

# Minerva Town Hall

## Spring 2024

Lili Gai, PhD, Director for High Performance Computing and Data  
April 16th, 2024



Icahn  
School of  
Medicine at  
**Mount  
Sinai**

# Outline

- ❑ 2023 User Survey Results
- ❑ Accomplishments & Updates
- ❑ 2024 Roadmap
- ❑ Appendix (Minerva Usage)



# 2023 Minerva User Survey Results

# 2023 Minerva Survey Results

## We asked five questions:

Q1: Overall, how satisfied are you with the LSF queue structure, compute and storage resources (GPUs, high-memory nodes, TSM, etc)?

Q2: Please rate current software environment (packages and services such as database, data transfer, container etc).

Q3: Please rate your satisfaction with operations (ticket system, responsiveness of staff, documentation, user support etc).

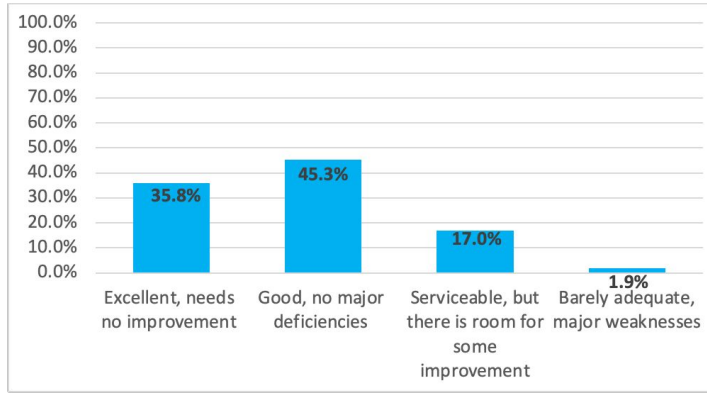
Q4: Which of the following would you most prefer for future Minerva expansion?

Q5: What suggestions do you have for improving our service?

**We received 54 responses and 39 comments from 1,100 active users in Jan 2024 (5.0% response rate).**

# 2023 Survey Results Question 1

**Q1: Overall, how satisfied are you with the LSF queue structure, compute and storage resources (GPUs, high-memory nodes, TSM, etc)?**



User satisfaction( $\geq$ Good)

2023: 81%

2022: 84%

2021: 85%

2020: 81%

2019: 65%

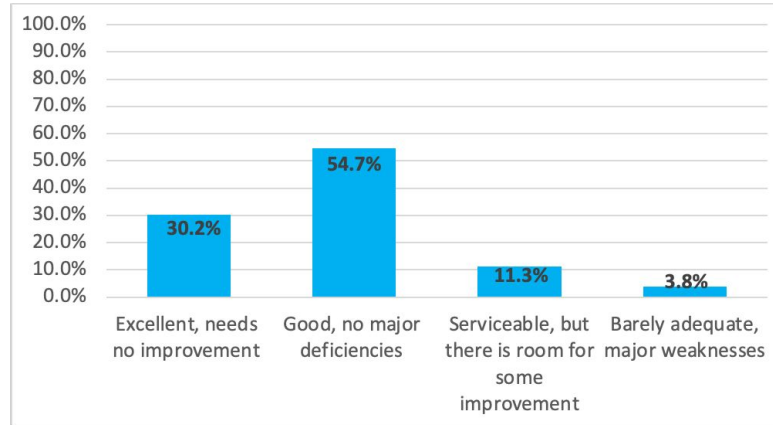
2018: 54%

## Raw comments(9):

- **LSF queuing:** I find that sometimes I will submit jobs and they will be pending for a long while, and other times they will run right away
- **Nodes:** Wait times for GPU resources can be very long, to the point it is shorter to simply run programs on CPU and wait longer for them to run. Sometimes it's very hard to get high-mem nodes even if my "priority" is high (meaning that I didn't calculate anything for weeks). More nodes in general, jobs can take a long time in queue currently

# 2023 Survey Results Question 2

**Q2: Please rate the current software environment (packages and services such as database, data transfer, container etc)**



User satisfaction(>=Good)

2023: 85%

2022: 93%

2021: 81%

2020: 80%

2019: 80%

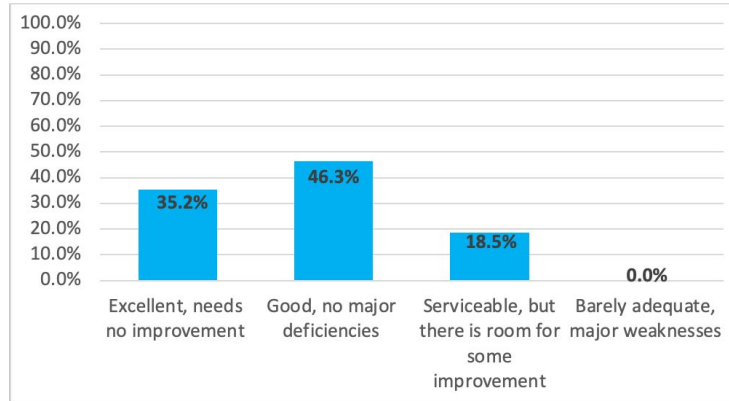
2018: 67%

## Raw comments(6):

- Notable limitations (e.g., docker); RStudio web server is extremely useful but it seems like only one person is managing it and it can only be used on li03c03.
- Files tab on OnDemand needs to have bookmarked folder locations so we can navigate faster. OnDemand needs a dark mode.

