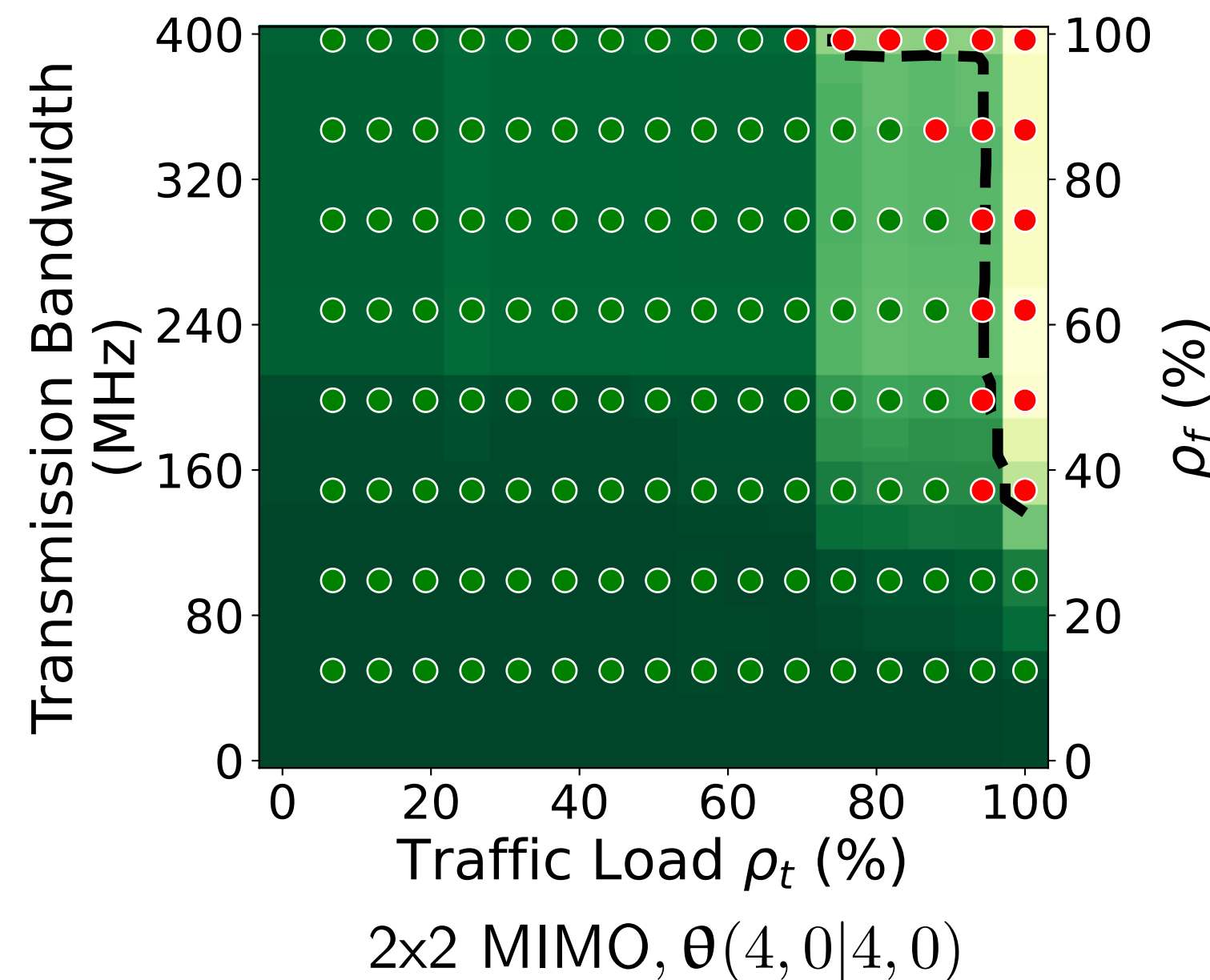
