

# Kubernetes as a Substrate for ATLAS Compute

Fernando Barreiro Megino (University of Texas at Arlington)

Lukas Heinrich (TU München)

KubeCon + CloudNativeCon Europe 2022, 18 May 2022, Valencia, Spain



UNIVERSITY OF  
TEXAS  
ARLINGTON



# Presenters

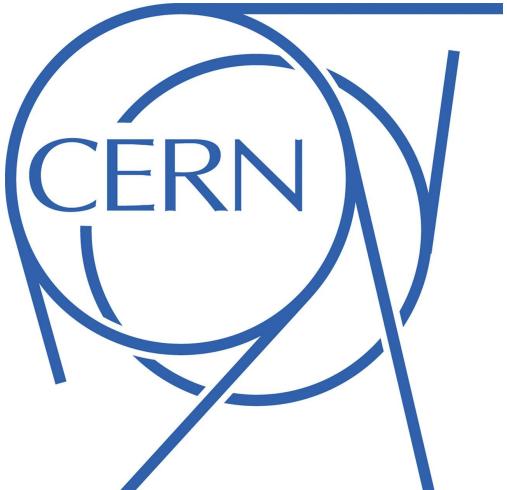


Fernando Barreiro works as a Computing Engineer at the ATLAS Experiment at CERN. He is part of ATLAS' core team for workload management. In recent years he has led the effort to integrate public and private cloud resources through native usage of Kubernetes



Lukas Heinrich is a particle physicist working on the ATLAS Experiment at the Large Hadron Collider. He focuses on introducing modern cloud computing and data science tools to more systematically search for phenomena beyond the Standard Model of Particle Physics.

# CERN: Where the Web was born



CERN DD/OC

Information Management: A Proposal

Tim Berners-Lee, CERN/DD

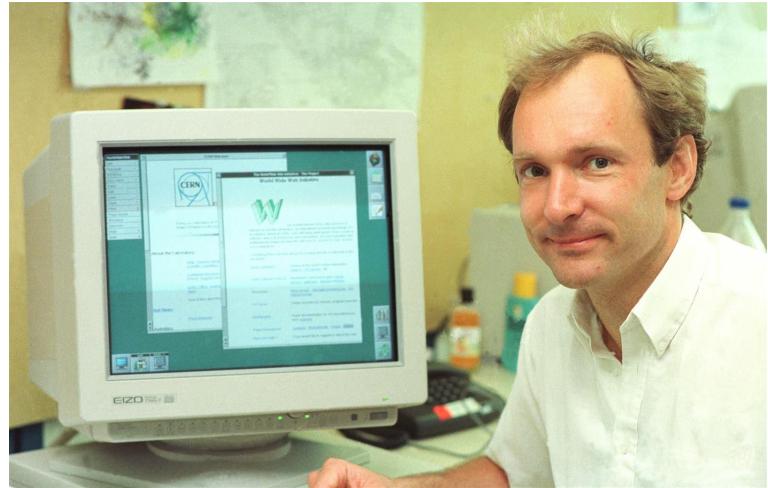
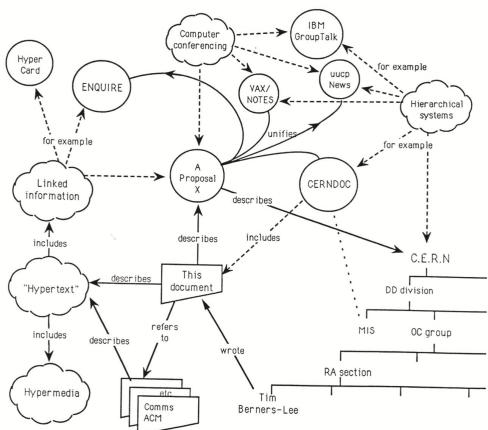
March 1989

