

Heterogeneous HPC to the rescue? Ways to improve the energy efficiency of climate simulations today and tomorrow.

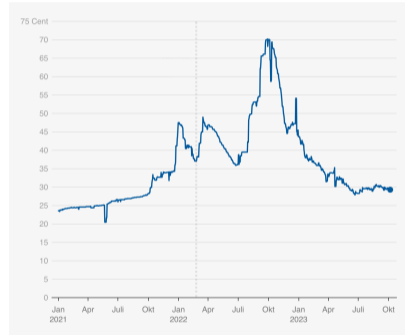
Jan Frederik Engels for the DKRZ Team



2023-10-11

Motivation

- Saving energy / using resources sustainably
- Mono-topical workload allows for good optimization
- Recent shifts in energy prices
- Increase in power efficiency between generations declines



Source: ndr.de / Verivox

Who?

Contributors:

- two BMBF projects (EECLIPS & EEHPC)
- Collaborators from various groups and departments in DKRZ
- Collaborators from the projects and elsewhere

DKRZ:

- 3MW, 3k Nodes, 130PB Storage, TOP 60
- AMD Milan & Nvidia A100
- Focus on German climate community, participation in EU-funded projects

Me:

- Applications and Services person
- Helping to put out fires where needed
- Leading the group working on EEHPC and EECLIPS

Outline

- ① The past: What has made a computing center energy-efficient?
- ② Today: What can we do with the existing machine to save energy?
- ③ The future: Can heterogeneity rescue us?

Key performance indicators for computing centers

$$\begin{aligned} \text{PUE} &= \text{Power usage efficiency} \\ &= \frac{\text{Total Energy}}{\text{IT-Equipment Energy}} \\ &\approx 1.1 \end{aligned}$$

$$\begin{aligned} \text{EUE} &= \text{Energy usage efficiency} \\ &= \frac{\text{Total Energy} - \text{Reused Energy}}{\text{IT-Equipment Energy}} \\ &\approx 0.8 \end{aligned}$$

How to optimize?

- Reduce usage of cooling machines
 - Warmwater cooling
 - Free cooling for cold water
- Heating of neighbouring building
- Hot aisle containment
- ...



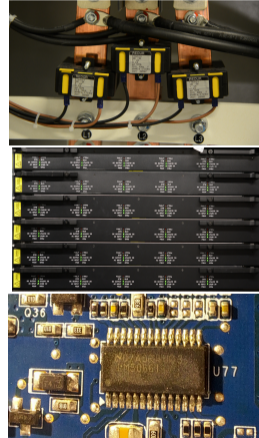
What does one aim for?

- Reducing power consumption?
- Reducing energy-to-solution?

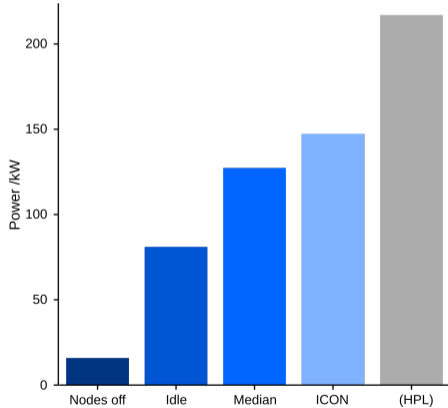


How much energy do we consume? And when?

- Measurement infrastructures
 - Building
 - Rack-PSUs
 - Node
 - Rapl / nvidia-smi
- Energy vs Power: $E = \int P(t)dT$

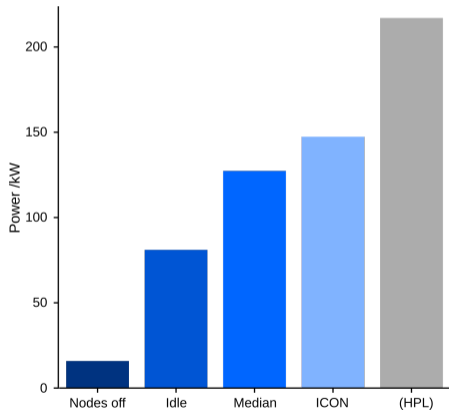


How much can we save?



