

Tutorial: *Scalable Infrastructures for Compute- Intensive Cognitive Neuroscience*

Workshop Materials:

tinyurl.com/CogSci-ScalableComputing

Please complete the Pre-Workshop Survey!

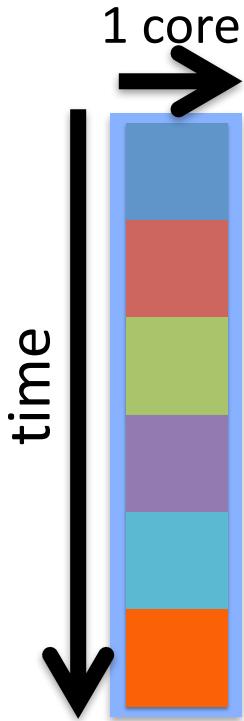
What is “large-scale” computing?

larger than a single computer

Computer Component:	Disk	Memory	Cores
Role of Component:	“long-term memory”	quick-access “working memory”	“math processing”
When You’ve Outgrown:	I don’t have room for all of the data.	The computer crashes.	It takes too long.

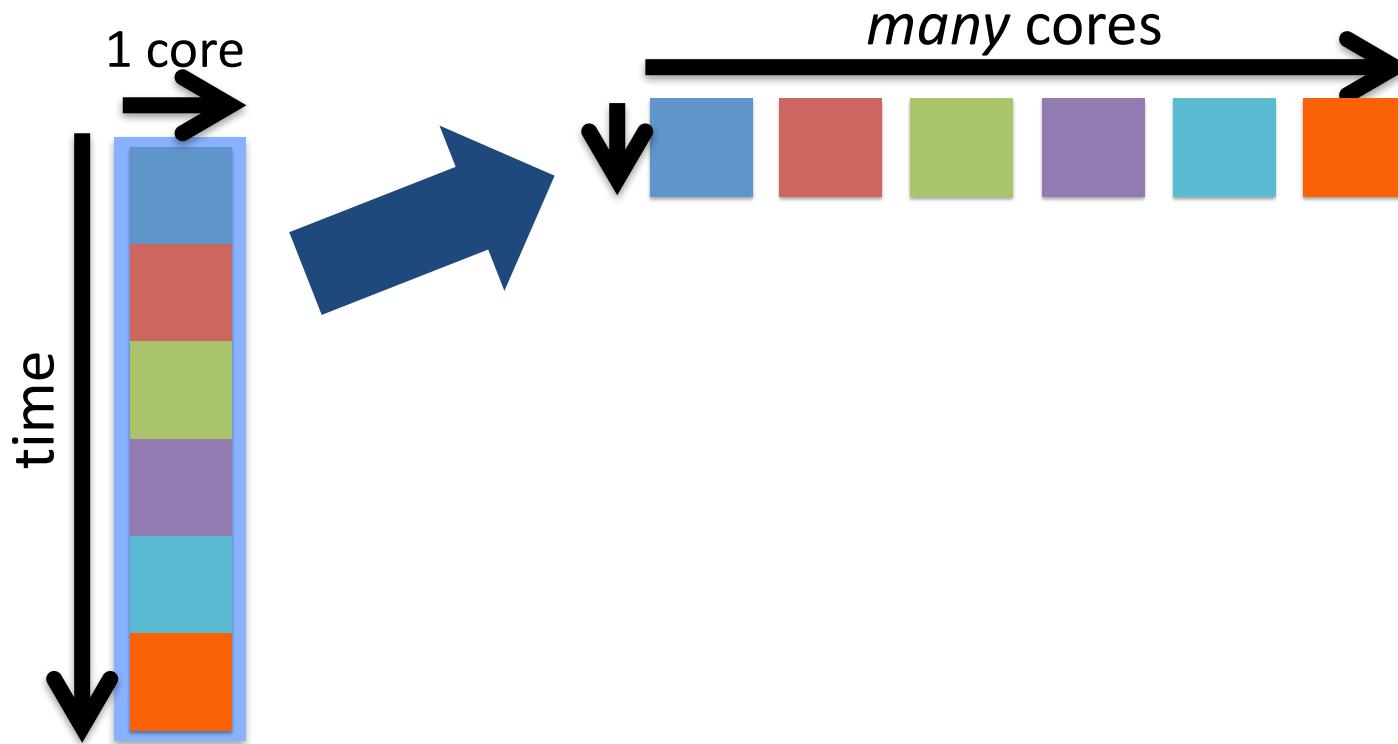
“It takes too long ...”