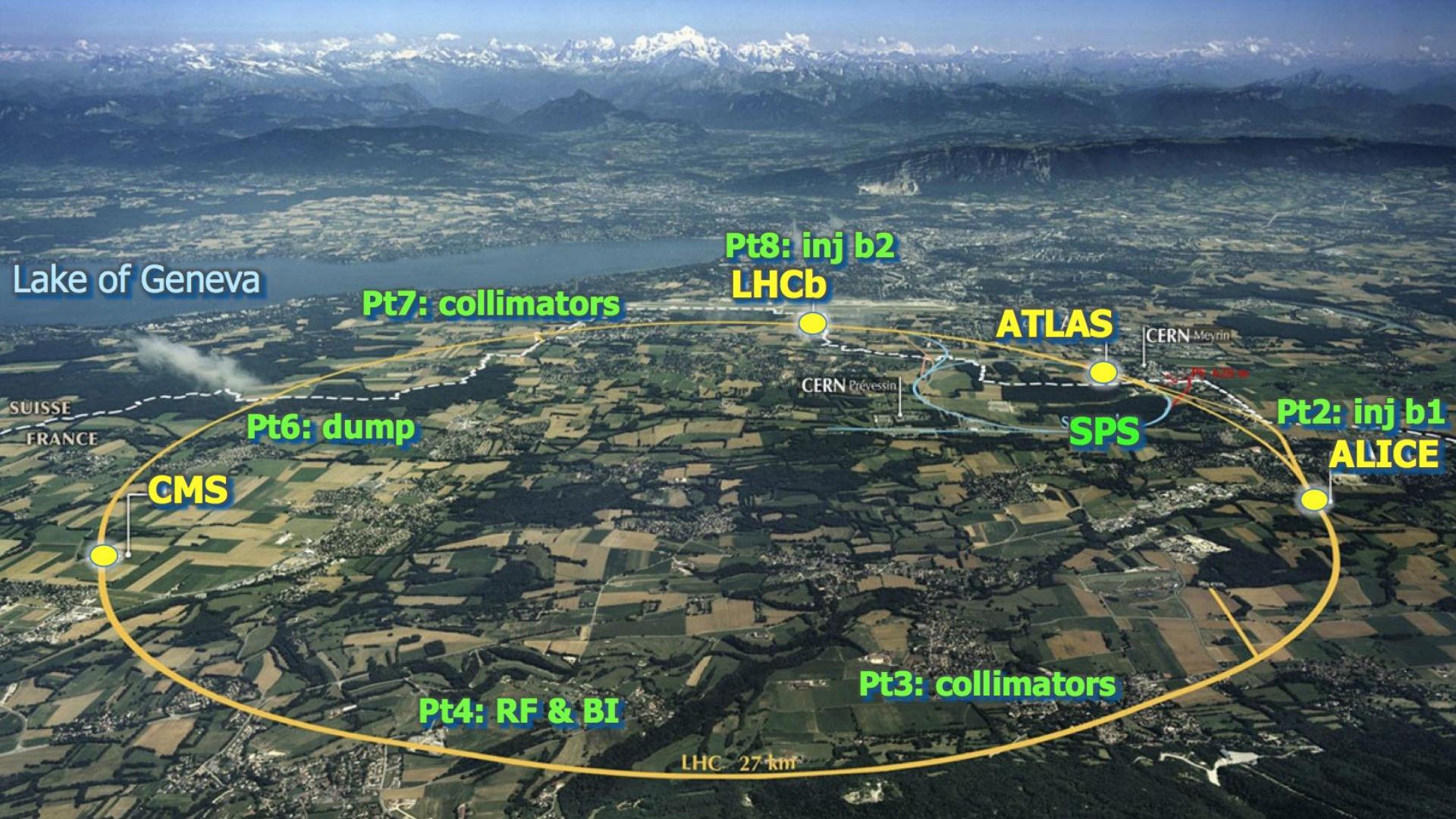
## Information Management: A Proposal

### Abstract

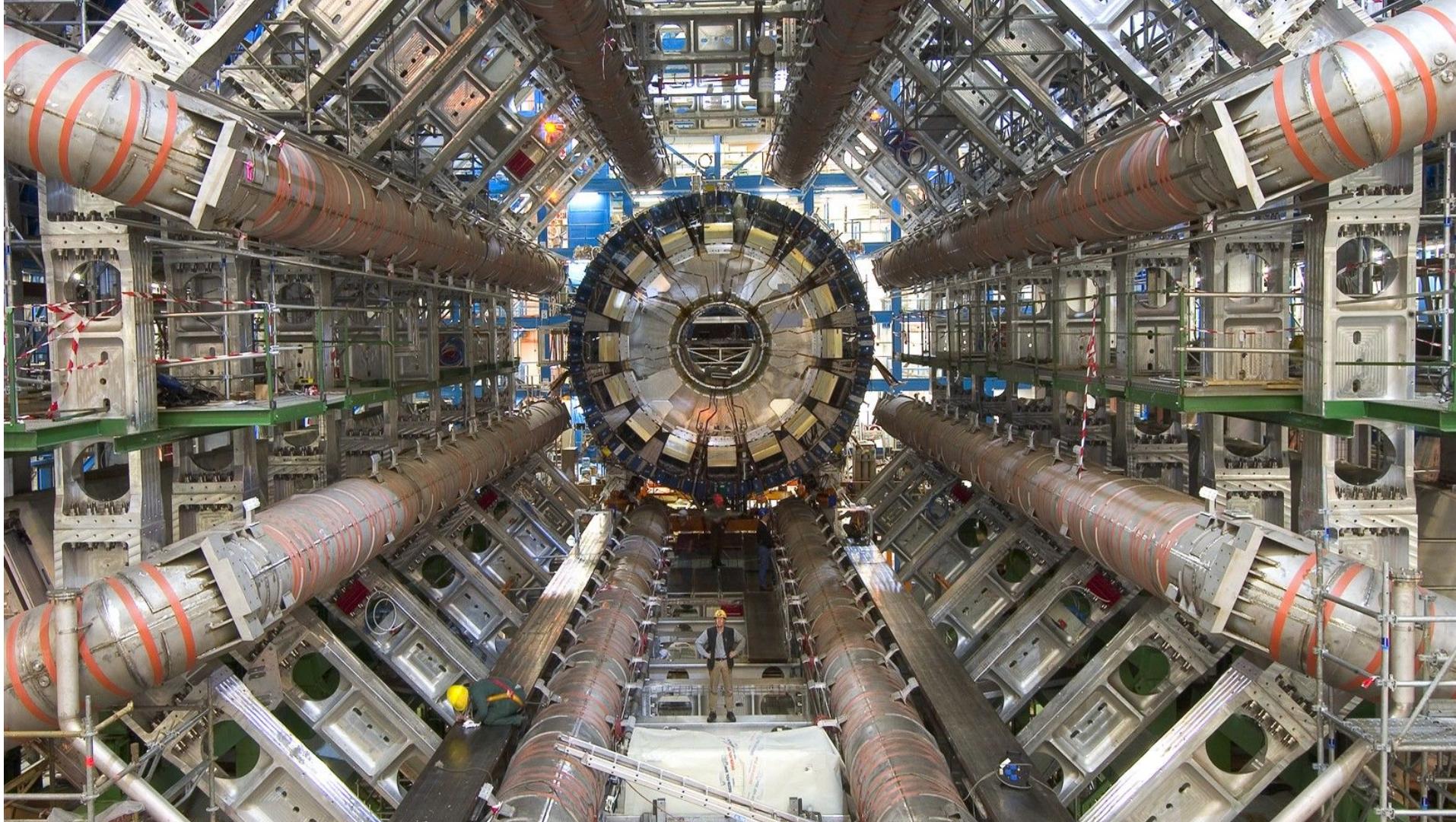
This proposal concerns the management of general information about accelerators and experiments at CERN. It discusses the problems of loss of information about complex evolving systems and derives a solution based on a distributed hypertext system.

Keywords: Hypertext, Computer conferencing, Document retrieval, Information management, Project control









# Physicists Find Elusive Particle Seen as Key to Universe

 Give this article    122



Scientists in Geneva on Wednesday applauded the discovery of a subatomic particle that looks like the Higgs boson. Pool photo by Denis Balibouse

By Dennis Overbye  
July 4, 2012

← 10 years ago...

