

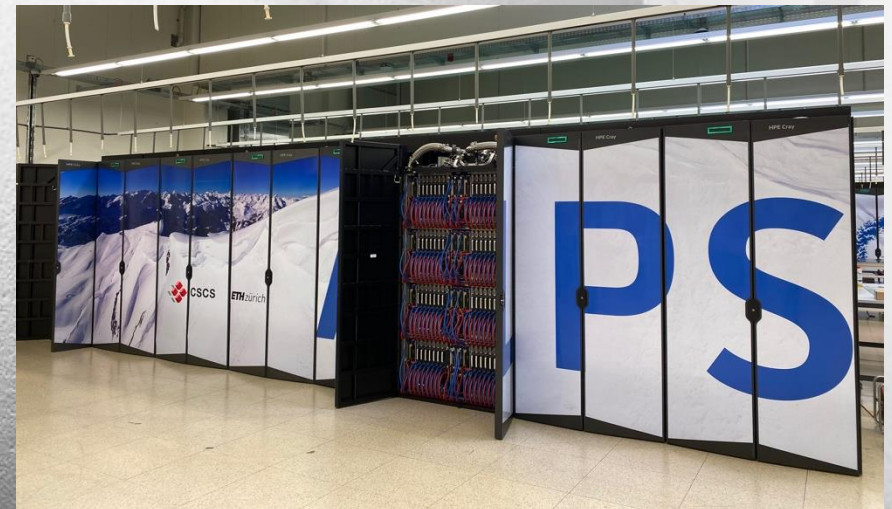


**Hewlett Packard
Enterprise**

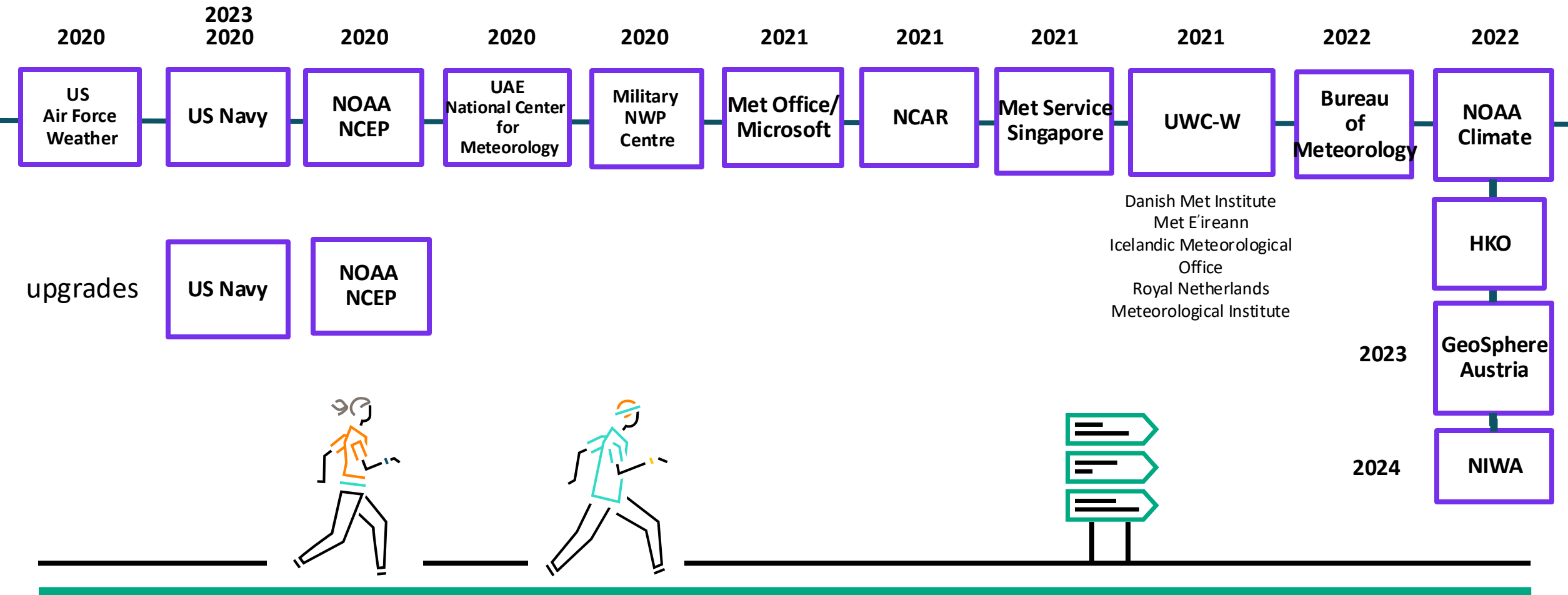
Future Technologies and NWP

Presented by: Edward Habjan

**Based on Details from Ilene Carpenter
Earth Sciences Segment Manager
ilene.carpenter@hpe.com**



New HPE Cray Supercomputers for WEATHER, CLIMATE and OCEANOGRAPHY



HPE Leadership in HPC systems for Weather Prediction

Supporting Operational NWP

AFW HPC11 at ORNL

- AMD CPUs
- NVIDIA A100 GPUs
- ClusterStor
- Slingshot interconnect

Bureau of Meteorology

- Intel SPR
- Slingshot interconnect

Met Office in Azure

- AMD CPUs
- ClusterStor
- Slingshot interconnect

NOAA WCOSS

- (Cactus and Dogwood)
- AMD CPUs
 - ClusterStor
 - Slingshot interconnect

UAE NCMS

- AMD CPUs
- NVIDIA A100 GPUs
- ClusterStor
- Slingshot interconnect

NIWA

- HPE, NVIDIA
- VAST Data

UWC-W

- AMD CPUs
- ClusterStor
- Slingshot interconnect

Met Service Singapore

- AMD CPUs
- ClusterStor
- Slingshot interconnect

US Navy

- AMD Genoa CPUs
- 128 MI300A APUs
- 24 NVIDIA L40 GPUs
- ClusterStor
- Slingshot interconnect

