

ECE 571 – Advanced Microprocessor-Based Design Lecture 6

Vince Weaver

`http://web.eece.maine.edu/~vweaver`

`vincent.weaver@maine.edu`

15 September 2022

Announcements

- HW#1 grades out
- HW#3 will be posted, RAPL
- Note, the earthquake benchmark takes a while to run (up to a few minutes). Don't give up on it.
- Some notes on Intel uops, `uops.info`



Paper Discussion

Producing Wrong Data Without Doing Anything Obviously Wrong! by Mytkowicz, Diwan, Hauswirth and Sweeney, ASPLOS'09.



Some Background on Academic Publishing

- In grad school you might read and write a lot of papers like this
- Especially if you are a PhD student
- Rundown of how “peer-review” works

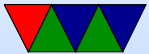


Paper Discussion Topics

- What is measurement bias?
- How did they set up their experiments?
- Did they find measurement bias in common computer engineering experimental methods?
- What were some of the sources of bias they found?
- What can be done to avoid it?
- In the years since this was published, was there any change in behavior by those publishing papers in Computer Engineering?



Power and Energy



Definitions and Units

People often say Power when they mean Energy

- Energy – Joules, kWh (3.6MJ), Therm (105.5MJ), 1 Ton TNT (4.2GJ), eV (1.6×10^{-19} J), BTU (1055 J), horsepower-hour (2.68 MJ), calorie (4.184 J)
- Power – Energy/Time – Watts (1 J/s), Horsepower (746W), Ton of Refrigeration (12,000 Btu/h)
- Volt-Amps (for A/C) – same units as Watts, but not same thing
- Charge – mAh (batteries) – need V to convert to Energy



Power and Energy in a Computer System

Power Consumption Breakdown on a Modern Laptop, A. Mahersi and V. Vardhan, PACS'04.

- Old, but hard to find thorough breakdowns like this
- Thinkpad Laptop, 1.3GHz Pentium M, 256M, 14" disp
- Oscilloscope, voltage probe and clamp-on current probe
- Measured V and Current. $P=IIR$. $V=IR$ $P=IV$,
subtractive for things w/o wires
- Total System Power 14-30W
- Old: no LED backlight, no SDD, etc.



Modern results are from CUGR/REU student research.

| | Laptop (2004) | Modern | Server? |
|--------------|---------------|--------|---------|
| Hard Drive | 0.5-2W | 5W | |
| LCD | 1W | | |
| Backlight | 1-4W | | |
| CPU | 2-15W | 60+W | |
| GPU | 1-5W | 50+W | |
| Memory | 0.5-1.5W | 1-5W | |
| Power Supply | 0.65W | | |
| Wireless | 0.1 - 3W | | |
| CD-ROM | 3-5W | | |
| USB | (max 2.5W) | | |
| USB keyboard | | 0.04W | |
| USB mouse | | 0.03W | |
| USB flash | | 0.5W | |
| USB wifi | | 0.5W | |



Skip Ahead

- We skipped ahead to cover RAPL first so you can start on the homework



Metrics to Optimize

- Power
- Energy
- MIPS/W, FLOPS/W (don't handle quadratic V well)
- $Energy * Delay$
- $Energy * Delay^2$



Power Optimization

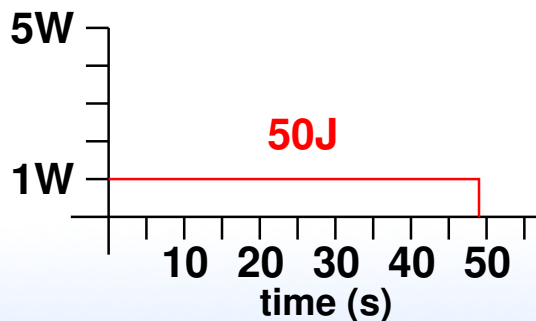
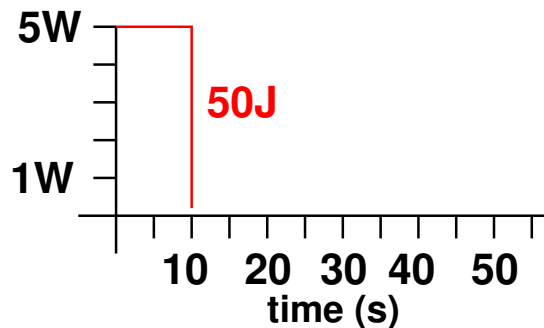
- Does not take into account time. Lowering power does no good if it increases runtime.



Energy Optimization

- Lowering energy can affect time too, as parts can run slower at lower voltages

Which is better?



