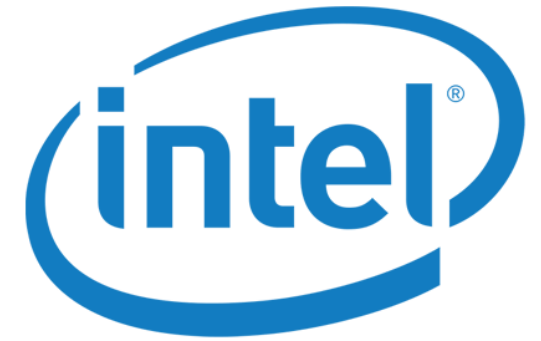


# Energy Efficiency Optimization with GEOPM

---

<http://geopm.github.io/geopm>



Jonathan Eastep [jonathan.m.eastep@intel.com]

Principal Engineer, PhD

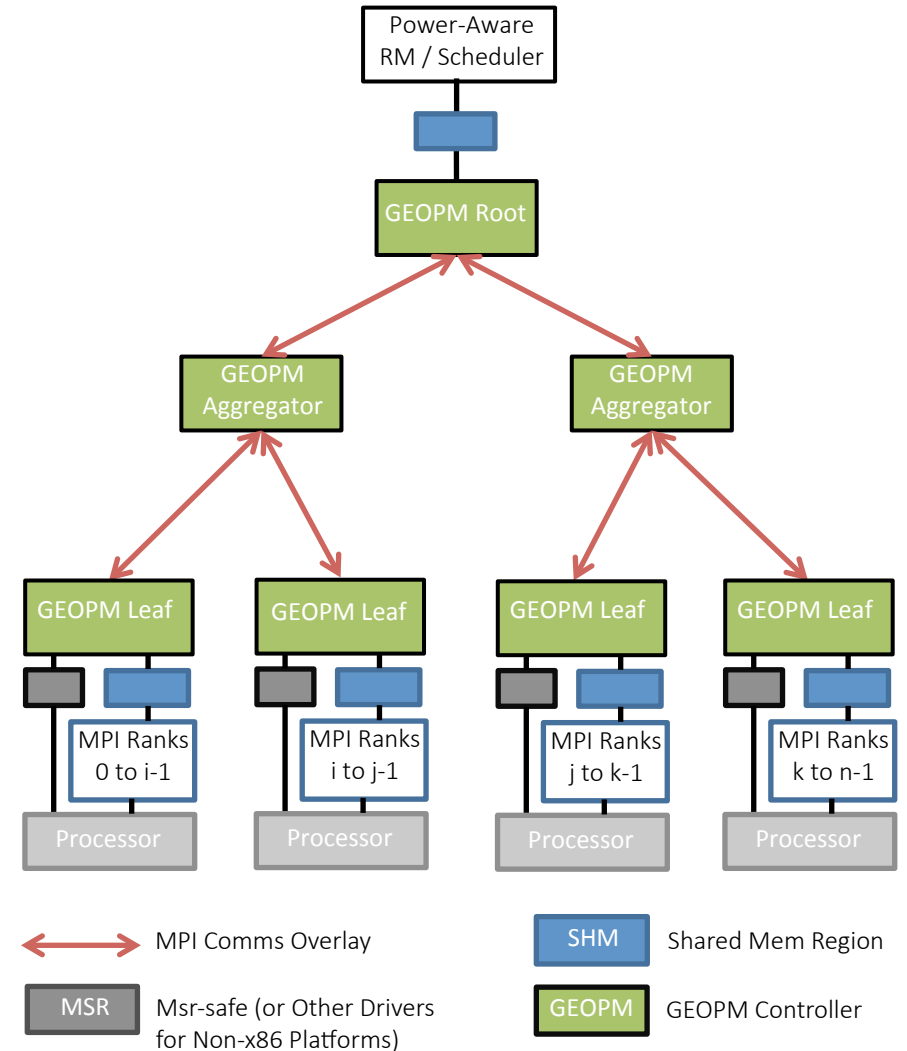
12 November 2017

# Outline

- GEOPM overview, use-cases, and status
- Intel, LRZ, LLNL, Argonne collaboration
- New experimental results
- Collaboration next steps

# G lobal E xtensible O pen P ower M anager

- Runtime for in-band power management and optimization
  - On-the-fly monitoring of HW counters & application profiling
  - Feedback-guided optimization of HW control knob settings
- Open source software (flexible BSD three clause license)
- Extensible through plugin architecture
  - Add new energy optimization strategies
  - Add support for new architectures beyond x86 (truly open)
- Designed for holistic optimization
  - Job-wide global optimization of HW control knob settings
  - Application-awareness for max speedup or energy savings
- Scalable via distributed tree-hierarchical design, algorithms



Project url: <http://geopm.github.io/geopm>

Contact: [jonathan.m.eastep@intel.com](mailto:jonathan.m.eastep@intel.com)

# GEOPM Use Cases

- Turn-key (requires no app annotation):
  - Automatic online job profiling
    - Node-level: trace samples of processor counters and correlate HW activity to each OpenMP parallel region
    - Job-level: aggregate the energy counters across all job compute nodes to monitor overall job power or energy
  - Automatic offline or online optimization
    - Will talk more about this today
  - Offline visualization of profile data
    - Python scripts leveraging pandas for data analysis
    - Helpful for debugging new plugins or understanding how they optimize energy or runtime
    - Plot trace of plugin decisions and data they're based on
- Advanced (requires using GEOPM profiling API for app annotation):
  - Automatic online rebalancing of power & perf among nodes
    - Purpose: accelerate critical path nodes in MPI bulk-synchronous applications
    - Refer to ISC'17 paper on GEOPM by Eastep et al. for more info
    - Note: work in progress to make the annotation automatic / turn-key too

# GEOPM Community (1)

Institution	Principal Investigator	Project Name	Project Scope	Contribution Type	Time Span	Quality Level	Funded?
Argonne	Kalyan Kumaran Vitali Morozov	CORAL	1. GEOPM 1.0 product development	Sponsor	Q2'15 – Q4'17	Product	Yes
* IBM STFC – Hartree	Vadim Elisseev Milos Puzovic Neil Morgan		1. GEOPM port to Power8 + NVLink 2. Integration of GEOPM with EAS	Contributor	Q4'16 – TBD	Research	Yes
LLNL	Barry Rountree Aniruddha Marathe	CRADA	1. Integration of GEOPM and Conductor runtime tech 2. Studies to motivate GEOPM/HW codesign	Contributor	Q3'13 – TBD	Research	Yes
* LLNL U. of Arizona Argonne	Tapasya Patki Dave Lowenthal Pete Beckman	ECP PS ECP Argo- GRM	1. Exascale power stack leveraging GEOPM 2. Integration of GEOPM + Caliper framework 3. Integration of GEOPM with EAS 4. Port of GEOPM to non-x86 architecture	Contributor	Q1'17 – Q4'19	Near-Product	Yes
LRZ	Dieter Kranzlmüller Herbert Huber Torsten Wilde		1. Energy optimization plugin for GEOPM 1.0 2. Power ramp limiting plugin for GEOPM 1.x	Contributor	Q3'17 – Q4'20	Near-Product	Yes
Sandia	James Laros Ryan Grant	Power API	1. GEOPM and Power API xface compatibility 2. Power API community WG kickoff at Intel	User	Q4'14 - TBD	Industry Standard	Yes