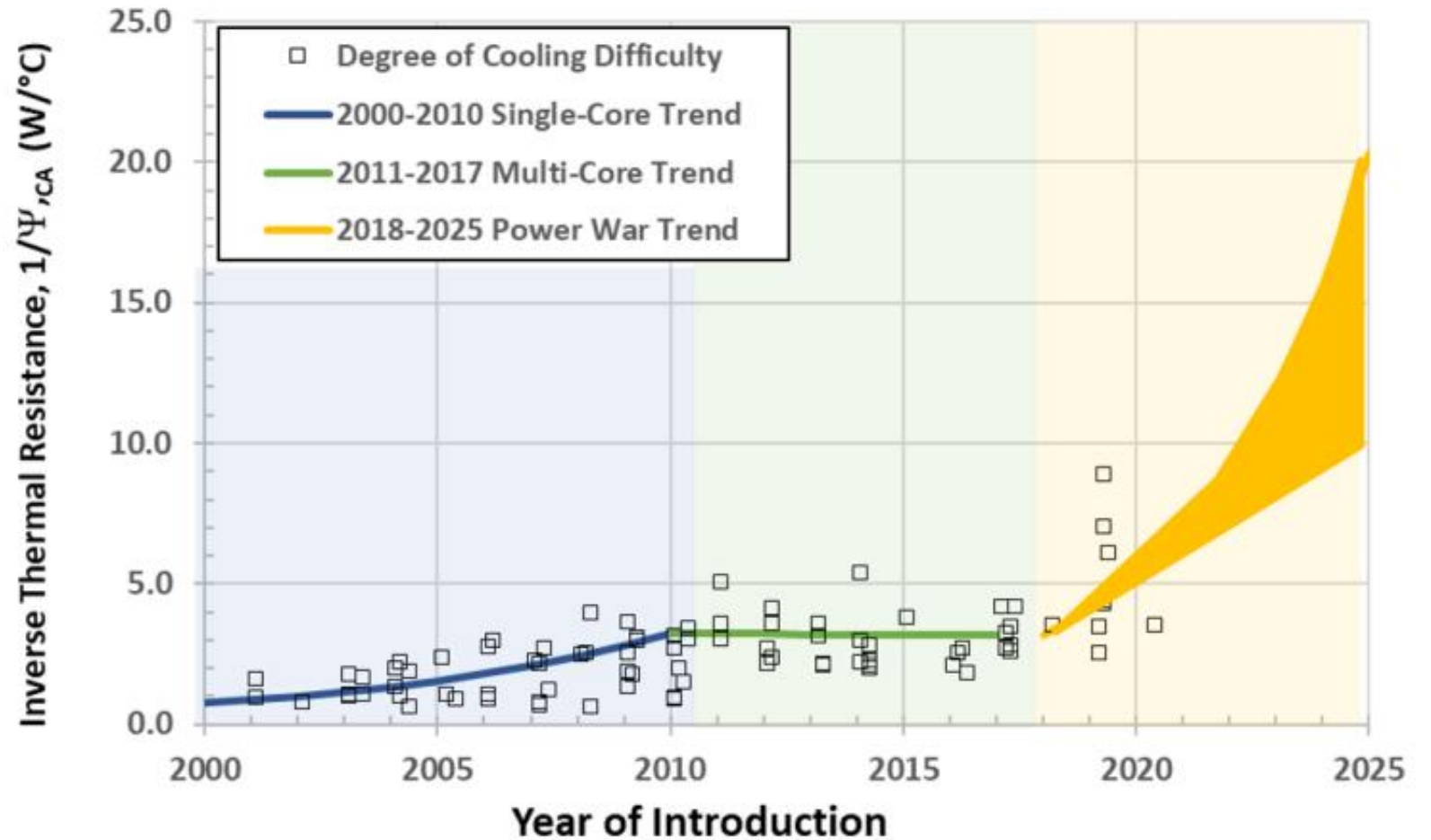
***Anyone can build a fast CPU.
The trick is to build a fast system.***

Seymour Cray



Changing Chip Designs

**COOLING MATTERS
MORE THAN EVER**



ASHRAE Technical Committee (TC) 9.9, Mission Critical Facilities, Data Centers, Technology Spaces, and Electronic Equipment White Paper: Emergence and Expansion of Liquid Cooling in Mainstream Data Centers

CPU thermals are dramatically increasing

Increasing power levels will stretch the limits of existing datacenters, and require long lead time to retrofit Power & Cooling

HPC/AI customers will need to need to Sacrifice Performance or Adapt

This is not even considering the uptake of GPUs !

