



High Performance  
Computing &  
Big Data Services

-  [hpc.uni.lu](http://hpc.uni.lu)
-  [hpc@uni.lu](mailto:hpc@uni.lu)
-  [@ULHPC](https://twitter.com/ULHPC)



# NETCOM

Optimizing the Energy Consumption of Computation

Tobias Fischbach  
University of Luxembourg

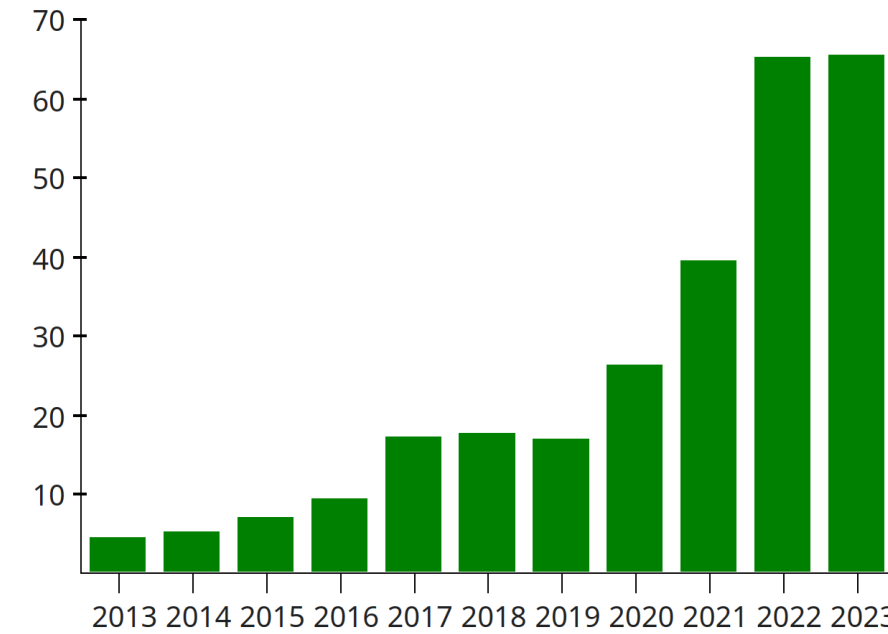
Pascal Bouvry  
University of Luxembourg

# Introduction

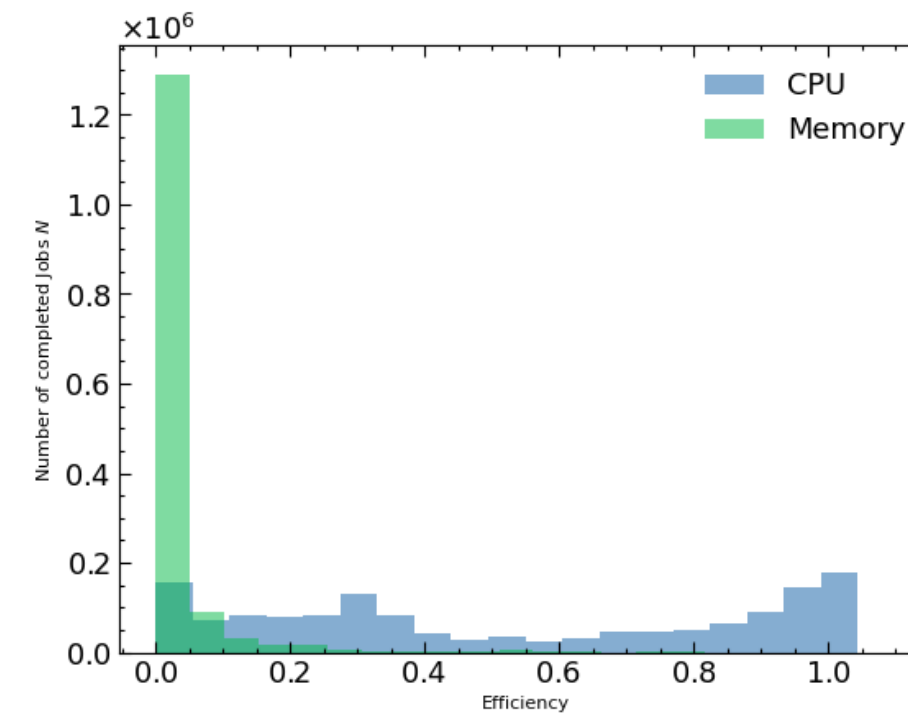
# Motivation

## Context

- ICT sector responsible for  $\approx 3-4\%$  of world wide power consumption<sup>1</sup>
- $\approx 20\%$  due to data centers
- continues research to improve energy efficiency of hardware (eg. Green500)<sup>2</sup>
- concerns are raised regarding the power demand of HPC-code in the era exascale computing
- focus on accelerators and mpi applications [<sup>3</sup>]<sup>4</sup>
- energy-aware scheduling successful in retrospects or tight bounds [<sup>5</sup>]<sup>6</sup>
- Software optimization is another tool to improve energy efficiency
- Average memory efficiency:  $\approx 2.675\%$
- Average cpu efficiency:  $\approx 52.460\%$



$\left[\frac{\text{GFlops}}{\text{W}}\right]$  of the Green500 winners per year<sup>2</sup>.



CPU and memory efficiency of all completed jobs on Iris (Feb. 2018 - Sep. 2023).

