



JAWS - Scientific Workflow Execution Service

Seung-Jin Sul
ssul@lbl.gov
June 10, 2026

About us

- High-throughput AI-centric genomic science user facility located at Lawrence Berkeley National Laboratory
- Provides the genomic capabilities, data, and expertise that supports the global research community in studying complex biological and environmental systems

Advanced genomic capabilities

Large-scale sequencing and synthesis

AI-ready scientific data

Curated data products for reuse

Professional expertise

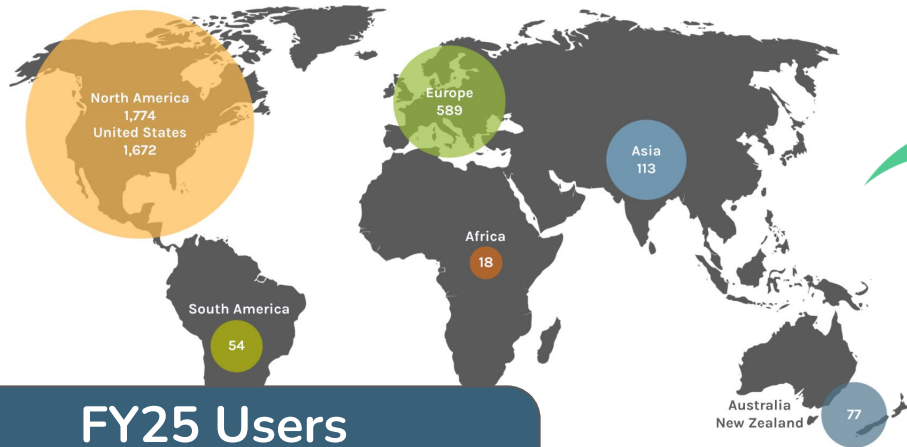
Support from domain and workflow experts



This mural at our site in Berkeley, California, shows some of the species that JGI users have studied.

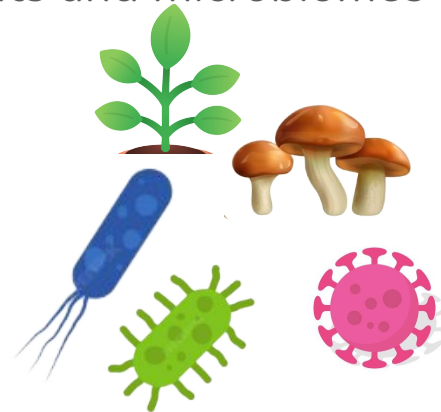


Big Data, Bigger Science, and Workflows



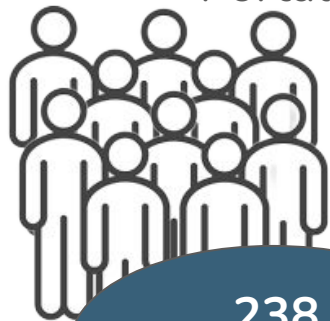
FY25 Users
Primary: 2,627
Secondary: > 18,000

Primary Users provide unique samples, from fungi, plants and microbiomes as part of their studies



Primary and Secondary Users leverage data through JGI

Portals



238 publications

~22M Files

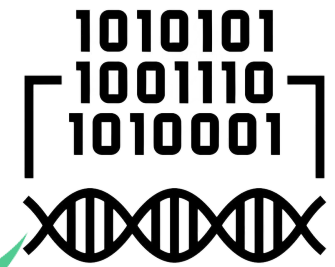
~20 PB Data

FY25

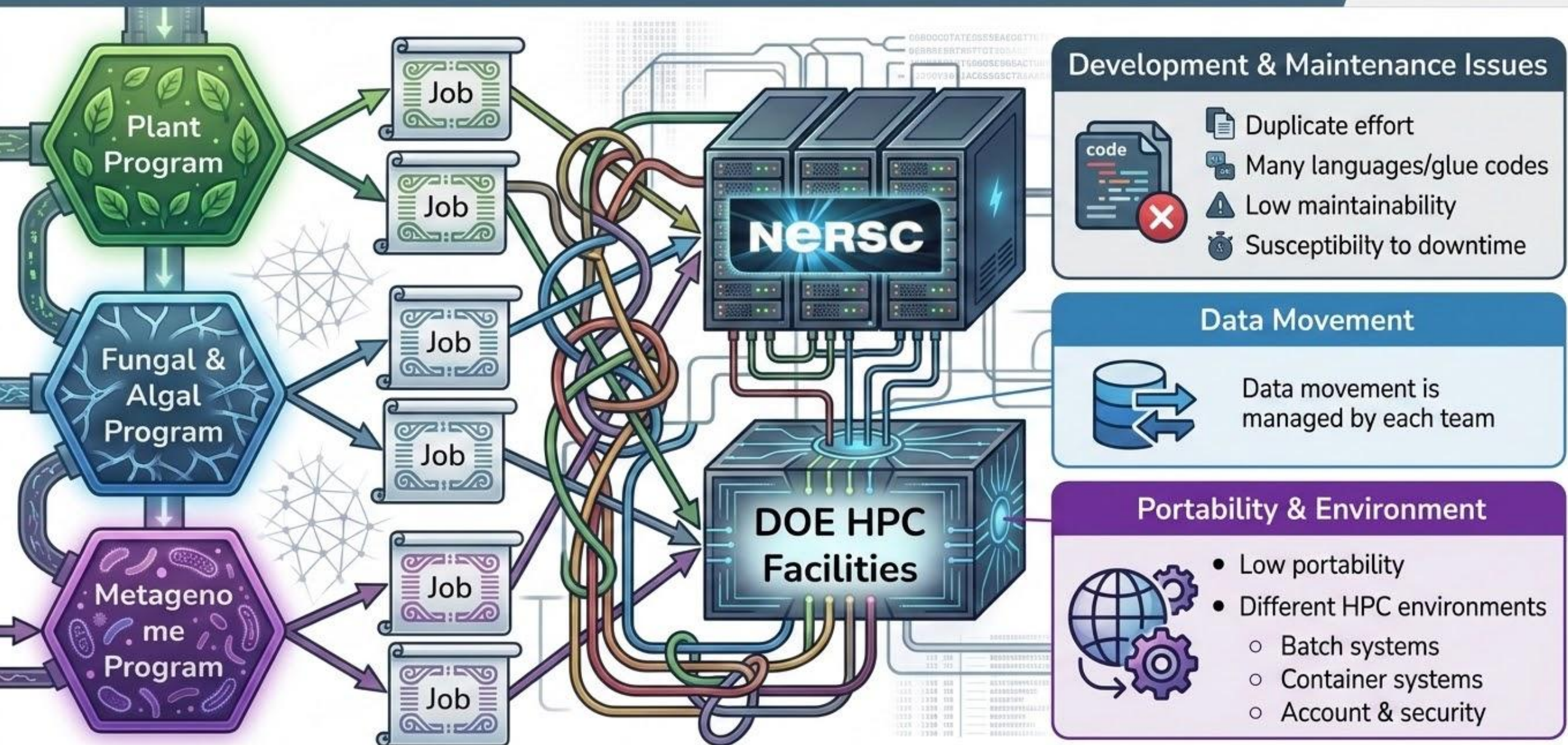
- Generated 1.14 Petabases of sequence
- Synthesized 8.53 Megabases of DNA
- Processed 17,179 metabolomics samples

