



# Parallel Programming and Heterogeneous Computing

## E1 - Energy-Aware Computing

Sven Köhler, Lukas Wenzel, Max Plauth, and Andreas Polze  
Operating Systems and Middleware Group

# About this Lecture

This unit is a highly condensed version of the *Energy-Aware Computing Systems (EASY)* lecture by Prof. Dr.-Ing. Timo Hönig (RUB, formerly FAU).

If you are interested in more content, check out the FAU-CS4 website or convince us to offer an entire semester-spanning lecture, here at HPI.

CS 4 / Lehre / SS 2020 / Energy-Aware Computing Systems

## Energy-Aware Computing Systems (EASY) im SS 2020

Home

News

Lecture

Overview

Slides

### Lecture Content

- Introduction
  - Overview
  - Organisation
- Fundamentals
  - Power, energy, and performance
    - causality (interdependencies, dimensions)

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart 2

[https://www4.cs.fau.de/Lehre/SS20/V\\_EASY/](https://www4.cs.fau.de/Lehre/SS20/V_EASY/)

# 1 Background

**ParProg21 E1  
Energy-Aware  
Computing**

Sven Köhler

Chart 3

# Our Computing Systems Use Massive Amounts of Energy

## SUPERCOMPUTER FUGAKU - SUPERCOMPUTER FUGAKU, A64FX 48C 2.2GHZ, TOFU INTERCONNECT D

[1]

Site:	RIKEN Center for Computational Science
System URL:	<a href="https://www.r-ccs.riken.jp/en/fugaku/project">https://www.r-ccs.riken.jp/en/fugaku/project</a>
Manufacturer:	Fujitsu
Cores:	7,630,848
Memory:	5,087,232 GB
Processor:	A64FX 48C 2.2GHz
Theoretical Peak (Rpeak)	537,212 TFlop/s
Nmax	21,288,960
Power Consumption	
Power:	29,899.23 kW (Optimized: <b>26248.36 kW</b> )

Power Measurement Level: 2



Energy Research & Social Science  
Volume 38, April 2018, Pages 128-137



[2]

Original research article

### Digitalisation, energy and data demand: The impact of Internet traffic on overall and peak electricity consumption

Janine Morley <sup>a</sup>, Kelly Widdicks <sup>b</sup>, Mike Hazas <sup>b</sup>

[Show more](#)

<https://doi.org/10.1016/j.erss.2018.01.018>

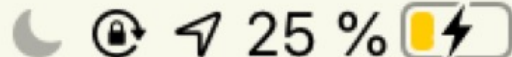
Under a Creative Commons license

[Get rights and content](#)

[open access](#)

### Abstract

Over the last decade, concerns have been raised about increases in the electricity used by information technologies, other consumer electronic devices, data centres



**ParProg21 E1**  
**Energy-Aware**  
**Computing**

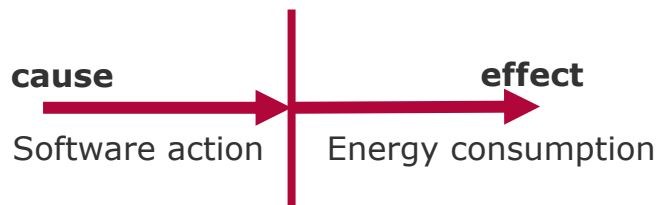
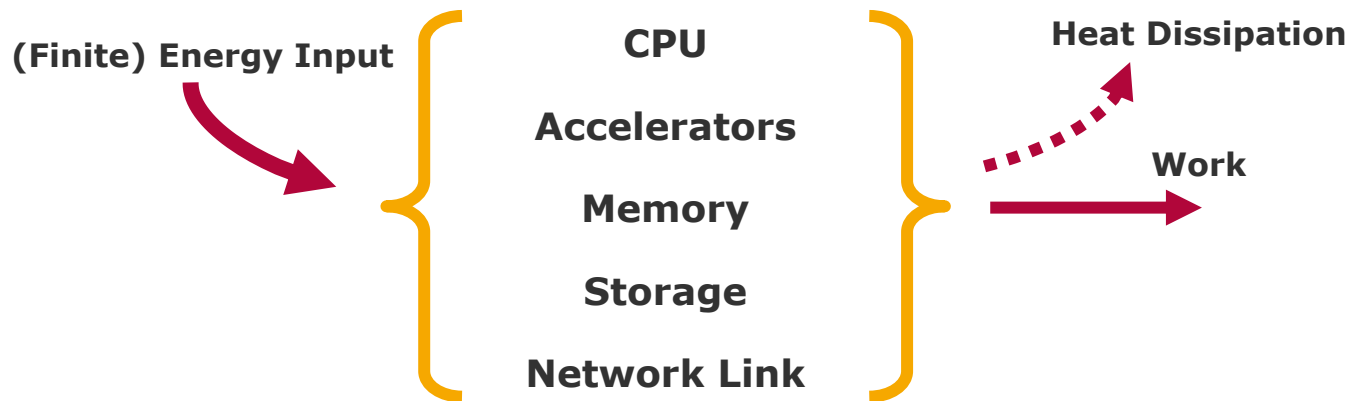
Sven Köhler

[1] Fugaku Supercomputer, Top500 List, Acc. 2021-06-22. <https://www.top500.org/system/179807>

[2] Morley, J., Widdicks, K., & Hazas, M. (2018). Digitalisation, energy and data demand: The impact of Internet traffic on overall and peak electricity consumption. Energy Research & Social Science, 38, 128-137.

