# THE UNIVERSITY OF CHICAGO

## BIOLOGICAL SCIENCES

Computing with the CRI: Storage, HPC, and Virtual Servers

November 2nd, 2017

# Introduction

- IT Operations & Infrastructure team
- Three main service lines
  - Storage (Lab shares)
  - Virtual servers and system administration
  - HPC and other compute services (e.g. Stats, WinStats)

- Procedures following the BSD's Information Security Policies to be able to store and process ePHI (HIPAA) data for research use

- All primary resources are behind firewalls in 6045 Kenwood datacenter

- Data is backed up to tape in the 1155 (E 60th St.) datacenter

# Resource Access

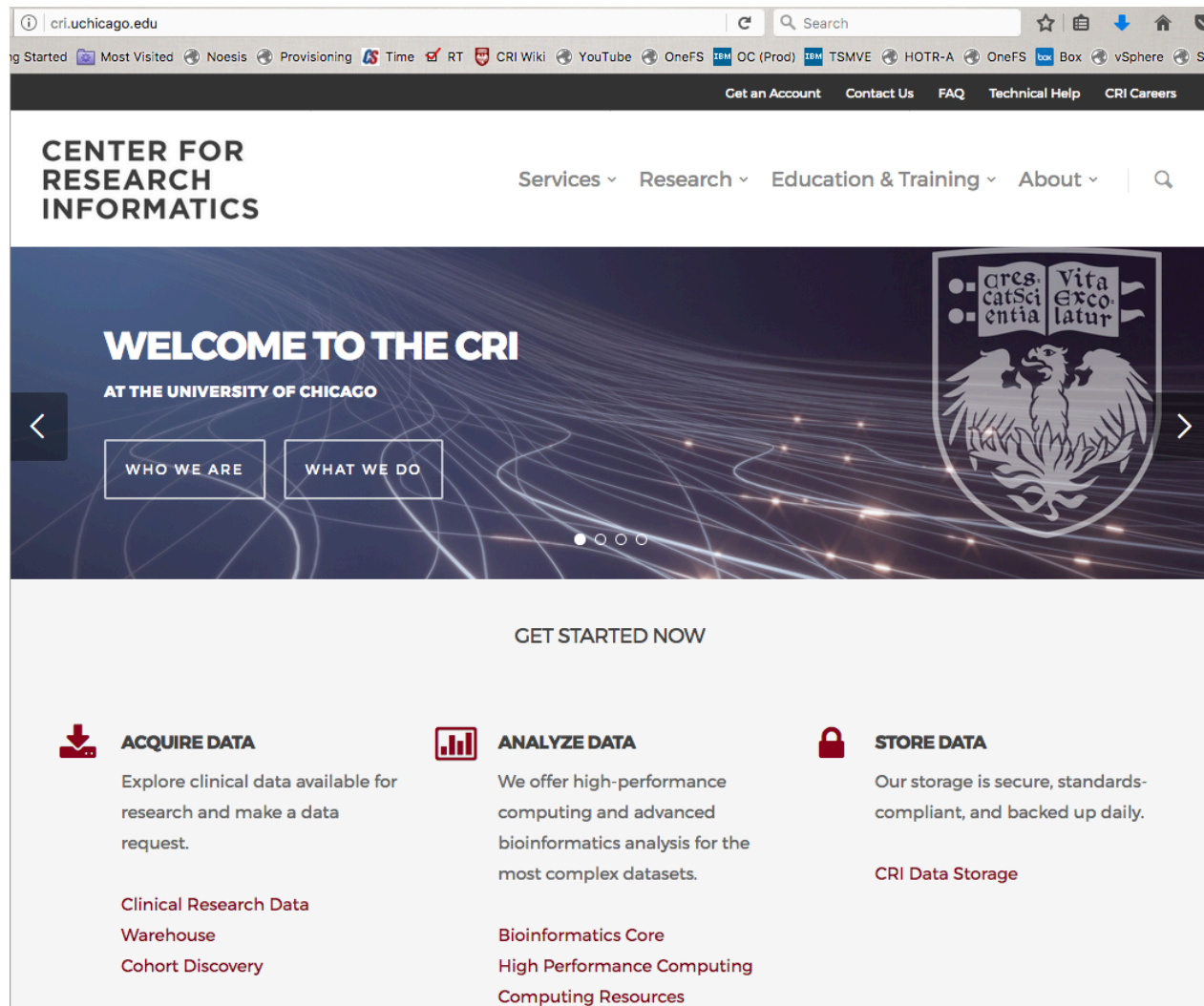| REDCap | HPC Cluster | Lab Shares | VMs | WinStats |
|--------|-------------|------------|-----|----------|

