Smarter technology for all

From high performance computing to artificial intelligence - my journey

Nils Smeds | 2023-06-20

2023 Lenovo Internal. All rights reserved.

Smarter technology for all

Topics of my talk

Energy efficient HPC and AI – Why and how
 Tuning applications for (energy) efficience
 My journey in HPC (so far)



Smarter technology for al

Energy Efficie

Dr. Thomas Alrutz Senior Systems Architect High Performance Computing S NO NO

2022 Lenovo Internal. All rights reserved.

The Present -> KiT HOREKA

- 610 SD650V2 46.360 CPU Cores
- 167 SD650-NV2 668 A100 GPUs
- 3+18+24 PB Bulk storage on HDD
- 250 TB Flash storage on NVMe
- Total power consumption ~940 kW*
- Custom CDU DeltaT of 5K
- Heat Recovery to WW ~ 85%



KIT – HOREKA CPU

(119#, Nov 2021) Lenovo SD650 V2 w/ Xeon 8368 38C 2,4GHz 270W



- Rack: 72 Nodes, 5472 Cores
- SpecFP2017 Rate: 275.114
- HPCG: ~35 TFlops
- Rack Power: ~55kW
- Flow Rate: 60 l/min

KIT – HOREKA GPU

(53#, Jun 2021- Green #13) Lenovo SD650-N V2 w/ Xeon 8368 38C 2,4GHz 270W + 4 A100 SXM4



- Rack: 36 Nodes, 2736 CPU Cores / 144 A100 GPUs
- HPCG: ~181 TFlops
- RackPower: ~60 80kW
- FlowRate: 108 l/min

Lenovo, ThinkSystem SD650 V2 (SilverSurfer/Toomie)

Feature	ThinkSystem SD650 V2 / SD650-N V2		
Form factor	 1U full wide double (CPU only) / single (GPU acc.) node tray in 6U6T Chassis (DW612) for 19inch rack cabinets 		
Processors	NPI: 2x Intel "Ice Lake" up to 270W TDP		
Memory	 16x DDR4 3200 R/3DS DIMM ECC Capacities: 16GB/32GB/64GB/128GB 		
Storage	 Up to 2x 7mm or 1x 15mm U.2/SATA, no drive/bp choice 2x M.2 SATA SSD Intel VROC RAID / SW RAID only 		
NIC	 1x SFP28 25GbE LOM, NCSI (10Gb capable) 1x RJ45 1GbE LOM, NCSI 		
PCIe	 2x x16 PCIe Gen4 LP (each in place of half the storage) Internal: 3x (+1 from above) x16 Gen 4 for GPU Side Car 		
Acceleration	4x Nvidia A100 400W SXM4 GPUs with Nvlink 3.0 (-N)		
Front Access	 All IO in front Power LED Button, ID and System Health LED KVM breakout connector, Pong 		
Rear Access	 2x RJ45 on SMM for XCC/Daisy Chain USB 2.0 dedicated to SMM LED for ID, Error, Power and Heartbeat 		
Mgmt/TPM	xClarity Controller (XCC), TPM 2.0		
Power	 Up to 9x HS Air CFF v4 (1800W PT, 2400W PT) N+1 redundancy (only air-cooled / without acceleration on DWC) 		
Cooling	 Up to 50°C warm water for component level cooling Up to 85% cooling efficiency at 45°C inlet temperature 		

blue are changes over SD650











Thinksystem DW612S - Chassis for Next Gen Silicon

DWC High Heat Recovery

6/9 Air Cool PSUs Heat to Water 85%-92%*

DWC Extreme Heat Recovery

3 Direct Watercooled PSUs Heat to Water Recovery 95%-98%*



Why are we so obsessed with temperatures?

- T_{Junction} T_{Case} drives the heat from the IC interior to the environment
- Increasing TDP (Thermal Design Point) 200W...400W...700W...1kW(?)
- The case temperature must drop further to prevent damage to the CPU
- Heat capacity (per volume) of water is much higher than air
- The advantage of water will increase further
- Once heat is trapped in water keep it there!