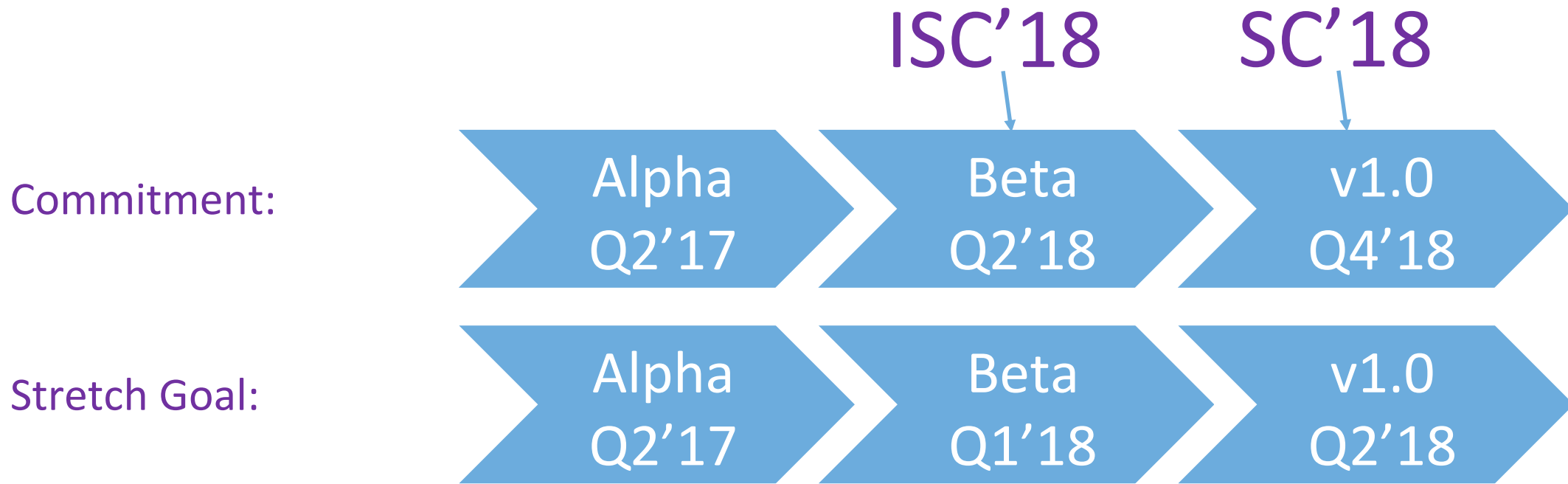
\* = collaborator will be sharing their GEOPM usages and experiences at SC17: BoF on Power API, GEOPM, and Redfish

# GEOPM Community (2)

Institution	Principal Investigator	Project Name	Project Scope	Contribution Type	Time Span	Quality Level	Funded?
Argonne	Kalyan Kumaran Vitali Morozov Kevin Harms		<ol style="list-style-type: none"> <li>1. GEOPM &gt;1.0 feature development</li> <li>2. GEOPM enablement for system power capping + EAS</li> <li>3. Studies to motivate GEOPM / hardware codesign</li> </ol>	Sponsor	Q1'18 – Q4'21	Product	WIP
CINECA	Carlo Cavazzoni		<ol style="list-style-type: none"> <li>1. System level runtime for power capping and power ramp limiting leveraging GEOPM</li> </ol>	Contributor	Q2'18 – Q1'21	Near-Product	WIP <sup>†</sup>
IT4I	Lubomir Riha		<ol style="list-style-type: none"> <li>1. GEOPM ports to OpenPOWER and ARM</li> <li>2. Extensions to GEOPM application profiler</li> <li>3. Integration of GEOPM with EAS</li> </ol>	Contributor	Q2'18 – Q1'21	Near-Product	WIP <sup>†</sup>
E4	Fabrizio Magugliani		<ol style="list-style-type: none"> <li>1. GEOPM port to OpenPOWER</li> </ol>	Contributor	Q2'18 – Q1'21	Near-Product	WIP <sup>†</sup>
PNNL	Leon Song		<ol style="list-style-type: none"> <li>1. GEOPM extensions to tune new HW control knob settings</li> <li>2. GEOPM extensions for coordinated tuning of SW params and HW control knob settings</li> </ol>	Contributor	Q1'19 – Q4'20	Research	WIP <sup>†</sup>

† = letter of intent or equivalent in-hand (non-binding)

# GEOPM Release Schedule



**Announcement:** OpenHPC application has been submitted. Under consideration.



**TOSS 3.x**

# Intel, LRZ, LLNL, and Argonne Collab

- Present focus:
  - Develop new techniques in GEOPM to improve energy-to-solution w/ modest impact to time-to-solution
  - Work toward integrating GEOPM and these techniques into production systems at LRZ, LLNL, and Argonne
- Energy optimization approach:
  - Automatically adapt processor core frequency based on characteristics of individual applications
  - Run the cores slower to save energy if the app is bottlenecked by the memory or network subsystems
  - Adapt frequency at a fine-grained timescale: different frequency for each computational phase in the app
  - It's critical to adapt to phases since each phase's runtime can have wildly different frequency sensitivity
- Innovation vs prior art:
  - Novel per-phase-adaptation enables bigger energy savings and lower impact to time-to-solution

Big wins are possible: up to **16.5% energy savings** at **0.3% increase in time-to-solution**

Status: trending to complete this work in time for GEOPM Beta release



# Two Techniques Under Development

- Start simple: **offline automatic** frequency optimization via scripts and Decider plugin
  - GEOPM plugin takes in each phase's frequency as an input, applies the inputted frequency upon phase entry
  - Offline, for a given application, scripts sweep over plugin frequency configurations to characterize the app's phases, identify *best* frequencies (min energy @ < 10% execution time impact)
  - Each time app is launched, it runs with best phase frequency configurations identified by the scripts
  - Offline approaches like this have known limitations:
    - They break down when phase execution time vs frequency scaling depends on runtime factors
- Get fancy: **online automatic** frequency optimization via Decider plugin
  - GEOPM plugin characterizes the application online and tunes each phase's frequency during an initial "learning" period
  - When finished learning, plugin uses best frequency for each phase and reaps energy benefits for the rest of execution
  - Online approaches like this have tradeoffs:
    - They're robust against cases when phase execution time vs frequency scaling depends on runtime factors and the best phase frequency can't be determined reliably offline
    - They break down when learning overhead is a non-trivial % of total execution time (not common case in long-running HPC apps)

