









#### Post-PC Era: Late 2000s - Present



Personal Mobile Devices (PMD): Relying on wireless networking, Apple, Nokia, ... build \$500 smartphone and tablet computers for individuals → Objective C, Java, Android OS + iOS

Cloud Computing: Using Local Area Networks, Amazon, Google, ... build \$200M Warehouse Scale Computers with 100,000 servers for Internet Services for PMDs

➔ MapReduce/Spark, Ruby on Rails





Some Videos... Amazon <sup>o</sup> Behind the Scenes | Tour of a High Tech Amazon Data Center | Infinite Innovations Tech <sup>o</sup> https://www.youtube.com/watch?v=yZINQZez0oU (2024, 11 min) ° Retired Data Center Hardware Gets a Second Chance? <sup>o</sup> https://www.youtube.com/watch?v=cstnNg1 oRo (2023, 3 min) Microsoft ° What runs ChatGPT? Inside Microsoft's AI supercomputer <u>https://www.youtube.com/watch?v=Rk3nTUfRZmo</u> (2023, 16 min) • Google: ° Google Data Center Security: 6 Layers Deep https://www.youtube.com/watch?v=kd33UVZhnAA (2023, 6 min) • NVIDIA, Facebook, IBM, HP, Dropbox, ... 15 BIGGEST Data Centers on Earth ° https://www.youtube.com/watch?v=1LmFmCVTppo (2024, 30 min) IBM, HP, Dropbox, ... 8





























#### **Distributed Programming Models and Workloads for Warehouse-Scale Computers**

- WSCs run public-facing Internet services such as search, video sharing, and social networking, as well as *batch applications*, such as converting videos into new formats or creating search indexes from Web crawls
- One of the most popular frameworks for batch processing in a WSC is MapReduce and its open-source twin Hadoop
  - ° E.g., Annual MapReduce usage at Google over time
  - ° Facebook runs Hadoop on 2000 batch-processing servers of the 60,000 servers it is estimated to have in 2011

	Aug-04	Mar-06	Sep-07	Sep-09
Number of MapReduce jobs	29,000	171,000	2,217,000	3,467,000
Average completion time (seconds)	634	874	395	475
Server years used	217	2002	11,081	25,562
Input data read (terabytes)	3288	52,254	403,152	544,130
Intermediate data (terabytes)	758	6743	34,774	90,120
Output data written (terabytes)	193	2970	14,018	57,520
Average number of servers per job	157	268	394	488

23

#### MapReduce Programming model and runtime for processing large data-sets E.g., Google's search algorithms Gool: make it easy to use 1000s of CPUs and TBs of data Inspiration: functional programming languages Programmer specifies only "what" System determines "how" Schedule parallelism, locality, communication,... Ingredients Automatic parallelization and distribution Fault-tolerance I/O scheduling Status and monitoring























# Recent Trends: Improving Energy Efficiency in WSC •Run WSC at higher temperature (71F or 22C) •Alternating cold/hot aisles, separated by thin plastic sheets

- Water-to-water intercooler
  - $^\circ$  Google's WSC in Belgium takes cold water from an industrial canal to chill the warm water from inside the WSC
- •Airflow simulation to carefully plan cooling



#### **Airflow within Container (for container based WSC)**



- Two racks (attached to ceiling) on each side of the container
- Cold air blows into the aisle in the middle of the container (from below) and is then sucked into the servers
  - "cold" air is kept 81°F (27°C)
  - Careful control of airflow allows this high temperature vs. most datacenters
- Warm air returns at the edges of the container
- Design isolates cold and warm airflows

















# **Memory Hierarchy of a WSC**

•Servers can access DRAM and disks on other servers using a NUMA-style interface

	Local	Rack	Array
DRAM latency (microseconds)	0.1	100	300
Disk latency (microseconds)	10,000	11,000	12,000
DRAM bandwidth (MB/sec)	20,000	100	10
Disk bandwidth (MB/sec)	200	100	10
DRAM capacity (GB)	16	1040	31,200
Disk capacity (GB)	2000	160,000	4,800,000





