- JGI Data Portal
- JGI MycoCosm THE FUNGAL GENOMICS RESOURCE
- JGI PhycoCosm THE ALGAL GENOMICS RESOURCE
- JGI Phytozome 14 THE PLANT GENOMICS RESOURCE
- JGI IMG/M INTEGRATED MICROBIAL GENOMES & MICROBIOMES
- JGI SMC SECONDARY METABOLISM COLLABORATORY

The sample becomes data



Fragmented Science: Porting is Hard, Availability is Necessary



* NERSC: DOE HPC center at LBNL

- Developed a workflow management service called JGI Analysis Workflow Service (JAWS)
- Unifies workflows across JGI groups and DOE facilities
- Users write their workflows in Workflow Description Language (WDL)
- **JAWS** takes care of moving data, handling compute backends of multiple facilities, and executing workflows under the hood
- Write once, run anywhere



User writes WDL; JAWS handles the rest

JAWS Client

- CLI / API / REST
- WDL workflow submit, status, cancel, logs, ...

JAWS Central & Sites

- File Transfer: Globus
- Protocols: HTTP and AMQP

Cromwell

- WDL execution engine
- Creates task DAG and HTCondor submit files
- Submits jobs to HTCondor

* DAG: directed acyclic graph

HTCondor

- Job scheduling and execution system between Cromwell and HPC clusters

JAWS Pool Manager

- HTCondor glidein pool manager
- Dynamically scales site compute pools
- Communicates with SLURM, PBS (AWS, k8s)

```
version 1.0
```

```
workflow hello {  
  call say_hello { }  
  output {  
    File greeting = say_hello.out  
  }  
}
```

```
task say_hello {  
  input {  
    String message= "Hello!"  
  }  
  command <<<  
    echo ~{message} > final_message.txt  
>>>  
  output {  
    File out = "final_message.txt"  
  }  
  runtime {  
    docker: "ubuntu:22.04"  
    cpu: 1  
    memory: "2G"  
    runtime_minutes: 3  
  }  
}
```



Workflow Description Language "Recipe"

