# How to do Large Scale Data Research on a Slurm HPC Cluster with OmniSci

## OmniSci Virtual Summit, May 2020

Devika Kakkar and Ben Lewis
*Harvard Center for Geographic Analysis*

Raminder Singh
*Harvard Research Computing*

# NSF Collaboration between OmniSci and CGA

Industry & University Cooperative Research Center (I/UCRC)
**Center for Spatiotemporal Thinking, Computing, and Applications**
Harvard University
George Mason University

**FAS RC**

FAS Research Computing
Harvard University
Faculty of Arts & Sciences

Q Search www.rc ✕

FASRC Cluster     Services     Training     News     About FASRC     ⌕ Documentation

# FASRC Research Software Engineering (RSE) Services

- Design, development, optimization, deployment, maintenance of scientific software packages and data services
- Development of
  - Data Science/Machine Learning/Deep Learning/AI apps and platforms
  - data-intensive and big data platforms
  - scientific packages (Python, R, C++, Fortran, Julia, MATLAB ...)
  - data acquisition and analysis automation platform
  - functional and robust UI/UX
  - microcontroller programs (Arduino, Teensy)
- Add critical features to existing codebases
- Performance turning of existing software packages
- Maintenance of the current codebases developed by researchers
- RSE Training
- RSE Consultation

RSE Service Request : rse@rc.fas.harvard.edu

**Building well-engineered software & data services for Harvard University researchers that support and enrich research productivity and reliability**

# Geotweet Archive:
## A global social media record spanning time, geography, and language:

- Developed in collaboration with the University of Salzburg Department of Geoinformatics
- Extends from 2010 to the present and updated daily
- Geotagged by GPS or user designated place name

for more information:

https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi%3A10.7910%2FDVN%2F3NCMB6

# Introduction to Harvard's FAS Research Computing Cluster

- Compute: 100,000 compute nodes, 8-64 cores/node, 12Gb to 512Gb memory/node, 2,500,000 NVIDIA GPU cores
- Software: CentOS 7 operating system, Slurm job manager, Singularity, 1000+ scientific tools and programs
- Storage: 100 GB (Home dir), 4TB+ (Lab storage), 70Gb/node (Local scratch), 2.4PB (Global scratch), 3PB (Persistent Research data)
- #144 in TOP500 Supercomputers in world

**FASRC CANNON**
HARVARD'S LARGEST CLUSTER

100,000 CPU CORES
3,000+ NODES

500 TB RAM
40PB STORAGE
2.5M CUDA CORES

29 MILLION JOBS/YR
300 MILLION CPU HR/YR

3 DATA CENTERS @ 10K+ FT²
BOSTON, CAMBRIDGE, & LEED PLATINUM
GREEN DATA CENTER IN HOLYOKE, MA

500+ LAB GROUPS
OVER 5500 USERS

CANNON: THE FASRC CLUSTER IS NAMED IN HONOR
OF ANNIE JUMP CANNON, A PIONEER IN ASTRONOMY.

# FASRC Services

Software:

- Operating System CentOS 7
- 1,000+ scientific tools and programs
  - https://portal.rc.fas.harvard.edu/apps/modules
- C, C++, Fortran and Intel compilers available
- Languages like Python, R and Julia etc. can be used
- Databases like MySQL, Postgres and MongoDB

**Accounts**
Internal Users
External Collaborators
Industry Partners
Lab Setup

**Infrastructure**
Shared Infrastructure
Storage Lease
Customized Infrastructure

**Consultation**
Infrastructure Setup
Secure Environment
Data Use Agreement
Software Development
Complex Workflows

# Data Science Virtual Desktop Apps

# GIS Databases for Big Data

- **PostGreSQL**: Powerful, open source object-relational database system
- **PostGIS**: Provides spatial objects for the PostgreSQL database, allowing storage and query of information about location and mapping
- **OmniSci**:
  - Designed to overcome the scalability and performance limitations of legacy analytics tools
  - Super fast queries/analytics (including machine learning) of unindexed data (open source)
  - Super fast interactive rendering (free for educational use) of millions or billions of features, on-the-fly on a map
  - Leverages the massively parallel processing of GPUs alongside traditional CPU compute

