

SUPERCOMPUTERS AND SUPERINTELLIGENCE

Horst Simon, <u>hdsimon@lbl.gov</u>

17th SIAM Conference on Parallel Processing for Scientific Computing, Paris, April 12 – 15, 2016



"It's not a human move ..."

- We just experienced another milestone in machine intelligence:
- Alpha Go of Deep Mind (Google) winning Go against Lee Sedol, one of the world's top go players.

March 11, 2016



+= +=

http://www.wired.com/2016/03/sadness-beautywatching-googles-ai-play-go/?mbid=social_fb



Recent concerns about the "machines" taking over



"Those disposed to dismiss an 'Al takeover' as science fiction may think again after reading this original and well-argued book." – Martin Rees, Past President, Royal Society



"Worth reading Superintelligence by Bostrom. We need to be super careful with AI. Potentially more dangerous than nukes." – Elon Musk



"If our own extinction is a likely, or even possible, outcome of our technological development, shouldn't we proceed with great caution?" – *Bill Joy* NICK BOSTROM



"Success in creating AI would be the biggest event in human history. Unfortunately, it might also be the last, unless we learn how to avoid the risks."

- Steven Hawking



Slide adapted from Larry Smarr, UCSD



- 1 Current Trends in Supercomputing (The Path to Exascale)
- 2 Computing and the Brain (Hype and reality)
- 3 What Computers Still Can't Do

(and what I think the real dangers are)



And what about supercomputers?

- At any given time one of the most powerful computers to solve scientific and engineering problems
- Supercomputers and HPC are largely absent from the public discussion about progress in AI



Exascale initiatives are advancing the computational power of supercomputers

- NSCI (National Strategic Computing Initiative) announced by President Obama in June 2015
- Exascale Computing Project ECP started by DOE in the US
- Similar initiatives in Europe, Japan, and China



The TOP500 Project (by Meuer, Strohmaier, Dongarra, Simon)



Listing of the 500 most powerful computers in the world Yardstick: Rmax of Linpack

- Solve Ax=b, dense problem, matrix is random
- Dominated by dense matrix-matrix multiply
- Updated twice a year:
- ISC'xy in June in Germany
- SCxy in November in the U.S.

All information available from the TOP500 web site at: www.top500.org

See also Strohmaier, Meuer, Dongarra, and Simon, IEEE Computer, Nov. 2015, pp. 32-39.



Hans Meuer (1936 - 2014)

PERFORMANCE DEVELOPMENT



500



SIAM PP 2016, Paris | April 12 - 15, 2016

PROJECTED PERFORMANCE DEVELOPMENT





1994 1996 1998 2000 2002 2004 2006 2008 2010 2012 2014 2016 2018 2020



- 1 Current Trends in Supercomputing (The Path to Exascale)
- 2 Computing and the Brain (Hype and reality)
- 3 What Computers Still Can't Do (and what I think the real dangers are)



TOP 500 Performance Projection -The Old Picture From 2007





The Exponential Growth of Computing



Adapted from Kurzweil, The Age of Spiritual Machines.



Growth of Computing Power & "Mental Power"





Hans Moravec, CACM 10, 2003, pp 90-97.

Arguments against this simplistic view

- Naïve extrapolation of current performance trends
 - transition to Post Moore's Law computing not clear
 - HPC systems are focused on excelling in scientific computing
- Scaling of AI and machine learning to millions of cores
 - most powerful computers today are not being used for cognitive tasks
- Successes of machine learning are accomplished through progress in algorithms



History Lesson: 1997

- IBM Deep Blue beats Gary Kasparov (May 1997)
- One of the biggest success stories of machine intelligence,
- However, the chess computer "Deep Blue", did not teach us anything about how a chess grandmaster thinks
- No further analysis or further developments
- 19 years later the story repeats itself with Go





Today's Supercomputers and the Brain



Transistors	Memory	Clock (GHz)	Power (W)	Weight (kg)	Size (ℓ)
9x10 ¹²	2x10 ¹⁵	2	1.2x10 ⁶	18,800	36,000

9x10 ¹⁰	1x10 ¹⁵	10 ⁻⁹ – 10 ⁻⁵	20	1.5	
Neurons	Syn. Conn.				A STAN

Slide from Peter Denes

Modha Group at IBM Almade



	-9	5	100
(h)-	1.	YY	H-N+
(The	TI	all	20
HAC'S	1	PG.	- A
000	J.	F.	S. 1

Mouse	Rat	Cat	Monkey	Human
N: 16 x 10 ⁶	56 x 10 ⁶	763 x 10 ⁶	2 x 10 ⁹	22 x 10 ⁹
S: 128 x 10 ⁹	448 x 10 ⁹	6.1 x 10 ¹²	20 x 10 ¹²	220 x 10 ¹²



Latest simulations in 2012 achieve unprecedented scale of 65*10⁹ neurons and 16*10¹² synapses

Real time simulation at the exascale





1994 1996 1998 2000 2002 2004 2006 2008 2010 2012 2014 2016 2018 2020

Ananathanarayanan et al., "The Cat is out of the Bag: Cortical Simulations with 10⁹ Neurons and 10¹⁵ Synapses", Proceedings of SC09.