func Workflow Call



Define task



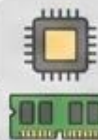
Define inputs type



Series of commands

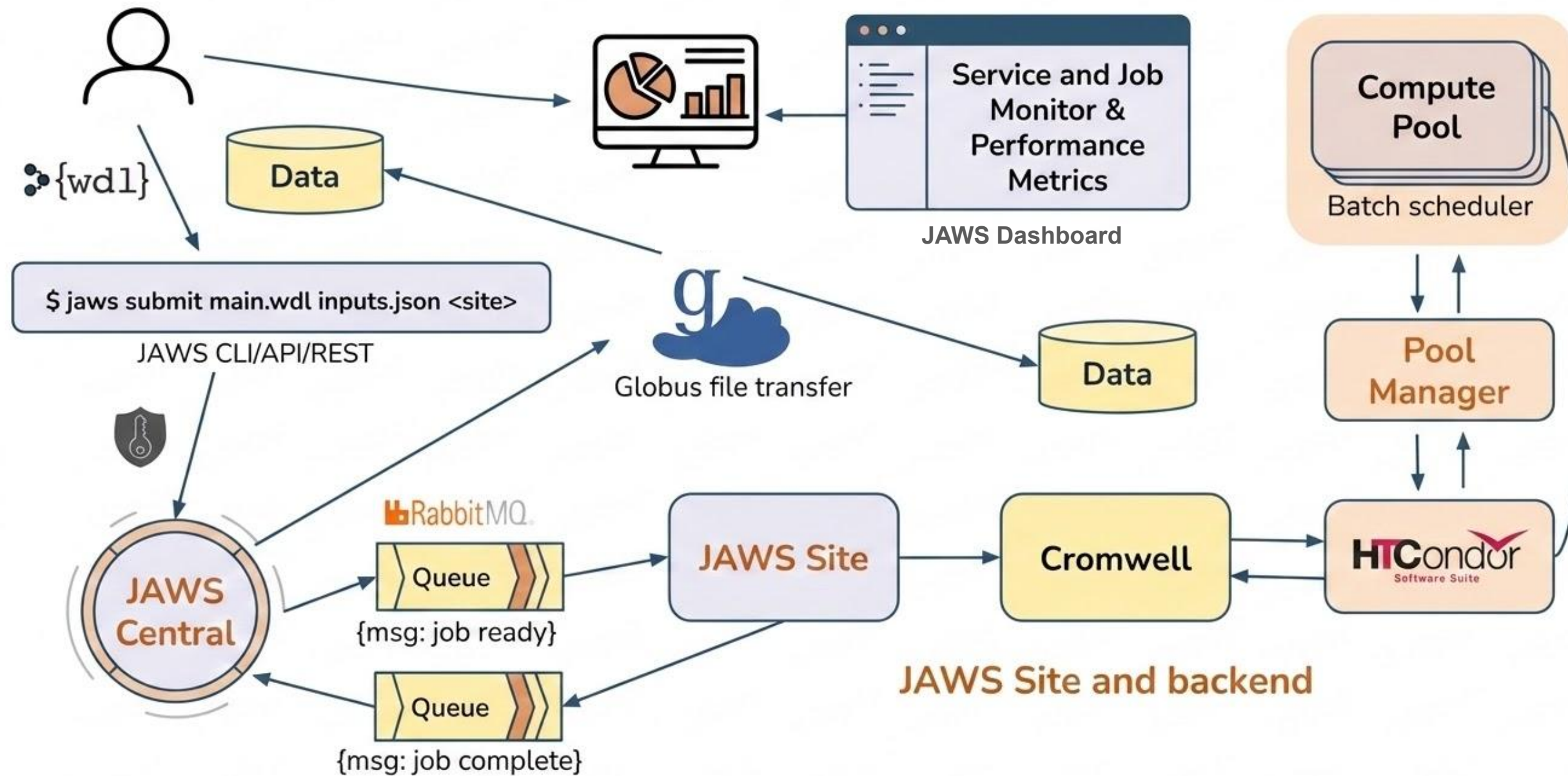


Save results → outputs

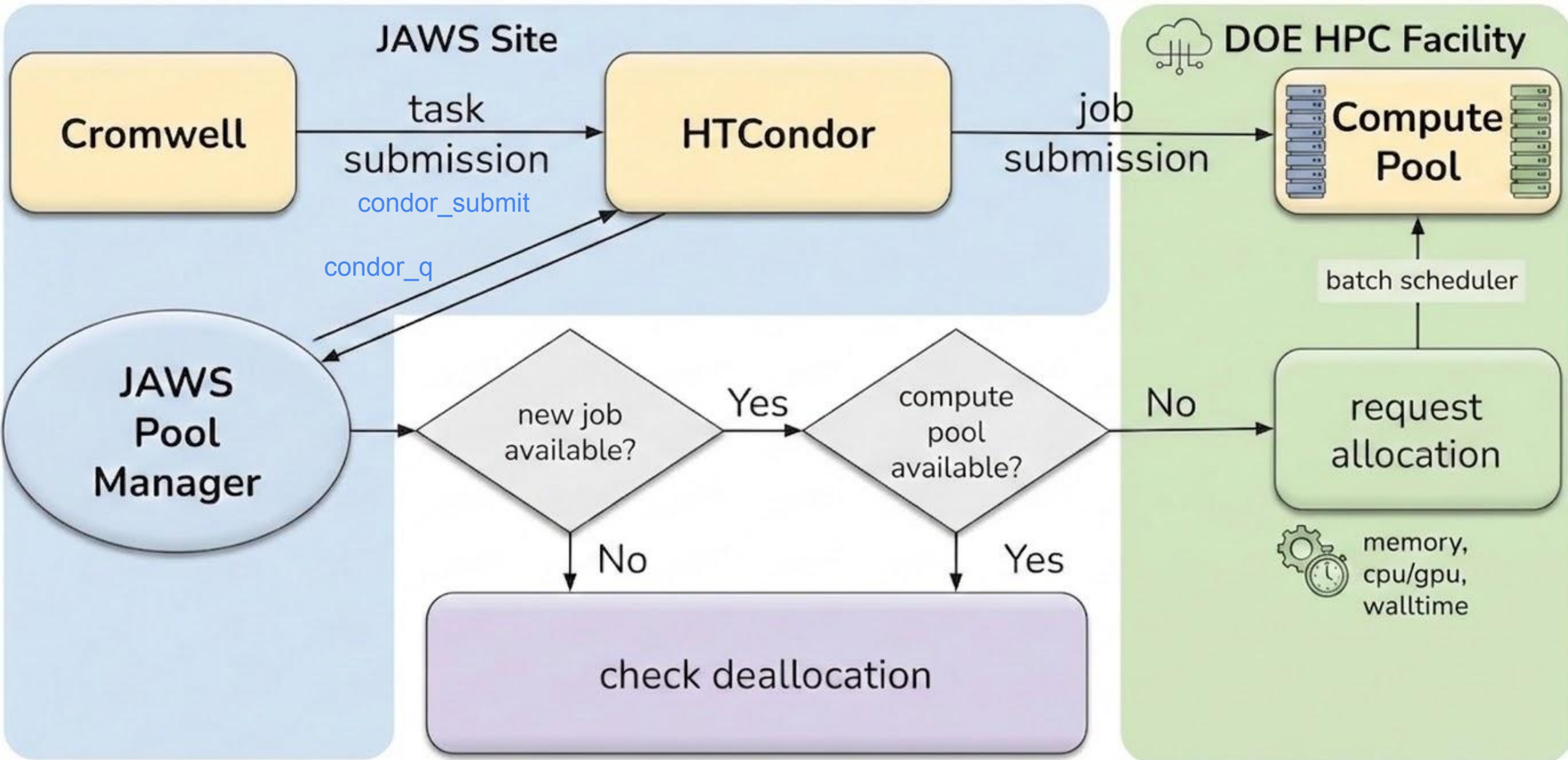


Compute resources
requirement

JAWS: You Bring the Science, We Manage the Complexity



JAWS Pool Manager and HTCondor Glidein



- **Pool types**

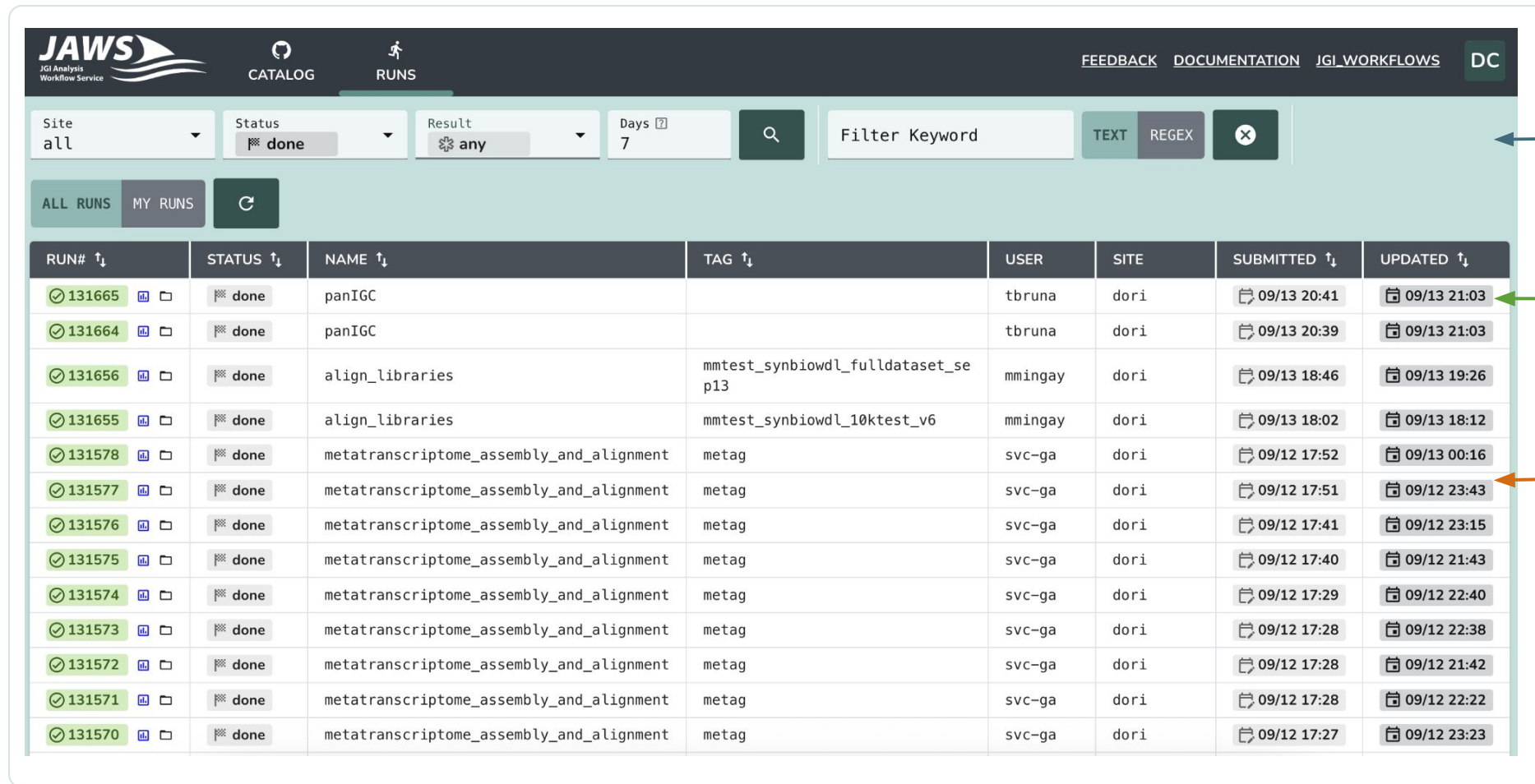
- Five pool types based on the cpu/gpu & memory size
- Each site has different set of node types
 - ex) Perlmutter: gpu and large (~500GB)
Dori: small (debug), large and xlarge (~1.5TB)
- Min/max pool sizes per pool type and site
 - ex) dori: 10-node large pool

- **Scheduling Policies**

- Walltime limit for job scheduling
- Policies for housekeeping
- Slot defragmentation



Operational visibility helps users, support teams, and workflow developers keep large runs moving