# 2023 Survey Results Question 3

**Q3: Please rate your satisfaction with operations (ticket system, responsiveness of staff, documentation, user support etc)**



User satisfaction(>=Good)

2023: 82%

2022: 88%

2021: 86%

2020: 91%

2019: 73%

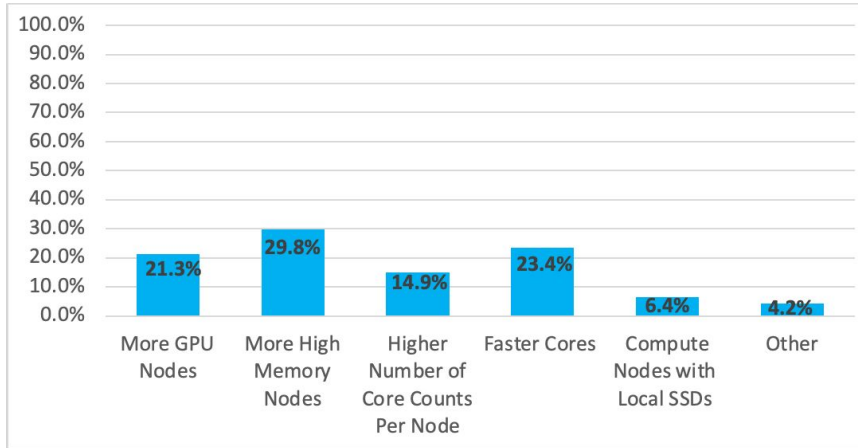
2018: 80%

## Raw comments (12):

- **Documentation:** Better and more up-to-date documentation (e.g., for Jupyter server, VS Code connection, etc.) needed; a lot of the documented methods don't work anymore
- **Tickets/Staff:** *There is a clear need for more staff. They are very helpful and responsive but take a lot of time. It is not reasonably possible for only 2-3 people to manage the whole system and be responsible for HPC tech support at the same time. More staff!!!* Any useful expansion will take time to integrate and as of now I don't think they can handle what we already have.

# 2023 Survey Results Question 4

Q4: Which of the following would you most prefer for future Minerva expansion?



With the response, we keep this mind with this upcoming Minerva refreshment

- More GPUs ✓
- More High memory nodes ✓
- Higher # of cores per node ✓
- Faster cores ✓
- Some compute nodes with local SSD ✓



# 2023 Survey Results: More Raw Comments

Thank you for your feedback! We are posting all the responses on our website.

## Positive Comments:

- I would like to express my sincere appreciation for the exceptional service provided by your team. Your accommodating nature, especially when it comes to handling out-of-the-ordinary requests, has been incredibly impressive and has significantly enhanced my experience.
- My concerns are always answered promptly and with professionalism.
- Staff is great, got back on all my issues with a very fast turnaround and helpful information
- with the addition of OnDemand, data transfer has become so much easier.
- No major suggestions. The Minerva HPC team continue to be excellent to work with. The team's responsiveness and ability to accommodate our team has been great. Thank you for all you do.

## Other comments:

- Due to staffing shortages, tickets take a little longer to clear, but this is completely understandable at the moment. • Sometimes helpdesk is understaffed and it takes a while to get an answer. Usually staff responds to inquiries, but there are also instances where no responses are given more than 24h later, which has an impact on working progress.
- Please more GPUs, they are constantly being used and it is limiting the research that I am doing. Does not have to be NVIDIA, other companies like AMD and Intel have respectable deep learning processors. Worth looking at these if it is more cost efficient.
- The only thing is that I often need working interactively with internet connection and, sometimes, it is difficult to find available interactive nodes. Compute nodes with internet access.

# 2023 Data Ark User Survey Results

# Data Ark Survey Results

Four questions we asked:

Q1: How would you rate your satisfaction with the Data Ark: Data Commons data quality and availability at Mount Sinai?

Q2: How would you rate your satisfaction with Data Ark support?

Q3: What barriers exist preventing your usage of Data Ark?

Q4: Please enter any other feedback for improving our services on Mount Sinai Data Ark.

**We received 19 responses with 9 comments out of 183 all Data Ark users and PIs in Jan 2024 (10% response rate)**

# Data Ark User Survey Results

Q1: How would you rate your satisfaction with the Data Ark: Data Commons data quality and availability at Mount Sinai?