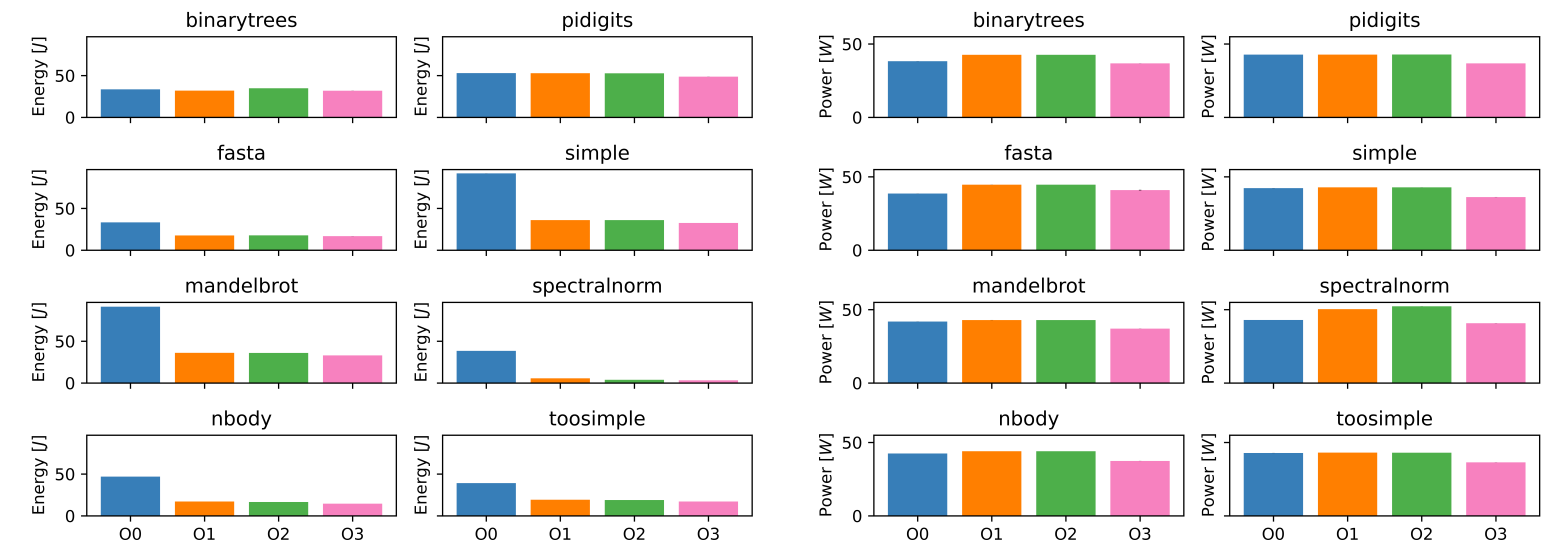


# The case for Software Optimization

## State of the Art

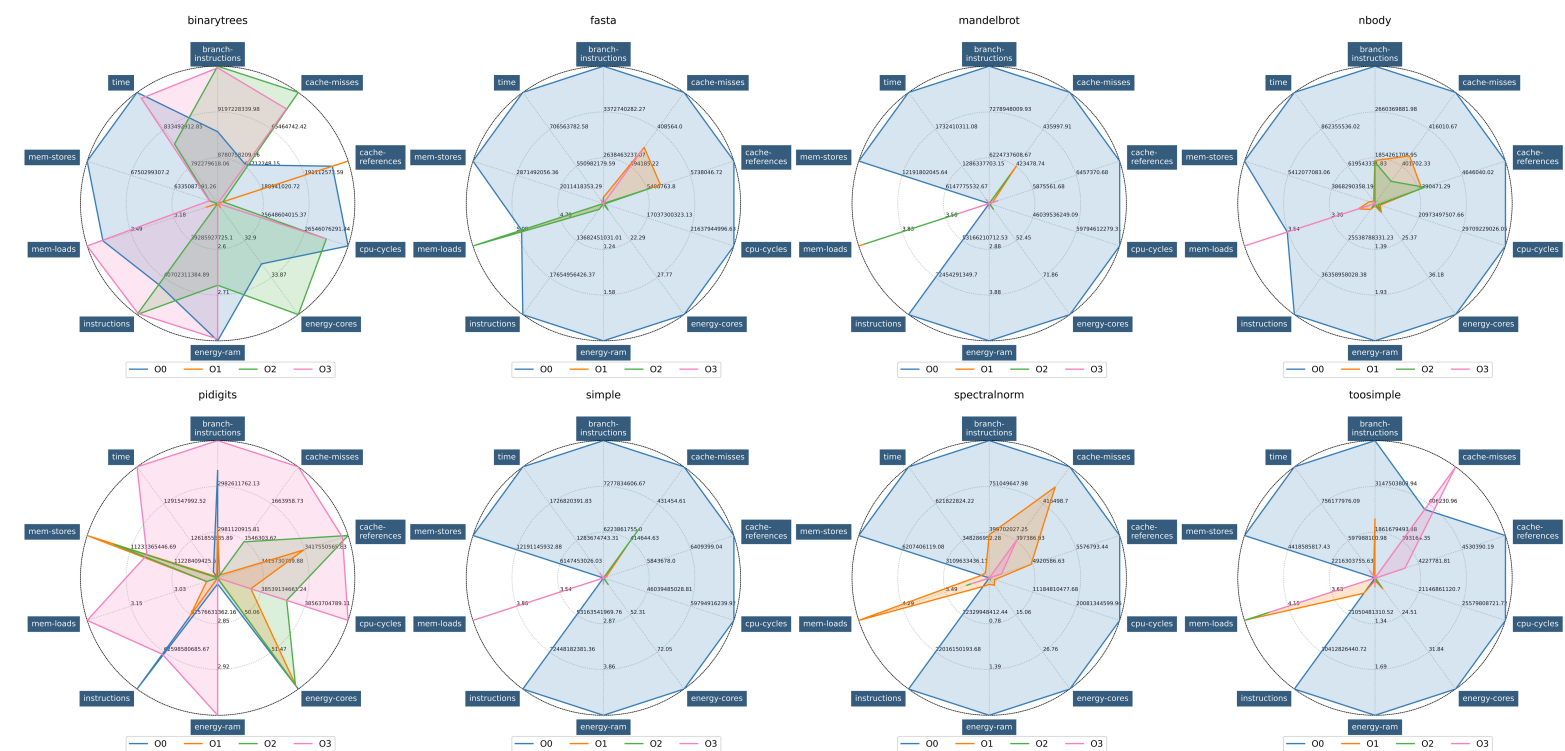
- faster software not necessarily more energy efficient
- cache misses and branching seem to be dominant drivers
- less optimized code can consume less power and / or less energy
- different runtime characteristics  $\rightarrow$  software categories

$\rightarrow$  What is the optimization pipeline that that minimizes the total energy consumption?



Total energy consumption of Clang for different opt. levels.

Average power consumption of Clang for different opt. levels.



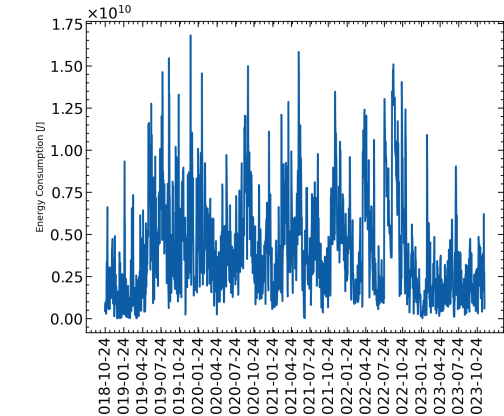
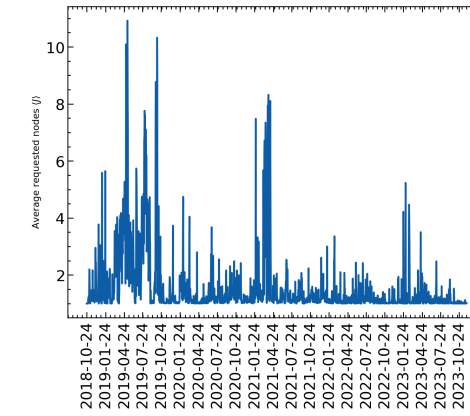
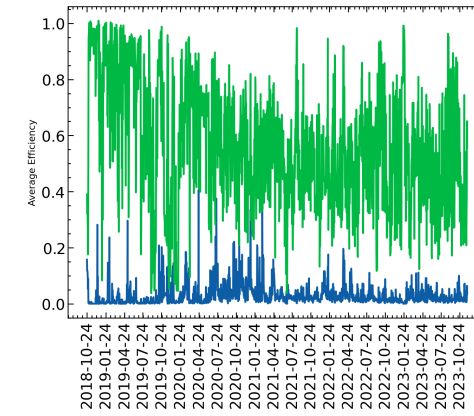
Runtime characteristics of Clang based opt. levels.



# Roadmap - PhD Sandwich

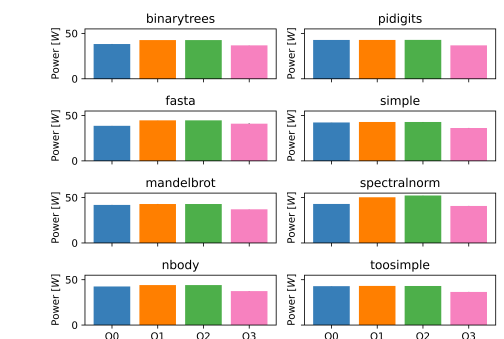
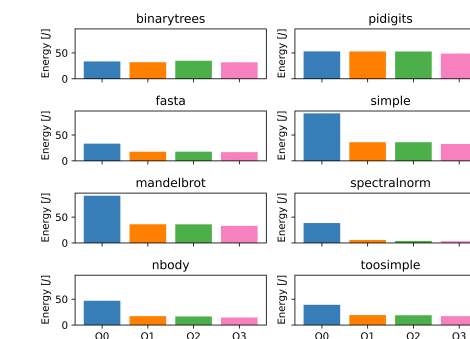
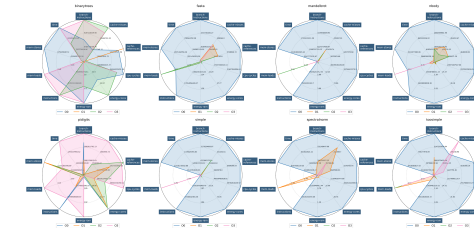
## 1. Top View: HPC Job Analysis

- environmental impact of HPC scheduling (PhD course project)
- clustering (done, but not promising)
- identify inefficient jobs upon submission



## 2. Middle View: Surrogate Software Optimization

- generate real-world and reproducible model
- provide an energy efficient pass pipeline for different software categories
- write own LLVM-Pass
- compare different CPU architecture



## 3. Bottom View:

- generate explanations
- guidelines for green code
- CPU architecture comparison



# Methodology

# Power measurements

## Software

- IPMI
- Intel RAPL (Linux, x86):
  - Likwid
  - PowerTOP
  - PowerStat
  - perf
- Syspower (macOS)

## Hardware

- Wattmeter
- PDUs
- Current sensors

## Challenges

- ARM and RISC-V require vendor support or extra hardware to monitor the power consumption
- Intel RAPL only supported by x86 and educated guessing

## Usage

- access counters with perf:

```
1 /proc/sys/kernel/perf_event_paranoid < -1
2 perf --event="power/energy-pkg/,power/energy-ram/" \
3     -j ./hal-9000.sh
```

- alternative direct access at “/dev/cpu/??/msr”
- available counters<sup>7</sup>:
  - energy-pkg := complete CPU socket
  - energy-cores := all CPU cores
  - energy-ram := DRAM
  - energy-psys := eDRAM and I/O Handling
- $f_{\text{update}} = 1 \text{ kHz}$

