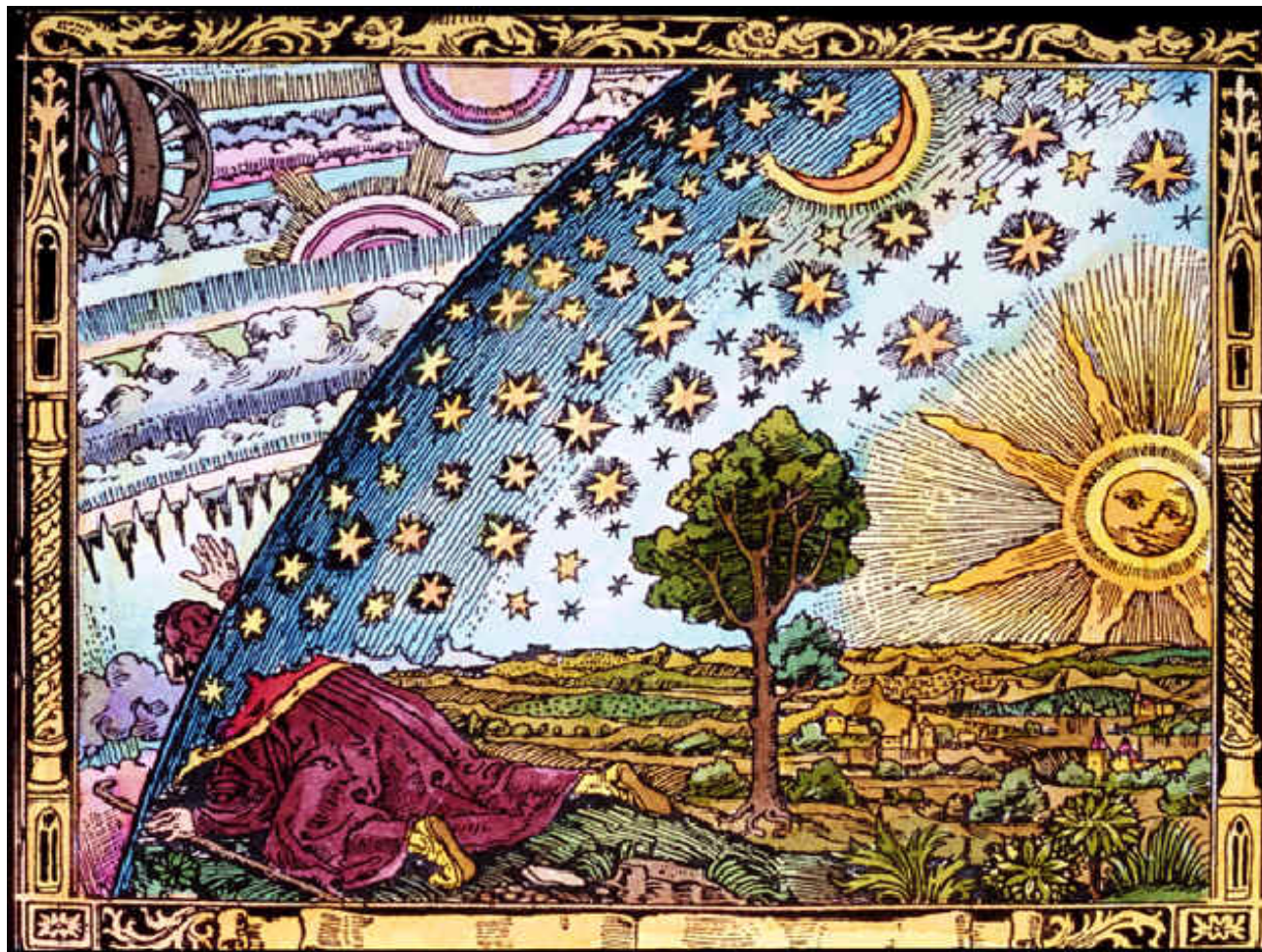


Computing at the Edge: Sensors, Learning, and Adaptation



Dan Reed
University of Utah

dan.reed@utah.edu
www.hpcdan.org

Quantitative creates qualitative and it changes everything

Viability determined by

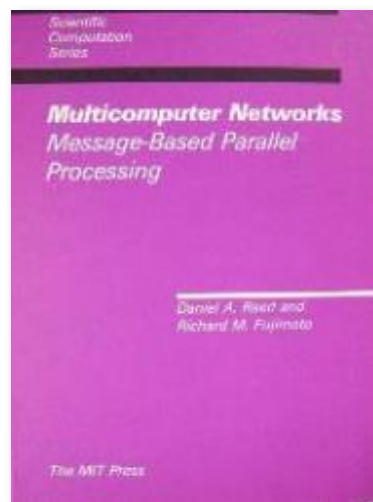
- *speed, capacity, cost, market scale*

... and their ratios

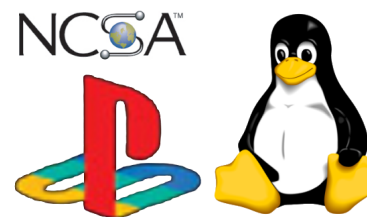


NCSA Cray X-MP (1985)

\$8,000,000 and *56 Kb/s NSFnet*
800 megaflops (peak)
(~\$18M in 2018 dollars)



Intel 80386 (1985)
16 MHz @ 2.3 watts
1.5 micron CMOS



TeraGrid™



iPhone 8 (2017)
~\$700 and LTE wireless
~3000 megaflops

Generation after generation: paradigm shifts

Success breeds complacency.
Complacency breeds failure.
Only the paranoid survive.