parallel computing approaches



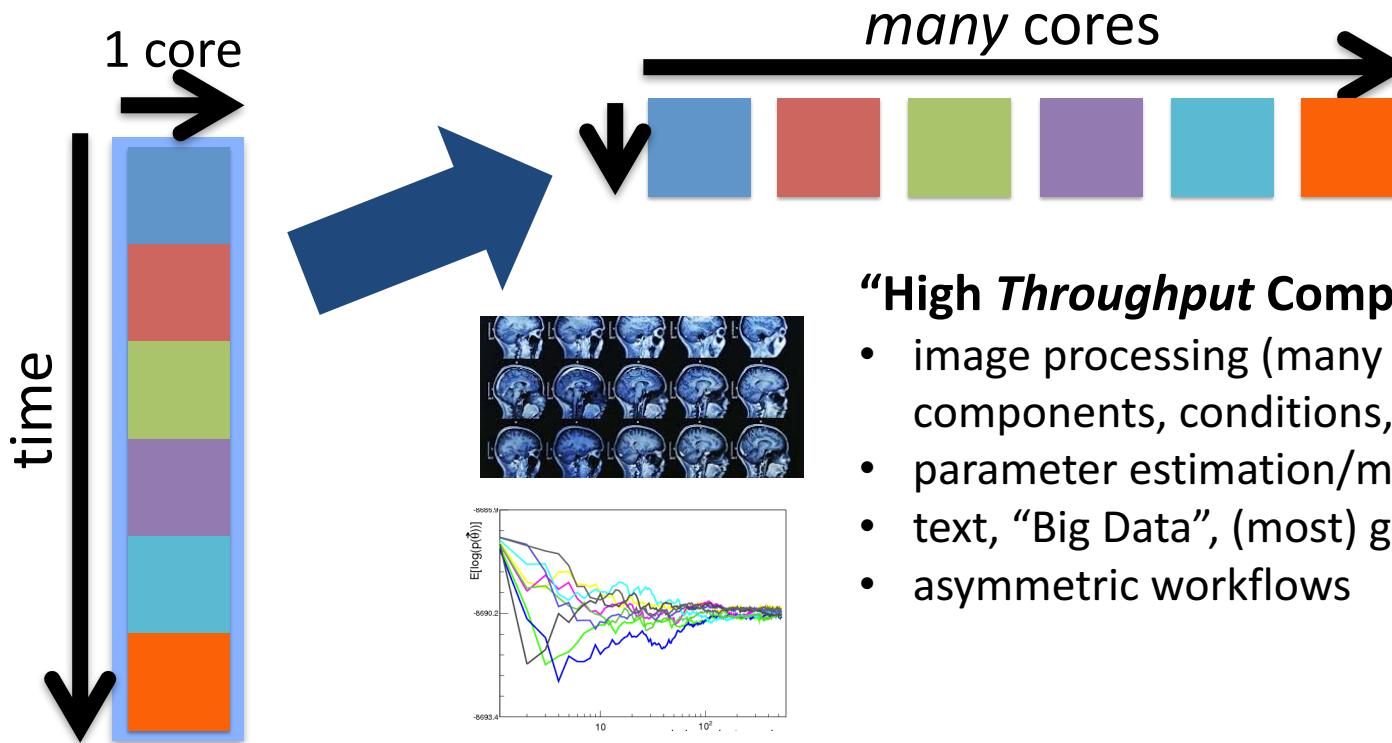
If order doesn't matter...

natural parallelization



If order doesn't matter...

