

Minerva Town Hall

Dec 2020

Patricia Kovatch
Dansha Jiang, PhD
Wayne Westerhold
Wei Guo, PhD
Eugene Fluder, PhD
Hyung Min Cho, PhD
Lili Gai, PhD

Dec 16, 2020



**Mount
Sinai**

Outline

Welcome and general comments

- 2020 Minerva Usage
- 2020 Accomplishments
 - Staffing
 - Minerva expansion
 - Deprecated older storage hardware (Hydra, Orga)
 - Less PM on Weekends
 - COVID and Translational Science (CATS) S10 proposal
 - HIPAA compliance
 - New service on Rstudio/Jupyter
 - Mount Sinai Data Commons - Data Ark
 - Training and documentation

2021 Initiatives and Roadmap

- Globus
- SLURM
- Visualization portal on web browser
- CATS
- New service on database
- Deploy new Electronic Health Record (EHR) OMOP data warehouse
- User survey

Q&A



2020 Minerva Usage (Jan - Sep)

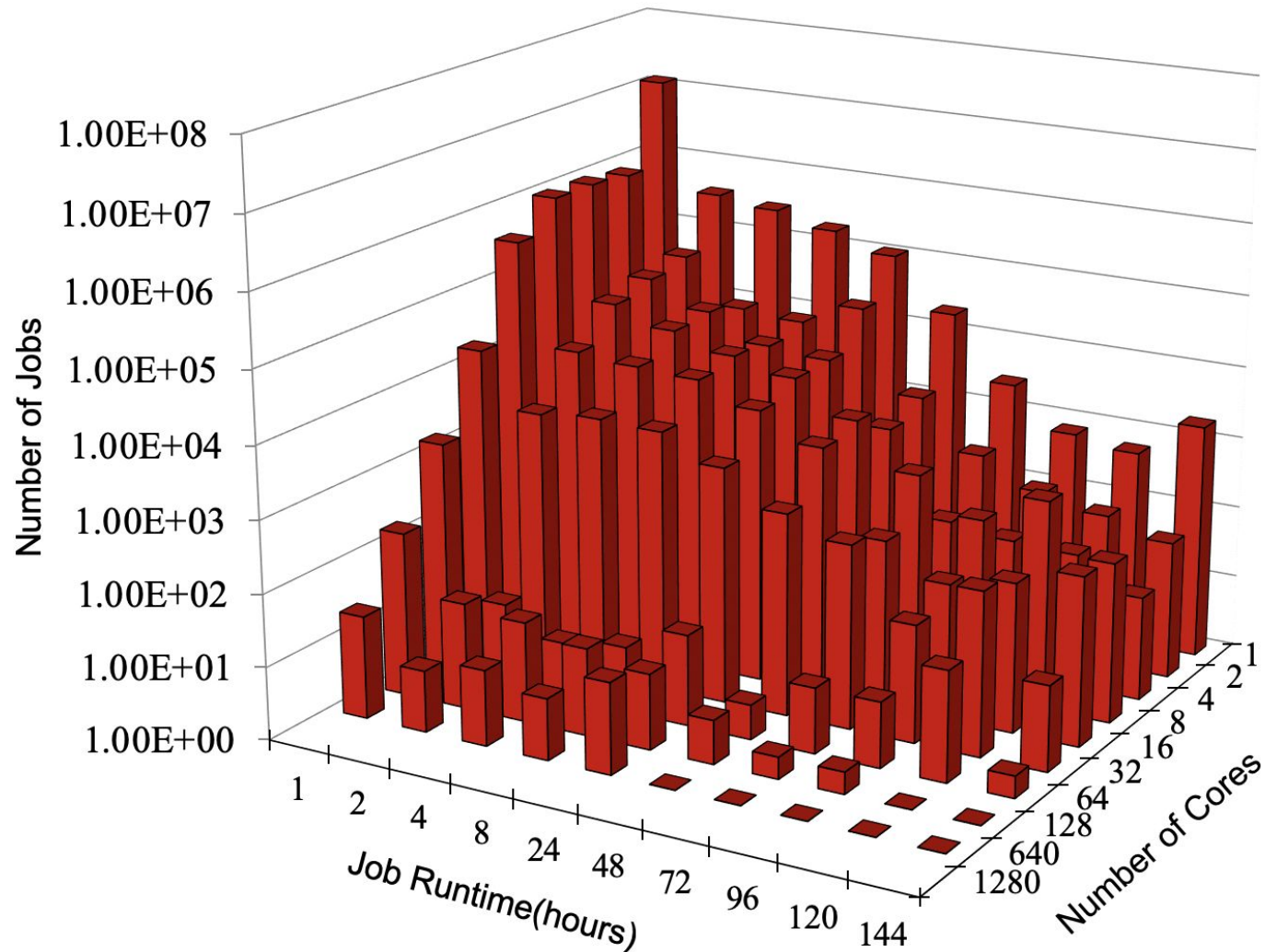
2020 Minerva usage summary (Jan-Sep)