# Experimental Setup: Measurement

- Common measurement methodology used for evaluating both techniques
- Programmers instrument application phase entry/exit with a profiling API provided with GEOPM
  - API is designed to be lightweight and easy for programmers to use
  - Nonetheless, work underway to automate phase instrumentation using OMPT
- GEOPM runtime performs energy and execution time measurements
  - Uses the phase instrumentation to track phase entry and exit timestamps
  - Periodically samples energy using processor counters (or other means)
  - Incorporates accounting logic to track total phase energy and execution time based on above
  - GEOPM report feature outputs this per-phase data to text file on each node
- All data points in charts shared today represent an average of at least 7 trials
- For each trial, results averaged across nodes since results vary for different nodes

# Experimental Setup: Workloads

- Studied workloads including:
  - **Proxy app**: model bulk-synchronous application with configurable balance of DGEMM and STREAM computation (from GEOPM tutorials and integration tests)
  - **FT**: a discrete 3-d FFT kernel (from NAS Parallel Benchmark suite)
  - **miniFE**: finite element code (from CORAL procurement benchmarks)
  - **Nekbone**: thermal hydraulics code (from CORAL procurement benchmarks)
- Applied the following conventions when configuring workloads:
  - Sized the problem to fit within available DRAM on the node
  - Tested several configs of MPI ranks and OpenMP threads per rank; used the config with lowest time-to-solution
  - Set application affinity masks to stay off of CPU0 to minimize OS jitter effects
  - Set GEOPM controller process affinity to CPU1 (application affinity set to stay off of this CPU)
  - Did not use hyperthreads

# Experimental Setup: Hardware

	JLSE Broadwell Xeon Cluster (Argonne)	Quartz Broadwell Xeon System (LLNL)
Workload configuration under study	proxy app: used 1 instance per node with no inter-node communication	FT, miniFE, Nekbone, Gadget: all used multi-node configurations with MPI communication
Processor and memory specs	4x 44-core nodes (dual-socket), Broadwell server processors, 128 GB of DRAM	Up to 8x 36-core nodes (dual-socket), Broadwell server processors, 128 GB of DRAM
Processor core frequency range	1.2 GHz to 2.2 GHz sticker, Turbo is enabled: frequency may exceed sticker	1.2 GHz to 2.1 GHz sticker, Turbo is enabled: frequency may exceed sticker
Network hardware	N/A	Intel OmniPath HFI and switches
Software environment	RHEL 7.4 Linux distro, Intel P-state driver running set to 'performance' with min=max=sticker, Intel compiler toolchain, Intel MPI implementation	RHEL 7.3 Linux distro with Intel P-state driver disabled, legacy governor set to 'performance', Intel compiler toolchain, MVAPICH2 MPI implementation

# Experimental Setup: 3 Investigations

## 1. Opportunity Analysis

- Use proxy app (parameterized model application) to determine envelope of energy-to-solution and time-to-solution impact we'll see over the landscape of BSP applications
- Measure energy-to-solution decrease and time-to-solution tradeoff **relative to running at sticker** on the JLSE cluster at Argonne
- Compare two different use-cases for the offline technique we developed:
  - 'Offline automatic **application** best-fit:' all phases run at common frequency (best-fit across all)
  - 'Offline automatic **per-phase** best fit:' each phase runs at the best frequency for it

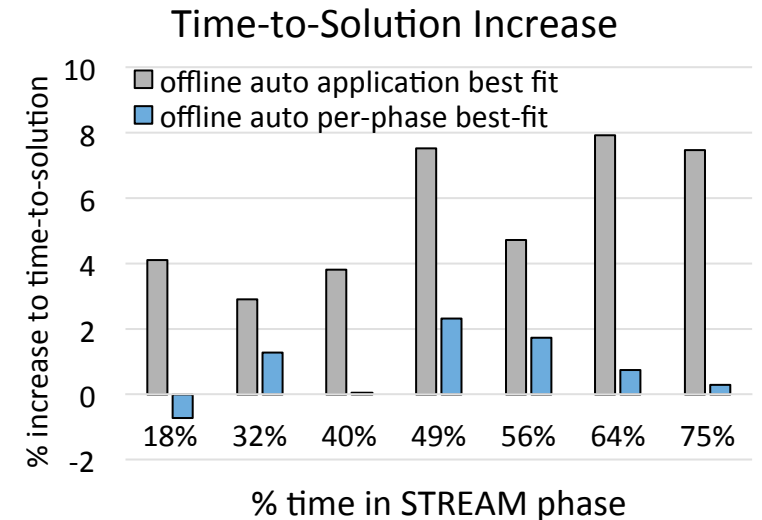
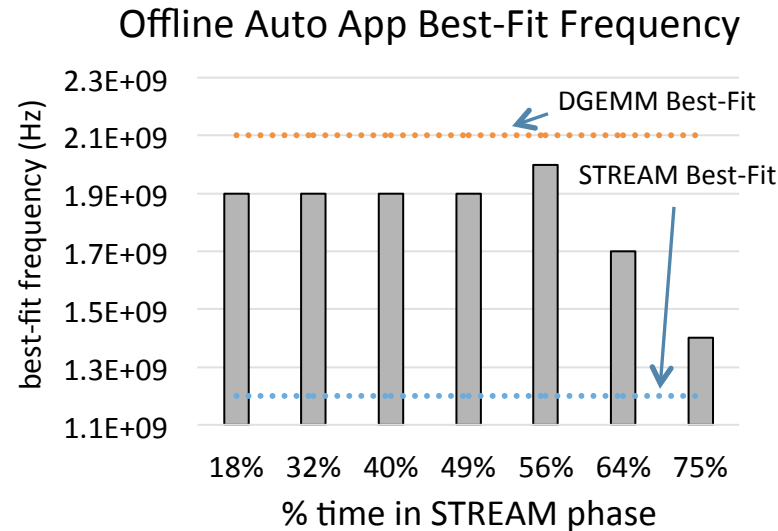
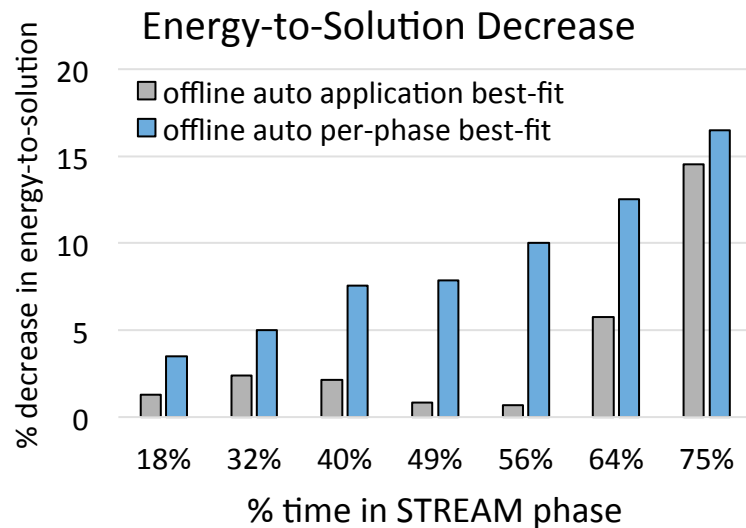
## 2. Benchmark **offline** energy optimization technique