Many Datacentres are perpetually behind

2013 - Intel Haswell/Broadwell
Top Bin = **145W**
~30kW - 72 Nodes/rack w/rich config

2015 - Intel Skylake/Cascade Lake
Top Bin = **205W**
~42kW - 72 Nodes/rack w/rich config

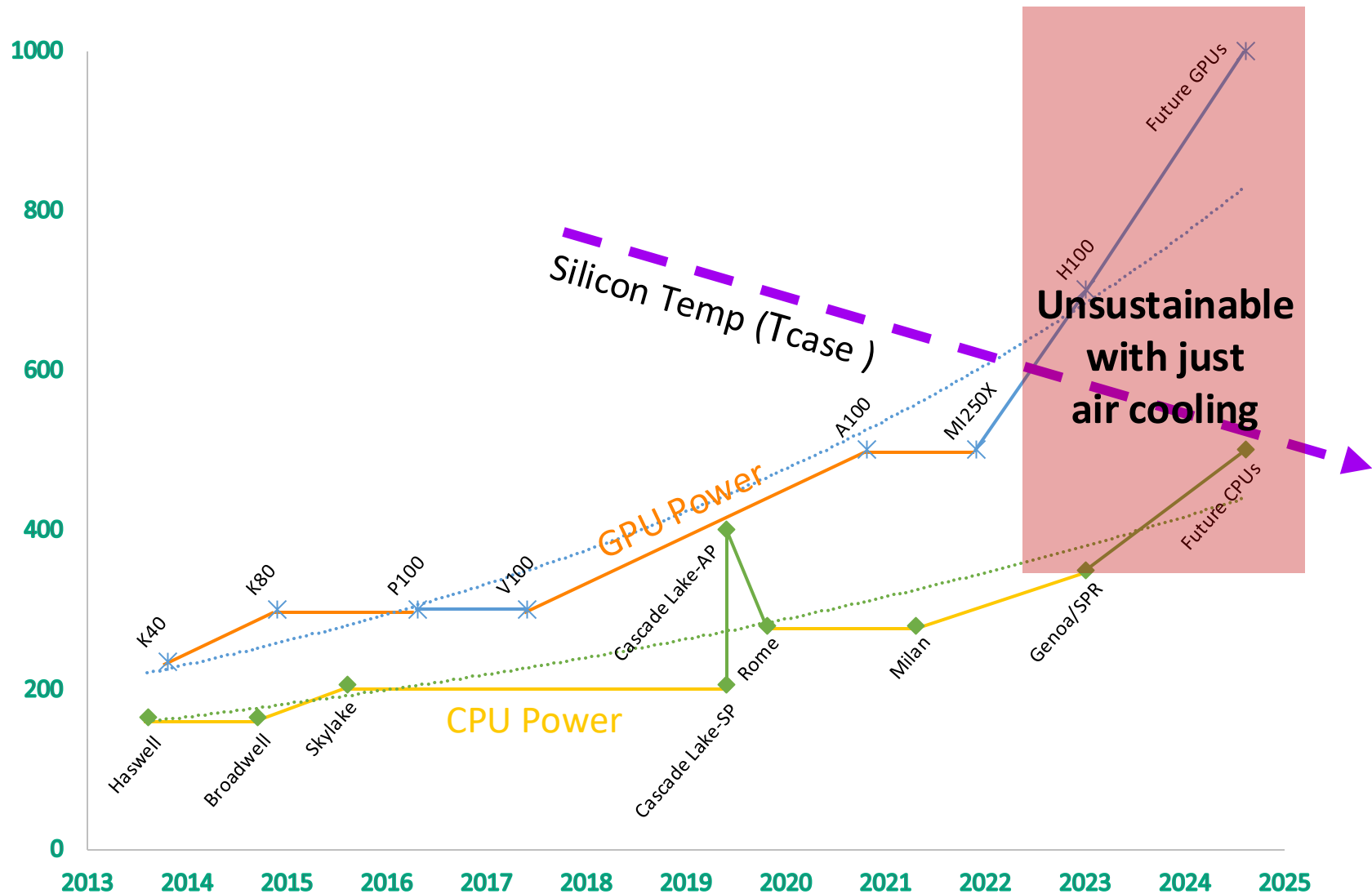
2019 - AMD Rome/Milan
Top Bin = **280W**
>50kW - 72 Nodes/rack of HPC

2023 - AMD Genoa / Intel Sapphire Rapids
Top Bin = **350W** or greater
>60kW - 64 Node/rack w/Slingshot