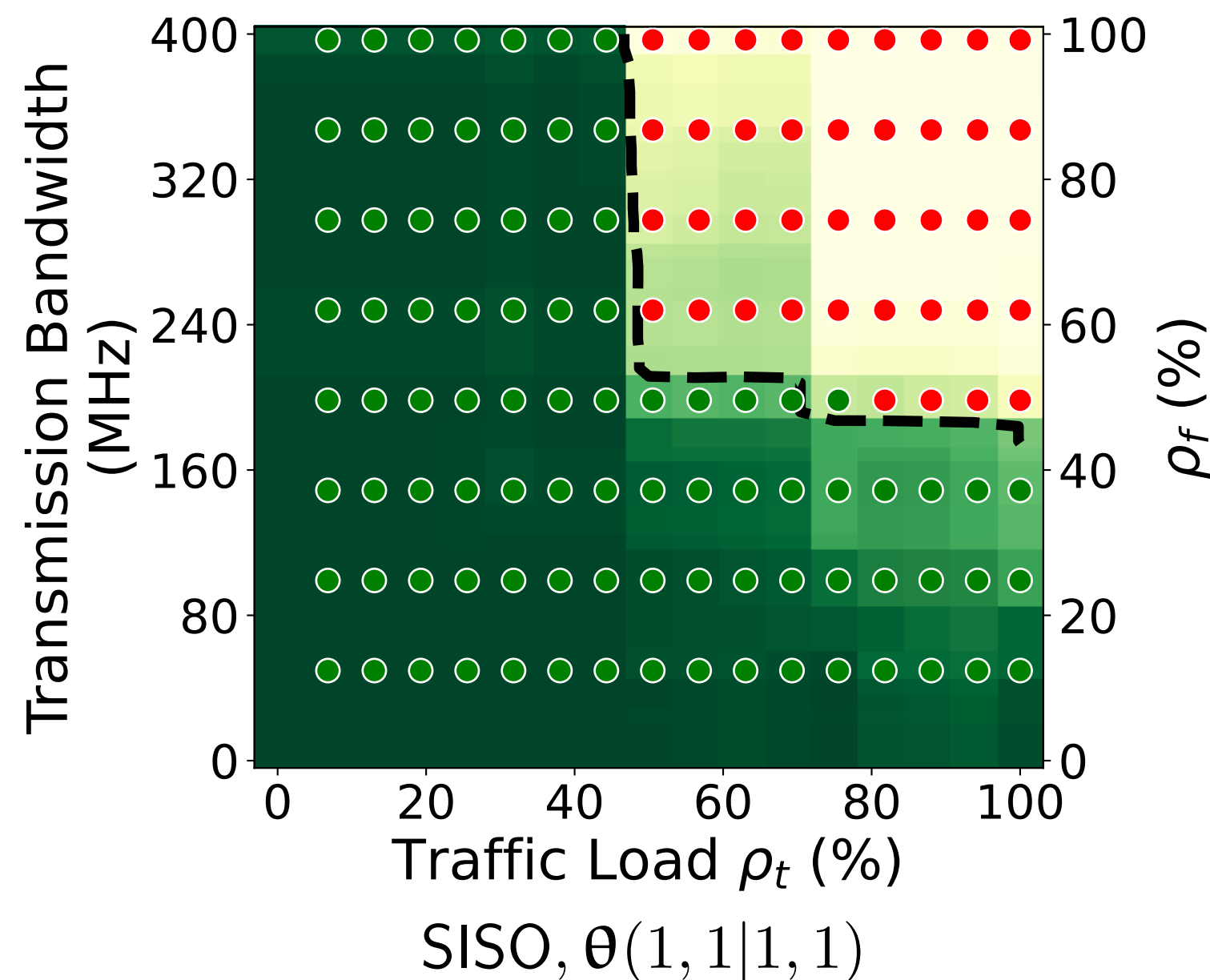