- Target **FT**, **miniFE**, and **Nekbone** workloads
- Same as above but targets less synthetic workloads and performs experiments on LLNL Quartz system

## 3. Benchmark **online** energy optimization technique

- Target the proxy app and perform experiments on JLSE cluster at Argonne
- Compare the online and offline techniques we developed:
  - '**Offline** automatic per-phase best-fit:' scripts identify best frequency via offline characterization
  - '**Online** automatic per-phase best fit:' GEOPM plugin performs characterization/tuning online

# Results: Opportunity Analysis



Big energy savings are possible with frequency optimization in GEOPM vs running workloads at sticker: up to **16.5% energy savings** at **0.3% increase in time-to-solution**

With per-phase optimization, energy savings increase with increase in % time in memory-limited phase

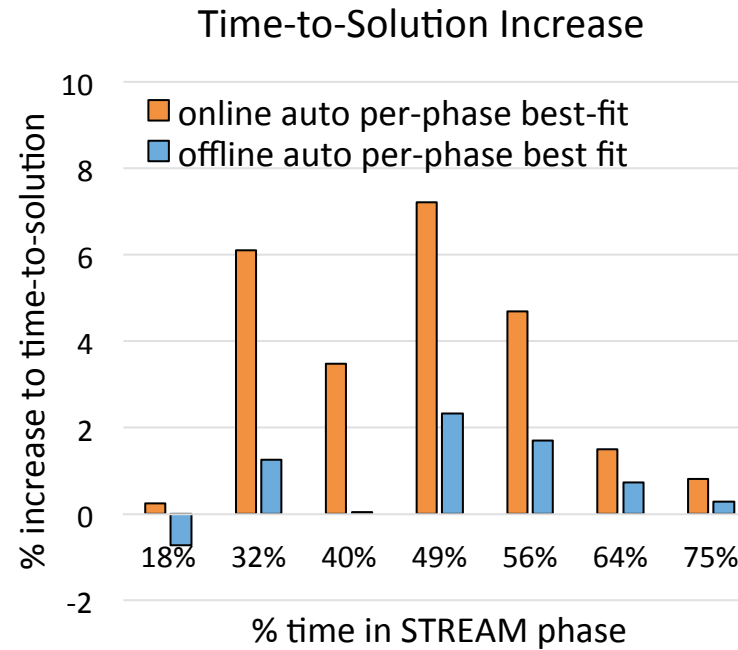
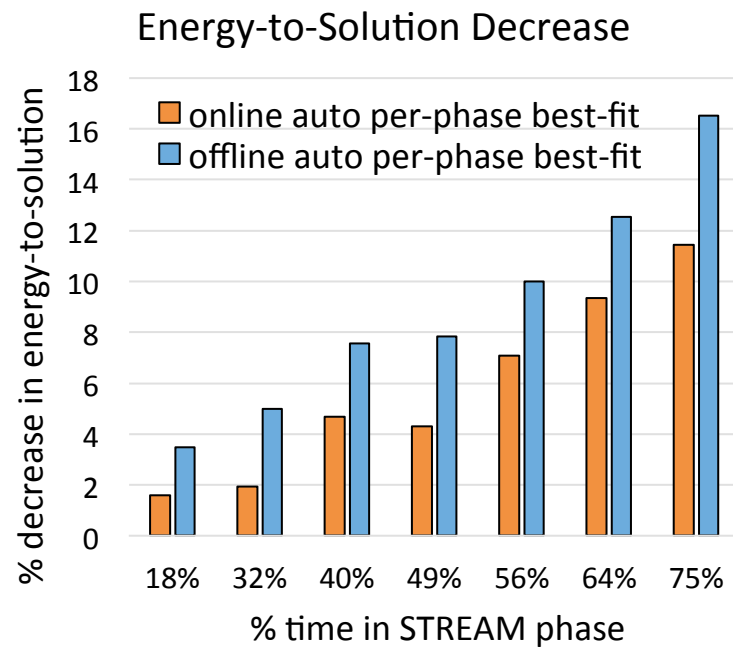
Per-phase optimization simultaneously offers better energy-to-solution AND time-to-solution versus optimizing frequency across the blended characteristics of all application phases

# Results: Offline App vs Per-Phase Best-Fit

Energy-to-Solution and Time-to-Solution Comparison on Quartz				
	Offline Automatic <b>Application</b> Best-Fit		Offline Automatic <b>Per-Phase</b> Best-Fit	
Workload	EtS Decrease vs Sticker	TtS Increase vs Sticker	EtS Decrease vs Sticker	TtS Increase vs Sticker
FT	9.5%	6.8%	15.8%	4.8%
miniFE	8.5%	5.8%	Collecting data now	Collecting data now
Nekbone	7.9%	2.4%	Collecting data now	Collecting data now

Results starting to confirm that GEOPM provides benefits for a number of workloads beyond our proxy app  
More data on the way, but data starting to suggest per-phase frequency optimization simultaneously offers better energy-to-solution AND time-to-solution vs optimizing frequency across blended characteristics of whole app

# Results: Online vs Offline Technique



## Explanation of EtS and TtS gaps:

- Runs were shorter than real apps -> noticeable “learning” overhead
- Reduced # samples in learning period to reduce overhead -> more noise-related control errors
- Observed latency between frequency change requests and enactment (10s of milliseconds) -> not running at desired frequency immediately, confusing algorithm

Remember, offline approach is brittle. The goal: same (or better) results via more robust online approach  
We think much of the EtS and TtS gap can be closed via addressing frequency latency & doing longer runs  
Fine-tuning needed, but already seeing promising decreases in energy-to-solution with online approach



# Next Steps

1. Extend online auto per-phase energy optimization plugin
  - Leverage GEOPM feature for automatic detection of OpenMP parallel region entry/exit to remove need for programmer to instrument phases with API
2. Expand evaluations to include more benchmarks (e.g. Gadget @ LRZ) and more architectures (e.g. Theta Knights Landing system @ Argonne)
3. Polish the energy optimization techniques demonstrated today and include them in GEOPM Beta
4. Work with LRZ, LLNL, and Argonne toward deployment on their production systems

# GEOPM Core Team Acknowledgements

## Lead Architect:

- Jonathan Eastep, Principal Engineer

## Hardware Team:

- Processor Firmware
  - Revathy Rajasree
- Hardware Architecture and Design
  - Fede Ardanaz
  - Fuat Keceli
  - Kelly Livingston
  - Lowren Lawson

## Software Team:

- GEOPM Development
  - Chris Cantalupo
  - Diana Guttman
  - Brad Geltz
  - Brandon Baker
- Research
  - Sid Jana
  - Asma Al-Rawi
  - Matthias Maiterth

