

# RUTGERS

# Abstract

As scientific applications target extreme scales, energy-related challenges are becoming dominating concerns. As a result, it is critical to explore emerging architectures (e.g., with multiple cores and deep memory hierarchies) and applications (e.g., coupled simulation workflows) from an energy perspective and investigate associated overheads and tradeoffs. For example, energy/power-efficiency have to be addressed in combination with quality of solution, performance and reliability, and other objectives, and achieving the desired levels of reduction in power consumptions requires a comprehensive cross-layer and application-aware strategy. In this talk I will explore these issues and will describe recent related research efforts at the Rutgers Discovery Informatics Institute (RDI2).



# RUTGERS

## **Rutgers Discovery Informatics Institute** (RDI<sup>2</sup>) Driving Innovation through Advanced Computing

- Established in March 2012 as New Jersey's Center for Advanced Computation with an overarching goal to create a world-class institute focused on computational and data sciences
- Fundamentally integrate research, education, ACI and industry partnerships to address core CDS&E / BigData challenges
- Broaden industry and academic access to state-of-the-art computing technology and expertise
  - NSF Cloud and Autonomic Computing IUCRC
- Integrate multidisciplinary research with ACI and industry partnerships

## http://rdi2.rutgers.edu



**RDI<sup>2</sup>** 







## **RUTGERS** Managing Energy/Power - A Wide Range of Techniques Aggressive component-level power management - CPU (e.g., DVFS), memory, disk, NIC, etc. Run-time power management - System level, exploiting slack in MPI programs, etc. Application/workload power management - Workload profiling/characterization, consolidation, etc. - Application, algorithm adaptations Autonomic policy adaptation - Characterize operational state and adjust management goals · Cooling and thermal management React to thermal hotspots, proactive workload placement ...., etc. Can we use these together in a cross-layer, coordinated, and consistent manner?



| RUTGERS RUTGERS  |
|--|
| Green HPC: Landscape (II/III)  |
| Dynamic scaling (contd.)   |
| <ul> <li>Scaling other components/subsystems</li> </ul>  |
| <ul> <li>Dynamic memory frequency</li> </ul>   |
| <ul> <li>Delaluz et al. [IEEE TC'01, DAC'02], Fradj et al. [DSD'06],<br/>Bianchini et al. [ASPLOS'11]</li> </ul>               |
| <ul> <li>Storage subsystem Multi-speed disks and RAIDS</li> <li>– Rotem et al. [IPDPS'09], Pinheiro et al. [ICS'04]</li> </ul> |
| <ul> <li>NVRAM – PCM/STTM, memristors, flash-based</li> </ul>  |
| <ul> <li>Caulfield et al. [ASPLOS'09, MICRO'10], Bianchini et al.<br/>[SC'11]</li> </ul>                                       |
| <ul> <li>Flash-based SSD</li> </ul>  |
| <ul> <li>Urgaonkar et al. [OSDI'08]</li> </ul>   |
|  |
|  |











# Proactive cluse application behaviors to enable aggressive and proactive power management Reactive power management is not always optimal and can result in large overheads Proactive adaptation of resources based on application behavior and requirements (e.g., subsystems demand over time) Maintain performance to the extent possible Overall approach [HiPC10, HiPC11] Application characterization (subsystem demand) Empirical quantification of possible power savings (upper bound) Proactive subsystem power management at system level









|           |        |                    |                        |              |         | Energy  | Savings |                                  |
|-----------|--------|--------------------|------------------------|--------------|---------|---------|---------|----------------------------------|
| Benchmark | DVFS   | Run Time (s)       | Energy (J)             | CPU          | Memory  | Disk    | NIC     | Total (J) / %                    |
| HPL       | ×      | 1,382 s<br>1,383 s | 298,546 J<br>292,824 J | -<br>5,722 J | 1,380 J | 5,338 J | 240 J   | 6,958 J 2.33%<br>12,680 J 4.33%  |
| b_eff_io  | ×      | 1,206 s<br>1,212 s | 164,224 J<br>161,460 J | -<br>2764 J  | 5,297 J | -       | 1,124 J | 6,421 J 3.90%<br>9,185 J 5.69%   |
| bonnie++  | ×      | 1,247 s<br>1,248 s | 190,613 J<br>187,533 J | -<br>3,080 J | 3,263 J | 574 J   | 3,841 J | 7,678 J 4.03%<br>10,758 J 5.73%  |
| TauBench  | ×      | 1,134 s<br>1,136 s | 251,904 J<br>244,473 J | -<br>7,431 J | 3,377 J | 7,297 J | 1,979 J | 12,653 J 5.02%<br>20,084 J 8.21% |
| FFTW      | ×<br>√ | 1,052 s<br>1,055 s | 198,621 J<br>193,146 J | -<br>5,475 J | 2,112 J | 6,927 J | 297 J   | 9,336 J 4.70%<br>14,811 J 7.67%  |
|           |        |                    |                        |              |         |         |         |                                  |