Energy Delay – Watt/t*t

- Horowitz, Indermaur, Gonzalez (Low Power Electronics, 1994)
- Need to account for delay, so that lowering Energy does not made delay (time) worse
- Voltage Scaling – in general scaling low makes transistors slower
- Transistor Sizing – reduces Capacitance, also makes transistors slower

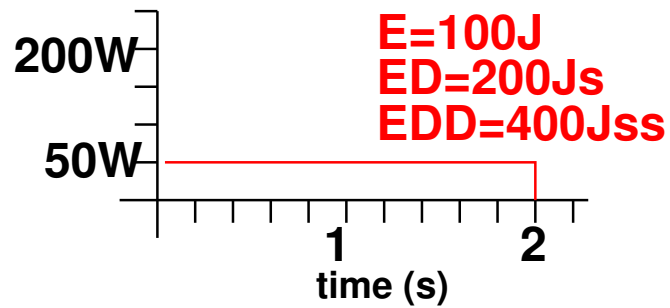
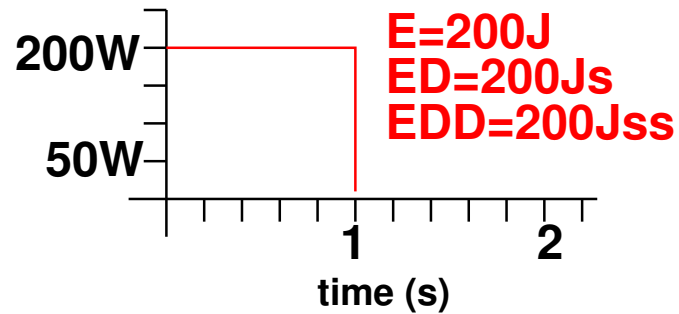


- Technology Scaling – reduces V and power.
- Transition Reduction – better logic design, have fewer transitions
Get rid of clocks? Asynchronous? Clock-gating?



ED Optimization

Which is better?



Energy Delay Squared– $E*t*t$

- Martin, Nyström, Péntzes – Power Aware Computing, 2002
- Independent of Voltage in CMOS
- $E*t$ can be misleading
 $E_a=2E_b$, $t_a=t_b/2$
Reduce voltage by half, $E_a=E_a/4$, $t_a=2t_a$, $E_a=E_b/2$,
 $t_a=t_b$
- Can have arbitrary large number of delay terms in Energy product, squared seems to be good enough



Energy Delay / Energy Delay Squared

Lower is better.

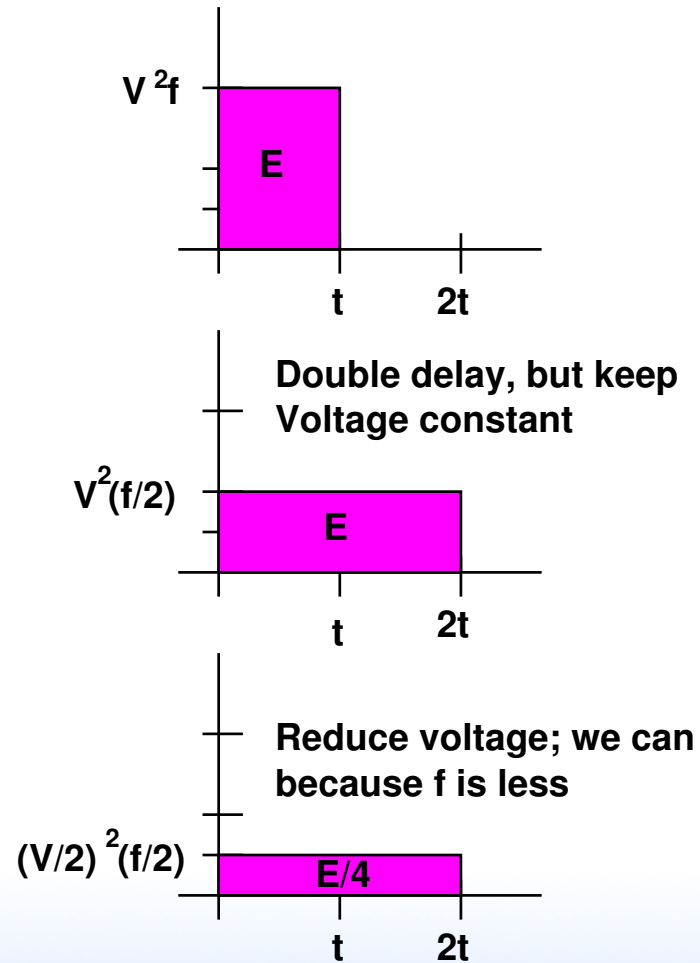
| Energy | Delay | ED | ED^2 |
|--------|-------|--------|----------|
| 5J | 2s | $10Js$ | $20Js^2$ |
| 5J | 3s | $15Js$ | $45Js^2$ |

