



TECHNISCHE
UNIVERSITÄT
DRESDEN



HAEC DFG

READEX
Runtime Exploitation of Application Dynamism
for Energy-efficient eXascale computing

DRESDEN
concept



Powernightmares: The Challenge of Efficiently Using Sleep States on Multi-Core Systems

Thomas Ilsche, Marcus Hähnel, Robert Schöne, Mario Bielert,
and Daniel Hackenberg

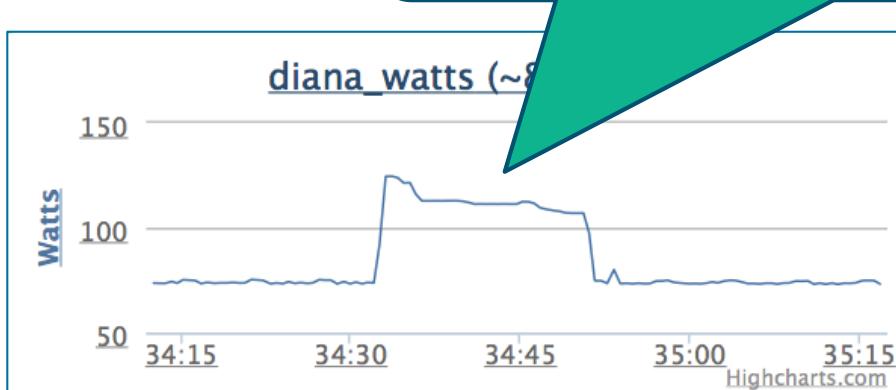
Technische Universität Dresden

29.08.17

5th Workshop on Runtime and Operating Systems for the Many-core Era

Observation

- Systems with continuous energy measurement
- Tuned for low idle power consumption
- Prolonged high power consumption
consumption → “*Powernightmare*”



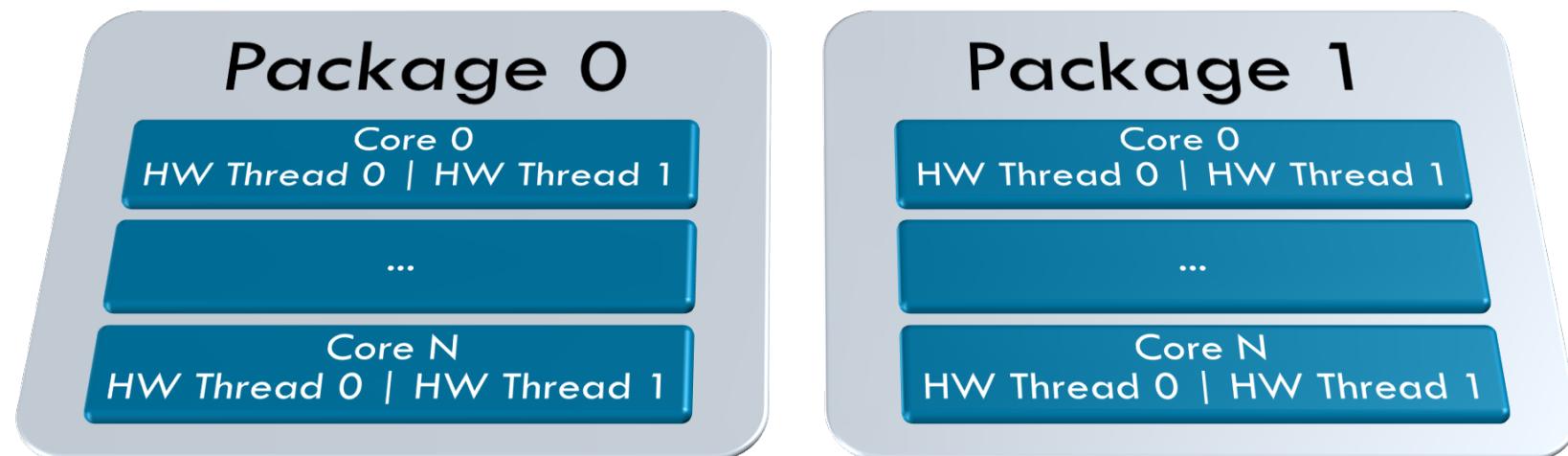
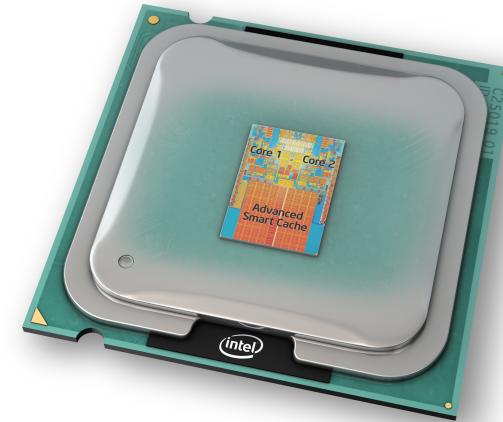
Background – Processor