The screenshot shows the JAWS Dashboard interface. At the top, there is a navigation bar with the JAWS logo, 'CATALOG', and 'RUNS' tabs. On the right, there are links for 'FEEDBACK', 'DOCUMENTATION', 'JGI_WORKFLOWS', and a 'DC' button. Below the navigation bar, there is a search and filter section with dropdown menus for 'Site' (set to 'all'), 'Status' (set to 'done'), and 'Result' (set to 'any'). There is also a 'Days' field set to '7' and a search input field labeled 'Filter Keyword'. Below this, there are buttons for 'ALL RUNS' and 'MY RUNS', along with a refresh icon. The main content is a table of workflow runs with the following columns: RUN# (with up/down arrows), STATUS (with up/down arrows), NAME (with up/down arrows), TAG (with up/down arrows), USER, SITE, SUBMITTED (with up/down arrows), and UPDATED (with up/down arrows). The table contains 13 rows of data, all with a status of 'done'. The first two rows are for 'panIGC' runs, and the remaining 11 rows are for 'align_libraries' and 'metatranscriptome_assembly_and_alignment' runs. The 'Submitted' and 'Updated' columns show timestamps in YYYY-MM-DD HH:MM:SS format.

RUN# ↑↓	STATUS ↑↓	NAME ↑↓	TAG ↑↓	USER	SITE	SUBMITTED ↑↓	UPDATED ↑↓
131665	done	panIGC		tbruna	dori	09/13 20:41	09/13 21:03
131664	done	panIGC		tbruna	dori	09/13 20:39	09/13 21:03
131656	done	align_libraries	mmtest_synbiowdl_fulldataset_se p13	mmingay	dori	09/13 18:46	09/13 19:26
131655	done	align_libraries	mmtest_synbiowdl_10ktest_v6	mmingay	dori	09/13 18:02	09/13 18:12
131578	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:52	09/13 00:16
131577	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:51	09/12 23:43
131576	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:41	09/12 23:15
131575	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:40	09/12 21:43
131574	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:29	09/12 22:40
131573	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:28	09/12 22:38
131572	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:28	09/12 21:42
131571	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:28	09/12 22:22
131570	done	metatranscriptome_assembly_and_alignment	metag	svc-ga	dori	09/12 17:27	09/12 23:23

Search & filter by site, status, tag

Track run status / logs

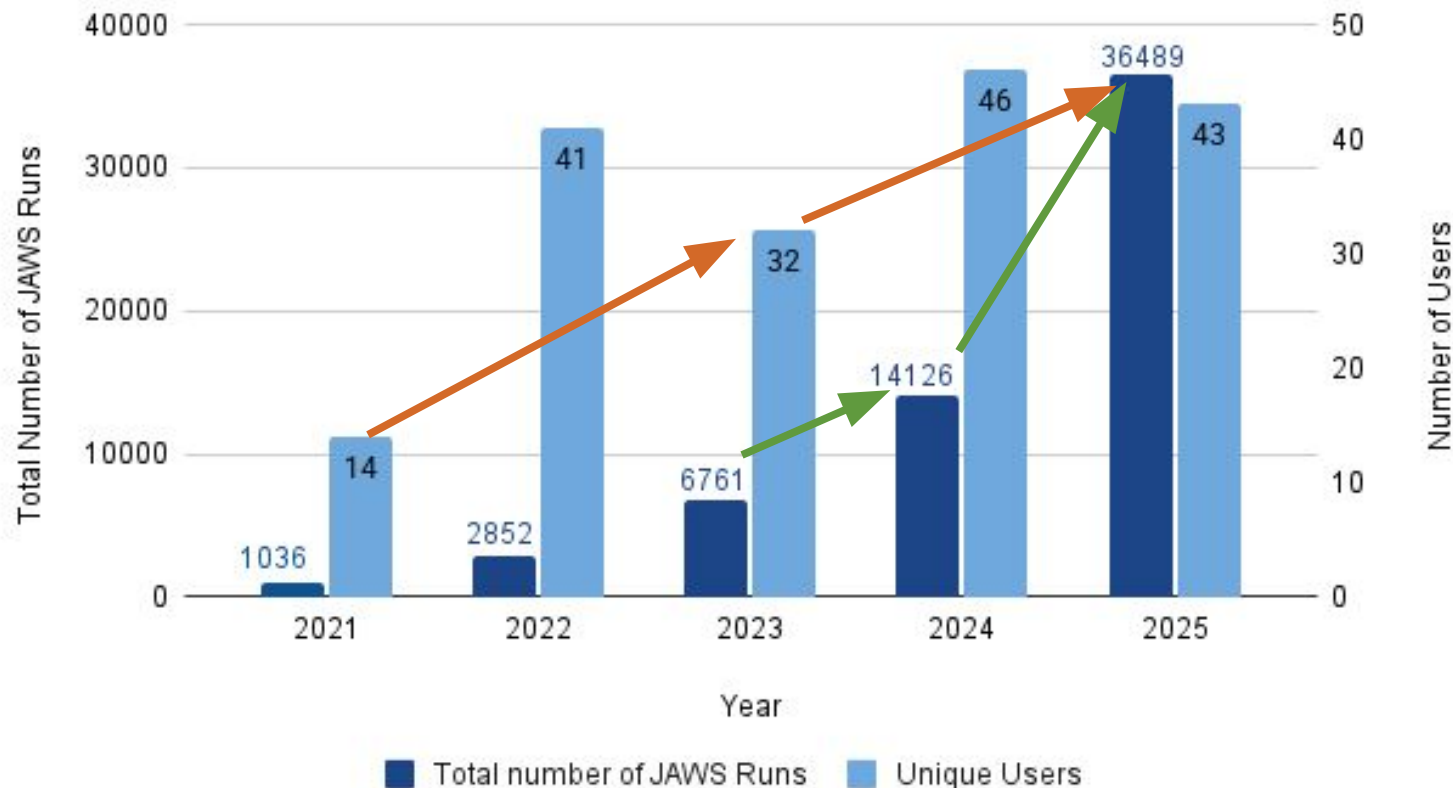
Expose performance metrics such as cpu/memory usage

The dashboard makes workflow status visible at scale

JAWS Usage Is Doubling Every Year!