Accounts	
Number of new users	227
Number of active users	714 (61 external users)
Number of total users	2,700
Number of project groups	344
Storage	
High-speed storage used (Hydra)	2.1 PiB (61% utilization)
3.5 PiB Usable in total	427,805,539 files
High-speed storage used (Arion)	6.2 PiB (65% utilization)
9.6 PiB Usable in total	1,616,501,226 files
Archival storage used	10.9 PiB
Compute	
Number of jobs run	16,969,087
Core-hours utilized	73,584,777 hrs
System	
Number of maintenance sessions	2 preventative maintenance (99% uptime)

Jobs and compute core hours by partition

Compute	# Jobs	CPU-hours	Utilization
Chimera	13,114,132	56,866,011	67.5 %
BODE2	3,567,153	14,523,779	60.6 %
Hi-memory nodes	14,458	871,538	71.1 %
GPU nodes	75,740	1,043,191	42.7 %
Total:	16,771,483	73,304,519	65.5 %
		(Jan - Jun)	57.1 %
		(Jul - Sep)	83.5 %

Job Mix



Top 10 users by core hours

PI	Department	# Core-hours	# Jobs
Panagiotis Roussos	Psychiatry	7,818,881	4,141,179
Gaurav Pandey	Genetics and Genomic Sciences	6,801,132	1,361,455
Zhongyang Zhang	Genetics and Genomic Sciences	6,347,886	1,361,455
Andrew Sharp	Genetics and Genomic Sciences	5,601,767	1,234,072
Bin Zhang	Genetics and Genomic Sciences	5,053,139	147,724
Joseph Buxbaum	Psychiatry	3,802,053	106,898
Gang Fang	Genetics and Genomic Sciences	2,763,749	202,632
Samir Parekh	Medicine	2,460,706	59,636
Marta Filizola	Structural and Chemical Biology	1,691,761	348,782
Judy Cho	Genetics and Genomic Sciences	1,673,491	502,965

Top 10 users by GPFS high speed storage

PI	Department	Storage usage
Panagiotis Roussos	Psychiatry	1,300 TiB
Bin Zhang	Genetics and Genomic Sciences	1,000 TiB
Robert Sebra	Genetics and Genomic Sciences	805 TiB
Zhongyang Zhang	Genetics and Genomic Sciences	463 TiB
Judy Cho	Genetics and Genomic Sciences	362 TiB
Samir Parekh	Oncological Sciences	293 TiB
Alison Goate	Neurosciences	251 TiB
Dalila Pinto	Genetics and Genomic Sciences	203 TiB
Towfique Raj	Neurosciences	152 TiB

Compute and storage usage by department/institute

Department/Institute	Compute Core Hours
Genetics and Genomic Sciences	40,547,333
Psychiatry	13,065,085
Oncological Sciences	3,436,530
Structural and Chemical Biology	2,908,027
Neurosciences	1,222,064
Medicine	813,695
Microbiology	693,592
Preventive Medicine	502,948
Pharmacology	381,971
Neurology	323,931

Department/Institute	Storage (TiB)
Genetics and Genomic Sciences	4,426
Psychiatry	1,500
Oncological Sciences	525
Neurosciences	515
Neurology	142
Medicine	118
Surgery	113
Structural and Chemical Biology	111
Mindich Child Health and Development	77
Microbiology	65

TSM Archival Storage Usage

Current archive storage usage	
Archived data	10.9 PiB
Total data with offsite copy	21.8 PiB
Number of tapes used	15,629

Statistics of 2020 Jan-Sep			
Amount of archived data	1.2 PiB	Amount of retrieved data	115 TiB
# of archive operations	27,552	# of retrieve operations	3,996
# of users who have issued archive commands	103	# of users who have issued retrieve operations	60

2020 Minerva Accomplishments

2020- Accomplishments Summary

Thank you very much for the feedback from user survey!

