

# The 200 Gbps Challenge: Imagining HL-LHC analysis facilities

Sam Albin <sup>1</sup>, Garhan Attebury <sup>1</sup>, Ken Bloom <sup>1</sup>, Brian Paul Bockelman <sup>2</sup>, Lincoln Bryant <sup>3</sup>,  
Kyungeon Choi <sup>4</sup>, Kyle Cranmer <sup>5</sup>, Peter Elmer <sup>6</sup>, Matthew Feickert <sup>5</sup>, Rob Gardner <sup>3</sup>, Lindsey Gray <sup>7</sup>,  
**Alexander Held** <sup>5</sup>, Fengping Hu <sup>3</sup>, David Lange <sup>6</sup>, Carl Lundstedt <sup>1</sup>, Peter Onyisi <sup>4</sup>, Jim Pivarski <sup>6</sup>,  
Oksana Shadura <sup>1</sup>, Nick Smith <sup>7</sup>, John Thiltges <sup>1</sup>, Ben Tovar <sup>8</sup>, Ilija Vukotic <sup>3</sup>, Gordon Watts <sup>9</sup>,  
Derek Weitzel <sup>1</sup>, Andrew Wightman <sup>1</sup>

<sup>1</sup> University of Nebraska–Lincoln, <sup>2</sup> Morgridge Institute for Research, <sup>3</sup> University of Chicago,

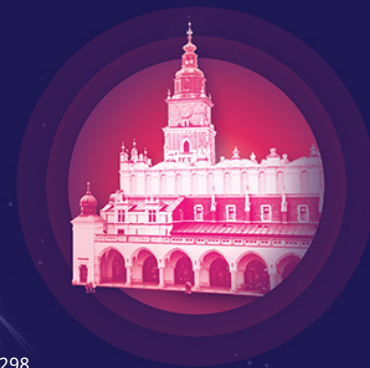
<sup>4</sup> University of Texas at Austin, <sup>5</sup> University of Wisconsin–Madison, <sup>6</sup> Princeton University,

<sup>7</sup> Fermilab, <sup>8</sup> University of Notre Dame, <sup>9</sup> University of Washington

CHEP 2024

<https://indico.cern.ch/event/1338689/>

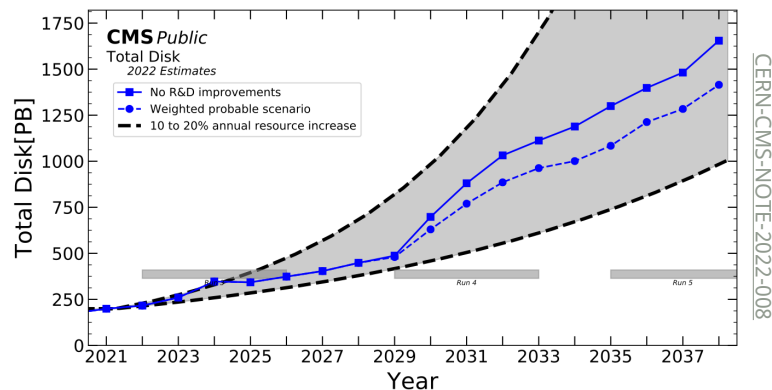
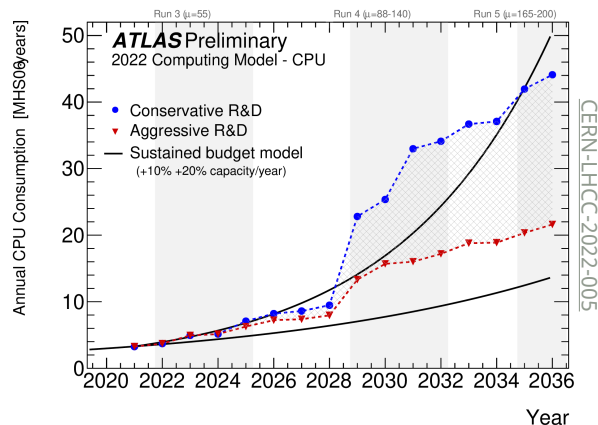
Oct 21, 2024



# CHEP 2004 and 2024

- Until the mid-1980s HEP's "computing problem" was often thought to be about **obtaining enough processor power**
- Then we worried about **storage capacity**
- The real problem has always been, in my opinion, **getting people to collaborate on a solution**

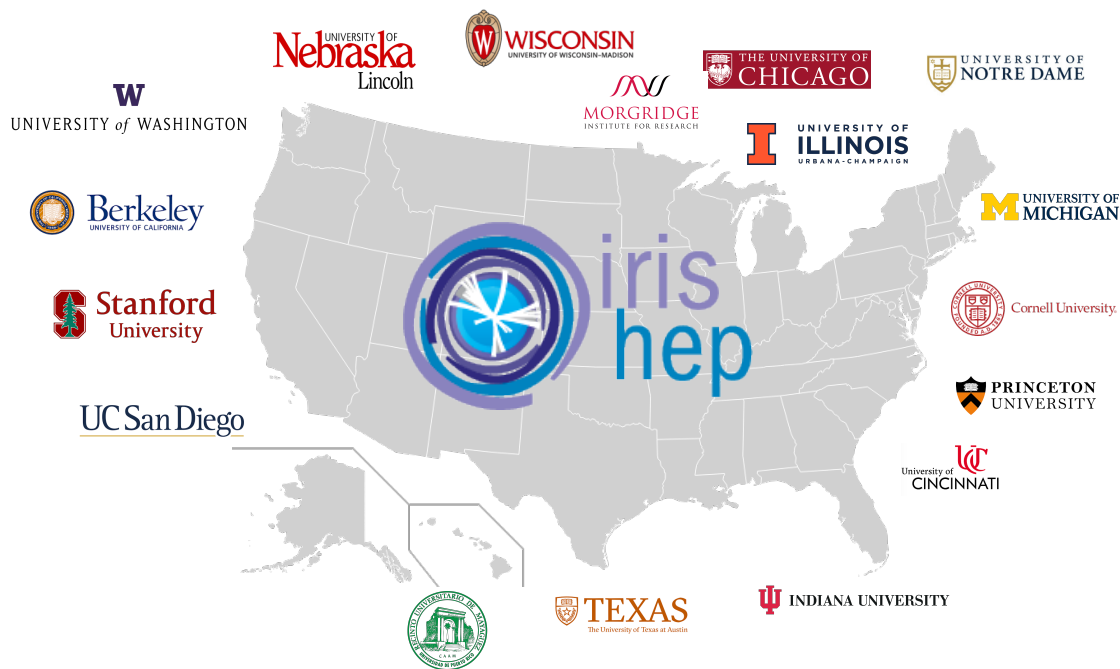
[David Williams: "50 years of Computing at CERN", CHEP 2004]



# IRIS-HEP and a HL-LHC vision

# The IRIS-HEP software institute

- **IRIS-HEP**: “Institute for Research and Innovation in Software for High Energy Physics”
  - a software institute funded by the US National Science Foundation, running 2018–2028
  - working in close collaboration with LHC experiments and facilities



[2024 IRIS-HEP retreat]



# R&D for the HL-LHC

- IRIS-HEP is working on **computing and software R&D for the HL-LHC**
  - a **software upgrade** accompanying detector hardware upgrades
  - focusing on a subset of **activity areas** today, connected through “**challenge**” format