Same ED , Different ED^2

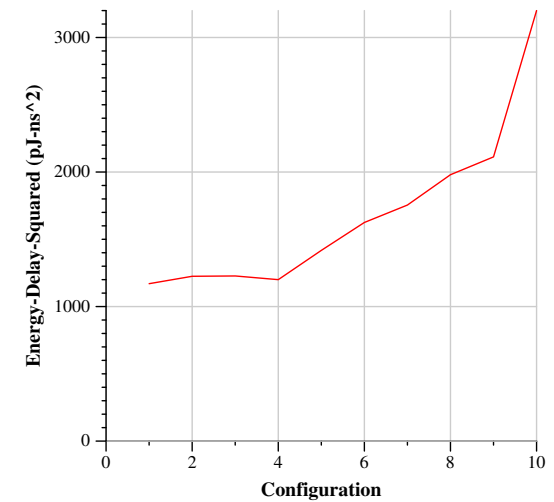
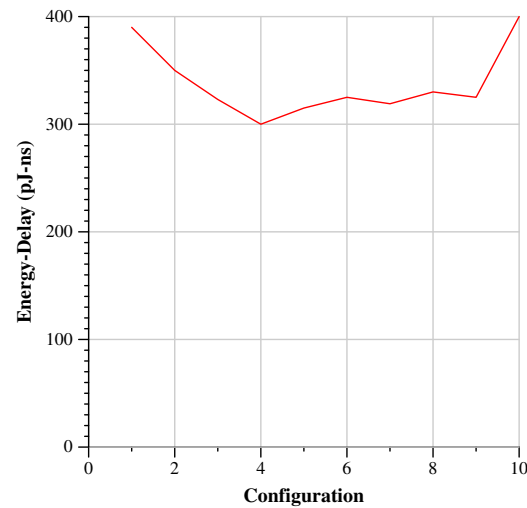
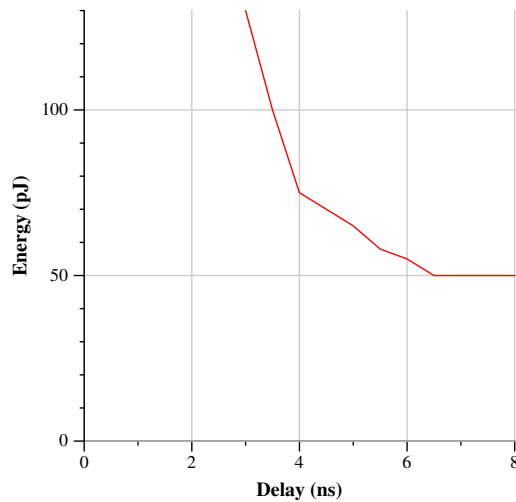
| Energy | Delay | ED | ED^2 |
|--------|-------|--------|----------|
| 5J | 2s | $10Js$ | $20Js^2$ |
| 2J | 5s | $10Js$ | $50Js^2$ |



Energy Example



Energy-Delay Product Redux



Roughly based on data from “Energy-Delay Tradeoffs in CMOS Multipliers” by Brown et al.



Raw Data

| Delay | Energy | ED | ED^2 |
|----------|-----------|------------|-------------|
| 3 | 130 | 390 | 1170 |
| 3.5 | 100 | 350 | 1225 |
| 3.8 | 85 | 323 | 1227 |
| 4 | 75 | 300 | 1200 |
| 4.5 | 70 | 315 | 1418 |
| 5 | 65 | 325 | 1625 |
| 5.5 | 58 | 319 | 1755 |
| 6 | 55 | 330 | 1980 |
| 6.5 | 50 | 390 | 2535 |
| 8 | 50 | 400 | 3200 |



Other Metrics

- $Energy - Delay^n$ – choose appropriate factor
- $Energy - Delay - Area^2$ – takes into account cost (die area) [McPAT]
- Power-Delay – units of Energy – used to measure switching
- Energy Delay Diagram – [SWEEP]
- Energy-Delay-FIT (reliability?)



Measuring Power and Energy



Why?

- New, massive, HPC machines use impressive amounts of power
- When you have 100k+ cores, saving a few Joules per core quickly adds up
- To improve power/energy draw, you need some way of measuring it



Energy/Power Measurement is Already Possible

Three common ways of doing this:

- Hand-instrumenting a system by tapping all power inputs to CPU, memory, disk, etc., and using a data logger
- Using a pass-through power meter that you plug your server into. Often these will log over USB
- Estimating power/energy with a software model based on system behavior



Measuring Power and Energy

- Sense resistor or Hall Effect sensor gives you the current
- Sense resistor is small resistor. Measure voltage drop.
Current $V=IR$ Ohm's Law, so $V/R=I$
- Voltage drops are often small (why?) so you may need to amplify with instrumentation amplifier
- Then you need to measure with A/D converter
- $P = IV$ and you know the voltage
- How to get Energy from Power?