# Backup Slides

# Per-Phase Best-Fit Frequency Table for FT

FT Phase	Best Frequency	% Total Exec Time at Sticker
Transpose_1	1.2 GHz	06.7%
Transpose_2	1.2 GHz	06.8%
MPI_All2All	1.2 GHz	35.2%
EVOLVE	1.2 GHz	09.9%
FFT_1_2	2.0 GHz	13.1%
FFT_1	2.0 GHz	13.1%
FFT_2	2.1 GHz	16.0%
INDEX_MAP	2.1 GHz	0.03%

Wide range in best-fit frequency across phases

Code inspection confirms the phases that tolerate 1.2 GHz are memory-limited or MPI phases, as expected

Phases tolerating 1.2 GHz use roughly the same % of total execution time as phases needing 2.0-2.1 GHz

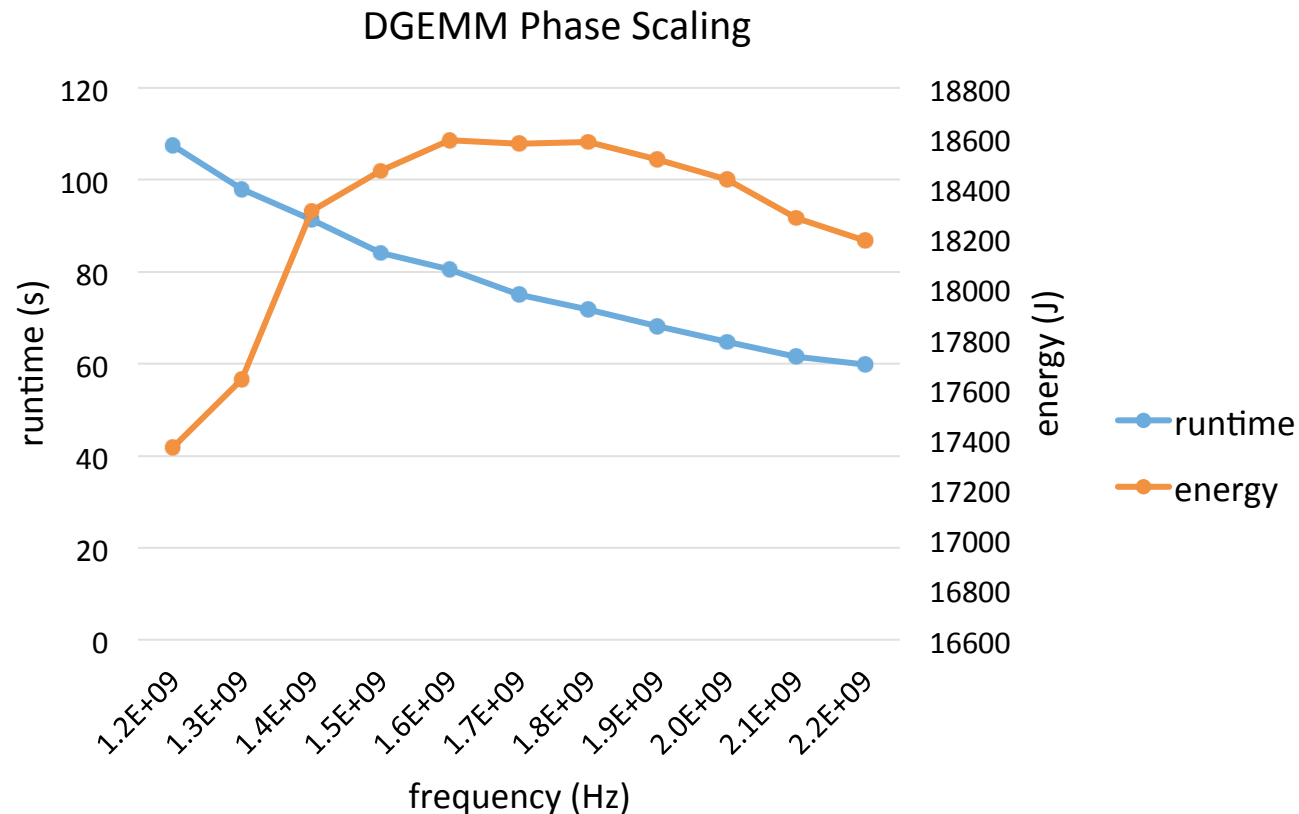
Running all phases at 2.0-2.1 GHz would waste energy in TRANSPOSE, MPI, and EVOLVE phases with little benefit to execution time

Running all phases at 1.2 GHz would harm execution time of FFT and INDEX\_MAP phases by much more than 10%

These facts illustrate why it is sub-optimal to apply the same frequency across all phases

# Energy May Not Monotonically Increase

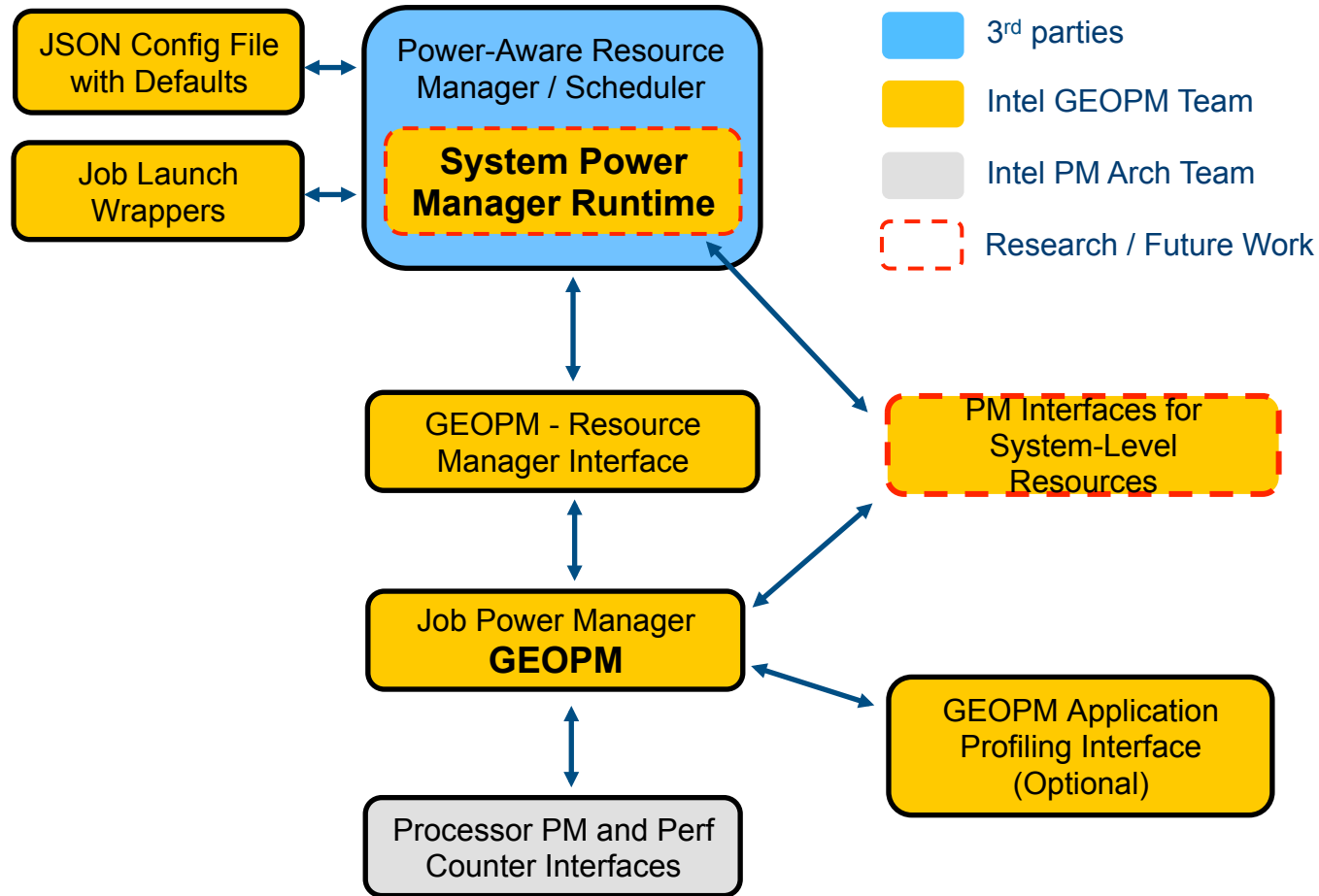
- Online algorithm can't simply walk down frequency until there's a 10% TtS increase (relative to sticker)
- Reason: reducing frequency below sticker *increases* (not decreases) energy for DGEMM phase; naïve algorithms would choose lower frequency for DGEMM than is optimal from an energy perspective!