Actions we took in 2020 (in response to the user survey and our last roadmap):

- ✓ Surpassed over **1,000** publications that utilized Minerva!!
- ✓ Hired 2 system admins: Wayne Westerhold and Wei Guo, PhD
- ✓ Expanded Minerva with more high memory nodes, GPUs, and more storage
- ✓ Deprecated older storage hardware (Hydra, Orga)
- ✓ Less PMs (5 in total) on weekends
- ✓ Submitted \$2M NIH S10 proposal for COVID and Translational Science supercomputer (CATS)
- ✓ Received HIPAA-compliance for Minerva starting Oct. 1 2020
- ✓ Implemented new services for Rstudio IDE, Rstudio-connect/Shiny, Jupyter
- ✓ Mount Sinai Data Commons - Data Ark on Minerva
- ✓ Data transfer node implemented: **data-xfer.hpc.mssm.edu**
- ✓ Updated the documentation and presented four tutorial sessions
- ✓ User support: Continued to support Minerva users through ticketing system (closed **3,290** tickets in 2020) and in-person meetings

Details will be presented in the following slides.

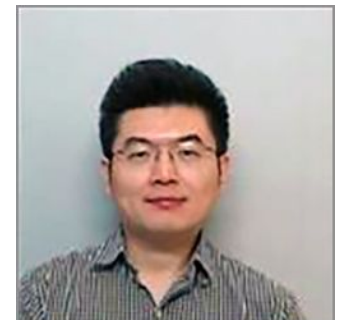
Staff

The HPC team consists of three computational scientists...

- Lili Gai, PhD
- Eugene Fluder, PhD
- Hyung Min Cho, PhD

...and three HPC admins:

- Dansha Jiang, PhD (Dansha is leaving Jan 2021. We have the position open and is looking for candidate)
- **Two vacant positions filled !**
 - **Wayne Westerhold**, joined this Sept. as an HPC Admin with an expertise in parallel file systems and performance. He is coming to us from Lenovo.
 - **Wei Guo, PhD**, joined as an HPC Architect this Oct. He most recently worked at St Jude Children's Hospital and has a PhD in Materials Science and Engineering.



Expanded Minerva with funds from a donor

All in production

33 high memory nodes (1.5 TB of memory per node)

- 48 Intel Xeon Platinum 8268 2.9 GHz Processors per node, for a total of 1,440 cores

32 A100 GPUs in 8 nodes

- 48 Intel Xeon Platinum 8268 2.9 GHz Processors per node, 384 GB memory per node, **1.92 TB SSD per node**, 4 A100 GPUs per node, 40 GB of memory on GPU, for a total of 384 cores, 32 A100 GPUs
- Cuda 11.x supported, local SSD available, and [more instructions here](#)

11.6 petabytes of usable storage (16 PB raw)

- Has been integrated into Arion file system, for a total of 21.2 PB of usable storage

Minerva has a total of 437 nodes for a total of 20,976 cores, and 21 petabytes of usable storage currently.

Minerva file system has migrated to Arion

Orga (/sc/orga/)

- Data migrated to Arion
- Retired in April 2020

Hydra (/sc/hydra/)