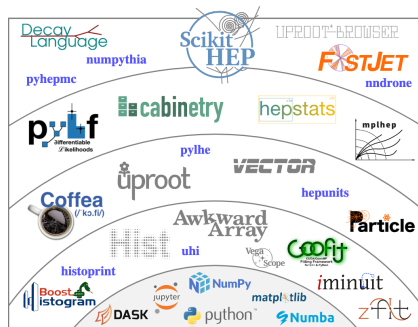
## DOMA

*data organisation and management*



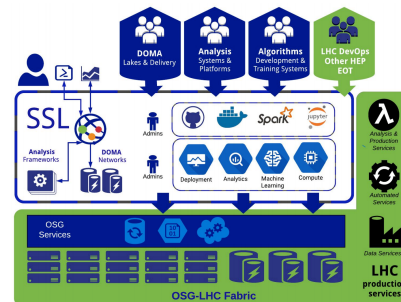
## AS

*analysis systems and tools*



## SSL and OSG-LHC

*deployment techniques and resources*



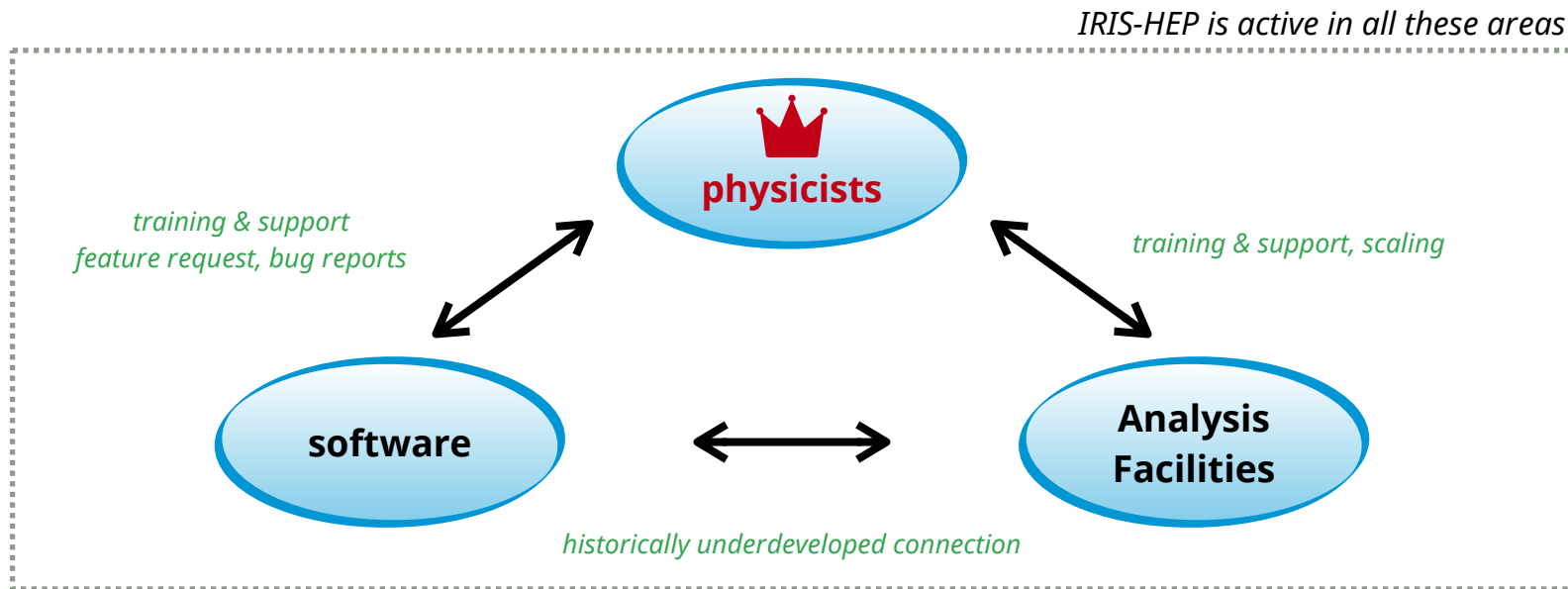
## SSC

*training*



# Empowering physicists, today and tomorrow

- This work is driven by the desire to **minimize time-to-insight** and **maximize the HL-LHC physics reach**
  - let **physicists** spend **more time doing physics, less time debugging, bookkeeping, waiting, ...**
  - **tighten feedback & support cycles** by connecting communities together
- **Physics is the end goal**: strive to find ways to overcome computing challenges



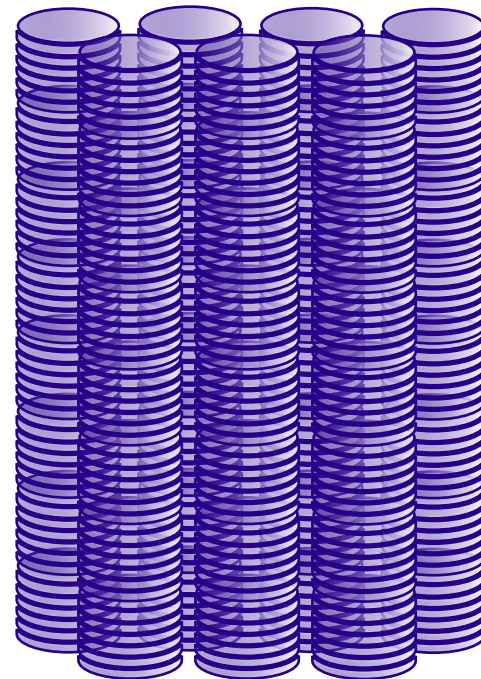
# Our end-user analysis vision

- **Analyze  $O(1000)$  TB of data within a few hours**
- Interactive analysis turnaround: “coffee break” timescale
- Fully integrated Analysis Facilities (AFs)
- UX to empower big & small teams
- Easy access to state-of-the-art ML + techniques
- Reproducibility, preservation, reuse



**today:**

create  $\mathcal{O}(1 - 10)$  TB ntuples  
on the grid  
in  $\mathcal{O}(\text{days} - \text{weeks})$ ,  
analyze on Tier-3 in  $\mathcal{O}(h - \text{days})$

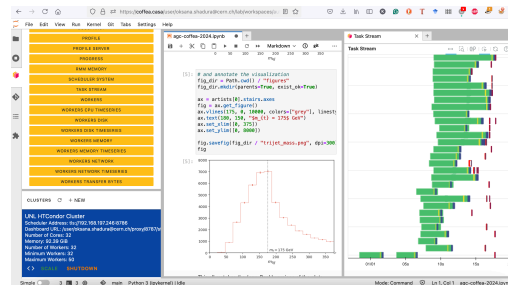


**HL-LHC:**

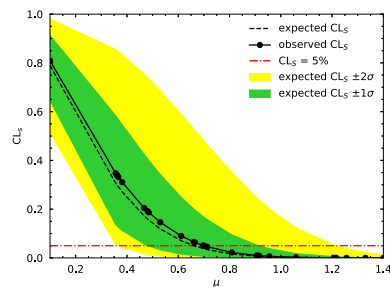
analyze  $\mathcal{O}(1000)$  TB of data  
straight out of central  
PHYSLITE / NanoAOD files in  $\mathcal{O}(h)$

# Our end-user analysis vision

- Analyze  $O(1000)$  TB of data within a few hours
- **Interactive analysis turnaround: “coffee break” timescale**
- Fully integrated Analysis Facilities (AFs)
- UX to empower big & small teams
- Easy access to state-of-the-art ML + techniques
- Reproducibility, preservation, reuse



*meaningful analysis iterations on timescale of a coffee break, interactive analysis design*

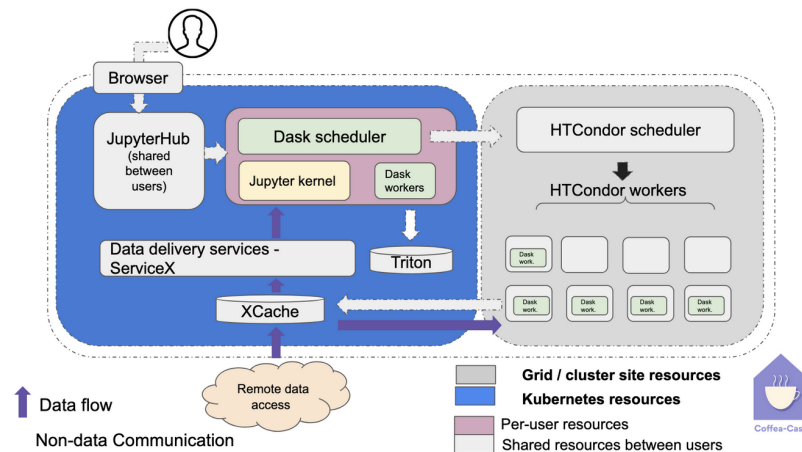


# Our end-user analysis vision

- Analyze  $O(1000)$  TB of data within a few hours
- Interactive analysis turnaround: “coffee break” timescale

