Our mission is to…

increase the efficiency and effectiveness of researchers engaged in data-driven science and scholarship through *sustainable* software.
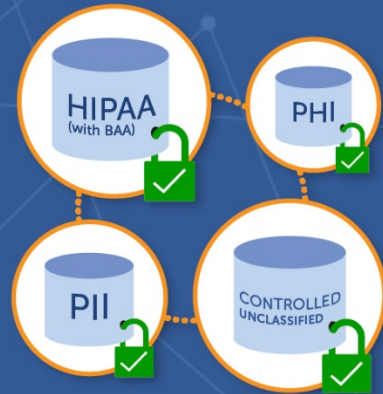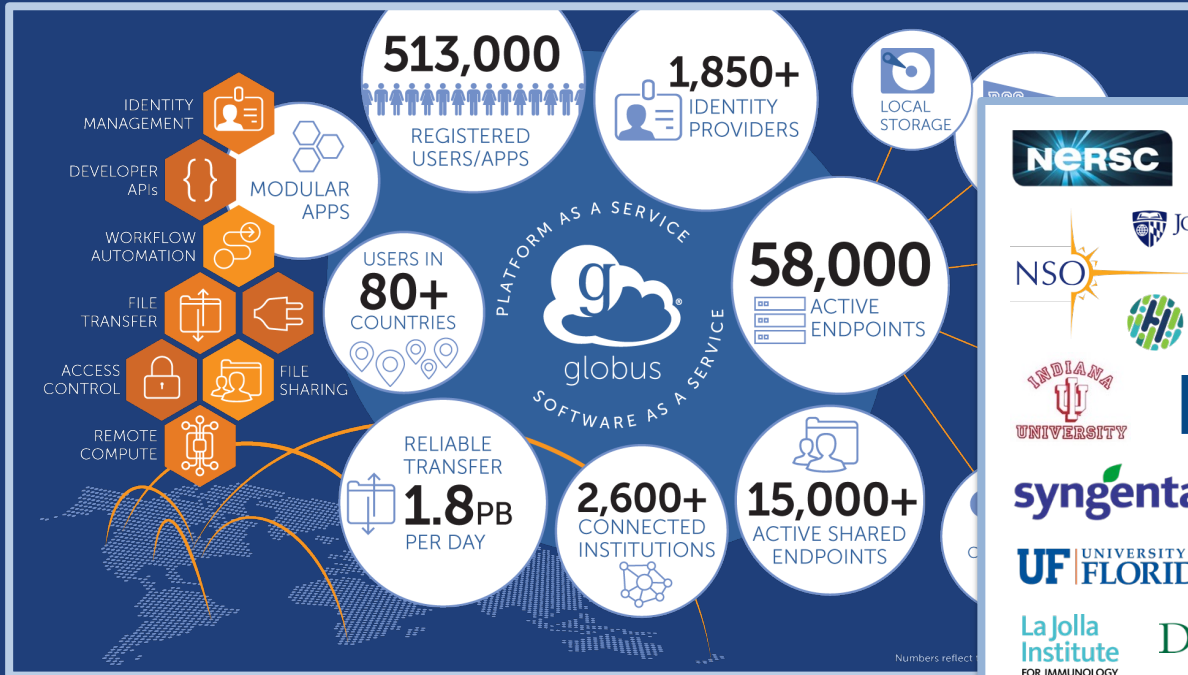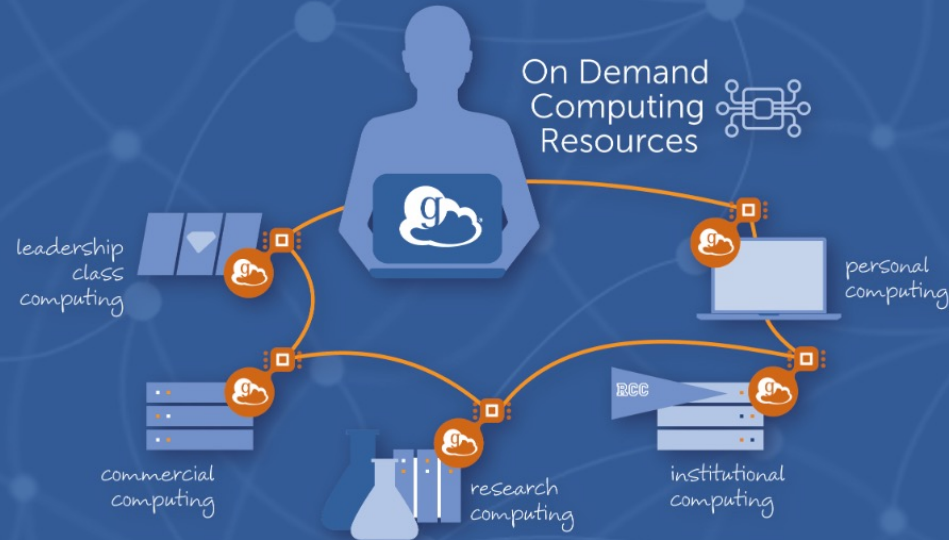
Development is partly funded by...

