# Coding the CONTINUUM

# Ian Foster

Argonne Training Program on Extreme-Scale Computing August 7, 2019



Now (2019)

#### Supercomputing

1.6 x 10<sup>8</sup> flop/s (Cray-1) 1.5 x 10<sup>17</sup> flop/s (Summit)



Now (2019)

#### Supercomputing

1.6 x 10<sup>8</sup> flop/s (Cray-1) 1.5 x 10<sup>17</sup> flop/s (Summit)

# Wide area networking104 bit/s (dial-up)1011 bit/s (optical)

3

"When the network is as fast as the computer's internal links, the machine disintegrates across the net into a set of specialpurpose appliances."

-- George Gilder, 2000



"When the network is as fast as the computer's internal links, the machine disintegrates across the net into a set of specialpurpose appliances."

-- George Gilder, 2000



### "network is as fast as the computer's internal links"



#### Innovation continues in the lab

Hollow core fiber: 99.7% speed of light (1.46x faster than fiber) 73.7 terabits per second





### But are we really free?

Time = 
$$T_{compute}$$
 + 2  $T_{latency}$ 



300 000 km per second\*. It's not just a good idea... It's the law!

\*207 000 km/sec in optical fiber: 5 x 10<sup>-6</sup> sec/km



"When the network is as fast as the computer's internal links, the machine disintegrates across the net into a set of specialpurpose appliances."

-- George Gilder, 2001



# "a set of special-purpose appliances"

More Flexible			Perf/W						
		CPUs 1X		Today's standard, most programmable, good for services changing rapidly		+			
neou		Manycore		Many simple cores (10s to 100s per chip), useful if	Conventional programming	C/C+			
noge			CPUs	3X	software can be fine-grain parallel, difficult to maintain				
НОН			GPUs	5-30X	Good for data parallelism by merged threads (SIMD), High memory bandwidth, power hungry		(UDA		
				Mast valies fully and menople settion. Cool for	Alternative	0			
		FPGAs	5-30X	Most radical fully programmable option. Good for	programming	log			
Specialized			J-20X	currently need to program in H/W languages.		Veri			
		Structured	20-100X	Lower-NRE ASICs with lower performance/efficiency.					
							ASICS		includes domain-specific (programmable) accelerators.
		Custom	1001	Highest efficiency. Highest NRE costs. Requires high	functionality	eril			
Mor	e	ASICs	> 100X	volume. Good for functions in very widespread use that	t	>			
Efficient		L]		are stable for many years. Source:	http://bit.ly/2SDC	GHzT			



# "a set of special-purpose appliances"

#### "Cloud computing 5x to 10x improved price point [relative to Enterprise]" — James Hamilton, http://bit.ly/2E78Wi1

#### Why?

- Improved utilization
- Economies of scale in operations
- More power efficient
- Optimized software



LBNL-1005775



# Zero-carbon cloud: Reduce energy cost and energy carbon footprint to 0



Andrew Chien DOI 10.1109/IPDPS.2016.96

Dr. Andrew A. Chien is the William Eckhardi Professor in Computer Science and Director of the CERES Center for Unstoppable Computing at the University of Chicago and a Senior Computer Scientist at Argonne National Laboratory. From 2005 to 2010, Chien was Vice President of Research of Intel Corporation, and from 1998 to 2005, the SAIC Chair Professor in Computer Science and founder of the Center for Networked Sustems at UCSD. His research has been recognized for excellence with numerous awards, and also supported by the NSF, DARPA, DOE, ONR, NASA, and industry. Dr. Chien currently serves as co-chair of the editorial board of the Communications of the ACM. From 1990 to 1998, he was a Professor of Computer Science at University of Illinois. Dr. Chien earned B.S., M.S., and Ph.D. degrees at the Massachusetts Institute of Technology, and is a Fellow of the ACM, IEEE, and AAAS.

Dr. Rich Wolski is a Professor of Computer Science at the University of California, Santa Barbara, and co-founder of Eucalyptic Systems Inc. Having received his M.S. and Ph.D. degrees from the University of California at Davis (while a research scientist at Lawrence Livermore National Laboratory) he has also held positions at the University of California, San Diego, and the University of California, San Diego, and the University of Tennessee, the San Diego Supercomputer Center, and Lawrence Berkeley National Laboratory. Dr. Wolski has led several national-scale research efforts in the area of distributed systems and is the progenitor of the Eucalyptus open source cloud project.