natural parallelization

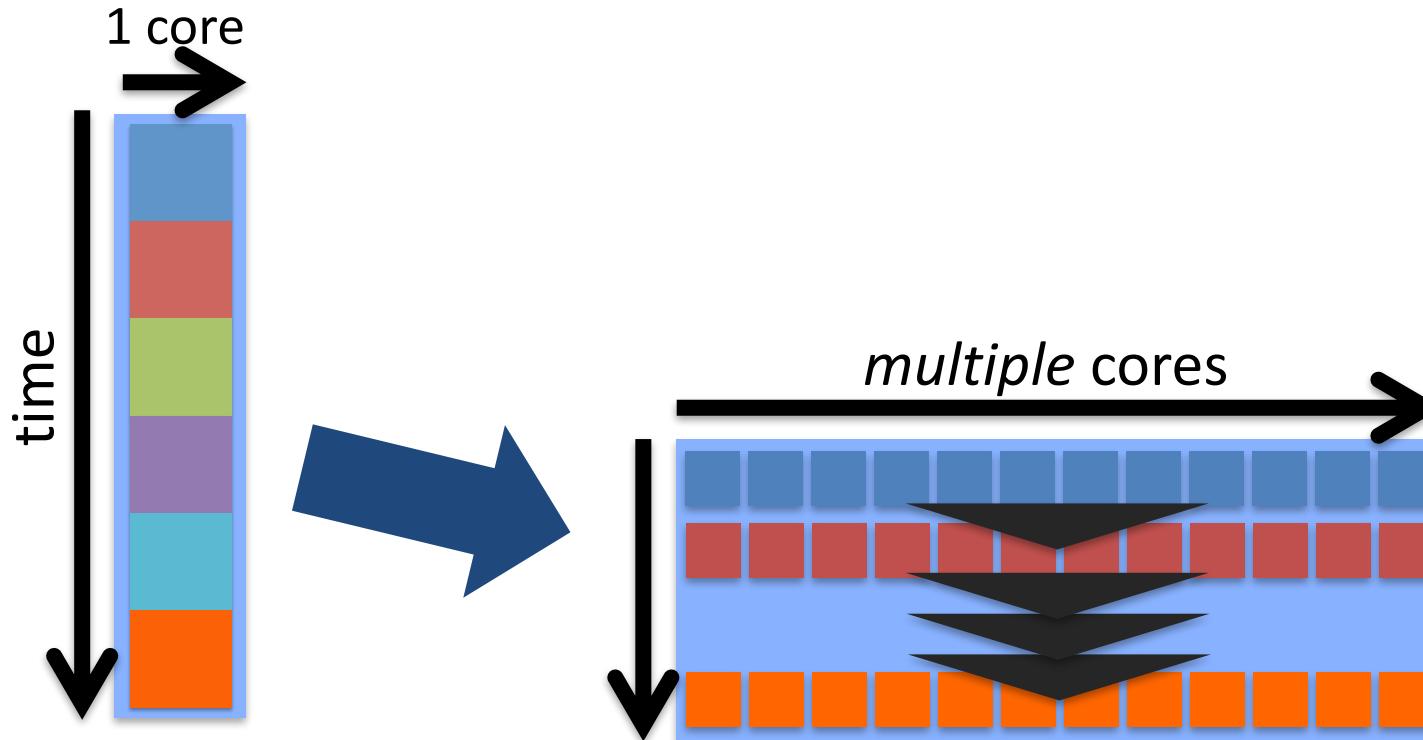


"High *Throughput* Computing (HTC)"

- image processing (many subjects, components, conditions, time frames)
- parameter estimation/model fitting
- text, "Big Data", (most) genomics
- asymmetric workflows

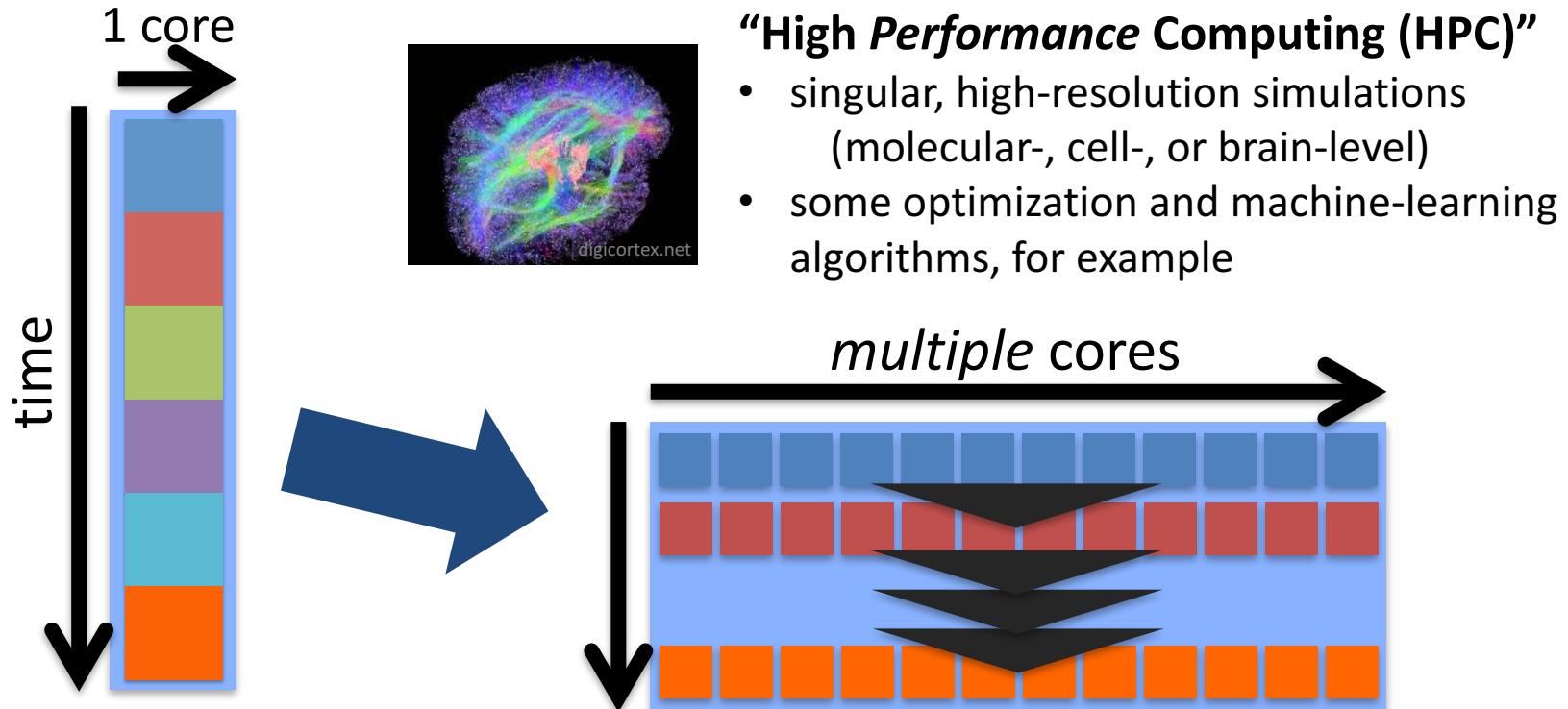
If order matters...