3

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Each processor is a package
- A package comprises multiple cores
- Each core has two hardware threads
- A hardware thread is called CPU



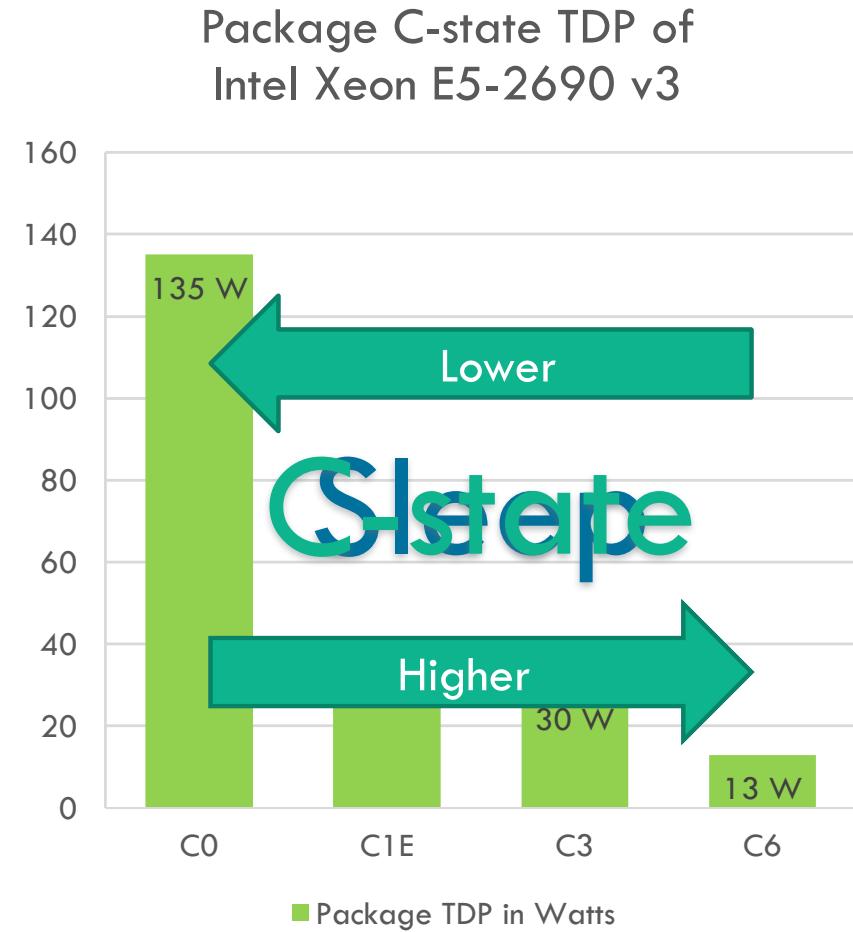
Background – C-states



4

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Idle power conservation
- Increasing latency
- Controllable per CPU,
but applied per core
- Package C-state
determined by lowest
core C-state
- Effective use is essential
for low idle power



Background – Linux idle governor



5

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Selects C-state for CPU
- ladder_governor gradually changes C-state
- menu_governor is based on a heuristic
- Heuristic used to predict idle time
 - ▣ Next timer event with correction factor
 - ▣ Repeatable interval detector (up to 8 data points)
 - ▣ Latency requirement