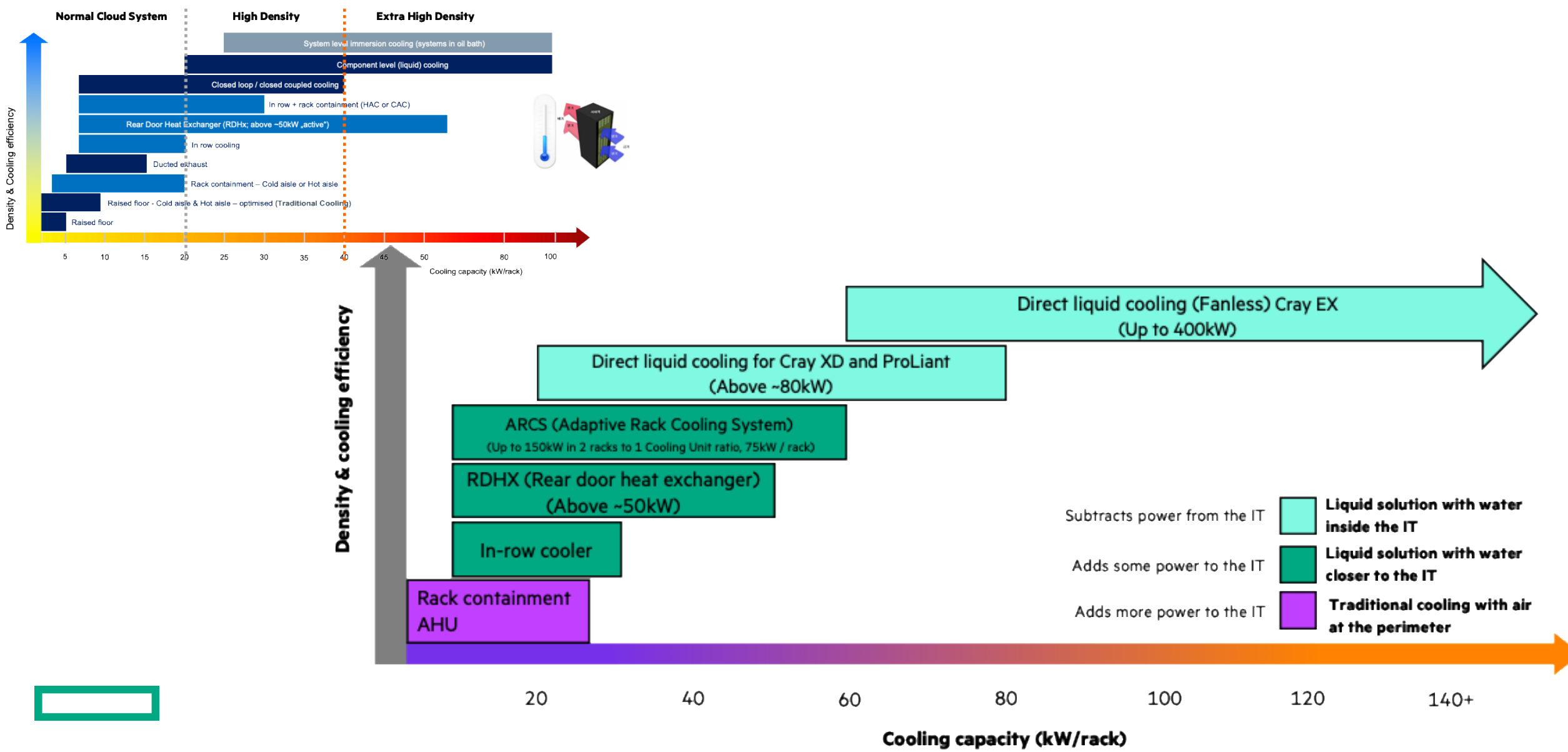
Many customers still stuck with **10-15kW**/rack datacenters, some have upgraded to **25-40kW**/rack

30-40kW/rack datacenters will struggle with air cooling

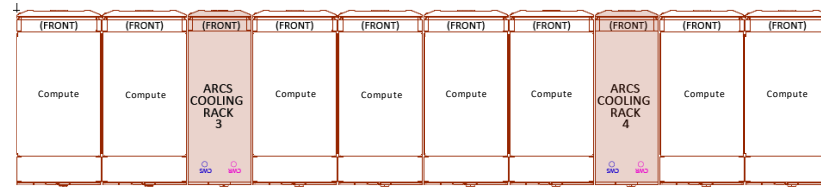
Processor Temperature – Cooling Dilemma



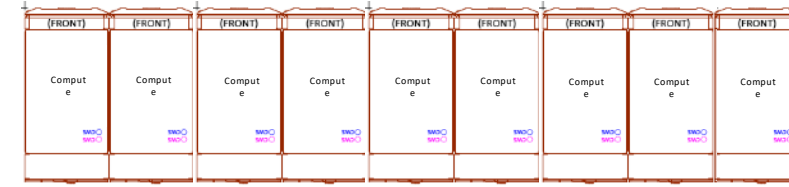
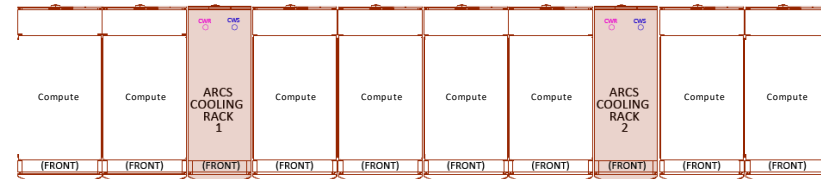
HPE Cooling Solutions