internal/interdependent parallelization



If order matters...

internal/interdependent parallelization



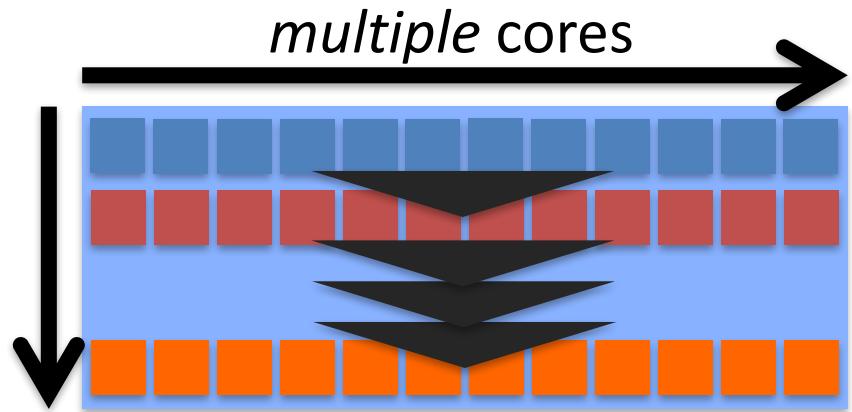
What is “large-scale” computing?

larger than a single computer

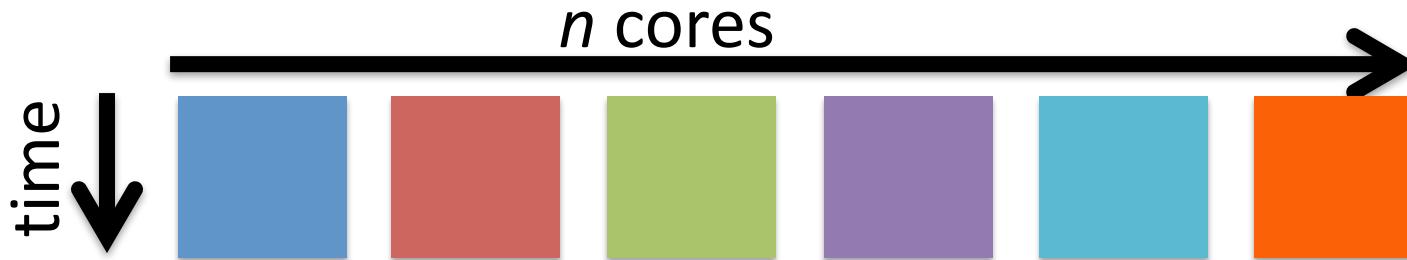
Computer Component:	Disk	Memory	Cores
Role of Component:	“long-term memory”	quick-access “working memory”	“math processing”
When You’ve Outgrown:	I don’t have room for all of the data.	The computer crashes.	It takes too long.

High Performance Computing (HPC)

- Benefits greatly from:
 - CPU speed + homogeneity
 - Shared filesystems
 - Fast, expensive networking and servers co-located
- Scheduling: **Must wait until all processors are available, at the same time and for the full duration**
- Requires special programming (MP/MPI)
- ***What happens if one core or server fails or runs slower than the others?***

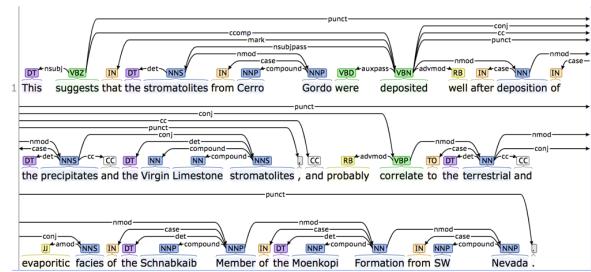


High Throughput Computing (HTC)