Chart 4

# What Consumes Energy?



**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart **5**

# Energy vs. Power

The energy demand  $E$  that is required to execute an operation is the integral over the system's power demand from start ( $t_s$ ) to end ( $t_e$ ) of the operation.

$$E = \int_{t_s}^{t_e} P(t) dt$$

Energy  $E$  (unit: J or Ws) is the ability to do work.

$E$  is a suitable metric for:

- your battery life
- your electricity bill
- your carbon footprint

Power  $P$  (unit: W or J/s) is the rate of doing work.

$P$  is a suitable metric for:

- power supply constraints (peak power)
- prediction of heat dissipation (cooling facilities)

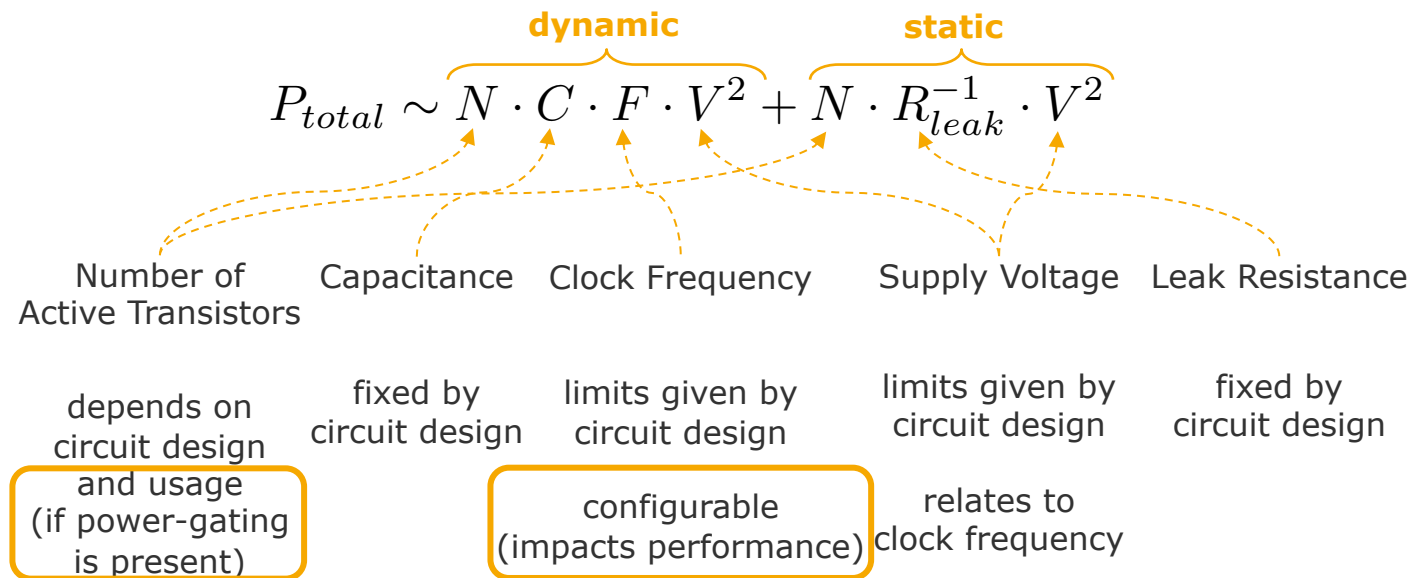
**Reducing the energy demand requires to reduce the run-time or the power demand.**

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart 6

# Power Demand of Computing Circuits



**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

**Reducing the power demand requires to shut off transistors or reduce the clock frequency.**

Chart 7

# Energy Management

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart 8



# Dynamic Voltage and Frequency Scaling

---

Modern compute architectures allow developers to actively regulate voltage and clock frequency at a fine time granularity (tens of milliseconds).

Examples:

- Intel CPUs: RAPL using e.g., `powergov` or direct control register access
- IBM POWER CPUs: EnergyScale via CIM or HMC
- ARM: Plenty of tools and libraries, usually by SOC/board vendor
- NVidia GPUs: `nvidia-smi` or NVidia Management Library
- AMD GPUs: In the Linux sysfs at `/sys/class/drm/.../pp_od_clk_voltage`

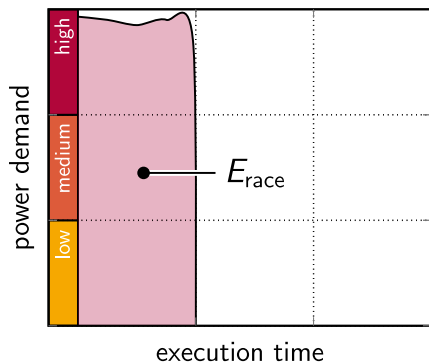
Proper power-gating is tricky. Without, your **idling** core **is wasting energy**. Thus, **minimize** the **idle time!** Put your cores to sleep, when you can.

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart 9

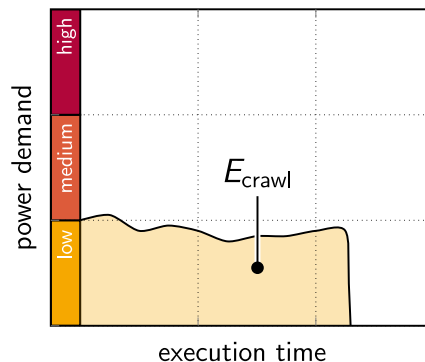
# Race or Crawl to Sleep?



## race-to-sleep

Maximize sleep time using a blocking management method after finishing pending work.

Suits especially compute-intensive processes



## crawl-to-sleep

Configure system at minimum voltage and clock rate, aiming for low average/peak power.