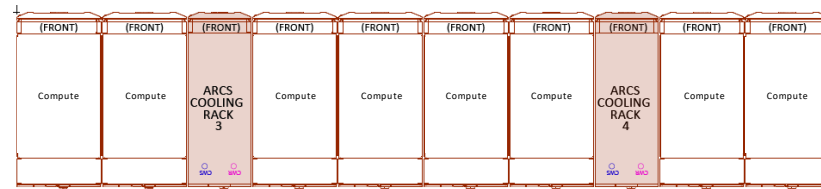
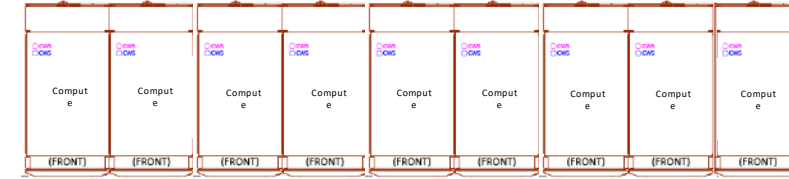
What difference does Packaging make – 4.26m x 7.9m comparison



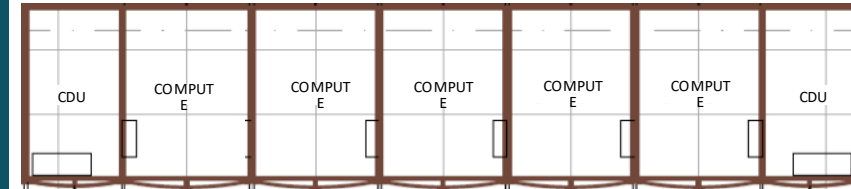
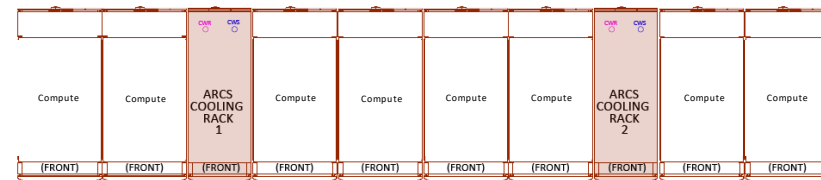
DL360 with ARCS 512 nodes



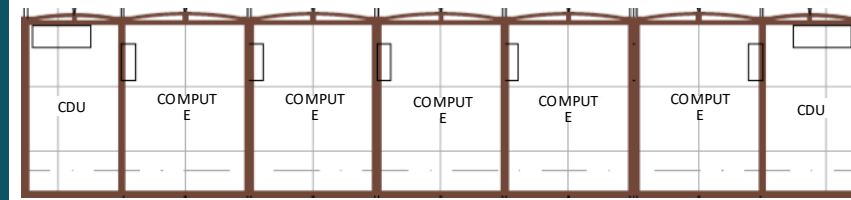
EX2500 1728 nodes



XD2000 with ARCS 1024 nodes



EX4000 2560 nodes



Trends and conditions affecting HPC for NWP

Processor diversity

- Many CPU choices, **lots of SKUs**
- Huge range of core counts
- dense/classic, P-cores/E-cores

GPU diversity

- Architectures
 - AI-optimized
 - FP64-optimized
- Software stacks

Workload evolution

- Rapid increase in AI/ML work
- Explosion in full forecast system emulators, new neuralGCMs

Performance/socket
driven by power (TDP)

Costs & supply chain
Performance/\$?

Will emulators replace
GPU-accelerated
loops?

- Lots of memory channels
- Lots of cores
- cache/core

Performance
Portability for Fortran-
based NWP and ES
models on GPUs

- OpenACC challenges
- DSLs

Big Uncertainties
Rapid change
Lots of choices



Trends and conditions affecting HPC for NWP - CPUs

Processor diversity

- Many CPU choices, lots of SKUs
- Huge range of core counts
- “dense” vs. “classic”, “P-cores vs. E-cores”
- X86 and Arm

Performance/socket driven by power (TDP)

- Lots of memory channels
- Lots of cores
- cache/core

As number of cores/socket increases, the number of choices increases:

Processors in 2025 will range from 8 to 192 cores with many variants

- Memory system differences
- Clock speed differences
- Impossible benchmark all options

Use performance modeling to understand application characteristics and trim choices.

NWP and ESM codes

- scale well & can use more cores to deliver performance
- need memory bandwidth and cache, not clock speed
- have no license costs (can use more cores)

How much memory/core do you need?

- Beware searching for the optimal SKU for a fixed workload
 - Complex optimization problem
 - Costs are unknown until last minute
 - Lifespan of optimal solution can be short – workload changes
- Aim for generally good

Beware of price and power increases in next-generation systems!

AMD Genoa – LOTS
of choices
(Turin has even
more!)