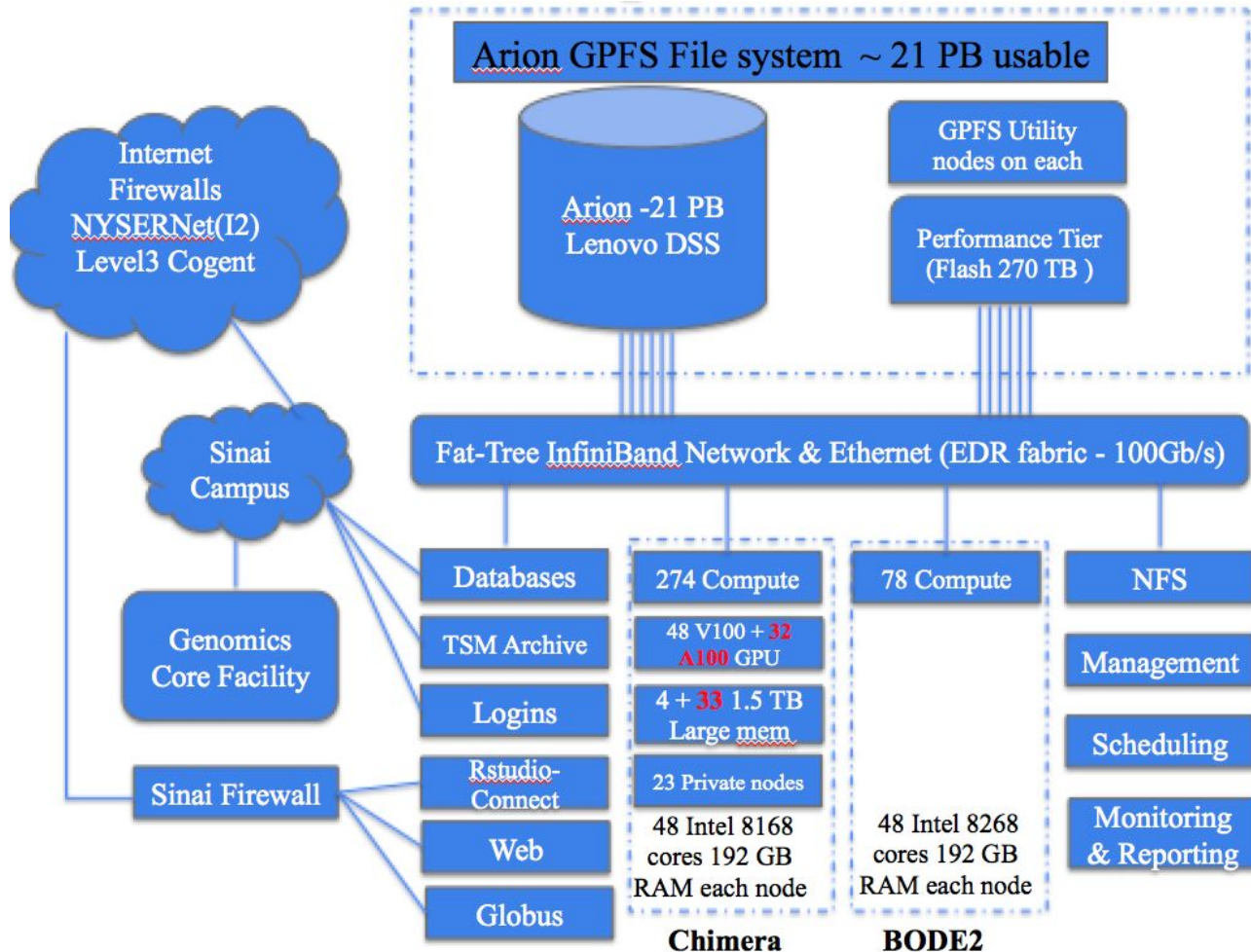
- Data migrated to Arion
- Will be retired by the end of year

Arion (/sc/arion/)

- The high speed online storage on Minerva, with a total of 21 PB of usable storage
- Apply for storage allocation [here](#)

Updated Minerva Diagram

- We now have a total of 437 compute nodes with a total of 20,976 compute cores with 21 PB of high speed online storage and have deprecated almost all out of warranty storage