# Geospatial on Harvard VDI

# OmniSci on FASRC

# OmniSci on FASRC

# Demo of the HPC interface - showing how easy it is to create a large instance

https://www.youtube.com/watch?v=TvqqikT_V58

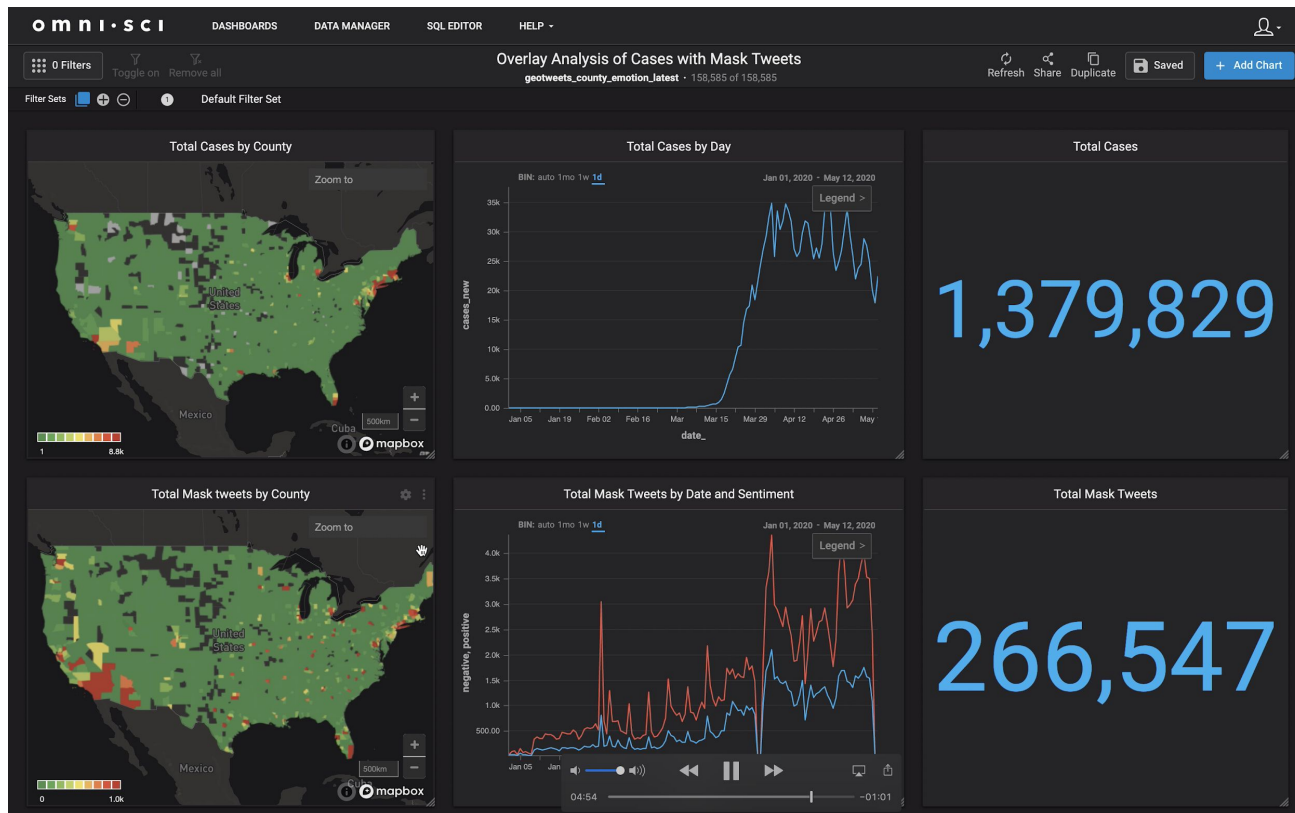# Demo of a Harvard project using OmniSci, running on the cluster

# Demo of a Harvard project using OmniSci, running on the cluster

# Demo of a Harvard project using OmniSci, running on the cluster

# Demo of a Harvard project using OmniSci, running on the cluster

https://www.youtube.com/watch?v=35nm5I__W-c

# Review of the scripts required to run OmniSci on the slurm cluster

Github Repo:
https://github.com/cga-harvard/GIS_Apps_on_HPC/tree/master/dev/OmniSci

The various script in the repo are run to achieve the following processes:

- Develop a User Interface to request the job parameter(partition, memory, CPUs, GPUs etc)
- Run a slurm job with the parameters specified by the user on Launch
- Install Omnisci using Singularity (Finding port, setting data directory, setting passwords etc.)
- Pass connection parameter to the user (if successful) else display the error

# How other slurm clusters can do the same, and get involved in the open source project

# Future Directions - Distributed OmniSci on FASRC



Distributed Configuration OmniSci [5]

# Q and A

# References

- Geospatial tools on Harvard Cluster Computing
  https://gis.harvard.edu/geospatial-data-science-tools-and-data-harvards-high-performance-computing-infrastructure

- Harvard Center for Geographic Analysis  / OmniSci Collaboration
  https://www.omnisci.com/blog/rejoining-forces-a-new-old-partnership-with-the-harvard-center-for-geographic-analysis

- Harvard Center for Geographic Analysis  https://gis.harvard.edu/

- Harvard Research Computing (FASRC)  https://www.rc.fas.harvard.edu/

- FASRC Cluster Architecture https://www.rc.fas.harvard.edu/about/cluster-architecture/

- Introduction to Cluster Computing:
  https://www.rc.fas.harvard.edu/wp-content/uploads/2019/12/Intro-to-Cannon.pdf

# Thank you

Devika Kakkar (kakkar@fas.harvard.edu)
Ben Lewis (blewis@cga.harvard.edu)
Raminder Singh (r_singh@g.harvard.edu)

# Outline

- Introduction to the project
- Harvard's Center for Geographic Analysis
- Intro to Harvard's Computation Cluster
- Overview of challenges researchers are facing with data which HPC is designed to address
- Installing OmniSci on Harvard Cluster
- Demo of the HPC interface - showing how easy it is to create a large instance
- Demo of a Harvard project using OmniSci, running on the cluster
- Review of the scripts required to run OmniSci on the slurm cluster
- How other slurm clusters can do the same, and get involved in the open source project
- Q and A