ASPEN, Colo. — Signaling a likely end to one of the longest, most expensive searches in the history of science, physicists said Wednesday that they had discovered a new subatomic particle that looks for all the world like the Higgs boson, a key to understanding



## NEWS

[Home](#) | [War in Ukraine](#) | [Coronavirus](#) | [Climate](#) | [Video](#) | [World](#) | [UK](#) | [Business](#) | [Tech](#) | [Science](#) | [Sports](#)

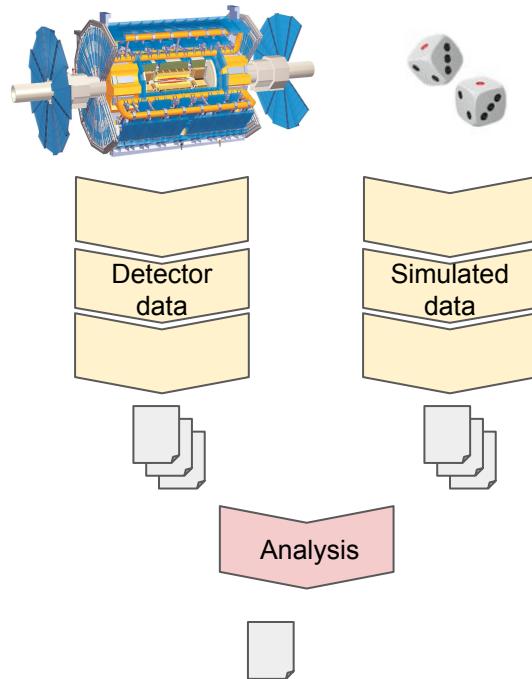
[Science](#)

## Higgs boson scientists win Nobel prize in physics

By James Morgan  
Science reporter, BBC News

© 8 October 2013 |  [Comments](#)

