JÜLICH
FORSCHUNGSZENTRUM

# Mapping Fine-Grained Power Measurements to HPC Application Runtime Characteristics on IBM POWER7

September 2nd, 2013 | Michael Knobloch

# The Exascale Innovation Center

## EIC - Exascale Innovation Center

- Project partners: IBM Germany R&D and JSC
- Goal: Co-Design for next-gen of Supercomputers
- One work-package on power and energy-efficiency
  - Investigation of power consumption on Blue Gene (EnA-HPC'11)
  - Fine-grained power measurements on POWER7 (this work)
  - Energy modelling on POWER7 (to be published)

# Test system – IBM Power 720

- 4-Core 3.0 GHz processor (Pseries, 8202-E4B)
  - 96 GFLOPS peak
  - 4 SMT threads per core
  - 64 kB L1 cache per core
  - 256 kB L2 cache per core
  - 16 MB L3 cache (shared)

- 16 GB memory, 2x 300 GB 10K RPM SAS disk
- **TPMD (Thermal Power Management Device)**
- **External power distribution unit (Raritan DPXS12A-16)**
  - 3 s measurement interval, 1 W resolution

# Amester

## Amester

IBM Automated Measurement of Systems for Temperature and Energy Reporting software.
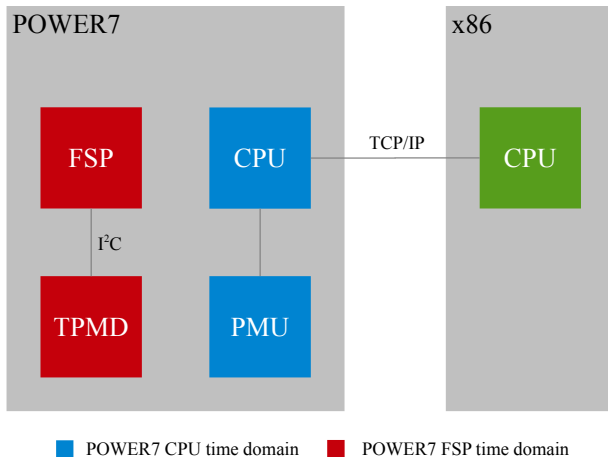
- Tool for monitoring and controlling power consumption of (IBM) servers – x86 and POWER
- Developed by Charles Lefurgy, IBM Research, Autisn, TX
- Histograms, traces for any sensor
- Scripting
  - Tcl command line
  - Send any IPMI command to measured system (ipmicmd)
  - On-line (50 ms interval) and off-line (buffered, 16 MB, 1 ms sampling) modes
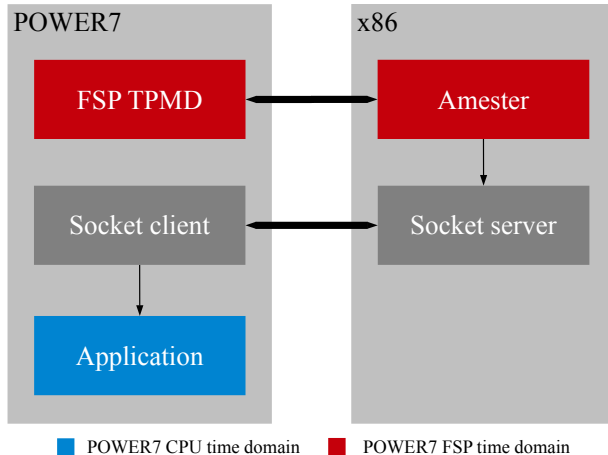
# Amester (condt.)

- Sensor data collection
  - Whole system power data collection (CPU, Memory, Fans, IO, Storage)
  - CPU temperature, processor speed, system utilization, instructions per second, memory bandwidth