## Limitations

- no on-die circuits since Haswell-E  $\rightarrow$  application of CPU-specific offsets to historic data
- power management dynamics not properly captured
- package based readings (core based readings in new AMD Power Management driver)
- systemwide readings  $\rightarrow$  measurements include os and other running software (**measurement script**)
- scheduling effects
- thermal effects
- discrepancy in:
  - low DRAM usage<sup>8</sup>
  - dynamic CPU usage<sup>9</sup>
  - short Applications<sup>9</sup>



# Reducing the noise

## System changes

- power management  $\rightarrow$  force c/p-state
- CPU scheduler  $\rightarrow$  exclude all cores but one
- thermal effects  $\rightarrow$  limit CPU frequency
- measurement script effects  $\rightarrow$  limit perf to  $(f_{\text{update}}=100 \text{ Hz})$

$\rightarrow$  Reproducible measurements (within  $<5\%$ ) but lack of real-world connection

## Reproducible vs Real-World

- reproducible configuration for:
  - CPU Architecture comparisons
  - surrogate model used for explanation generation
- real-world configuration:
  - software optimization using a surrogate model

## Example Configuration

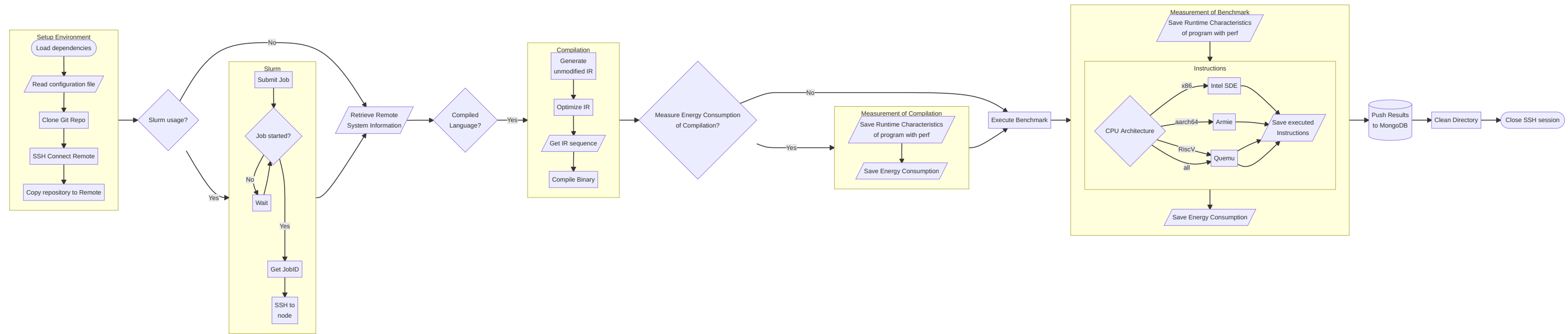
```
1 boot = {
2     kernelParams = [
3         "processor.max_cstate=1"
4         "intel_idle.max_cstate=0"
5         "intel_pstate=passive"
6         "intel_pstate=no_hwp"
7         "intel_pstate.no_turbo=1"
8         "isolcpus=1,2,3,4,5,6,7,8,9,10,11"
9         "mitigations=off"
10    ];
11    kernel.sysctl."kernel.perf_event_paranoid" = -1;
12    kernel.sysctl."kernel.kptr_restrict" = 0;
13 };
14 systemd.extraConfig = ''
15     CPUAffinity=0
16 '';
17 systemd.services.set-cpu-freq = {
18     wantedBy = [ "multi-user.target" ];
```

## Example Task Spawner

```
1 perf stat -j --event="power/energy-cores/" \
2     taskset -c 1-11 <command>
3 # or
4 perf stat -j --event="power/energy-cores/" \
5     numactl -physcpubind=1-11 <command>
```



# Experimental Pipeline - Overview



- reduce measurement overhead by separating experiment and data processing
- nix is used to create a reproducible software environment
- mongoDB for data management
- extract executed instructions Qemu, Intel SDE and Armie
- saves LLVM-IR by default
- modular (replace measurement method, integrate with optimization pipeline, ...)





# Experimental Pipeline - Usage

## Usage

1. Define building, jumping (opt.) and destination machine
2. Setup language specifics, including custom compiler passes
3. Specify how environments should be started (here nix)
4. mongoDB database and column setup
5. Define how programs should be called and how to measure

## Limitations

- Retrieval of executed Instructions through emulation slow
- Capturing runtime Characteristics with perf:
  - introduces overhead
  - depends on the argument order
- Relying on a running SSH service (extra overhead)

## Example Configuration

```
1 {
2     "dic_dest": {
3         "server": "<destination-adress>",
4         "username": "<username>",
5         "password": "<password-or-key>",
6         "port": <port>
7     },
8     "dic_jump": {
9         <jumphost-config>
10    },
11    "dic_local": {
12        <localhost-config>
13    },
14    "remote_dir": "<remote-dir>",
15    "language": "<language>",
16    "opt": "<custom-pass-pipeline>",
17    "env_loader_build": "<shell-script-or-command-to-load-build-env>",
18    "env_loader_exp": "nix-shell -Q <path-to-experiment-env>/shell.
19    "..."
20 }
```

## Example Call

```
1 $ python driver.py my-python-experiment.conf
2 [2023-11-11 17:54:22.678456]: Setup /home/tobias/src/python/bin/r
3 [2023-11-11 17:54:23.007983]: Enqueue runs.
4 [2023-11-11 17:54:23.008093]: Execute experiment.
5 [elapsed time: 0:02:13]:|*****
6 [2023-11-11 17:56:36.868175]: Setup experiment for language pyth
7 [2023-11-11 17:56:38.224519]: Start remote build.
8 [2023-11-11 17:57:12.296006]: Copy /tmp/build to /home/tobias on
```

How can we trust our measurements? \(\rightarrow\) Validation Experiment



# Programming Language Benchmark

## Game Rules

- no external libraries (eg. no numpy)
- code from the computer language benchmarks game<sup>10</sup>
- same input parameters
- consider consumed power by compilation and subtract baseline
- $$\begin{aligned} E^{\{\text{bench}\}} &= \& \left< E^{\{\text{bench}\}}_{i\right> + \left< E^{\{\text{comp}\}}_{i\right> - 2 * \left< E^{\{\text{base}\}}_{i\right> \text{ , } \left[\text{J} \right] \\ P^{\{\text{bench}\}} &= \& \frac{\left< E^{\{\text{bench}\}}_{i\right> + \left< E^{\{\text{comp}\}}_{i\right> - 2 * \left< E^{\{\text{base}\}}_{i\right>}}{\left< t_{i\right\}^{\{\text{bench}\}} \right>} \text{ , } \left[\frac{\{\text{J}\}}{\{\text{s}\}} \right] \end{aligned}$$

$\rightarrow$  reproduces<sup>11</sup> and<sup>12</sup>  $\rightarrow$   
**experimental pipeline is sound**

### The Computer Language 23.03 Benchmarks Game

“Which programming language is fastest?”