User satisfaction(>=Satisfied)

2023: **79%**

2022: **73%**

	N %
Very Satisfied	42%
Satisfied	37%
Neither Satisfied Nor Dissatisfied	16%
Dissatisfied	0%
Very Dissatisfied	5%

Q2: How would you rate your satisfaction with Data Ark support?

User satisfaction(>=Satisfied)

2023: **84%**

2022: **73%**

	N %
Very Satisfied	42%
Satisfied	42%
Neither Satisfied Nor Dissatisfied	5%
Dissatisfied	5%
Very Dissatisfied	5%

# Raw Comments from the Survey

Thank you for your feedback! We are posting all the responses on our website.

## Positive comments:

- Great support team! MSDW dataset is very useful.
- I've appreciated the improvements to the Data Ark website over the past 6 months or so.

## Other comments:

- Provide a platform or a upload form for people to upload publicly available files and datasets. Provide guidelines in a website for data sharing and upload by users, so everyone can contribute.
- A lack of knowledge of how to query data, and a lot of red tape to get access.
- Have people who've used the system (or designed it) for complex work in the Data Commons to explain how to use it better. Going for help and getting the response "I don't know" or "Just pay us and we'll get the data for you" is very frustrating.
- But further improvements in clarity and specificity regarding data use agreement requirements, IRB requirements for specific data sets, procedures for use, and guidelines around acknowledging the CTSA would be valuable.

# Accomplishments & Updates

# Staffing

The HPC team consists of four computational scientists/bioinformaticians

- Hyung Min Cho, PhD
- Jielin Yu, PhD
- Yiyuan Liu, PhD (part time Data Ark)
- One GPU/AI expert joining this May

...and six HPC architects/admins positions

- Wei Guo, PhD
- Kali McLennan
- Eric Rosenberg (part time on Minerva)
- **Three positions are just reopened Q1**
  - Lead HPC Architect
  - HPC Architect
  - HPC admin

**We have got some good candidates!!**



# Minerva Refreshment and Expansion: Processing PO

Budget finally “Approved” for new equipment in HPC+ facility cost

## Hardware Purchased:

- **146 high memory** compute nodes (1.5TB DDR5 memory) - 14,016 cores
  - **Latest** Intel latest 5th Xeon(R) Emerald Rapids 8568Y+ 48C, 2.3GHz
- **210 GPUs in total**
  - **188 H100-80GB (SXM5)** NVlinked GPUs in 47 nodes
  - **32 L40S** GPUs in 4 nodes
- 3.84 TB Local NVME SSD per node
- NDR400 IB networking (400Gb/s)
- 4 login nodes and 6 service node
- 300 TB memory in total
- Direct water-cooling solution





# Minerva Refreshment and Expansion: Timeline

In production targeted at Dec 2024 - Jan 2025 ???

## The Refreshment Processes:

- RFQ issued for bidding on **Feb. 23** ✓
- Quotes from all the vendors received on **Mar. 11** ✓
- PO submitted on **March 28th** ✓
- HPC equipment arrives by **mid Sep**
- Construction work in Hess done by **Nov. 16th** ? - working with facility to expedite
- Installation and acceptance test
- Soft open to friendly users
- Open to public



# LSF Operation: Changes on GPU Queues

## Goal:

- Improve the throughput and enhance utilization efficiency of the limited GPU resources

## Action took (02/23/2024):

- The half of the A100 and H100 nodes, and three V100 nodes were moved to the gpuexpress queue, with the maximum wallclock time as 15 hours.
- Moved one dedicated interactive node back to GPU queue
- Walltime for GPU queue remains 144 hours
- Further tune as we learn more

# Annual Agreement to Cite Sinai's CTSA: New Form Launched in Jan 2024

**Goal:** To improve the acknowledgment of NIH S10 and CTSA

## **Status:**

- **On Feb 2024**, 254 user accounts disabled; 58 PIs NOT signed with their 81 Projects locked
- The template for the acknowledgement is accessible on your form page, website and Minerva message of the day.

## **Actions needed:**

A requirement for your continued access to Minerva is an annual *HIPAA Compliance Agreement* and an annual *NIH acknowledgement Agreement* (started in 2024). Please review and submit both agreements at <https://forms.hpc.mssm.edu/> by **Jan. 31, 2024**. After signing in, you will see the list of forms pending on your submission

## **Who?**

**All users and PIs on Minerva are required to complete both Agreements** irrespective of individual plans to use Protected Health Information (PHI) on Minerva and funding status.

## **Impact:**

Any PI on Minerva who has not signed the form by then will have **their projects and accounts locked** until the form is signed. Any non-PI user who has not signed the form by then will have **their accounts locked** until the form is signed.