Name	# of CPU Cores	# of Threads	Max. Boost Clock ¹	Base Clock	L3 Cache	Default TDP
AMD EPYC™ 9754S	128	128	Up to 3.1 GHz	2.25 GHz	256 MB	360W
AMD EPYC™ 9754	128	256	Up to 3.1 GHz	2.25 GHz	256 MB	360W
AMD EPYC™ 9734	112	224	Up to 3 GHz	2.2 GHz	256 MB	340W
AMD EPYC™ 9684X	96	192	Up to 3.7 GHz	2.55 GHz	1152 MB	400W
AMD EPYC™ 9654P	96	192	Up to 3.7 GHz	2.4 GHz	384 MB	360W
AMD EPYC™ 9654	96	192	Up to 3.7 GHz	2.4 GHz	384 MB	360W
AMD EPYC™ 9634	84	168	Up to 3.7 GHz	2.25 GHz	384 MB	290W
AMD EPYC™ 9554P	64	128	Up to 3.75 GHz	3.1 GHz	256 MB	360W
AMD EPYC™ 9554	64	128	Up to 3.75 GHz	3.1 GHz	256 MB	360W
AMD EPYC™ 9534	64	128	Up to 3.7 GHz	2.45 GHz	256 MB	280W
AMD EPYC™ 9474F	48	96	Up to 4.1 GHz	3.6 GHz	256 MB	360W
AMD EPYC™ 9454P	48	96	Up to 3.8 GHz	2.75 GHz	256 MB	290W
AMD EPYC™ 9454	48	96	Up to 3.8 GHz	2.75 GHz	256 MB	290W
AMD EPYC™ 9384X	32	64	Up to 3.9 GHz	3.1 GHz	768 MB	320W
AMD EPYC™ 9374F	32	64	Up to 4.3 GHz	3.85 GHz	256 MB	320W
AMD EPYC™ 9354P	32	64	Up to 3.8 GHz	3.25 GHz	256 MB	280W
AMD EPYC™ 9354	32	64	Up to 3.8 GHz	3.25 GHz	256 MB	280W
AMD EPYC™ 9334	32	64	Up to 3.9 GHz	2.7 GHz	128 MB	210W
AMD EPYC™ 9274F	24	48	Up to 4.3 GHz	4.05 GHz	256 MB	320W
AMD EPYC™ 9254	24	48	Up to 4.15 GHz	2.9 GHz	128 MB	200W

Name	# of CPU Cores	# of Threads	Max. Boost Clock ¹	Base Clock	L3 Cache	Default TDP ¹
AMD EPYC™ 9684X	96	192	Up to 3.7 GHz	2.55 GHz	1152 MB	400W
AMD EPYC™ 9384X	32	64	Up to 3.9 GHz	3.1 GHz	768 MB	320W
AMD EPYC™ 9184X	16	32	Up to 4.2 GHz	3.55 GHz	768 MB	320W



HPE Cray EX Supercomputer Blades – Direct Liquid Cooling

x86 Compute Blade
Cray EX4252 (Genoa)



Accelerated Blade
Cray EX254n (GH200)



HPE Slingshot
Interconnect



x86 Compute Blades
Cray EX420 (Sapphire Rapids)

Accelerated Blade
Cray EX255a (MI300A)



GPUs for LLMs is changing the landscape of HPC for Scientific Computing

Impact of HUGE commercial purchases

- Cost of data center space
- Availability of power for data centers
- Cost of GPUs
- Lead time for GPUs

Elon Musk wants to purchase 300,000 Blackwell B200 Nvidia AI GPUs — Hardware upgrades to improve X's Grok AI bot

News

By Aaron Klotz published June 3, 2024

Nvidia's B200 GPUs will be used to boost the X platform's AI capabilities

Zuckerberg's Meta Is Spending Billions to Buy 350,000 Nvidia H100 GPUs

In total, Meta will have the compute power equivalent to 600,000 Nvidia H100 GPUs to help it develop next-generation AI, says CEO Mark Zuckerberg.



By [Michael Kan](#) January 18, 2024



AI

Microsoft has a target to amass 1.8 million AI chips by the end of the year, internal document shows

[Ashley Stewart](#) Apr 17, 2024, 4:54 PM MDT

[Share](#) | [Save](#)

Trends and conditions affecting HPC for NWP and ESM - GPUs

GPU diversity

- Architectures
 - AI-optimized
 - FP64-optimized?
- Software stacks

Performance

Portability for Fortran-based NWP and ESMs on GPUs

- OpenACC
- DSLs

Rapidly changing roadmaps for GPUs, driven by AI investments

- Designed for transformer models
- What GPU/CPU ratio?
- How much HBM do you need?
- Costs and TDPs are eye-watering!
 - Performance/\$?
 - Performance/Watt?

Coupled ESMs have some components that aren't likely to be efficient on GPUs

- need enough CPU resources to run them
- should they share node? If so, how much memory do they need?
- Will they run optimally on nodes designed for LLM training?



Trends and conditions affecting HPC for NWP - ML

Rapid Workload evolution driven by advances in ML methods

- A few years ago
 - Emulators trained offline for specific parameterization
 - Post-processing/bias correction
- 2022-2023
 - Explosion in full forecast system emulators
 - Trained on ERA5 data
 - Inference starting from the analysis – extremely fast
- 2024 – stochastic models, diffusion models
 - neuralGCMs with dynamics solver combined with ML physics with online physics training
 - Mesoscale models
 - ML for DA &/or inferencing directly from observations
 - Foundation models

Collaborations between meteorologists at Met services with companies like Google, Microsoft, NVIDIA

Increasing value of high temporal and spatial resolution authoritative data from weather services.

