# Global Survey of Energy and Power-aware
# Job Scheduling and Resource Management in Supercomputing Centers

**Siddhartha Jana**
siddhartha.jana@intel.com
Intel Corporation

**Gregory A. Koenig**
koenig@acm.org
Energy Efficient HPC Working Group

**Matthias Maiterth**
matthias.maiterth@intel.com
Intel Corporation

**Kevin T. Pedretti**
ktpedre@sandia.gov
Sandia National Laboratory

**Andrea Borghesi**
andrea.borghesi3@unibo.it
University of Bologna

**Andrea Bartolini**
andrea.bartolini@iis.ee.ethz.ch
IIS, ETH Zurich

**Bilel Hadri**
bilel.hadri@kaust.edu.sa
KAUST Supercomputing Lab

**Natalie J. Bates**
natalie.jean.bates@gmail.com
Energy Efficient HPC Working Group

## System Design Challenges:
- Building systems for HPC under a Power Budget
- Peak power demands for future Exascale systems > 20MW
- Instantaneous power fluctuations: 8MW
- Microarchitecture improvements and high degree of parallelization not sufficient

## Eight Survey Questions for the sites:
1. Motivation behind investing in Energy and Power Aware Job Scheduling and Runtime Management (EPA-JSRM)
2. Target infrastructure (e.g. site-wide power budget, cooling capacity, etc.)
3. Workload characteristics
4. Adopted design for EPA-JSRM
5. Implementation details for EPA-JSRM
6. Application/task level and topology-aware solutions
7. Results and challenges
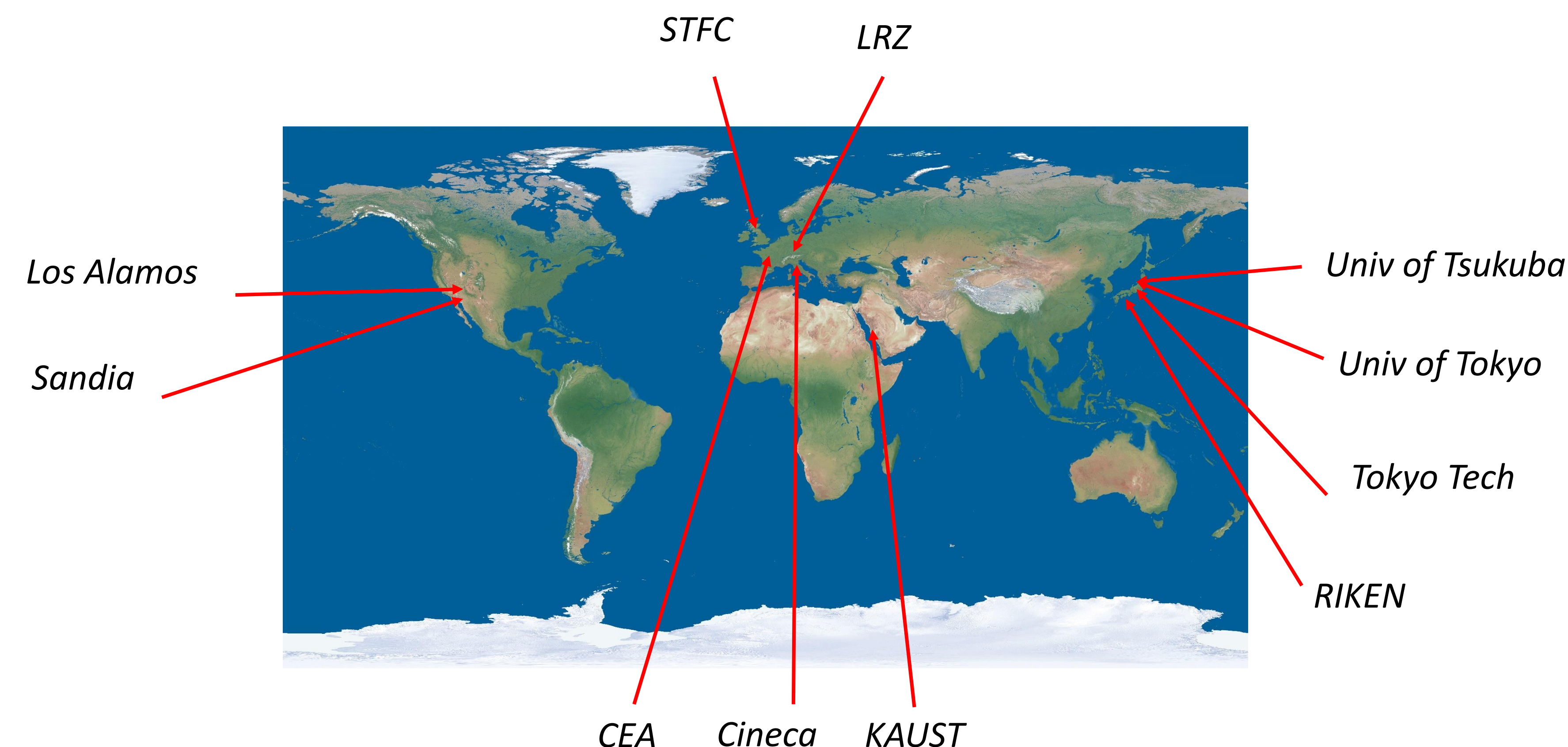8. Next steps including system procurement