Top 5+ program performance comparisons –

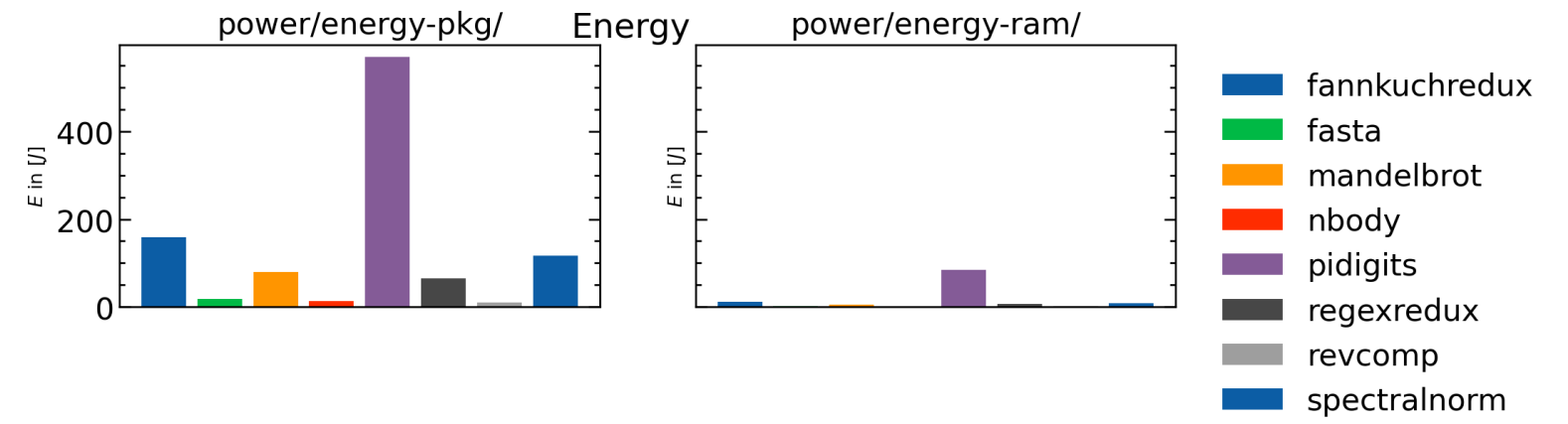
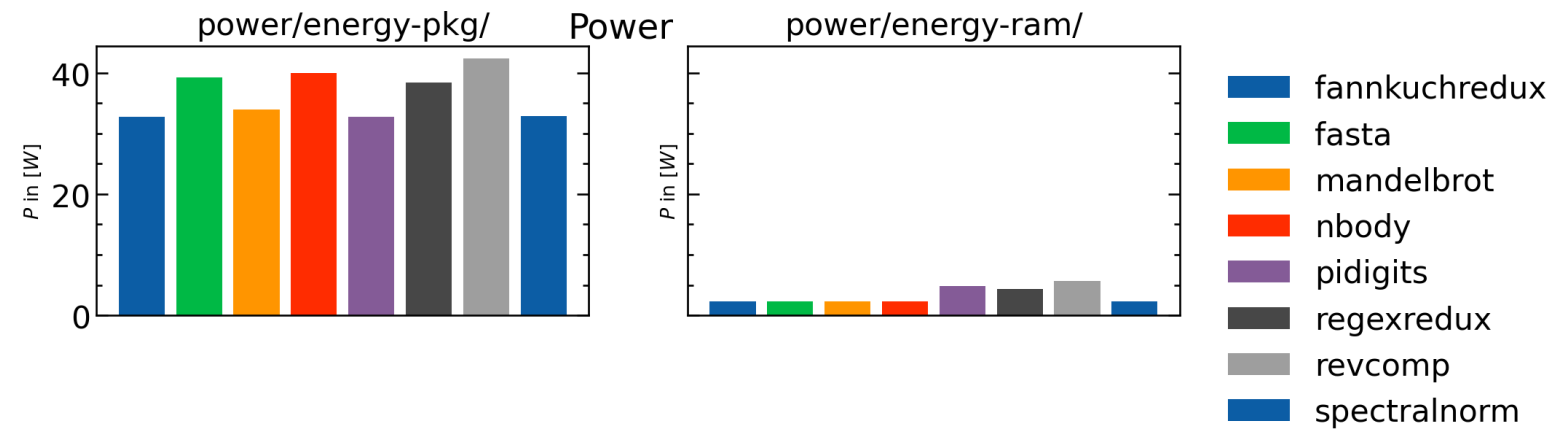
C# vs Java    Go versus Java  
Ruby vs Python    Rust versus C++  
Rust vs Go

Compare measurements of a transliterated program –

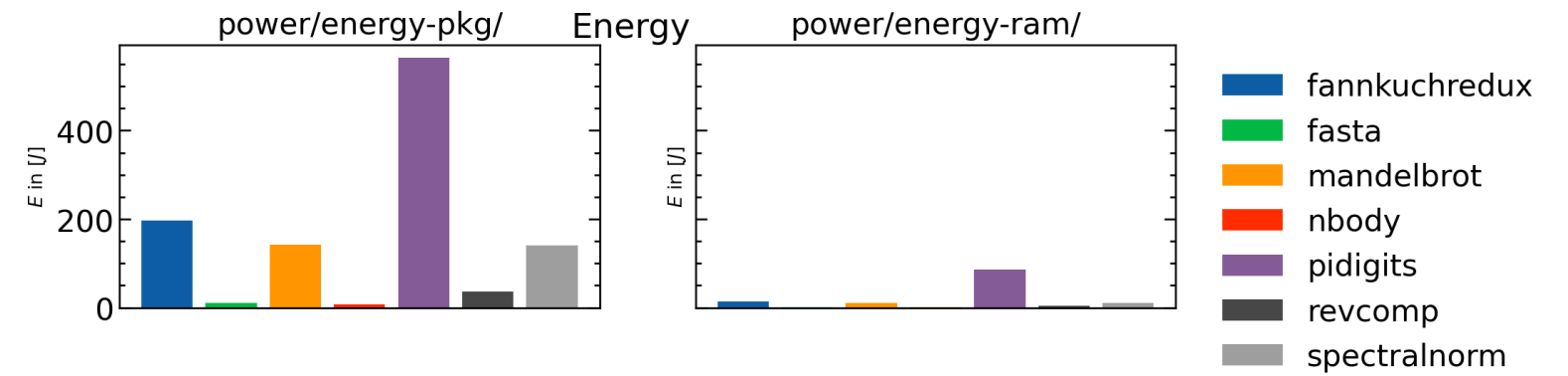
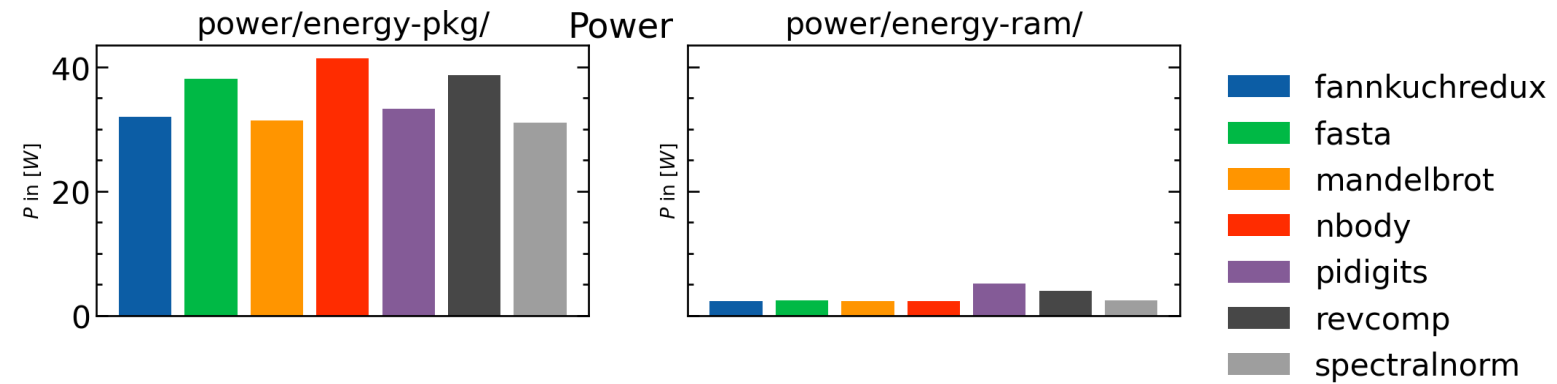
too simple    simple