# Expand Open OnDemand with more User-Friendly Features

**Goal:** Provide easy graphical access to Minerva without Linux command needed

**Status:** 215 users (25% of active users) are using it since its launch on Aug 2 2023!!

- We are supporting Chimera Desktop, GUI (Matlab, MarketScan, SAS, Stata), servers( VS code, Rstudio, Jupyter )

This product offers a fully-compliant job management and desktop portal requiring minimal knowledge of Linux high-performance computing (HPC) environments with no end-user installation requirements other than an up-to-date web browser (Chrome or Firefox recommended)

The service portal is accessed at URL: <https://ondemand.hpc.mssm.edu>

Documentation is available at <https://hpc.mssm.edu> Documentation>Open OnDemand or [here](#)

# New Encrypted Server for Hosting CMS Data

## Goal:

- Centralize the CMS data within a secured and compliant Minerva Ecosystem

**Status:** New encrypted storage server for CMS data in production in Jan. 2024!

- Cost: \$100 per TB per year
- Current capacity: total storage 100TB with 30TB used
- We are facilitating the SAQ (Self-Assessment Questionnaire) with Minerva CMS-compliant environment

# Minerva Shared-Web Server Updates for Improved Security

Updated procedures for getting public website for your research

1. Please fill out the request form at <https://redcap.link/g08ytzki>
2. PI review and sign the “Comprehensive Webserver Security Policy for DMZ Server v1.0” from email notification
3. Cybersecurity will scan the web application. This may take **one or two weeks.**
4. If no critical/high vulnerabilities reported, we will config the website for public access

**02/22/24** Comprehensive Webserver Security Policy for DMZ Server in effect required by Cybersecurity

- **File Restrictions:** The directory must not contain, at any time, files with non-public information, malware, spyware, viruses, unlicensed material or any form of malicious software or data.
- **File Screening:** All files uploaded to the \$HOME/www/ directory must undergo a preliminary screening process to check for PHI, malware, and other prohibited content.
- **Regular Audits:** Periodic audits of the contents of the \$HOME/www/ directory will be conducted to ensure compliance with this policy. Prohibited content shall be removed from the directory.
- **Best Practices for Webserver Security on File Encryption:** Where possible, files should be stored in an encrypted format to provide an additional layer of security.

22

**03/28/24** Disabled Directory Listing required by Cybersecurity for improved security



# TSM Archival Storage LTO-5 Tape Solution

**Updates: 35% of all archived files have been migrated from old LTO-5 tapes to new LTO-9 tapes (13 PB in total) since June 2023.**

**We are migrating everyone's files to the new tapes for maximal user convenience...Extra effort from our system admin.**

## No Cluster-Wide PM in last six months

System admins try to minimize the system-wide downtime

- Some short windows on specific servers and TSM
- Well-prepared worksheet by system admin before changes made on system

**There will be some PMs for the new installation.**





# Minerva Seven Training Sessions Spring 2024

Seven training sessions in person/Zoom this Spring with more info [here](#)

- Four additional training sessions on GPU/AI and an office hour session to lead the AI initiative at Mount Sinai, jointly presented by Minerva staff and NVIDIA domain experts

Session	Topic	Date	Time	# of Attendees
1	Minerva Intro	27-Mar	1PM	48
2	LSF Job Scheduler	3-Apr	1PM	22
3	Intro to GPU resources on Minerva	10-April	1PM	23
4	5 Ways to Get Started with GPUs	12-Apr	1PM	22
5	Accelerated general data science in medicine with RAPIDS, CuPy and Numba	17-Apr	1PM	TBD
6	Data Ark	24-Apr	1PM	TBD
7	Introduction to accelerated genomics analysis	1-May	1PM	TBD
NA	Office Hour on GPU/AI	15-May	1PM	TBD

# Documentation and Others

- Documentation updated
  - Our website at <https://labs.icaahn.mssm.edu/minervalab/>
  - We provided additional training material (including slides & recording) online
- Digital Concierge - Weekly Weds
- HPC Town Hall in person/Zoom - Twice yearly
- For most recent announcement and updates:
  - Join our mail-list: [hpcusers@mssm.edu](mailto:hpcusers@mssm.edu)
  - Minerva user group meetings will be scheduled as needed
  - Message Of The Day on Minerva

# Data Ark Data Commons Datasets

There are 17 datasets hosted under Data Ark currently

## Access within 24 hours after DUA signed

### Public Data Sets

- GTE<sub>x</sub>
- GWAS Summary Stats
- gnomAD
- eQTLGen
- BLAST
- Reference Genome
- Genebass
- 1,000 Genomes Project
- UKBB-LD
- Partial of the Cancer Genome Atlas (TCGA)
- LDSCORE

### Mount Sinai Generated Data

- The CBIPM-BioMe Data Set (coming soon)
- MSDW De-identified OMOP Data set (pending on approval)
- MSDW COVID-19 EHR Data Set
- Mount Sinai COVID-19 Biobank
- The Living Brain Project
- STOP COVID NYC Cohort