Suits especially I/O- or memory-bound processes

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart **10**

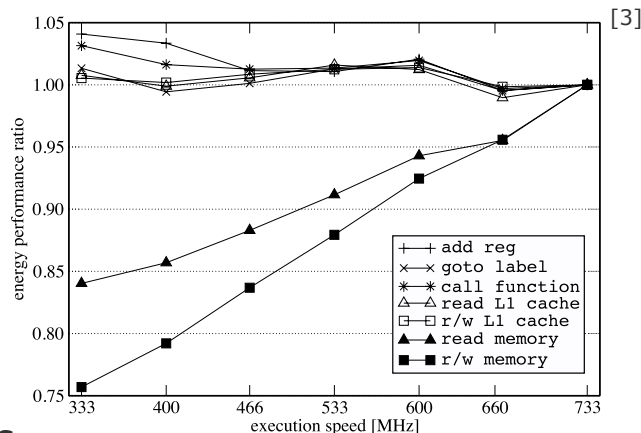
# Data Processing And Computing

A **naïve approach** to energy-aware computing:

*Run memory-bound and CPU-bound threads with low and high clock speed, respectively.*

Problems of this approach:

- dynamic characteristics of workloads
- simple system model (#cores, interlocked voltages, cache size)
- input-dependent, variable size of working set
- costs for frequency switching



**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart 11

# Memory-aware Scheduling (Combining) I

## Observation:

Contention between cores due to resource demand (caches, memory) leads to run-time penalties (depending on thread characteristics).

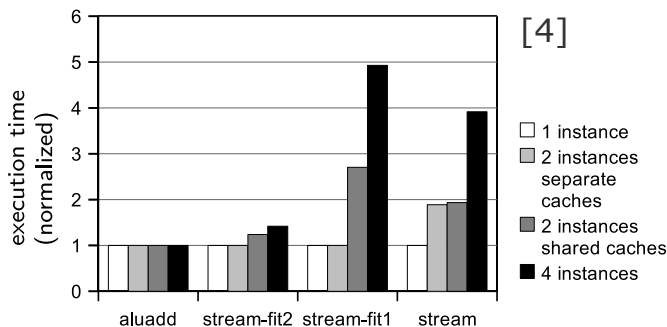


Figure 1. Normalized runtime of microbenchmarks running on the Core2 Quad

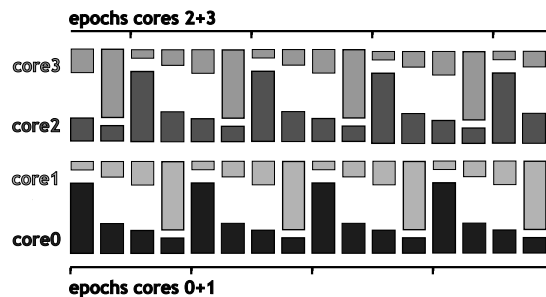


Figure 4. Sorted scheduling. Bars correspond to memory intensity.

## Proposed strategy:

Combine and co-locate compute-bound and memory-bound threads to reduce contention (Gang scheduling<sup>[5]</sup>)

**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
Sven Köhler

[4] Merkel, A., Bellosa, F.: Memory-aware Scheduling for Energy Efficiency on Multicore Processors. In: Proceedings of the Workshop on Power Aware Computing and Systems (HotPower'08), 2008, S. 123–130

[5] Ousterhout, J. K. et. al.: Scheduling Techniques for Concurrent Systems. In: Proceedings of the 1982 International Conference on Distributed Computing Systems (ICDCS'82) Bd. 82, 1982, S. 22–30

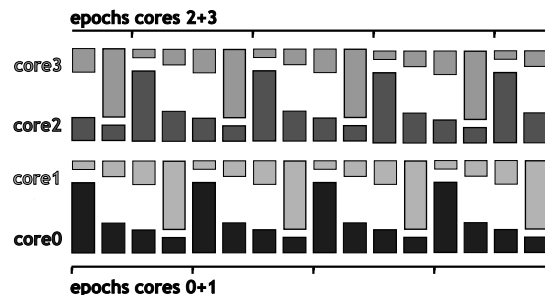
# Memory-aware Scheduling (Combining) II

## Implementation:

- group CPU cores into pairs of two
- Run threads with complementary resource demands on each pair
- Scale to **lowest** frequency if **no** compute-bound threads are ready (only memory-bound threads ready)
- Scale to **highest** frequency if **at least one** compute-bound thread is ready

## Limitations and Considerations:

- inferences with kernel scheduling strategy (risks priority inversion)
- scheduling policy only effective for specific working set sizes
- memory hierarchy and cache sizes must be considered



**Figure 4.** Sorted scheduling. Bars correspond to memory intensity.

# Access and Execute (Sequencing)

Sequenced execution the extend phases of homogenous operations.

Reorder your instructions into two streams operations of the same kind



- |   |   |
|---|---|
| ▪ prefetch data into caches, write intermediate results to memory | ▪ Execute operations on data in hot caches (i.e., computations) |
| ▪ run with low clock speed  | ▪ run with high clock speed                                     |

Eliminates unnecessary CPU stalling and memory waits, but requires some compiler support and might cause additional synchronization efforts.

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

[6] Smith, J. E.: Decoupled Access/Execute Computer Architectures. In: Proceedings of the 9th Annual Symposium on Computer Architecture (ISCA'82), 1982, pp 112–119

[7] Koukos, K., Black-Schaffer, D., Spiliopoulos, V., Kaxiras, S.: Towards More Efficient Execution: A Decoupled Access-execute Approach. In: Proceedings of the 27th International ACM Conference on International Conference on Supercomputing (ICS'13), 2013

## What else?

---

- Pick a more energy-efficient system (e.g., FPGA over CPU or high-efficiency cores like on ARM big.LITTLE machines)
- Optimize your algorithm!
- Optimize your implementation for performance, go to sleep
- But: Fast systems may use more energy than they save in time<sup>[8]</sup>

*You will never know if your algorithm, implementation or management strategy is more energy efficient than another, unless you measure ...*

**ParProg21 E1  
Energy-Aware  
Computing**  
Sven Köhler

[8] Hönig, T., Janker, H., Eibel, C., Mihelic, O., & Kapitza, R. (2014). Proactive Energy-Aware Programming with PEEK. In 2014 Conference on Timely Results in Operating Systems (TRIOS 14).

