

# **EPI TUTORIAL: POWER ASPECT**

CO-LOCATED WITH THE HIPEAC2020

BOLOGNA, ITALY

22 JANUARY 2020



**European  
Processor  
Initiative**



# FRAMEWORK PARTNERSHIP AGREEMENT IN EUROPEAN LOW-POWER MICROPROCESSOR TECHNOLOGIES

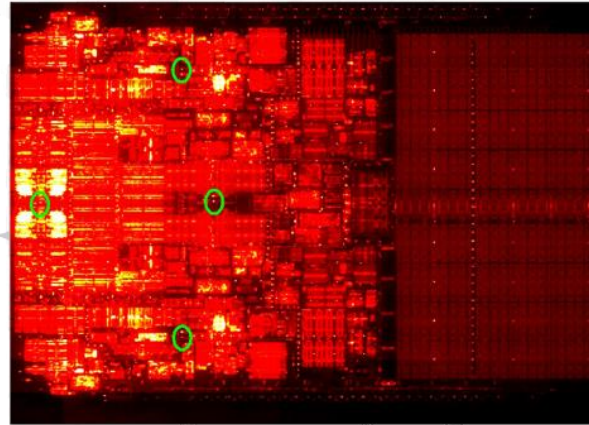
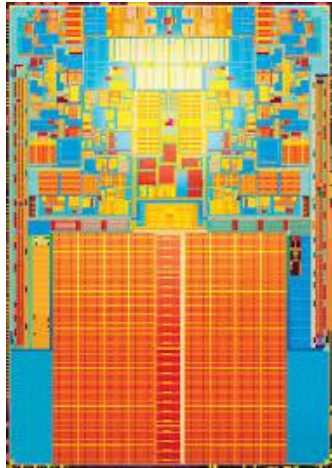


THIS PROJECT HAS RECEIVED FUNDING FROM THE EUROPEAN UNION'S HORIZON 2020 RESEARCH AND INNOVATION  
PROGRAMME UNDER GRANT AGREEMENT NO 826647



**European  
Processor  
Initiative**

Intel Pentium - 1 core



○ = Thermal Sensor

End of Dennard's Scaling  
=> Power density increases  
=> Thermal issues

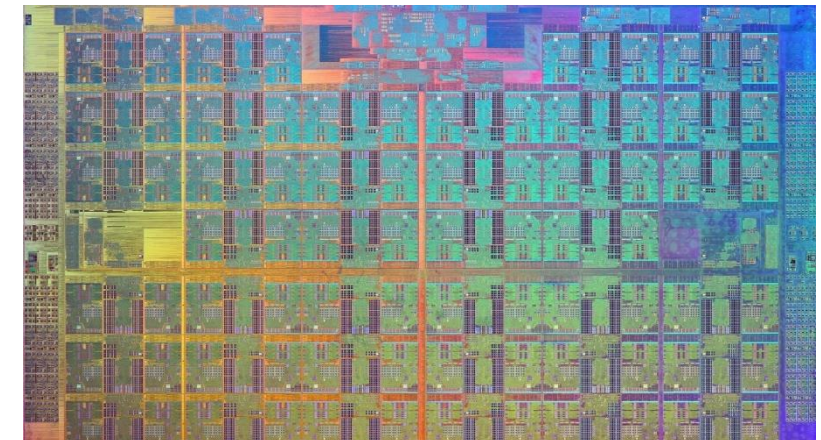
Moore's Law: "The number of transistors and resistors on a chip doubles every 24 months"

↑ theoretical max frequency  
↑ #cores

Same Power Budget !