## Restricted Access

### Public Data Sets

- UK Biobank

### User Group-Acquired Data Sets

- MarketScan®



# Data Ark Dataset Updates

- ▶ Reinstated the DUA for accessing public datasets to comply with terms provided by Data Owner
  - Data Ark requested and got approval from Data Owners to redistribute those datasets in Nov. 2023
- ▶ Migrated UK Biobank phenotype data to PIs' project repositories in December 2023
- ▶ Data Ark and CBIPM IRB protocols approved: BioMe data can now be released via Data Ark
  - Announcement coming soon
- ▶ Updated MarketScan DUA with new terms for data access cost and collaborative support for users
  - After first 90 days of data access, users are required to contact Dr. Parul Agarwal on cost for continued data access
  - Integrating the dataset within Minerva environment for better computational resources
- ▶ De-identified Digital Pathology Slides
  - 1.2 million de-identified pathology images (2002-2023) are on Minerva
  - Cohort building functionality is available in MSDW Leaf
  - Images are being linked to patients' EHR data at MSDW



# Data Ark: Researcher Engagement Highlights Oct 2023 - Mar 2024

Events	Date	Attendees	Delivered Service
<b>Data Ark Training</b>	10/18/23	32 trainees	Data Ark datasets; UK Biobank data access
<b>Data Ark Town Hall</b>	10/24/23	5 trainees	Latest updates on Data Ark and roadmap ahead
<b>GGG Faculty Meeting</b>	10/25/23	20 GGS faculty	Offered Data Ark services
<b>Workshop at TCI Shared Resources Fair</b>	10/26/23	30 TCI researchers	Cancer-centric Data Ark resources
<b>Sinai Course 'AI/Data Learning in the Clinic'</b>	11/15/23	25 Graduate students and postdoctoral fellows	Briefing on datasets offered by Data Ark
<b>GGG Retreat</b>	11/17/23	100 GGS researchers	Poster presentation of Data Ark services
<b>MarketScan Data Analysis Training</b>	01/31/24	25 researchers	Introduction of MarketScan databases and considerations in data analysis
<b>Digital Concierge</b>	Weekly Wednesday	5 researchers	Real-time interactive support for Data Ark dataset access

# 2024 Roadmap

# 2024 Roadmap

## Q2

- **Deploy new TSM infrastructure with improved network performance and reliability**
- **Submit S10 proposal to expand GPUs and storage resource**
- **Continue migrating data from archival storage LTO-5 tape solution to LTO-9**
- **Implement audit logging and encryption for security compliance**
- **Update 1000 Genomes and gnomAD with new versions**
- **More Data Ark user engagement through training and research activities**
  - ◆ SCD Lunch and Learn 05/09/24
  - ◆ CTSA Research Day 06/11/24
- **Issuing a new survey on datasets @Sinai on behalf of our CTSA**
  - ◆ To better understand the diversity of data types in internal and external data sets, and the willingness of PIs to share information (metadata) about these data sets.

# 2024 Roadmap

## Q3

- **Minerva OS and software upgrade**
- **Deploy JupyterHub as a job proxy to increase users' experience on IDE software**
- **Integrate genomic data from somatic testing**
  - ◆ Results include both structured and raw genomic data
  - ◆ File formats include BAM, VCF, PDF, XML, JSON and CSV but vary by vendor and IRB required
- **De-identified Digital Pathology Slides linked to Image-EHR data**
  - ◆ User training session(s) on accessing Digital Pathology Slides by Sharon Nirenberg, MD and Cyrus Hevat, MD, PhD

## Q4

- **Deploy Minerva refreshment and expansion**
  - ◆ Preparation for the deployment of new machines including new network and management servers including NFS servers, xcat, mail server, proxy server, LSF, LDAP, etc
- **Deploy new GPFS storage for expansion**



