



Modernizing Data Workflows from the Research Lab to the Data Center

Scott Yockel

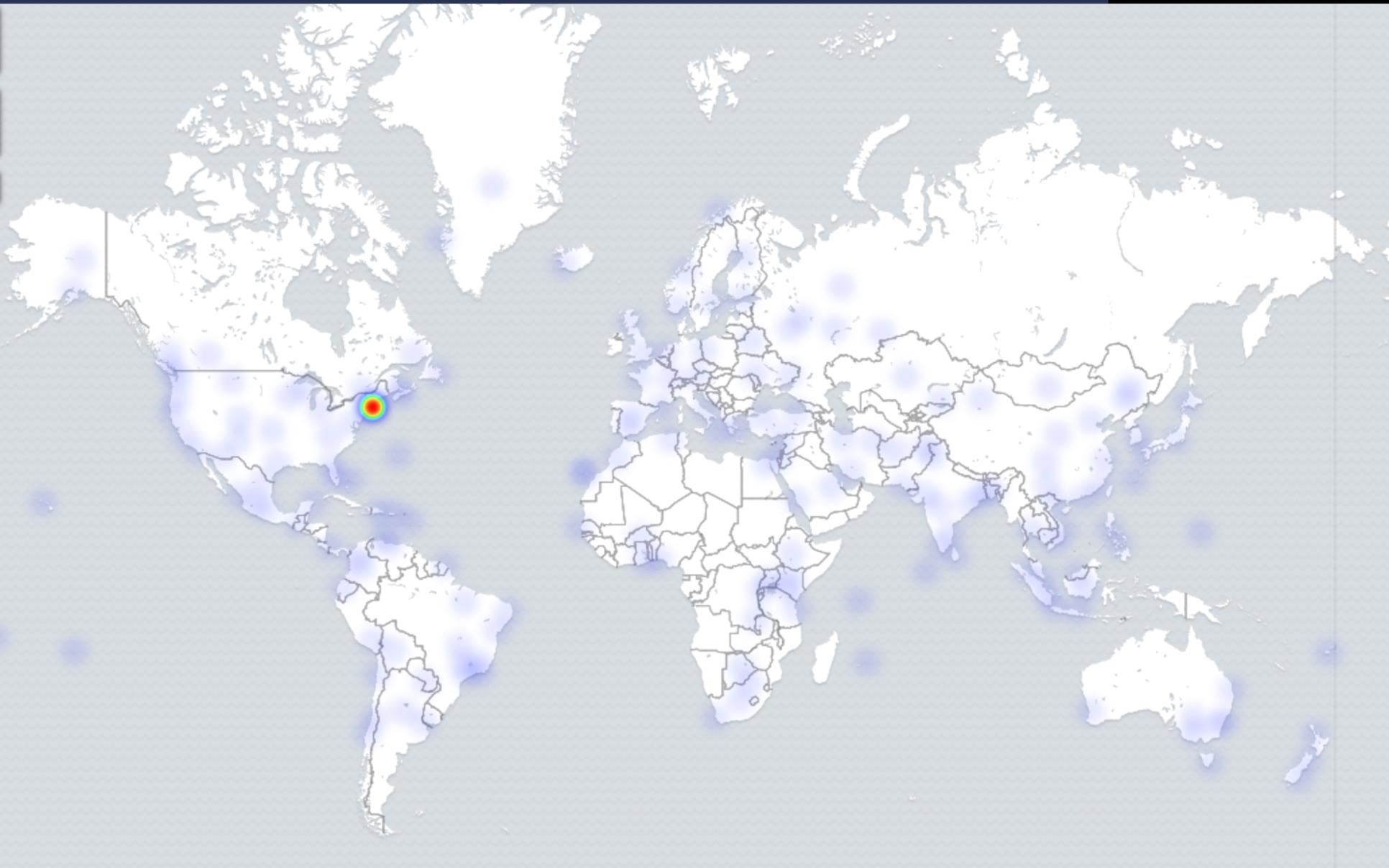
Director of Research Computing

Harvard University



Research Computing Resources

- 79k cores HPC
- Over 2 Petaflops GPUs (NVidia)
- 38.0PB of storage (Isilon, gluster, Lustre, XFS/NFS, Ceph, GPFS)
- 400+ virtual machines (KVM/OpenNebula)
- 2MW of research computing equipment in 3 data centers
- 23 FTE in 4 groups
 - ARCS: Advanced Research Computing Support
 - Software as Infrastructure (OpenNebula/VMs, Containers)
 - Systems Engineering & Data Center Operations
 - **Research Software Engineering (POSITION OPEN)**
- Supporting 600 research groups and 3500+ users