Recent Evidence for Petabyte Size Memory

- Bartol et al., "Nanoconnectomic Upper Bound on the Variability of Synaptic Plasticity", Jan. 2016, eLife.
- 26 sizes of synapses corresponds to about 4.7 bits per synapse and thus about 4 5 Petabytes



from http://www.salk.edu/news-release/memory-capacity-ofbrain-is-10-times-more-than-previously-thought/



Compute Power of the Human Brain

- Estimate of compute power for the human brain is about 1-10 Exaflops and 4-5 Petabytes
- Three different paths lead to about the same estimate
- A digital computer with this performance might be available in about 2024 with a power consumption of at best 20–30 MW (goal of the Exascale project)
- The human brain takes 20 W
- A digital exaflops computer using CMOS technology will still be a factor of a million away from brain power





Dimensions of Intelligence

- 1. Verbal-Linguistic
 - ability to think in words and to use language to express and appreciate complex concepts
- 2. Logical-Mathematical
 - makes it possible to calculate, quantify, consider propositions and hypotheses, and carry out complex mathematical operations
- 3. Pattern Recognition
 - capacity to recognize and think about common pattern in our fourdimensional environment
- 4. Bodily-Kinesthetic
 - ability to manipulate objects and fine-tune physical skills
- 5. Musical
 - sensitivity to pitch, melody, rhythm, and tone
- 6. Interpersonal
 - capacity to understand and interact effectively with others

After Howard Gardner. *Frames of Mind: The Theory of Multiple Intelligences*. New York: Basic Books, 1983, 1993.







A Comment about the Turing Test

- Test only one dimension - verbal linguistic

- Need additional tests that explore other dimensions of human intelligence
- My favorite is the "Ikea test" for pattern recognition and 3d spatial thinking: have a robot build furniture from the schematic drawing

http://io9.gizmodo.com/8-possible-alternatives-to-theturing-test-1697983985





Performance Development



Source: TOP500 November 2015.

Historical Perspective, or "Why We Are Here Today"



John von Neumann and Robert Oppenheimer, Princeton, IAS, 1952

- 60 years of large scale
 computational physics
 applications driving computer
 development
- Remarkable longevity
- Architecture matched to application: Von Neumann architecture and focus on floating point performance



From the Bomb to the Cloud in Sixty Years



- Von Neumann architecture is ideally suited for large-scale, floating point intense, 2D or 3D grid-based, computational physics applications
- Today we are using the same basic architecture for social networking, Web searches, music, photography, etc.





From A. Merolla et al., *Science*, Aug. 8, 2014



Recent Neuromorphic Projects

Five complementary approaches to Neuromorphic Computing (massively parallel, asynchronous communication, configurable):

- Commodity microprocessors (SpiNNaker, HBP)
- Custom fully digital (IBM Almaden)
- Custom mixed-signal (BrainScaleS, HBP)
- Custom subthreshold analog cells (Stanford, ETHZ)
- Custom yybrid (Qualcomm)





IBM SyNAPSE Project (D. Modha)

- 5.4 B transistor chip TrueNorth, 4096 neurosynaptic cores, 1 M spiking neurons, 256 M configurable synapses
- 63 mW power per chip, significantly less energy per event (176,000) when compared to a simulator



- Scalability to large system
- Corelet programming model







- 1 Current Trends in Supercomputing (The Path to Exascale)
- 2 Computing and the Brain (Hype and reality)
- 3 What Computers Still Can't Do (and what I think the real dangers are)



About 1967



Kosmos LOGIKUS "Spielcomputer"

My first computer:

- 10 electric lamps
- 10 switches
- Programming by wiring

http://www.logikus.info/

What I learned: computers are just machines, wires, batteries, and light bulbs. They cannot have a mental state or experience. Many years later I realized that this position is equivalent to that I don't believe in strong AI.











Strong and Weak AI

Strong AI: A physical symbol system can have a mind and mental states.

Does my play computer get upset when I win?

Does the Edison computer at NERSC know what a hurricane IS?