# Our story back then…

- **Oct. 1998: Globus Toolkit v1 released—The Grid**

- **Nov. 2010: Globus Online released—The Cloud**

- Nov. 2013: Freemium model launched

- Jan. 2019 - 100th subscriber signed, >50% sustainable

- ??? – Globus becomes fully self-sustaining

| | | | |
|---|---|---|---|
| **1,042** most shared endpoints at a single institution | **450 PB** moved | **75 billion** files processed | **1,700** active server endpoints |
| | **90** subscribers | **100,000** users | **3 months** longest running transfer |
| **18,000** active personal endpoints | **500** identity providers | **1 PB** largest single transfer to date | **8,000** active shared endpoints / **99.9%** availability |

# What a difference 5 years make!



- **Network effects driving growth**

- **Asymptotically approaching sustainability :-)**

# What's driving the need to reimagine research IT?

# Circa 1980…

# Circa 2020: The Instruments are Coming!

New York Genome Center

Advanced Light Source

NISC — NIH Intramural Sequencing Center

Purdue Science — Purdue Cryo-EM Facility

High Resolution Cryo-EM — National Cryo-Electron Microscopy Facility

# More (and larger) collaborations



**Active Guest Collections**
(a Globus mechanism by which researchers can share data with collaborators)

Legend:
- Shared endpoints
- Guest collections (GCS)
- Guest collections (GCP)
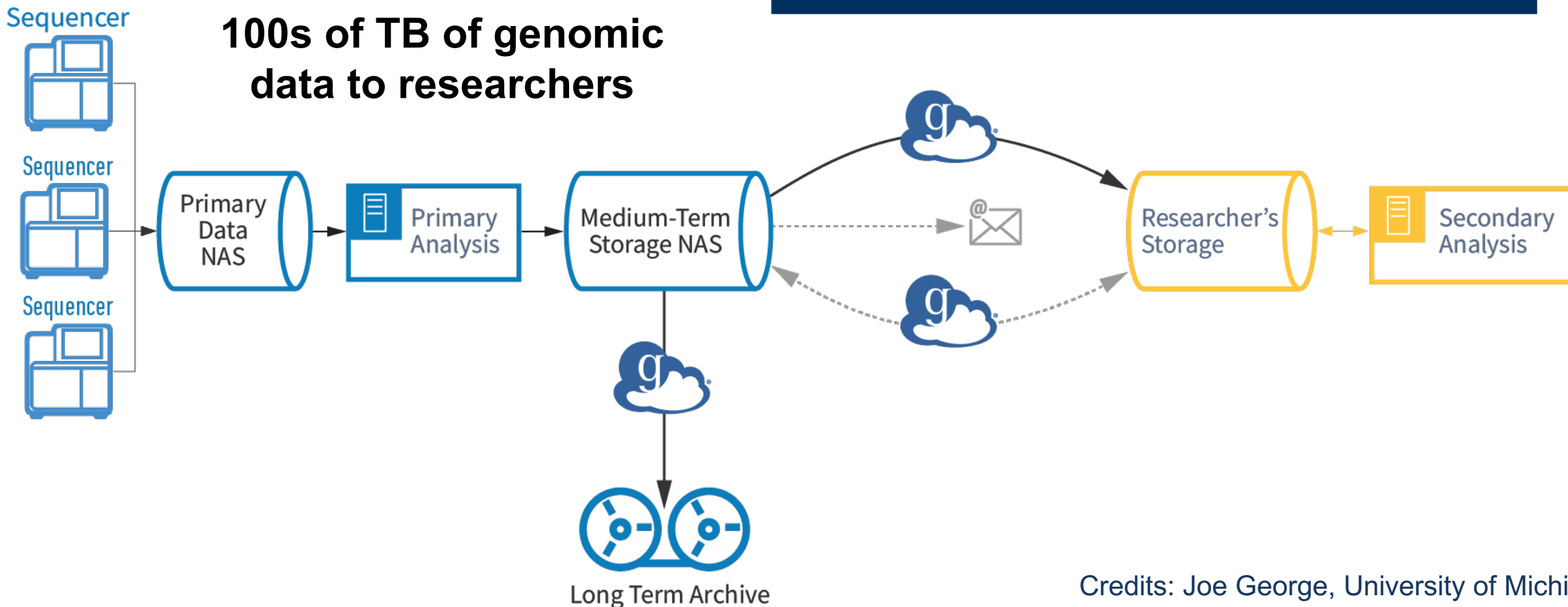
**Dramatic changes in sensor data rates and broader collaborations demand large-scale automation, from capture through publication**

