

# Run-Time DFS/DCT Optimization for Power-Constrained HPC Systems\*

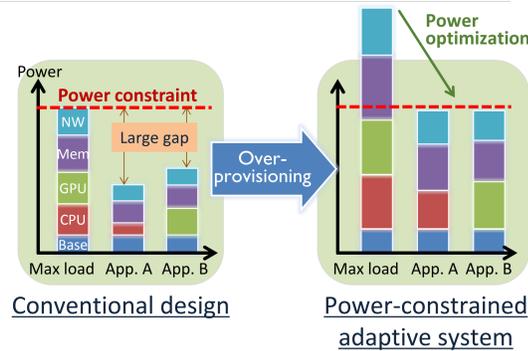
Ikuo Miyoshi<sup>1</sup>, Shinobu Miwa<sup>2</sup>, Koji Inoue<sup>3</sup> and Masaaki Kondo<sup>4</sup>

<sup>1</sup>Fujitsu Limited, <sup>2</sup>The University of Electro-Communications, <sup>3</sup>Kyushu University, <sup>4</sup>The University of Tokyo

\*This work was partially supported by JST under CREST.

## Background

- For development of exascale HPC systems, power consumption is one of the major design constraints. HW-overprovisioning is an effective approach to build systems under such constraints.
- Since job throughput of HW-overprovisioned systems is restricted by power limit, “performance-per-watt” has become a more important indicator of efficient system use.



## Motivation

- Because application developers are not expected to pay much attention to “performance-per-watt”, power optimizations should be done by system providers and system administrators.
- Although optimization techniques for power and energy, such as DFS and DCT, have been studied well, practical use is limited to a DB-based approach, which records and reuses the relationship between a job’s execution time and its locked CPU frequency.
- Conducting DFS/DCT optimization during the execution of a job is simple to use, adaptive to the environment of executing systems, and possibly robust in terms of CPU’s manufacturing variability.

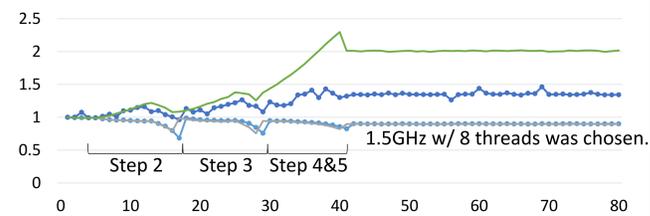
## Optimization Algorithm

- Assumptions
  - ✓ Target application does iterative computation, such as time integration and iterative solver, and can be divided into several regions based on computational characteristics.
  - ✓ From the application users’ viewpoint, performance fluctuation by ~10% is ignored or acceptable.
- The optimization target is a DFS/DCT configuration of max “performance-per-watt” for each region in the range of acceptable performance degradation.
- Optimization steps
  - For use as reference performance, observe the performance during initial iteration(s)
  - Under the “Turbo” frequency, search the # of threads in decreasing order for the best efficiency
  - Under the base frequency, search the # of threads in the same way as Step 2
  - With the # of threads chosen by Step 2&3 as the best efficiency in the range of acceptable performance degradation, search CPU frequency in decreasing order for the best efficiency
  - When performance degradation exceeds the acceptable range, continue the search after adding one more thread
- Calculation of efficiency
  - When measuring by RAPL, efficiency = “measured energy for the reference performance” / “measured energy for a configuration”.
  - When estimating without RAPL, “measured energy” is replaced with “estimated energy”; “estimated energy” = “CPU cycles” \* (1 + (“# of threads” - 1) \* “correction factor”).

The “correction factor”, which is 0.06 for Haswell and 0.13 for Sandy Bridge, was derived from the results of the STREAM benchmark by regression analysis.