Andy Grove, Intel



Mainframes
IBM S/360



Minicomputers
VAX 11/780



Workstations
SUN 3/50



Personal Computers
IBM PC



Windows PCs



Smartphones



Edge + IoT

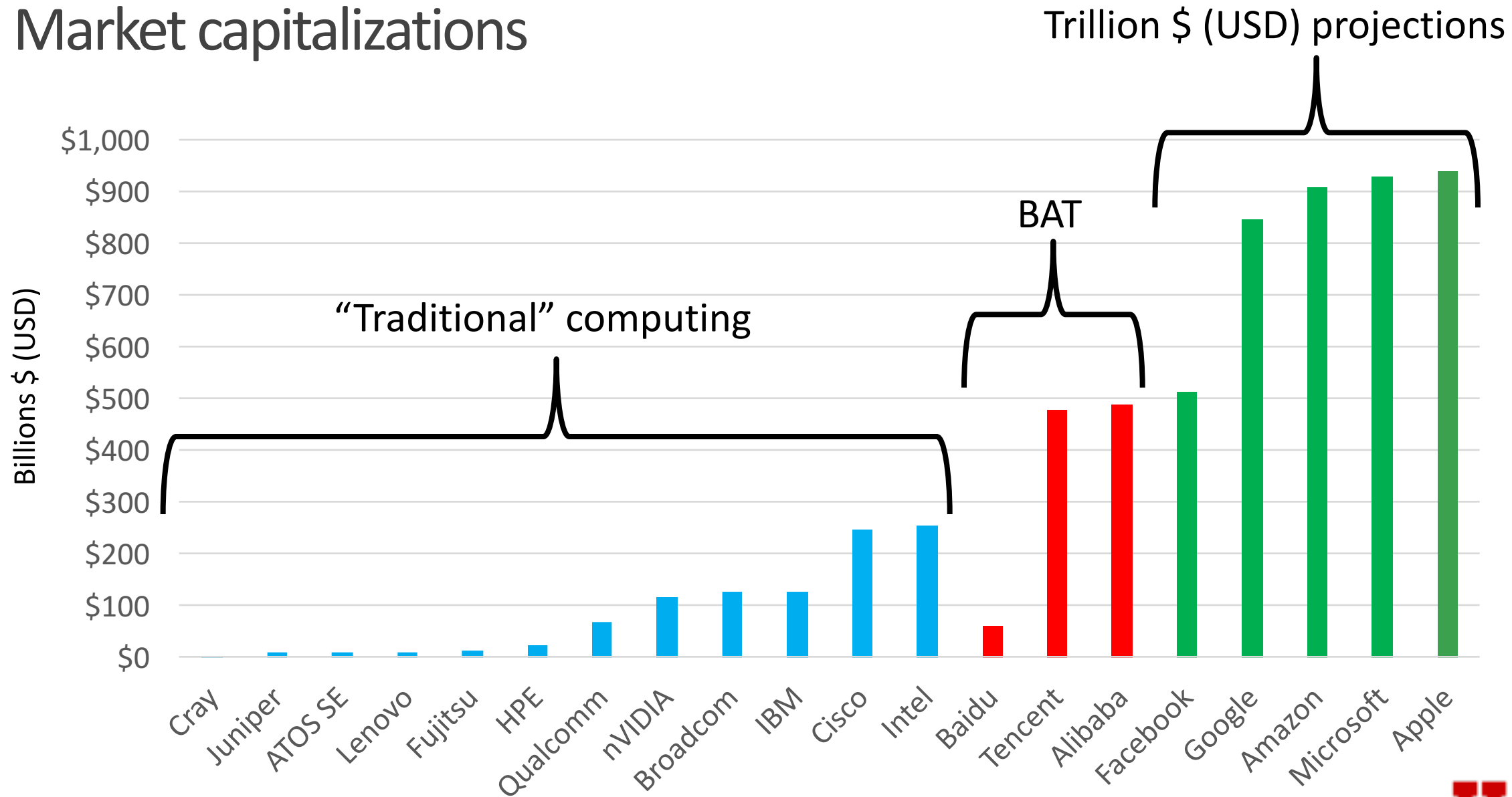


Cloud Data Centers

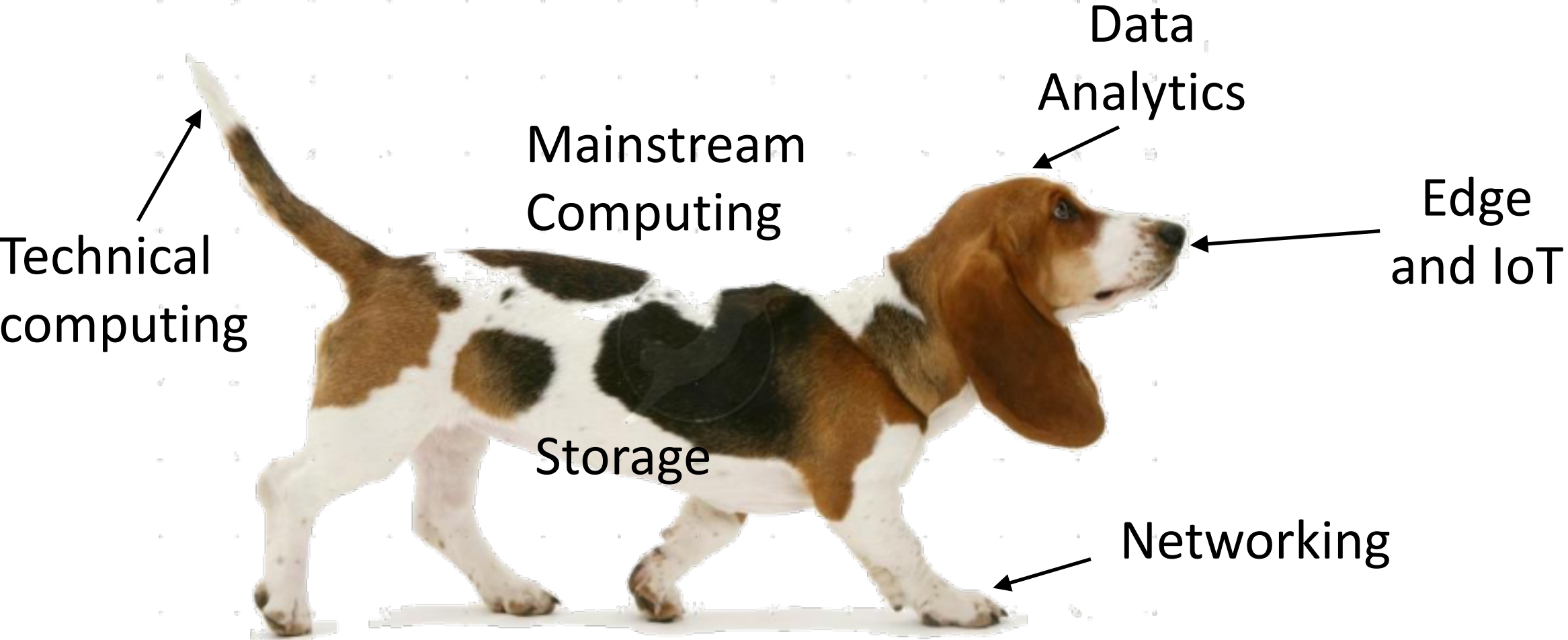


Our world has changed

Market capitalizations



Follow the money and the users ...



From petascale to the global race to exascale



ORNL Summit (200 PF)



From K to Post-K Computer



Meanwhile, everything gets smart and the trolley paradox gets real



"CAN I INTEREST YOU IN A FIREWALL FOR YOUR TOASTER?"



Oura Ring



Oxford Nanopore Mobile DNA sequencer



IFTTT



The computing continuum: holistic thinking needed



Size	Nano	Micro	Milli	Server	Fog	Campus	Facility
Example	Adafruit Trinket	Particle.io Boron	Array of Things	Linux Box	Co-located Blades	1000-node cluster	Datacenter & Exascale
Memory	0.5 KB	256 KB	8 GB	32 GB	256 GB	32 TB	16 PB
Network	BLE	WiFi/LTE	WiFi/LTE	1 GigE	10GigE	40GigE	N*100GigE
Cost	\$5	\$30	\$600	\$3K	\$50K	\$2M	\$1000M

Count = 10^9
Size = 10^1



Count x Complexity = ~Constant



Count = 10^1
Size = 10^9

Stateful vs. Stateless

Three computing revolutions ...