## Acknowledgements

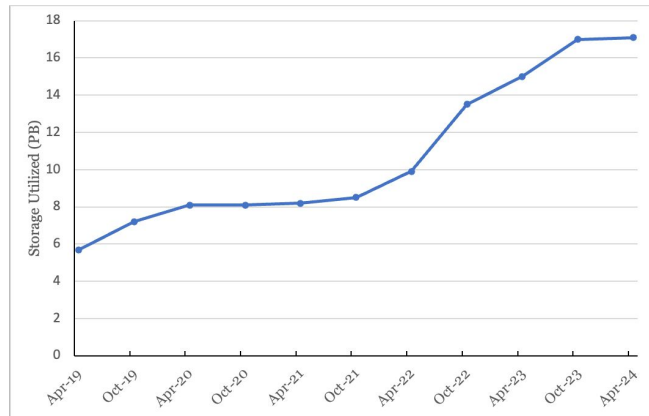
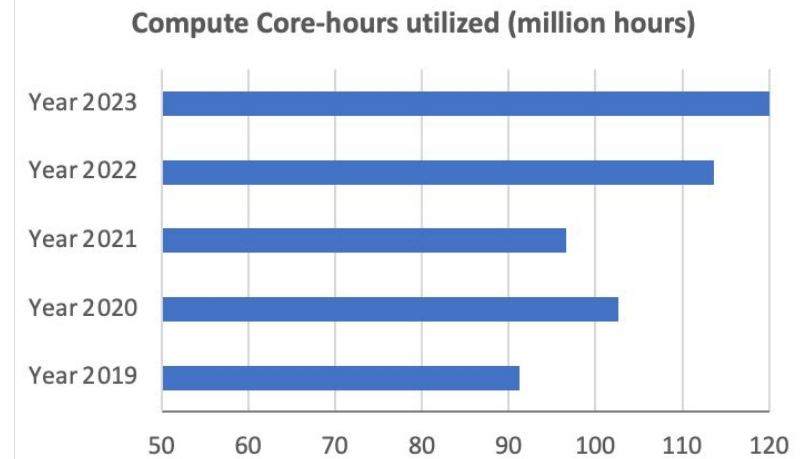
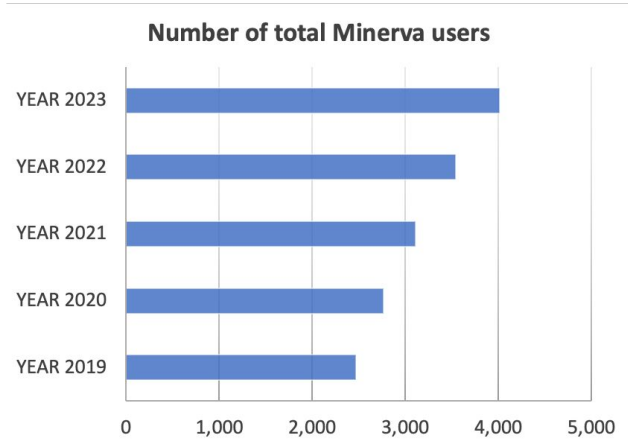
- ▶ Supported by the Clinical and Translational Science Awards (CTSA) grant UL1TR004419 from the National Center for Advancing Translational Sciences, National Institutes of Health.

**CTSA** Clinical & Translational<sup>®</sup>  
Science Awards

**Thank you!**

# **Appendix: Minerva Usage (Oct. 2023 - Mar. 2024)**

# Minerva Usage Over Years



# of users increase by 468 in 2023  
Storage usage doubled since 2022

# Minerva Usage Summary (Oct. 2023- March. 2024)

Accounts and tickets	
Number of active users	859
Number of total users	4,154
Number of project groups	557 (403 active)
Number of support tickets	2,404
Storage	
High-speed storage used (Arion)	17.1 PB (52% utilization) 7,464,377,393 Files
Archival storage used	17.9 PB
Compute	
Number of jobs run	24,119,818
Core-hours utilized	74,280,165 hrs
System	
Number of maintenance sessions	No preventative maintenance (99.6% uptime)

# Minerva publications > 1,700 since 2012!!

**We collect publications twice a year (Jan & June). Thank you!!!**

We sent email to PIs and delegate

Year	# pubs
2012	55
2013	59
2014	62
2015	115
2016	149
2017	165
2018	133
2019	178
2020	146
2021	234
2022	174
2023	239
2024	31