Globus Usage This Year

62	Unique Users
1119	Transfers
99,537,224	Files
574 TB	Data transferred



Outline

- Where does all this data come from?
 - Traditional HPC : Numerically Intensive
 - External Repositories : Data Intensive
 - Modern Instruments : Data + Numerically Intensive
- How do we deal with it?
 - Traditional Globus Transfer
 - Making use of User Plus Personal Endpoints



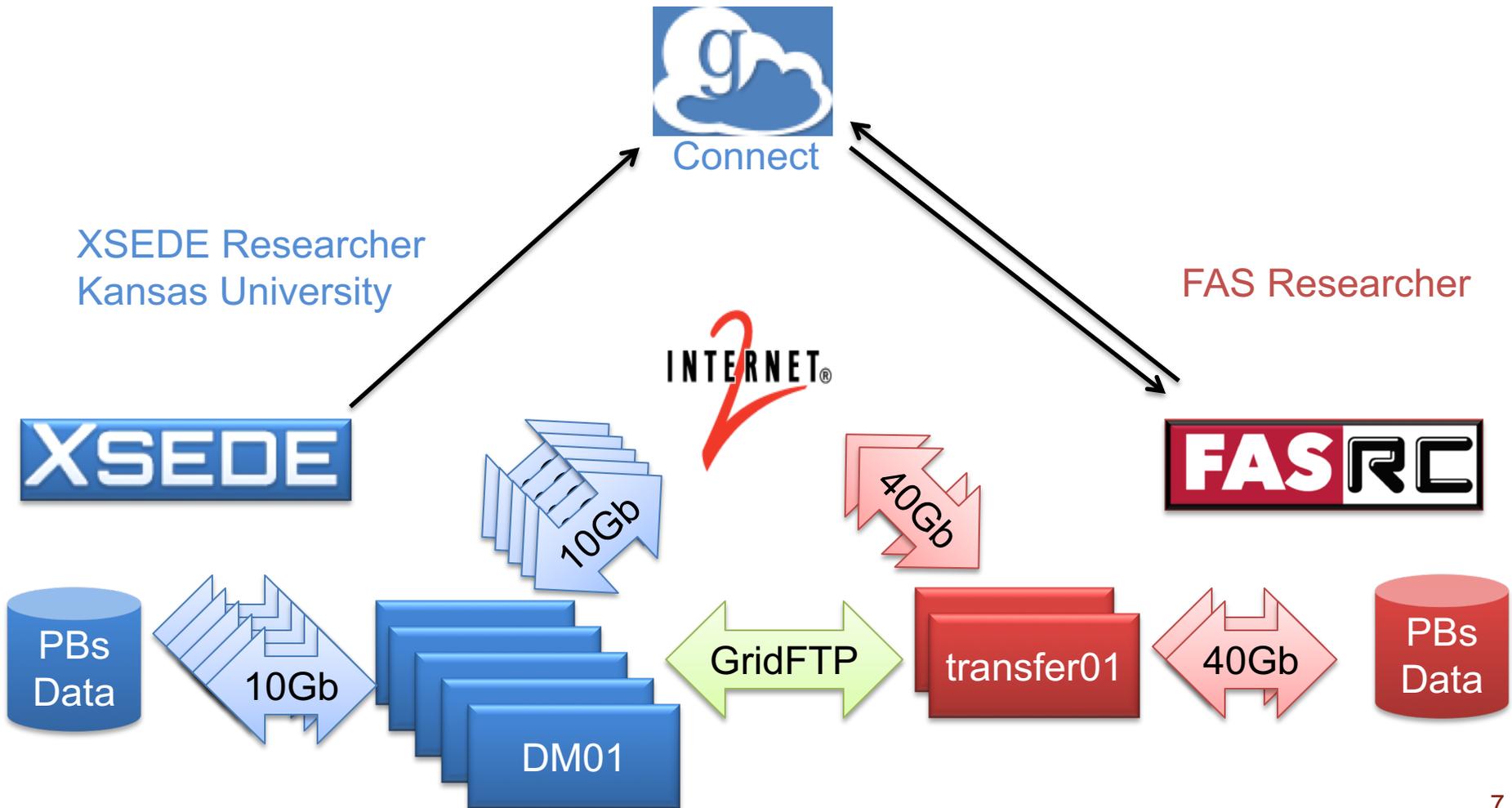
One Researcher's Ask

Also, our storage estimates have gone up. We have 6 different models we want to run for each of 12.5K species. For each model, each species outputs ~1GB