# SUPER Energy-efficiency HPC Research

**SUPER**
INSTITUTE FOR SUSTAINED PERFORMANCE, ENERGY, AND RESILIENCE

*Laura Carrington (lead), Ananta Tiwari*
*UCSD/PMaC*

*Rob Fowler*
*RENCI*

*Dan Terpstra*
*UTK*

<u>Abstract:</u> SUPER's Energy thrust is charged with understanding how computation and communication patterns affect the overall energy requirements of HPC applications. We then leverage this understanding to design software- and hardware-aware optimization techniques that reduce the DOE's HPC energy footprint. Two focus areas have emerged within this thrust: software solutions that provide fine-grained access to the power measurements and energy efficiency research that utilizes these measurements to develop green optimization strategies. We highlight recent accomplishments in each area and present empirical results that illustrate SUPER's contributions in minimizing DOE's HPC energy requirements.

## Energy-Constrained Computation:   Measurement and Adaptation

### Technological and Commercial Imperatives Are Driving the Problem.

Moore's law → The number of gates/chip grows.
- Vendor compete to build and sell increasingly complex high-end chips.

Denard scaling → If geometry, voltages, and clocks scale co-linearly, power density remains constant.
- While these scaling properties worked, high-end chips could run at increasing speed and still be adequately cooled. (This is why Moore's law translated into faster, not just bigger, computers.)
- In the past decade,  gate insulator thickness, voltages stopped scaling.  Chips become power/cooling-limited.
- Instruction-level parallel designs became dominated by "housekeeping" overhead.

The multicore (and "System on a Chip") response.
- Add more power-efficient cores and other units (memory controllers, NICS, GPUs) to chips.
- Cores can still run fast, but sustained speed is limited power/cooling of chip package → Sell chips using "de-rated" specs.
  - Intel E5-2680 is  nominally 2.7GHz but individual cores run easily at  TurboBoost speed of  3.5GHz.
  - AMD "computational sprinting" is similar.
- Entering an era of "dark" or "dim" silicon in which parts of chips are shut down or run slowly.

Impact on HPC of power/cooling constrained computing.
- On chip controllers automatically react to thermal state by adjusting speeds and feeds of major components.  Example, Intel's controller uses a "running average power limit" (RAPL) model.
- The performance advantages of hardware adaptation are significant, but
  → Performance is becoming non-deterministic, varies by : socket,  core,  location in rack, time of day, other programs.
- Measurement, monitoring, and analysis tools are needed to assess the impact.
  - Adaptation requires real time feedback tools.
- System and application code need to interact flexibly with adaptive hardware.

### Challenges and Questions

How should power and thermal effects be incorporated in performance experiments?  Used in tuning exercises?