# What Problems Does GEOPM Address?

- 1. At-scale load imbalance due to manufacturing variation in power-capped systems.** This problem is deemed one of the key Exascale-era power challenges. Developing GEOPM and techniques to address this problem over the past 6 years made me a Principal Engineer at Intel.
- 2. Gap in community energy management research tools.** There was previously no platform for energy management research that was open, scalable, robust, flexible, portable (truly open), and backed by serious engineering resources. Now the community is using GEOPM, porting to non-x86 architectures, integrating their optimization techniques into it, and integrating it with other software components.
- 3. Gap in industry server power management roadmaps and technical directions.** Power management was previously done *node-locally*. Techniques were *oblivious to application-level information* such as bottlenecks on remote nodes that could limit overall performance and were *unable to forecast* what computation was going to happen in the future and optimize power-performance policy accordingly. GEOPM adds a critical layer of global optimization across nodes, application and application phase awareness, and forecasting capabilities. See ISC'17 paper for demo of benefits (up to 30% speedup).

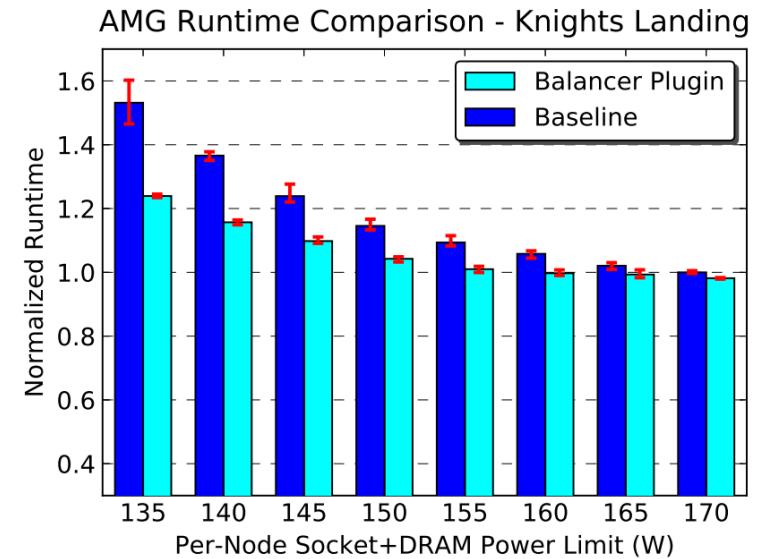
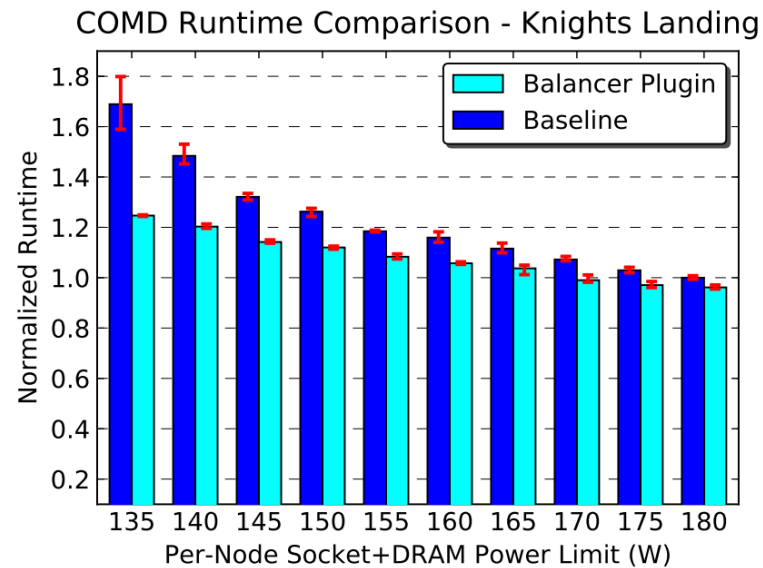
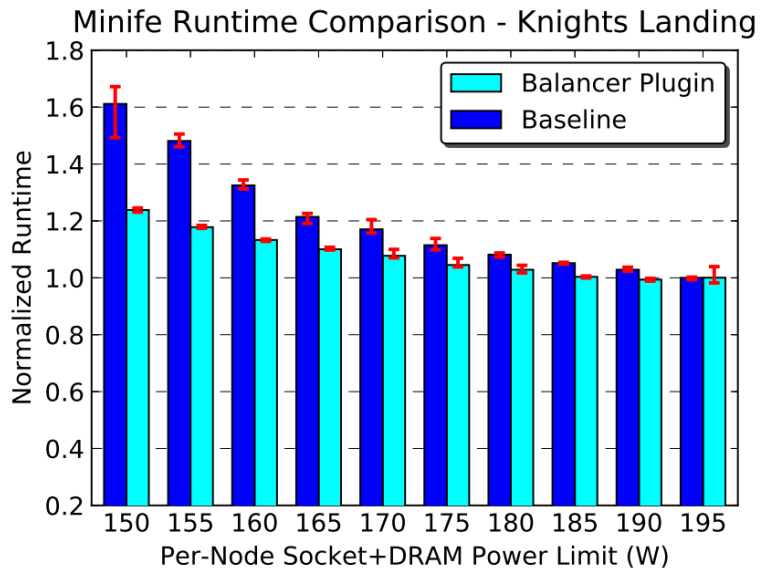
# GEOPM Interfaces and HPC Stack Integration