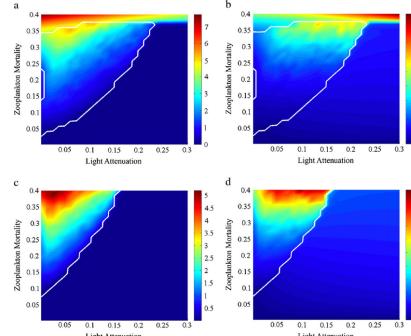


- Scheduling: only need **1 CPU core for each** (shorter wait)
- Easier recovery from failure
- No special programming required
- Number of concurrently running jobs is *more* important
- CPU speed and homogeneity are *less* important

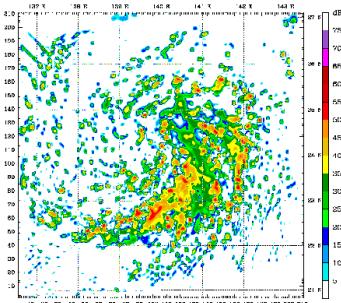
HTC-able Research Examples



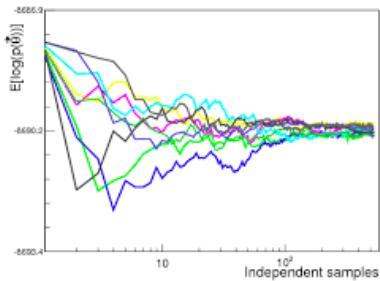
text analysis (most genomics ...)



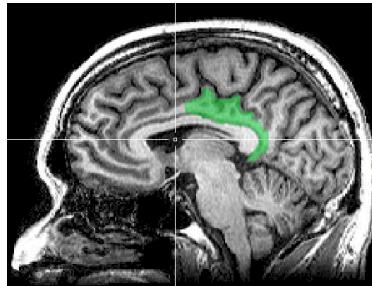
parameter sweeps



multi-start simulations



statistical model optimization
(MCMC, numerical methods, etc.)



multi-image and
mult-sample analysis

Scope of Computing Infrastructures

Accessibility and Scale

- per research group, per institution, national, or commercial
- free, by-proposal, or fee-based

Supported computing modes

- HTC- vs. HPC-optimized; availability of specialized hardware: GPUs, high-memory, etc.

Interface graphical/pipeline (specialized) or command-line (generalizable)

Human Support

- local vs. remote personnel; domain-specificity?
- learning materials



National-Scale Resources

XSEDE (HPC)

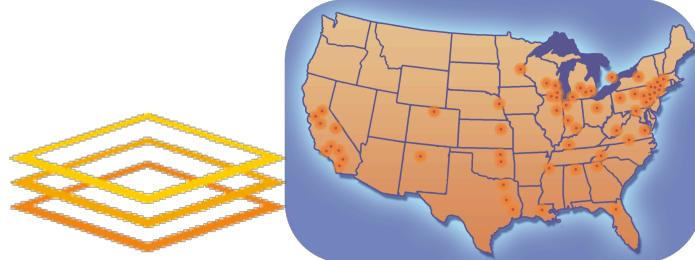
- multiple, large HPC-style clusters
- proposal-based allocations

XSEDE

Extreme Science and Engineering
Discovery Environment

Open Science Grid (HTC)

- extreme-scale HTC sharing
- **OSG Connect:** free for academic, govt., and non-profit researchers



Open Science Grid

Other HPC Clusters via Funding or Collaboration:



Commercial Cloud Providers



Google
Cloud Platform



Computation

- run your own HPC or HTC (more frequently) system
- some ready-to-use tools (map-reduce, machine-learning, visualization, workflows, genomics, etc.)
- some free, proposal-based allocations via NIH, NSF, etc.