1. ↓ frequency
2. ↓ # active core at the same time
3. ↑ cooling cost

Intel KNL – 72 cores

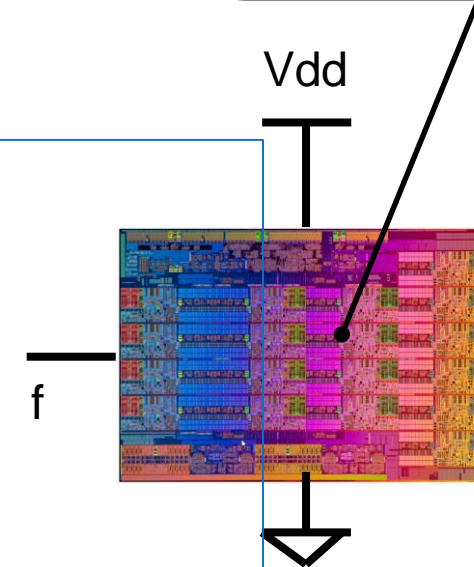
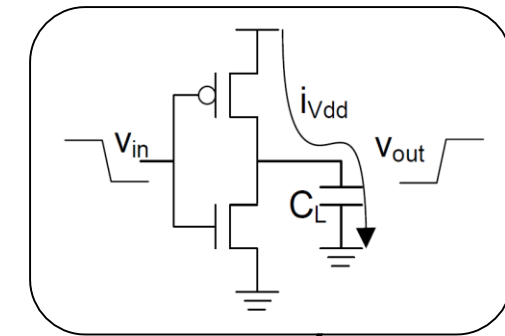


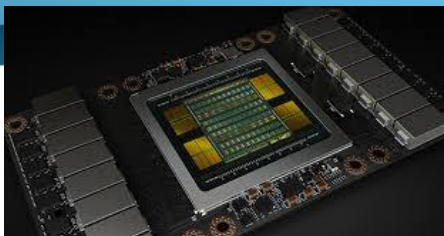
# DYNAMIC POWER

$$P_{dynamic} = \underbrace{a \times C_L}_{\text{"effective capacitance" } (C_{Effective})} \times V_{dd}^2 \times f$$

activity factor      load capacitance      supply voltage      clock frequency

- **Linear** ↓ with ↓  $C_{Effective}$
- **Linear** ↓ with ↓  $f$
- **Quadratic** ↓ with ↓  $V_{dd}$
- **Cubic** ↓ with ↓ both  $V_{dd}$  and  $f$



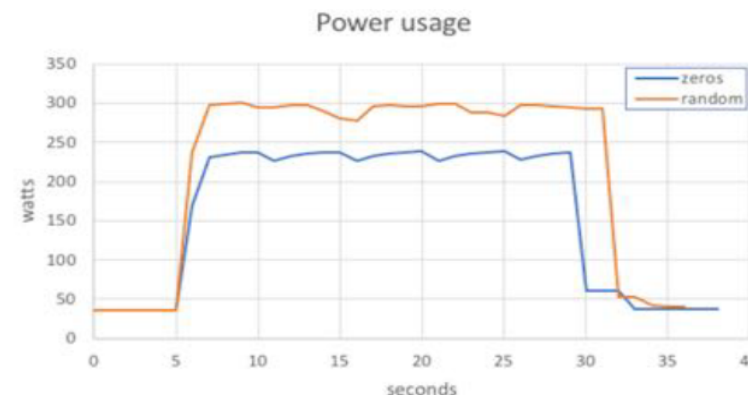
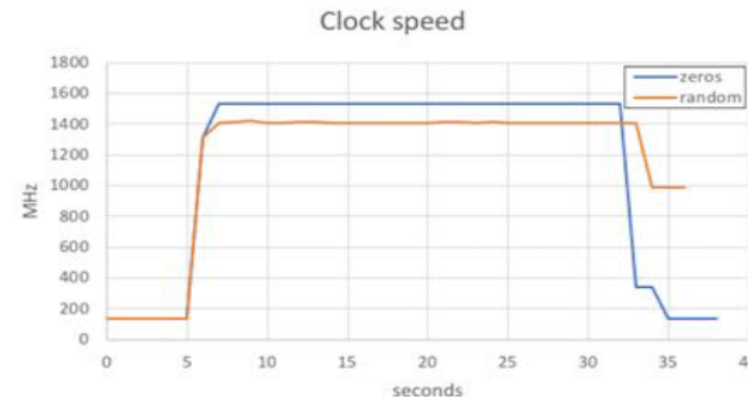


## Ceff EFFECT

- TC/HGEMM has surprising data-dependent performance: **125 TF** theoretical peak, **113 TF** achievable on zero-filled matrices, **105 TF** peak on random CCC matrices, **~95 TF** peak on matrices with fully random FP16 entries

### Issues

- Measurements on 1 Summit GPU using `nvidia-smi`
- Data-dependent performance of Tensor Cores is due to 300W power/frequency throttling of Voltas on Summit
- Baidu DeepBench GEMM benchmark has a bug (reported), incorrectly fills FP16 matrices with zeros instead of the intended random values, thus miscalculates GPU performance