# Measuring Power and Energy

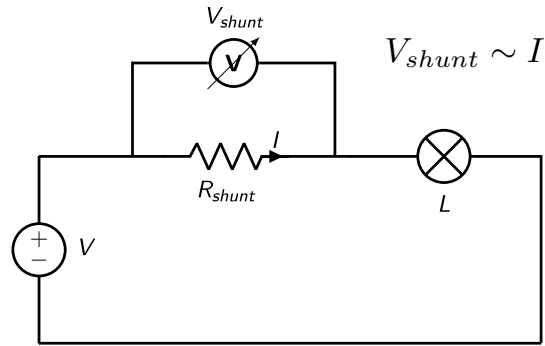
**ParProg21 E1  
Energy-Aware  
Computing**

Sven Köhler

Chart 16



# Measurement Methods



## physical measurements

Direct or indirect physical method, like measuring the voltage drop across a resistor.

Quite accurate, little overhead, requires setup alteration

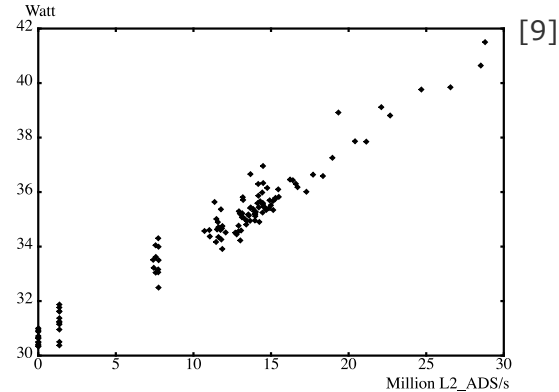


FIG. 3. Correlation of L2 Cache references and energy consumption

## logical measurements

Based on a software power model, initially build upon physical measurements.

No additional circuits required, but model might be error-prone

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

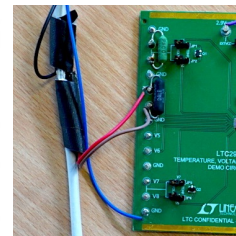
Sven Köhler

Chart **17**

# Measurement Facilities

## External

i.e., standalone devices intercepting the supply between power source and measured device



## On-Board

Part of the mainboard or SOC, often allow for distinction of separate power rails

IPMI, BMC,  
Jetson counters

## On-Chip

Integrated with the individual hardware platform, allows for most details

RAPL, PowerOCC  
Apple M1 counters

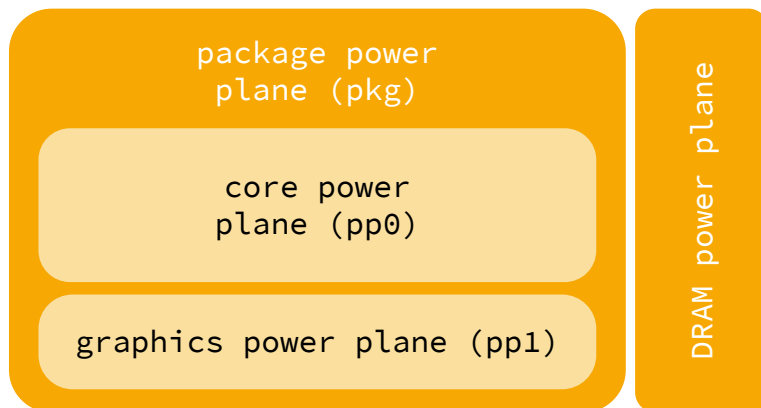
**ParProg21 E1  
Energy-Aware  
Computing**

Sven Köhler

Chart **18**

# Running Average Power Limit (RAPL)

- Available for Intel platforms, since Sandy Bridge
- Registers capture cumulative energy consumption (not power draw), at ~1 ms resolution (wrap around after ~60s)
- Accessible via control registers, Linux sysfs, or perf\_event\_open
- Semi-compatible AMD implementation since Ryzen Gen 3



**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

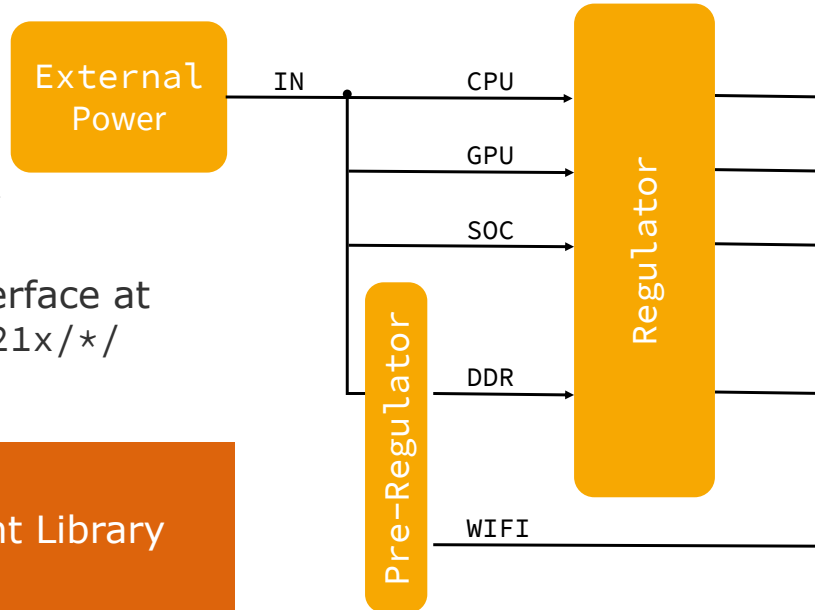
Chart **19**

# NVidia Jetson TX2 Boards

Two triple-channel INA3221 power monitors:

- report averaged power draw, voltage and current
- estimated 5% sample accuracy, 20 Hz sampling frequency
- I<sup>2</sup>C exposed via Linux sysfs-interface at `/sys/bus/i2c/drivers/ina3221x/*/iio_device/in_power`

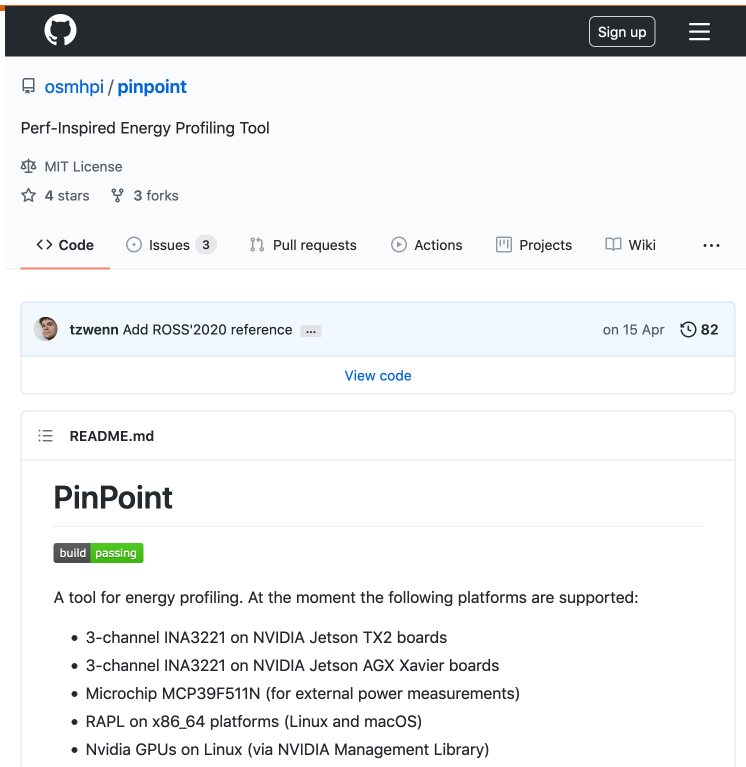
For all other NVidia GPUs:  
Check out the NVidia Management Library (nvmml) or the nvidia-smi tool.



**ParProg21 E1  
Energy-Aware  
Computing**  
Sven Köhler

Chart **20**

# Shameless Self-Plug: PinPoint



osmhpi / pinpoint

Perf-Inspired Energy Profiling Tool

MIT License

4 stars 3 forks

Code Issues 3 Pull requests Actions Projects Wiki

tzwenn Add ROSS'2020 reference on 15 Apr 82

View code

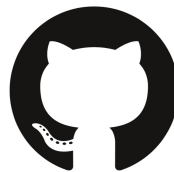
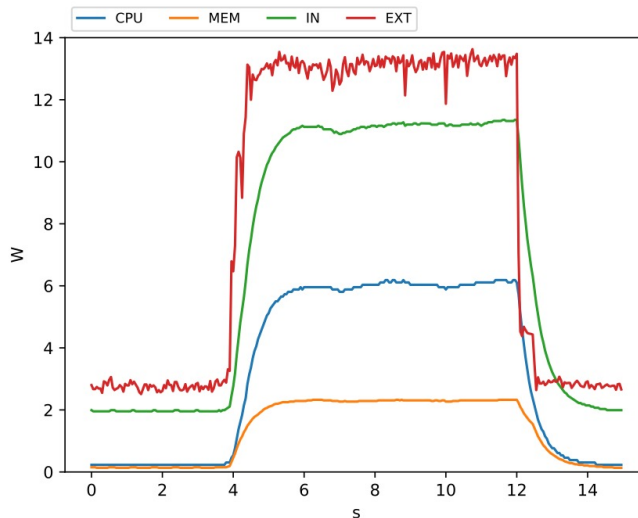
README.md

## PinPoint

build passing

A tool for energy profiling. At the moment the following platforms are supported:

- 3-channel INA3221 on NVIDIA Jetson TX2 boards
- 3-channel INA3221 on NVIDIA Jetson AGX Xavier boards
- Microchip MCP39F511N (for external power measurements)
- RAPL on x86\_64 platforms (Linux and macOS)
- Nvidia GPUs on Linux (via NVIDIA Management Library)