HDS: This is an interesting philosophical question, but we can leave it aside for this discussion.

Weak AI: A physical symbol system can act intelligently.

Is it possible to develop a computer system that performs indistinguishable from a human?

HDS: Yes, in principle, but it will be very, very hard.



Some challenges ahead for modeling the brain

- Unsuitability of current architectures
 - HPC systems are focused on excelling in computing; only one of the six (or eight) dimensions of human intelligence
- Fundamental lack of mathematical models for cognitive processes
 - That's why we are not using the most powerful computers today for cognitive tasks
- Lack of data, standards. That's what the BRAIN initiative in the US should address
- We should be able to model the brain as a complex physical system, but we have barely started



What we should be really concerned about: (1) large scale complex systems controlled by algorithms

Example:

High Frequency Trading (HFT) is now accounting for over 60% of the volume in US equity markets.

The interaction of multiple algorithms create a complex system that we don't understand any longer, yet our prosperity depends on it.



See also Center for Innovative Financial Technologies (CIFT) at LBNL: http://crd.lbl.gov/departments/data-science-and-technology/sdm/current-projects/cift





What we should be really concerned about: (2) confluence of pattern recognition machines + image analysis + big data + behavioral prediction

Deployment of neuromorphic processors that excel at pattern recognition inexpensively in large scale (IoT)

Collection of huge amounts of data and image in the cloud

Capability of real time streaming data analysis

Behavioral prediction

For friends of SF (Philip K. Dick): I am more concerned about "Minority Report Future" than a "Blade Runner Future"







SIAM PP 2016, Paris | April 12 - 15, 2016

Towards the pattern recognition machine





From D. Modha, IBM Research

Summary – Key Messages

- Current HPC technology is about a factor of 10⁹ away from the real time performance of the human brain (10⁶ in power, 10³ in computation, <10 in memory)
- Both technology and architectural innovation are needed to close the gap.
- Given what we know today, an artificial, sentient, superintelligence is unlikey (strong AI).
- Given what we know today, a realistic, highly accurate simulation of brain functions will require major advances in systems, algorithms, and mathematical modeling, in parallel with progress in neuroscience (weak AI).
- The real danger is turning over decisions to systems that we don't understand and control.



Thank You for Contributions

Erich Strohmaier (LBNL) Jack Dongarra (UTK) John Shalf (LBNL) Peter Denes (LBNL) Michael Wehner and team (LBNL)

Dharmendra Modha and team (IBM) Karlheinz Meier (Univ. Heidelberg) and Wikipedia







41st List: The TOP10

#	Site	Manufacturer	Computer	Country	Cores	Rmax [Pflops]	Power [MW]
1	National University of Defense Technology	NUDT	Tianhe-2 NUDT TH-IVB-FEP, Xeon 12C 2.2GHz, IntelXeon Phi	China	3,120,000	33.9	17.8
2	Oak Ridge National Laboratory	Cray	Titan Cray XK7, Opteron 16C 2.2GHz, Gemini, NVIDIA K20x	USA	560,640	17.6	8.21
3	Lawrence Livermore National Laboratory	IBM	Sequoia BlueGene/Q, Power BQC 16C 1.6GHz, Custom	USA	1,572,864	17.2	7.89
4	RIKEN Advanced Institute for Computational Science	Fujitsu	K Computer SPARC64 VIIIfx 2.0GHz, Tofu Interconnect	Japan	795,024	10.5	12.7
5	Argonne National Laboratory	IBM	Mira BlueGene/Q, Power BQC 16C 1.6GHz, Custom	USA	786,432	8.59	3.95
6	Los Alamos NL / Sandia NL	Cray	Trinity Cray XC40, Xeon E5 16C 2.3GHz, Aries	USA	301,0564	8.10	
7	Swiss National Supercomputing Centre (CSCS)	Cray	Piz Daint Cray XC30, Xeon E5 8C 2.6GHz, Aries, NVIDIA K20x	Switzerland	115,984	6.27	2.33
8	HLRS – Stuttgart	Cray	Hazel Hen Cray XC40, Xeon E5 12C 2.5GHz, Aries	Germany	185,088	5.64	
9	King Abdullah University of Science and Technology	Cray	Shaheen II Cray XC40, Xeon E5 16C 2.3GHz, Aries	Saudi Arabia	196,608	5.54	2.83
10	Texas Advanced Computing Center/UT	Dell	Stampede PowerEdge C8220, Xeon E5 8C 2.7GHz, Intel Xeon Phi	USA	462,462	5.17	4.51