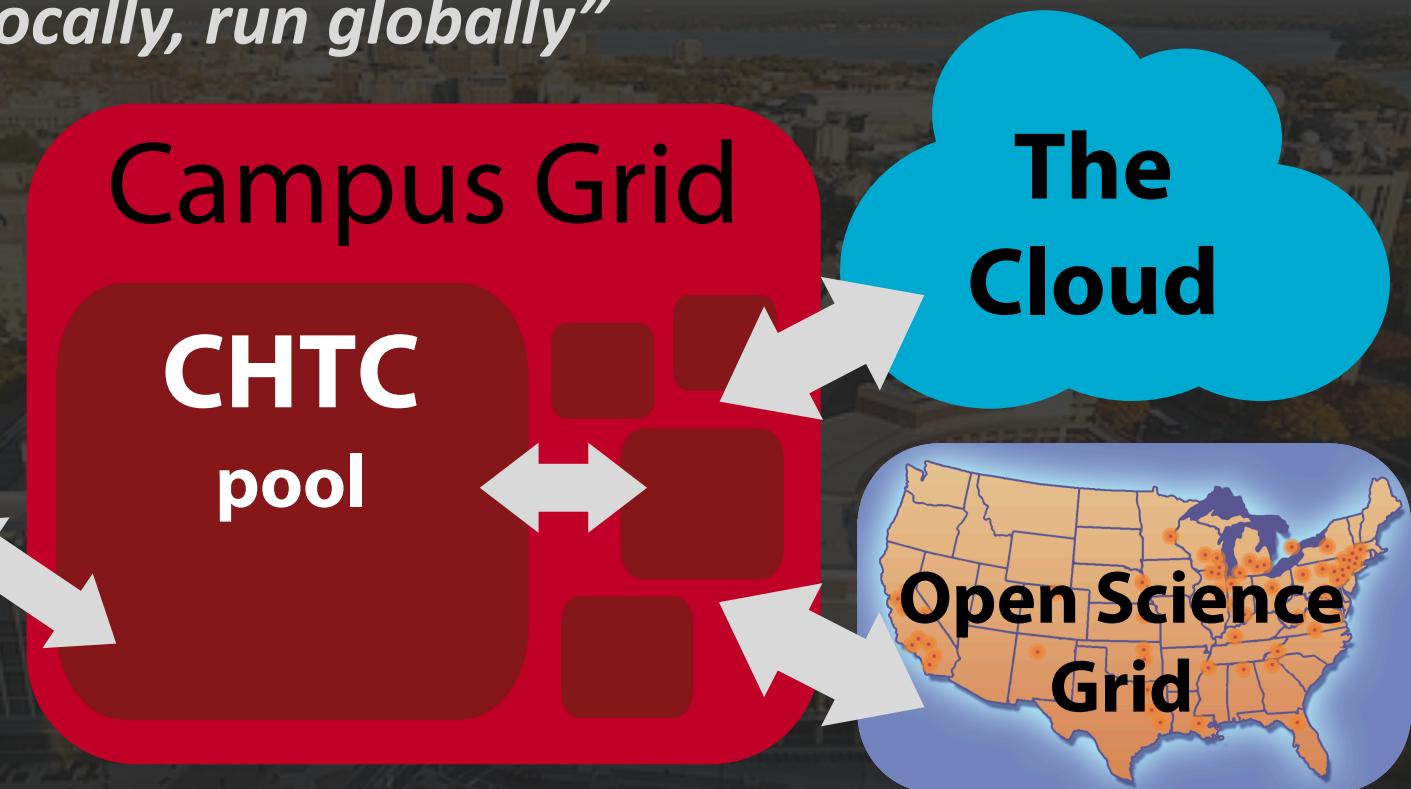
Data Management

- GoogleDrive, Box, databases, etc.
- Large-scale data transfer



Agile, Shared Compute Systems

“submit locally, run globally”



Science Gateways

Easy-to-use, Method-/Domain-Specific Computing

The screenshot shows the homepage of the Neuroscience Gateway (NSG) portal. At the top, there is a browser header with navigation icons, a URL bar showing <https://www.nsgportal.org>, and a set of small icons. Below the header is the NSG logo, which consists of a stylized green 'NSG' monogram next to the text 'NEUROSCIENCE GATEWAY' and 'A Portal for Computational Neuroscience'. To the right of the logo are a search bar and a menu icon. The main content area features a background image of several blue neurons on a black background. On the left side of this area, the heading 'Available Tools' is displayed above a list of software names: NEURON, GENESIS3, MOOSE, NEST, PyNN, Brian, Freesurfer, The Virtual Brain Pipeline, and BluePyOpt.

Available Tools

- NEURON
- GENESIS3
- MOOSE
- NEST
- PyNN
- Brian
- Freesurfer
- The Virtual Brain Pipeline
- BluePyOpt

[Register Account »](#)

[Access NSG Portal »](#)

[Join Mailing List »](#)

Scope of Computing Infrastructures

Accessibility and Scale

- per research group, per institution, national, or commercial
- free, by-proposal, or fee-based

Supported computing modes

- HTC- vs. HPC-optimized; availability of specialized hardware: GPUs, high-memory, etc.