Investigation – lo2s



6

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

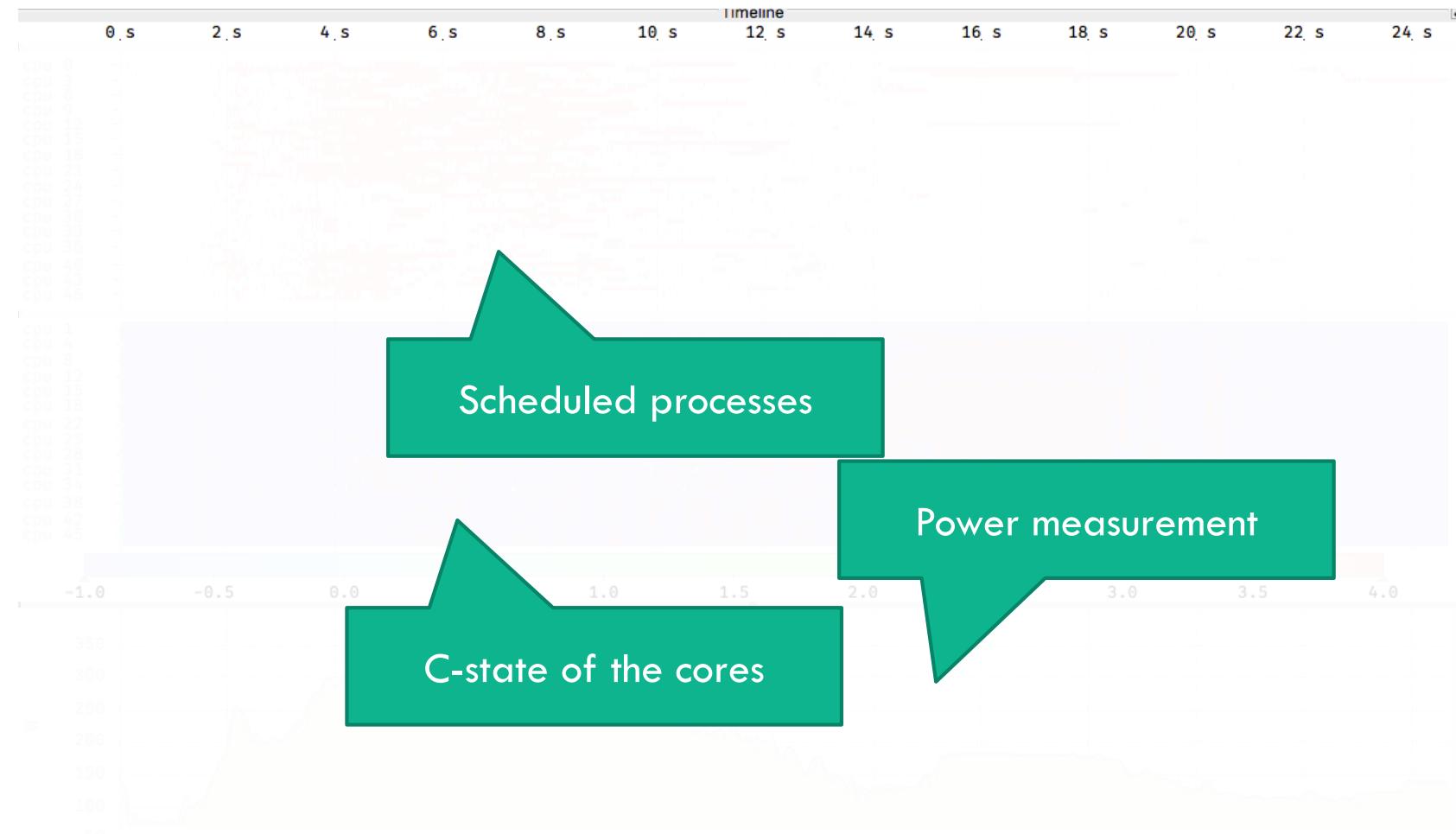
- Uses Linux' perf infrastructure
- Create a trace combining
 - ▣ Active processes using the trace point sched_switch
 - ▣ Selected C-state using the cpu_idle trace point
 - ▣ External power measurements
 - ▣ C-state residency using x86_adapt
- Available at <https://github.com/tud-zih-energy/lo2s>

Investigation – lo2s



7

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing



Vampir showing a lo2s trace of a parallel build using make
5th Workshop on Runtime and Operating Systems for the Many-core Era

29.08.17

Investigation – Powernightmare



8

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Up to 3 wakeups needed for correction after a misprediction by the heuristic



Zoomed Full elaboration Power Right nightmare Sched Credates system states socket power

Triggering the issue

9

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

□ Code to reliably trigger a Powernightmare

```
int main() {
#pragma omp parallel
{
    #pragma omp barrier
    while (1) {
        for (int i = 0; i < 8; i++) {
            #pragma omp barrier
            usleep (10);
        }
        sleep (10);
    }
}
```

Approaching the problem



10

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Changing task behavior
- Improving the idle time prediction
- Biasing the prediction error
- C-state selection by hardware
- Mitigating the impact