	AV	VS I	Price	25		
Viewing 758 available instances						
٩			<	<b>1</b> 2 3 4	5 6 7 38	>
Instance name	On-Demand hourly rate  ⊽	vCPU ⊽	Memory 🔻	Storage ♥	Network performance	▽
t4g.nano	\$0.0042	2	0.5 GiB	EBS Only	Up to 5 Gigabit	
t4g.micro	\$0.0084	2	1 GiB	EBS Only	Up to 5 Gigabit	
t4g.small	\$0.0168	2	2 GiB	EBS Only	Up to 5 Gigabit	
t4g.medium	\$0.0336	2	4 GiB	EBS Only	Up to 5 Gigabit	
t4g.large	\$0.0672	2	8 GiB	EBS Only	Up to 5 Gigabit	
t4g.xlarge	\$0.1344	4	16 GiB	EBS Only	Up to 5 Gigabit	
t4g.2xlarge	\$0.2688	8	32 GiB	EBS Only	Up to 5 Gigabit	
t3.nano	\$0.0052	2	0.5 GiB	EBS Only	Up to 5 Gigabit	



![](_page_29_Figure_0.jpeg)

![](_page_29_Figure_1.jpeg)

![](_page_30_Picture_0.jpeg)

![](_page_30_Picture_1.jpeg)

![](_page_31_Figure_0.jpeg)

![](_page_31_Figure_1.jpeg)

WSC Summary
<ul> <li>WSC not the same as traditional Datacenters</li> <li>Scale, datacenter-is-computer, RLP, costs</li> </ul>
<ul> <li>Total cost of ownership</li> <li>Capex-amortized (facility, P&amp;C, server, networking), Opex (P&amp;C, people, bandwidth)</li> </ul>
<ul> <li>Architecture and key building blocks</li> <li>Warehouse, container, computer, storage, network, power delivery design, cooling design</li> <li>Balance: HW/SW co-design</li> </ul>
<ul> <li>Reliability, availability and serviceability (RAS), Energy</li> <li>Scale changes everything: non-traditional models for RAS, more aggressive energy efficiency focus</li> </ul>
<ul> <li>WSC Software         <ul> <li>Cluster-level infrastructure software</li> <li>Resource management (e.g., cluster scheduler), Hardware abstraction and other basic services (e.g., GFS), programming frameworks (e.g., MapReduce/Spark)</li> </ul> </li> <li>Deployment and maintenance         <ul> <li>Service-level dashboards, performance debugging tools, platform-level monitoring (Google Health Infrastructure)</li> <li>Application-level software</li> </ul> </li> </ul>

# <section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item>

![](_page_33_Figure_0.jpeg)

![](_page_33_Figure_1.jpeg)

![](_page_34_Picture_0.jpeg)

# **Exascale Computing: The Future of Supercomputing**

![](_page_34_Figure_3.jpeg)

Science at Scale	
<ul> <li>Combustion simulations improve future designs</li> <li><sup>°</sup> Model fluid flow, burning and chemistry</li> <li><sup>°</sup> Uses advanced math algorithms</li> <li><sup>°</sup> Requires petascale systems today</li> </ul>	
<ul> <li>Need exascale (10<sup>18</sup> operations per second) computing to design for alternative fuels, new devices</li> </ul>	Simulations reveal features not visible
<ul> <li>Impacts of Climate Change</li> </ul>	In lab experiments
<ul> <li>Warming ocean and Antarctic ice sheet key to sea level rise</li> <li>Previous climate models inadequate</li> <li>Adaptive Mesh Refinement (AMR) to resolve ice-ocean</li> </ul>	X
<ul> <li>interface</li> <li>Dynamics very fine resolution (AMR)</li> <li>Antarctica still very large (scalability)</li> </ul>	
<ul> <li>Exascale machines needed to improve detail in models, including ice and clouds</li> </ul>	

![](_page_35_Picture_1.jpeg)

#### • From Simulation to Image Analysis

- $^{\circ}\,$  Computing on Data key in 4 of 10 Breakthroughs of the decade
  - 3 Genomics problems (better DNA, microbe, ancestry analysis) + CMB (cosmic microwave background; to understand origin of universe)
- ° Data rates from experimental devices will require exascale volume computing

#### Image Analysis in Astronomy

- Data Analysis in 2006 Nobel Prize
   Measurement of temperature patterns
- ° Simulations used in 2011 Prize
  - Discovery of the accelerating expansion of the universe through observations of distant supernovae
- <sup>o</sup> More recently: astrophysics discover early nearby supernova.
  - Rare glimpse of a supernova within hours of explosion, 20M light years away
  - Telescopes world-wide redirected to catch images

![](_page_35_Picture_13.jpeg)

![](_page_35_Picture_14.jpeg)

#### **Science through Volume: Screening Drugs to Batteries**

• Large number of simulations covering a variety of related materials, chemicals, proteins,...