# C/C++

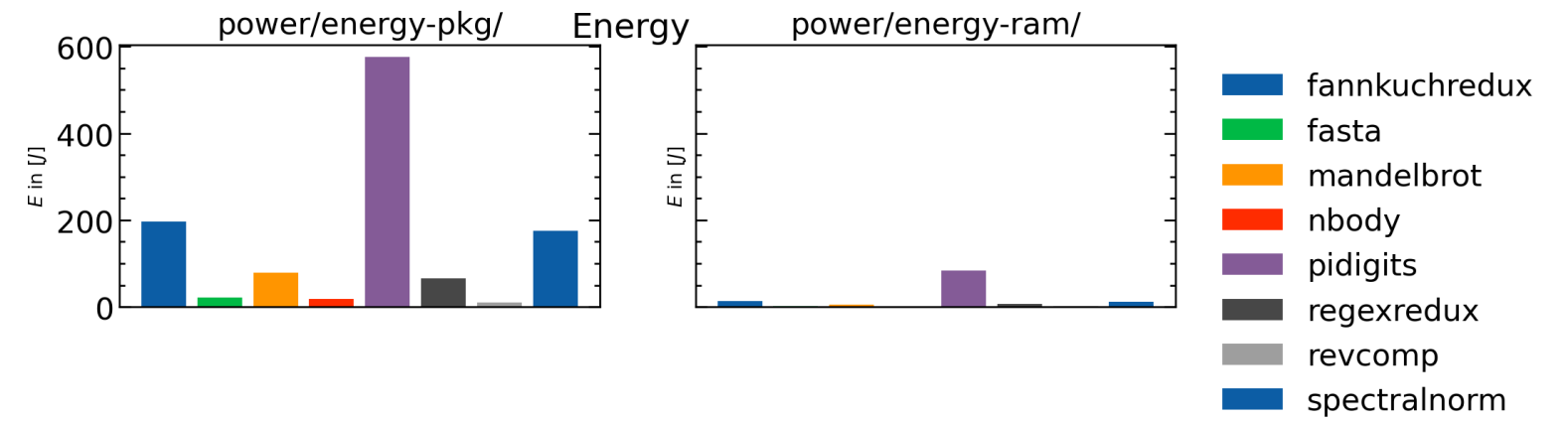
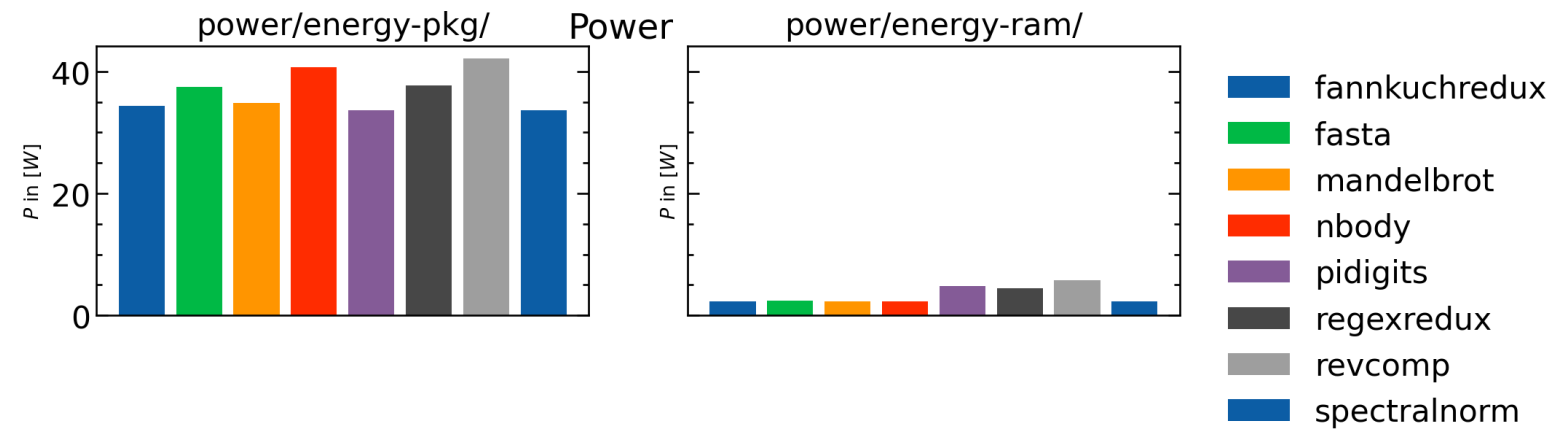


## Clang 16

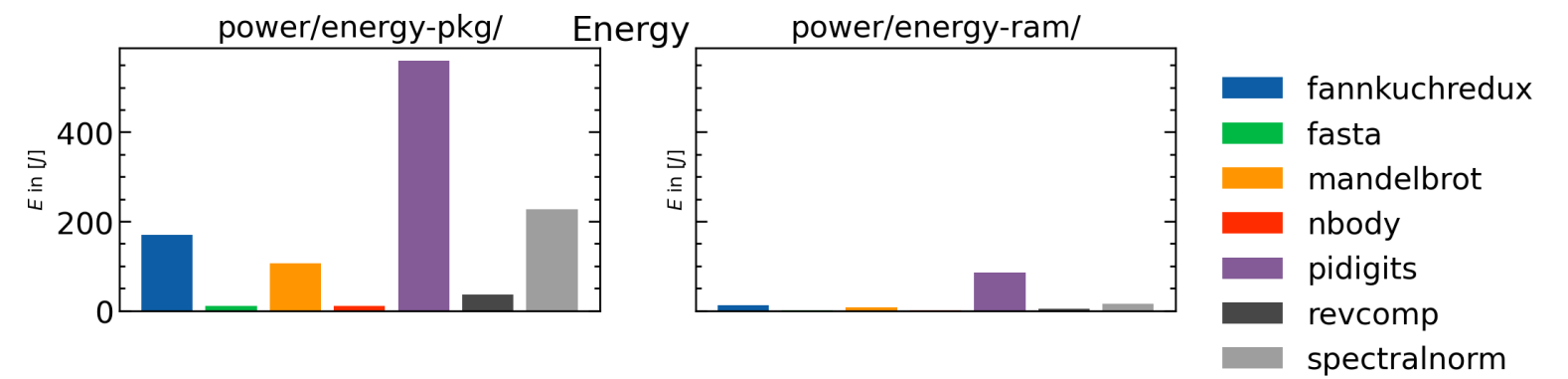
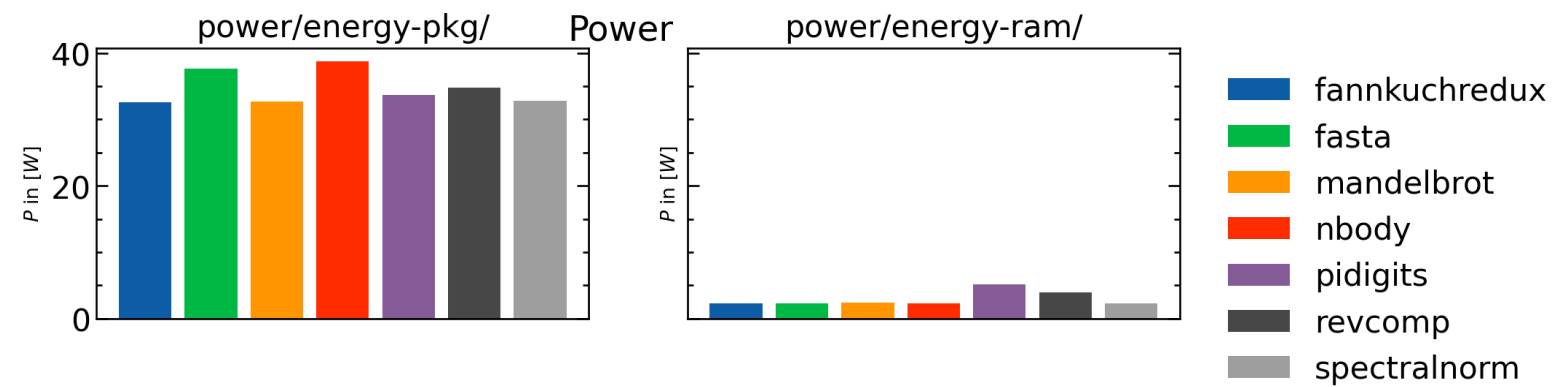