- **Fully integrated Analysis Facilities (AFs)**

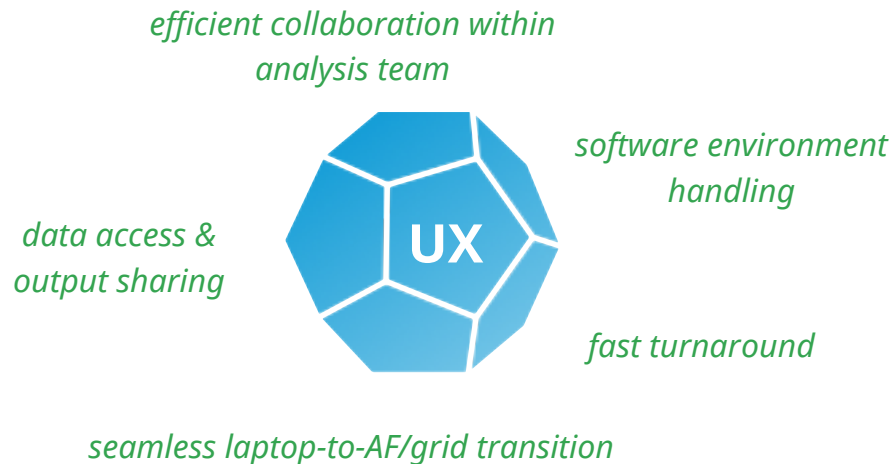
- UX to empower big & small teams
- Easy access to state-of-the-art ML + techniques
- Reproducibility, preservation, reuse



*required services available,  
convenient interfaces,  
access to powerful resources*

# Our end-user analysis vision

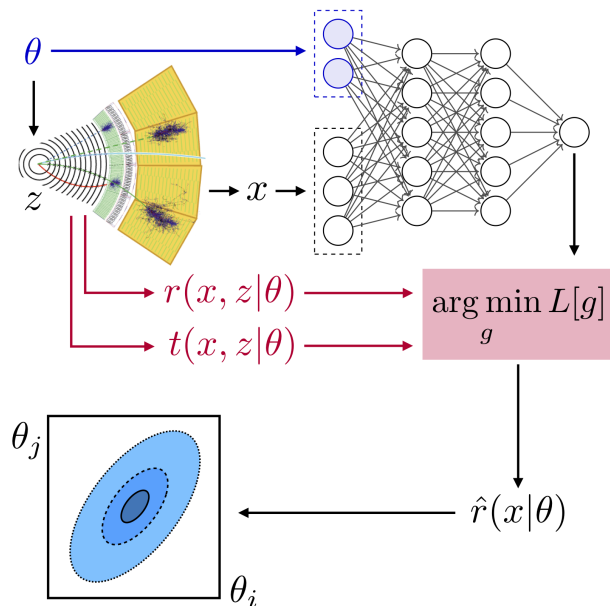
- Analyze  $O(1000)$  TB of data within a few hours
- Interactive analysis turnaround: “coffee break” timescale
- Fully integrated Analysis Facilities (AFs)
- **UX to empower big & small teams**
- Easy access to state-of-the-art ML + techniques
- Reproducibility, preservation, reuse



see also: [HSF AF White Paper](https://arxiv.org/abs/2404.02100)  
<https://arxiv.org/abs/2404.02100>

# Our end-user analysis vision

- Analyze  $O(1000)$  TB of data within a few hours
- Interactive analysis turnaround: “coffee break” timescale
- Fully integrated Analysis Facilities (AFs)
- UX to empower big & small teams
- **Easy access to state-of-the-art ML + techniques**
- Reproducibility, preservation, reuse



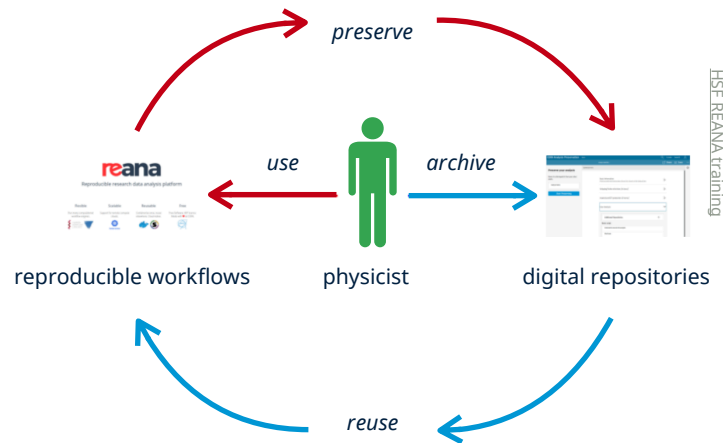
from MadMiner tutorial

*simulation-based inference techniques use different workflows from traditional histogram-based approaches*

# Our end-user analysis vision

- Analyze  $O(1000)$  TB of data within a few hours
- Interactive analysis turnaround: “coffee break” timescale
- Fully integrated Analysis Facilities (AFs)
- UX to empower big & small teams
- Easy access to state-of-the-art ML + techniques

• **Reproducibility, preservation, reuse**



*sustainable research  
maximizing long-term impact and legacy*

*see also:*

[Nature 533, 452–454 \(2016\)](#)

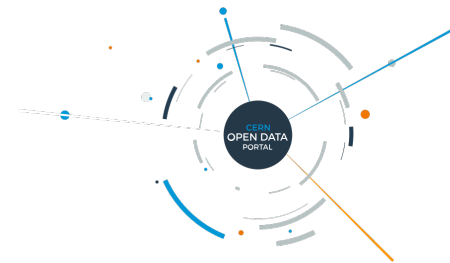
[arXiv:2203.10057 \[hep-ph\]](#)





# The Analysis Grand Challenge (AGC)

# A test case for HL-LHC analysis

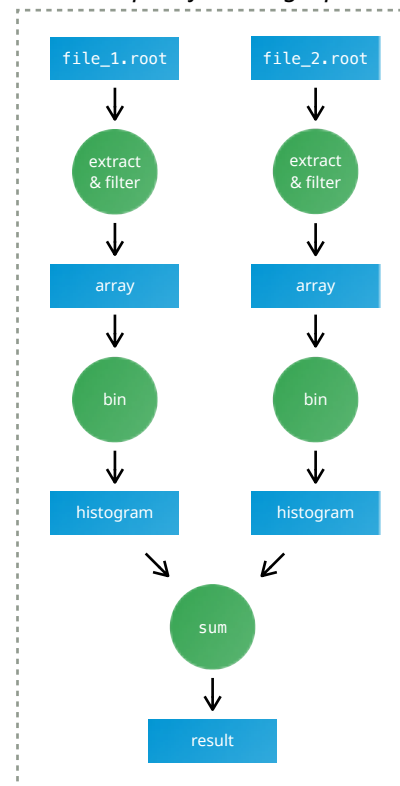
- The **Analysis Grand Challenge (AGC)** defines a **physics analysis task** with **Open Data** to test **HL-LHC workflows**
  - **columnar data extraction** from large datasets & data processing into **histograms**
  - statistical model construction and **statistical inference**, relevant **visualizations**
  - **ML training & inference**



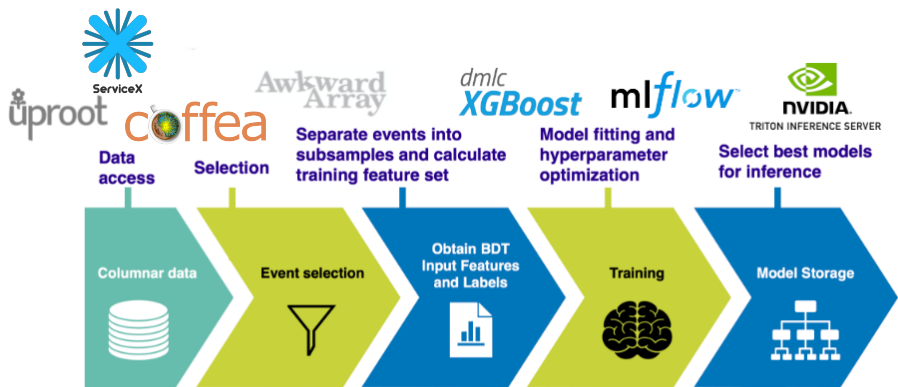
# Analysis with task graphs and dask

- We employ **task graphs** to **express & execute** data analysis operations
- This relies on  **dask**, a **Python library** providing
  - an interface to describe **manipulations of data via task graphs**
  - a **task scheduler** to execute task graphs
- **Deep integration of Dask** and existing **Python HEP toolset** with minimal API changes
  - arrays via  **dask-awkward**, histograms, **coffea** etc.
  - Dask emerging as **common feature in Analysis Facilities**