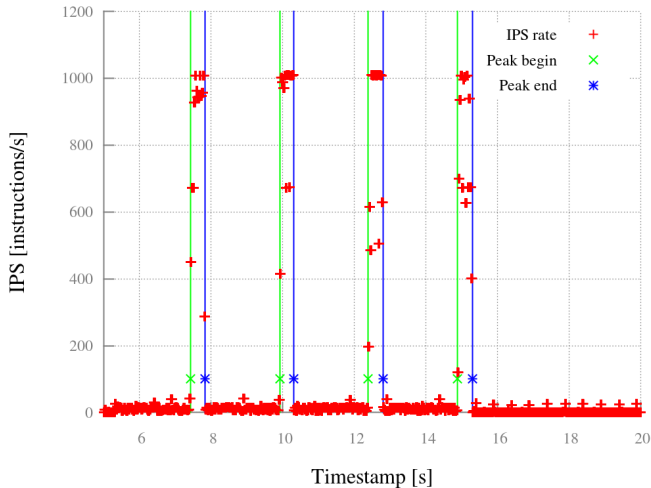| Sensor name | Units | Time scale | Description |
|---|---|---|---|
| PWR1MS | W | Instantaneous | Node power consumption |
| PWR1MSP0 | W | Instantaneous | Processor power consumption |
| PWR1MSMEM0 | W | Instantaneous | Memory power consumption |
| PWR32MS | W | avg. over last 32 ms | Node power consumption |
| PWR32MSP0 | W | avg. over last 32 ms | Processor power consumption |
| PWR32MSMEM0 | W | avg. over last 32 ms | Memory power consumption |
| IPS32MS | Mips | Every 32 ms | Instructions per second rate |

# Hardware Setup



POWER7 CPU time domain    POWER7 FSP time domain

# Software Setup



POWER7

x86

FSP TPMD

Amester

Socket client

Socket server

Application

■ POWER7 CPU time domain   ■ POWER7 FSP time domain
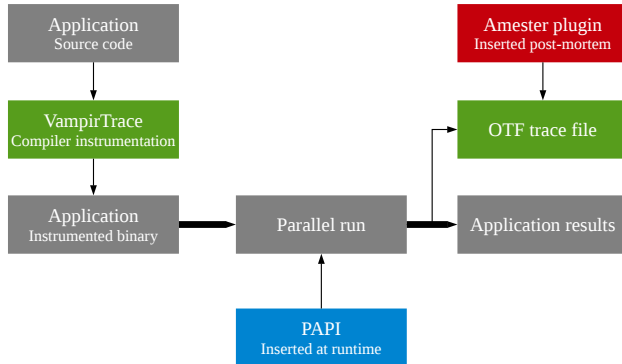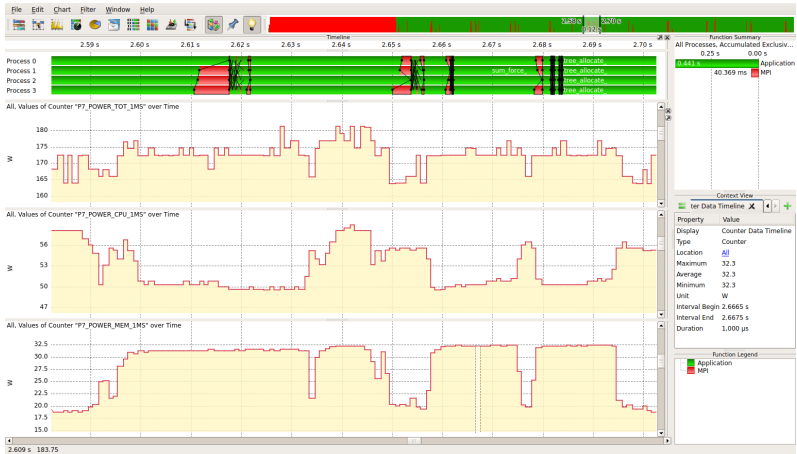
# Time-stamp synchronization
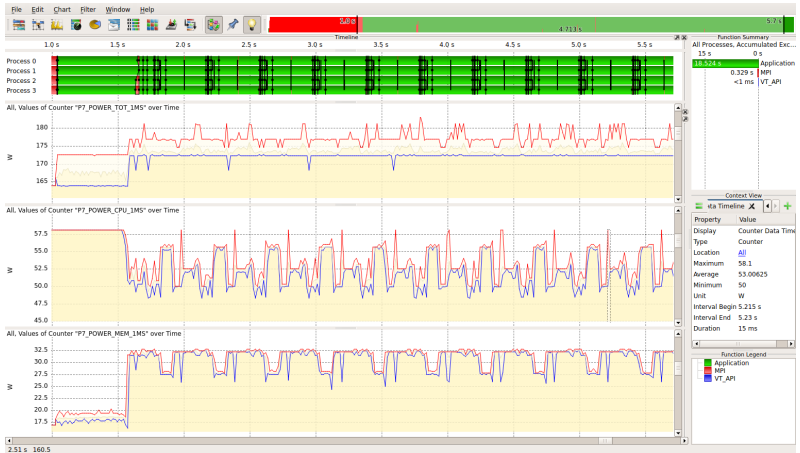
# VampirTrace Workflow
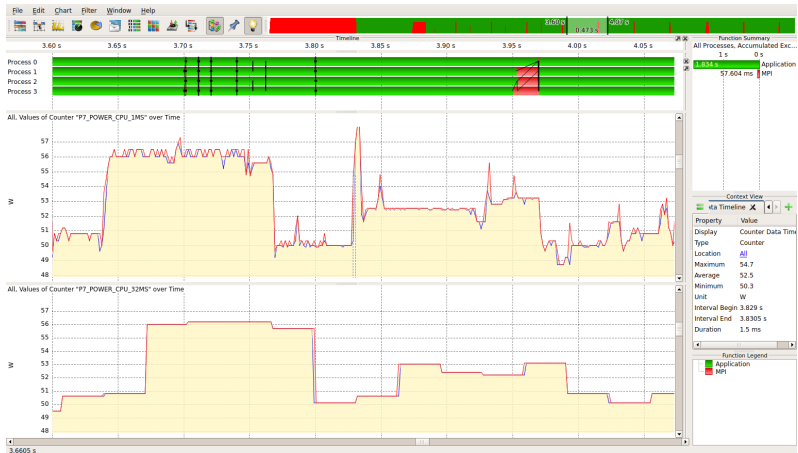
# PEPC: Full run (without initialization)
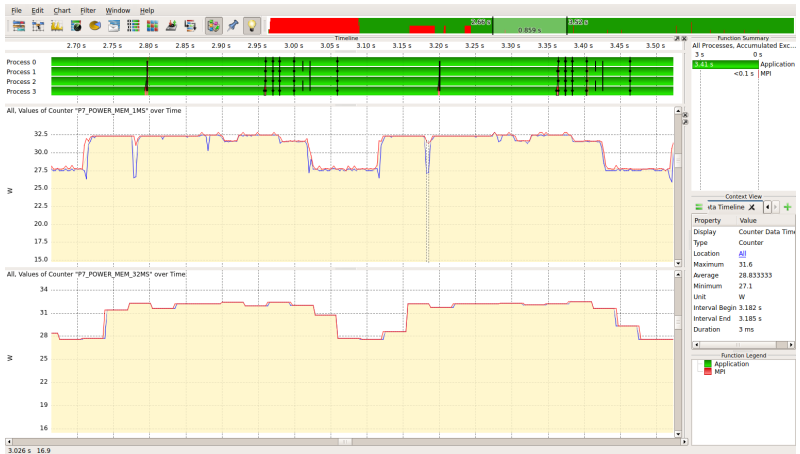
# PEPC: 1 Iteration