How can we use the interfaces to the package controllers to better manage power states?
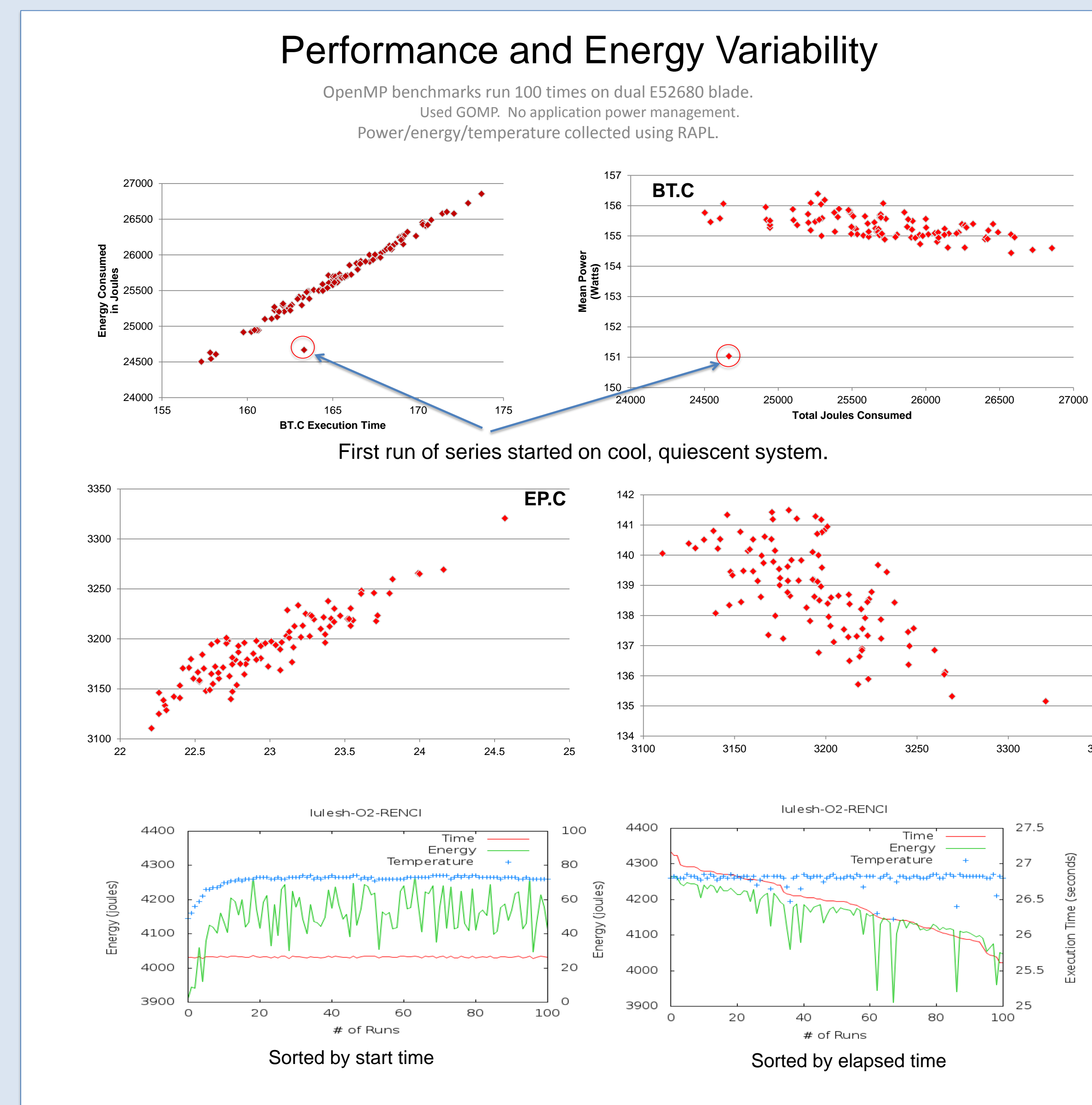- Can this be made portable?

Application and library support:
- Thread management mechanisms and policies to improve throughput while saving power?

The V part of DVFS is becoming less effective and is package-wide.  Core idling and "clock modulation" (in Intel-speak) is effective on a per core basis.  Can we  promote this?

Does better (water) cooling at the package level make the problem easier?

### Performance and Energy Variability

OpenMP benchmarks run 100 times on dual ES2680 blade.
Used GOMP.  No application power management.
Power/energy/temperature collected using RAPL.



First run of series started on cool, quiescent system.