*example of a task graph*

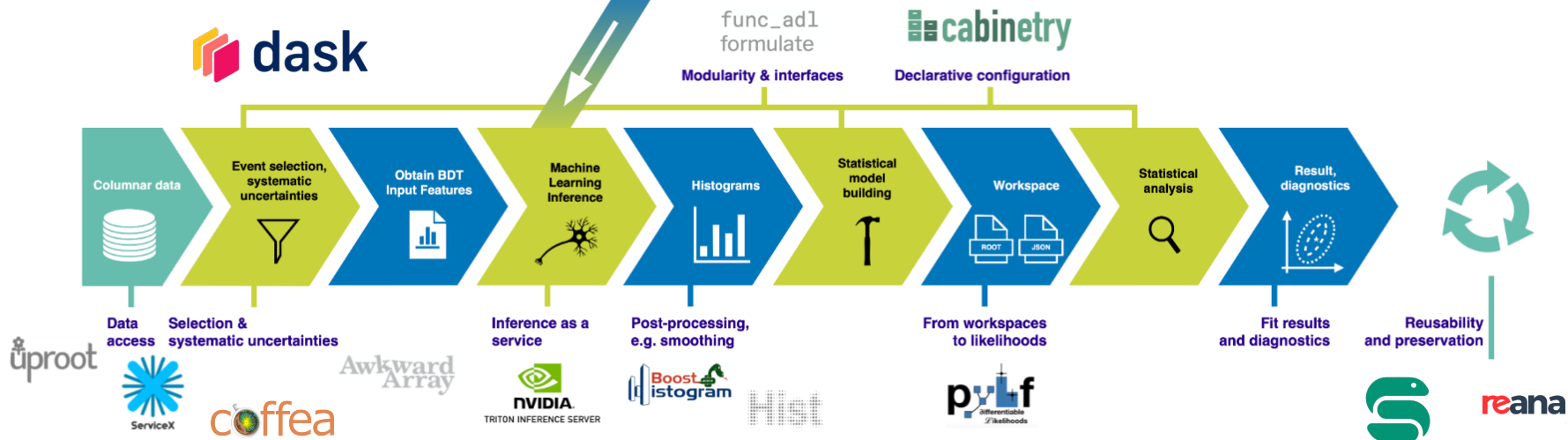


# The IRIS-HEP AGC reference implementation



The **IRIS-HEP reference implementation** employs the **Scikit-HEP/ PyHEP ecosystem** and serves as **ideal environment** to test our **latest R&D**.

find it all on [GitHub](#) and <https://agc.readthedocs.io/>



# A community project

- **Variety of AGC implementations** have been developed, more are welcome!

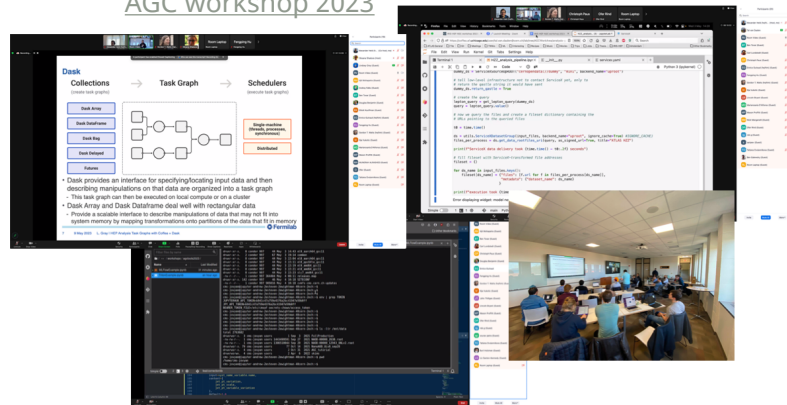


- **Regularly hosting related events**

- ▶ AGC workshops (#1, #2, #3, #4)
- ▶ “Demo day” event series

- **AGC keeps evolving:** wishlist for future additions

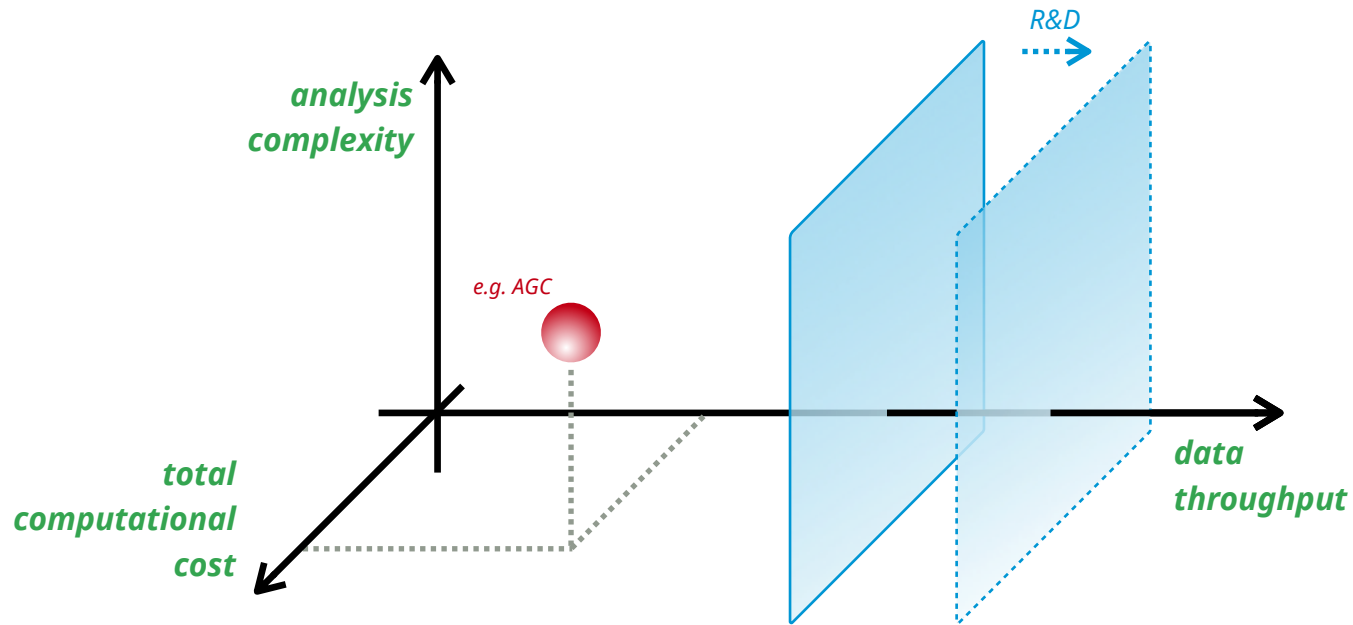
AGC workshop 2023



# The 200 Gbps Challenge

# Breaking down physics analyses

- Currently **limited agreement** on how “**HL-LHC physics analyses**” will look
  - **factorize into independent challenges**, push the boundaries in all directions



# Defining the 200 Gbps Challenge



Reaching these scales poses significant challenges. We set ourselves an ambitious goal.  
*... and had only 8 weeks to reach it.*

## CMS NanoAOD example

With **2 kB / event**, this means

- **90 B events**,
- **50 MHz event rate**,

or 1k cores with 50 kHz and 25 MB/s each.

## ATLAS PHYSLITE example

With **10 kB / event**, this means

- **18 B events**,
- **10 MHz event rate**,

or 2k cores with 5 kHz and 12.5 MB/s each.