Deep neural networks at the edge

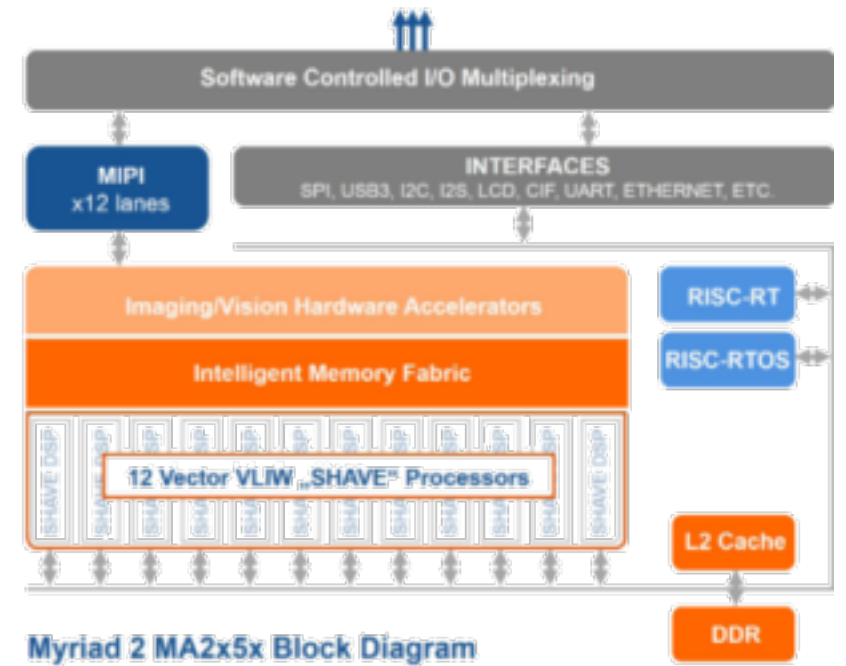
Most data *never* leaves the device



\$249 (USD)

Amazon Deep Lens

- Integrated camera
- Intel Atom and Ubuntu 16.04
- Intel Gen9 graphics engine
- TensorFlow and Caffe support
- AWS integration (obviously)



Myriad 2 MA2x5x Block Diagram



~\$80 (USD)

Movidius (Intel) Neural Compute Stick

- Custom Vision Processing Unit (VPU)
- TensorFlow and Caffe support

Edge ferment



Coral USB



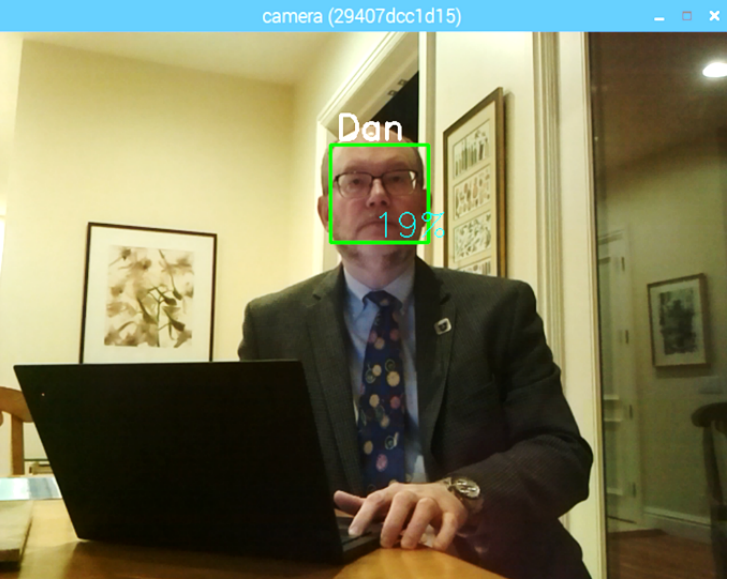
Google Coral (\$149)

- Quad-core Cortex-A53 + Cortex-M4F
- *Google Edge TPU ML accelerator*
- 8 GB eMMC
- 1 GB LPDDR4
- GigE, HDMI, USB-C, USB-3
- MicroSD slot

Jetson Nano (\$99)

- Quad-core A57
- *128 core nVidia Maxwell*
- 4 GB LPDDR4
- GigE, HDMI, USB 3
- MicroSD slot

Face recognition on the cheap



Hardware Item	Cost
Raspberry Pi 2/3 Case (optional)	\$9.99
AZ Delivery Raspberry Pi Camera	\$7.99
Raspberry Pi Model 3	\$34.99
SanDisk 64 GB microSDXC UHS-I card	\$11.59
5V 2.5A Micro USB AC Adapter	\$10.99
Qubo Phone Tripod Camera Stand (optional)	\$20.99
TOTAL	\$96.54



Serious action

Now at the two extremes

The very small (edge/fog computing and sensors)

The very large (clouds, exascale, and big data)