Interface graphical/pipeline (specialized) or command-line (generalizable)

Human Support

- local vs. remote personnel; domain-specificity?
- learning materials



Research Computing Facilitators

accelerating research transformations

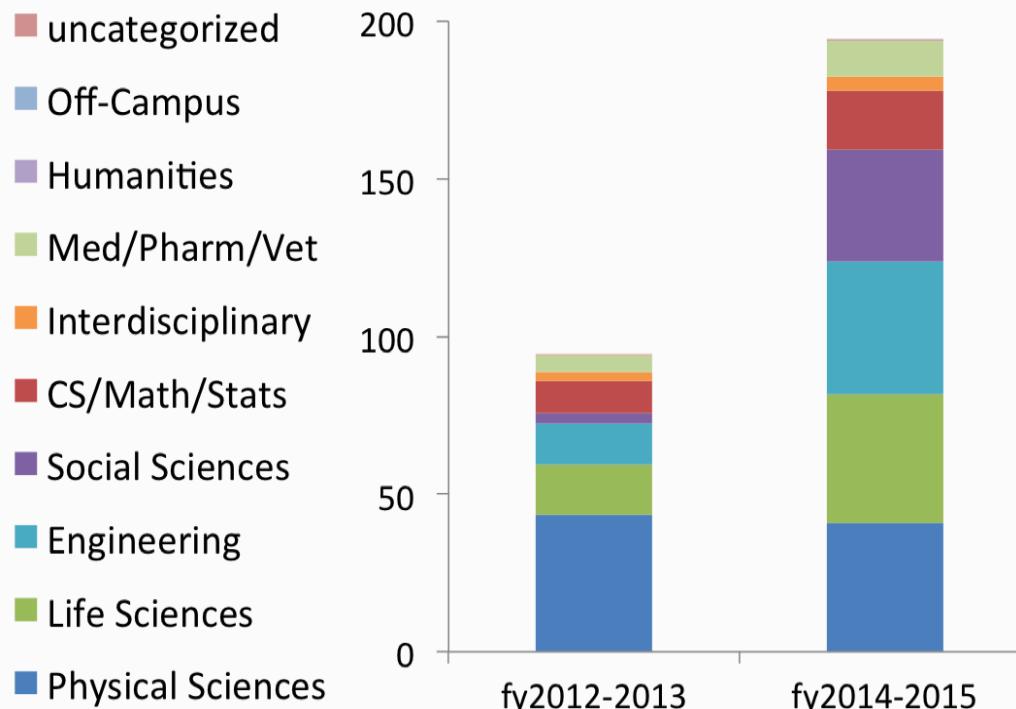
proactive engagement
personalized guidance
teach-to-fish training
technology agnostic
collaboration liaising
upward advocacy





Impact Across Domains

Millions of CPU Hours via CHTC

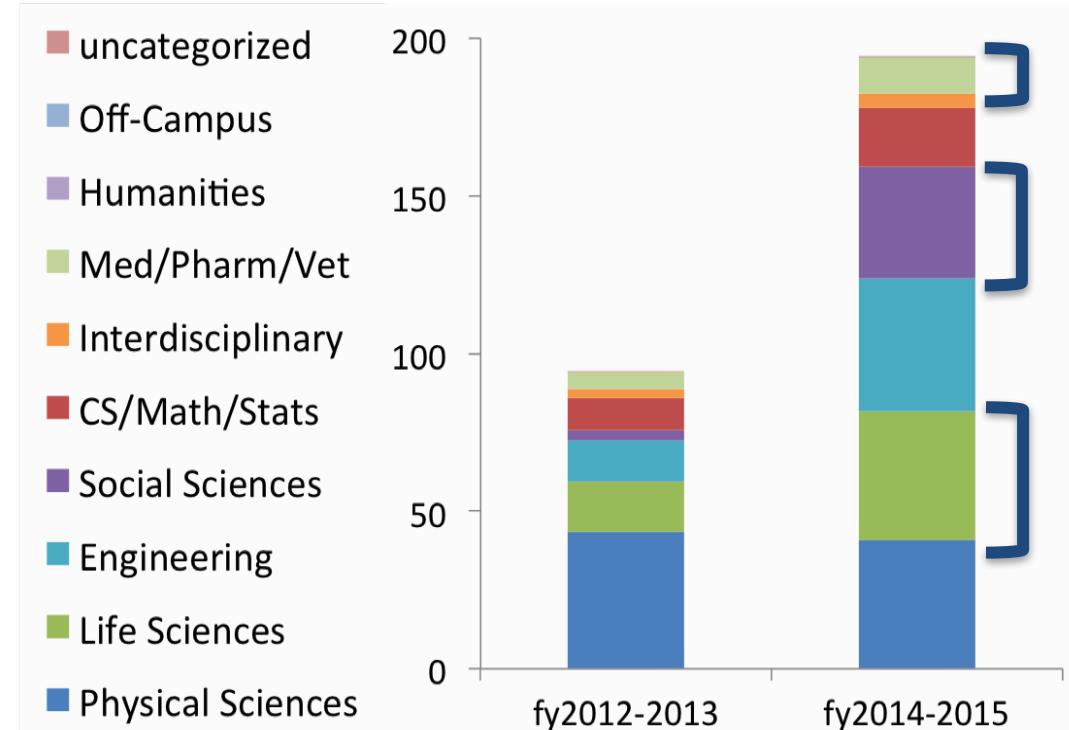


Facilitators hired:
Jan 2013, Nov 2014



Impact Across Domains

Millions of CPU Hours via CHTC



>95% high
throughput
computing



Impact Across Domains

Future Directions for
**NSF ADVANCED
COMPUTING
INFRASTRUCTURE**
to Support U.S. Science
and Engineering
in 2017–2020