**Use Globus to deliver 100s of TB of genomic data to researchers**

Credits: Joe George, University of Michigan

# Automating NGS

**Globus Flows**

High performance storage

Archive storage

**PromethION ~10TB/run**

**MinION ~50GB/run**

**Transfer**
Transfer raw files

**Transfer**
Backup raw files

**Compute**
Base calling

Analysis pipelines

**Share**
Grant lab access

**Transfer**
Results to lab servers

**Compute**
Variant calling

# Cryo-EM will eat (all) your resources



Cryo-EM resolution, angstroms

- Year's average reported resolution
- Year's highest reported resolution

Their work is advancing new treatments against a whole range of diseases, such as Covid-19, chlamydia, and some cancers.

Ramos, S. GlobusWorld 2023
https://www.youtube.com/watch?v=q-o6VKIEPj4

# Serial crystallography at the Argonne APS

- Serially image chips with thousands of embedded crystals

- Quality control first 1,000 to report failures

- Analyze batches of images as they are collected

- Report statistics and images during experiment

- Return crystal structure to scientist



Darren Sherrell, Gyorgy Babnigg, Andrzej Joachimiak

# How Globus enables SSX automation at scale

**Data capture**

**Globus Flows**

**Compute** — Launch QA job

**Carbon!** — Check threshold

**Transfer** — Transfer raw files

**Compute** — Analyze images

**Image processing**

**Data publication**

**Search** — Ingest to index

**Share** — Set access controls

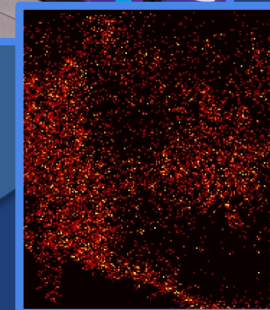**Transfer** — Move results to repo

**Compute** — Gather metadata

**Compute** — Visualize
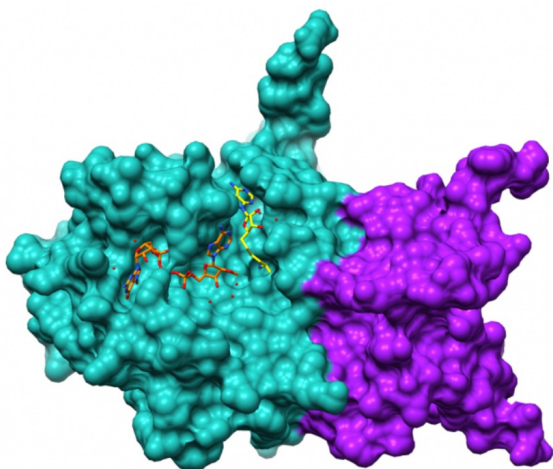
# Outcomes and impact



"These data services have taken the time to solve a structure from weeks to days and now to hours"