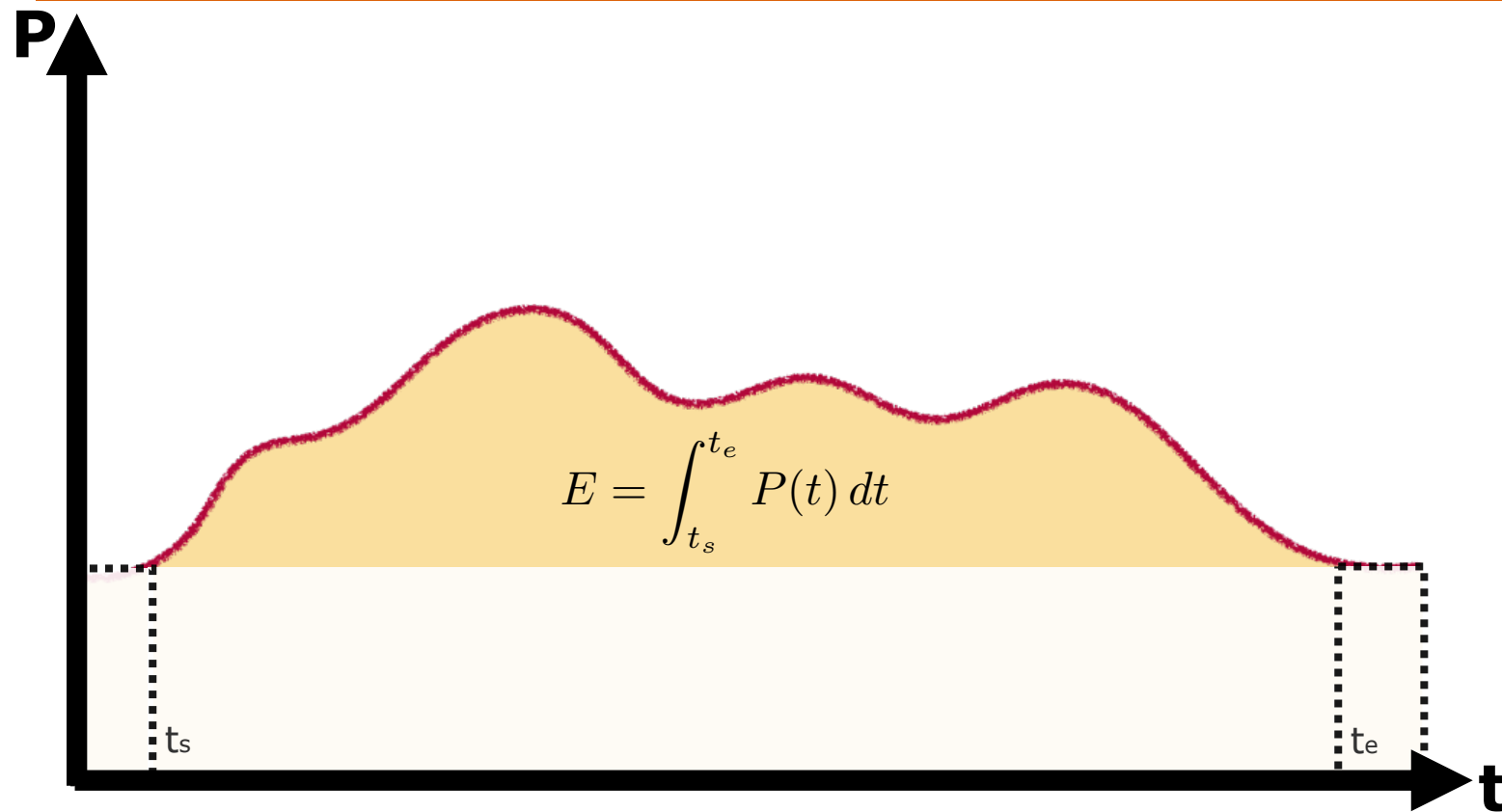
<https://github.com/osmhpi/pinpoint>

**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
 Sven Köhler

Chart 21

[12] Köhler, S., Herzog, B., Hönig, T., Wenzel, L., Plauth, M., Nolte, J., Polze, A., & Schröder-Preikschat, W. (2020, November). Pinpoint the Joules: Unifying Runtime-Support for Energy Measurements on Heterogeneous Systems. In *2020 IEEE/ACM International Workshop on Runtime and Operating Systems for Supercomputers (ROSS)* (pp. 31-40). IEEE.