of data, so each model generates about 12.5 TBs of data. 4 of these models are essential, which puts us at the low end of our storage requirements at ~50 TB, but if we can get the XSEDE resources (and storage) it would be great to have all 6 models, which would mean we'd need ~75 TB. Is that possible? Again, it also depends on getting the resources from XSEDE, but I just wanted to double-check with you about the storage.

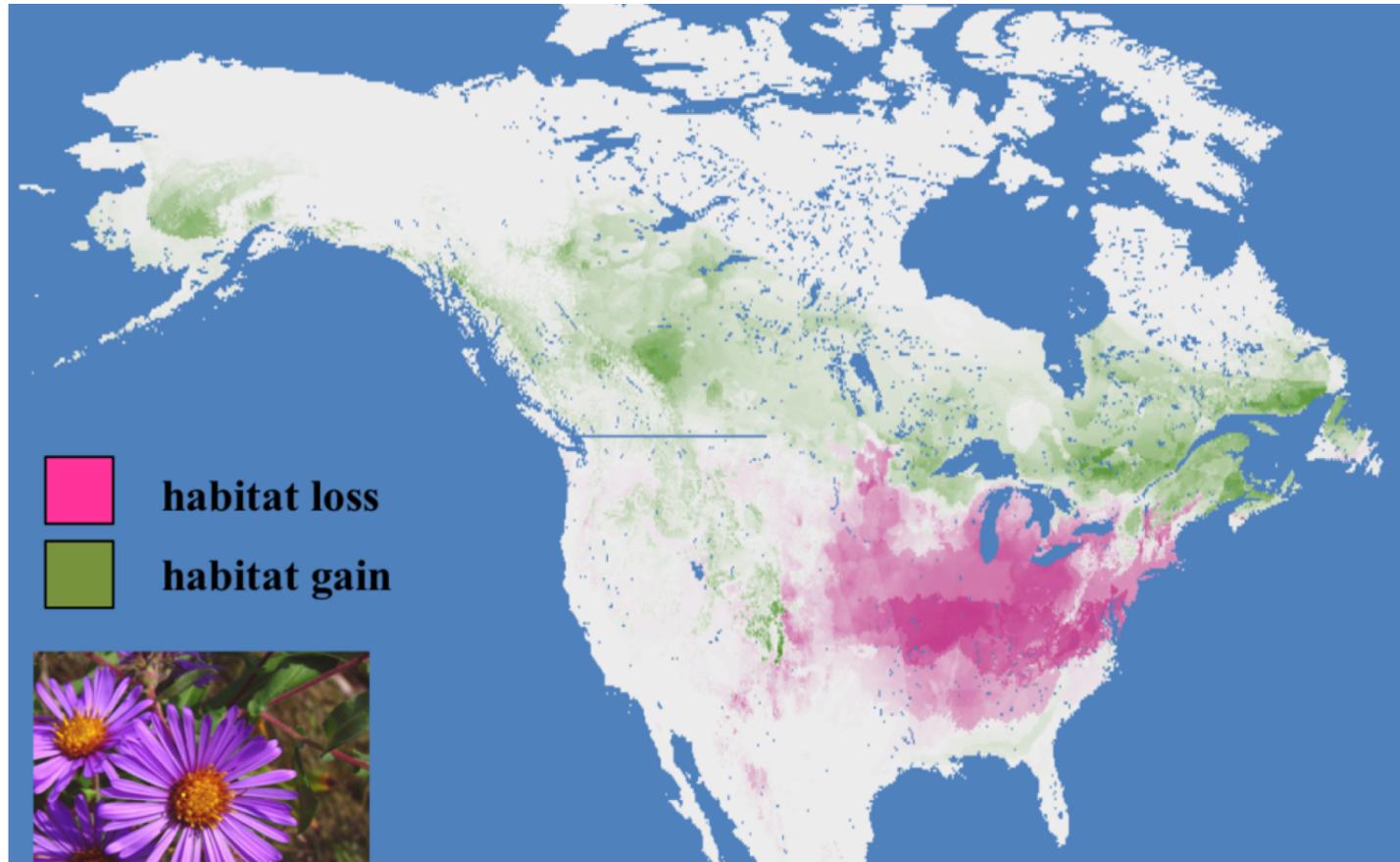
Cheers,
Charlie

NSF Postdoctoral Fellow
Harvard University Herbaria

Getting Data In/Out



Species Migration (Davis - OEB/Herbaria)