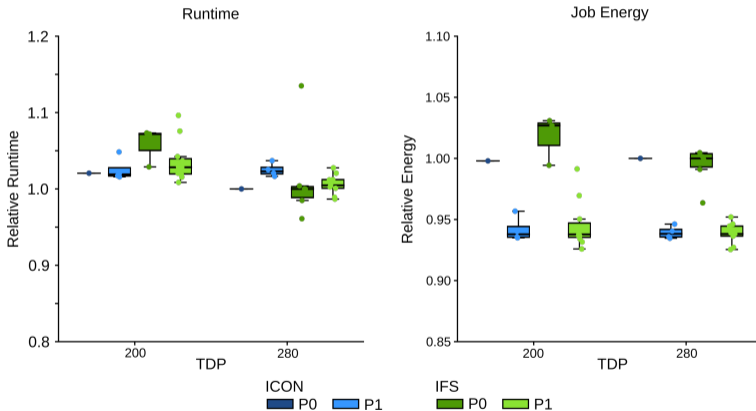
Nodes off	19%
Idle	100%
Median	158%
ICON	181%

How much can we save?

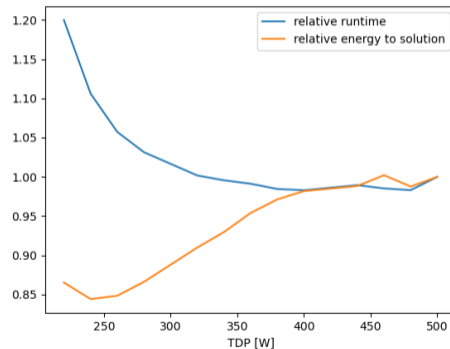
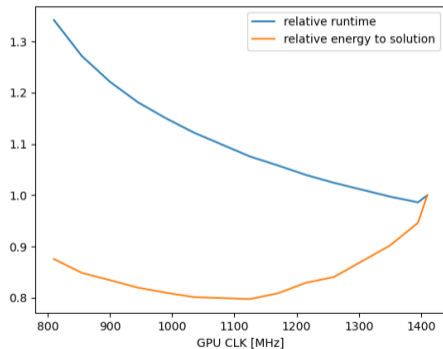


- Wiggle room for optimization is low
 - First order approximation:
Reducing runtime is saving energy
- ⇒ Trivial saving: Switch off idle nodes.

CPU knobs (AMD Milan)



GPU knobs (NVIDIA A100, preliminary)



Raising awareness

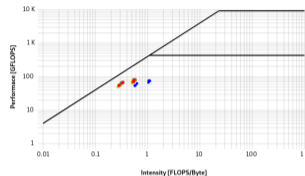
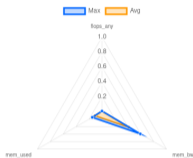
But what about inefficient resource usage?

[redacted] (levante)

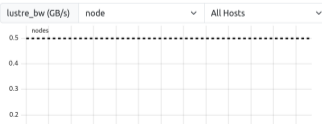
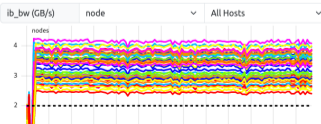
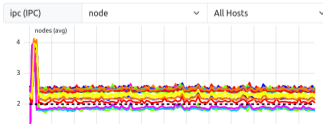
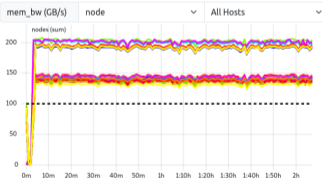
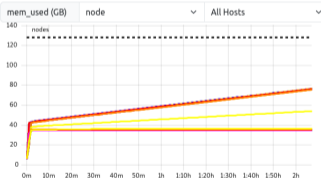
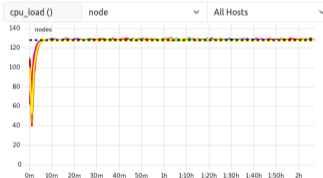
[redacted]
 [redacted]