# SUB-THRESHOLD LEAKAGE CURRENT

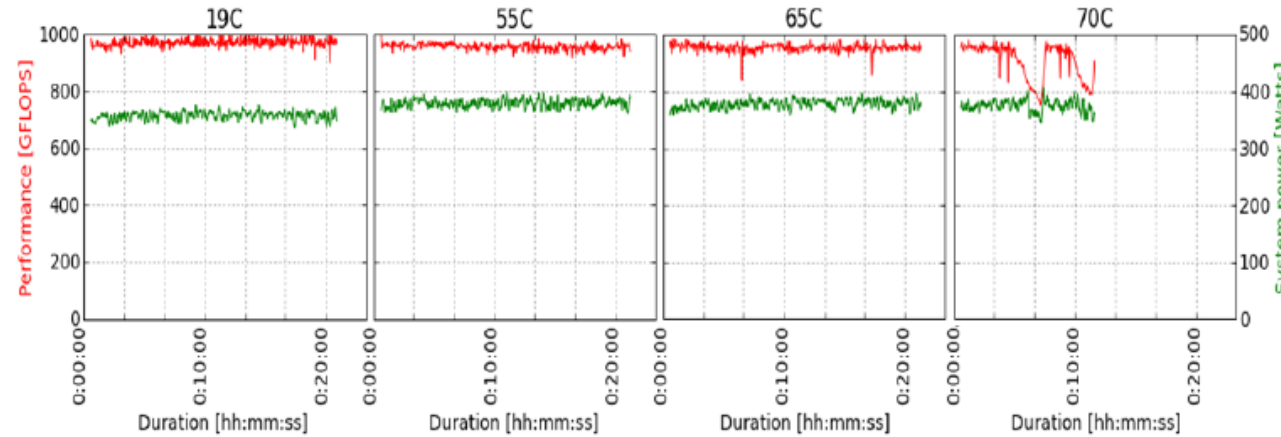
$$I_{leakage} = k_1 \times \left(1 - e^{-k_2 \times V_{ds} / T}\right) \times e^{k_3 \times (V_{gs} - V_{TH} - V_{off}) / T}$$

Diagram illustrating the components of the sub-threshold leakage current equation:

- constants** (points to  $k_1$  and  $k_2$ )
- temperature** (points to  $T$ )
- drain to source voltage** (points to  $V_{ds}$ )
- gate to source voltage** (points to  $V_{gs}$ )
- threshold voltage** (points to  $V_{TH}$ )
- empirical parameter** (points to  $V_{off}$ )

- **Exponential ↓ with ↓  $V_{gs}$  ( $\sim V_{dd}$ )**
- **Exponential ↓ with ↑  $V_{TH}$**
- **Exponential ↓ with ↓  $T$**

# LEAKAGE IN SUPERCOMPUTING



A. Moskovsky et al., "Server level liquid cooling: Do higher system temperatures improve energy efficiency?" Supercomputing frontiers and innovations, 3(1):67-74, 2016.

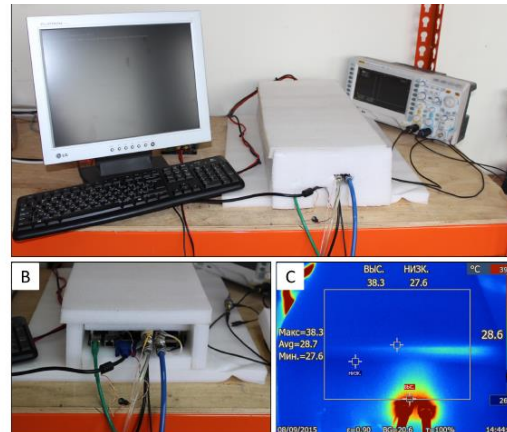
**Figure 3.** Illustration of system performance and power profiles at  $T_{LC} = 19, 55, 65$  and  $70\text{ }^{\circ}\text{C}$

in Temperatures Improve Energy Efficiency?