- 12.5k species @ ~1GB/species = 12.5 TB
- 6 different models = 75 TB



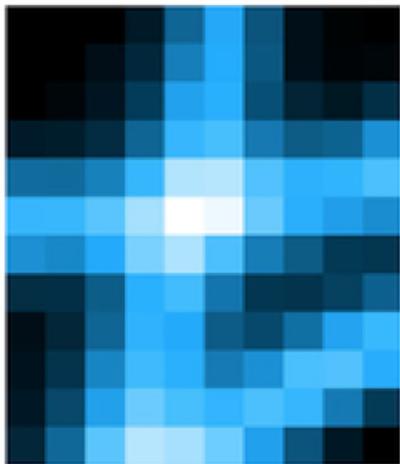
Modern Instrument Data

- Processing raw data and statistical analysis used to be done by high-end workstation when:
 - Data sets were smaller: 1-2 GB
 - Statistical algorithms were less rigorous

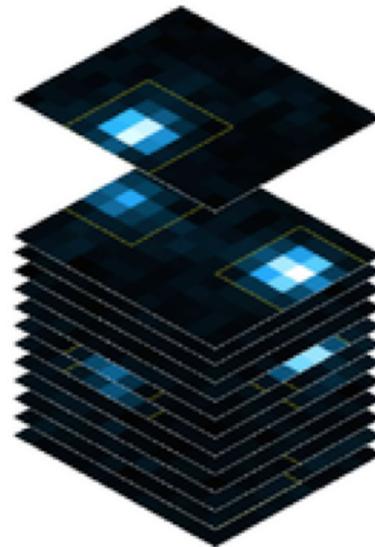
STORM Microscopy (Zhuang - CCB)

- Sub-diffraction-limit imaging by **Stochastic Optical Reconstruction Microscopy (STORM)**
- Using photo-switchable fluorescent probes to temporally separate the spatially overlapping images of individual molecule.
- a super-resolution image to be reconstructed from the positions of the localized probes