In order to access CRI resources:
- You need to have an active BSD account
- You need to request access using our Web Provisioning Forms
- A ticket will be created and you will be sent access instructions

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Requesting Access to CRI Resources

THE UNIVERSITY OF

**CHICAGO**

CENTER FOR
**RESEARCH**
INFORMATICS

**Web Provisioning**

crescat scientia; vita excolatur

Home      Admin      Help

[ Admin Login ]

## CRI Web Provisioning - New Service Request

**Request Access to CRI Resources**
Request access to the High Performance Computing Cluster or Stats Server. You must have a BSDAD account in order for access to be granted. You can submit a request for a BSDAD account by visiting the BSDAD Account Request page.

**Request Lab Share Creation**
Request creation of a Lab Share to store research data in a highly secure clustered storage system with tape library backup.

**Request User Access to Lab Share**
Request user access to an existing Lab Share

**Request Access to Gardner HPC**
Request access to the Gardner High Performance Computing Cluster

**Request Software Installation on the HPC**
Request installation of software on the High Performance Cluster

**Request Software for Commercial Statistical Analysis on Windows**
Request Software for Commercial Statistical Analysis on Windows

**Request Virtual Machine Creation**
Request creation of a virtual machine instance

**Request Firewall Access**
Request firewall access

**Revoke Resource Access**
Request removal of user access to the HPC Cluster, Stats Server, REDCap, or Lab Share

**Request Collaborator Account Creation**
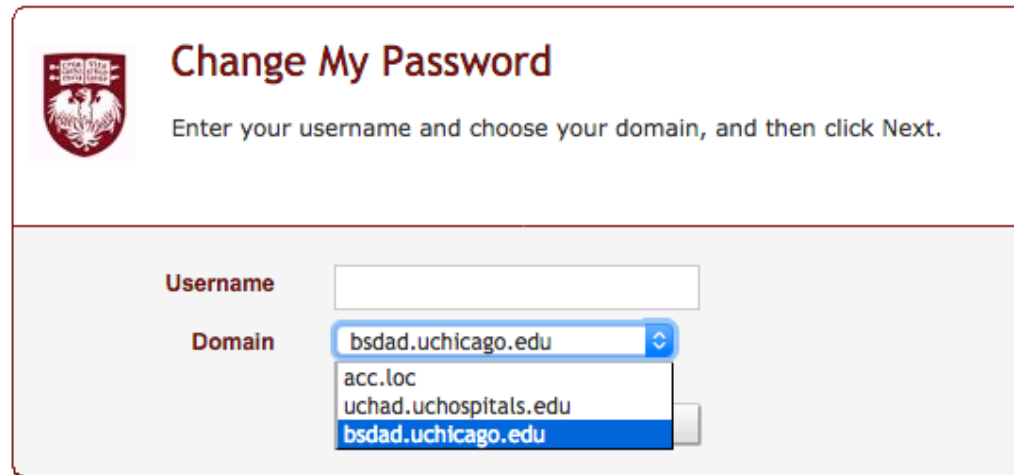Request creation of a collaborator account for REDCap, HPC or Storage access

THE UNIVERSITY OF
**CHICAGO MEDICINE &**
**BIOLOGICAL SCIENCES**
Center for Research Informatics

# Changing your Temporary Password

https://mail.uchospitals.edu/changemypassword/



- Make sure you change your temporary password quickly.
- Do not write down your password; you can write down a clue to help you remember the password instead and keep in a secure location.
- Do not share your password or account access with anyone else.
- When choosing a new password, you will be required to use at least one capital letter, one lowercase, one number, and one symbol.
- In addition to these requirements, it is good practice to refrain from using personal information like a birthday or name.