*Darren Sherrell, SBC beamline scientist APS Sector 19*



**SCIENCE**

## Argonne researchers use Theta for real-time analysis of COVID-19 proteins

**AUTHOR** NILS HEINONEN
**PUBLISHED** 07/28/2020
**DOMAIN** BIOLOGICAL SCIENCES
**SYSTEMS** THETA

- **Automation pipeline** collects data, analyzes and visualizes the data, solves protein structure and loads results into a searchable portal for real-time feedback
- **Achieved over 10-100x speed up** in time to solution of protein structures at APS

# Where does Globus fit in this picture?

**Managed transfer & sync**

fast
secure transfer
reliable

A → B

**Publication & discovery**

DOI 456

**Platform-as-a-Service**

your app

**Collaborative data sharing**

collaborators

colleagues

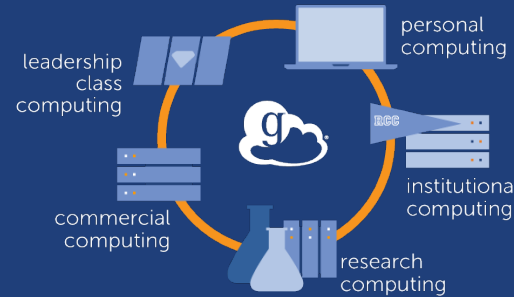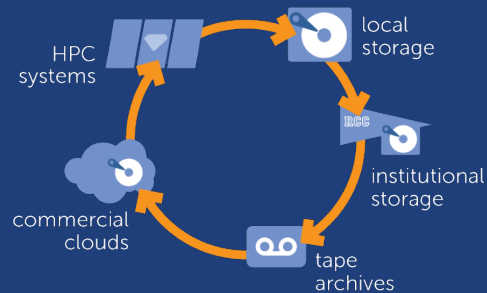**Managed remote execution**

leadership class computing

personal computing

commercial computing

institutional computing

research computing

**Unified data access**

HPC systems

local storage

institutional storage

tape archives

commercial clouds

**Reliable automation**

Auth    Transfer    Search

GET CREDENTIALS    TRANSFER DATA    INGEST

**Software-as-a-Service**

File Manager    Panels

Collection    Search    Search
Path

Start    Transfer & Timer Options    Start

Search for a collection to begin

Get started by taking a short tour.

Search for a collection to begin

Get started by taking a short tour.

# Making data FAIR by default

**Managed transfer & sync**

fast
secure transfer
reliable

A → B

**Collaborative data sharing**

colleagues
collaborators

**Unified data access**

HPC systems
local storage
institutional storage
commercial clouds
tape archives

## Publication & discovery

DOI 456
DOI 1...

**Reliable automation**

GET CREDENTIALS — TRANSFER DATA — INGEST

**Software-as-a-Service**

# Globus Search: Data description and discovery

- **Metadata store with fine-grained visibility controls**

- **Schema agnostic**

- **DOI minting via DataCite**

- **Simple search using URL queries**

- **Complex search using filters and facets**

**Search Index**

globus

**docs.globus.org/api/search**

# Making data more FAIR by default



Collaborative data sharing

colleagues

collaborators

Unified data access          Reliable automation

your app

Platform-as-a-Service

File Manager

Software-as-a-Service

# Secure data sharing ...from any storage

Select files to share, select user or group, and set access permissions

**1**

**Globally accessible multi-tenant service**

- **Fine-grained access control**
- **Storage system "overlay"**
- **Share with identity/email/group**
- **No data staging required**

Globus controls access to shared files on existing storage

**globus connect**

Collaborator logs into Globus and accesses shared files; no local account required; download via Globus

**2**

**On-prem or public cloud storage**

Google Cloud

**globus connect**

**Laptop, server, compute facility**

# Distinct access policies may be applied to

# Data *and* Metadata

# Support for managing protected data

**Security controls**
→ **NIST 800-53**
→ **NIST 800-171**

**Restricted data handling**
→ **PHI, PII, CUI**
→ **Compliant data sharing**

**BAA w/Uchicago**
→ **UChicago BAA with Amazon**

# Making data more FAIR by default



HPC systems

local storage

institutional storage

tape archives

commercial clouds

**Unified data access**

your app

**Platform-as-a-Service**

File Manager
Panels
Collection    Search                Search
Path

Start ▷        Transfer & Timer Options ⌄        Start

Search for a collection to begin

Get started by taking a short tour.