Increasing asks for AI-optimized storage in addition to traditional PFS



Hybrid Cloud and agility for AI

On-prem/dedicated single-tenant

Best match for well-understood, stable applications and workflows

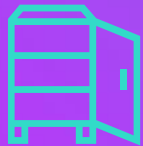
- Optimize the architecture and configuration for apps
 - Traditional (lengthy) procurement processes, benchmarks
- Include some generally applicable AI technology for in-house development and exploration, especially for training with large data produced by your own models and online training

Cost-savings + full control

Multi-tenant commercial cloud

- Low commitment
- Quick access to the most diverse set of architectures and configurations
- Enables exploration of AI methods using latest technologies
- Expensive to store and/or egress large volumes of data
 - Scientists like to save everything “just in case”
 - Massive, potentially ephemeral data -> reduce in time & space -> expand to create derived variables
 - Only curated open datasets get cheap/subsidized storage rates

Procure an optimized hybrid cloud solution



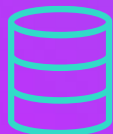
Mod/Sim



AI



Private cloud



Open Source vs. Proprietary Software

Open Science Requirements

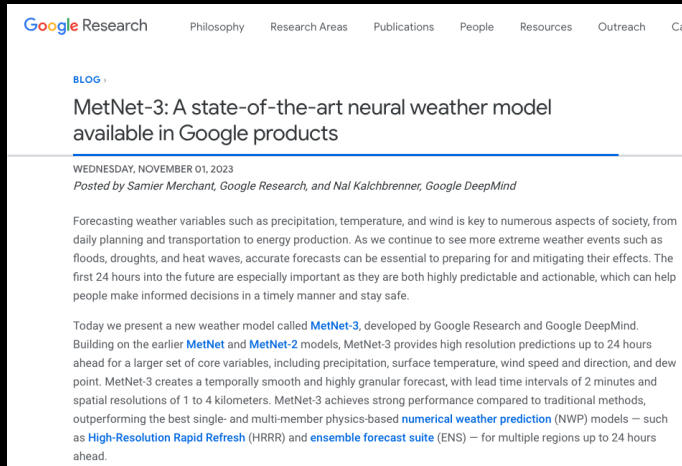
- Journals have moved towards demanding that results are reproducible (statistical, not bit-wise)
 - Suggests proprietary components must have an open-source alternative or be freely distributable
- Major AI frameworks are already open-source with highly permissive licences
- Some popular software moving to closed licenses (MongoDB, Redis, etc.), historically followed by development of open source alternatives

Collaboration Requirements

- Containerized applications for multiple platforms
- Do compilers need to be multi-platform?

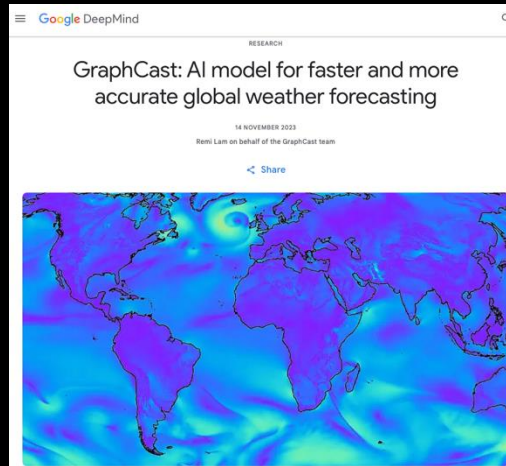
Proprietary for Performance, Open-source for sharing

Global Weather Forecast - GOOGLE



MetNet-3

- Regional forecast
- Traditional neural network
- High resolution and accuracy
- Operational



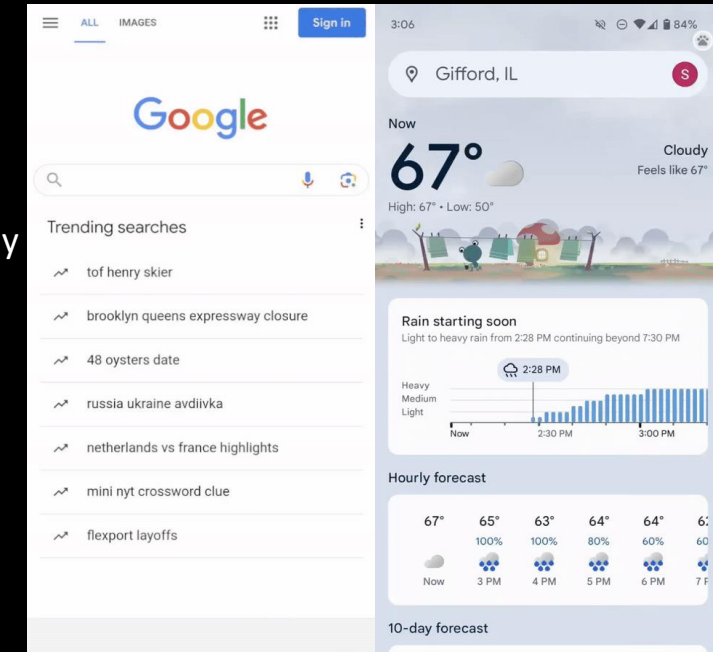
GraphCast

- Global forecast
- Graph Neural Network
- ~4 weeks training on 32 Cloud TPUv4
- 6 variables at 37 levels plus 5 surface variables (0.25 degree, 28x28km at eq.)