# ATLAS computing flow



## “Production”

Exabyte scale

$O(10^5)$ - $O(10^6)$  vCPUs

$O(10)$  groups following organised campaigns

Turnaround expectation: weeks

Batch processing

## “Analysis”

Petabyte scale

$O(10^5)$  vCPUs

$O(100)$  simultaneous individual users

Turnaround expectation: ASAP / days

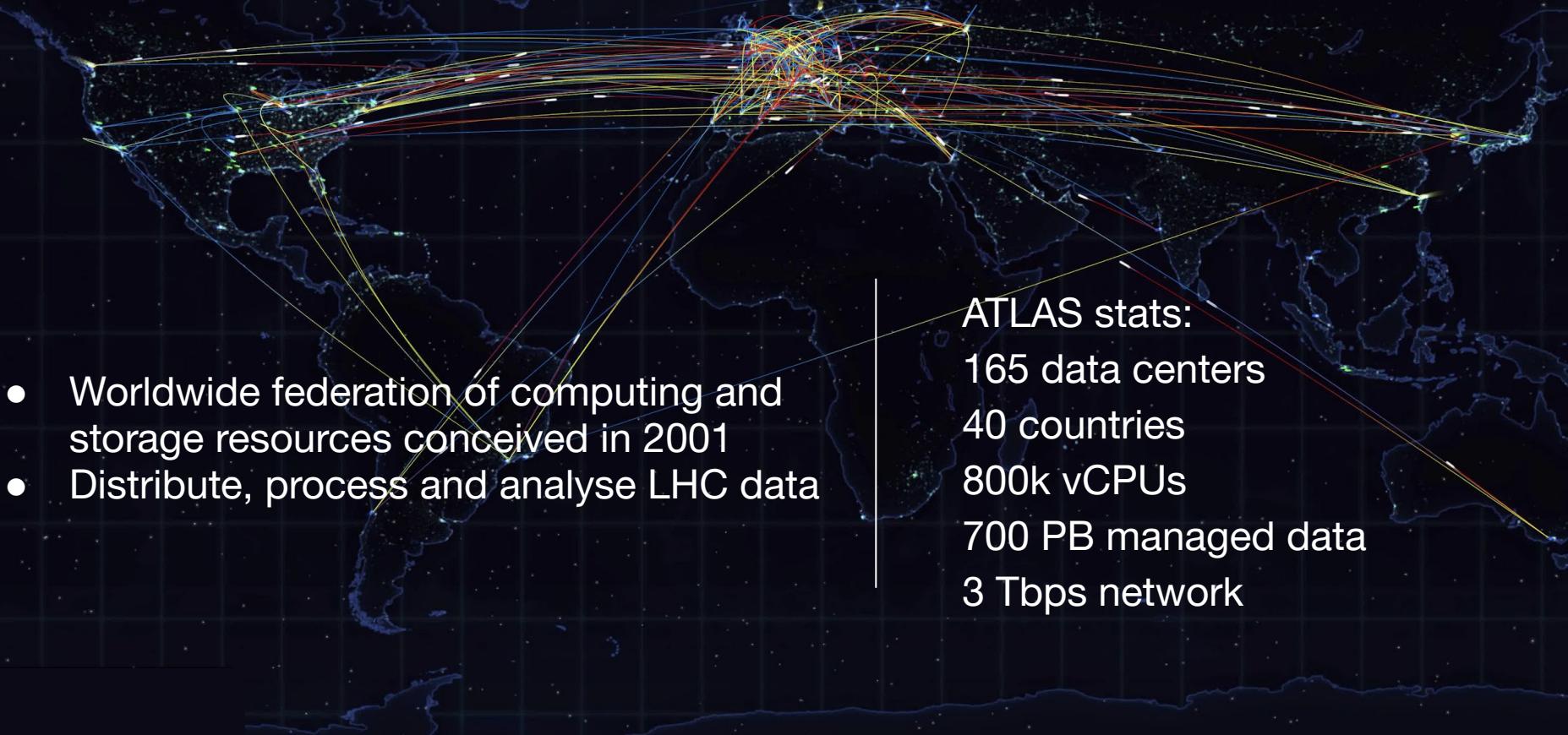
Batch and interactive processing



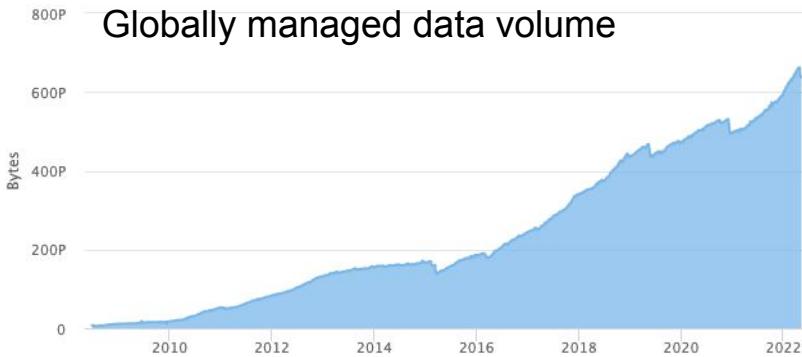
# Worldwide LHC Computing Grid



**WLCG**  
Worldwide LHC Computing Grid

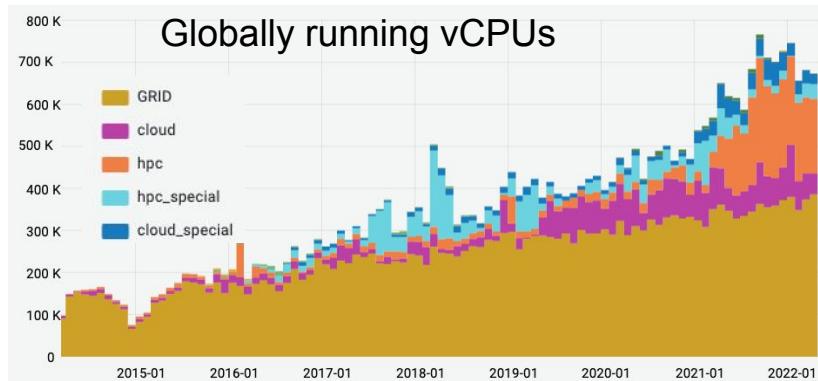