Minerva PM on weekends

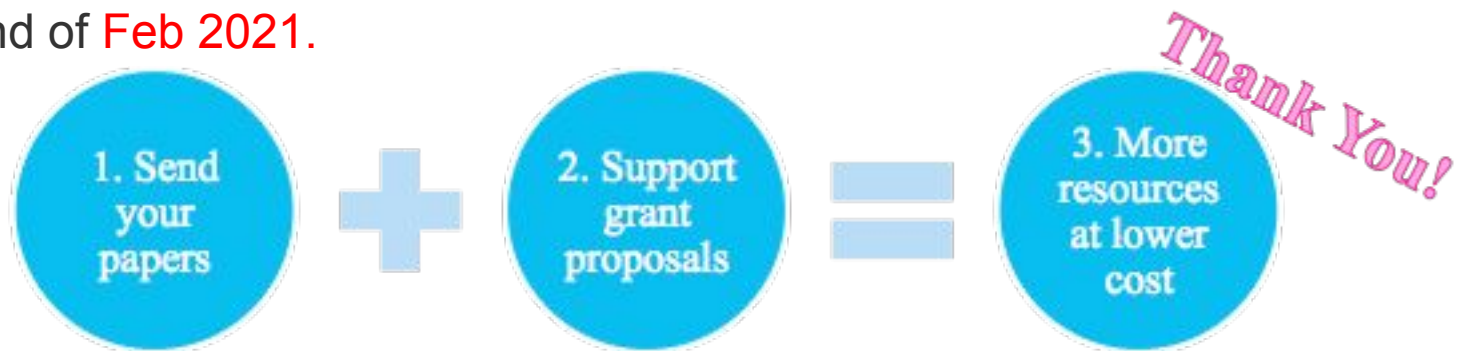
5 full PMs this year on weekends

- Thanks to our system admins Dansha, Wei and Wayne, we performed all the PMs this year on weekends to accommodate users' request.
- We will try our best to accommodate this in 2021, but not guaranteed.



\$2M COVID and Translational Science (CATS) S10 proposal

- To enable new kinds of scientific discovery and translation related to the newly emerged COVID-19 pandemic, we requested a new high-performance instrument with large shared memory nodes, and more storage to reduce your cost
- If awarded, CATS will contain with 2,640 Intel cores, 82 TB of memory and 16 PB of raw storage. **And your storage cost will be reduced again!**
- We received a score of **22** from the reviewers, noting that the proposal was “outstanding” and will have a “high impact”. Final decision will be made at the end of **Feb 2021**.



Please keep sending us your publications! Thank you

Minerva is HIPAA-compliant starting Oct 1 2020!

- Minerva is HIPAA compliant on October 1st, 2020, i.e., Protected Health Information (PHI) data will be allowed to be stored and processed on Minerva.
- All users have to read the HIPAA policy and complete Minerva HIPAA Agreement Form at <https://labs.ica hn.mssm.edu/minervalab/hipaa/> annually.
- Any user who has not signed the agreement gets their accounts locked until the agreement is signed.

Rstudio Connect server

- Rstudio Connect server deployed this Aug, with a license for 20 users
- You can publish **Shiny, R Markdown, Jupyter Notebook** for collaborators or others to Minerva Rstudio Connect server (More details on rstudio connect [here](#))
- If interested, please contact us at hpchelp@hpc.mssm.edu.

Content / 012-datatables

Columns in diamonds to show:

- ☒ carat
- ☒ cut
- ☒ color
- ☒ clarity
- ☒ depth
- ☒ table
- ☒ price
- ☒ x
- ☒ y
- ☒ z

diamonds mtcars iris

Show 10 entries Search:

	carat	cut	color	clarity	depth	table	price
1	1.24	Premium	D	SI1	62.4	59	7486
2	1.2	Premium	G	VS2	62.1	61	7728
3	0.73	Very Good	F	SI1	59.7	60	2473
4	1.53	Premium	I	SI1	61.5	59	8911
5	0.3	Premium	D	SI1	62.1	59	515
6	0.58	Ideal	H	VS1	61.2	55	1671
7	0.51	Ideal	E	SI2	61	56	1098
8	1.5	Ideal	G	VVS2	61.3	56	17176
9	2.66	Good	H	SI2	63.8	57	16239
10	0.3	Premium	F	VVS2	61.4	59	737