| CPU frequency | # of threads |      |      |      |      |      |      |      |      |      |      |      |      |      |      |      |
|---------------|--------------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
|               | 1            | 2    | 3    | 4    | 5    | 6    | 7    | 8    | 9    | 10   | 11   | 12   | 13   | 14   | 15   | 16   |
| 1.2           | 0.24         | 0.44 | 0.61 | 0.77 | 0.91 | 1.04 | 1.15 | 1.26 | 1.34 | 1.39 | 1.42 | 1.43 | 1.45 | 1.44 | 1.43 | 1.42 |
| 1.3           | 0.26         | 0.46 | 0.64 | 0.80 | 0.94 | 1.06 | 1.18 | 1.27 | 1.34 | 1.39 | 1.41 | 1.42 | 1.41 | 1.40 | 1.40 | 1.38 |
| 1.4           | 0.27         | 0.48 | 0.66 | 0.83 | 0.96 | 1.10 | 1.21 | 1.30 | 1.36 | 1.38 | 1.40 | 1.41 | 1.40 | 1.38 | 1.37 | 1.36 |
| 1.5           | 0.28         | 0.50 | 0.69 | 0.85 | 1.00 | 1.12 | 1.23 | 1.32 | 1.36 | 1.38 | 1.38 | 1.37 | 1.37 | 1.35 | 1.34 | 1.32 |
| 1.6           | 0.29         | 0.52 | 0.71 | 0.87 | 1.01 | 1.15 | 1.25 | 1.33 | 1.36 | 1.37 | 1.37 | 1.36 | 1.35 | 1.33 | 1.32 | 1.30 |
| 1.7           | 0.30         | 0.54 | 0.73 | 0.89 | 1.03 | 1.16 | 1.26 | 1.33 | 1.35 | 1.35 | 1.34 | 1.34 | 1.32 | 1.31 | 1.29 | 1.27 |
| 1.8           | 0.31         | 0.55 | 0.73 | 0.90 | 1.03 | 1.15 | 1.25 | 1.30 | 1.32 | 1.31 | 1.30 | 1.29 | 1.27 | 1.25 | 1.23 | 1.20 |
| 1.9           | 0.32         | 0.56 | 0.75 | 0.90 | 1.04 | 1.16 | 1.23 | 1.28 | 1.29 | 1.28 | 1.27 | 1.25 | 1.24 | 1.21 | 1.19 | 1.17 |
| 2.0           | 0.33         | 0.57 | 0.76 | 0.91 | 1.05 | 1.16 | 1.23 | 1.27 | 1.26 | 1.25 | 1.22 | 1.21 | 1.19 | 1.16 | 1.14 | 1.11 |
| 2.1           | 0.34         | 0.58 | 0.77 | 0.92 | 1.05 | 1.17 | 1.22 | 1.25 | 1.24 | 1.23 | 1.21 | 1.20 | 1.18 | 1.15 | 1.13 | 1.11 |
| 2.2           | 0.34         | 0.59 | 0.77 | 0.93 | 1.06 | 1.16 | 1.21 | 1.22 | 1.21 | 1.20 | 1.19 | 1.17 | 1.15 | 1.12 | 1.10 | 1.08 |
| 2.3           | 0.35         | 0.59 | 0.78 | 0.93 | 1.06 | 1.16 | 1.20 | 1.21 | 1.21 | 1.19 | 1.17 | 1.15 | 1.12 | 1.10 | 1.08 | 1.05 |
| Turbo         | 0.38         | 0.59 | 0.74 | 0.86 | 0.97 | 1.04 | 1.08 | 1.14 | 1.15 | 1.13 | 1.11 | 1.09 | 1.06 | 1.04 | 1.01 | 1.00 |