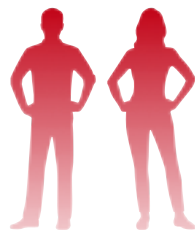


# Defining the 200 Gbps Challenge



- Targeting **“HL-LHC scale” analysis**, including **decompression** & **data in memory** as arrays
- **Two different setups, targeting realism**, all code on GitHub
  - **Nebraska**: analyze **Run-3 NanoAOD** CMS data ([iris-hep/idap-200gbps](https://github.com/iris-hep/idap-200gbps))
  - **UChicago**: analyze **Run-2 PHYSLITE** ATLAS data ([iris-hep/idap-200gbps-atlas](https://github.com/iris-hep/idap-200gbps-atlas))
  - similar tasks broadly, **important differences**: facilities, event sizes, object types, compression, ...

# Ingredients for 200 Gbps throughout



*team of experts from IRIS-HEP and beyond,  
rallied behind a shared vision*

*planning, structure, schedule*

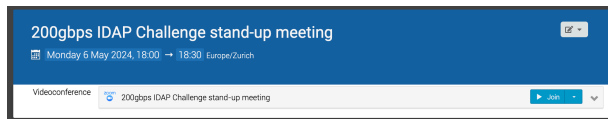


*100s of messages per day,  
dedicated communication channels*



*challenging timeline: **8 weeks**  
from first idea to WLCG/HSF workshop*

*dedicated meetings in multiple formats*



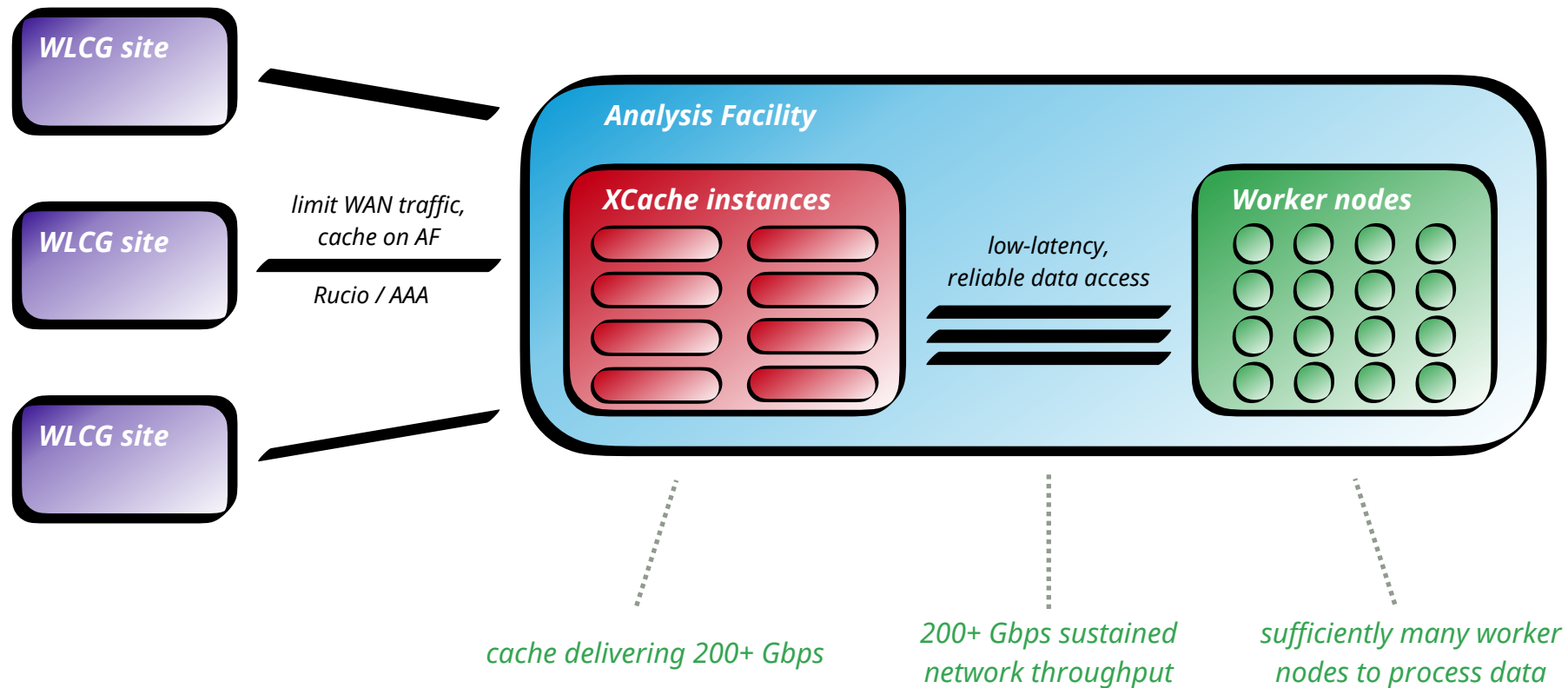
[image by DALL-E 3]



Demonstrator Analysis 200 Gb/s  
Hoersaal, DESY

Brian Paul Bockelman  
16:00 - 16:20

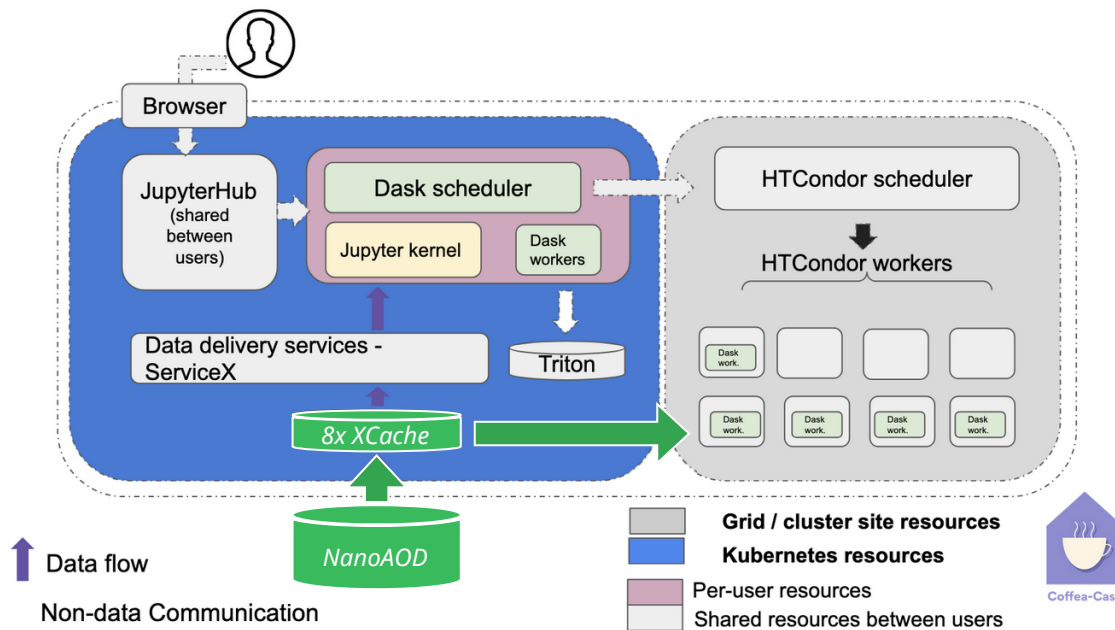
# Key Analysis Facility elements for 200 Gbps



# Coffea-casa at Nebraska: 200 Gbps setup

- **R&D prototype of a future Analysis Facility**

- designed as **hybrid setup** including **Kubernetes** and **Nebraska CMS Tier-2 / Tier-3 resources**