Fan Yang is currently a third-year Ph.D. student in the Department of Computer Science, University of Chicago. She joined the Large-Scale System Group at UChicago in 2013. She has broad research interests in parallel and distributed systems, data-intensive

computing, parallel programming and optimization, large-scale system software, and computer architecture. In the past year, she worked on graph computing and benchmarking of graph-processing systems. She

currently focuses on green datacenter design and stranded power in renexable electricity generation. She received her B E.R.g. in computer science from National University of Defense Technology, Changsha, Hunan

Province, China, in 2012. This work was supported in part by the Office of Advanced Scientific Computing Research, Office of Science, L.S. Department of Energy, under Award DE-SC0008080 and Contract DE-AC02-06C1H1357, as well as the National Science Foundation under Anoards CNS-1403959, STCI-0751315, CNS-

0905237, and CNS-1218808. The authors also gratefully acknowledge generous support from Hewlett-Packard, Keysight, Huawei, Nvidia, and the Semonur Goodman Foundation.

CrossMark

#### The Zero-Carbon Cloud: High-Value, Dispatchable Demand for Renewable Power Generators

Variability is an ongoing challenge to growth of largescale renewable power generation, posing challenges for the power grid and ambitious renewable portfolio standards. The authors propose Zero-Carbon Cloud (ZCCloud), a new high-value, dispatchable demand for renewables that improves their economic viability. Initial studies show that ZCCloud can create high-value computing resources with payback periods of just a few years.

Andrew A. Chien, Richard Wolski and Fan Yang

#### I. Rising Renewable Power Standards: The Stranded Power Opportunity

Generation costs for renewable power are an ongoing challenge to growing adoption of higher targets for renewable portfolio standards (RPS). Among leading RPS states, California produced 20 percent of its power from renewable sources in 2010 and has an ambitious target of 33 percent by 2020, and even Midwestern states such as Illinois<sup>1</sup> have adopted RPS targets that increase renewables in the electrical utilities fuel mix from 10 percent in 2015 to 25 percent in 2025. The United States federal government has established a goal of 20 percent by 2020 for all federal agencies, a dramatic increase

110 1040-6190/ © 2015 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license The Electricity Journal (http://creativecommons.org/licenses/by/4.0/)., http://dx.doi.org/10.1016/j.tej.2015.09.010

### The performance landscape becomes peculiar

A program can run on two computers On  $C_1$ , it takes 0.01 seconds On  $C_2$ , it takes 0.005 seconds Which is faster?



# The performance landscape becomes peculiar

A program can run on two computers On  $C_1$ , it takes 0.01 seconds On  $C_2$ , it takes 0.005 seconds Which is faster?

The answer depends on their location. Say C<sub>1</sub> is adjacent and C<sub>2</sub> is 500 km distant  $t(C_1) = T_1 = 0.01$  sec  $t(C_2) = T_2 + 2 \times 500 \times 5 \times 10^{-6} = 0.01$  sec



# The performance landscape becomes peculiar

A program can run on two computers On  $C_1$ , it takes 0.01 seconds On  $C_2$ , it takes 0.005 seconds Which is faster?

The answer depends on their location. Say C<sub>1</sub> is adjacent and C<sub>2</sub> is 500 km distant  $t(C_1) = T_1 = 0.01$  sec  $t(C_2) = T_2 + 2 \times 500 \times 5 \times 10^{-6} = 0.01$  sec



The apparent **speed** of a computer depends on its **location**; the apparent **location** of a computer depends on its **speed** 

### Continuum

A set of elements such that between any two of them there is a third element [dictionary.com]

#### For example, the **computing continuum**:

IoT/Edge	е		Fog			HPC/Cloud	
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
Memory	0.5K	256K	8GB	32GB	256G	32TB	16PB
Network	BLE	WiFi/LTE	WiFi/LTE	1 GigE	10GigE	40GigE	N*100GigE
Cost	\$5	\$30	\$600	\$3K	\$50K	\$2M	\$1000M
Count = 10 Size = 10 <sup>1</sup>	) <sup>9</sup>		ARRAY-THING				Count = 10 Size = 10

Credit: Pete Beckman, beckman@anl.gov

# The space-time continuum

"space by itself, and time by itself, are doomed to fade away into mere shadows, and only a kind of union of the two will preserve an independent reality ..."