## Clang++ 16

# C/C++

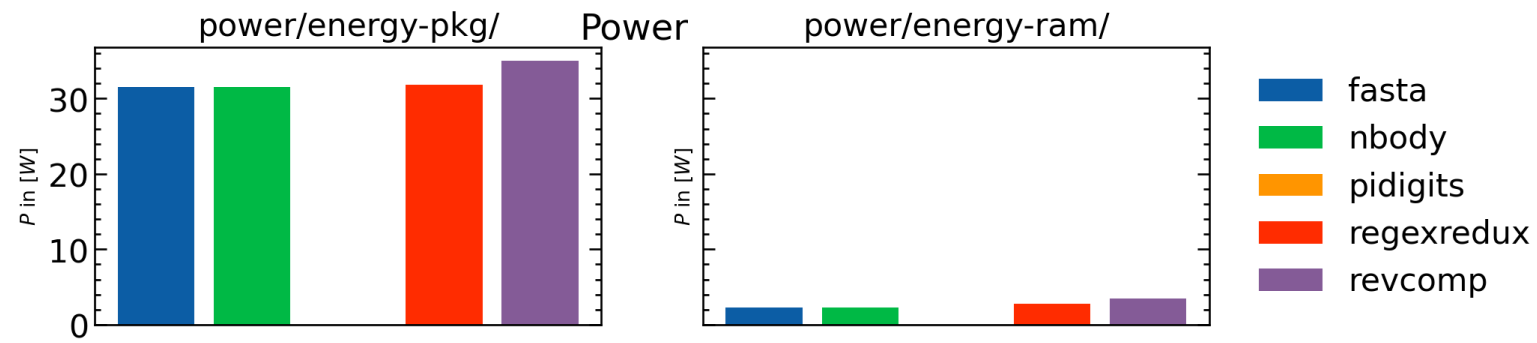


## GCC 12

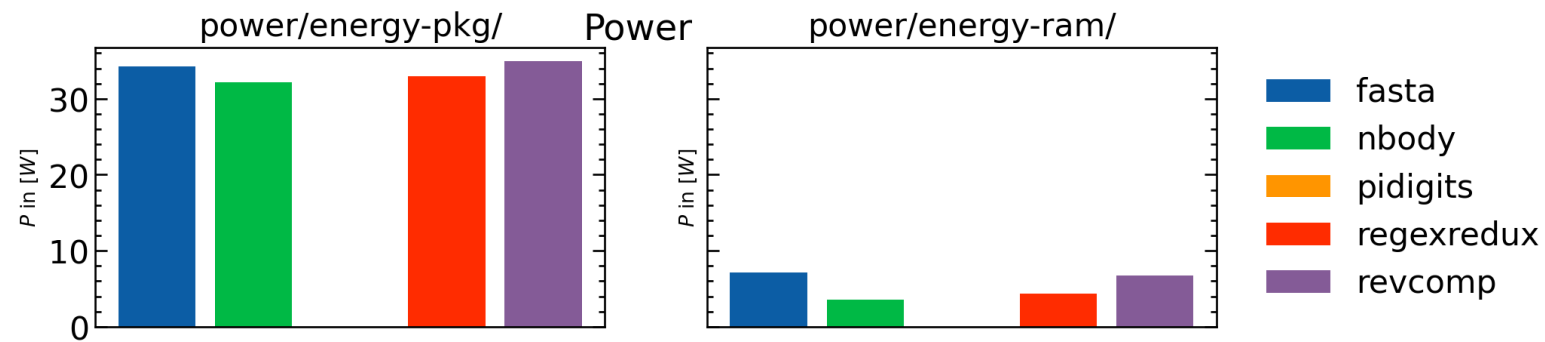


## G++ 12

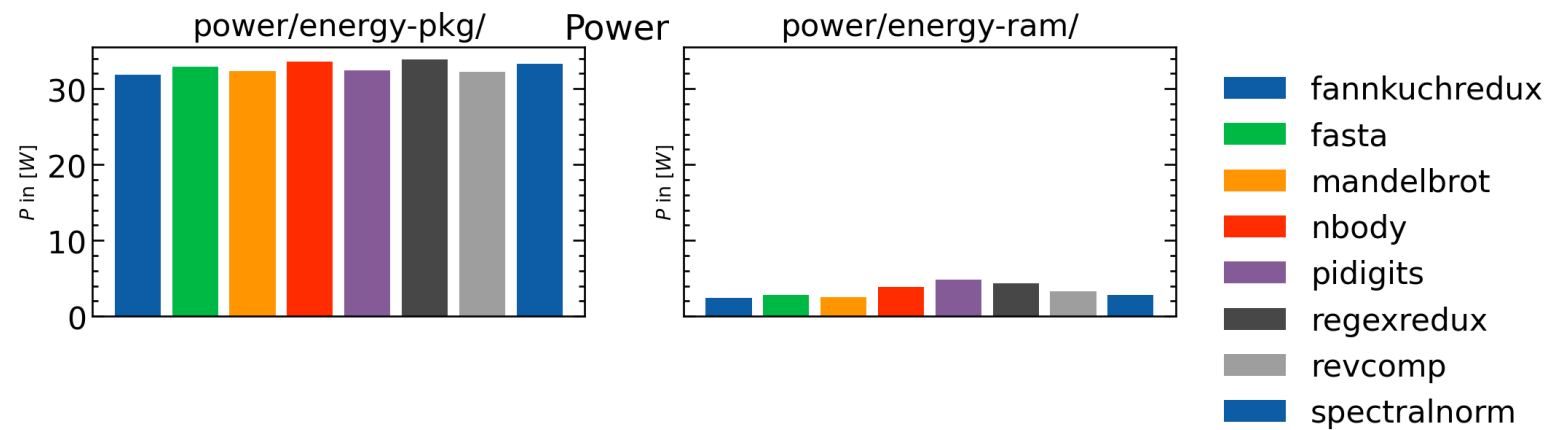
# Interpreted Languages



## Python 3.10



## PyPy 7.3.8



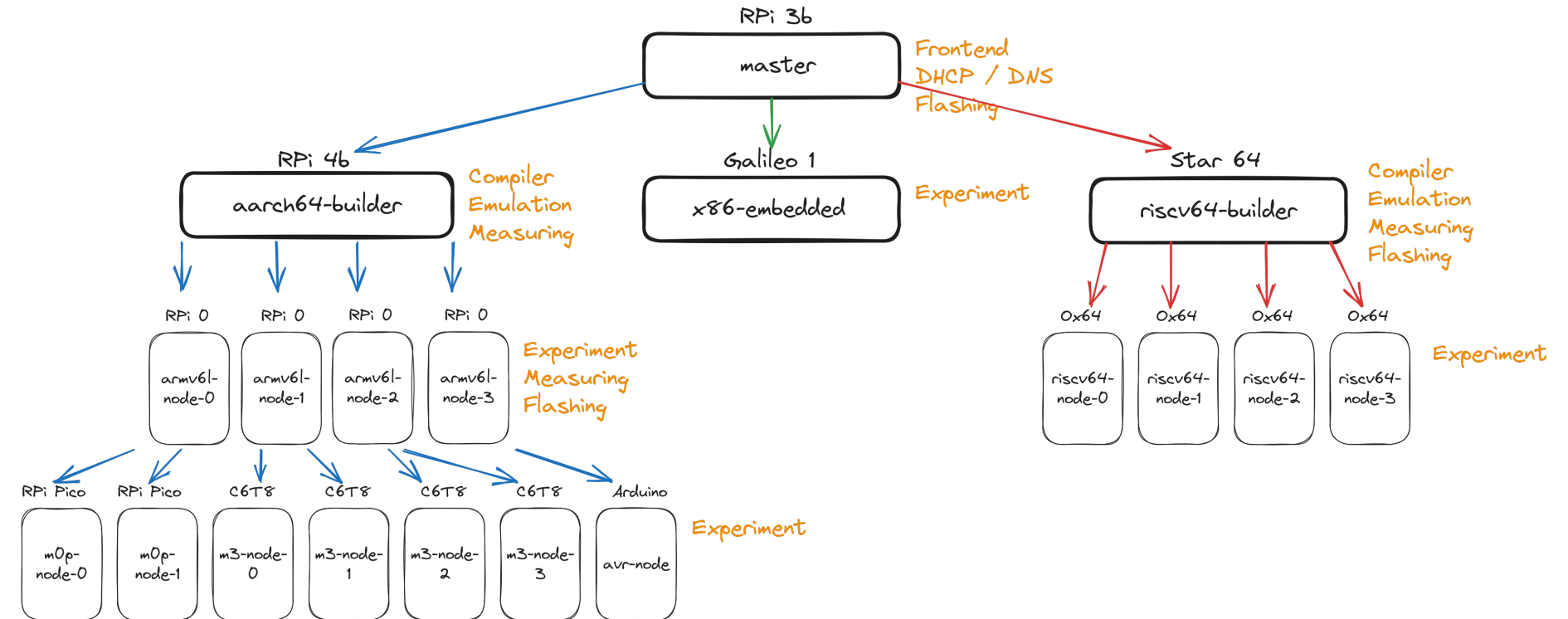
## Julia 1.9



# Programming Language Benchmark: Whats next?

## Computer Architecture

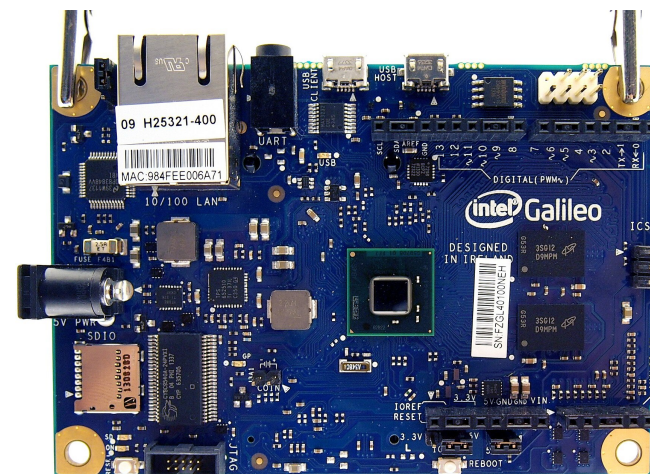
- compare energy consumption for different cpu architectures
- currently missing:
  - Slurm setup on master
  - small code to interact with current sensor via GPIO



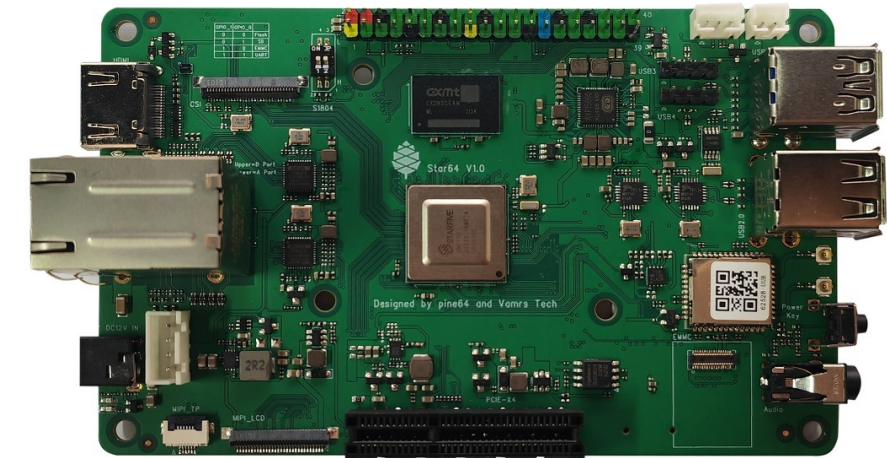
## Parallel Computing

- compare the energy efficiency scaling from single-core to multi-core to multi-nodes
- using matrix-multiplication
- currently implementing

Tree Diagram of a cluster that is used for several experiments. Different x86, ARM, RISC-V and microcontrollers are supported. The energy consumption is measured with external current sensors and is always system-wide.