Real-time update
Personalization
Accuracy and Reliability



Data collection +
Data Processing +
Data Visualization



Global Weather Forecast - Microsoft



Weather from Microsoft Start

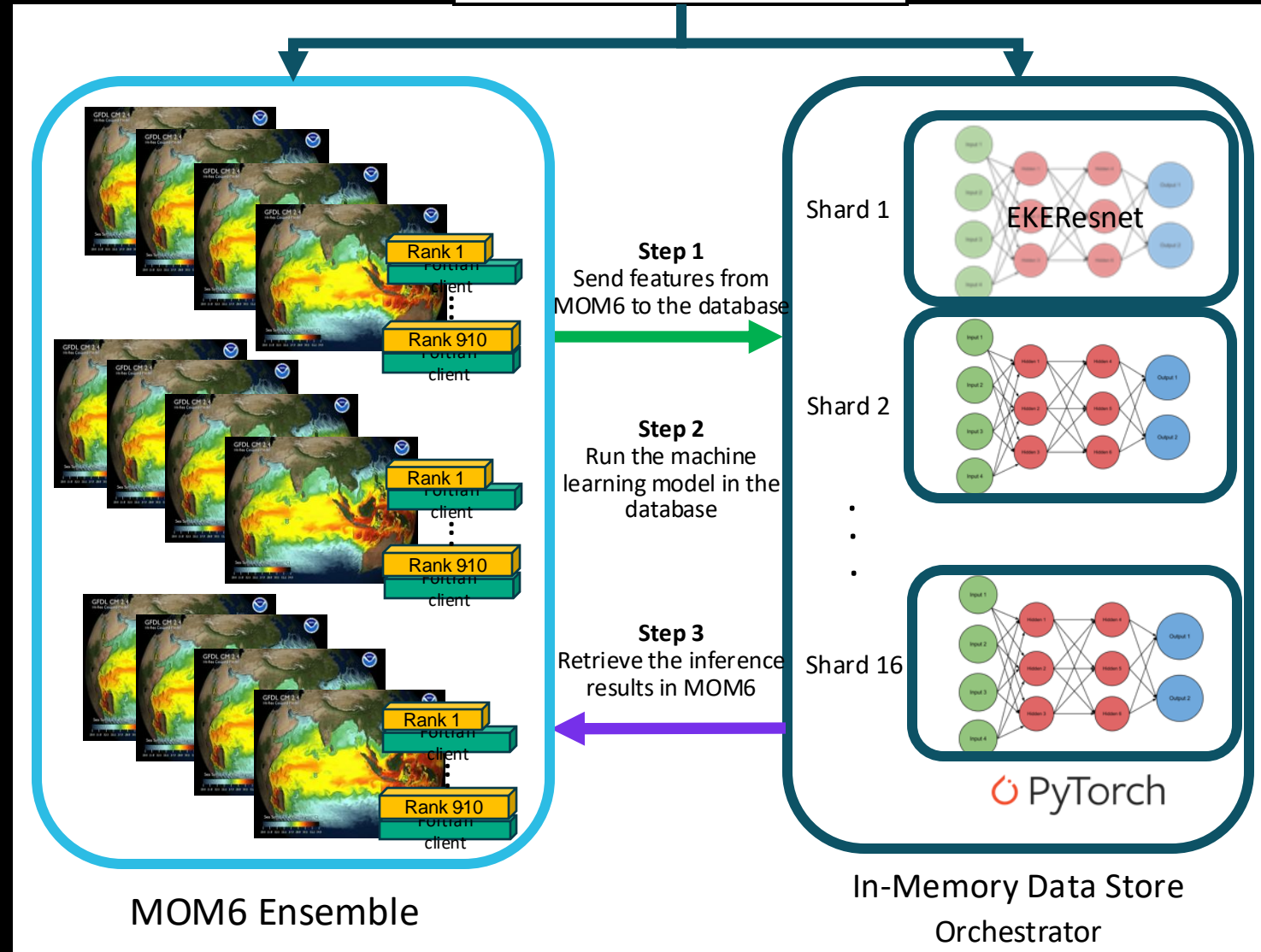
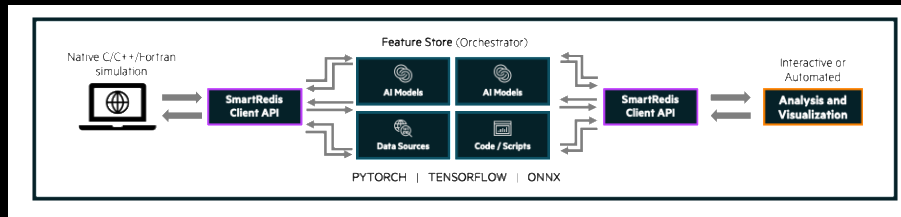
- Ambassador of NOAA's Weather Ready Nation
- A part of UN's Early Warnings for ALL

ClimaX

- Developed by Microsoft
- Built as a foundation model, trained on CMPI6 climate model output, can be fine tuned with other data at different scales
- Source code in github

Hybrid AI Models

- “Cognitive Simulation” or “CogSim” → Model Physic + AI
 - Emulators for complex physical parameterizations
 - Challenging to integrate one ML parametrization into a larger model and maintain stability
- Bias correctors inside time-stepping (ECMWF IFS and 4DVAR)
- SmartSim from HPE (open source) enables Fortran + data science in one application



Thank you



Edward.Habjan@hpe.com

Ilene.Carpenter@hpe.com



CONFIDENTIAL | AUTHORIZED

© 2024 Hewlett Packard Enterprise Development LP