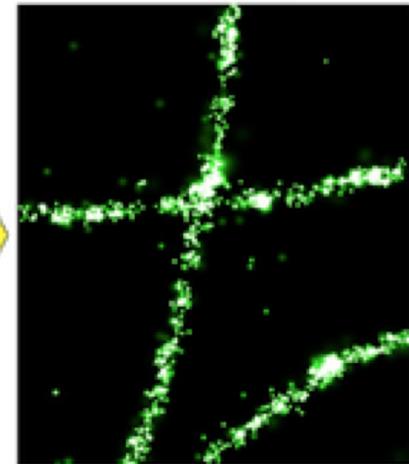
Diffraction-limited image



Stochastic activation of single molecules over many frames



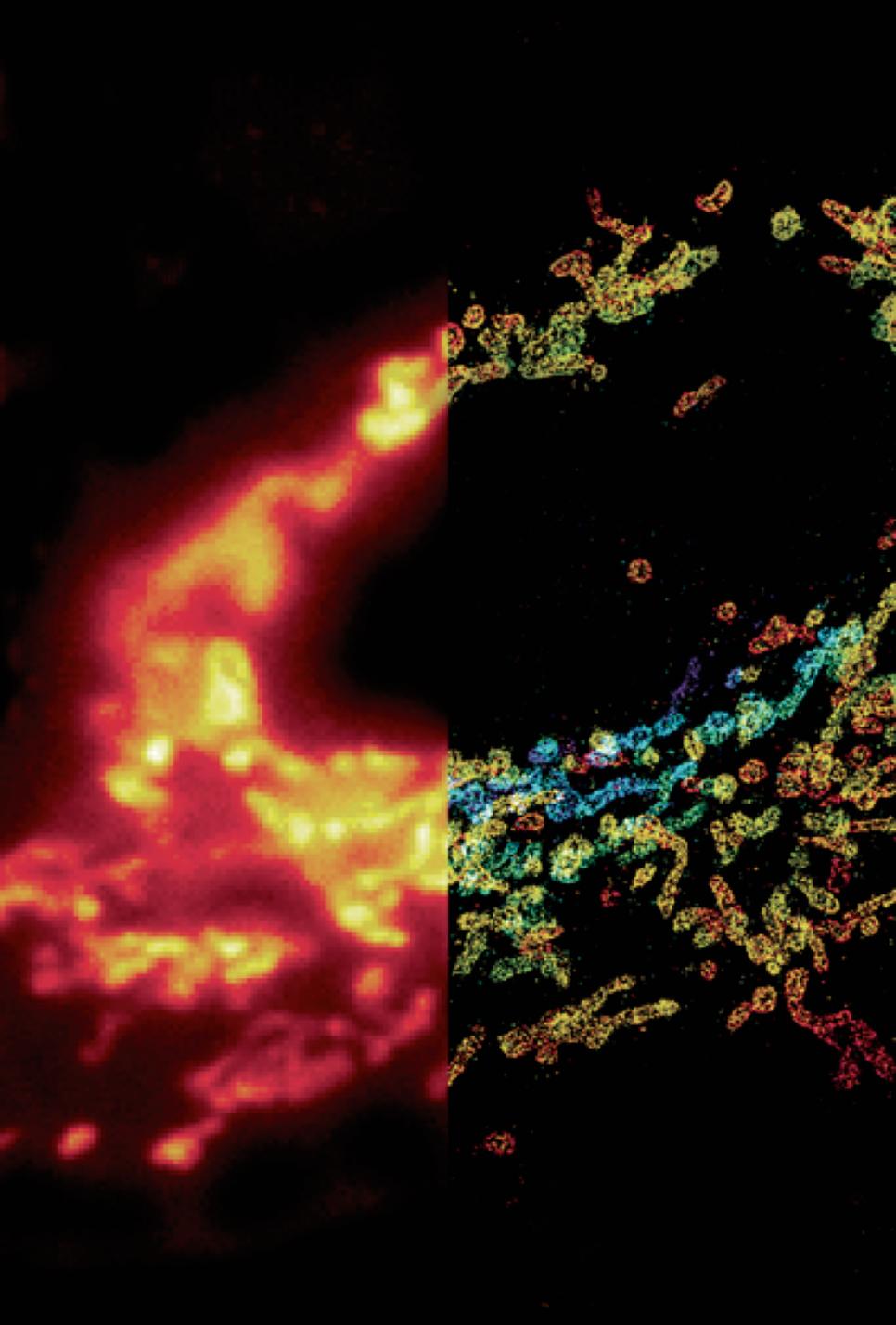
STORM image



- 2 TB created (100MB tiffs), 2 TB to transfer, 2 TB to process on the cluster every single day, 2 instruments. ~ 4TB/day or ~1PB / year

3D STORM

- Mitochondria in a cell
- Left: conventional image
- Right: 3D STORM image colored topology