Galileo is the x86 X1000 32-bit “microcontroller” from Intel (no SIMD like MMX or SSE). Comparable to i586.



Star64 is a RISC-V board based on JH7110. No vector extension support. Qemu is supports this sock.



# Software Optimization using Surrogate models

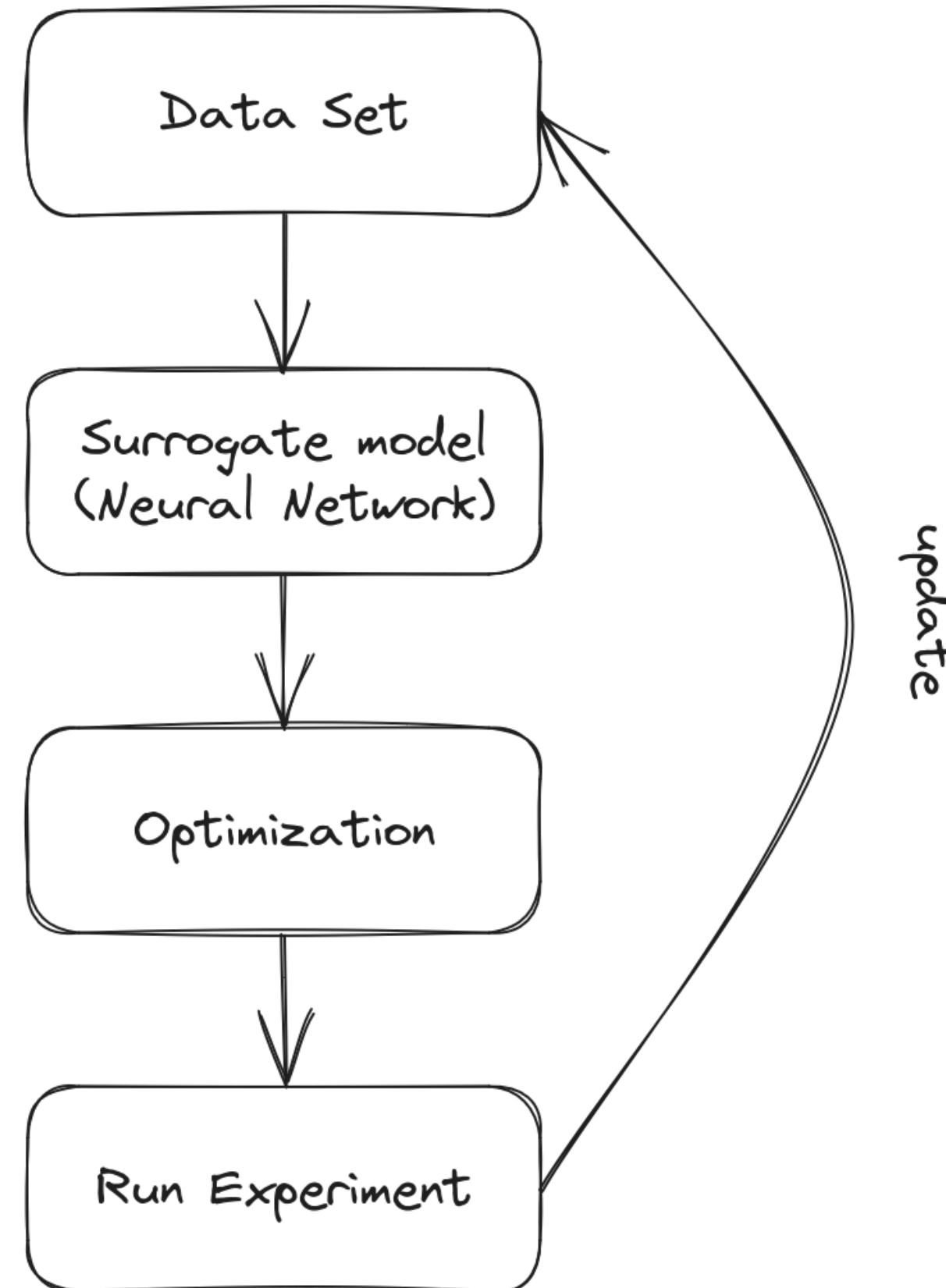
## Goal

- find the transformation pass pipeline that minimizes the energy consumption for different software types
- target IR instead of individual programming languages
- Approach detailed in: “Challenges in Automatic Software Optimization: the Energy Efficiency Case”, Fischbach, Kieffer, Bouvry (2023)

## Why Surrogate Optimization?

- time consuming evaluation
- blackbox (neural network based)
- model can be used to generate explanations