![](_page_36_Figure_2.jpeg)

Dynameomics Database

Improve understanding of disease and drug design, e.g., 11,000 protein unfolding simulations stored in a public database.

![](_page_36_Figure_5.jpeg)

#### Materials Genome

Cut in half the 18 years from design to manufacturing, e.g., 20,000 potential battery materials stored in a database

### **Many Other Domains Need Exascale**

![](_page_36_Picture_9.jpeg)

![](_page_36_Picture_10.jpeg)

Energy Storage Understanding the storage and flow of energy in nextgeneration nanostructured carbon tube supercapacitors

![](_page_36_Picture_12.jpeg)

Source: Steven E. Koonin, DOE

Turbulence Understanding the statistical geometry of turbulent dispersion of pollutants in the environment.

![](_page_36_Picture_15.jpeg)

#### Biofuels

A comprehensive simulation model of lignocellulosic biomass to understand the bottleneck to sustainable and economical ethanol production.

![](_page_36_Picture_18.jpeg)

#### Nuclear Energy High-fidelity predictive simulation tools for the design of pext-operation nuclear

of next-generation nuclear reactors to safely increase operating margins.

#### Smart Truck

Aerodynamic forces account for ~53% of long haul truck fuel use. ORNL's Jaguar predicted 12% drag reduction and yielded EPA-certified 6.9% increase in fuel efficiency.

![](_page_36_Picture_23.jpeg)

Nano Science Understanding the atomic and electronic properties of nanostructures in nextgeneration photovoltaic solar cell materials.

## **Summary: Need for Supercomputing**

#### • Strategic importance for supercomputing

- Essential for scientific research
- Critical for national security
- Fundamental contributor to the economy and competitiveness through use in engineering and manufacturing
- Supercomputers are the tools for solving the most challenging problems through Simulations!

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	DOE/SC/Oak Ridge National Laboratory United States	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11 HPE	8,699,904	1,206.00	1,714.81	22,786
2	DOE/SC/Argonne National Laboratory United States	Aurora - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Stingshot-11 Intel	9,264,128	1,012.00	1,980.01	38,698
3	Microsoft Azure United States	Eagle - Microsoft NDv5, Xeon Platinum 8:480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR Microsoft Azure	2,073,600	561.20	846.84	
4	RIKEN Center for Computational Science Japan	Supercomputer Fugaku - Supercomputer Fugaku, A6&FX 48C 2.2GHz, Tofu interconnect D Fujitsu	7,630,848	442.01	537.21	29,899
5	EuroHPC/CSC Finland	LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11	2,752,704	379.70	531.51	7,107

# **Challenges of Exascale Computing**

- 10-100 Million processing elements (cores or mini-cores) with chips perhaps as dense as 1,000 cores per socket?
- Energy costs?
  - ° At ~\$1M per MW, energy costs are substantial
  - ° 1 petaflop in 2010 used 3 MW
  - $^\circ~$  1 exaflop in 2018 possible in 200 MW with "usual" scaling
- 3D packaging?
- Large-scale optics/photonics based interconnects?
- 10-100 PB of aggregate memory?
- Hardware and software-based fault management?
- Heterogeneous cores?
- Performance per watt?
- Power, area and capital costs will be significantly higher?

![](_page_38_Figure_13.jpeg)

• Ex	cascale Computing Project ° <u>https://www.exascaleproject.org/about/</u>
• Th	1e Energy Exascale Earth System Model (E3SM) Project ° <u>https://e3sm.org/about/</u>
• M	odeling and Simulation at the Exascale for Energy and the Environment ° <u>https://science.osti.gov/-/media/ascr/pdf/program-documents/docs/Townhall.pdf</u>
• Cr	<ul> <li>ossCut Report – Exascale Requirements Reviews</li> <li><u>https://science.osti.gov/-/media/ascr/pdf/programdocuments/docs/2018/DOE-ExascaleReport- CrossCut.pdf</u></li> </ul>
• N'	VIDIA at Supercomputing <sup>o</sup> https://www.nvidia.com/en-us/events/supercomputing/?ncid=pa-srch-goog- <u>32907&amp;gclid=EAIaIQobChMI4pyVpE7QIVTfDACh0n2QIeEAAYAiAAEgIzyPD_BwE#cid=hpc06_pa-srch</u> <u>goog_en-us</u>
• Th	ie Convergence of AI and HPC <sup>o</sup> <u>https://www.intel.com/content/www/us/en/high-performance-computing/supercomputing/exascale computing.html</u>
•	