Table 1: Performance and power characteristics at different coolant temperatures

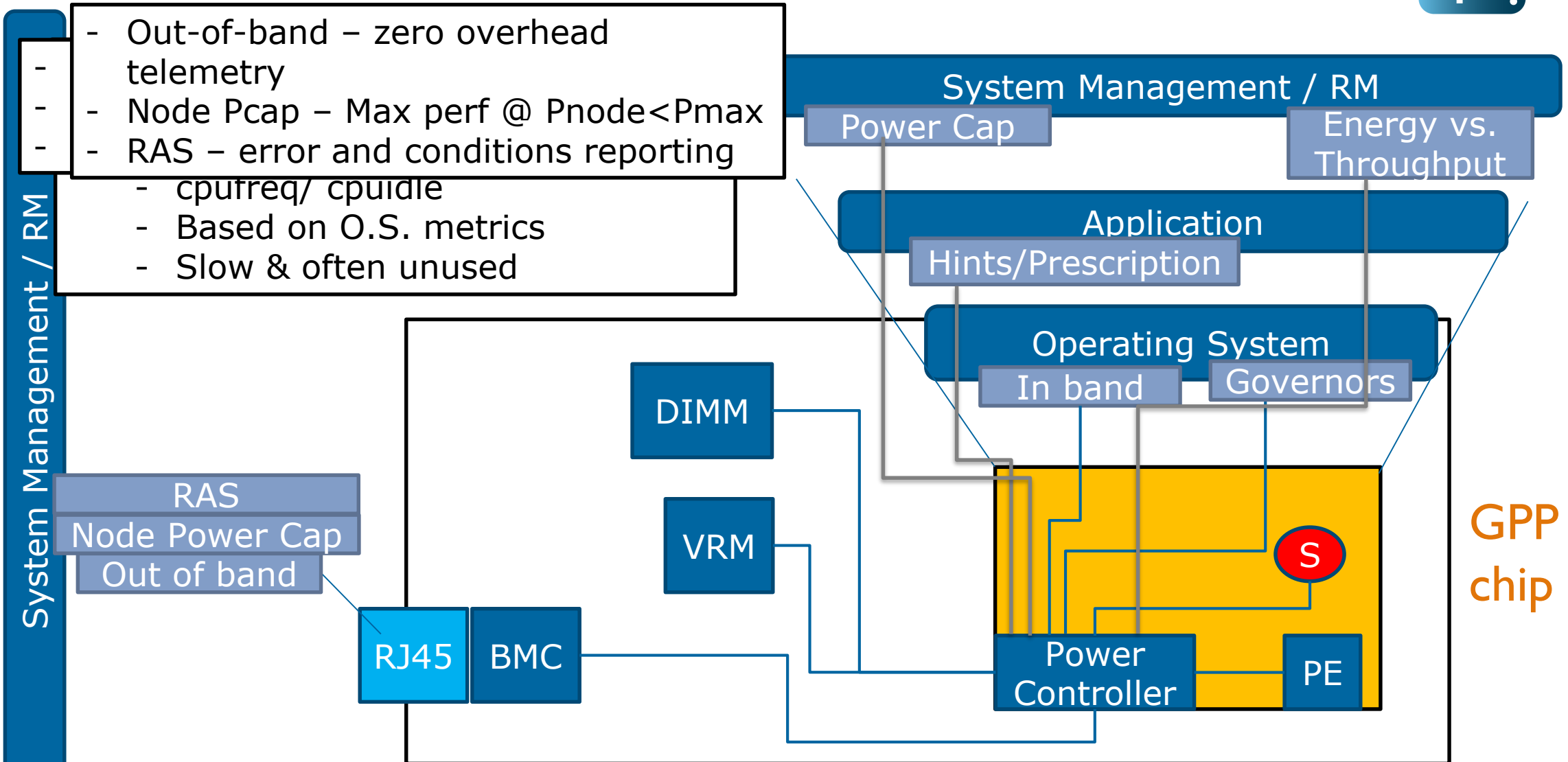
$T_{LC}$ ( $^{\circ}\text{C}$ )	Performance (GFLOPS)	System power (W)	Efficiency (GFLOPS/W)	Efficiency decrease (%)	Possible energy reuse <sup>1</sup> (%)
19	974	359	2.72	–	-1.8
45	985	388	2.54	7.0	7.8
50	985	390	2.53	7.5	9.7
55	981	395	2.49	9.4	10.8
60	976	396	2.46	10.3	12.3
65	969	398	2.44	11.5	14.9

<sup>1</sup> assuming environment temperature  $T_0 = 25\text{ }^{\circ}\text{C}$



**Figure 1.** a) Overview of the experimental setup; b) liquid cooling and system connection  
Ti32 view ( $T_{\max} = 38.3\text{ }^{\circ}\text{C}$ ,  $T_{\min} = 27.6\text{ }^{\circ}\text{C}$ )