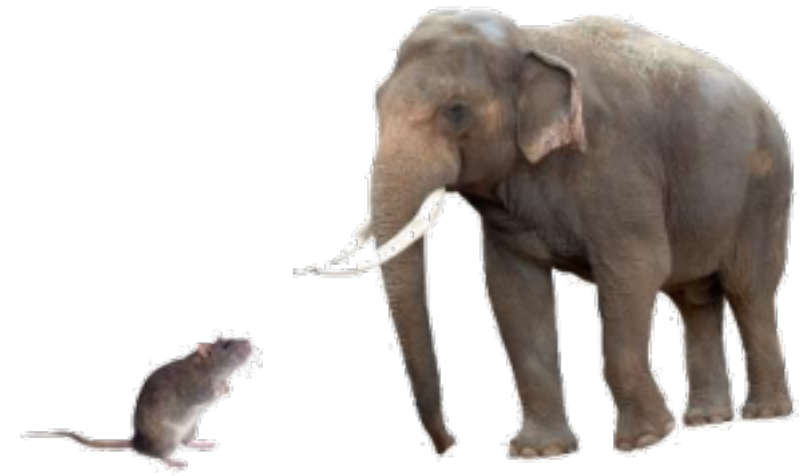
But with real constraints on each

Technical implications

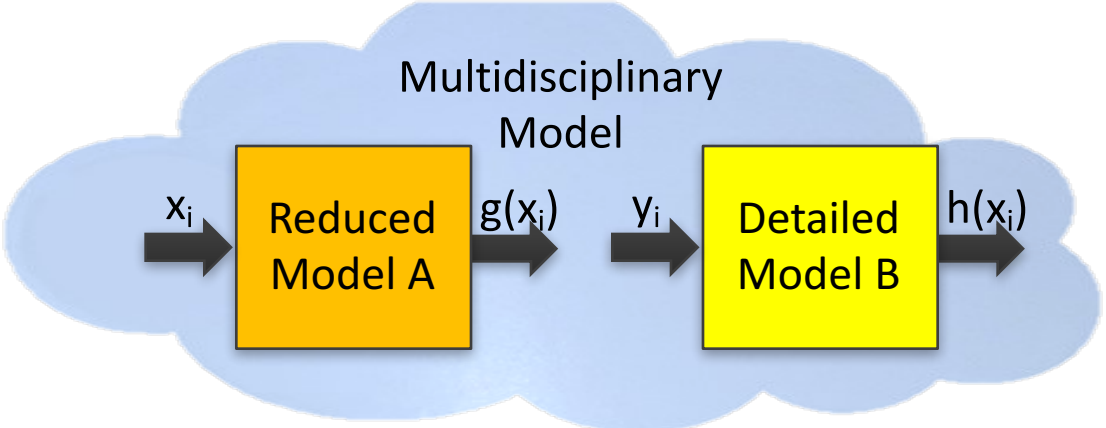
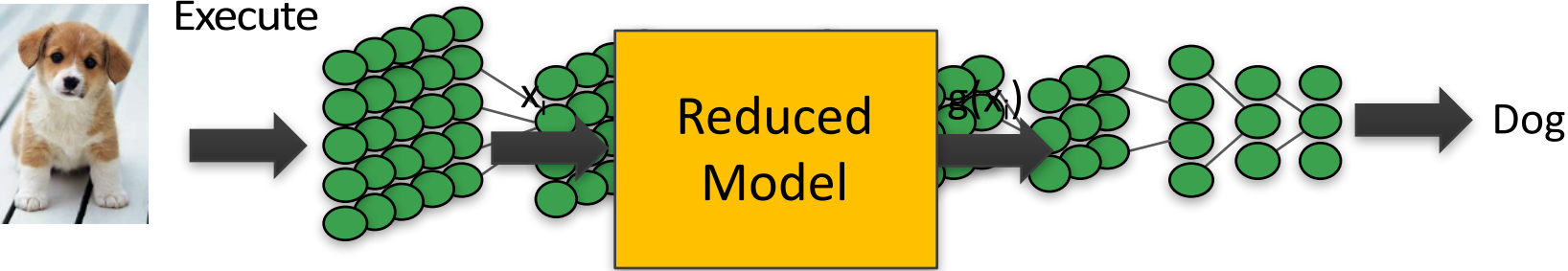
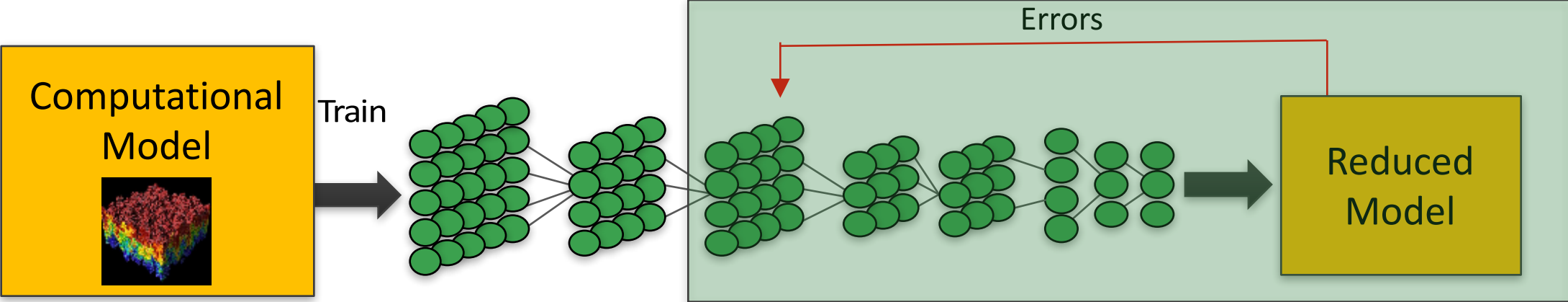
- Fluid end-to-end cyberinfrastructure
- Interdisciplinary data and infrastructure sharing

Cultural implications

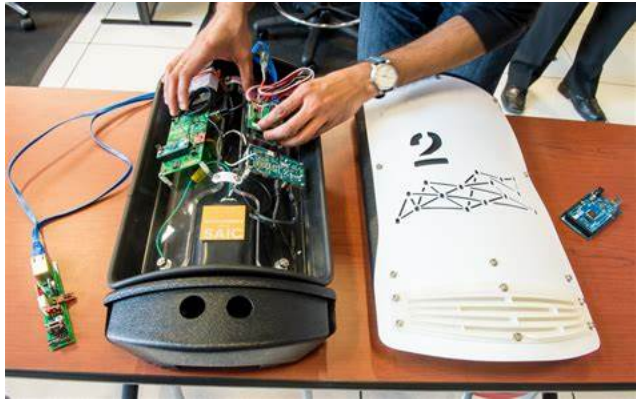
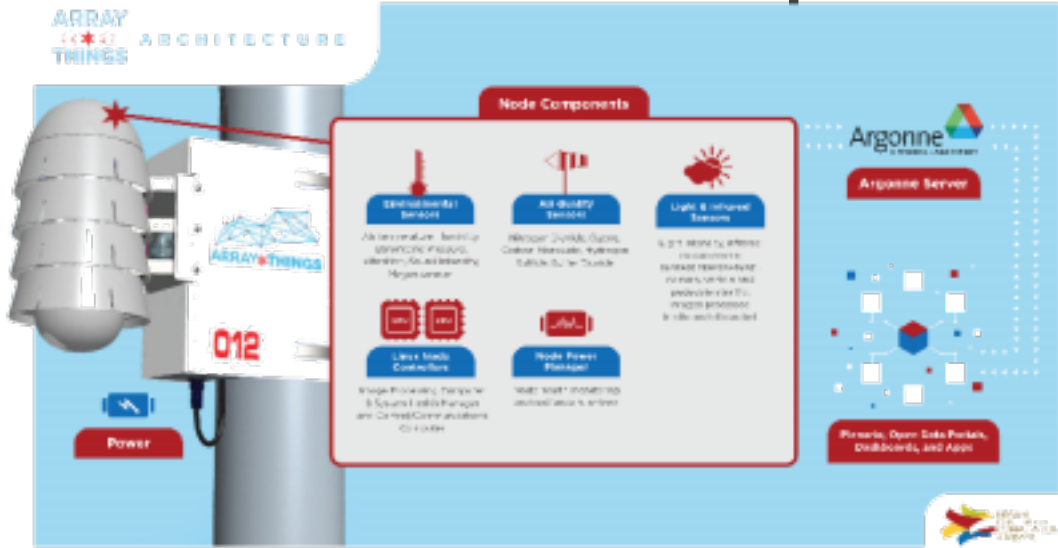
- Change management and strategic planning
- Community collaboration
- “Run forever” applications



Deep nets and complex models at the edge



Smart cities and personalized medicine



Charlie Catlett



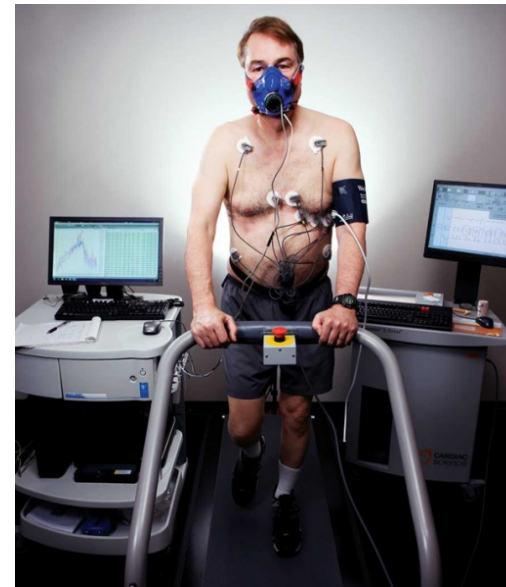
Connected devices + The Quantified Self

What problem are you solving for?