Impact mitigation approach



11

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Set a wakeup timer if huge difference between next known timer and predicted idle time

Prediction correct

- Wakeup event in predicted time interval
- Cancel timer

Prediction incorrect

- Timer triggers wakeup
- Ignore recent residency
- Enter high C-state
- ✓ Misprediction corrected

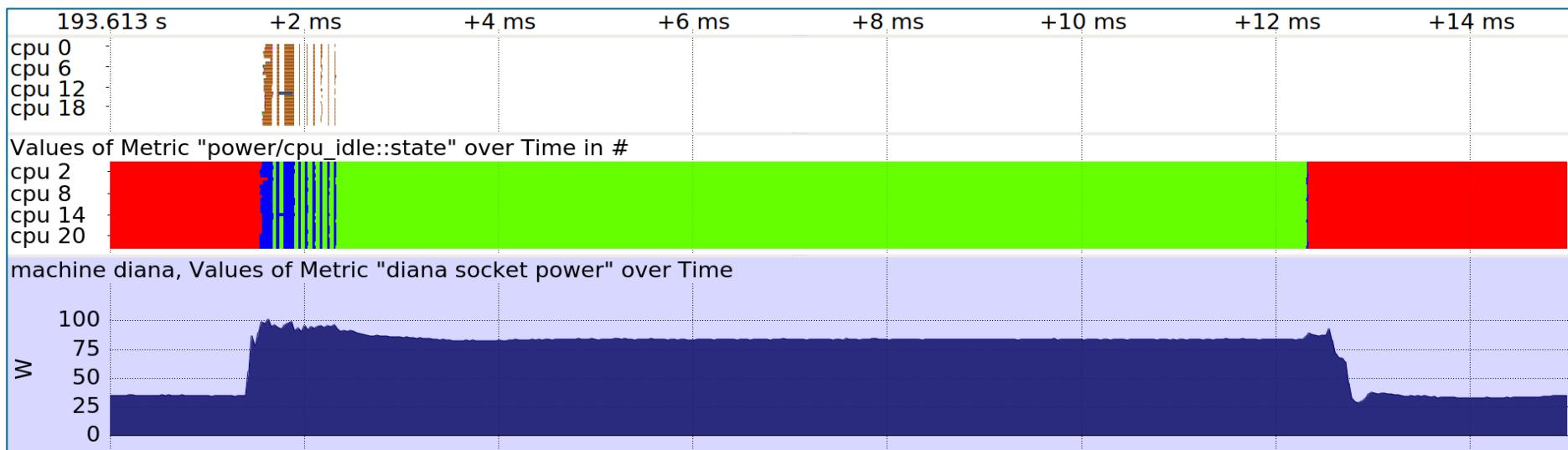
Powernightmare with timer



12

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- ❑ Fallback timer corrects wrong C-state selection
- ❑ Only 10 ms of shallow sleep