*using 8 XCache instances behind 2x100 Gbps uplink each*

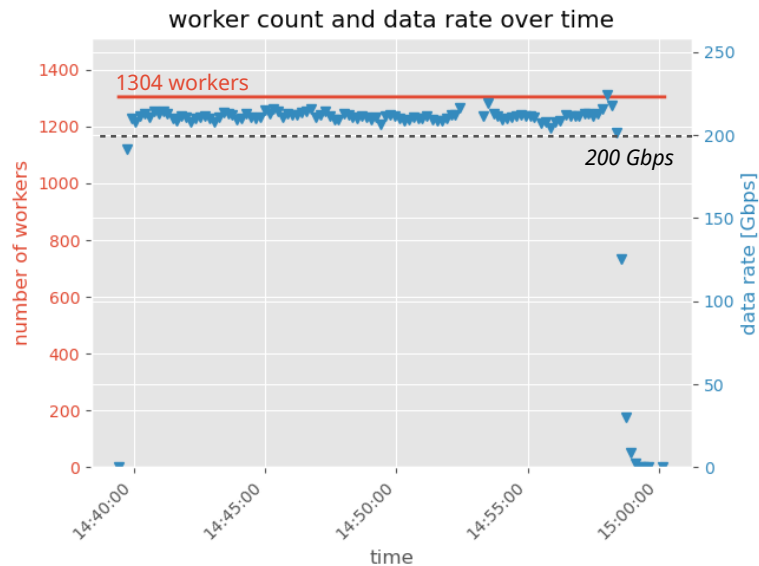
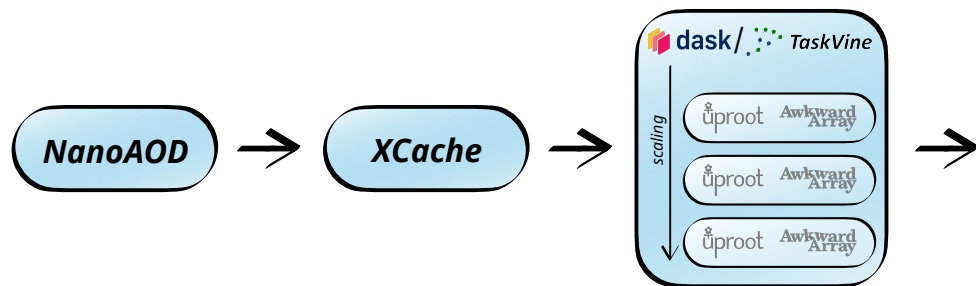
# Coffea-casa at Nebraska: measurements

- **200 Gbps sustained throughput of data for physics**

- scheduling with **Dask & TaskVine**, scaling with **HTCondor & Kubernetes**
- **re-compressed NanoAOD** (LZMA → ZSTD) for 2.5x event rate increase

*details for this example:*

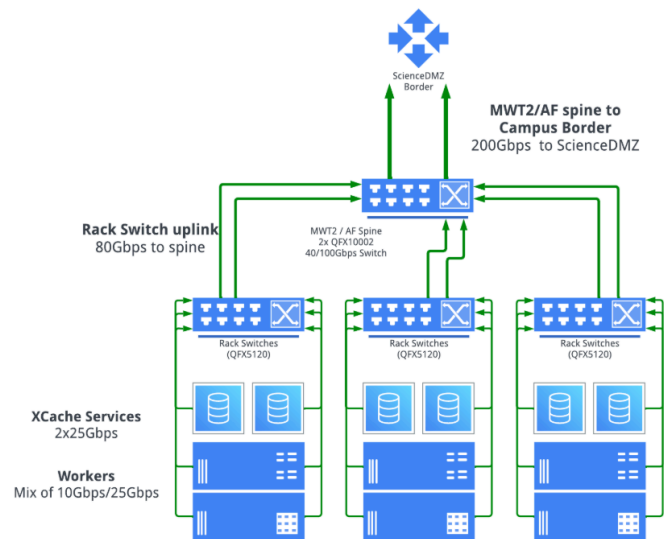
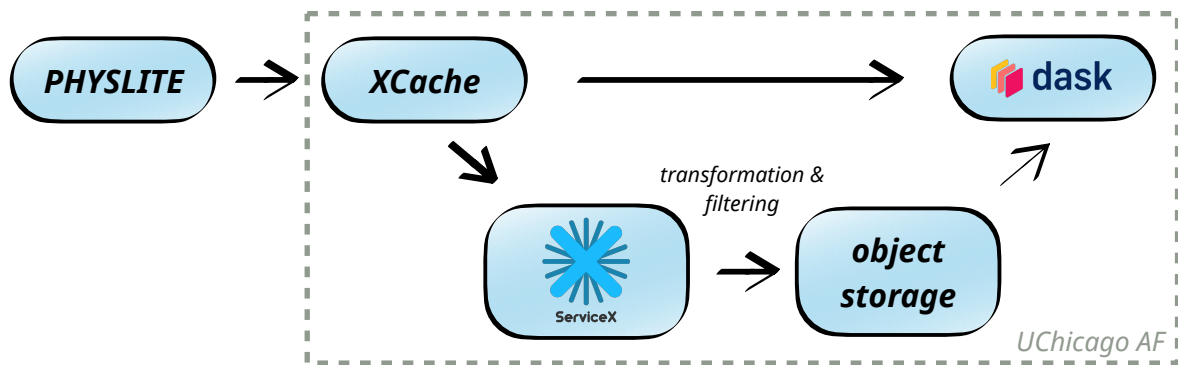
- 40 B events, 64k files
- 1304 workers
- 32 MHz event rate
- data processed (compressed): 30 TB
- data processed (uncompressed): 71 TB



*200+ Gbps with Dask + HTCondor*

# UChicago AF: 200 Gbps setup

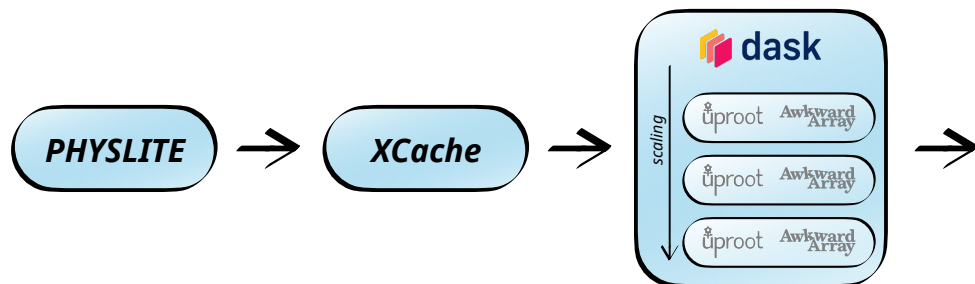
- **Production Analysis Facility for ATLAS**
  - built on **Kubernetes**, **partially reconfigured** for needs of challenge
- **Two configurations explored** with Kubernetes scaling (HTCondor available)
  - **uproot** scaled with **Dask** reading from **XCache**
  - **ServiceX** as data delivery service writing to object storage



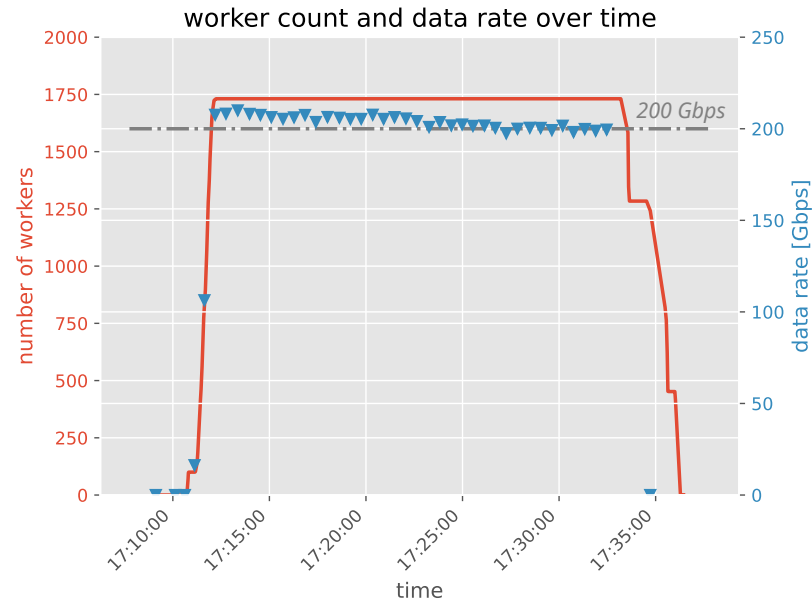
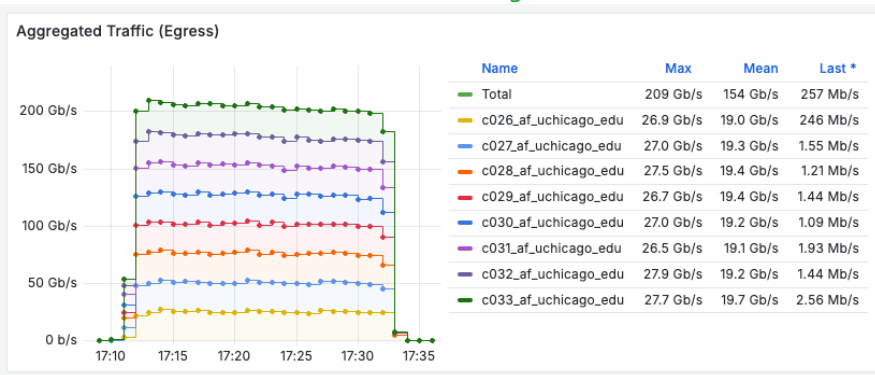
*8 XCache instances total,  
distributed to maximize bandwidth*