48
compute

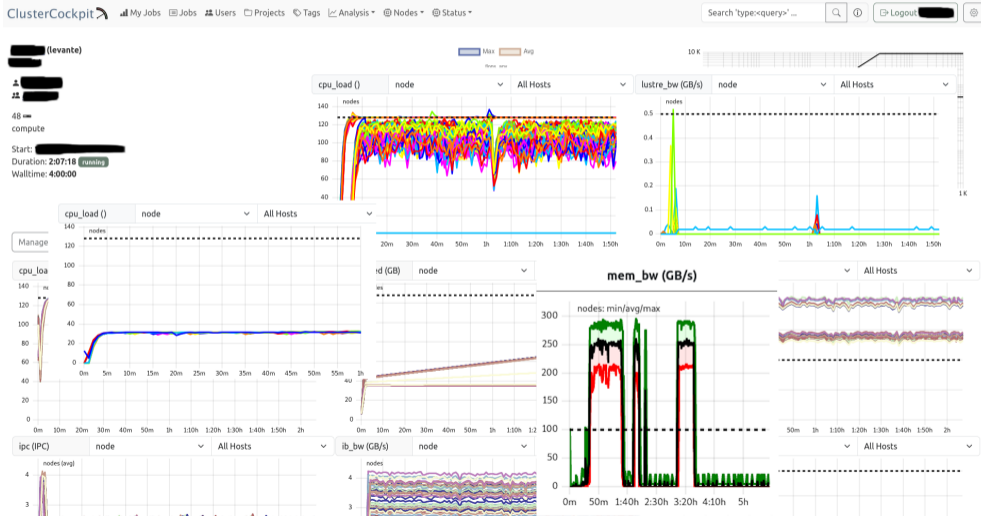
Start: [redacted]
Duration: 2:07:18 running
Walltime: 4:00:00



Manage Tags



Raising awareness



A coupled ICON simulation

	CPU 224 SDPD		GPU / CPU 231 SDPD	
	Nodes	Energy	Nodes	Energy
Atmosphere	130	37850 Wh	16	12000 Wh*
Ocean	16	4200 Wh	16	4250 Wh
Output	7	1150 Wh	7	1180 Wh

* *preliminary*

Heterogeneity beyond CPU-GPU-coupling?

Research questions:

- **Can we setup an ICON simulation where each component runs on its most energy-efficient architecture?**
- What is the most energy-efficient architecture for a component?
- How to build a test-cluster for this?
- How much energy can this approach save?



First step: Measuring (sub-) components

time per iteration	nh_solver	icefem_solver	ocean_vel_diffusion	...
SPR DDR	60.2	10.7	5.27	WIP
SPR HBM	32.8	10.2	2.55	WIP
Genoa	58.0	7.6	3.84	WIP
Aurora 1	86.4	WIP	WIP	WIP
A64FX		Compiler licence issues		
Grace Hopper		Awaiting delivery		
Grace Grace		Awaiting delivery		
(Levante A100-80)		Code parts missing / WIP		
(Aurora 3)		"2.4x faster than Aurora 1"		
(MI300?)		OpenACC support missing		

(preliminary)

First step: Measuring (sub-) components

energy per iteration	nh_solver	icefem_solver	ocean_vel_diffusion	...
SPR DDR	7.79	1.04	0.56	WIP
SPR HBM	4.3	1.04	0.37	WIP
Genoa	5.43	0.15	0.36	WIP
Aurora 1	3.13	WIP	WIP	WIP
A64FX		Compiler licence issues		
Grace Hopper		Awaiting delivery		
Grace Grace		Awaiting delivery		
(Levante A100-80)		Code parts missing / WIP		
(Aurora 3)		"40% less than Aurora 1"		
(MI300?)		OpenACC support missing		

(preliminary)

Next Steps

- Integrate the remaining architectures
- Define more (sub-)components and measure those
- Work out the optimal placement of components onto architectures
- Set up the ideal simulation

Ideas for components

- Atmosphere
- Atmospheric chemistry
- Radiation
- Ocean
- Ocean Biogeochemistry
- ...

Thanks!

A big thanks for contributing to the work behind this talk:

- Pay Giesselmann, Julius Plehn, Erik Pfister
- Stephan Jaure (Eviden)
- ... our projects partners

Disclosure: AMD and Intel both sponsored CPUs for the heterogeneous test-cluster.

Aliasing