- GEOPM = job-level power manager
  - Coordinates optimization of hardware control knob settings across all compute nodes in job
  - Userspace software; access to control knobs facilitated via OS driver like msr-safe from LLNL
  - Supported control knobs: node power caps, processor frequency controls, more coming
- Run GEOPM via job launch wrappers
  - Includes wrappers for srun and aprun so far
  - Same syntax but with added flags to configure geopm; e.g.: power budget and plugin
  - Admin can provide defaults via JSON config file
- Integrates with RM and scheduler
  - Near-term: GEOPM can stand alone
  - Long-term: integrates with emerging System Power Manager (SPM) runtime component
  - GEOPM interface to system power manager
    - Feedback: GEOPM reports job power consumption and other job characteristics
    - Control: SPM dynamically reconfigures job power budget and/or GEOPM plugin

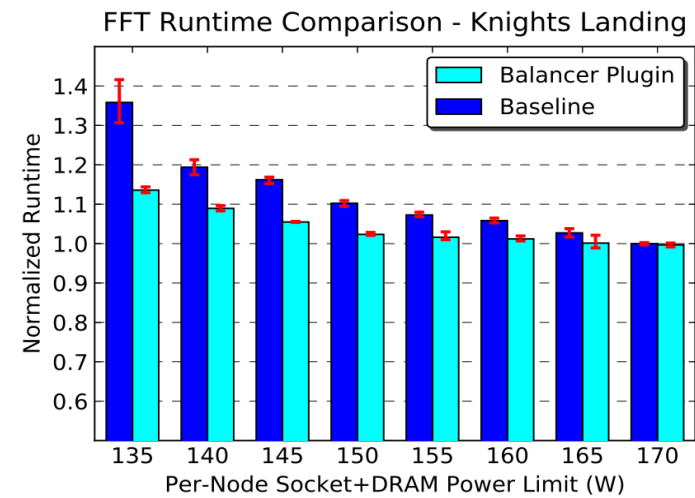
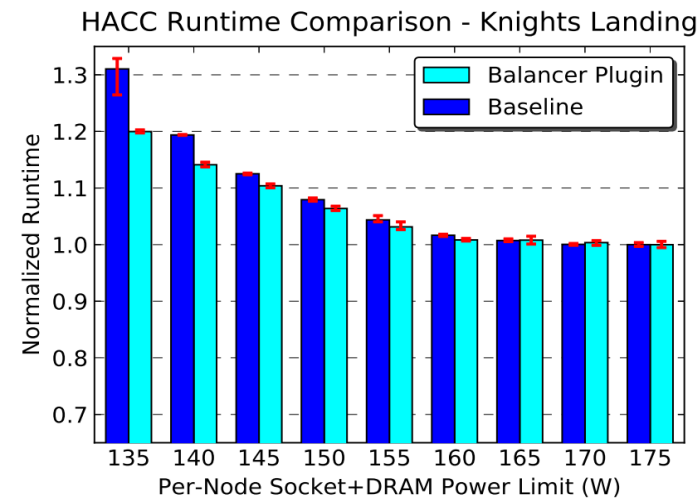
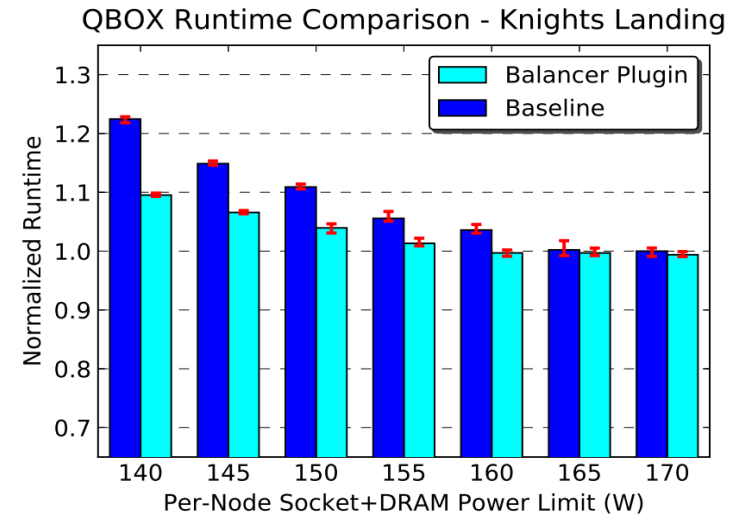
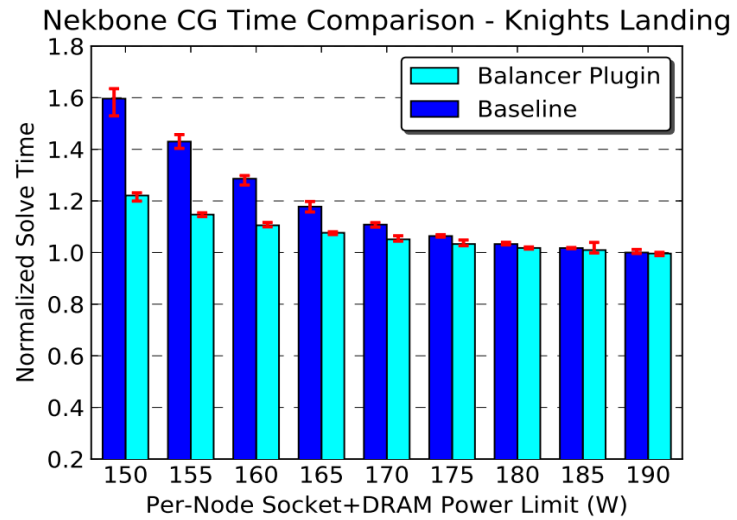
# Results: Inter-Node Power Balancing Use Case

- See GEOPM ISC'17 [paper](#) by Eastep et al. for details of experimental setup and further analysis
- Paper demonstrates power balancing plugin: it leverages annotation of application's outer synchronization loop to detect critical path nodes and then reallocates power among nodes in order to equalize their time to complete a loop iteration
- Compared overall time-to-solution when capping job power on 12-node KNL cluster with power balancer plug-in vs. static uniform power division (baseline); swept over a range of different job power caps
- Region of interest in job power caps: low-end of job power caps was selected to avoid inefficient clock throttling and the high-end of the job power caps equals the unconstrained power consumption of the workload
- Main result: **up to 30% improvement** in time-to-solution at low end of caps (miniFE, CoMD, AMG), with **up to 9-23% for the rest**. Improvement generally increases as power is more constrained