# UChicago AF: measurements

- **200 Gbps sustained throughput of data for physics**



network monitoring

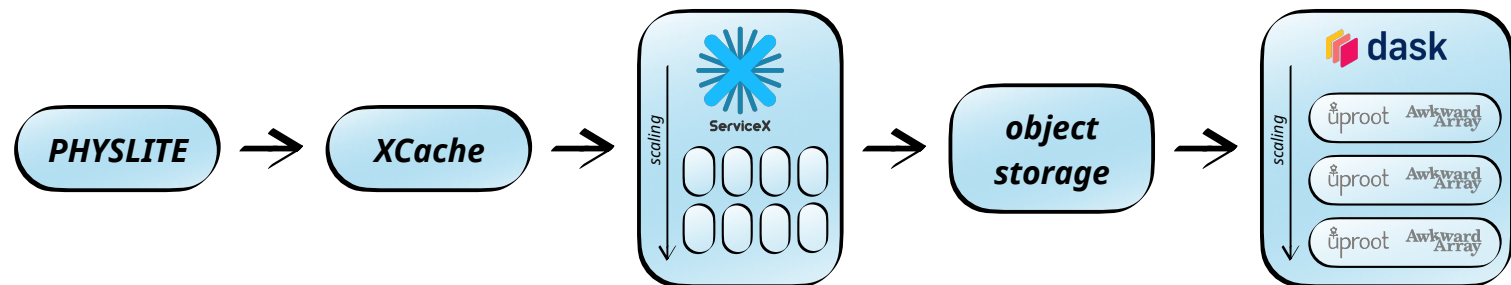


more details:

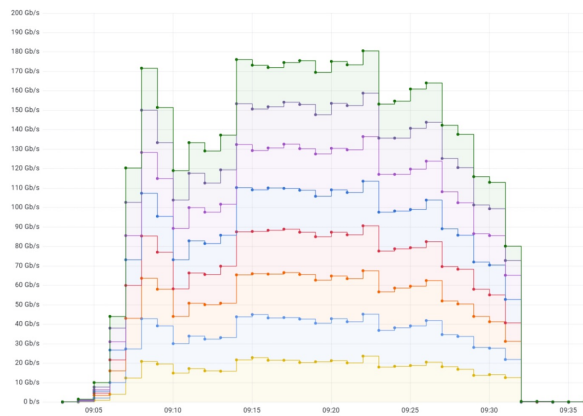
- 32 B events, 190 TB data, 218k files
- 1739 workers peak
- 15 MHz event rate, 5–20 kHz per core
- 200 Gbps throughput sustained
- data processed (compressed): 32 TB
- data processed (uncompressed): 80 TB

# ServiceX as data delivery service

- **Idea:** filter data with **ServiceX**, then further process output with Dask
  - rapid turnaround from **cached ServiceX output**



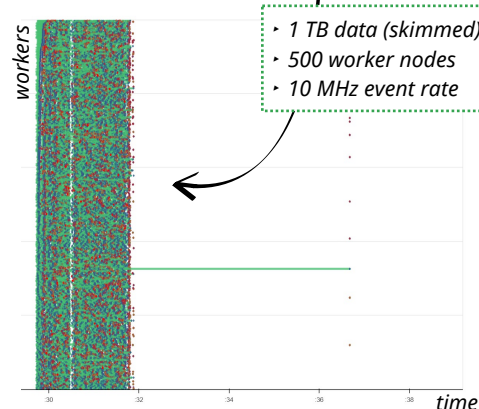
170 Gbps parallel data processing with ServiceX



- 19 B events, 146 TB data
- up to 1k pods
- 10 MHz event rate

multi-stage processing schema,  
transparent to users

Dask tasks





# Towards multi-user interactivity

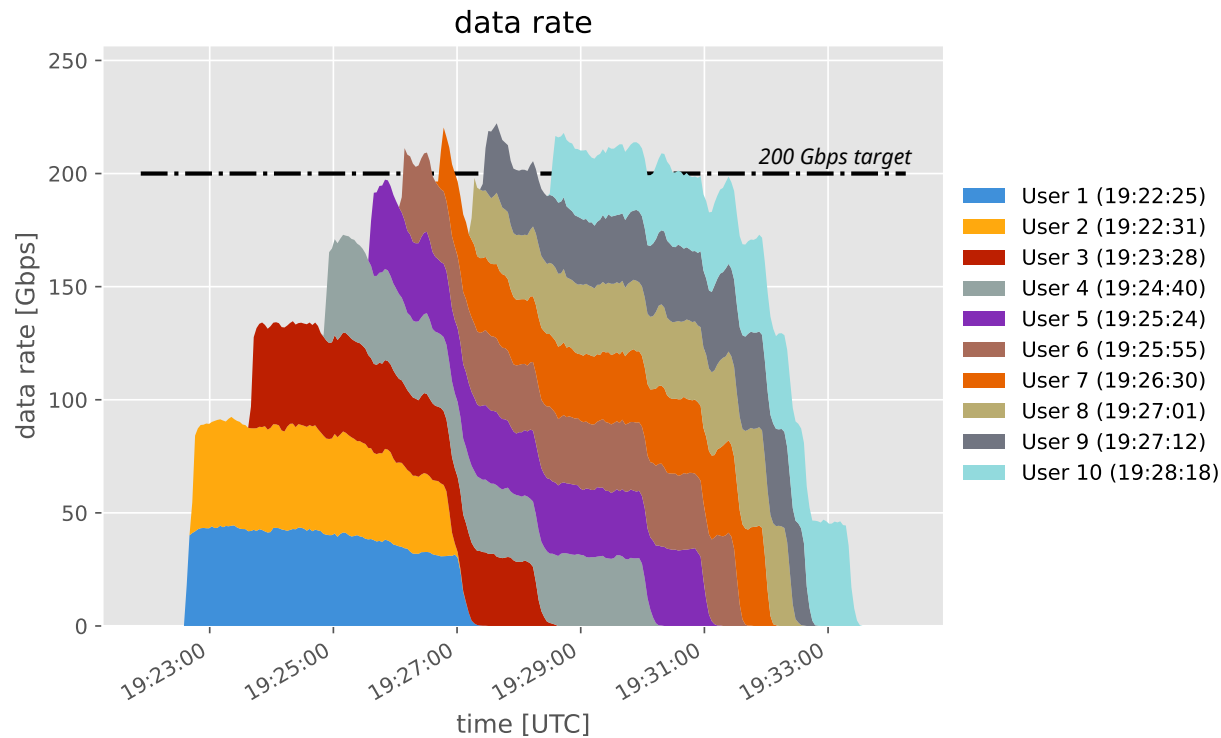
- Analysis Facilities will host **many users** looking to achieve **interactive turnaround simultaneously**
  - ensure **scale testing** includes number of users
- **Live exercise** at 2024 IRIS-HEP retreat: **200 Gbps setup with 16 participants**
  - **automatic CPU scaling** with Dask
  - limited **maximum number of cores** per user
- Intended as part of a bigger **discussion about fair-share & interactivity**