I want to control my own WELL BEING.



Oura Ring



Oxford Nanopore Mobile DNA sequencer



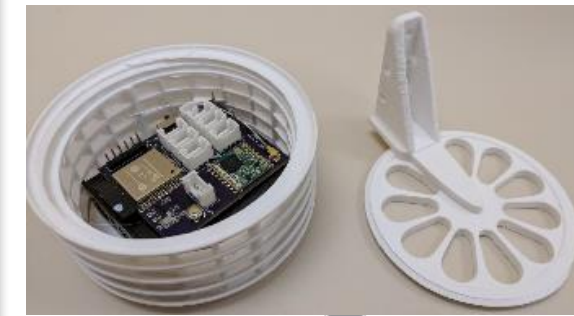
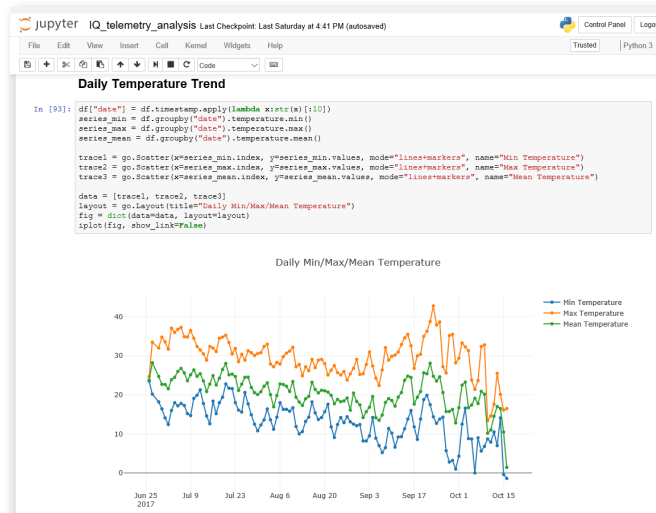
A framing question ...

How would you use

- Hundreds of ~\$50 wireless sensors?
- Streaming environmental data to understand watersheds and irrigation?



Solar powered Raspberry PI gateway



Custom Arduino Sensor



mongoDB



docker



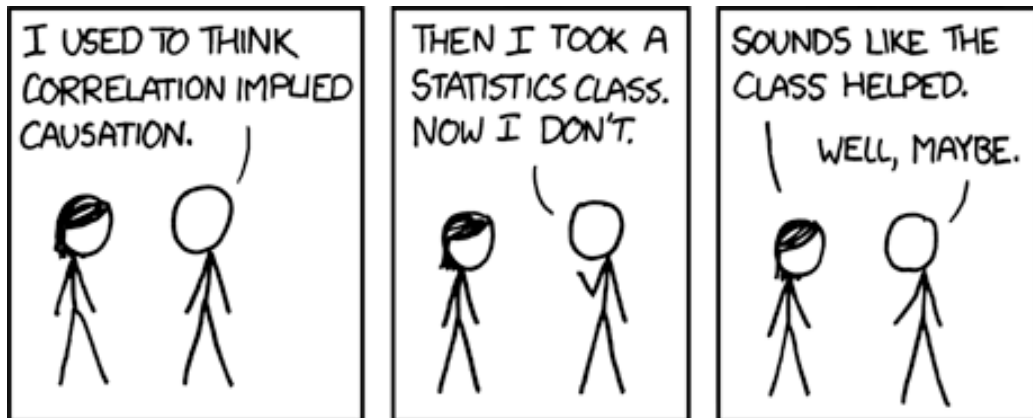
python™



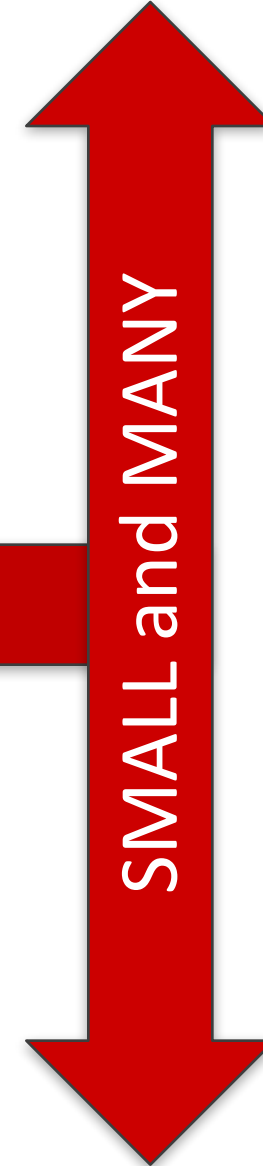
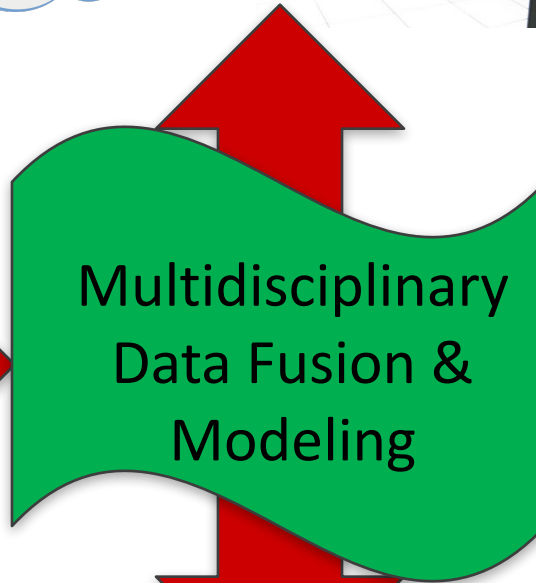
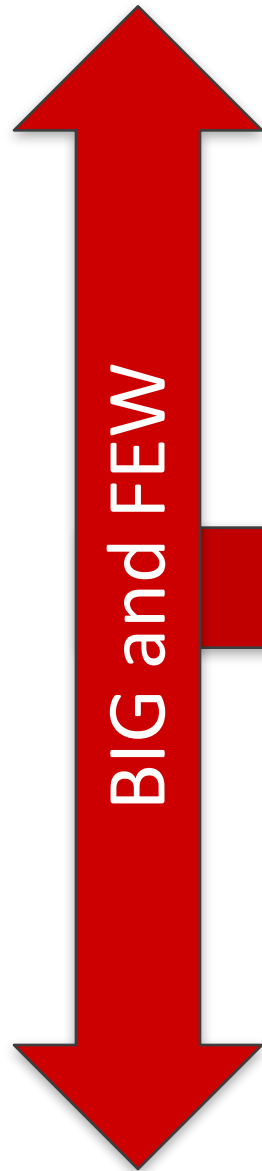
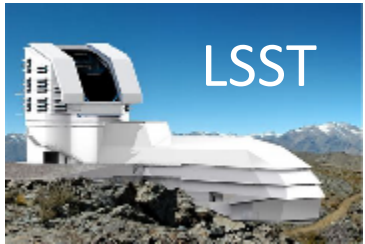
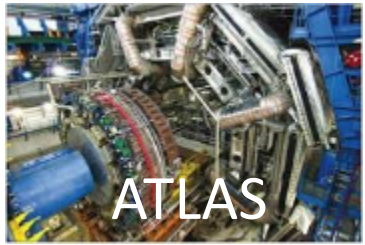
Big data: changing perceptions, shifting challenges

What information consumes is rather obvious: it consumes the attention of its recipients. Hence *a wealth of information creates a poverty of attention*, and a need to allocate that attention efficiently among the overabundance of information sources that might consume it.