Sorted by start time
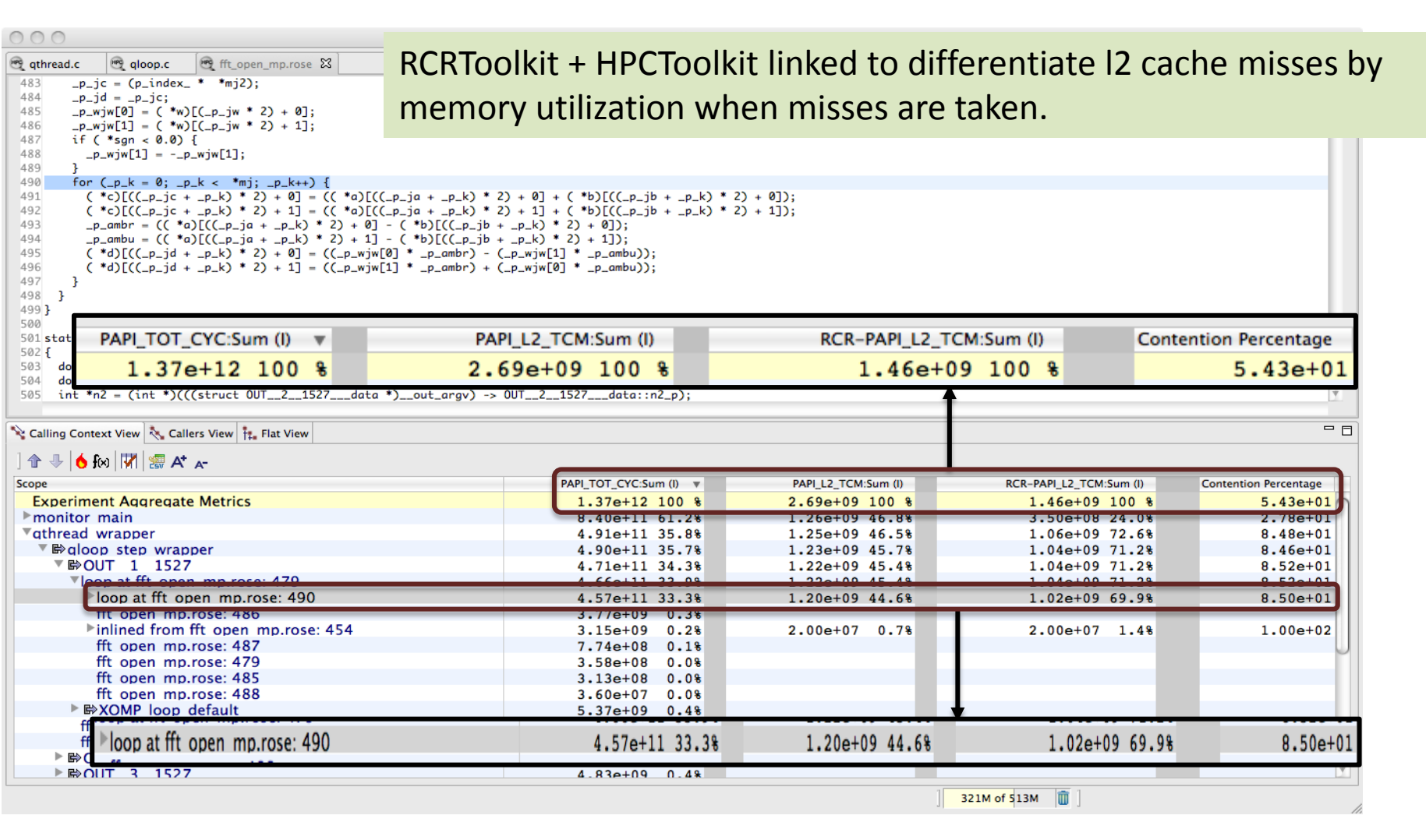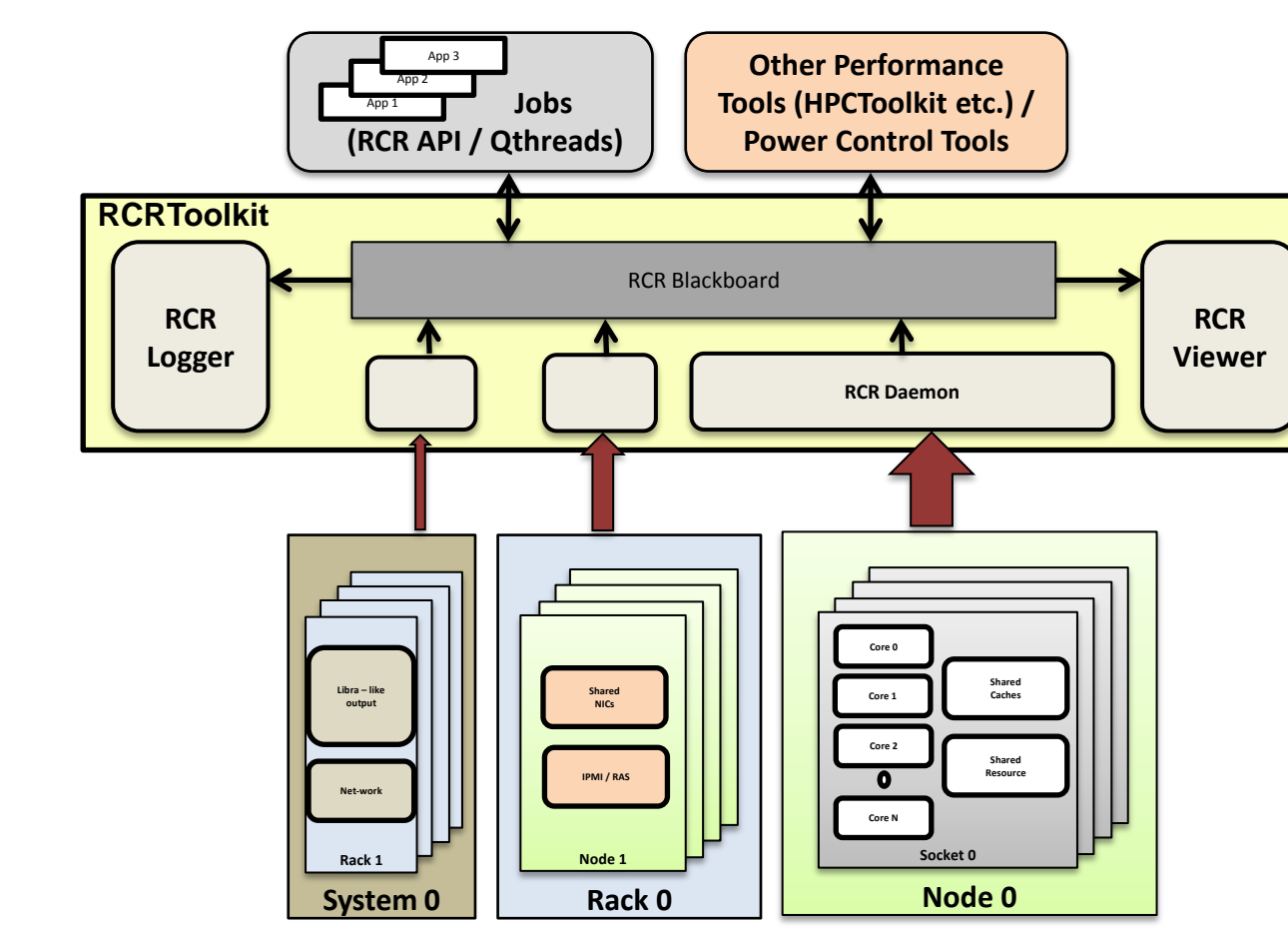
Sorted by elapsed time