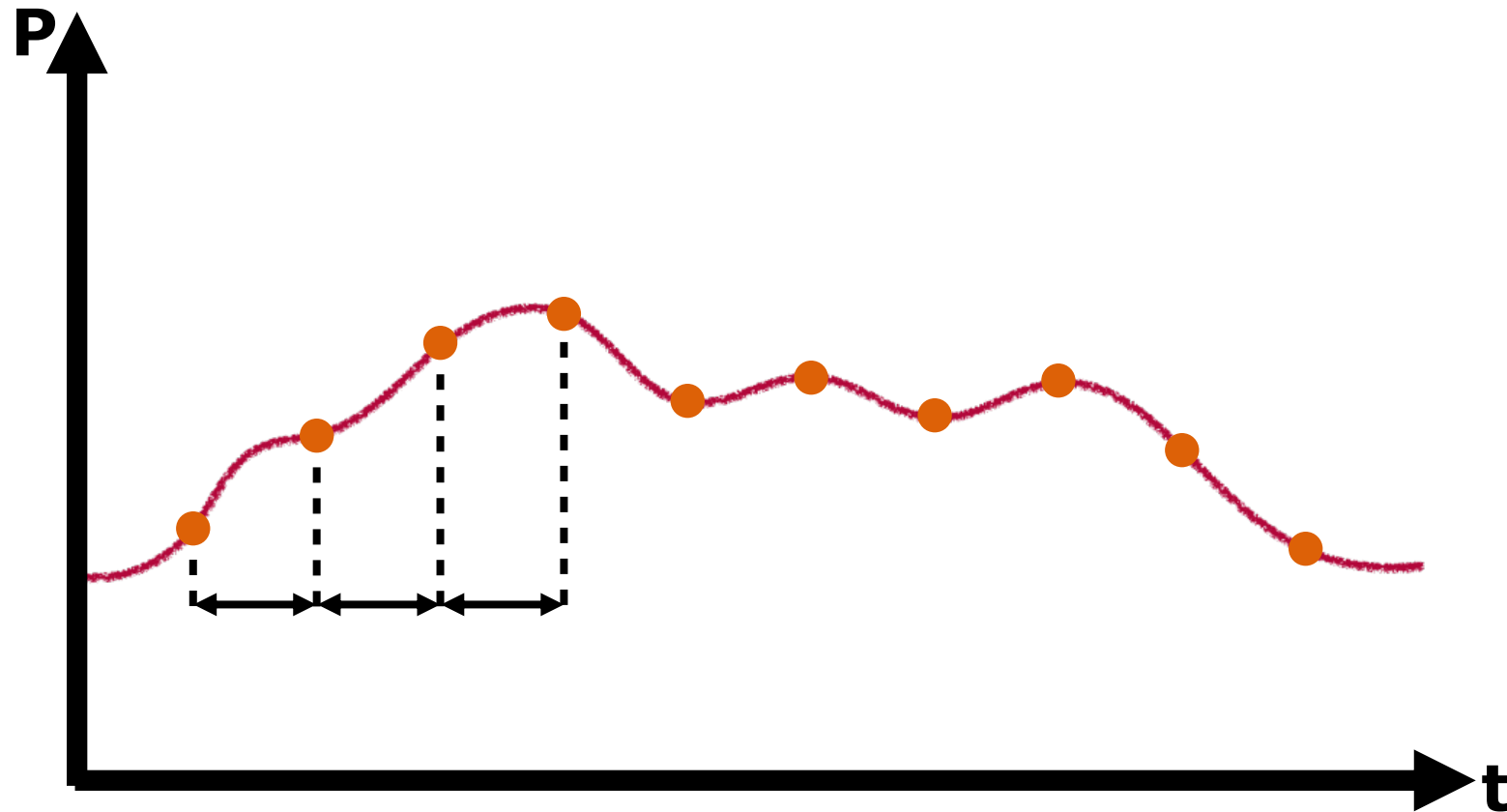
# Be Careful What You Measure



**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
Sven Köhler

Chart 22

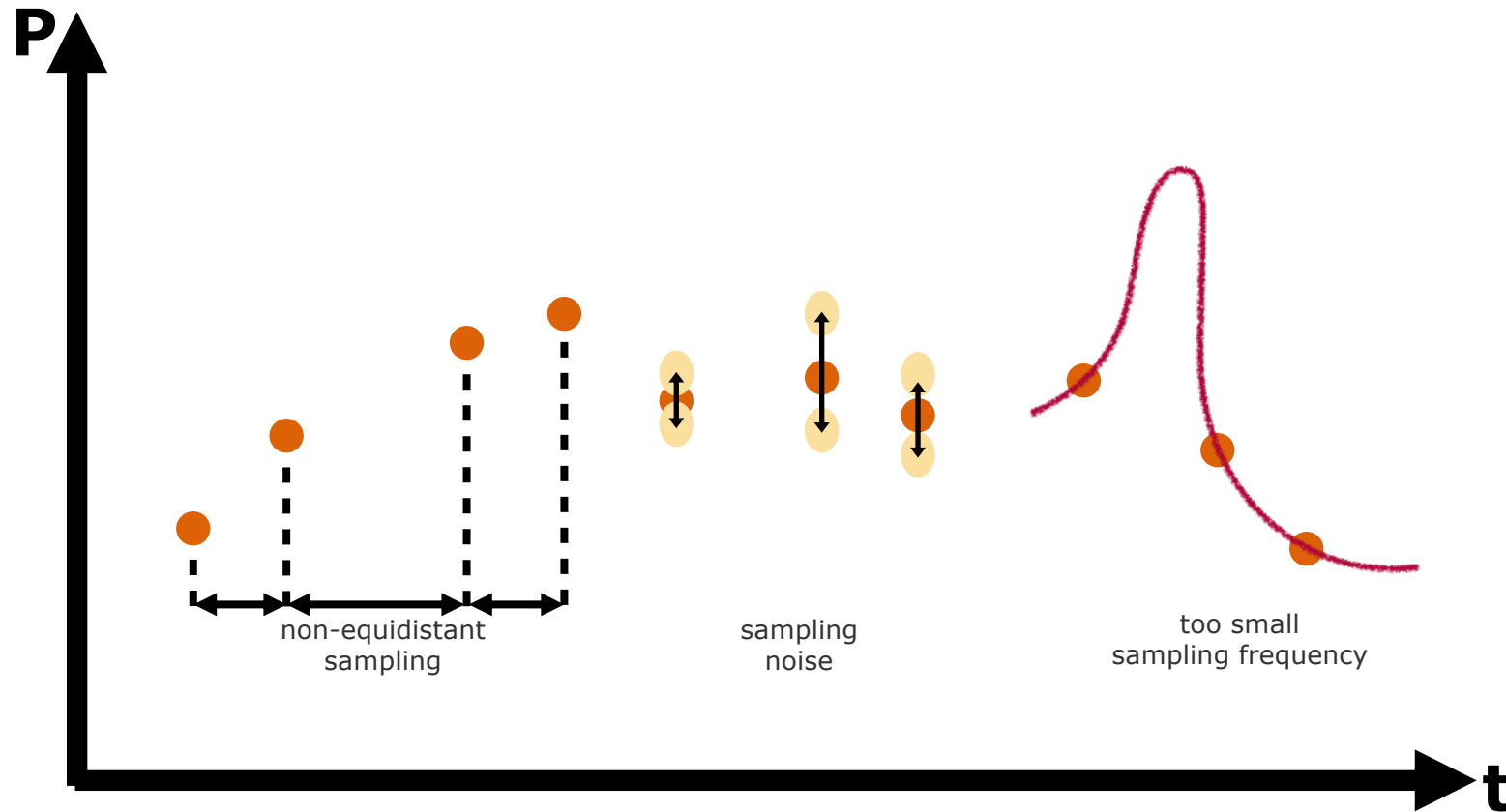
# Be Careful What You Measure



ParProg21 E1  
Energy-Aware  
Computing  
Sven Köhler

Chart 23

# Be Careful What You Measure



**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
 Sven Köhler

Chart 24



# Extended and Composite Metrics

- Power and energy demand are insufficient metrics
- Other system characteristics (e.g., **performance** or **latency**) may **differ** strongly, even though **power** or **energy** characteristics are **the same**
- Extended metrics combine basic metrics (power/energy demand) with additional system properties like execution time:
  - **Power-Delay Product (PDP)**:  $P_{\text{avg}} \cdot t$   
(approximates energy per switching event, good for fixed voltage)
  - **Energy-Delay Product (EDP)**:  $E \cdot t \simeq P_{\text{avg}} \cdot t \cdot t$   
(equal weight for changes of energy demand and performance, but misleading metric for systems with dynamic voltage scaling<sup>[13]</sup>)
  - **Energy-Delay-Squared Product (ED<sup>2</sup>P)**:  $\text{EDP} \cdot t$   
(good for fixed micro-architecture with dynamic voltage scaling<sup>[14]</sup>)

**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
Sven Köhler