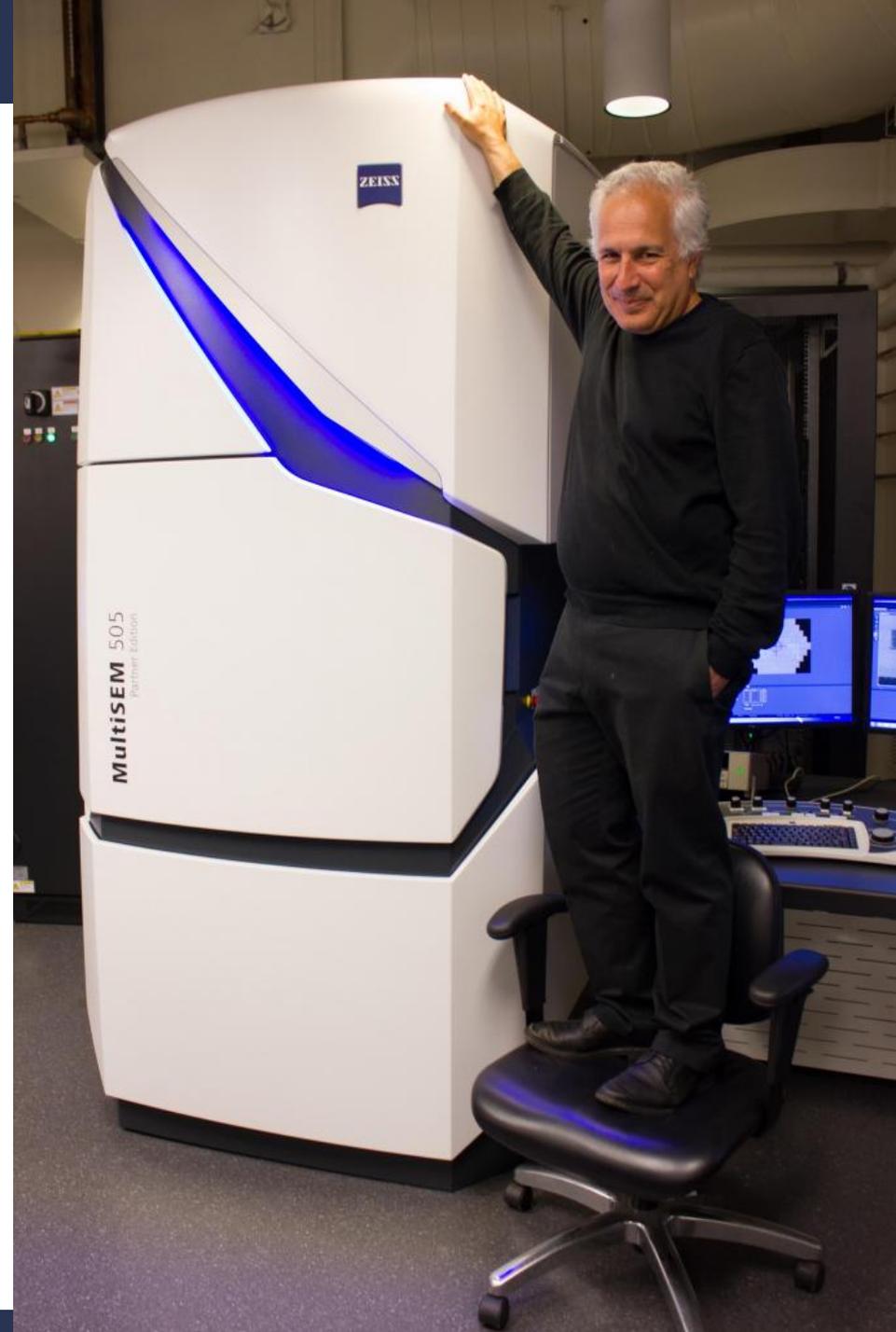
Multi-SEM (Lichtman - MCB)

3 TB

Every

Single

HOUR!