Training sparked adoption. Automation now drives scale.

Total number of JAWS Runs and Users



Early Growth (2021–2023):

→ Internal training workshops

Scaling Up (2023–2025):

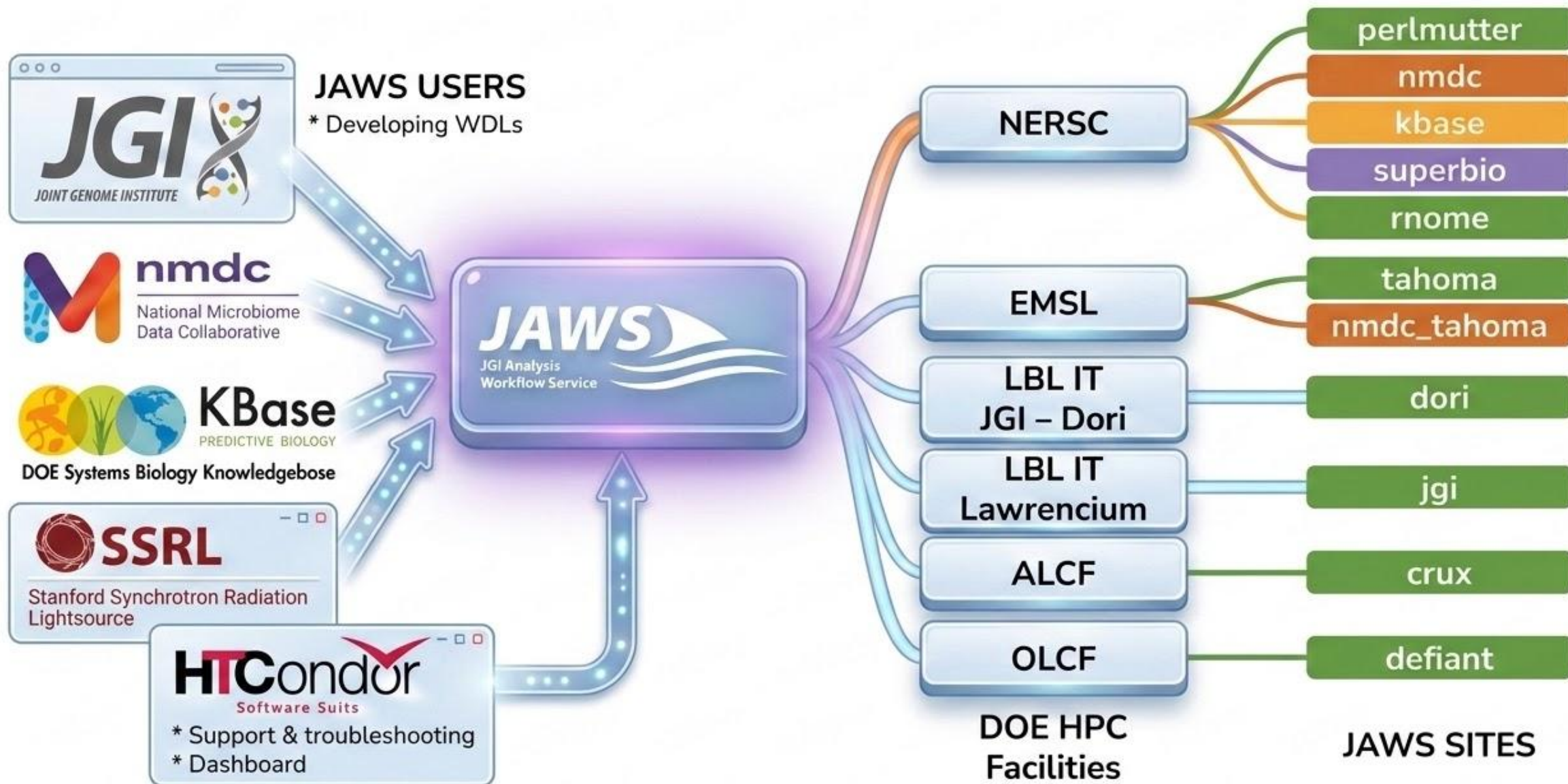
→ Usage more than doubled

Today:

→ The larger share of runs now comes from automation accounts.

“The JAWS team is a model of what good teams can be! Super responsive, transparent, always seeking and acting on feedback.”
(JAWS User Survey 2025)