H. Minkowski, 1908

aktion der Elektr Space-time diagram https://en.wikipedia.org/wiki/Spacetime

### The spacetime continuum in computational systems



#### The spacetime continuum in computational systems



Misquoting Minkowski: "Henceforth, **location** for itself, and **speed** for itself shall completely reduce to a mere shadow, and only some sort of union of the two shall preserve independence." <sub>21</sub>

### A real example: High energy physics trigger analysis



# Reasoning about the computing continuum (a) Assumptions

- A1: N consumers, distributed uniformly, X secs apart
- A2: Each consumer requests compute units at 1 Hz
- A3: Infinite bandwidth: i.e., only compute time and latency are of concern
- A4: An individual computer takes T secs to complete a compute unit
- A5: A compute center containing Z computers is faster by a factor of  $\sqrt{Z}$



# Reasoning about the computing continuum (a) Assumptions

- A1: N consumers, distributed uniformly, X secs apart
- A2: Each consumer requests compute units at 1 Hz
- A3: Infinite bandwidth: i.e., only compute time and latency are of concern
- A4: An individual computer takes T secs to complete a compute unit
- A5: A compute center containing Z computers is faster by a factor of  $\sqrt{Z}^{\,*}$



\* Grosch's law (1953): computer power rises by the square of the price

# Reasoning about the computing continuum (b) Without response time bounds

Max time is:



# Reasoning about the computing continuum (b) Without response time bounds

Max time is:



# Reasoning about the computing continuum (b) Without response time bounds

Max time is:



# Reasoning about the computing continuum (c) With response time bound, B

We want to know D for which:

 $\frac{1}{\sqrt{size}} + 2D \le B$ As size is  $\pi D^2/X^2$ , we want to solve:  $\frac{T}{\sqrt{\pi D^2/X^2}} + 2D = B$  From A1, there are  $\pi D^2/X^2$  consumers within distance D of a compute center



# Reasoning about the computing continuum (c) With response time bound, B

We want to know D for which:

 $\frac{1}{\sqrt{size}} + 2D \le B$ As size is  $\pi D^2/X^2$ , we want to solve:  $\frac{T}{\sqrt{\pi D^2/X^2}} + 2D = B$ 

With B=0.01, T=0.001, X=0.0001 sec: D = 0.004964 sec (~1000 km)

Then:

Size =  $\pi D^2 / X^2$  = 7854

Max processing time is

2 × 0.004964 + 0.001/√7854

= 0.01 seconds

From A1, there are  $\pi D^2/X^2$  consumers within distance D of a compute center



# Reasoning about the computing continuum (d) Discussion

The model emphasizes the importance of **distribution** ("disintegration") of function and **aggregation** of capability

The model can be improved:

- Empirical data on scaling of cost and speed with size
- Data transfer costs
- Empirical data on workloads

Optimal solutions will involve compute centers of multiple sizes  $\rightarrow$  Not just "center" and "edge"

# Small and midsize data centers: Server intensity



Code: verb.

1) to arrange or enter in a code

Code: verb.

1) to arrange or enter in a code

#### 2) to write code for

Code: verb.

1) to arrange or enter in a code

#### 2) to write code for

Now that the machine has disintegrated across the net, how do we program it?

Code: verb.

1) to arrange or enter in a code

2) to write code for

Now that the machine has disintegrated across the net, how do we program it?





### An example: Serial synchrotron crystallography

#### For each sample:

- Image crystals at ~50 Hz:
  - Validate each image
  - After 1000, quality control
  - After 26000, full analysis
- If good:
  - Determine crystal structure
  - Return crystal structure



### Coding the continuum: Serial crystallography





#### Similar needs arise across modern (AI-enabled) science



### Learned Function Accelerators (LFAs)



### Coding the continuum: Closed solution



https://read.acloud.guru/aws-greengrass-the-missing-manual-2ac8df2fbdf4

#### Thanks to colleagues, especially:





Yadu Babuji Ananthakrishnan











globus 🛆

**Kyle Chard** 



labs

Ryan Chard



Zhuozhao Li











Logan Ward



DLHub Automate Write 🥻 Parsl programs SCRTMP funcX Auth

In [1]:	<pre>from funcx_sdk.client</pre>	t import FuncXClient						
	<pre>fxc = FuncXClient()</pre>							
In [2]:	<pre>func = """ def add(data):     sum_val = sum(dat     return sum_val """</pre>	ta['data'])						
In [3]:	<pre>fxc.register_function("add_func", func, description="Sum a list of numbers.")</pre>							
In [4]:	]: input_data = [1, 2, 3]							
	<pre>res = fxc.run(input_</pre>	data, "user#laptop", "add_func")	)					
In [5]:	print(res)							
	6							
Ро	rtable code	Any access	Any compute					
	Python	SSH, Globus,	Clusters,					