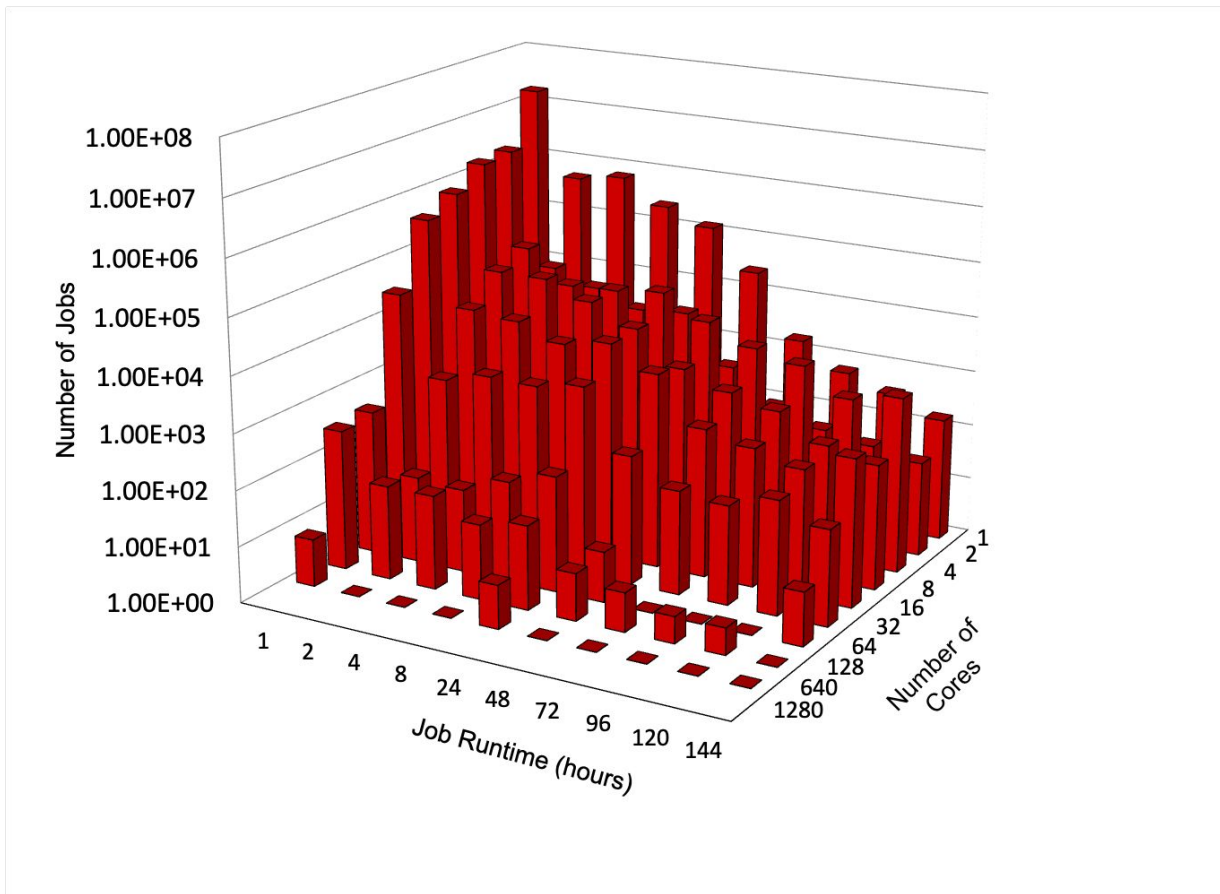
**Kovatch P**, Gai L, Cho H, Fluder E, Jiang D, Optimizing High-Performance Computing Systems for Biomedical Workloads, The 19th International Workshop on High Performance Computational Biology (HiCOMB), IPDPS, IEEE International Parallel and Distributed Processing Symposium, **May 2020**.

**Kovatch P**, Costa A, Giles Z, Fluder E, Cho H, and Mazurkova S, Big Omic Data Experience, SC'15: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, **November 2015**.

# Jobs and compute core hours by partition

Compute	# Jobs	CPU-hours	Utilization
Chimera	13,809,761	43,195,670	74.8 %
BODE2	5,849,282	10,888,357	66.7 %
Hi-memory nodes	1,569,186	6,426,001	92.4 %
CATS	2,730,022	10,823,993	67.6 %
GPU nodes	161,567	2,946,144	66.2 %
<b>Total:</b>	<b>24,119,818</b>	<b>74,280,165</b>	<b>73.1 %</b>

# Job mix





# Top 10 users compute core hours

PI	Department	# Core-hours	# Jobs
Roussos, Panos	Psychiatry	8,757,769	8,502,416
Buxbaum, Joseph	Genetics and Genomic Sciences	8,467,668	2,749,351
Pejaver, Vikas	Institute for Genomic Health	6,431,410	454,725
Goate, Alison	Genetics and Genomic Sciences	5,695,031	452,209
Raj, Towfique	Genetics and Genomic Sciences	4,571,777	1,106,970
Charney, Alexander	Genetics and Genomic Sciences	4,207,960	1,546,483
Sharp, Andrew	Genetics and Genomic Sciences	2,828,702	724,364
Luksza, Marta	Oncological Sciences	2,819,350	384,200
Schlessinger, Avner	Pharmacology	2,458,942	24,532
Zhang, Bin	Genetics and Genomic Sciences	2,340,084	159,902

# Top 10 PIs GPFS high speed storage

PI	Department	Storage usage
Zhang, Bin	Genetics and Genomic Sciences	1.3 PiB
Charney, Alexander	Genetics and Genomic Sciences	1.2 PiB
Roussos, Panagiotis	Psychiatry	1.2 PiB
Fuchs, Thomas	AI and Human Health	1.2 PiB
Raj, Towfique	Neuroscience	1.2 PiB
Sealfon, Stuart	Neurology	847 TiB
Sebra, Robert	Genetics and Genomic Sciences	834 TiB
Nadkarni, Girish Charney, Alexander	Genetics and Genomic Sciences	633 TiB
Buxbaum, Joseph	Psychiatry	615 TiB
Goate, Alison	Genetics and Genomic Sciences	560 TiB