Search for a collection to begin

Get started by taking a short tour.

**Software-as-a-Service**

# Globus supports diverse storage systems

# Enabling reuse via remote computation



**Managed transfer &**

**Collaborative data sh**

**Unified data access**

**Reliable automation**

**Managed remote execution**

leadership class computing

personal computing

commercial computing

research computing

institutional computing

**Software-as-a-Service**

fast secure transfer reliable

colleagues

collabor

HPC systems

local storag

commercial clouds

instit stora

tape archives

GET CREDENTIALS

TRANSFER DATA

INGEST

# Globus Compute: Function-as-a-Service …on any system

- **Fire and forget function execution**
- **Uniform interface to diverse compute resources**
- **…from a laptop to a supercomputer**

**Globally accessible multi-tenant service**

User submits a function to be run on compute endpoints

**1**

**3**

Results returned to the user

Globus manages the function execution on any endpoint

**2**

**2**

**Compute Facility**

**Laptop, server, compute facility**

# Scaling FAIRness in the research enterprise

**Managed transfer & sync**

**Collaborative data sharing**

**Unified data access**

Auth — GET CREDENTIALS

Transfer — TRANSFER DATA

Search — INGEST

**Reliable automation**

# Globus Flows: Reliable, secure task orchestration

- A platform for defining, executing, and sharing distributed research automation flows

