

# NSF/IUCRC CAC PROJECT

---

## MONITORING, VISUALIZING, AND PREDICTING HEALTH STATUS OF HPC CENTERS

Jie Li

Doctoral Student, TTU

06/12/2020

Advisors:

Mr. Jon Hass, SW Architect, Dell Inc.

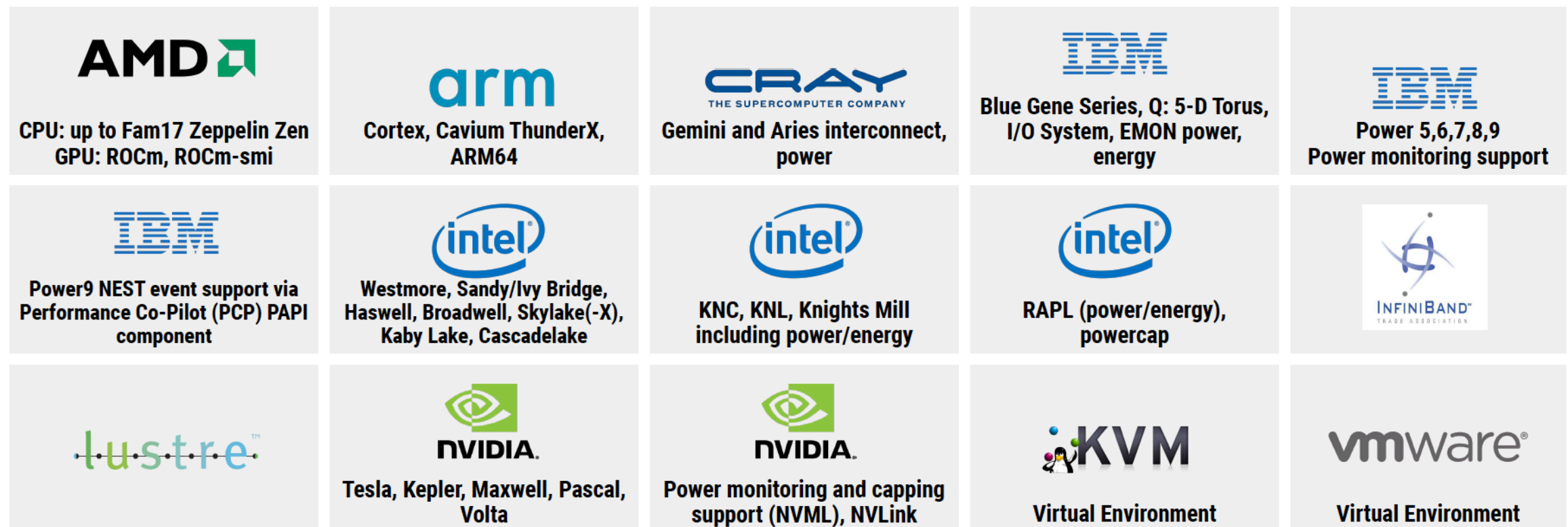
Dr. Alan Sill, Managing Director, HPCC, TTU