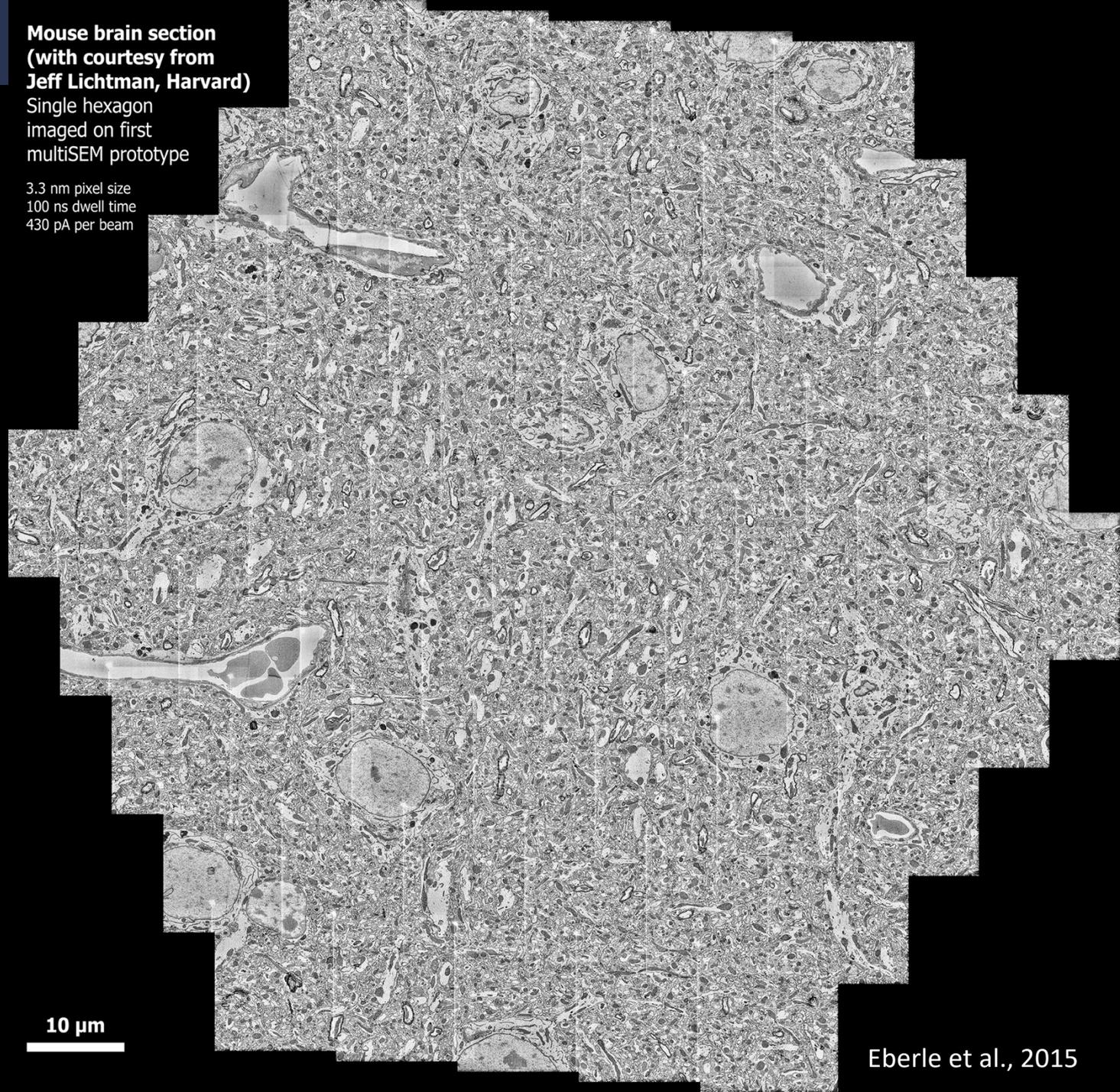


61 Beam
Hexagonal
Image 3.3
nm pixel
resolution

Mouse brain section
(with courtesy from
Jeff Lichtman, Harvard)

Single hexagon
imaged on first
multiSEM prototype

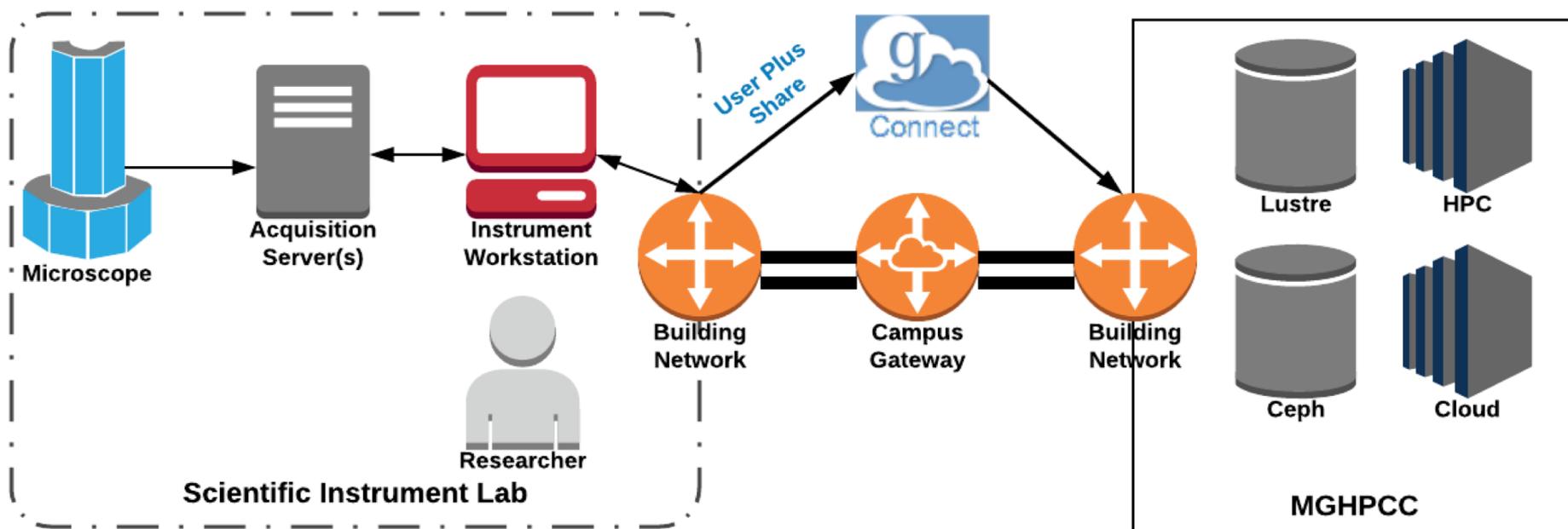
3.3 nm pixel size
100 ns dwell time
430 pA per beam