- Flows comprise **Actions**

- **Action Providers**: Called by Flows to perform tasks

- **Triggers**\*: Start flows based on events

\* Coming soon

# Extending FAIRness beyond the institution

**Managed transfer & sync**

fast
secure transfer
reliable

A B

**Collaborative data sharing**

colleagues
collaborators

**Publicatio**

DOI 1

leadership
class
computing

commercial
computing

**Managed re**

**Unified data access**

HPC
systems

local
storage

BCC

institutional
storage

commercial
clouds

tape
archives

**Reliable automation**

Auth

GET
CREDENTIALS

TRANSFER
DATA

INGEST

your
app

**Platform-as-a-Service**

**Software-as-a-Service**

# Portal frameworks to enable reusability

**acdc.alcf.anl.gov**

# A growing ecosystem



Registered applications

# Data FAIRness through data mobility



fast
secure transfer
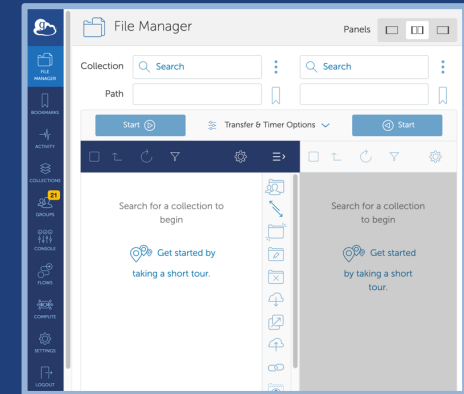reliable

**Managed transfer & sync**

Unified data access

Reliable automation

**Platform-as-a-Service**

**Software-as-a-Service**

# Readily accessible by researchers everywhere

- **Federated login for 1,800+ institutions**

- **Access via ORCID ID for millions of PIs**

- **Open to the world via Google, GitHub, …**

# Large-scale automation is increasingly enabling intelligent analysis and experiment steering

# What keeps you up at night?

# "Smart" instruments are emerging

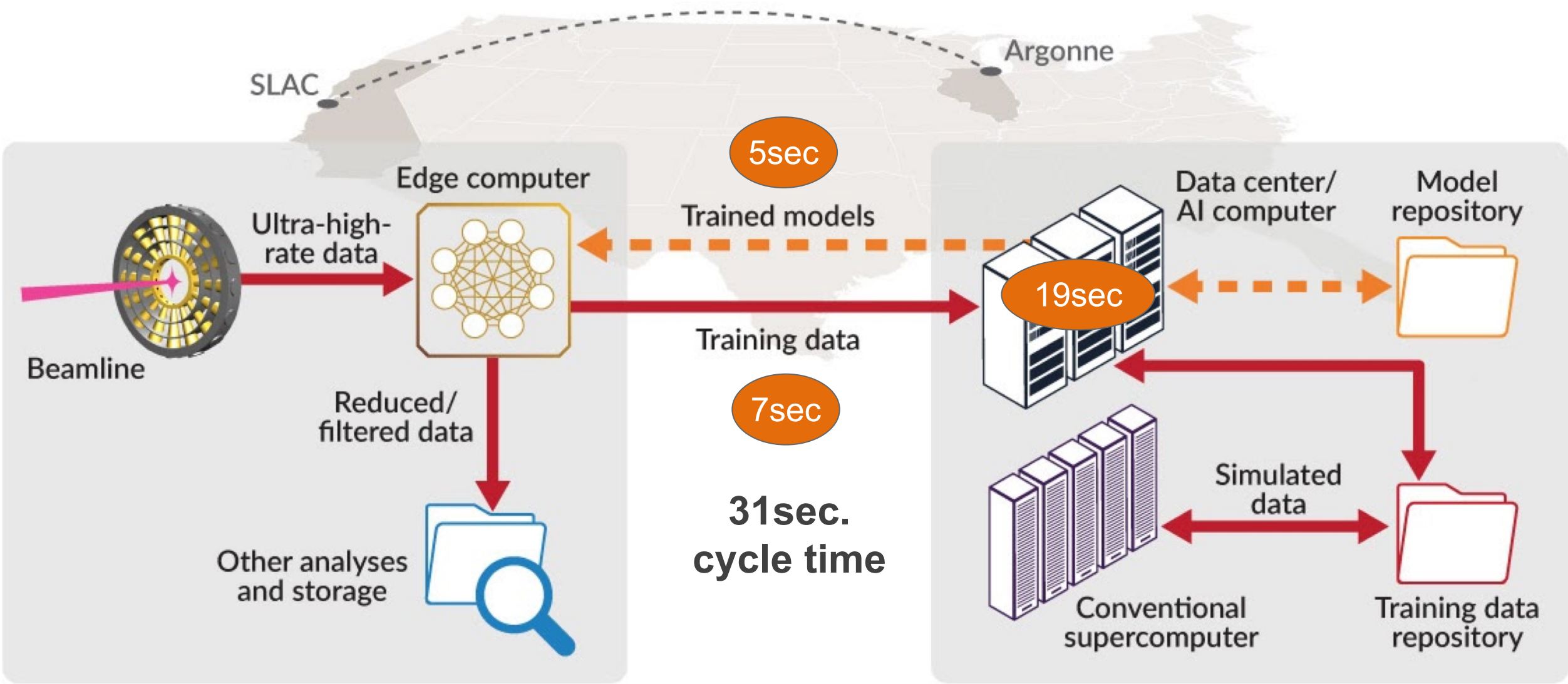# Current research will increase access to advanced instrument automation capabilities