# WP3 – POWER MANAGEMENT & CONTROLLER



# POWER MANAGEMENT SOA & REQUIREMENTS

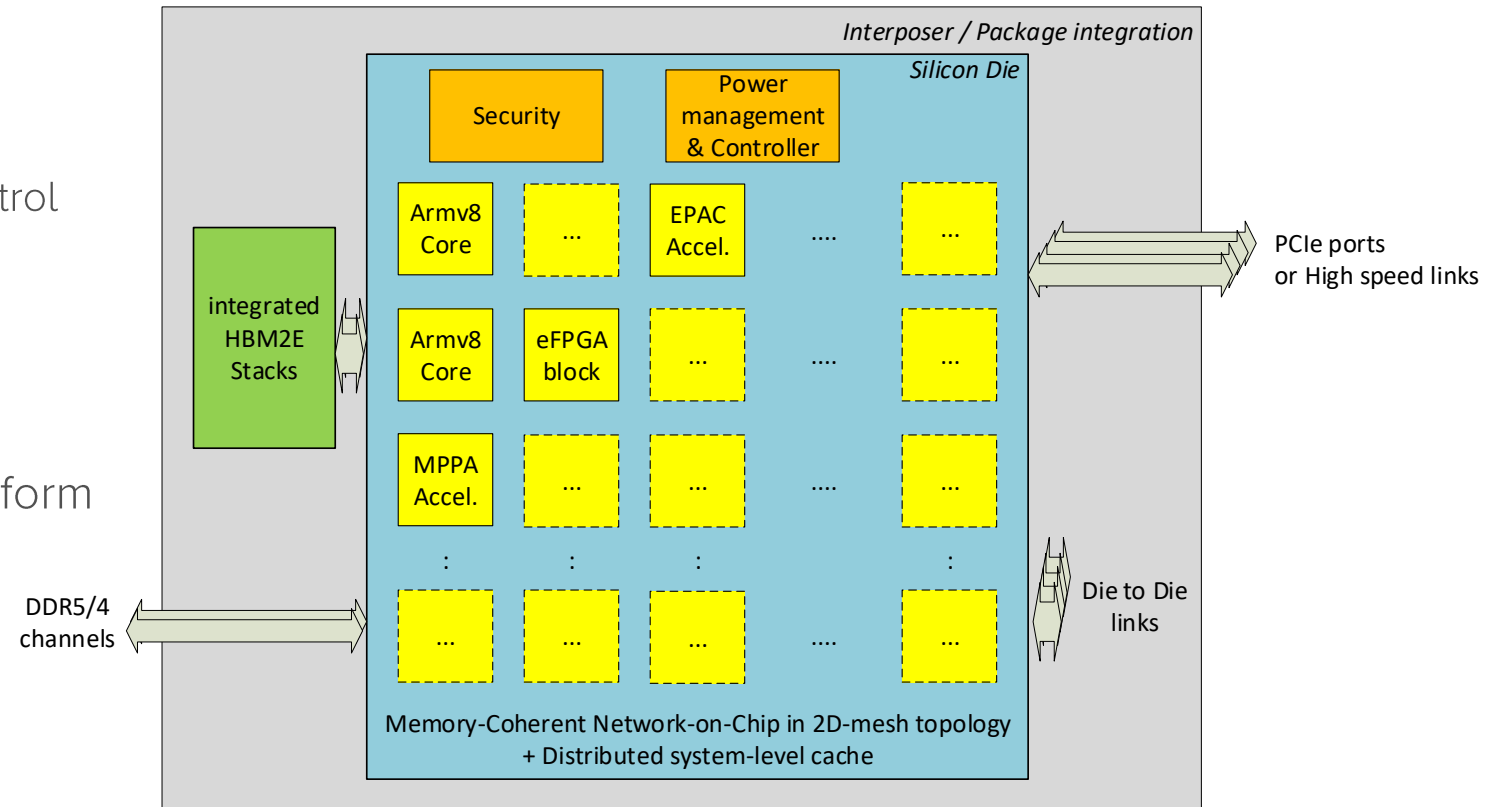
	Intel	IBM	ARM	AMD	Cray	Fujitsu
Monitor (Domain,Granularity)	S, M, A, T <b>1ms</b>	N, S, M, A, T, U <b>500us</b> ,10ms aggregation 16ms for T & U, 100ms aggregation	S, M, T 1-10KHz with SCP	N, S, M, A, T 1 sec (C ), 1ms (G)	N, S, M, A, N OOB (100ms)	N, S, C, M <b>1ms</b> (N), ~ns - model based (C)
Control (Domain,Granularity)	S, M RAPL <b>1ms</b> (in-band), DVFS <b>500us</b>	N, S, M, A 10-100ms	S, M 1-10KHz (100ms to 1s)	N, S, M, A ~secs	N, S, M, A DVFS, RAPL, min-max range, 10- 30s at job launch	S, C, M, <b>DVFS</b> , Decode Width, HBM2 B/W
Interfaces, Tools, etc	RAPL MSRS, msr-safe, libmsr, PAPI, likwid  <i>Source PowerStack / 9</i>	<b>OpenBMC</b> , amester, <b>Memory Map</b>	ACPI, SCP (sys ctrl proc), IPA (intelligent allocator), PAPI	Likwid, PAPI, <b>Memory Map</b>	CapMC, PAPI, Cray <b>BMC interfaces</b>	Power API, PAPI
Socket (S), Core (C), Memory (M), Accelerator (G), Node (N), Utilization (U), Temperature (T)						