Agora / Savannah-mc

When  $C = 1, V = 1$   
Savannah-sc

# NEXUS: Single-cell Optimization

- **Configuration space reduction** based on key observations from comprehensive system profiling
  - E.g., from 152 to 15 with  $(C, V) = (26, 16)$
- **Random forest (RAF) based feasibility prediction:**  $f : (\sigma, \theta) \mapsto [0, 1]$ 
  - Trained on >9K profiling experiments, each with 20K frames with detailed PHY processing timing



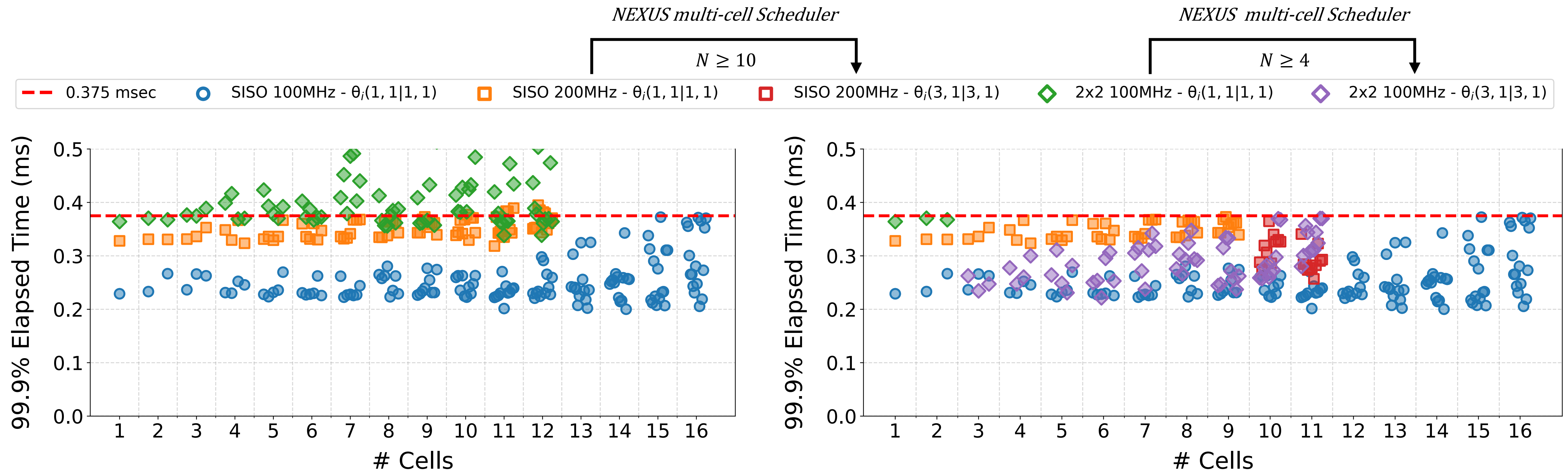
RAF-based feasibility prediction (dashed lines) vs. ground truth (green: feasible; red: infeasible), for an unseen MCS-18

	Accuracy	Inference Time
MLP	99.1%	582 usec
1D CNN	99.2%	615 usec
TinyNet	99.3%	570 usec
RAF (10 trees)	99.0%	1.4 usec
RAF (100 trees)	99.1%	4.3 usec