Showing 1 to 10 of 1,000 entries

Previous 1 2 3 4 5 ... 100 Next

Info Access Runtime Schedule Tags Vars Logs

Who can view this application

You

Who can change this application

Lili Gai gail01

Add collaborator

Who runs this content on the server

The default user rstudio-connect

Content URL

/hpcshowcase/

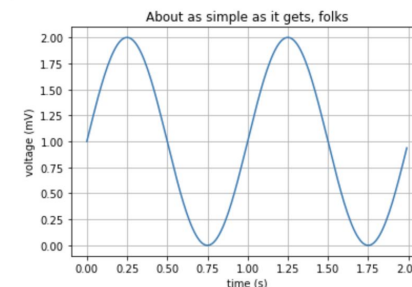
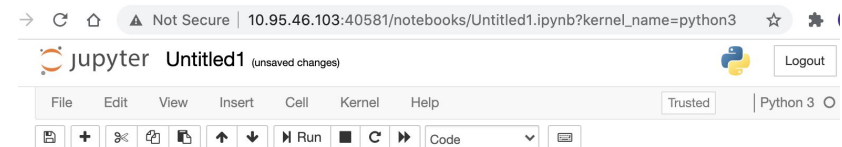
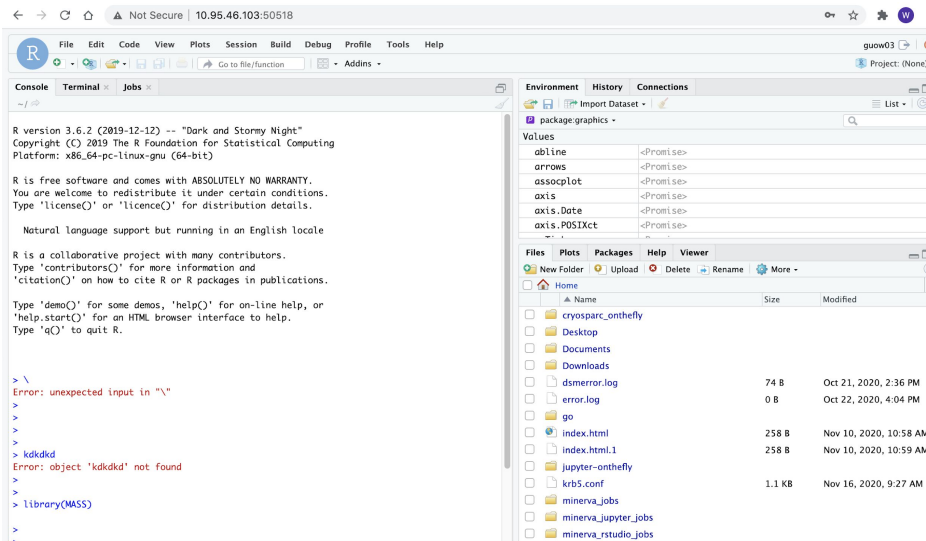
<https://rstudio-connect.hpc.mssm.edu> Copy

On-the-fly RStudio and Jupyter Notebook in a Minerva job

- One simple command to get interactive web sessions in a HPC job
- Containerized application for workflow reproducibility, packages installed in \$HOME
- Available on login nodes only
- See usage:

- `minerva-rstudio-web.sh -h`
- `minerva-jupyter-web.sh -h`

```
pip install numpy --user
```



Data Ark: Mount Sinai Data Commons

- ▶ Goal: Increasing **genomic data accessibility** and **reusability**
- ▶ Space on Minerva for Mount Sinai and publicly available data sets (version 1):
 - **UK Biobank:** 500k genotype-phenotype data, with 150k WES soon
 - **1000 Genomes:** WGS data on 1,000 individuals of European ancestry
 - **GTEx:** Gene expression data on hundreds of individuals across 50 tissues
 - **Post-mortem brain gene expression data:** combines three large Mt Sinai gene expression data sets, providing the largest available data set anywhere
 - **Mount Sinai COVID Electronic Health Record (EHR) data:** 350 phenotypic data elements from all Mount Sinai patients diagnosed with or under suspicion for COVID; updated daily
- ▶ Incentives for data sharing
 - Data-generating PIs and contributors will be offered **authorship**
 - Boost in data use/publication
 - **Credit** will be given by the faculty Appointment Promotions and Tenure Committee
 - **Free QC**, standardized processing, management and descriptive analysis of shared data
 - **Increased data quality** due to multiple analysts interrogating it
- ▶ Strict terms of use
 - Process for requesting access

Documentation updates and training sessions

- For most recent announcement and updates:
 - Join our mail-list: hpcusers@mssm.edu
 - Minerva user group meetings will be scheduled as needed
 - Message Of The Day on Minerva
- Offered four training sessions this year via Zoom:
 - Two sets of training sessions in spring and fall
 - Topics include “Introduction to Minerva” and “LSF job scheduler”
 - Spring sessions at March 11 and March 18 2pm -3pm
 - Fall sessions at September 16 and September 23 2pm -3pm
- Documentation updated on the website:
 - New website at <https://labs.icaahn.mssm.edu/minervalab/>
 - We will add newer pages/articles as needed
 - We provided additional training material (including slides) online

2021 Initiatives and Roadmap

What's next for 2021?

Deploy Globus under HIPAA+BAA subscription, making data share more easily

- Be able to share data with anyone using their identity or their email address
- Be able to upgrade users' endpoint to globus Plus as needed
- We will disable other forms of outgoing data transfer over time

Consider the SLURM job scheduler

- We will evaluate it to save budget--timeline for deployment
- We will support cloud bursting going forward

Visualization Portal

- Set up Open OnDemand as better visualization portal to access Minerva through modern web browsers

COVID And Translational Science Supercomputer - S10

- If awarded, we will install more resources in next spring!!