The National Academies of
SCIENCES • ENGINEERING • MEDICINE

*“Increased advanced computing capability has historically enabled new science, and many fields today increasingly rely on **high-throughput computing** for discovery. Modeling and simulation, the historical focus of **high-performance computing**, is a well-established peer of theory and experiment. Data-driven research, a complementary “fourth paradigm” for scientific discovery, needs data-intensive computing capabilities and resources.” - p. S1 (Summary)*

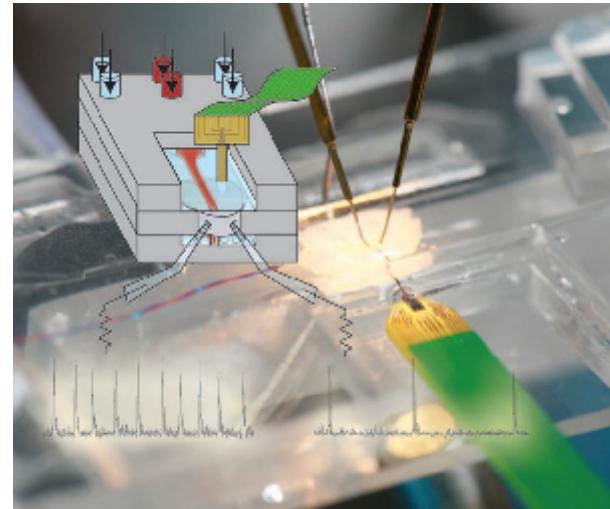
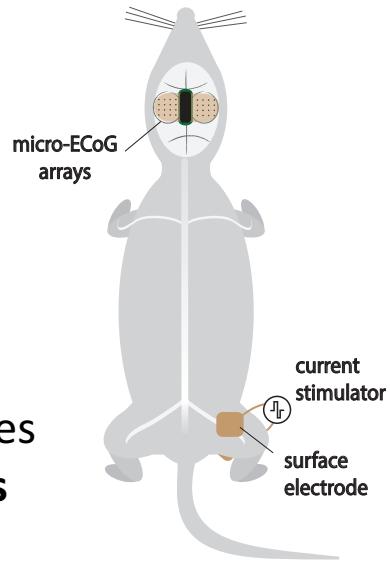
“For data-volume-driven workflows, scientific outcomes are best achieved when the advanced computing is configured for efficient, high-throughput processing, at scale ...” - p. 5-3



Modeling Sensory Networks for the Design of Brain-Computer Interfaces

Surgically-implanted arrays record potentials from the sensory cortex for estimation of the neural connection between regions.

Ricardo Pizarro used a *Discriminatory Category Matching* (DCM) algorithm to calculate network likelihood and features for **10,000 independent parameter sets** across **hundreds of experiments**.



4.3 million HTC hours in one year, ~1 millennium



Elucidating Cognition Pathways

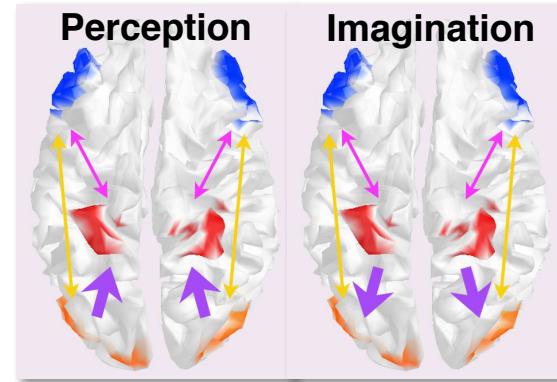
Expectation-maximization algorithm
Barry Van Veen predicts neural connectivity from EEG data.

Used by numerous on-campus clinical projects examining:

- **short-term memory**
- **imagination vs perception**
- **sleep versus waking states**
- and others



EEG data collection



Brain Regions and Connectivity

Per subject, per condition, per time point:
dozens trials X 20,000 Monte Carlo iterations

13 million CPU hours in one year, ~1.5 millennia

See Also:

Saturday Symposium, 3pm:

Data-intensive brain imaging: The state of the art

also presenting insights from...

NeuroComp17 Workshop, Aug 2017

<https://conferences.discovery.wisc.edu/neurocomp17/>

*What scalable computing modes and infrastructures are
most important to major challenges in neuroscience?*