Accuracy and inference time comparison between different methods

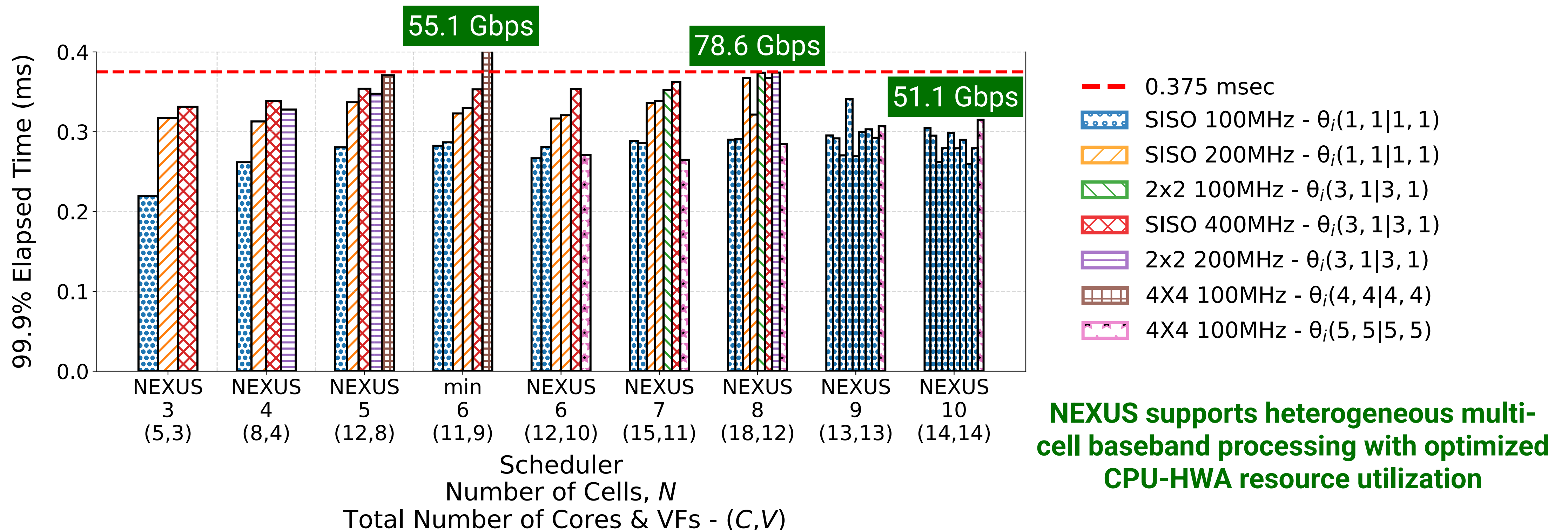
# NEXUS: Multi-cell Optimization

- Assume **non-overlapping allocation** of compute resources (CPU cores, ACC100 VFs) across  $N$  cells
- Leverage single-cell optimization to reduce configuration space across  $N$  cells
- **Contention-aware latency adjustment** of the RAF-based feasibility prediction:
  - If the predicted cell configuration cannot meet the latency requirements, go to the next best configuration



# NEXUS: Implementation and Results

- ~5K LOC (C++), ~9K configuration files, ACC100 VF management based on Intel's `bbdev` tool (`pf_bb_config`) and DPDK's threading model (`rte_thread`)
- Two servers emulating the set of RUs and DUs, connected via 100 GbE links (Mellanox CX6)
- Example: NEXUS with heterogeneous cell configurations and heterogeneous resource allocations



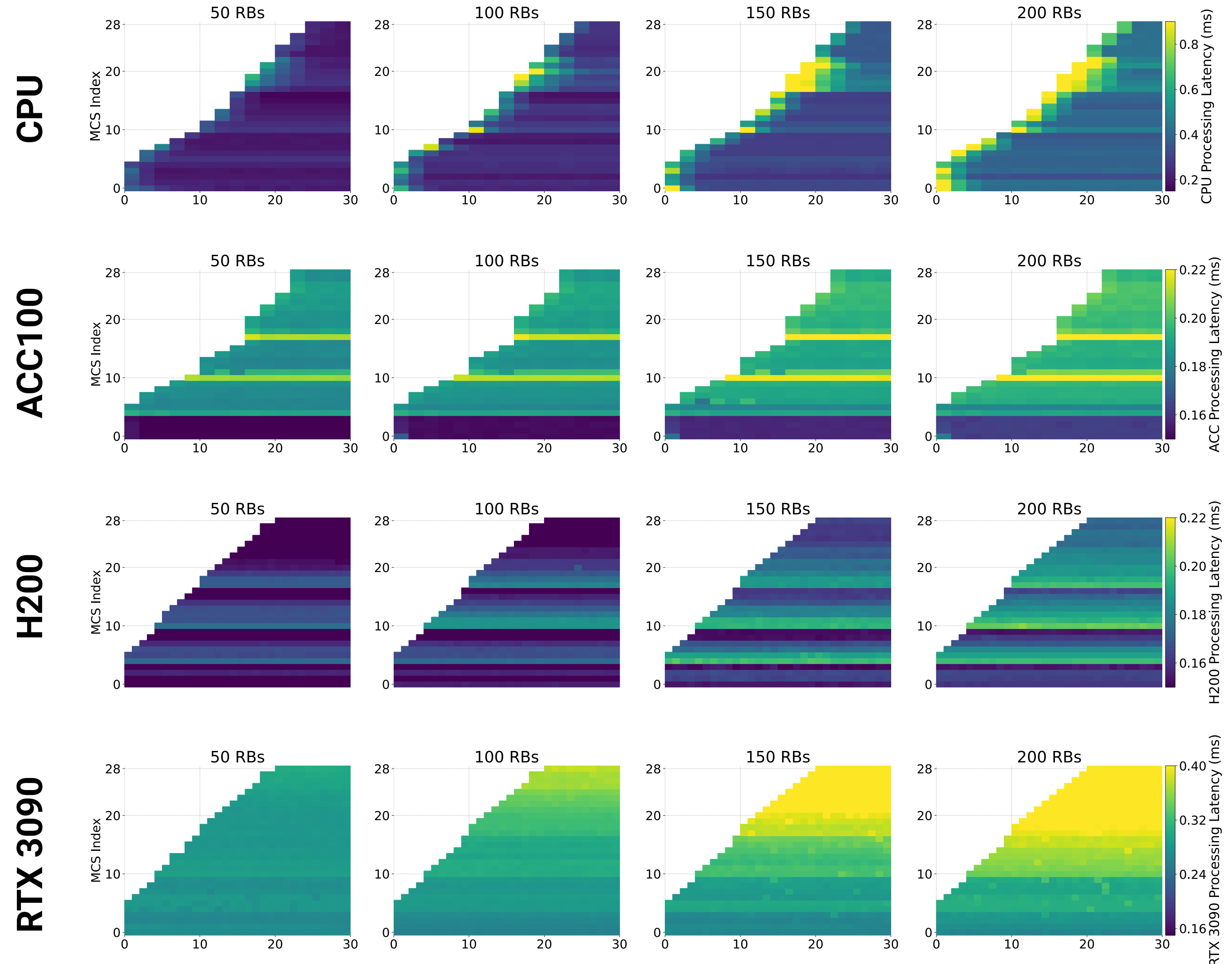
# DecodeX: LDPC Decoding Across CPU, GPU, and ASIC Platforms