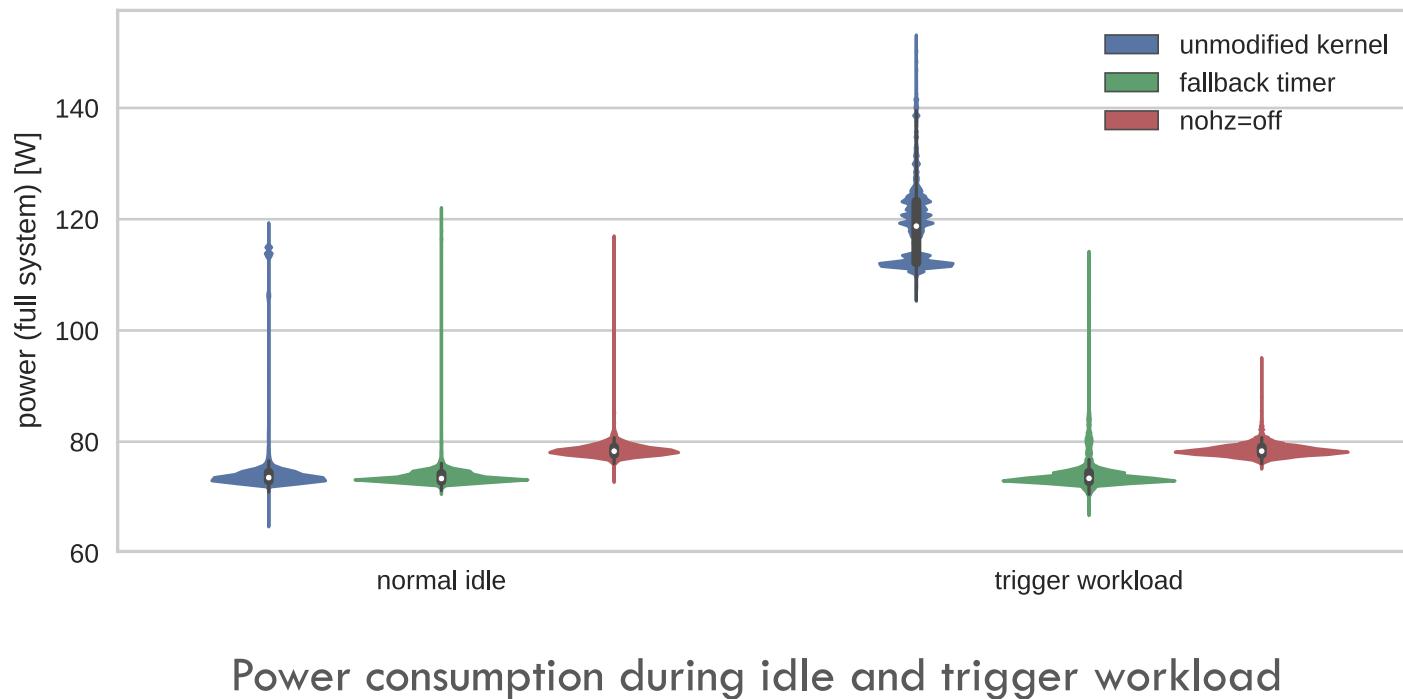
Reduced impact of Powernightmare with active fallback timer

Verification

13

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Measurements taken over 20 minutes
- Trigger workload every 10 seconds

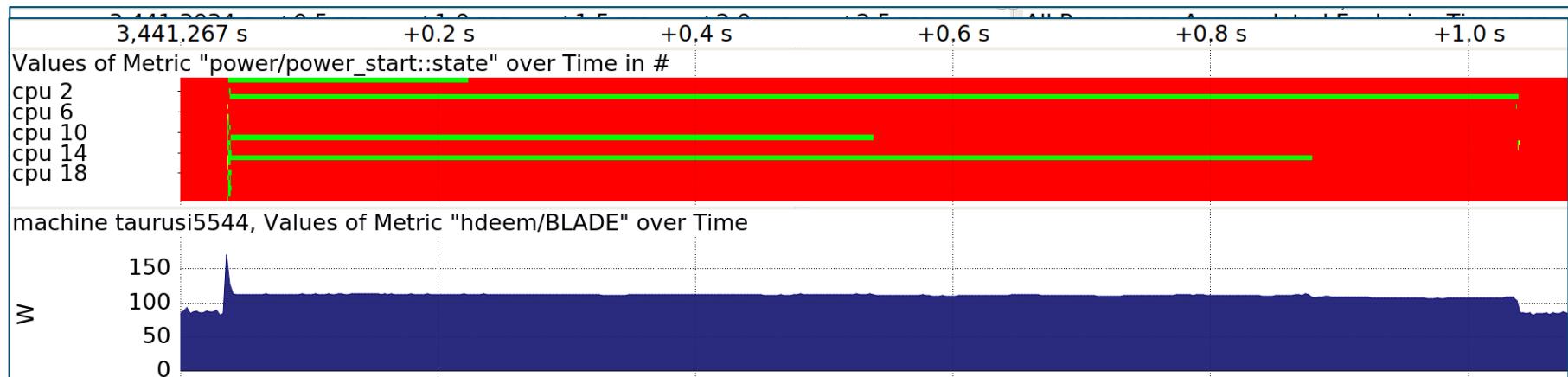


Production servers?

14

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Found on node of production HPC system “taurus”
- Lustre related pattern every 25 seconds
- Triggers one second Powernightmare



A short discussion point about Lustre related skew in task C1

Summary

15

Collaborative Research Center 912: HAEC – Highly Adaptive Energy-Efficient Computing

- Analyzed pattern of inefficient use of sleep states
- Developed a methodology and tools to observe
- Investigation shows misprediction in idle governor
- Proposed solution to mitigate effect
- Discussion with Linux community initiated
- Increasing probability with rising number of cores
- Effect not limited to HPC Systems

Any questions?