Hall Effect Current Sensors

- Output voltage varies based on magnetic field.
- Current in wire causes magnetic field
- Voltage output is linear proportional to current
- Ideally little to no resistance (unlike sense resistor)
- Can measure higher current. 5, 20, 30A
- Need that? 100W CPU at 3.3V is roughly 30A



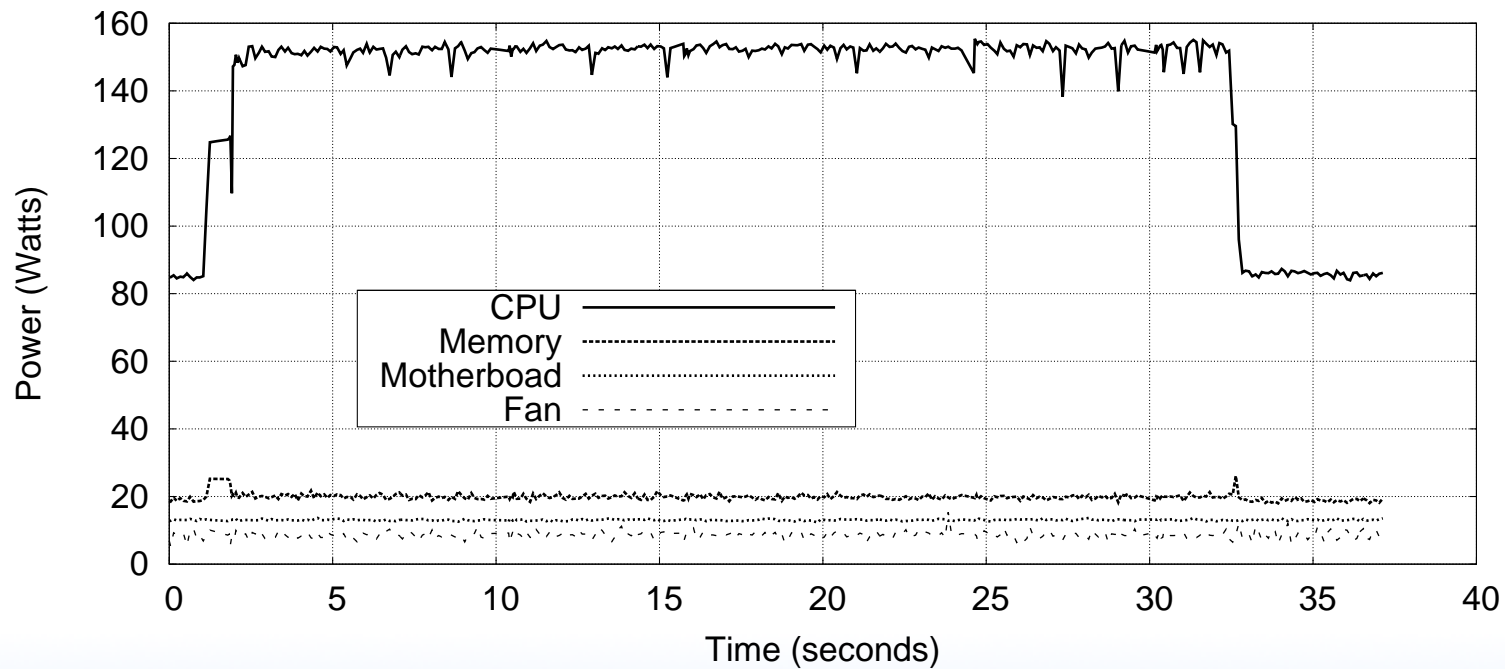
Other Issues

- Matching up internal and external measurements?
- Serial port? ntp? signal?
- Hard for small time intervals.



Existing Related Work

Plasma/dposv results with Virginia Tech's PowerPack



Powerpack

- Measure at Wall socket: WattsUp, ACPI-enabled power adapter, Data Acquisition System
- Measure all power pins to components (intercept ATX power connector?)
- CPU Power – CPU powered by four 12VDC pins.
- Disk power – measure 12 and 5VDC pins on disk power connector



- Memory Power – DIMMs powered by four 5VDC pins
- Motherboard Power – 3.3V pins. Claim NIC contribution is minimal, checked by varying workload
- System fans



PowerMon 2

- PowerMon 2 is a custom board from RENCI
- Plugs in-line with ATX power supply.
- Reports results over USB
- 8 channels, 1kHz sample rate
- We had hardware at UT, but managed to brick it