Tianhe-2 (TH-2) at NUDT, China – #1 on the TOP500 list



Summary of the Tianhe-2 (TH-2) Milkyway 2

Model	TH-IVB-FEP	Summary of t	
Nodes	16,000	Items	
Processor	Intel Xeon IvyBridge E5-2692		-
Speed	2.200 GHz		
Sockets per Node	2		
Cores per Socket	12	Processors	
Coprocessers	Intel Xeon Phi 31S1P	1100033013	
Coprocessors per Node	3		
Cores per Coprocessor	57		
Coprocessors total	48,000		
Operating System	Kylin Linux	Interconnect	
Primary Interconnect	Proprietary high-speed interconnecting		
network	(TH Express-2)	Momony	
Peak Power (MW)	17,8	wentory	
Size of Power Measurements (Co	ores) 3,120,000		
Memory per Node (GB)	64	Storage	
Summary of all components		Cabinata	-
	004.000	Cabinets	
CPU Cores	384,000		

48,00

1,024,000 GB

Summary of the Tianhe-2 (TH-2) or Milkyway-2				
Items	Configuration			
Processors	32,000 Intel Xeon CPU's + 48,000 Xeon Phi's (+4096 FT-1500 CPU's frontend) Peak Performance 54.9 PFlop/s (just Intel parts)			
Interconnect	Proprietary high-speed interconnection network, TH Express-2			
Memory	1 PB			
Storage	Global Shared parallel storage system, 12.4 PB			
Cabinets	125 + 13 + 24 = 162			
Power	17.8 MW			
Cooling	Closed air cooling system			



Accelerators/CP

Memory

Accelerator/CP Cores

Two Supercomputers in Berkeley



The Exponential Growth of Computing



Adapted from Kurzweil, The Age of Spiritual Machines



Replacement Rate





Performance Development



1994 1996 1998 2000 2002 2004 2006 2008 2010 2012 2014



Projected Performance Development





History Lesson: 1987

"Legendary" CM-2 by Thinking Machines

Architecture evolved into CM-5 (1992) built as MPP for scientific applications

Early history of AI applications on parallel platforms has been lost





History Lesson: 2011

- IBM Watson beats two best human players
- Still only a narrowly defined game, but Watson demonstrates significant progress in language processing
- Uses supercomputer
 architecture
- IBM is developing commercial applications based on Watson
- No impact on HPC







DOE / ASCR -> BRAIN Initiative?

Neuroscience is a vast field

- BRAIN Initiative is a small part
- Driven by measurement, technology and computing

DOE National Labs would bring

- Team science + interdisciplinary integration
- Systems engineering
- · Facilities and problems of scale

ASCR would bring







Computing for BRAIN is synergistic with ASCR strategic directions



Computing requirements for BRAIN are well aligned with DoE/ASCR data strategy. Investment in applied mathematics and computer science will be synergistic.



Data Management



Devise standardizations and models for curation, provenance-tracking, and fusion of multi-

Develop modality agnostic data analytics methods and visualizations to reveal structure of brain data.

Analysis

Methods

X

 J_d

 f_d^{-1}

Y

Theory and **Models**



HPC **Facilities**



Construct rigorous mathematical theories and simulations to bridge multiple spatiotemporal scales.

Provide community access to centralized data repositories and high-performance computing facilities.



modal brain data.

ASCR can uniquely contribute to BRAIN



ASCR can play a unique role in BRAIN computing through advances in applied mathematics and computer science together with HPC facilities.





Integration and Synthesis are Crucial



Persistent, open (Web-based) access to multimodal data to allow the neuroscience community to perform exploratory analysis for data driven discovery (**Data Superfacility**).



Function

dynamic data

Theory & Models

abstractions



Derived ****products

Integration & Synthesis

Structure

static data





BRAIN Initiative

#3. **The brain in action**: Produce a *dynamic picture* of the functioning brain by developing and applying improved methods for large-scale monitoring of neural activity. #5. **Identifying fundamental principles**: Produce conceptual foundations for understanding the biological basis of mental processes through *development of new theoretical and data analysis tools.*

#2. **Maps at multiple scales**: Generate circuit diagrams that *vary in resolution* from synapses to the whole brain.