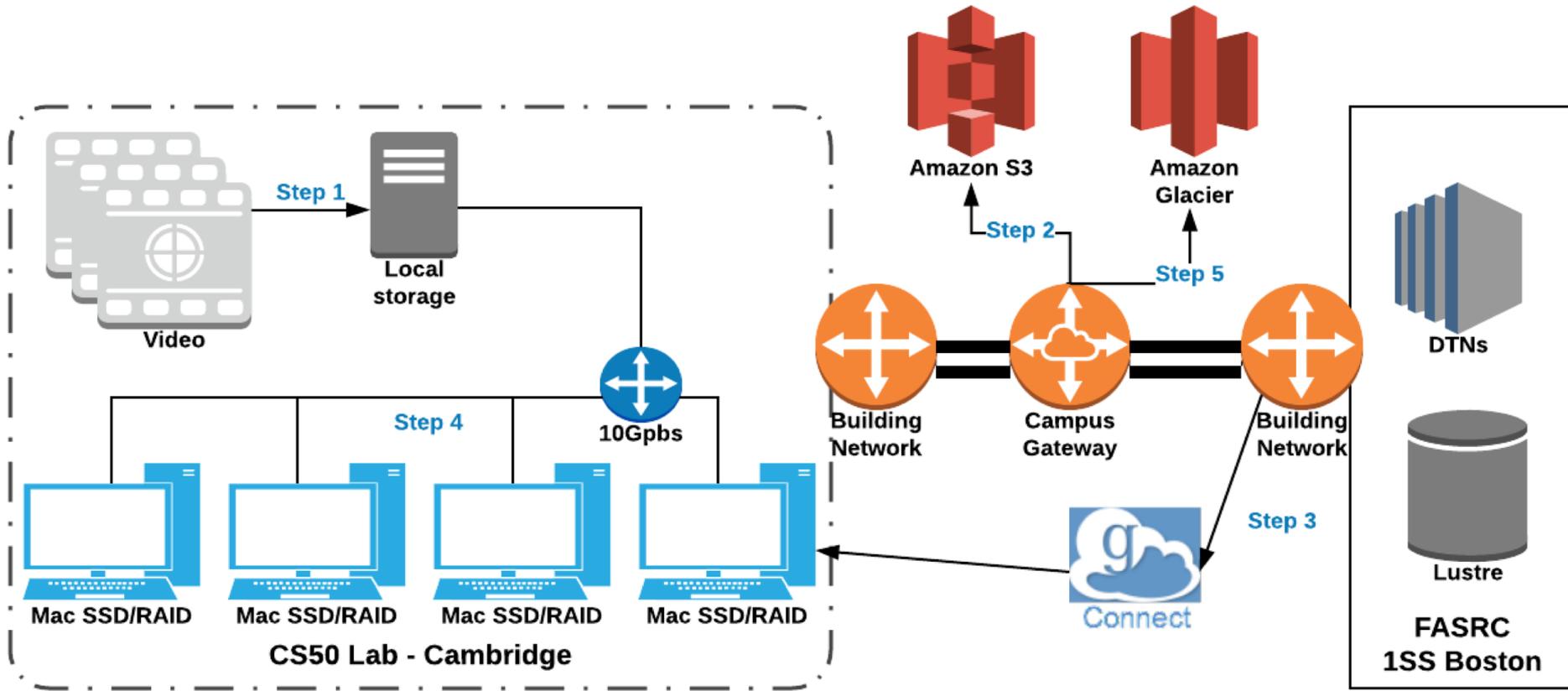
What took
5 hours
now takes
5 minutes!

10 μ m

Globus User Plus Endpoint



CS50 Example





Big Thanks ^^^

April 2018

GlobusWorld