Explaining power and the <u>cost</u> of cooling

Typical 700W 1U Server

The Data Center View





Our 30kW is really 50kW operationally



Lenovo's DWC Design Tenets?

High Heat Recovery	>90% at 45°C inlet temperature	
High Cost Effectiveness	As much common parts/design as possible.	
High Water Temperature	Up to 50°C inlet for efficient waste heat reuse.	PROFESSIONAL UV Leak DETECTOR KIT PROFESSIONAL UV Leak DETECTOR KIT Protector Part and Eav Protector Prote
High Material Reliability	Copper avoids leakage or microfractures.	
Environmental Safe	Use high efficiency water, no solvents/glycols	
High TDP *PU Support	Up to 700W TDP in the future – new challenges.	

Performance Benefits of Liquid Cooling

- Highest absolute performance
- Lower power efficiency out of the box
- Increased thermal efficiency with liquid cold plates
- Liquid cooled with Neptune
 - Stay well below throttling temperatures
 - No fan power required.
 - Provides reliable clock speeds





Best-in-class, drip-less quick connect shared between two nodes

Parallel loop ensuring <5°C ΔT between CPU1 and CPU2 for max perf.





Direct liquid cooling for standard 2.5' drives – including NVMe.

Direct liquid cooling for standard memory DIMMs – including Intel Optane.



ThinkSystem SD 650 V2



ThinkSystem SD650-N V2



Energy Aware Run time (EAR): Motivation

- Strom und Energieverbrauch ist mittlerweile eine kritische Größe für heutige HPC-Systeme
- Performance und Stromverbrauch von parallelen Applikationen hängt ab von:
 - Parametern der Architektur
 - Laufzeit der Knoten Konfiguration
 - Charakteristiken der Applikation
 - Eingangsdaten
- Manuelle "beste" Frequenz
 - Schwer manuell festzulegen und ein zeitaufwendiger Prozess (Ressourcen und dann Strom) und nicht wiederverwendbar
 - Kann sich im Laufe der Zeit ändern
 - Kann über Knoten hinweg unterschiedlich sein



Benchmarking and tuning



- Wall clock time is often reported by the application
- nvidia-smi dmon can report various metrics (not well documented, and only one of the usages described worked for me)
- Batch system tools may provide data such as total run-time and energy consumption
- nmon can produce system wide data in csv format (or InnoDB/Grafana)
- Profilers are great if you can isolate a small enough section of critical code

What can be controlled by the user?

- Very site dependent!
- Memory speed, CPU governors/frequencies
- No of GPUs used, power envelope
- Subdivision of GPU (eg GPU split/vGPU), multiple processes/MPI tasks per GPU
- CPU utilization (processes/threads/HT)

- Important to find out early
- Time to solution paper publication?
- Using my project resource efficiently?
- Maximizing completed runs vs. system TCO?
- Or just the beauty of knowing what is going on inside?



Optimizing energy efficiency



- Air cooled server, 8×NVidia A100 PCIe 40GiB
- Tensorflow benchmark, 4 GPUs used (2 on each socket)

My HPC journey



Introduction to computing









No real urge for anything particular

- Naturwissenschaft Gymnasium 1980-83
- USA College 1 yr. Comp Sci + other things
 - Found out I liked assistant teaching and computers
- Military service (obligatory at that time)
- Gap year working in radio communications
- KTH, Stockholm, Engineering Physics

Still not really any thing I really had a wish to become

- Extra teacher
 - My old gymnasium, at a prison, computer science, numerical methods and solid mechanics
- Did my master thesis in solid mechanics
- Started at the PhD program in solid mechanics (my biggest mistake)
- Didn't quit (second mistake)
- A second national supercomputer center started across the street at KTH (started to fix my mistakes)

- Don't be worried if you have not found "your thing"
- Try to be open to new openings
- Don't listen to closely to career advise
 - Most people over estimate their own influence on their success.
- Be kind, give credit to others, offer to help
 - But be aware that you will be used by some people
- Your closest manager is the most critical to your happiness at work
 - Your fellow co-workers come next in importance

Smarter technology for all

2022 Lenovo Internal. All rights reserved.