Herbert Simon



Science instrument continuum



Vehicles



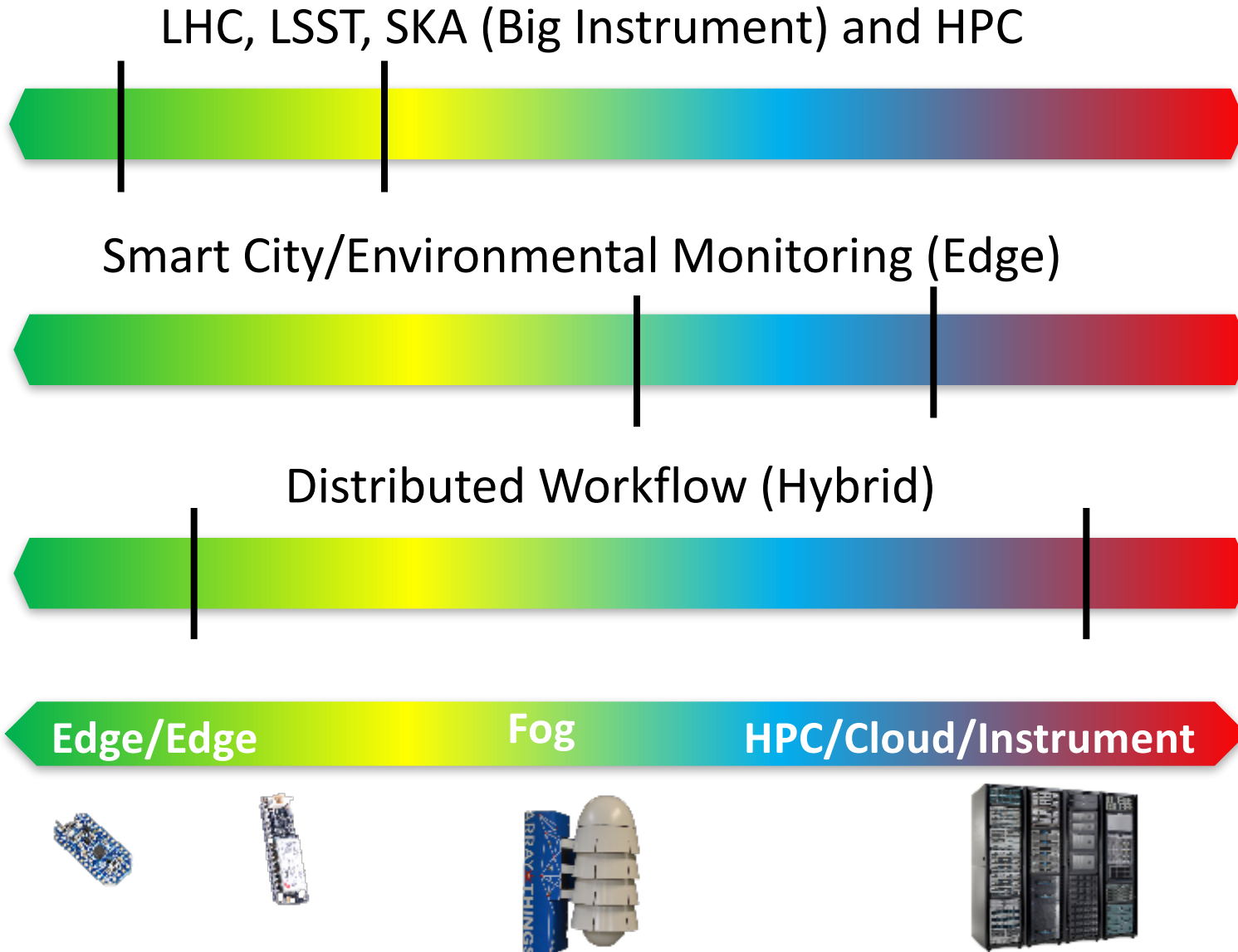
Array of
Things



Health



Building fluid capabilities: AND not OR

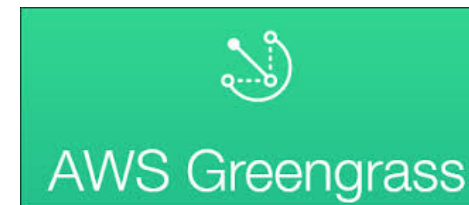


Mapping

- What
- Where
- When

Subject to

- Speed
- Capacity
- Latency
- Resilience
- Security



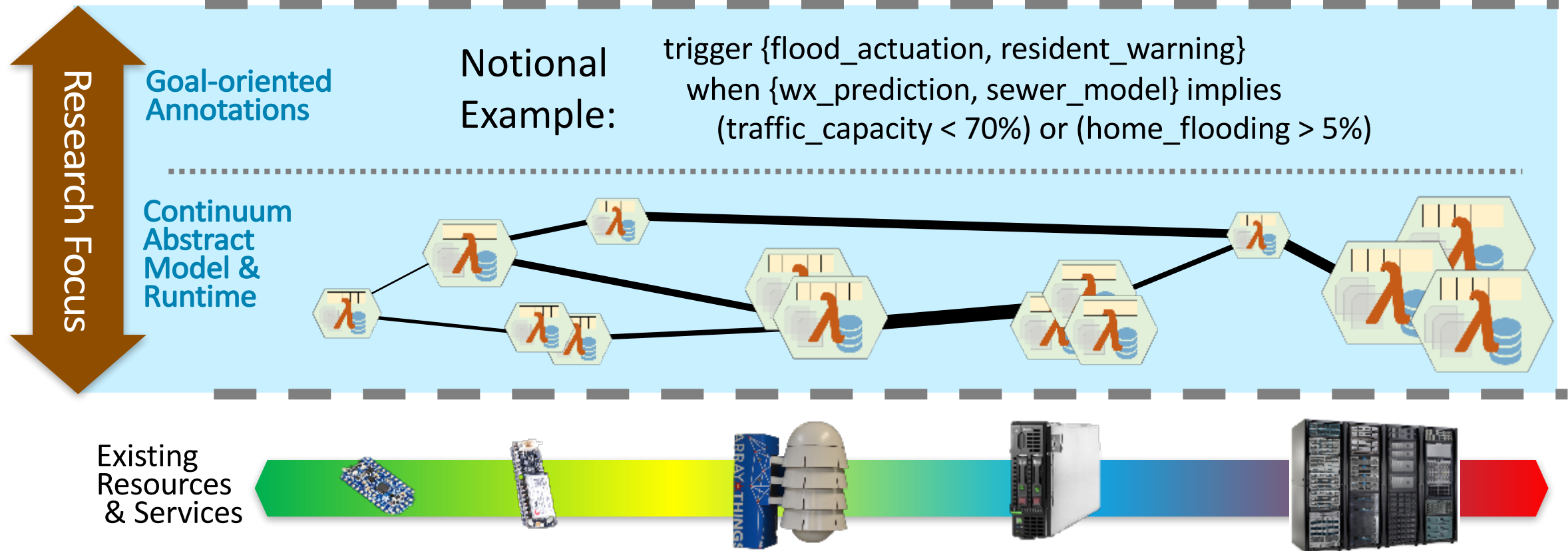
Harnessing the computing continuum