Shortcomings of current methods

- Each measurement platform has a different interface
- Typically data can only be recorded off-line, to a separate logging machine, and analysis is done after the fact
- Correlating energy/power with other performance metrics can be difficult
- How often can you measure (a lot happens on a CPU at 2GHz)



Watt's Up Pro Meter

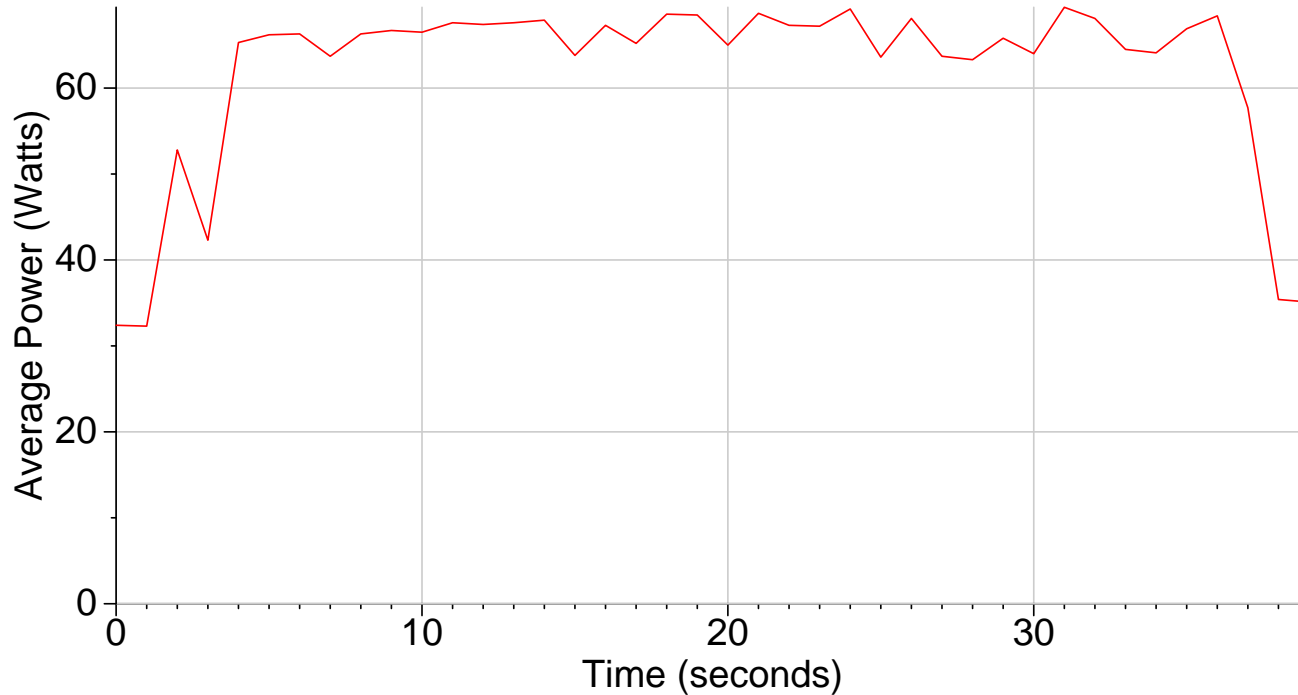


Watt's Up Pro Features

- Can measure 18 different values with 1 second resolution (Watts, Volts, Amps, Watt-hours, etc.)
- Values read over USB
- Joules can be derived from power and time
- Can only measure system-wide

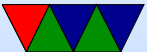


Watt's Up Pro Graph



PLASMA Cholesky Factorization N=10,000 threads=2

Measured on Core2 Laptop



Estimating Power

- Popular thing to do. One example: *Real Time Power Estimation and Thread Scheduling via Performance Counters* by Singh, Bhadauria and McKee.
- Have some sort of hardware measurement setup.
- Then measure lots of easy-to-measure things. Performance counters. Temperature. etc.
- Create a model (machine learning?) that can estimate
- Apparently using as few as 4 counters can give pretty good results



RAPL

- **R**unning **A**verage **P**ower **L**imit
- Part of an infrastructure to allow setting custom per-package hardware enforced power limits
- Also for TurboBoost
- User Accessible Energy/Power readings are a bonus feature of the interface



How RAPL Works

- RAPL is *not* an analog power meter (usually, Haswell-EP exception)
- RAPL uses a software power model, running on a helper controller on the main chip package
- Energy is estimated using various hardware performance counters, temperature, leakage models and I/O models
- The model is used for CPU throttling and turbo-boost, but the values are also exposed to users via a model-specific register (MSR)