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# VPN Access

- Needed if you want to access CRI resources from outside campus
- Collaborator accounts are enabled for BSDVPN access
- Download and install the Cisco AnyConnect Secure Mobility Client
  - https://cvpn.uchicago.edu
  - https://bsdvpn.uchicago.edu
- You can setup two factor authentication for your account
  - https://2fa.bsd.uchicago.edu



THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Getting Help

- Contact the CRI using our Support Portal:
  - cri.uchicago.edu > Technical Help > Support for Resources

  - This will create a ticket in the appropriate queue for faster processing
  - Ensures we collect information necessary to start processing the ticket

# Getting Help (CBIS)

If you have a non-CRI managed account, please contact the CBIS Service Desk at help@bsd.uchicago.edu or call 773-702-3456.

- *For BSD email issues*
- *BSD account password reset*



THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# CRI Data Storage Infrastructure

- 1.8PB Capacity

- High performance Nodes

- Archive Tier Nodes

- Backed up nightly

- Connected to the HPC cluster

- Easy and flexible share access

- Extended access controls

- Quota notifications



Windows

OS X

Linux

HPC Cluster

WinStats

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# CRI Storage Shares

- Home Directories

- Lab Shares

- Departmental Shares

- Group Shares

- Scratch Space

  – Accessible only from the HPC Cluster
  – Temporary space for staging input for analysis jobs and temp files from job execution.

# Storage Access Methods

- SMB/CIFS
- NFS (Servers Only)
- SCP/SFTP
  You can use graphical SFTP clients such as:
  - WinSCP, CoreFTP (Windows Only) or
  - FileZilla, Cyberduck (Windows & Mac)

- Rsync (use cri-syncmon.cri.uchicago.edu)

# Virtual Server Environment

- Current Environment
  - 2 servers for management and CRI systems
    - Currently using 8% of CPU
    - Currently using 37% of memory
  - 6 servers for all other servers
    - Currently using 6% of CPU
    - Currently using 49% of memory
  - Each server:
    - 2 sockets – each socket with 20 cores
    - 512GB Memory
    - No hard drive – only dual SD card
  - Data Storage:
    - Compellent
    - SSD layer – All writes happen here
    - 7K Layer – Less used data moved here overnight

# Request a Virtual Machine

- Fill out form

- We will read the answers and respond to the ticket with questions

- Our Process:

  1. Create VM
  2. Configure server with needed space
  3. Configure security settings
  4. Install main software
  5. Give over VM to customer to install/configure the rest of the server
  6. After configuration, we will run security scans
  7. Get firewall configured for final network
  8. Move server out of staging network
  9. Setup basic monitoring of system: disk, CPU, memory, known applications
  10. Setup backups – Nightly snapshot of VM and exports of known databases

# Request Form

- People Information

**Requester Info**

Requester First Name: [                    ]

Requester Last Name: [                    ]

Requester Email: [                    ]

Requester Phone Number: [                    ]

Requester Department: [                    ▼]

☐ Requester is the PI/Lab Manager

**PI/Lab Manager Info**

PI/Lab Manager First Name: [                    ]

PI/Lab Manager Last Name: [                    ]

PI/Lab Manager Email: [                    ]

PI/Lab Manager Phone Number: [                    ]

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Request Form

- Environment

**Virtual Machine Info**

Environment:

Network Access:

Operating System:

Server Configuration:

Purpose of VM:

- Network

Network Access:

Operating System:

Server Configuration:

- OS

Operating System:

Server Configuration:

✓
Test
Development
QA
UAT
Production

✓
Campus
Internal Only
External

✓
Redhat Linux 7
Windows 2012 R2

# Request Form

- Size of server

Server Configuration:

Purpose of VM:

✓

Basic – 2G MEM, 1CPU
Advanced – 4G MEM, 2 CPU
Enterprise – 8G MEM, 4 CPU

- Purpose

Purpose of VM:

- Data Sensitivity

Data Sensitivity:  ☐ PHI  ☐ IP  ☐ PII  ☐ PCI  ☐ Other  ☐ No Sensitive Data  glossary

- Definition of Types of Data

**Data Sensitivity: Glossary of Terms**

- PHI: Protected Health Information
- IP: Intellectual Property
- PII: Personally Identifiable Information
- PCI: Payment Card Industry

**Close**

# Request Form

- **C** onfidentiality

What type of data will reside on this system?
⚪ Any type of personal identifying data  ⚪ Private, but not protected data  ⚪ Public Data

Reason for 'type of data' selection:

- **A** vailability

What is the expected availability of the system?
⚪ Needs to be up 24/7 with scheduled maintenance windows
⚪ Only needed during business hours, can be rebooted after 7pm without notifying anyone
⚪ Does not matter if it goes down during the day

Reason for 'expected availability' selection:

- **I** ntegrity

How difficult would it be to recreate the data on this system if corrupted?
⚪ Impossible  ⚪ Not impossible, but difficult and time consuming  ⚪ Easy

Reason for 'data recreation difficulty' selection:

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Request Form

- ### Other Software

  What software will be running on the system?