Powered by Partnership. Ready for Collaboration

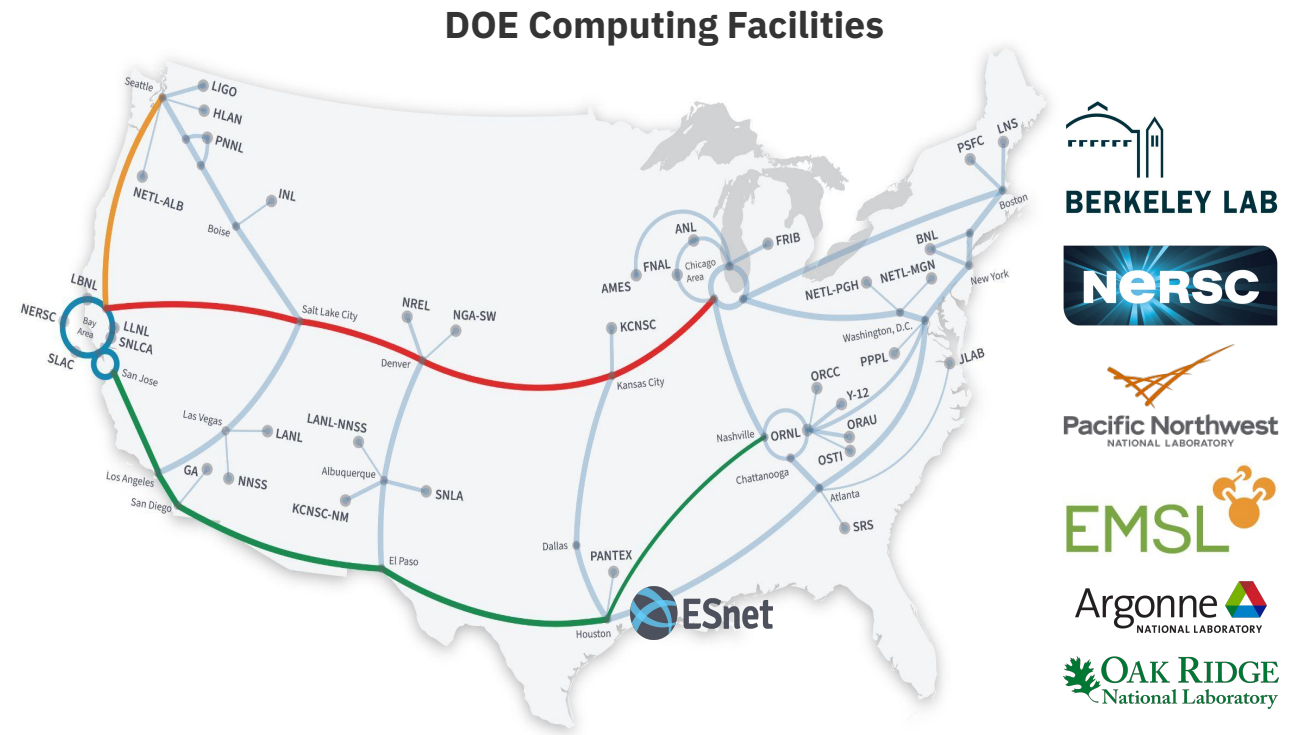


AI for Bio: Next-Gen Data Infrastructure

JGI is promptly moving forward to become an AI-centric facility that's Data Lakehouse enabled!

- Expanding sites
- Service resiliency
- Monitor and telemetrics

“Hey, JAWS!”
Powered by jaws.ai



HTCondor Team

Miron Livny, Todd Tannenbaum, Todd Miller

JAWS Team

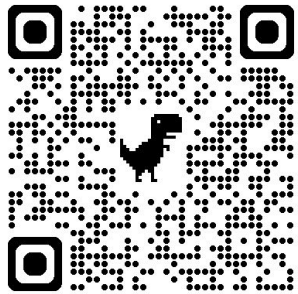
Seung-Jin Sul, Mario Melara, Ramani Kothadia, Ludovico Bianchi, Joshua Boverhof, Daniela Cassol, Mike Sneddon, Set Sarrafan, Kjersten Fagnan

NERSC

Nick Tyler

Thank you!

JAWS Docs



JAWS Code



Acknowledgements

The work conducted by the U.S. Department of Energy Joint Genome Institute (<https://jgi.doe.gov/>), a DOE Office of Science User Facility, is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-C02-05CH11231.



- **Cromwell**

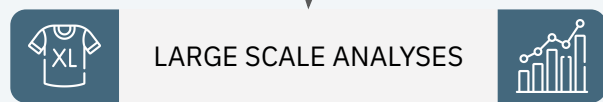
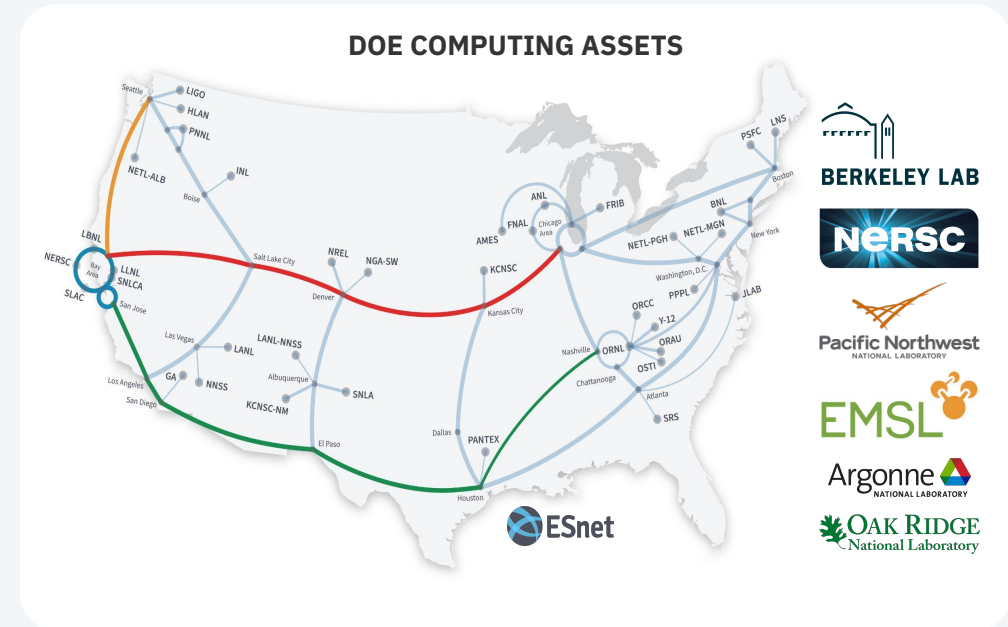
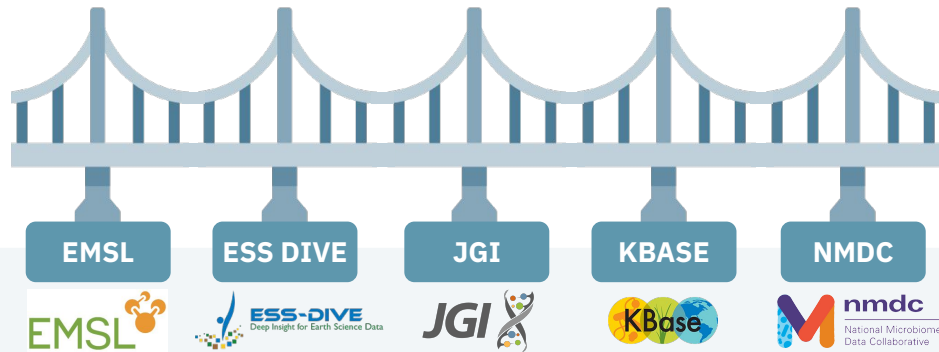
- Open-source workflow management software by the Broad Institute
- Executes WDL workflows with various compute backends
- Creates HTCondor job submit files
- Submits/manages jobs using HTCondor