| Benchmark | Configuration                 | Run Time (s)                  | %                     | Energy (J)                          | %                | EDP                                       | %                 |
|-----------|-------------------------------|-------------------------------|-----------------------|-------------------------------------|------------------|---|-------------------|
| HPL       | Reference<br>PAPM<br>PAPM+SSD | 1,383 s<br>1,385 s<br>1,385 s | -<br>+0.14%<br>+0.14% | 292,824 J<br>287,906 J<br>281,559 J | -1.67%<br>-3.84% | 404,975,592<br>398,749,810<br>389,959,215 | -1.53%<br>-3.70%  |
| b_eff_io  | Reference<br>PAPM<br>PAPM+SSD | 1,212 s<br>1,217 s<br>1,134 s | +0.41%                | 161,460 J<br>157,335 J<br>143,768 J | -2.55%           | 195,689,520<br>191,476,695<br>163,032,912 | -2.15%            |
| bonnie++  | Reference<br>PAPM<br>PAPM+SSD | 1,248 s<br>1,249 s<br>1,169 s | -<br>+0.08%<br>-6.33% | 187,533 J<br>182,904 J<br>168,606 J | -2.47%           | 234,041,184<br>228,447,096<br>197,100,414 | -2.39%<br>-15.78% |
| TauBench  | Reference<br>PAPM<br>PAPM+SSD | 1,136 s<br>1,139 s<br>1,137 s | -<br>+0.26%<br>+0.08% | 244,473 J<br>236,496 J<br>229,446 J | -3.26%<br>-6.14% | 277,721,328<br>269,368,944<br>260,880,102 | -3.00%<br>-6.06%  |
| FFTW      | Reference<br>PAPM<br>PAPM+SSD | 1,055 s<br>1,057 s<br>1,051 s | -<br>+0.19%<br>-0.38% | 193,146 J<br>187,677 J<br>177,071 J | -2.83%           | 203,769,030<br>198,374,589<br>186,101,621 | -2.64%            |



# Rutgers

## **Application-aware Power Management:** Lessons Learned

• Application-aware power control can lead to improved energy efficiency without significant performance penalty

RDI<sup>2</sup>

- Power management at the subsystem level can not be neglected
- Low power devices and upcoming technologies (e.g., SSD, storage-level memory) can significantly increase energy efficiency
- However, power management at a single level is not enough
- A cross-layer approach is required!















| RUTGERS   | VV       |              | E                      |            |             | Rutgers Discovery Informatics Institute |  |  |  |  |
|---|----------|--------------|------------------------|------------|-------------|---|--|--|--|--|
| <ul> <li>Use Case 2: Sobel Filter Application</li> <li>Sobel filter has two different phases</li> <li>Two studied strategies <ul> <li>Using the same policy during the whole application execution</li> <li>Using hints during the initial phase (e.g., PM_POWER)</li> </ul> </li> <li>Hints avoid time delay and maintains energy savings <ul> <li>The right power mode cannot be identified at runtime or system</li> </ul> </li> </ul> |          |              |                        |            |             |   |  |  |  |  |
|   |          |              |                        |            |             | Power over time (Sobel)                 |  |  |  |  |
| Language extension<br>(user hints)  | Sobel    | - fix policy | Sobel – dynamic policy |            |             | 100<br>80<br>60                         |  |  |  |  |
|   | Time (%) | Energy (%)   | Time (%)               | Energy (%) | 8 20<br>9 0 |   |  |  |  |  |
| PM_PERFORMANCE  | 100.0    | 100.0        |                        |            |             | (a) PM_POWER policy, base configuration |  |  |  |  |
| PM_POWER  | 147.6    | 75.0         | 101.8                  | 79.4       | er (W)      |   |  |  |  |  |
| PM_AGGRESSIVE_POWER   | 193.8    | 72.8         | 101.9                  | 73.7       | OC Powe     |   |  |  |  |  |
| PM_CONSERVATIVE   | 118.0    | 75.7         | 101.5                  | 73.7       | S           | 0 100 200 300 400 500 600 700 800       |  |  |  |  |
|   |          |              |                        |            |             | (b) PM_POWER policy, application-aware  |  |  |  |  |















### **RUTGERS Data Staging over Deep Memory Hierarchy Motivation** · Small DRAM capacity per core - even aggregated memory on dedicated nodes can hardly keep all coupled data (given the ratio of resource allocations for compute nodes and dedicated nodes) **Hybrid Staging** Simulation cores Primary resources · Spans horizontally across Data staging and processing cores compute nodes of both primary Secondary resources Data and secondary resources Spans vertically across the multi-Node-level processor cores level memory hierarchy, e.g. Node-leve DRAM/NVRAM/SSD, to extend memory storage Staging Abstraction the capacity of in-memory data staging HPC core computation resource







# RUTGERS

## **Power Behavior of In-situ Analytics Pipeline**

(with J. C. Bennett, H. Kolla, J.Chen, T. Bremer, A. G. Landge, A. Gyulassy, P. McCormick, S. Pakin, V. Pascucci)

- Motivation. Coupled simulation workflows use online data processing to reduce data movement. Need to explore energy/performance tradeoffs.
- Use case. Combustion simulation workflow with an in-situ data analytics pipeline.