cluster or HPC

scheduler

clouds, HPC,

accelerators

Docker, Shifter,

Singularity



# *func*X: Transform clouds, clusters, and supercomputers into high-performance function serving systems



# *func*X: Transform clouds, clusters, and supercomputers into high-performance function serving systems





### Common FaaS systems, compared

	Function Language	Intended Infrastructure	Virtualization	Triggers	Maximum Walltime (s)	Billing
Amazon Lambda	C#, Go, Java, Powershell, Ruby Python, Node.js	Public cloud, Edge (Greengrass)	Firecracker (KVM)	HTTP, AWS services	900	Requests, runtime, memory
Google Cloud Functions	BASH, Go, Node.js, Python	Public cloud	Undefined	HTTP, Pub/Sub, storage	540	Requests, runtime, memory
Azure Functions	BASH, Java, Python, Visual Studio	Public cloud, local	OS images	HTTP, APIM, MS services	600	Requests, runtime, SLA
OpenWhisk	Ballerina, Go, Java, Node.js, Python, Go	Kubernetes, Private cloud, Public cloud	Docker	HTTP, IBM Cloud OW-CLI	300	IBM Cloud: Requests, runtime Local: NA
Kubeless	Node.js, Python .NET, Ruby Ballerina, PHP	Kubernetes	Docker	HTTP, scheduled, Pub/Sub	Undefined	NA
SAND	C, Go, Java, Node.js, Python	Public cloud, Private cloud	Docker	HTTP, Internal event	Undefined	Triggers
Fn	Go, Java, Ruby, Node.js, Python	Public cloud, Kubernetes	Docker	HTTP, direct trigger	300	NA
Abaco	Container	TACC clusters	Docker	НТТР	Undefined	Undefined
funcX	Python	HPC, Public cloud	Singularity, Shifter, Docker	HTTP, Globus Automate	No limit	HPC SUs

#### Incremental construction of a personalized cost map

- Build black-box performance models from observed execution times for different codes on different platforms
- Transfer learning across codes, problem sizes, and hardware platforms
- Experiment design to choose experiments that maximize reduction in uncertainty
- Evolve models over time as codes and platforms change
- Use models for instance selection and scheduling



#### Example: A cost map for bioinformatics applications on different AWS instance types IndexBam performs better on compute-optimized instances. Poorly



On average, within 30% of final error after 4 experiments and within 2.3% after 6



Detect and respond to events

→ E.g., in HPC file systems: FSMon (Arnab Paul et al.) Invoke RESTful services, and accept user input Manage short- and long-lived activities



# Flow automation in a neuroanatomy automation





Cloud-hosted services support data lifecycle events
 → Cloud for high-reliability, modest-latency actions
 → Integrated OAuth-based security with delegation



DLHub		-0 >		ζ → Automate
Data and Learning Hub for Science				
A simple way to find, share, publish, and run m learning models and discover training data for s	achine	<b>○</b>	Write programs	Parsl
Python SDK CLI O Contribute			Cost	CODIMD
Get Started			map	SCRIMP
1	2	3	Function fabric	funcX
Describe 🗐	Publish	Run 📆		
<pre>m = KerasModel() m.create_model("p1b1-example.h5") m.set_title("CANDLE Pilot 1 - Benchmark 1") m.set_name("candle_p1b1") # short name m.set_domains("genomics","biology","HPC")</pre>	<pre>from dlhub_sdk.client import DLHubClient dl = DLHubClient() dl.publish_servable(m)</pre>	<pre>from dlhub_sdk.client import DLHubClient dl = DLHubClient() mid = dl.get_id_by_name("candle_p1b1") data = np.load("pilot1.npy") pred = dl.run(mid, {'data': [data.tolist()]})</pre>	Data fabric	Data services

Model

registry

DLHub

Aut

https://arxiv.org/abs/1811.11213 Paper @ Session 7, 1:30pm today

dlhub.org

# Coding the Continuum: Thanks for support



US Department of Energy



**US National Science Foundation** 



US National Institutes of Health





Amazon Web Services



Globus subscribers

"the machine disintegrates across the net into a set of special-purpose appliances"

#### Code: verb:

#### 1) to arrange or enter in a code

"Henceforth, **location** for itself, and **speed** for itself shall completely reduce to a mere shadow, and only some sort of union of the two shall preserve independence."

#### 2) to write code for

Distribute computational tasks across a heterogeneous computing fabric





# Coding the [locationspeed] continuum