### Real Time Measurement/Analysis to Support Adaptation.

Resource-Centric Reflection Tools

An infrastructure for real time collection and analysis of node- and system-wide performance information to guide runtime adaptation.
- Capture "uncore" counter information for shared resources.
  - This  includes power, energy, and temperature data.
- Capture OS and Network information.
- Share across the entire software stack  through a blackboard.



RCRToolkit + HPCToolkit linked to differentiate l2 cache misses by memory utilization when misses are taken.

### Capturing power usage via PAPI

Power measurement is  essential to all energy efficient research because it helps us understand, exploit and limit hardware policies for transparent and adaptive power control. Such understanding allows the implementation of more effective hardware-software co-managed energy optimizations.

PAPI provides access to hardware performance counters, enabling software engineers to measure relationships between software performance and processor or other hardware events.
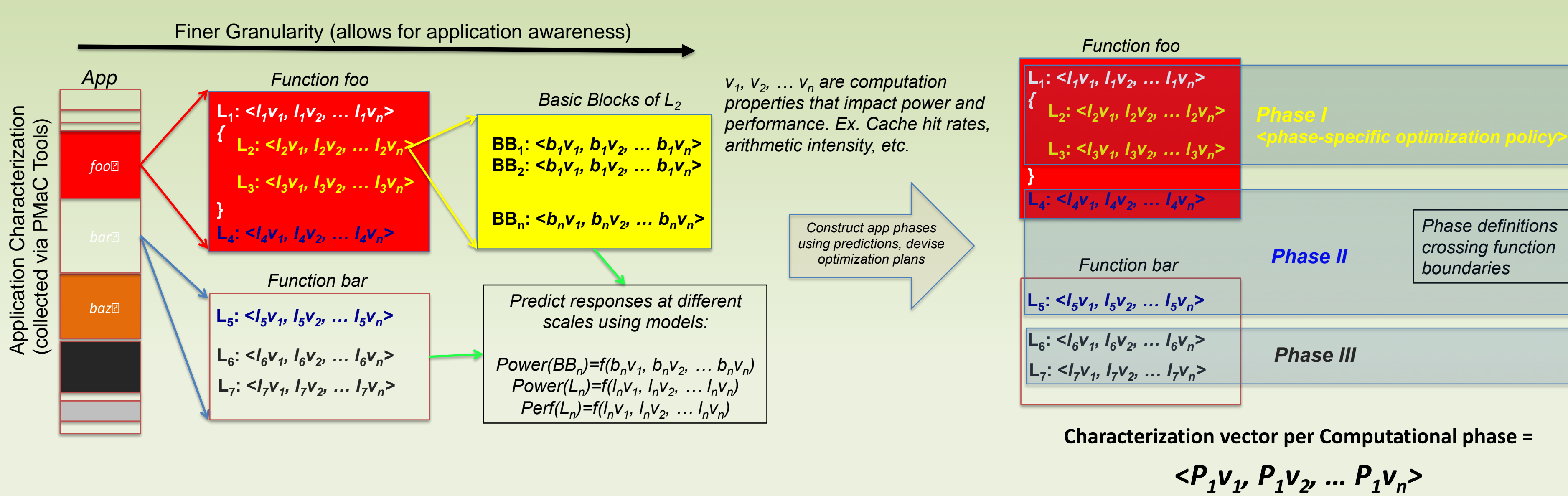
Component PAPI can be extended to allow energy and power measurements without any changes to core infrastructure or measurement tools. Current users of PAPI on HPC systems can now analyze power and energy with little additional effort. Recent and ongoing component development supports measuring energy and power usage via a variety of sources including:

- MIC-side power measurement component
- Host-side MIC power measurement component
- PowerMon2 : RENCI's real-time power monitoring card
- RAPL: Intel's Running Average Power Limit model
- NVML: NVIDIA's Management Library for Tesla cards
- Watts Up?: AC line voltage and current monitor
  - In discussion with IBM for power measurement via EMON2 & PAPI on BG/Q

## Optimizing energy usage via models for the  performance & power of computational phases of HPC applications

### PMaC's Green Queue Framework



Characterization vector per Computational phase = $<P_1v_1, P_1v_2, ... P_1v_n>$

*A fully automated framework that utilizes fine-grained application characterizations and power and performance models  to devise and deploy energy efficient policies*
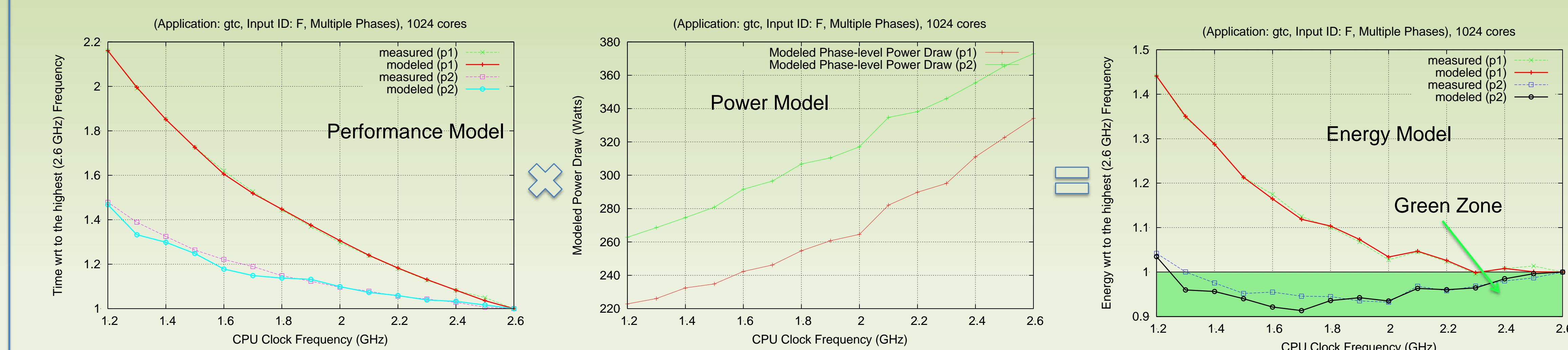
*An application's computational behavior is captured by a series of characterization vectors; these vectors are inputs to power and performance models.*

### Phase Characterization

- Collecting data at various scales (from function to loop to basic blocks) allows flexible definition and construction of phases. A computational phase is defined as a window in the execution of a program where some relevant characteristics of the computation remain fairly unchanged. Phase boundaries can change for different types of optimizations

- Identify phases based on some behavior of interest:
  - Data footprint, L3 Misses, Power draw, Accelerator performance, vectorization

### Case Study: Modeling Two Computational Phases of GTC (1024 core run; Gordon)



These plots show the measured and modeled behavior of two different computational phases of GTC on 1024 cores. We use the phase-level characterization data to predict performance and power responses of the two phases. Predictions for power and performance are then combined to predict energy. We note that for this graph, we normalize the energy required to run each phase at all available frequencies with respect to the energy required to run the phase at the highest frequency. A ratio of less than 1 for a given frequency/phase pair means that we can conserve energy for that phase by running it at that frequency compared to running that phase on the default system frequency. The green zone marks those frequency selections that provide energy savings and illustrate how the models enable fine-grained customized DVFS settings for an application's individual computational phases.