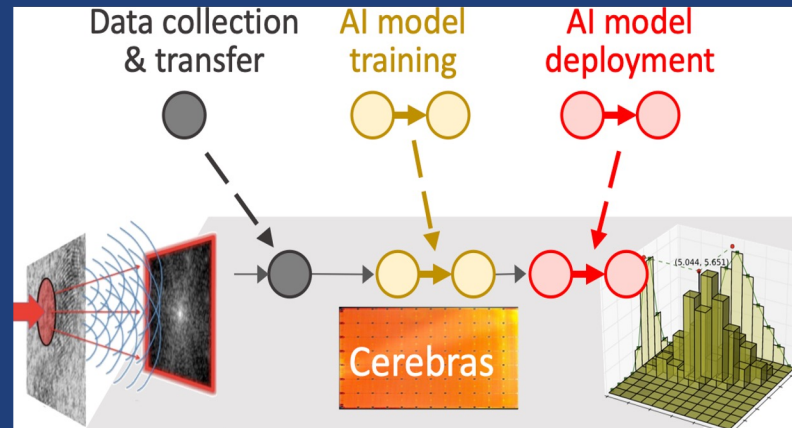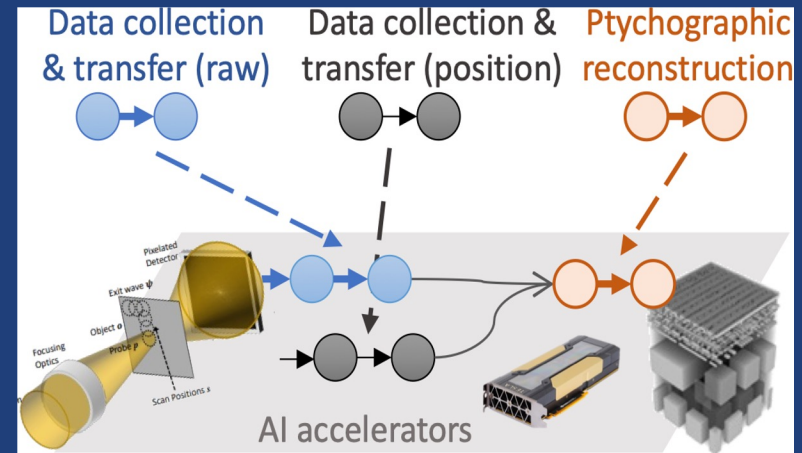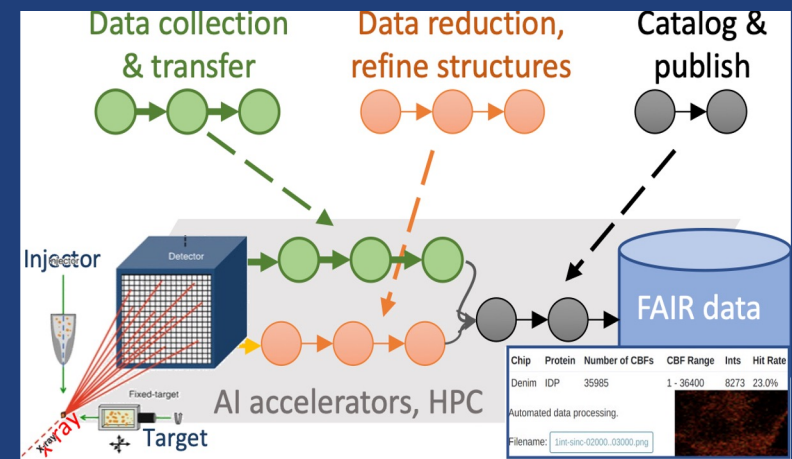
Braid (2020-2023) developed methods and tools for defining, running, and scheduling flows that link instruments, computers, data repositories, people

- **Method and tools:** Globus Flows for representing flows, Braid for scheduling sets of flows, HPC DL training algorithms, data pub methods, etc.

- **Deployment and application:** APS, LCLS, microscopes, CryoEM, etc.; ALCF and other compute

- **Online experimental data analysis, experiment steering, data ingest pipelines, climate data analysis, etc.**

- **Operational experience:** 10,000s of flow runs, 100s TB data processed, 10,000s node hours

- **Science impact:** Determination of Covid protein structures, online high energy diffraction microscopy analysis, etc.

https://doi.org/10.1016/j.patter.2022.100606

Blaiszik, Chard, Chard, Foster, Huerta, Nicolae, Vescovi, Wozniak

# Diaspora (2023-2028) will develop a hierarchical event fabric and resilience solutions that will be applied to distributed applications involving long-term campaigns, time-sensitive analysis, and distributed data integration

- **Argonne, ORNL, SLAC, Texas Tech University**

- **Event fabric supports eventing, monitoring, and resilience across institutions/facilities**

- **Resilient distributed data structures can be deployed across diverse resources**

https://diaspora-project.github.io

Foster, Rao, Thayer, Corsi