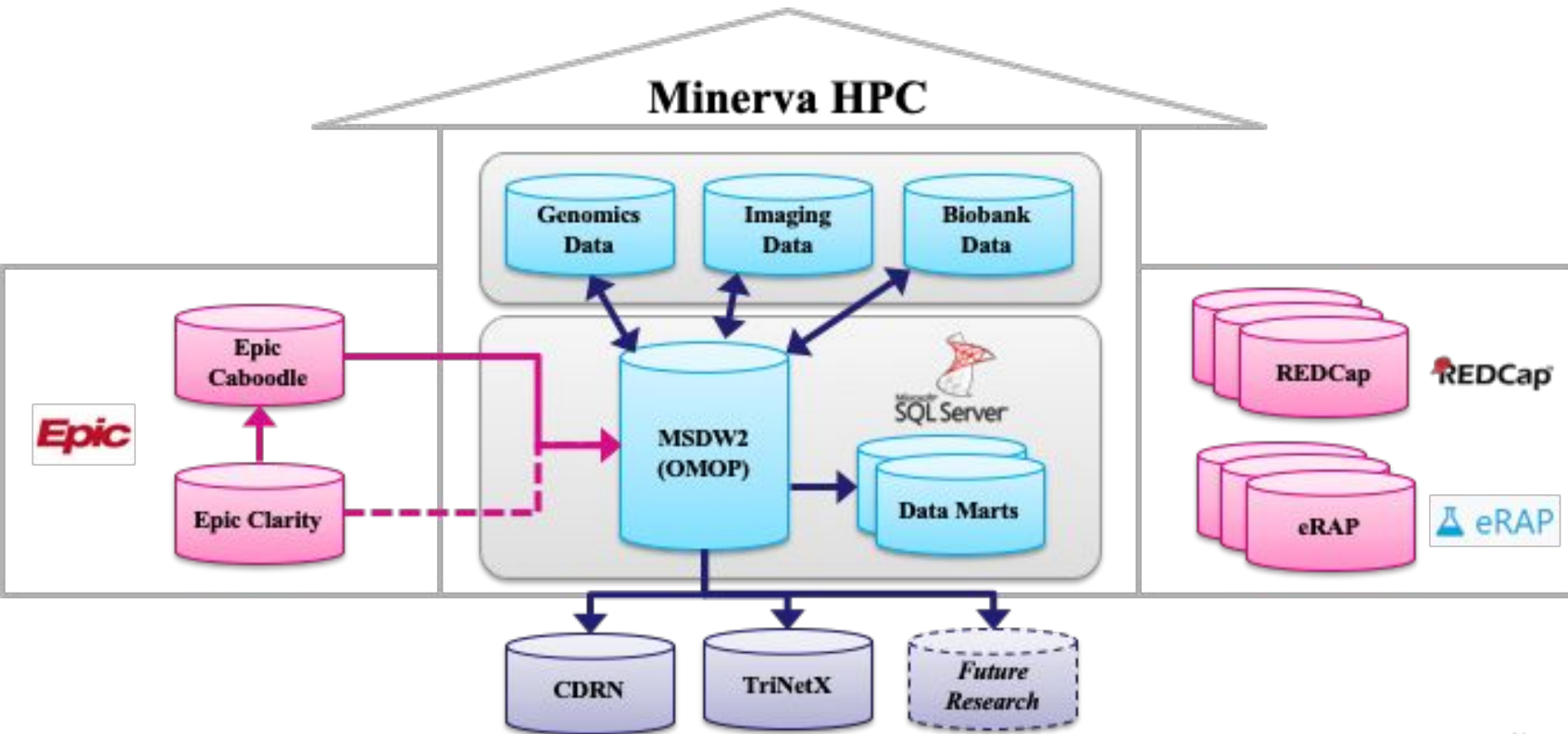
New services such as MongoDB user databases

Documentation

Deploy new Electronic Health Record (EHR) OMOP data warehouse

- Will contain de-identified and identified EHR data
- Will connect to image research warehouse

Minerva HPC Ecosystem & MSDW2 Database



Will issue 2020 user survey in Jan 2021

We will ask the following 4 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: General suggestions for service improvement.

Open-ended

Question and comments

Thank you!