Available RAPL Readings

- `PACKAGE_ENERGY`: total energy used by entire package
- `PP0_ENERGY`: energy used by “power plane 0” which includes all cores and caches
- `PP1_ENERGY`: on original Sandybridge this includes the on-chip Intel GPU
- `DRAM_ENERGY`: on Sandybridge EP this measures DRAM energy usage. It is unclear whether this is just the interface or if it includes all power used by all the DIMMs too

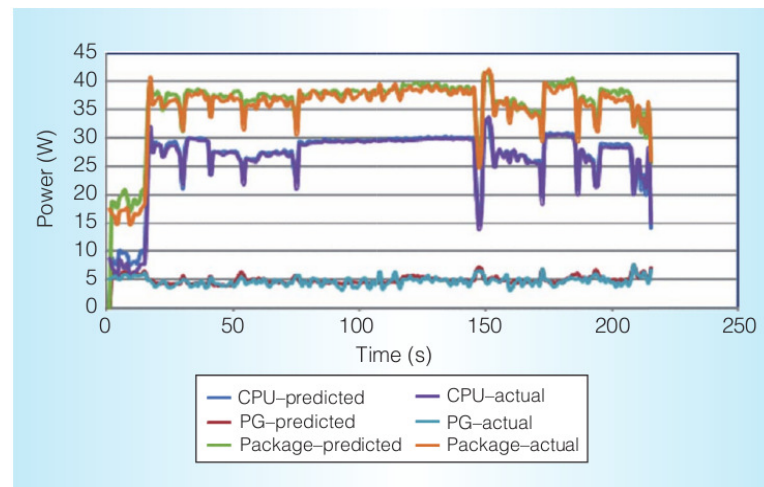


- SoC energy (skylake and newer?)



RAPL Measurement Accuracy

- Intel Documentation indicates Energy readings are updated roughly every millisecond (1kHz)
- Rotem et al. show results match actual hardware



Rotem et al. (IEEE Micro, Mar/Apr 2012)



RAPL Accuracy, Continued

- The hardware also reports minimum measurement quanta. This can vary among processor releases. On our Sandybridge EP machine all Energy measurements are in multiples of 15.2nJ
- Power and Energy can vary between identical packages on a system, even when running identical workloads. It is unclear whether this is due to process variation during manufacturing or else a calibration issue.