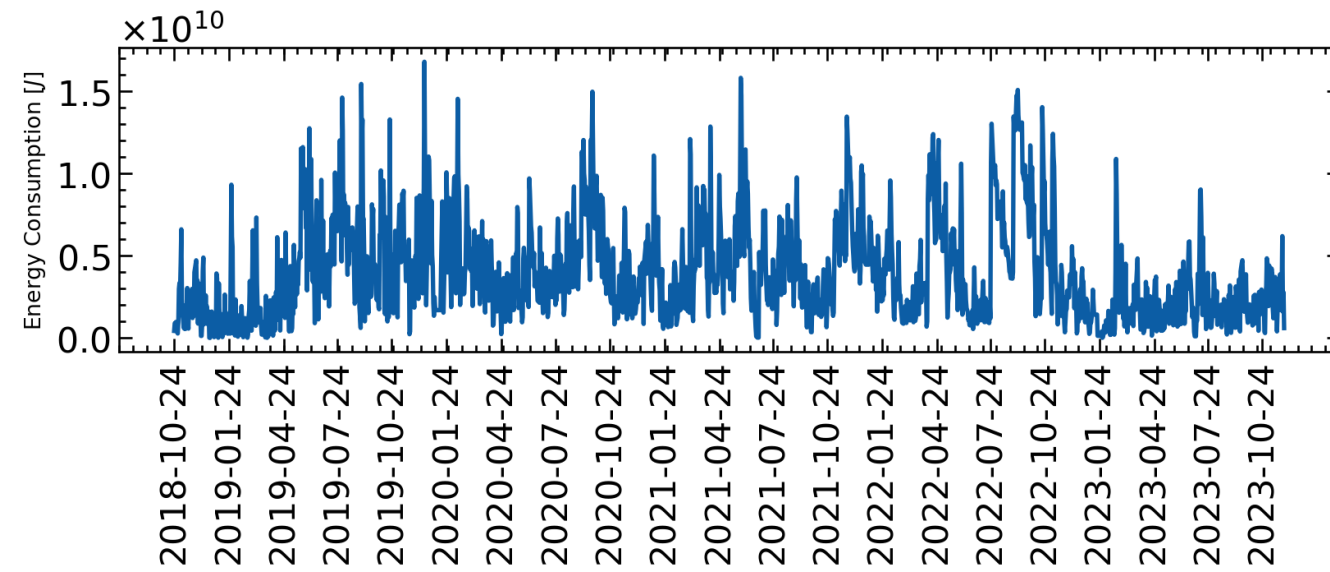
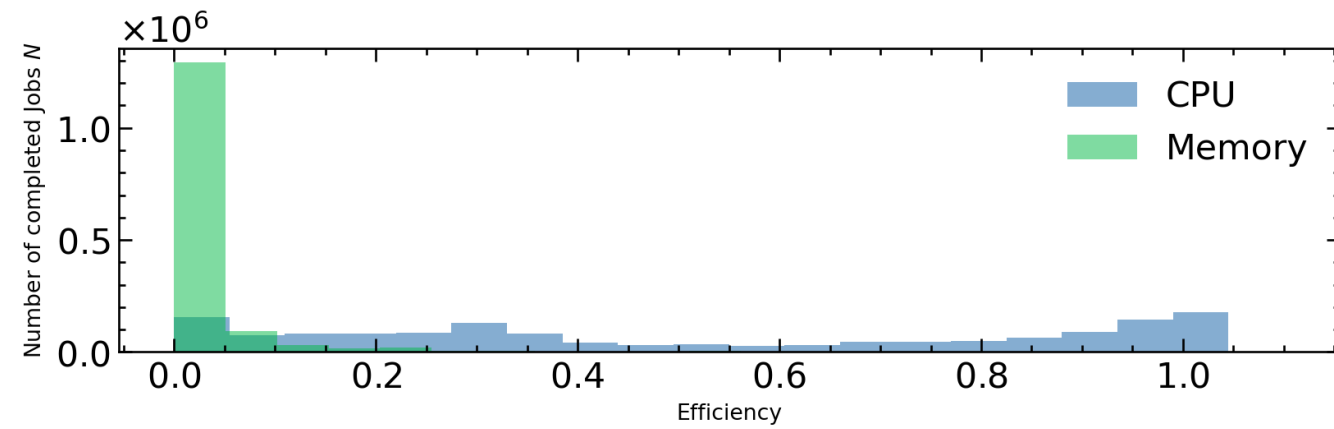
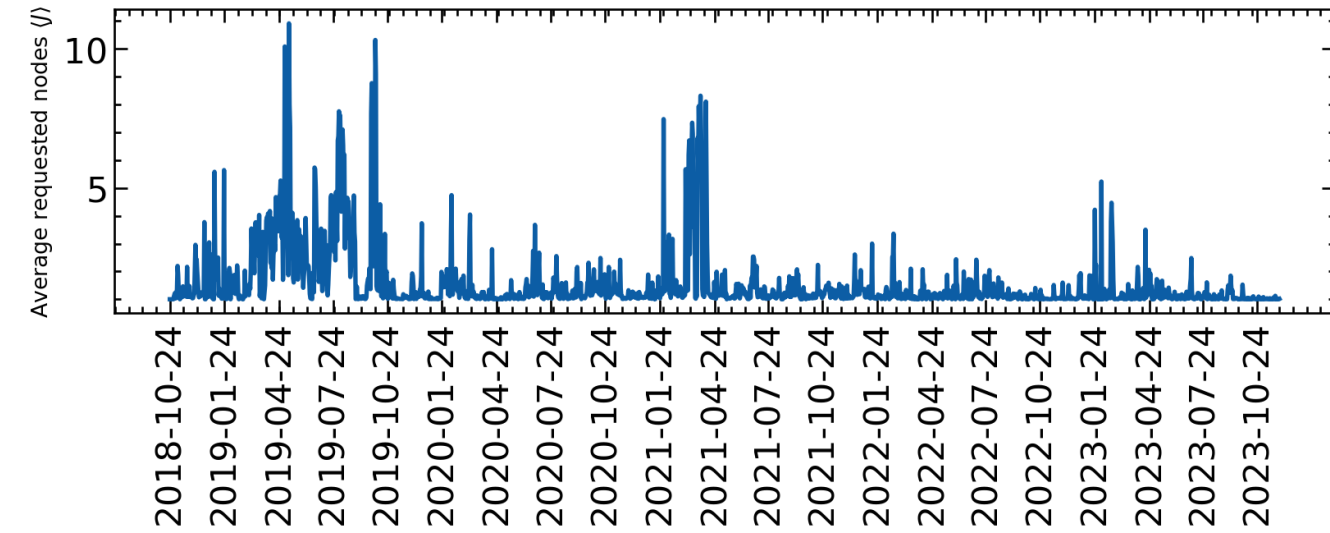
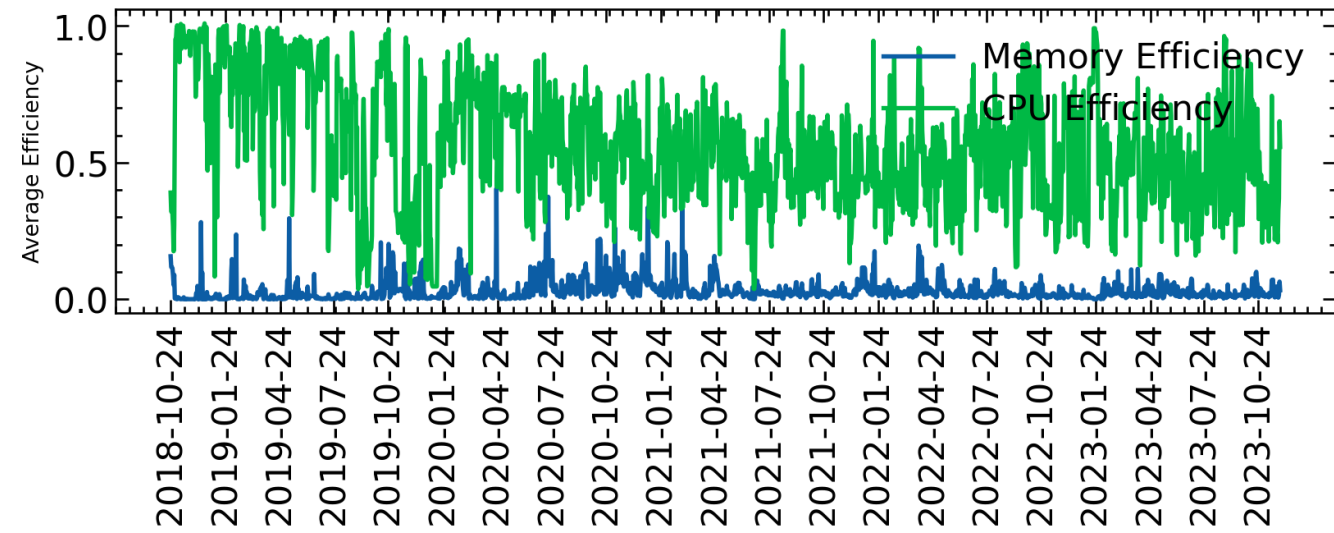
\(\rightarrow\) Next step is the integration with the experimental pipeline.





# Job Efficiency

# Iris



- CPU efficiency:  $(C_{\text{eff}} = \frac{t_{\text{TotalCPU}}}{t_{\text{Elapsed}} * N_{\text{requested}}})$
- Memory Efficiency  $(M_{\text{eff}} = \frac{\text{MaxRSS}}{\text{Mem}_{\text{requested}} * N_{\text{requested}}})$



# Outlook

# Summary

- HPC users are better in utilizing CPUs than memory
- most of the requested memory is unused
- modular framework to capture runtime dynamics and energy consumption
- introduction of programming language benchmark
- next surrogate optimization of LLVM pass-pipeline

Thank you for your attention



# Bibliography

1. [The Beating Heart of the World's First Exascale Supercomputer - IEEE Spectrum.](#)
2. [Green500 TOP500.](#)
3. Rountree, B. *et al.* Bounding energy consumption in large-scale MPI programs. in *Proceedings of the 2007 ACM/IEEE conference on Supercomputing* 1–9 (ACM, 2007). doi:[10.1145/1362622.1362688](#).
4. Mittal, S. & Vetter, J. S. [A Survey of Methods for Analyzing and Improving GPU Energy Efficiency.](#) *ACM Computing Surveys* **47**, 19:1–19:23 (2014).
5. Garg, S. K., Yeo, C. S., Anandasivam, A. & Buyya, R. [Energy-Efficient Scheduling of HPC Applications in Cloud Computing Environments.](#) *arXiv.org* (2009).
6. Kocot, B., Czarnul, P. & Proficz, J. [Energy-Aware Scheduling for High-Performance Computing Systems: A Survey.](#) *Energies* **16**, 890 (2023).
7. [Intel® 64 and IA-32 Architectures Software Developer Manuals.](#) *Intel*.
8. Desrochers, S., Paradis, C. & Weaver, V. M. A Validation of DRAM RAPL Power Measurements. in *Proceedings of the Second International Symposium on Memory Systems* 455–470 (Association for Computing Machinery, 2016). doi:[10.1145/2989081.2989088](#).
9. Zhang, H. & Hoffmann, H. A Quantitative Evaluation of the RAPL Power Control System.
10. [Which programming language is fastest? \(Benchmarks Game\).](#)
11. Pereira, R. *et al.* Energy efficiency across programming languages: How do energy, time, and memory relate? in *Proceedings of the 10th ACM SIGPLAN International Conference on Software Language Engineering* 256–267 (Association for Computing Machinery, 2017). doi:[10.1145/3136014.3136031](#).
12. Pereira, R. *et al.* [Ranking programming languages by energy efficiency.](#) *Science of Computer Programming* **205**, 102609 (2021).