Function



Theory & Models

abstractions



Structure

static data



Derived **U**products

Integration & Synthesis



Technology to the rescue



Tri-Institutional Partnership - a model



Tri-Institutional Partnership BRAIN R&D Initiative to support innovative neurotechnology



Home

NEWS & ANNOUNCEMENTS

New Partnership Launches with Meeting at Berkeley Lab April 1, 2014

National Lab

Universities

Facilities Capabilities

Neuroscience

Clinical data

UC San Francisco, the three al Partnership as a means to

ects in neurotechnology. The neement of a peer-reviewed o catalyze bold, potentially

s convened at Berkeley Lab for the opportunities provided by in research ideas and forge new



s institutions even more so," said Graham Fleming, Vice Chancellor for Research at Berkeley. tions and the energy in the room demonstrates the level of excitement they share in addressing

er's Day. Our scientists were joined by 23 scientists from UC SF and 28 from UC Berkeley," said me that the Bay Area scientific community is ready for the tri-institutional partnership in order to

technical excellence, innovation, and the substantive involvement of the collaborative partners ave a clear path from concept to the development of a competitive proposal for outside funding.

e BRAIN R&D seed-funding project to support innovative neurotechnology.



Challenges of **brain simulation**: link structure to function across scales

Disparate spatiotemporal scales

- nanometers to meters
- picoseconds to years
- Diverse data types
- genomics to (functional/structural) connectomics
- electrical, optical and other measurements
- behavior, sensory stimuli, perturbations
- Complex analysis issues
- · fusion of multi-modal data
- inference robust to noise
- provide insight into computations
- statistical prediction of future events
- extraction of important features

Lack of data and models

Human Brain Project, 2012





LBL BRAIN Related/Motivated Research

LBL LDRD

Computation

+Machine Learning (K. Bouchard)

+Data model (BRAINformat) (O. Ruebel)

+Graph Analytics (A. Buluc)

+Real-Time Processing & Vis (D. Donofrio & G. Weber)

+Neural Networks (K. Bouchard)

+Neuron Reconstruction at NERSC (Prabhat)

+fMRI analysis (D. Ushizima)

+Neuromorphic Computing (2 LDRDs, FY' 16)

Technology

+High-density electrophysiology (P. Denes)

+Up-Converting Nanoparticles (B. Cohen)

+Opto-acoustic waveguides (P. Schuck)

+Chemical Sensors (C. Chang)







Experimentation

+Neurodegeneration (C. McMurray)

+Neurological Aging (W. Jagust)

+Sensorimotor circuits (K. Bouchard)

+Neurocognitive Resilience (A. Wyrobek)

+Toxicant impact on host and microbiome (S. Celniker, A. Snijders, J-H Mao)

+Molecular Basis of Epilepsy (B. Brown)



National Institutes of Health

MOLECULAR FOUNDRY

Facilities

:17



Tri-Institutional Partnership





UNIVERSITY OF CALIFORNIA J Lawrence Berkeley National Laboratory



Electrical Sensors and Optical Emitters for BRAIN

Ultra high-density electrophysiology





Go from $10^2 \rightarrow 10^5$ electrodes

Data volume

- Today: at limit of "workstation"
 - ↑ 3-4 orders of magnitude



SIAM PP 2016, Paris | April 12 - 15, 2016

Advanced Computing for **BRAIN** at LBL

Common cycle in DOE computing





Summary

- Proponents of the rapid development of superintelligence use straightforward extrapolation of current computer performance. This ignores the end of Moore's Law, and the multidimensional nature of human intelligence.
- Simulating the human cortex in real time will require a system with 10 Exaflops, 5 Petabytes, and 20 MW.
- With the human brain taking only 20 W, current technology is at least a factor of a billion away from human brain performance.
- We must investigate new architectures if we want to close this gap.
- We must advance various "brain initiatives" to get the necessary data.



Technology Trends: Microprocessor Capability



2X transistors/chip every 1.5 years – called "Moore's Law"

Microprocessors have become smaller, denser, and more powerful



Gordon Moore (co-founder of Intel) predicted in 1965 that the transistor density of semiconductor chips would double roughly every 18 months



Sustained Growth of Technology Allowed Us to Ignore Architecture



- Two technology transitions since 1940s
- Moore's Law stated for integrated circuits only
- Kurzweil et al. claim accelerated exponential growth across technologies



Motivation for the Title of My Talk

"The computer model turns out not to be helpful in explaining what people actually do when they think and perceive."

Hubert Dreyfus, pg.189

Example: one of the biggest success stories of machine intelligence, the chess computer "Deep Blue", did not teach us anything about how a chess grandmaster thinks





Why This Is Important

- This could be the beginning of the development of a "right-brain" architecture for computers
- The future is a fast pattern recognition machine. It is not relevant if the chips are actually resembling the brain
- Energy efficiency has been demonstrated with conventional 26nm process
- Small systems could be made available easily to wide developer community (think NVIDIA and CUDA, or Raspberry Pi for pattern recognition)

67

- Potential widespread use in mobile
- Scale up and integrated into HPC.