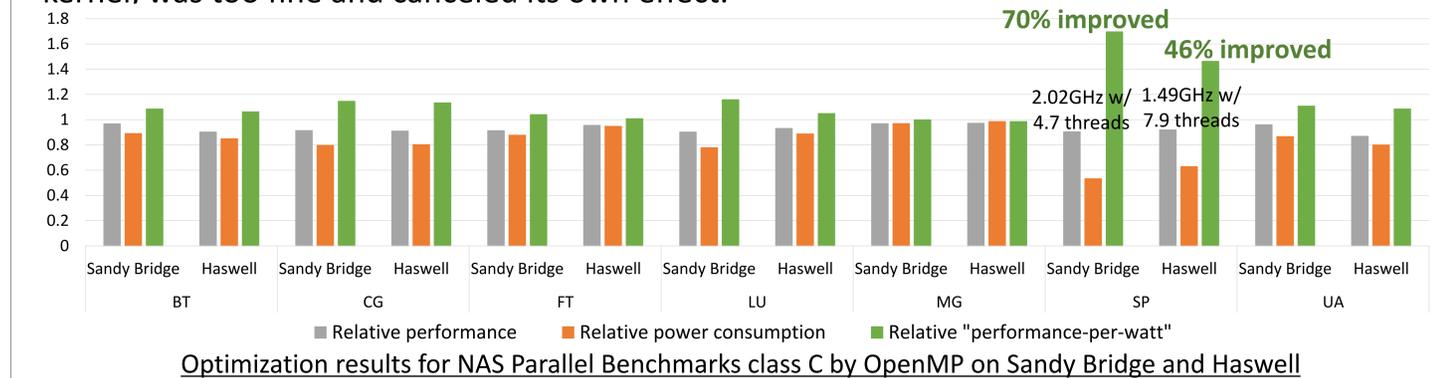
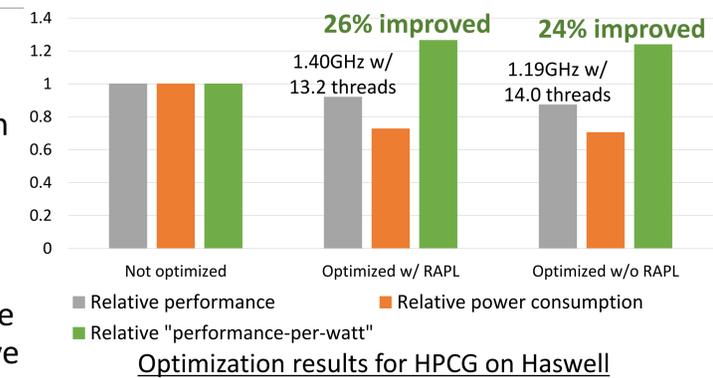
Relative “performance-per-watt” of HPCG on Haswell



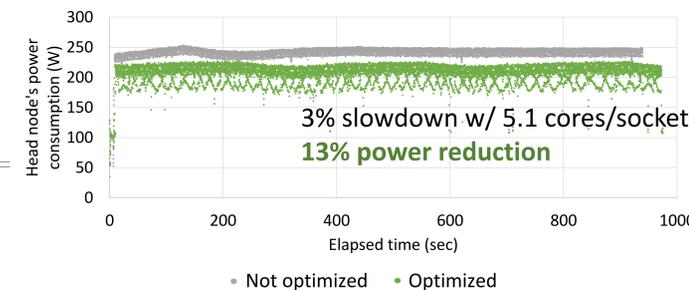
Optimization behavior for ComputeSPMV region of HPCG on Haswell

## Evaluation

- For the HPCG benchmark on Haswell, the optimization chose 1.19GHz with 14.0 threads on average without the RAPL measurement, which improved the “performance-per-watt” by 24%.
- Among 7 iterative kernels in the NAS Parallel Benchmarks, 6 kernels showed improvement in “performance-per-watt” up to 70%. However, the optimization for MG, which is a memory-intensive kernel, was too fine and canceled its own effect.



- Results on the 10 nodes of Sandy Bridge running NICAM-DC-MINI, which is a mini app for Post-K development, indicated the possibility to increase headroom for power-shifting between jobs.



## Related Work

- READEX[1] proposes a methodology for automatic tuning to improve energy efficiency. In contrast to our approach, pre-execution of applications with a representative dataset is necessary for this methodology.
- GEOPM[2] is an open source run-time framework for researching energy management solutions. By changing the RAPL settings, its power-balancing plug-in improves a job’s execution time, whose performance was originally degraded by constant power-capping by RAPL.

## Summary and Future Work

- DFS and DCT techniques were integrated into a run-time “performance-per-watt” optimization and its “on-the-fly” optimization demonstrated its effectiveness.
- Further studies are planned as follows: confirming the applicability to new generation Xeons, studying performance controllability by the degree of acceptable performance degradation, drafting of a concrete scenario of power-shifting between jobs, and so on.

## References

- Oleynik, Yury, et al. “Run-time exploitation of application dynamism for energy-efficient exascale computing (READEX).” *Computational Science and Engineering (CSE), 2015 IEEE 18th International Conference on*. IEEE, 2015.
- Eastep, Jonathan, et al. “Global Extensible Open Power Manager: A Vehicle for HPC Community Collaboration on Co-Designed Energy Management Solutions.” *International Supercomputing Conference*. Springer, Cham, 2017.