## Participating Sites:
- CEA ( Alternative Energies and Atomic Energy Commission), France
- Cineca, Italy
- KAUST (King Abdullah University of Science and Technology), Saudi Arabia
- LRZ (Leibniz Supercomputing Centre), Germany
- Riken, Japan
- STFC (Science and Technology Facilities Council), United Kingdom
- Tokyo Institute of Technology, Japan
- University of Tokyo and University of Tsukuba (JCAHPC), Japan
- Los Alamos and Sandia National Laboratories (Trinity), United States

## EPA-JSRM solutions depicted on the right have already been adopted - at least, in parts by these sites.

### Next phase of JSRM roadmap:
- Continue working on more stable designs of system-wide frameworks (e.g. job schedulers) that allocate resources in a power-aware manner
- Invest in robust energy/power predictors that rely on statistical modeling
- Leverage power-capping mechanisms exposed by vendors



Earth Map Credit: NASA's Earth Observatory

## Telemetry Monitoring solutions adopted:
- Sensors for monitoring energy and power
  - Both in-band as well as out-of-band
  - Direct real-time measurements
- Thermal-based sensors coupled with prediction models
- Model to indirectly derive power-based metrics
- Ongoing design of high-level APIs for end-users and resource managers
  - Energy and power monitoring
  - Feedback mechanisms
- Implementation of statistical approaches for prediction
  - Model based on job demand, size, length, etc.
  - Helps assign power budget to specific users

## JSRM solutions adopted:
- Dynamic shutdown of jobs in response to limited power budget (Reactive approach)
  - Job selection based on job size, job length, etc. to shut down
- Automated reduction of node availability by the resource manager (Proactive approach)
  - Reduces the theoretical maximum power that can be consumed
  - Drop in system utilization
- Use of power-capping mechanisms supported by CPU and system vendors
  - Attempts to keep total power consumed below a specific limit
  - Power cap applied over a specific time-window (in order of minutes)
- Use system interface to trigger specific p-states (operating frequencies) supported by the platform
  - Design of portable APIs
- Design of system-wide frameworks (like job schedulers) that use static prediction models
- "Tagging" applications based on their power usage characteristics (feedback-driven approach)
  - Mapping of "tags" to performance metrics
  - Storage of historical records attained over past job runs
  - Use of tag-values for future budget assignment

### System Characteristics

| Organization | Site Power Budget | Site Cooling capacity | Major HPC System | System Power Draw |
|---|---|---|---|---|
| RIKEN | Up to 25 MW | Up to 40 MW | K computer (83,944 nodes) | Up to 15 MW |
| Tokyo Tech | Up to 5 MW | Up to 5 MW | TSUBAME2.5 (1400 nodes) | Up to 5 MW |
| CEA | Up to 10 MW | Up to 10 MW | Anticipated 25PF System in 2017 | Up to 5 MW |
| KAUST | Up to 5 MW | Up to 5 MW | Shaheen 2 (6174 nodes) | Up to 5 MW |
| LRZ | Up to 10 MW | Up to 15 MW | SuperMUC Phase 1 / 2 | Up to 5 MW |
| STFC | Up to 5 MW | Up to 5 MW | 846 x dual SKX (128GB), 840 x KNL (96GB), 24x dual SKX (1TB) | ** |
| LANL + SNL (Trinity) | Up to 20 MW | Up to 30 MW | Trinity (9436 HSW nodes + 9984 KNL nodes) | Up to 10 MW |
| Cineca | Up to 10 MW | ** | Marconi (7500 nodes) | Up to 4 MW |
| JCAHPC | Up to 10 MW | ** | Oakforest PACS (8208 nodes) | Up to 5 MW |

** Information unavailable as of Oct 2017