[13] Horowitz, M., Indermaur, T., Gonzalez, R.: Low-power digital design. In: Proceedings of 1994 IEEE Symposium on Low Power Electronics, 1994, S. 8–11

[14] Brooks, D. M., Bose, P., Schuster, S. E., Jacobson, H., Kudva, P. N., Buyuktosunoglu, A., Wellman, J., Zyuban, V., Gupta, M., Cook, P. W.: Power-aware microarchitecture: design and modeling challenges for next-generation microprocessors. In: IEEE Micro 20 (2000), Nov, Nr. 6, S. 26–44

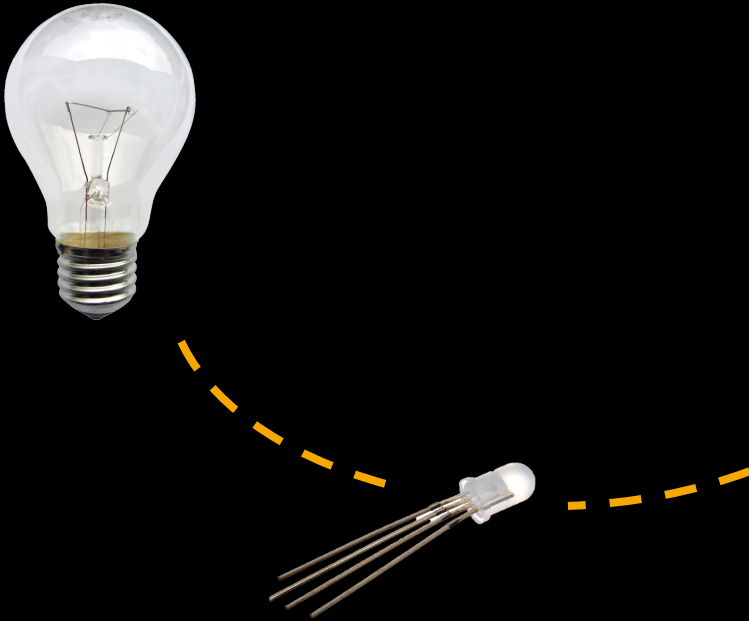
# 4 Closing Remarks

**ParProg21 E1**  
**Energy-Aware**  
**Computing**

Sven Köhler

Chart 26

# Rebound Effect

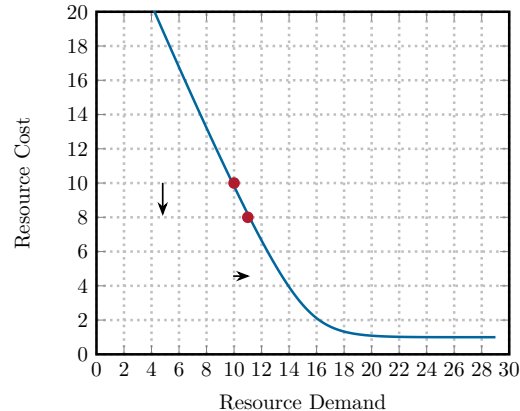


**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
Sven Köhler

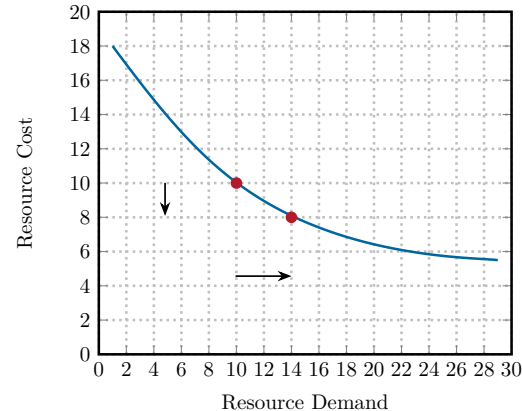
Chart **27**

# Jevons Paradox

*Technological progress that increases the efficiency with which a resource is used tends to increase (rather than decrease) the rate of consumption of that resource.*



increase efficiency by 20%  
 ⇒ increase demand by 10%  
 (Jevons Paradox does **not** apply)



increase efficiency by 20%  
 ⇒ increase demand by 40%  
 (Jevons Paradox **does** apply)

**ParProg21 E1**  
**Energy-Aware**  
**Computing**  
 Sven Köhler

Chart **28**

And now for a break and  
a glass of water\*.



\*or drink of your choice