Intentional specification

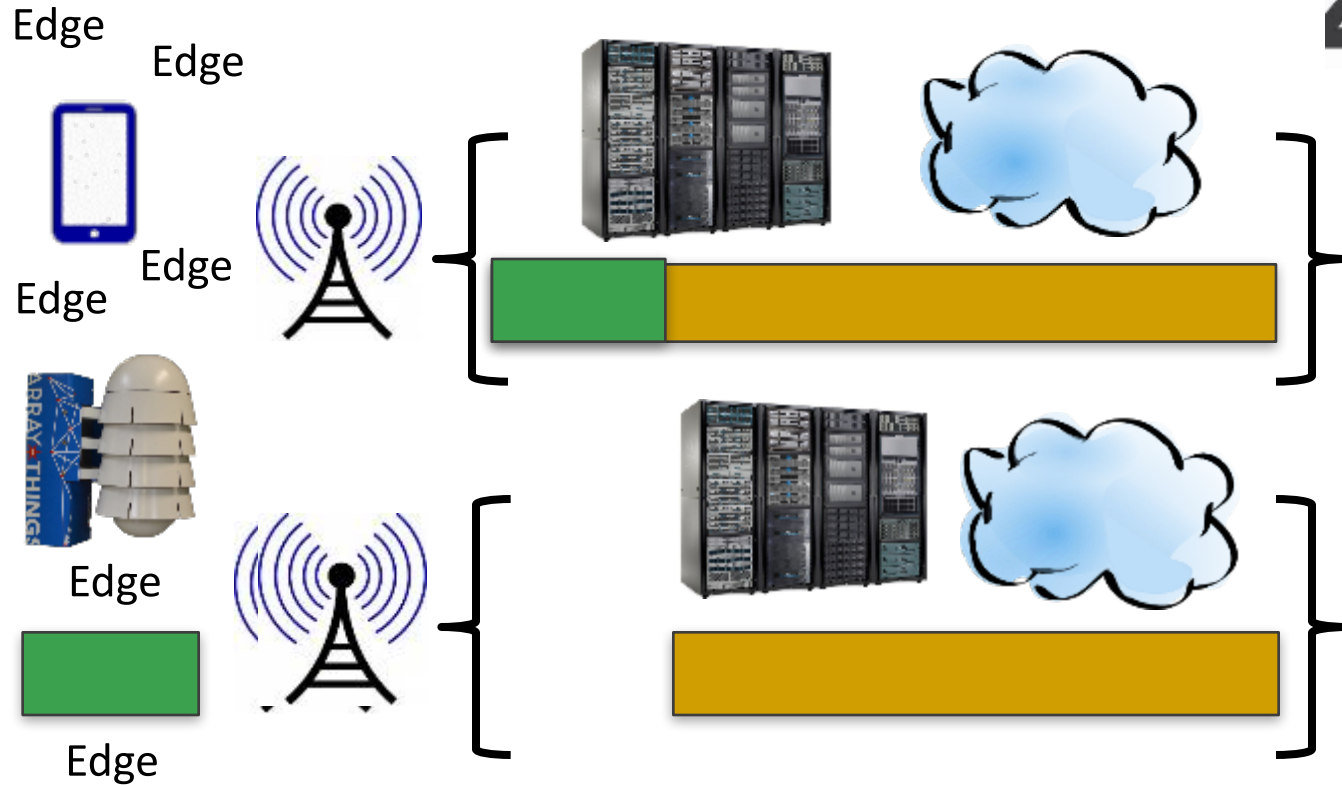
Science-driven Problems



e.g.: “Predict urban response to rainfall, trigger intelligent reaction...”



Building fluid capabilities



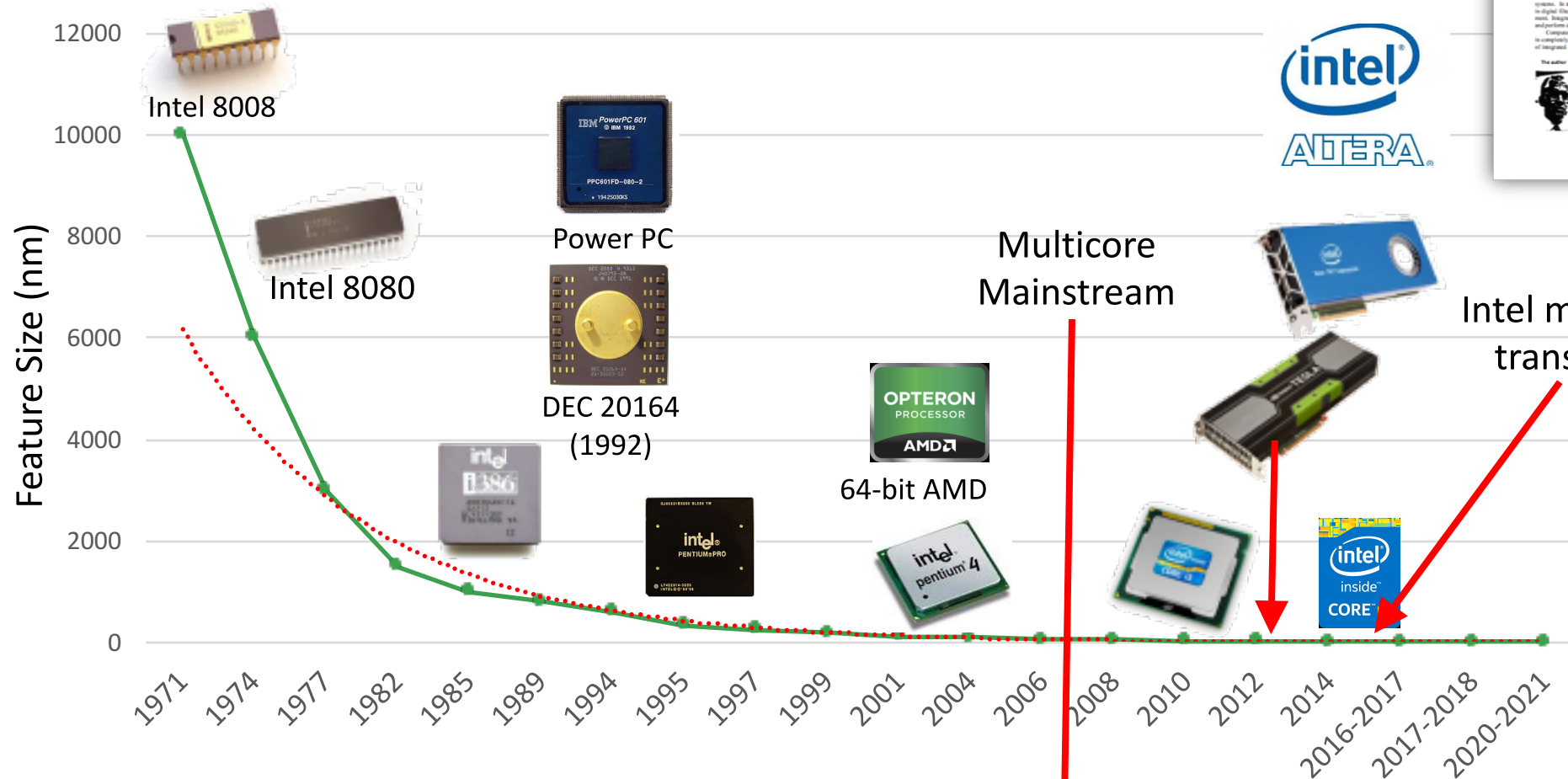
Latency
Bandwidth
Energy

Storage
Knowledge
Context



Chip feature sizes and Dennard scaling

No exponential is forever (except in the textbooks)