	Cromwell backend interface	HTCondor command
Task submission	“submit”	condor_submit
Task cancellation	“kill-docker”	condor_rm
Status check	“check-alive”	condor_q

AI for Bio: Next-Gen Data Infrastructure

JGI is promptly moving forward to become an AI-centric facility that's Data Lakehouse enabled!

BRIDGE - **B**iological and **e**nvi**R**onmental **I**nfrastructure for **D**ata management and **E**xploration



LARGE SCALE ANALYSES

“Hey, JAWS!”
Powered by jaws.ai



- BERKELEY LAB
- NERSC
- Pacific Northwest NATIONAL LABORATORY
- EMSL
- Argonne NATIONAL LABORATORY
- OAK RIDGE National Laboratory

The **BRIDGE** initiative is a biological data infrastructure program developed by the U.S. Department of Energy (DOE) Joint Genome Institute (JGI).

Its name stands for Biological and enviRonmental Infrastructure for Data manaGement and Exploration.

The initiative aims to create an AI-ready data architecture and unified ecosystem to make genomic and environmental research data highly accessible and interoperable.

Key Features of JGI's BRIDGE Initiative

- Data Lakehouse: Integrates raw data repositories and data warehouses for seamless exploration.
- JAMO Integration: Uses JGI's Archive and Metadata Organizer to securely store and transfer data files.
- AI Readiness: Structures vast biological datasets to fuel modern Artificial Intelligence models and tools.
- Cross-Facility Access: Connects with other DOE hubs—such as the Environmental Molecular Sciences Laboratory (EMSL) and KBase—to allow predictive ecosystem biology.

The Department of Energy (DOE) Joint Genome Institute (JGI) is building a **data lakehouse** to transition into an AI-centric user facility.

This architecture combines the flexibility of a data lake (for raw, unstructured data) with the structure of a data warehouse (for indexed, highly-queryable data).

Key Features of the JGI Data Lakehouse

- **AI-Ready Data:** Standardizes and organizes massive genomic datasets to train foundational AI models across all kingdoms of life.
- **Unified Ecosystem:** Integrates data from the JGI, KBase, Environmental Molecular Science Laboratory (EMSL), and National Microbiome Data Collaborative (NMDC).
- **Better Accessibility:** Replaces legacy, disjointed search systems to provide researchers with a single, highly responsive entry point for over 13 petabytes of data.
- **Community Driven:** Actively collects user feedback to ensure the new architecture addresses past data frustration and fills current use cases.

Sites and Compute Pools

Site	Type	#Nodes	Mem (GB)*	Minutes	#Threads	#GPUs
Perlmutter (NERSC)	Large	3072	492	2865	256	0
	GPU (4x NVIDIA A100 (40GB))	1536	256	2865	128	4
	GPU (4x NVIDIA A100 (80GB))	256	256	2865	128	4
JGI (Lab-IT)	Large	8	492	4305	32	0
Dori (Lab-IT)	Large	100	492	4305	64	0
	Xlarge	16	1980	20160	128	0
Tahoma (EMSL)	Medium	184	364	2865	36	0
	Xlarge	24	1480	2865	36	0
	GPU (2x NVIDIA Tesla V100 32GB)	24	1480	2865	36	2
Crux (ALCF)	Medium	256	256	1425	256	0
Defiant (OLCF)	Medium	20	492	4305	128	0

About us

- High-throughput AI-centric genomic science user facility located at Lawrence Berkeley National Laboratory
- Provides the genomic capabilities, data, and expertise that supports the global research community in studying complex biological and environmental systems

Advanced genomic capabilities

Large-scale sequencing and synthesis

AI-ready scientific data

Curated data products for reuse

Professional expertise

Support from domain and workflow experts



This mural at our site in Berkeley, California, shows some of the species that JGI users have studied.