# Top compute and storage usage department/institute

Department/Institute	Compute Core Hours
Genetics and Genomic Sciences	24,860,683
Psychiatry	19,084,279
Institute for Genomic Health	7,334,105
Neurosciences	5,233,447
Oncological Sciences	4,880,902
Neurology	3,091,138
Pharmacology	2,473,088
Medicine	1,739,956
Radiation Oncology	1,009,593
Precision Immunology Institute	886,126

Department/Institute	Storage (Tibibytes)
Genetics and Genomic Sciences	6,554
Psychiatry	1,946
AI and Human Health	1,256
Neurosciences	1,024
Oncological Sciences	998
Neurology	735
Institute for Genomic Health	288
Structural and Chemical Biology	217
Microbiology	216
Precision Immunology Institute	194

## Top 10 PIs - GPU usage hours

PI	Department	GPU hours	# Jobs
Raj, Towfique	Neurosciences	498,264	97,090
Fuchs, Thomas	AI and Human Health	463,617	8,562
Filizola, Marta	Structural and Chemical Biology	340,938	360
Schlessinger, Avner	Pharmacology	256,204	7,164
Beck, Erin	Neurology	171,306	2,093
Osman, Roman	Structural and Chemical Biology	156,176	3,096
Nadkarni, Girish	Medicine	125,666	1,529
Crary, John	Pathology	115,923	1,000
Davies, Terry	Medicine	115,309	113
Shen, Li	Neuroscience	82,536	1,690

# Total TSM archival storage usage (Oct 2023- Mar 2024)

Current archive storage usage	
Archived data	17.9 PB (LTO5: 8.4 PB, LTO9: 9.5 PB)
Total data with offsite copy	38.1 PB (LTO5: 19.1 PB, LTO9: 19.0 PB)
Number of tapes used	15,627 (14,390 LTO5 + 1,237 LTO9)

Statistics of Oct 2023 - Mar 2024			
Amount of archived data	2,061 TB	Amount of retrieved data	366 TB
# of users who have issued archive commands	64	# of users who have issued retrieve operations	39

LTO5 to LTO9 Migration			
Amount of data migrated	4.5 PB	Est data remaining max	8.4 PB
# tapes decommissioned	2,432	Est migration completion	Q2 2025

# Data Ark Usage Summary

# of active unique users: **52**

# of support tickets: **159**

- gnomAD has replaced TCGA as the #1 accessed dataset

Dataset	Size (GB)	# of unique users	# of times data accessed
gnomAD	8,628	19	11,925,250
TCGA	154	14	5,278,122
UK Biobank LD	2,866	12	157,956
Genebass	903	8	146,136
UK Biobank	12,695	15	84,326
GWAS Summary Statistics	6,826	11	78,868
1000 Genome	143	18	75,226
LD Score Regression	173	9	10,288
GTE <sub>x</sub>	1,888	13	6,746
Reference Genome	142	15	3,624
MSDW OMOP	3,076	3	2,441
MSDW Covid	1	2	1,993
Blast	1,116	10	1,974
eQTLGen	39	9	260

# Pathology Slide Cohort Builder Already Available

The screenshot displays the Leaf Pathology Slide Cohort Builder interface. At the top, the 'leaf' logo is on the left, and the status 'Unsaved Query' with '0 patients' is on the right. A search bar is located below the logo. The main content area is divided into a left sidebar and a right panel. The sidebar lists various categories: 'Lab Results & Measurements (LOINC)', 'Medications (ATC)', 'Patient Cohorts' (with a sub-count of 1,779,462), 'Procedures (CPT4)', and 'My Saved Cohorts'. Under 'Patient Cohorts', several cohorts are listed with their respective patient counts. The 'Digitized Pathology Slides Cohort' is highlighted with a red box and has a 'Learn More' button next to it. The right panel features a 'Run Query' button and a 'Limit to' section with three filter boxes. Each filter box contains the text 'Patients Who', 'Anytime', and 'At Least 1x'.

leaf

Unsaved Query  
0 patients

All Concepts... Search...

Lab Results & Measurements (LOINC)

Medications (ATC)

Patient Cohorts 1,779,462

- BioMe Biobank 45,212
- BioMe Biobank Global Diversity Array (Sema4) 15,570
- BioMe Biobank Global Screening Array (Regeneron) 23,463
- BioMe Biobank Whole Exome Sequencing (Regeneron) 22,778
- Cancer Institute Biorepository 13,456
- Cancer Patient Cohort 254,938
- Dental Patient Cohort 80,464
- Digitized Pathology Slides Cohort 107,936**
- Imaging Research Warehouse 1.0 466,836
- Imaging Research Warehouse 2.0 1,554,814

Learn More

Procedures (CPT4)

Need Help? Vitals 3,315,306

Run Query

Limit to

Patients Who  
Anytime  
At Least 1x

And  
Anytime  
At Least 1x

And  
Anytime  
At Least 1x