EPI power management design is powered UNIBO and targets:

- Support for fine grain power monitoring, and control
- An higher performance power controller capable of supporting advanced power control algorithms.

# GENERAL ARCHITECTURE

- Top level infrastructures
  - Power management & controller
  - Dedicated power management and control network
  - Security
- EPI Power Management Subsystem
- RISC-V ISA, Derived from the PULP platform
- Parallel processor w. DSP extensions
- Open-Source Design



# THE POWER CONTROLLER FIRMWARE

## T control

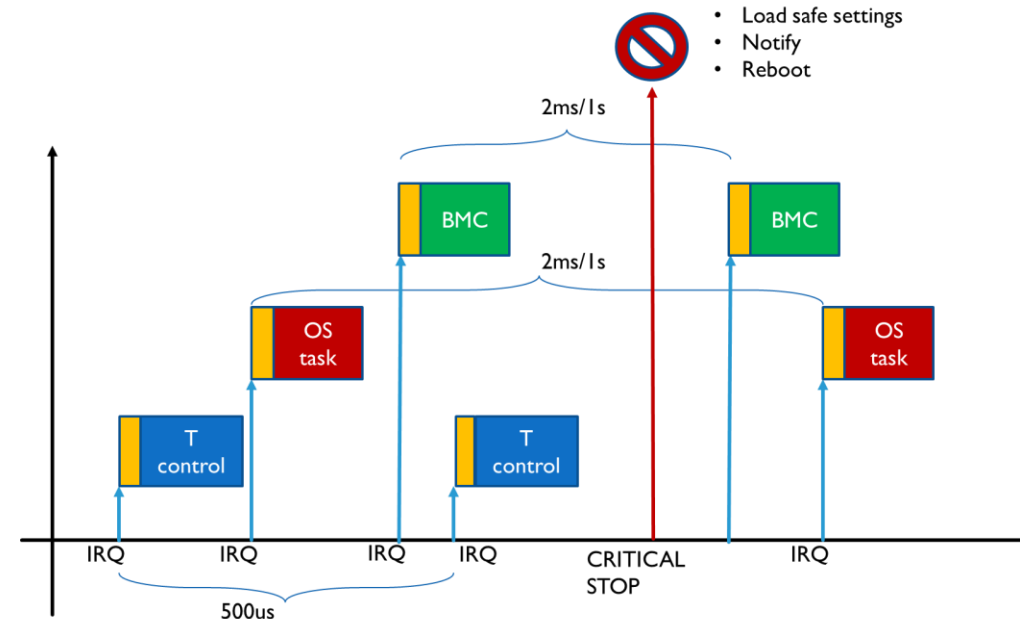
Watchdog reset  
 Write power controller settings  
 Write to internal memory telemetry data  
 Read PVT sensors  
 Read workload from O.S.  
 Read **target P/C state settings, power budget**  
 Read Pending BMC requests  
 Compute controller settings

## PM task

- Read voltage regulator, power, status (VR)
- Power model update

## BMC

- Read pending command queue
- Decode Command/data
- Perform action:
  - Change **target P/C state, power budget**
    - Set pending BMC
  - Ask telemetry data



# CONCLUSIONS

- Power management is a key aspect of HPC processors
- Implemented by mean of a embedded computing subsystem with extensions for interfacing with the power management IPs.
- EPI will leverage a best-in class power management subsystem based on parallel architecture with DSP extensions.