# Results: Four Additional Workloads

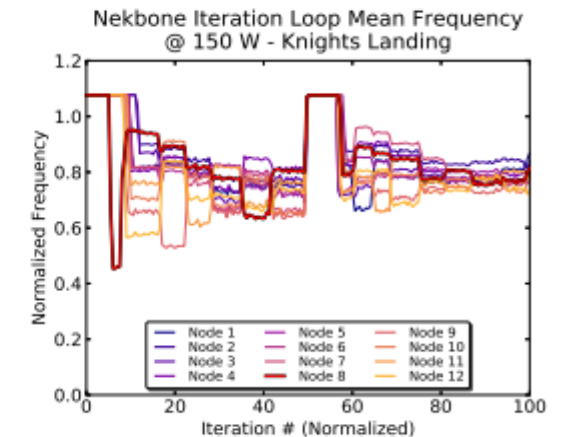
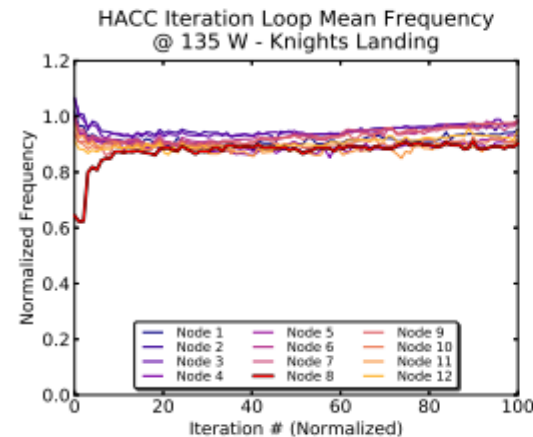
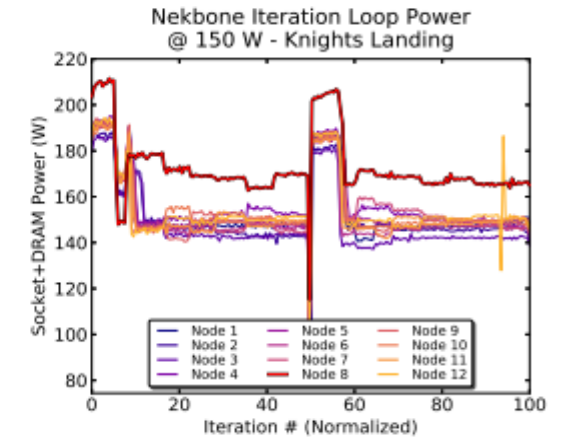
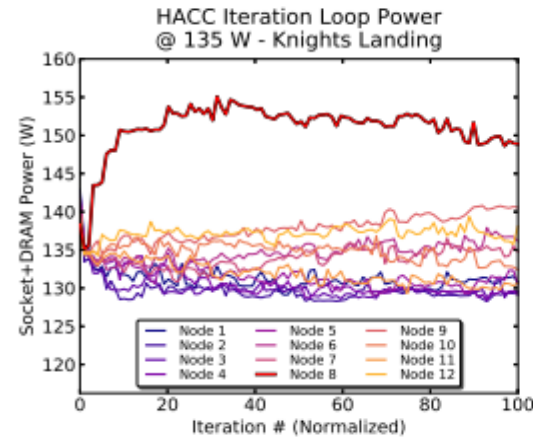
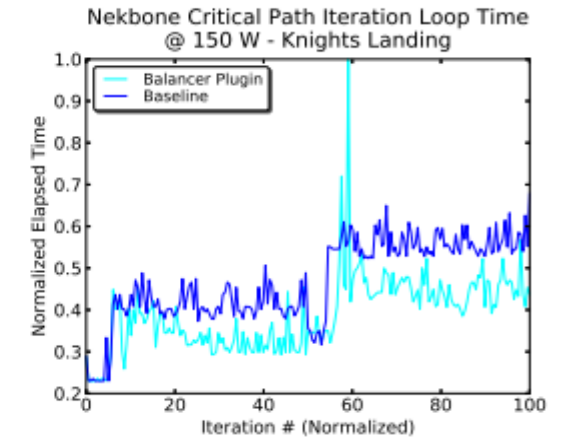
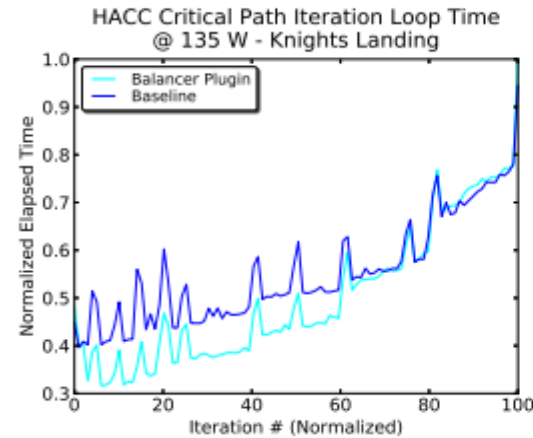


# GEOPM Speedup Analysis

(using included GEOPM Trace and Python Visualization Tools)

## Take-away points:

- Results demonstrate robustness of power balancing algorithm against time-varying amounts of work in the outer loop and sharp shifts in computational-intensity (top graphs)
- Node 8, with lowest power efficiency in our KNL cluster, is allocated more power (middle graphs)
- Power balancing algorithm improves critical path loop time by finding the power allocation that roughly equalizes the frequencies of all nodes (bottom graphs)



# Research on GEOPM/HW/FW Codesign

- GEOPM project is not just a software project. It also drives codesign of the features in Intel hardware for power-performance monitoring and control
- Goals are to significantly advance the state-of-the-art in HPC power management technology and to ensure GEOPM runs best on Intel
- Research areas:
  - Processor: improvements to granularity, reaction time, and interfaces for existing features
  - Processor: hooks for GEOPM to guide allocation of Turbo headroom among cores
  - Memory: hooks for GEOPM to hint to mem controller when it's best to enter low-power states
  - Network: hooks for GEOPM to estimate power, manage tradeoffs between power and bandwidth in HFI and switches, and hint to HFI when it's best to enter low-power states

# GEOPM New Business Opportunities

- GEOPM software package is open source, provides a rich feature set free of charge
- Intent is for Intel's future work on the software to be open source as well
- 3<sup>rd</sup> parties are able to make proprietary extensions of GEOPM (BSD 3-clause license)
  - Enables integrators like Dell/Cray/HPE to develop commercial for-profit plugins (i.e. add power management secret sauce to differentiate your systems vs the competition)
  - GEOPM team can help integrators with this in a consulting capacity
- Intel can explore developing custom processor firmware enhancements for customers
  - Enables processor power management firmware and GEOPM plugins to be co-optimized for individual customer needs
  - Enables management of hardware control knob settings which are not (yet) publically available
  - Providing GEOPM NRE funding in a system contract is a good way to establish such an engagement

Inquire with Jonathan Eastep for more information: [jonathan.m.eastep@intel.com](mailto:jonathan.m.eastep@intel.com)