[https://science.energy.gov/~media/ascr/ascac/pdf/meetings/201812/ASCAC BRN Microelectronics.pdf](https://science.energy.gov/~media/ascr/ascac/pdf/meetings/201812/ASCAC_BRN_Microelectronics.pdf)
[https://science.energy.gov/~media/bes/pdf/reports/2018/Microelectronics Brochure.pdf](https://science.energy.gov/~media/bes/pdf/reports/2018/Microelectronics_Brochure.pdf)



Why us? Why now? (channeling Jim Larus)

Why us?

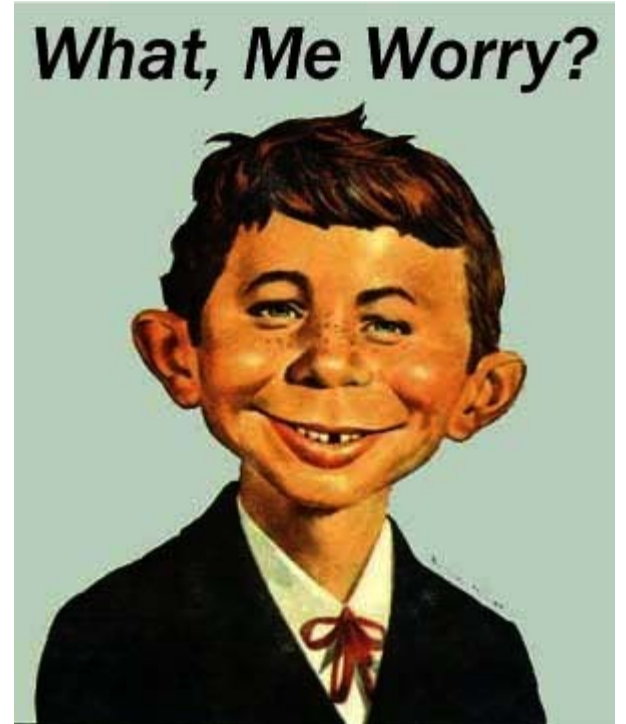
- Someone needs to think about higher level models

Why now?

- End-to-end, “run forever” services are the future
- Architectural specialization puts a premium on portability
 - Minimizing data movement at many levels
 - Maximizing operations/joule

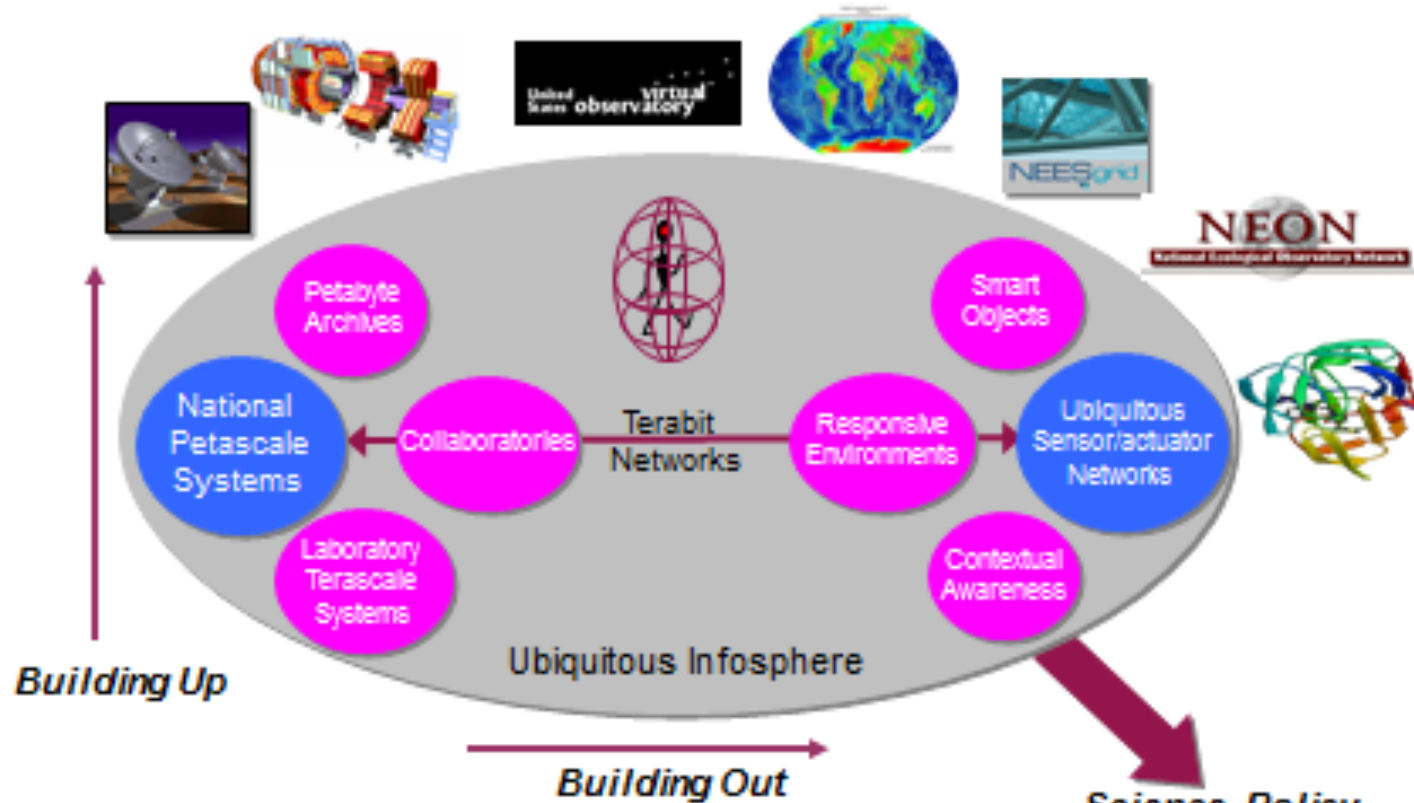
Why fusion?

- Integration will enable new capabilities
- It's more than workflows, containers, and libraries



Not yet fully realized ... my 2002 prediction

Futures: The Computing Continuum



The big questions don't change ...



Discussion

... but the approaches and answers do