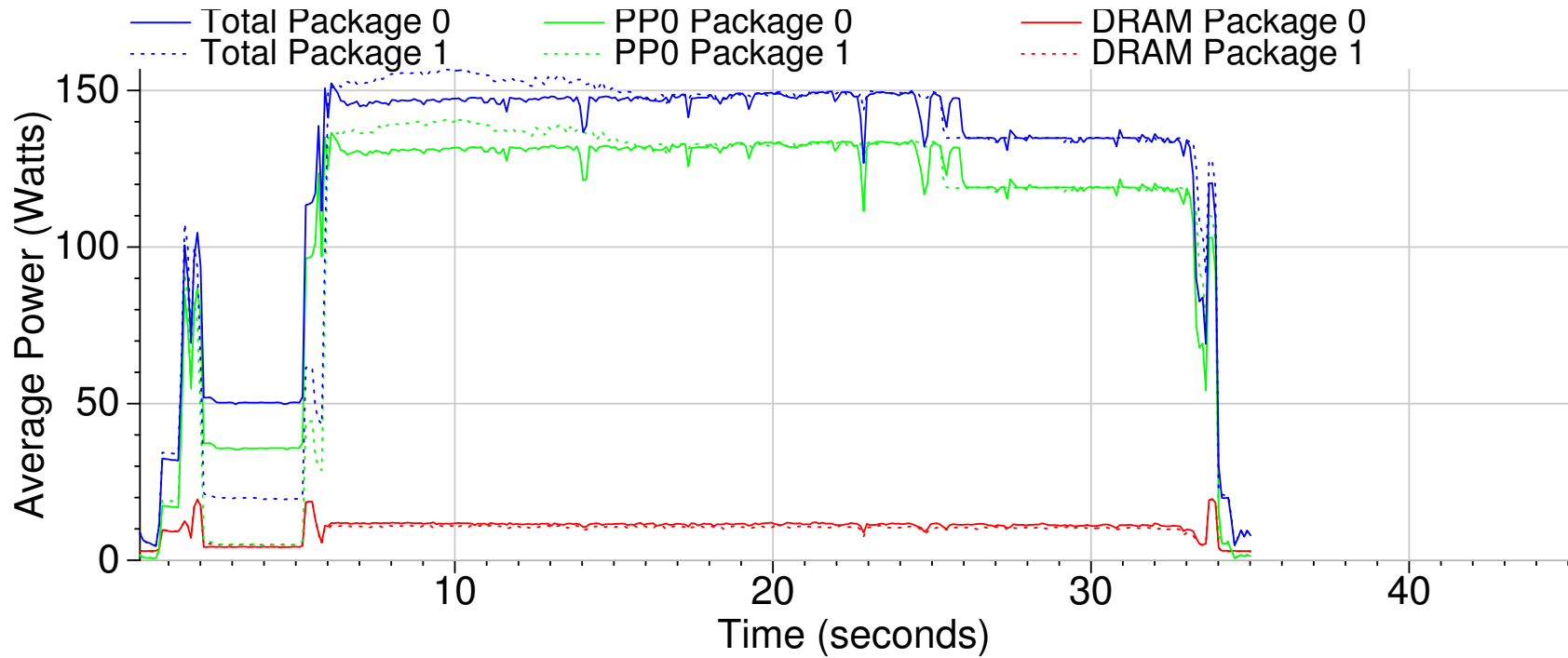


RAPL Validation

- The Dresden Paper
- My MEMSYS paper (include some plots?)