*live demonstration with retreat participants*



# Bandwidth shared between multiple users



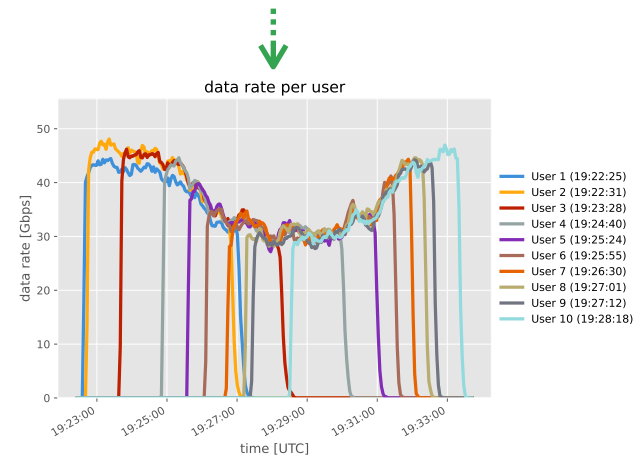
*task launch times were randomly distributed to simulate reality of random submissions*

- Test with **ten simultaneous users at UChicago**

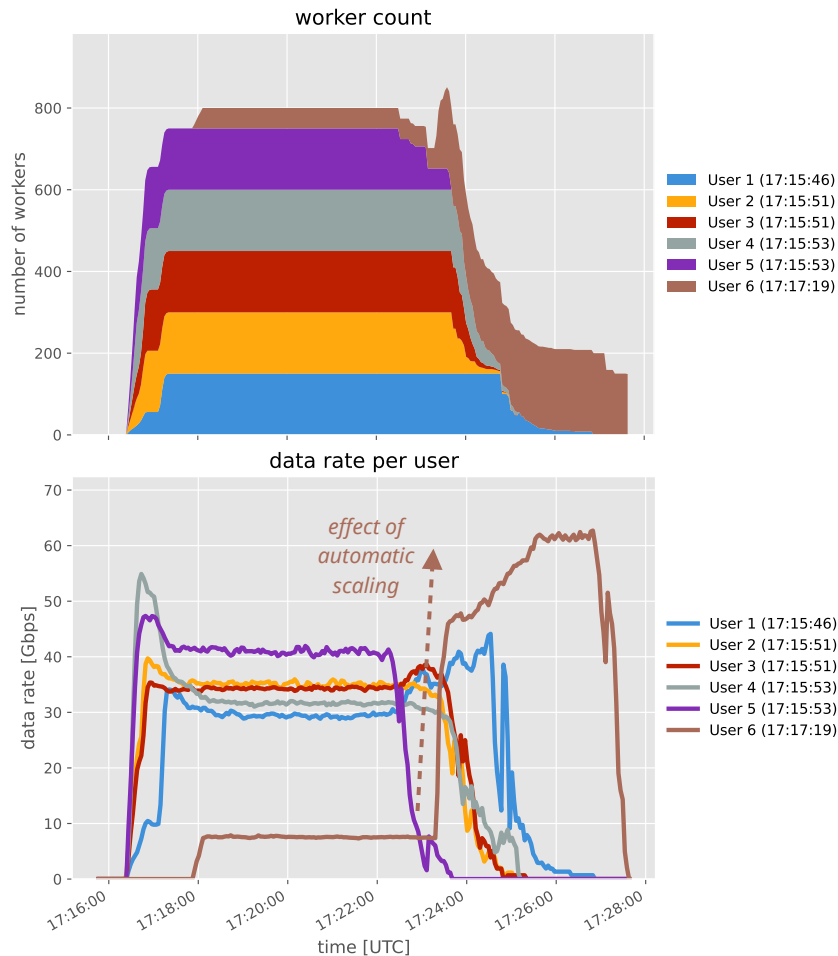
- users limited to **max 200 cores**

- Reached **200 Gbps collectively**

- **network saturation** effect visible



# Automatic worker scaling

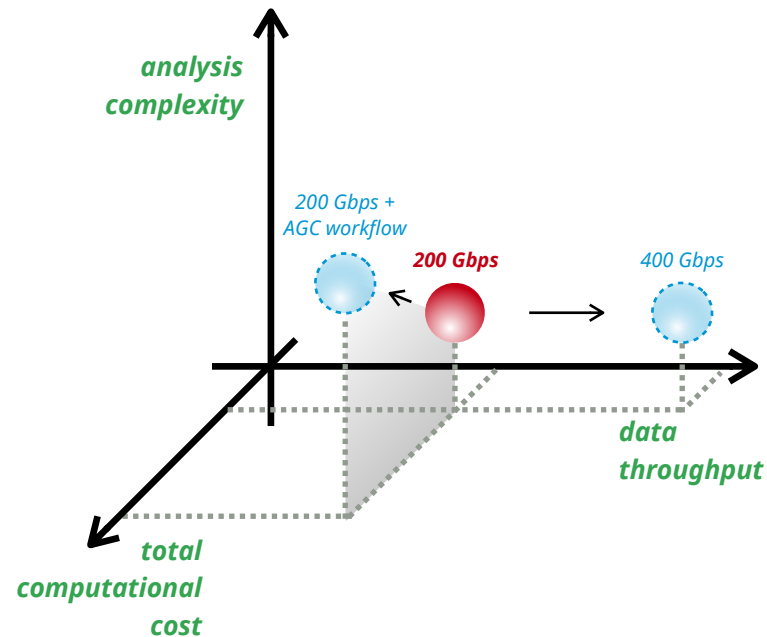


- Test with **six simultaneous users at Nebraska**
- First five users launch at the same time
  - **stable parallel processing**
- **User 6** receives last available cores, slower initially
  - rapid automatic scaling once more resources become available

Where to go from here?

# Next steps

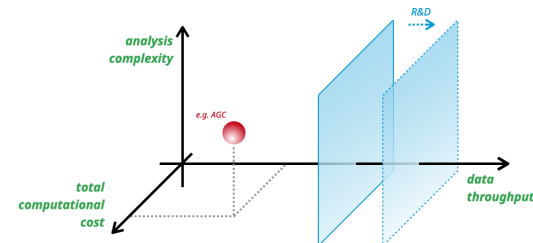
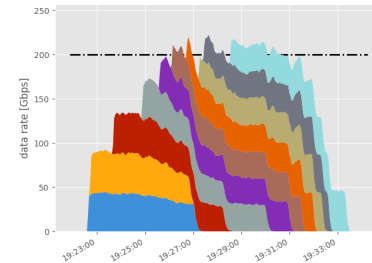
- **Further explore the parameter space** of HL-LHC analyses
  - extend 200 Gbps setup towards **full AGC-type workflow**
  - medium term: **400 Gbps exercise**
- Further **collaboration with community & knowledge sharing**
  - lessons learned **help analyses already today**
- Help **inform Analysis Facility evolution**



*axis closely connected to  
environmental sustainability  
→ [David Britton's talk](#)*

# Conclusion

- **Successful 200 Gbps Challenge** shows technology readiness, checkpoint towards HL-LHC
  - extremely valuable project to generate **feedback** and identify **potential bottlenecks**
  - planned **extensions for more realism** (Analysis Grand Challenge)
- **Difficult to predict future** of HEP end-user data analysis
  - **factorize challenges** & push boundaries, **remain open** to new ideas
- **Close collaboration remains crucial:** take advantage of CHEP to connect to colleagues from different areas!
  - **physics is the end goal**



Also check out these *related contributions* this week!

[Investigating Data Access Models for ATLAS: A Case Study with FABRIC Across Borders and ServiceX](#) [Oct 22, 15:18]

[GIL-free scaling of Uproot in Python 3.13](#) [Oct 23, 14:42]

[Building Scalable Analysis Infrastructure for ATLAS](#) [Oct 23, 17:09]

[Operating the 200 Gbps IRIS-HEP Demonstrator for ATLAS](#) [Oct 24, 16:33]

[Tuning the CMS Coffea-casa facility for 200 Gbps Challenge](#) [Oct 24, 16:51]

[Benchmarking massively-parallel Analysis Grand Challenge workflows using Snakemake and REANA](#) [Oct 24, 17:45]

- **Contact:** join us at [analysis-grand-challenge@iris-hep.org](mailto:analysis-grand-challenge@iris-hep.org) (sign up: [google group link](#))