# MP2C: Full run (without initialization)
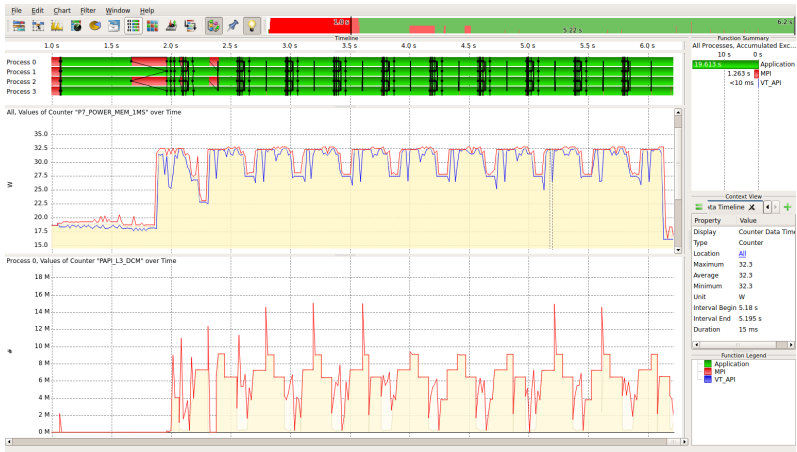
# MP2C: CPU Counter Resolution Comparison

# MP2C: MEM Counter Resolution Comparison

# MP2C: CPU Power and IPS

# MP2C: MEM Power and L3 Data Cache Misses

## Outlook

### Conclusions

- Fine-grained power measurements help to better understand application power consumption
- Amester requires complicated setup
- Mapping to other metrics can be difficult due to timing issues
- **Correlation does not imply causation**

### Future Work

- Energy modelling (to be published)
- Integration in Score-P and Scalasca