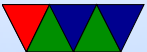


RAPL Power Plot

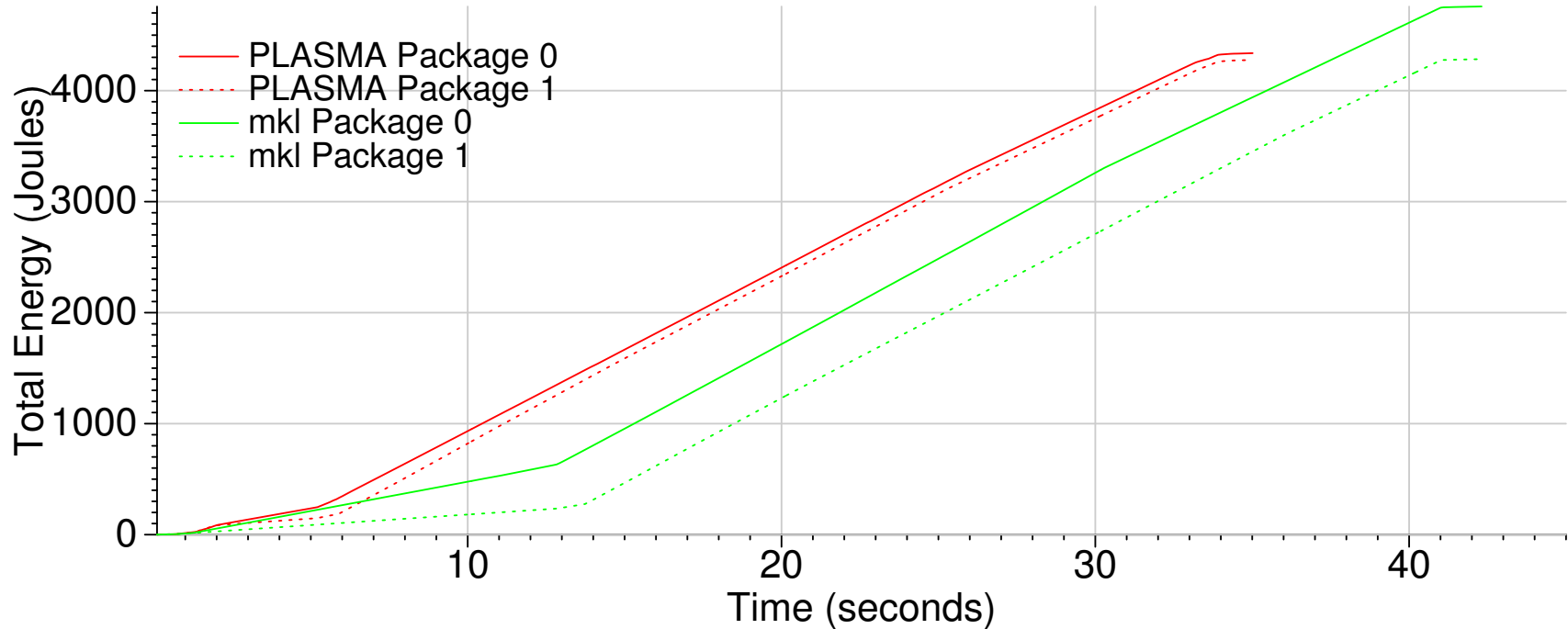


PLASMA Cholesky Factorization N=30,000 threads=16

Measured on SandyBridge EP

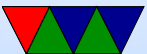


RAPL Energy Plot



Cholesky Factorization N=30,000 threads=16

Measured on SandyBridge EP

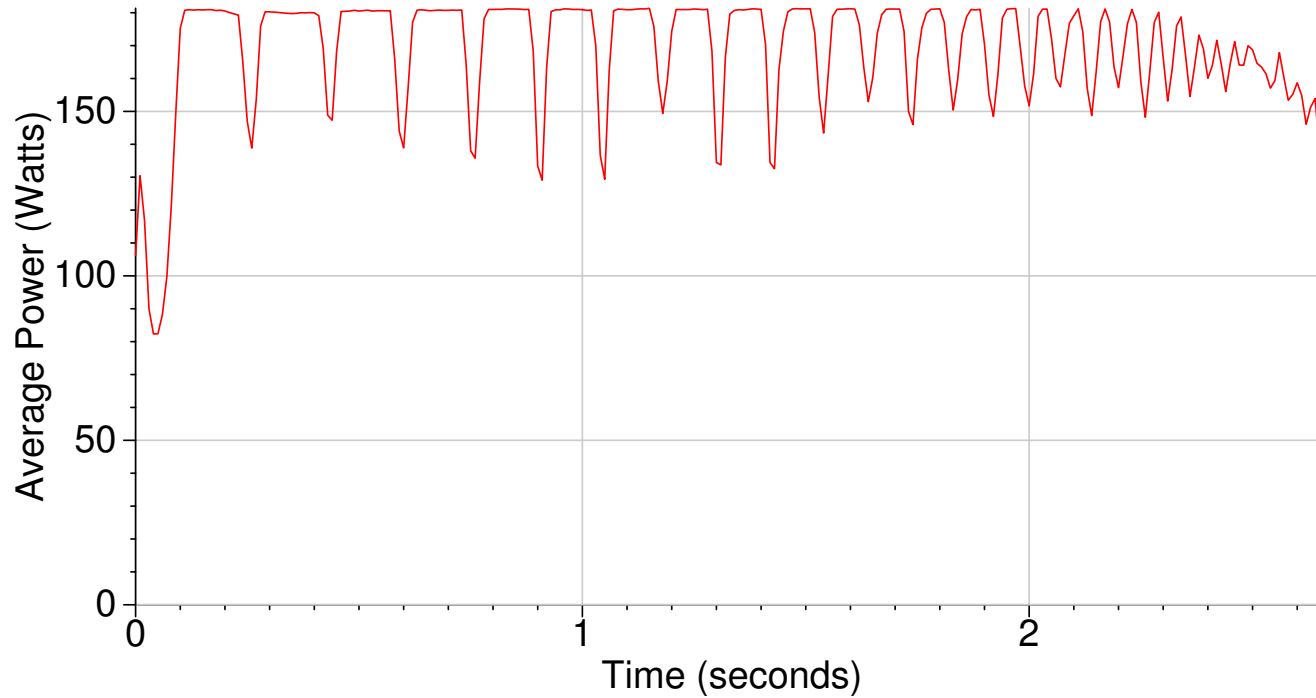


NVML

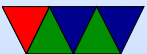
- Recent NVIDIA GPUs support reading power via the NVIDIA Management Library (**NVML**)
- On Fermi C2075 GPUs it has milliwatt resolution within $\pm 5W$ and is updated at roughly 60Hz
- The power reported is that for the entire board, including GPU and memory



NVML Power Graph



MAGMA LU 10,000, Nvidia Fermi C2075



AMD Application Power Management

- Recent AMD Family 15h processors also can report “Current Power In Watts” via the Processor Power in the TDP MSR
- Support for this can be provided similar to RAPL
- Have had bad luck getting accurate readings. Have found various chip errata on fam15h and fam16h hardware



Other ways to measure Power

- IPMI – many server machines have built in (low frequency) measurement of power supply values.
- Thermal? IR camera? Can see how much individual parts of chip use.
Overheat? Use IR transparent liquid to cool it?



Using RAPL

- On Linux, at least 4 ways to get these values
- Read msr directly, either with instruction or `/dev/msr`.
Need root as you can do bad things with msrs. “safemsr”
- `perf_event`
- `hwmon/powercap (/sys/class/powercap/)`
- `lmsensors`



Listing Events

```
$ perf list
...
power/energy-cores/           [Kernel F
power/energy-gpu/             [Kernel F
power/energy-pkg/             [Kernel F
power/energy-ram/             [Kernel F
...
```

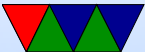


Measuring

```
$ perf stat -a -e power/energy-cores/,power/energy-ram/,instru
```

```
Performance counter stats for 'system wide':
```

```
        63.79 Joules power/energy-cores/
         2.34 Joules power/energy-ram/
21038123875      instructions          #    1.06
19782762541      cycles
3.407427702 seconds time elapsed
```



Measuring

- The key is -a which enables system-wide mode (needs root too if not configured as such)
- Why do you need system-wide?
- What does that do to the other metrics?