Dr. Yong Chen, Associate Professor, CS Dept, TTU

Dr. Tommy Dang, Assistant Professor, CS Dept, TTU

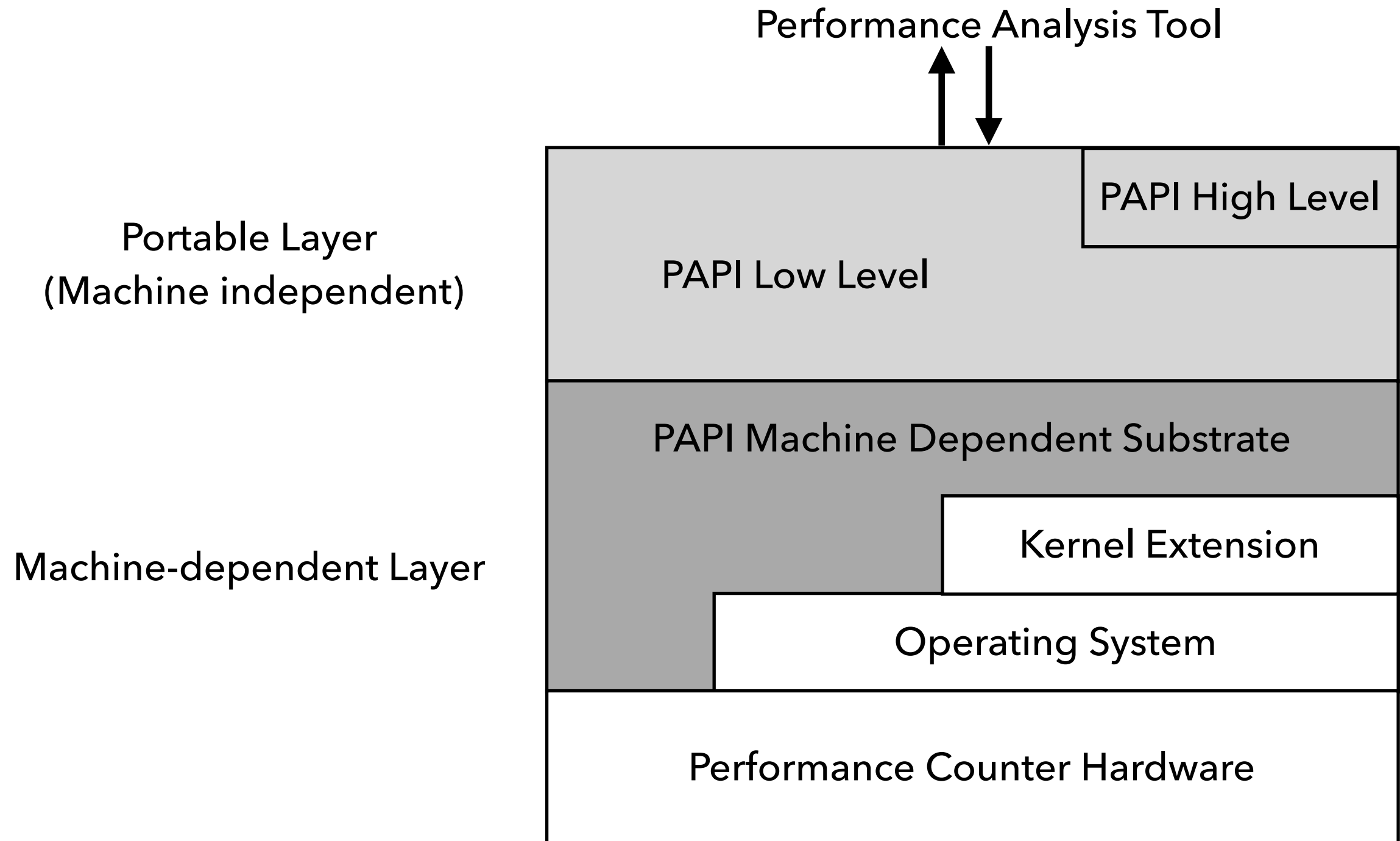
- On-chip performance monitoring (PM) exists on almost every microprocessor.
- PM hardware consists of a small number of registers with connections to various other parts of the chip.
- Two methods of using the PM hardware:
  - Aggregate (direct)
    - Involves reading the counters before and after the execution of a region of code and recording the difference.
    - Permits explicit, highly accurate, fine-grained measurements.
  - Statistical (indirect)
    - The PM hardware is set to generate an interrupt when a performance counter reaches a preset value. This interrupt includes the program counter (PC), the text address at which the interrupt occurred.
    - Facilitates a good high-level understanding of where and why the bottlenecks are occurring.

- Performance Application Programming Interface (PAPI) library.
- Provides **a consistent interface and methodology** for using low-level performance counters in CPUs, GPUs, on/off-chip memory, interconnects, I/O system, and energy/power management.
- Enables monitoring of both types of performance events— **hardware- and software-related events**— in a uniform way.
- Implemented on **a wide variety of architectures and operating systems.**



# PAPI-ARCHITECTURE

---



# PAPI-COMPONENTS

	Component Name	Description
CPU	perf_event	Linux perf_event CPU counters (default)
	perf_event_uncore	Linux perf_event CPU uncore and northbridge (default)
	...	
GPU	cuda	CUDA events and metrics via NVIDIA CuPTI interfaces
	nvml	NVIDIA hardware counters (usage, power, temperature, fan speed, etc)
	...	
Power	powercap	Linux powercap energy measurements
	rapl	Linux RAPL energy measurements
	...	

# PAPI-COMPONENTS

	Component Name	Description
Network	infiniband	Linux Infiniband statistics using the sysfs interface
	net	Linux network driver statistics
I/O	appio	Linux I/O system calls
	io	Linux I/O statistics from /proc/self/io
	lustre	Lustre filesystem statistics
	...	
Other	coretemp	Linux hwmon temperature and other info
	Vmware	Support for VMware (vmguest and pseudo counters)
	...	

### Sample events in the **perf\_event** component:

- **L1 and L2 cache hit rates**, indicate how cache-friendly a program is.
  - L1 cache hit rate:  $1.0 - (PAPI\_L1\_DCM / (PAPI\_LD\_INS + PAPI\_SR\_INS))$
  - L2 cache hit rate:  $1.0 - (PAPI\_L2\_DCM / PAPI\_L1\_DCM)$
- **TLB misses**:  $PAPI\_TLB\_DM$
- **Floating-point operations**:  $PAPI\_FP\_OPS$
- Completed operations per cycle (**IPC**):  $(PAPI\_TOT\_INS / PAPI\_TOT\_CYC)$


DEMO

---

DEMO



- Explore **other events** and understand the meaning.
- Investigate the **overhead** of PAPI.
- Build **metrics collector** utilizing PAPI interface.



**QUESTIONS?/COMMENTS?**