  ☐ Apache  ☐ IIS  ☐ MS SQL  ☐ MySQL  ☐ NGINX  ☐ Oracle  ☐ Other

- ### Firewall Rules

  What connections does this system have? Or what ports need to be opened (port 80, 3306, etc)

- ### Disk Space Needs

  How much data outside of the operating system will be loaded onto this system?

- ### Other Users

  Please use the search form to find and add inviduals who need to access this virtual machine.
  Open Search Form

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Other Software used by CRI

- SaltStack – Configuration Management

- IPA Server – Redhat's Identity Policy Audit Linux Domain

  - Trust with BSDAD to provide BSDAD Accounts to Linux servers and bulkstorage

- Device 42 – Inventory System with API

- Nagios – Monitoring/Alerting - Can have custom setup

- RequestTracker (RT) – Ticketing System from Best Practical – Can provide ticket queues for other groups

- GitLab – https://git.cri.uchicago.edu - Can provide git, and other development tools (code review, wiki, issue tracking), to other groups

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Commercial Software - Winstats

- Specifications

  - 2 CPUs: 2.4GHz Intel Xeon E7-8870

  - 10 cores per CPU

  - 512 GB of RAM

  - 1.18 TB of dedicated SSD storage

# Commercial Software – Stats (Linux)

- Specifications

  - 2 servers

    - 2 Intel Xeon E5-2698 v3 CPUs

    - 2.3 GHz

    - 16 ops/cycle

    - 16 cores per CPU

    - 768 GB RAM per server

# Commercial Software - Software

- SPSS (Winstats)

- Stata

- SAS

- R (not really commercial but...)

- Matlab (Linux)

- Mathematica

# High Performance Computing - Staff

- Mike Jarsulic (Sr. HPC Administrator)
    - Lived in Pittsburgh for about 32 years
    - Attended the University of Pittsburgh (at Johnstown)
    - Bettis Atomic Power Laboratory (2004-2012)
        - Scientific Programmer (Thermal/Hydraulic Design)
        - Thermal/Hydraulic Analyst - USS Gerald R. Ford
        - High Performance Computing
    - University of Chicago (2012-present)
    - Dislikes Mimes

- Qiannan Miao (Student HPC Administrator)

    - Current UChicago MPCS Student
    - Lots of coursework focusing on machine learning
    - Graduates 12/2017
    - Dislikes Mimes

# High Performance Computing – Martin Gardner



- Graduate of the University of Chicago (1936)

- Yeoman on the USS Pope during WWII

- Amateur Magician

- Mathematical Games

- Skepticism

- Literature

- Art

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# High Performance Computing – What is a HPC Cluster?

- Storage

- Login Nodes

- Scheduler Node

- Compute Nodes

- GPUs

- Xeon Phi

# High Performance Computing – Storage

- Isilon Storage Cluster

  – Home Directories (/home/<userid>)

  – Lab Shares (/group/<labid>)

  – Applications (/apps/software)

  – Total Size: 1.8 PB

  – Bandwidth: 1 Gb/s

  – Permanent, Quota'd, Backed Up

  – Available everywhere

- Scratch Space

  – Total Size: 175 TB

  – Bandwidth: 56 Gb/s

  – Purged, Private, Not Quota'd, Not Backed Up

  – Limited availability

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# High Performance Computing – Login Nodes

- Purpose

  – User interface for the HPC cluster

  – Composing editing jobs

  – Submitting jobs

  – Tracking/Managing jobs

  – Writing source code

  – Compiling source code

- DO NOT RUN ANALYSIS ON THE LOGIN NODES!!!