- A comprehensive suite of LDPC decoder implementations (kernels, APIs, and test vectors) for benchmarking and R&D:

- CPU (FlexRAN)
- GPU (NVIDIA Aerial and Sionna-RK)
- ASIC (ACC100)

- Allow for evaluation of L1 processing models with different HWAs:

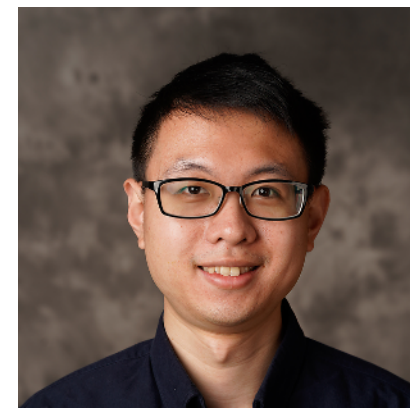
- **Inline**: CPU or GPU only (e.g., FlexRAN, Aerial)
- **Lookaside**: CPU + ACC100 (e.g., Savannah, NEXUS), CPU + GPU
- **Hybrid**: CPU-GPU with unified memory (e.g., Sionna-RK)



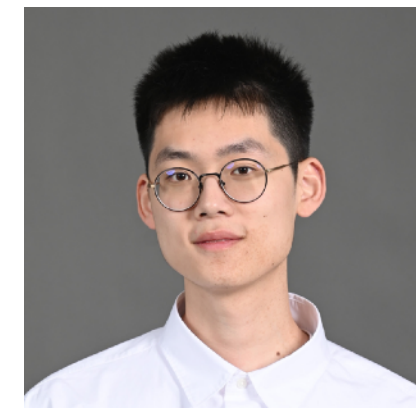
# Thank you! Questions?



Zhenzhou  
Qi



Chung-Hsuan  
Tung



Yuncheng  
Yao



Yiming  
Li



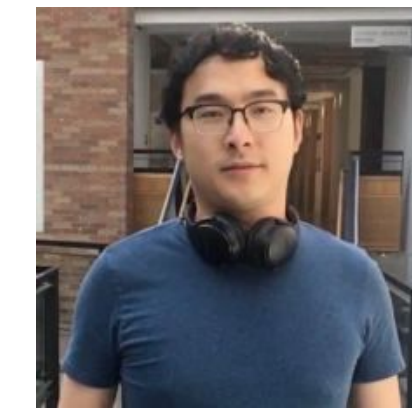
Yunyao  
Zheng



Zihui  
Gao



Anuj  
Kalia



Danyang  
Zhuo

[tingjun.chen@duke.edu](mailto:tingjun.chen@duke.edu); <https://tingjunchen.com>



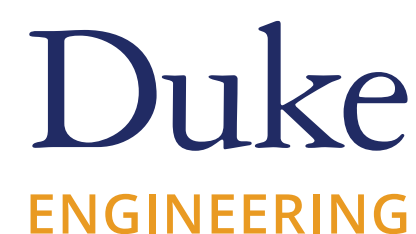
Platforms for Advanced  
Wireless Research



Semiconductor  
Research  
Corporation



NSF AI Institute for Edge Computing



Center for Ubiquitous Connectivity 16