# High Performance Computing – Scheduler Node

- Purpose

    – Keeps track of which resources are available on each compute node of the cluster
    – Schedules jobs based on available resources
    – Maintains historical metrics on jobs

# High Performance Computing – Compute Nodes

- Specifications

  - Nodes:
    - 88 standard
    - 28 mid-tier
    - 4 high-tier

  - Processors (2 per node):
    - Intel Xeon E5-2683 v3
    - 14 cores
    - 2.0 GHz
    - 16 ops/cycle

  - Memory
    - Standard: 4 GB/core
    - Mid-tier:  16 GB/core
    - High-tier:  45 GB/core

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# High Performance Computing – Accelerator Nodes

- GPUs

  - 5 Nodes
  - 1 NVidia Tesla K80 per node
    - Contains 2 NVidia Tesla GK210 GPUs
    - 2496 CUDA cores
    - 24 GB memory

- Xeon Phi

  - 1 node
  - 2 Intel Xeon Phi P5110 coprocessors
    - 60 Cores per coprocessor
    - 8 GB RAM

# High Performance Computing – Software

- Compilers

  - GNU
  - Intel
  - PGI
  - Java JDK

- Scripting Languages

  - Perl
  - Python
  - R

- Software Environment

  - LMod

# High Performance Computing – Obtaining an Account

- Prerequisites: BSD Account

- Sign up for and account

  - http://cri.uchicago.edu
  - Experience Level
  - Software Requests
  - Email Address for Job Output
  - Emergency Phone Number

- Collaborator Accounts

# High Performance Computing - Support

- How to get help?

  - Email: hpc@rt.cri.uchicago.edu
  - Support form on the CRI Website
  - Documentation: Coming Soon
  - Ask a friend
  - Office Hours (Tuesday/Thursday)
  - Bioinformatics Core
  - CRI Seminar Series

- How not to get help?

  - Calling me (unless it's an emergency)
  - Email me directly

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Future: New Storage System

- Problems with Isilon Cluster

  - Enterprise solution; not designed for research environments
  - High expansion cost
  - High maintenance cost
  - Ethernet requirement
  - Not POSIX-compliant (parallel file access)
  - Low read iops

- Problems with Scratch Storage

  - Glusterfs
  - Poor performance on small files
  - Poor metadata performance
  - Stability

# Future: New Storage System

- Goals

  - Combine scratch and permanent storage into one system
  - Lower costs
  - Increase performance
  - Utilize Infiniband network
  - POSIX-compliant
  - Maintain stability and security of the Isilon

- Options

  - Lustre (Intel)
  - GPFS (IBM)

# Future: New Storage System – Metadata Ops

Values are in operations per second

All operations are file operations

| Operation | Isilon | Scratch | Lustre | GPFS |
|---|---|---|---|---|
| Creation | 3188 | 884 | 22025 | 30665 |
| Stat | 80940 | 5684 | 96346 | 39247 |
| Read | 6474 | 6176 | 50557 | 88245 |
| Removal | 210 | 209 | 16521 | 11736 |

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Future: New Storage System – Multiple Files

All Values in MB/s

8 processes; 1 file per process

| API | Operation | Isilon | Scratch | Lustre | GPFS |
|---|---|---|---|---|---|
| POSIX | Read | 110 | 347 | 31949 | 22430 |
| POSIX | Write | 108 | 106 | 3459 | 11292 |
| MPI-IO | Read | 9 | 370 | N/A | 22212 |
| MPI-IO | Write | 10 | 138 | N/A | 7522 |

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Future: New Storage System – Single File

All Values in MB/s

8 processes; single shared file

| API | Operation | Isilon | Scratch | Lustre | GPFS |
|---|---|---|---|---|---|
| POSIX | Read | 109 | 119 | 24980 | 18447 |
| POSIX | Write | 94 | 112 | 4503 | 9176 |
| MPI-IO | Read | 5 | 110 | N/A | 17439 |
| MPI-IO | Write | 3 | 65 | N/A | 8467 |

THE UNIVERSITY OF
CHICAGO MEDICINE &
BIOLOGICAL SCIENCES
Center for Research Informatics

# Future: Current Projects

- Globus Online for data transfers (being set up now)

- Deep Learning / Machine Learning (exploring needs and options)

- Container-based computing for both HPC and applications (exploring needs, options w/ tests)

- Visualization (exploring and testing)

# Questions?