# ATLAS Distributed Computing



## Data management: Rucio

- Manage experiment data
- Interact with storage systems

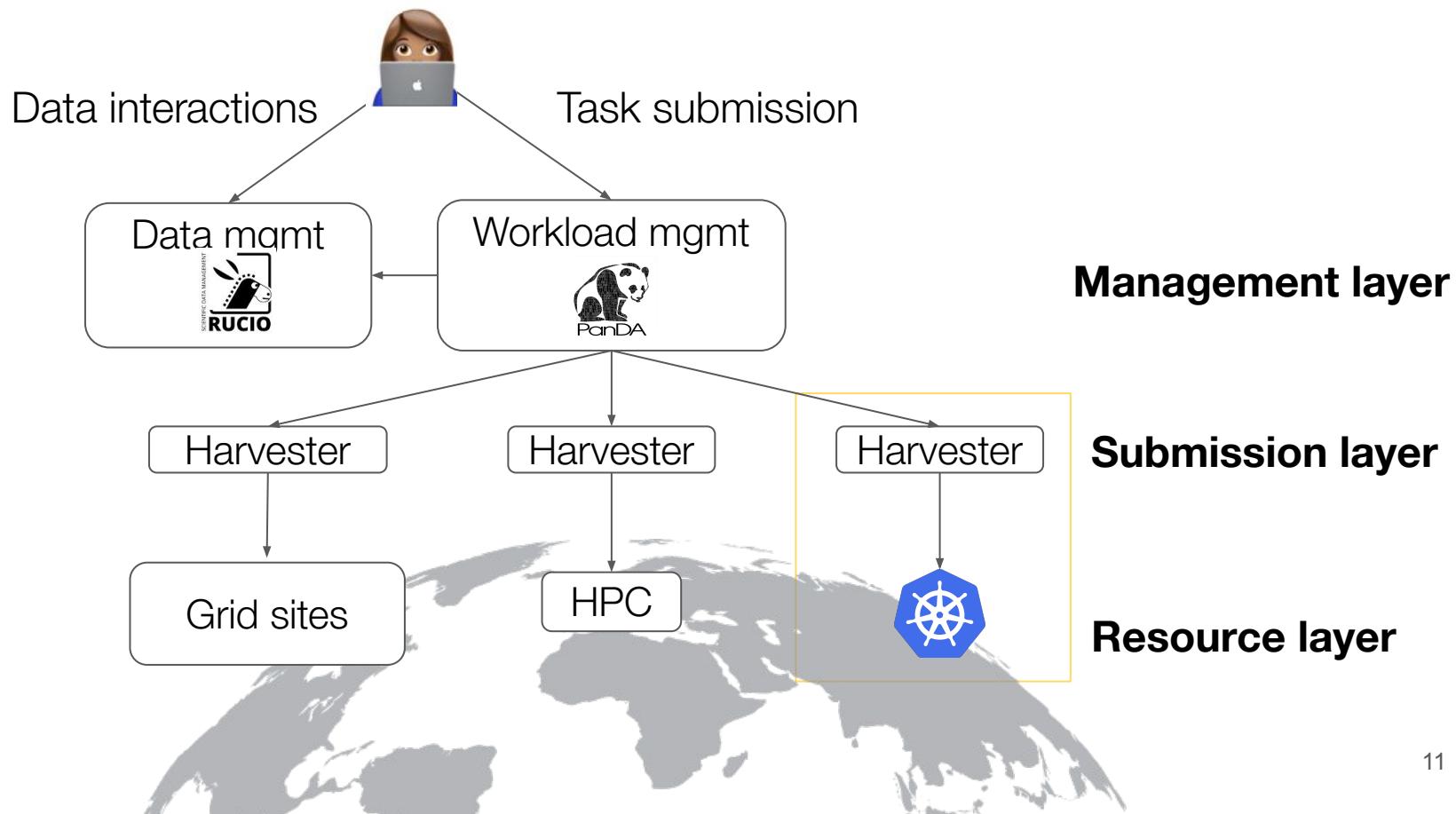


## Workload management: PanDA

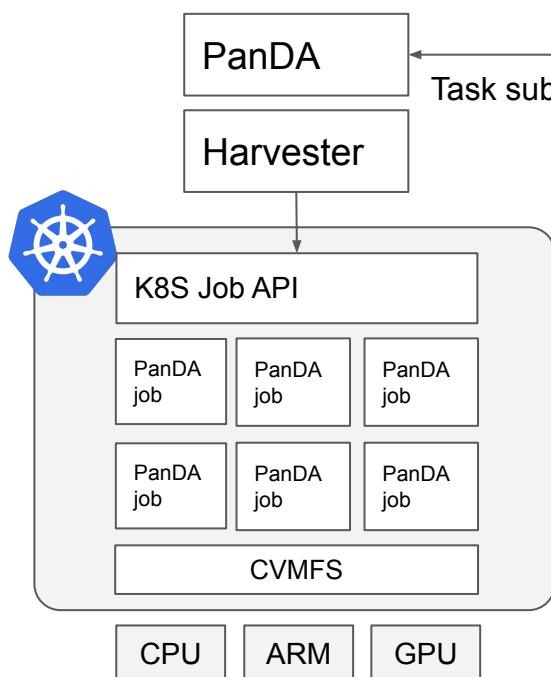
- Broker computational tasks
- Interact with compute systems (Harvester)



# ATLAS job submission through Harvester



# Grid processing: Harvester-Kubernetes integration

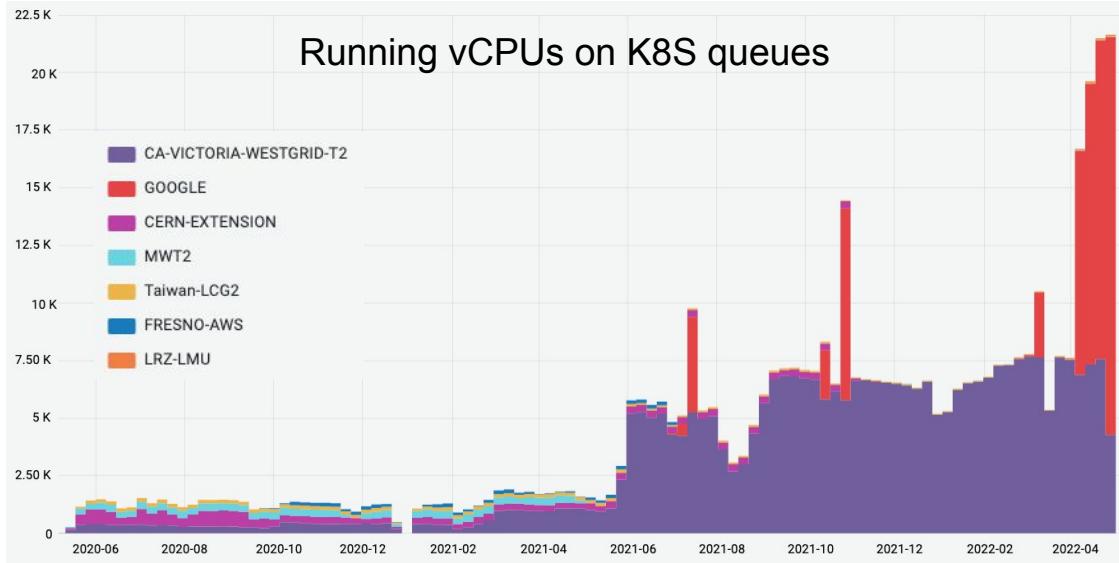


- Integrate K8S to work as any Grid site, while keeping setup as simple as possible
- Purely based on native K8S Job controller
- Harvester plug-ins for job submission & monitoring
- Usage of K8S common options
  - Limits for CPU, memory, disk, duration
  - Pod affinity to pack nodes
  - Priority classes
- High Energy Physics filesystem: CVMFS daemonset



Batch processing

# Onward & upward: on prem and cloud K8S queues



Mini K8S Grid scaled from few hundred vCPUs to many thousands vCPUs  
Provider independent integration



University  
of Victoria



arbutus  
cloud



Google Cloud



THE UNIVERSITY OF  
CHICAGO

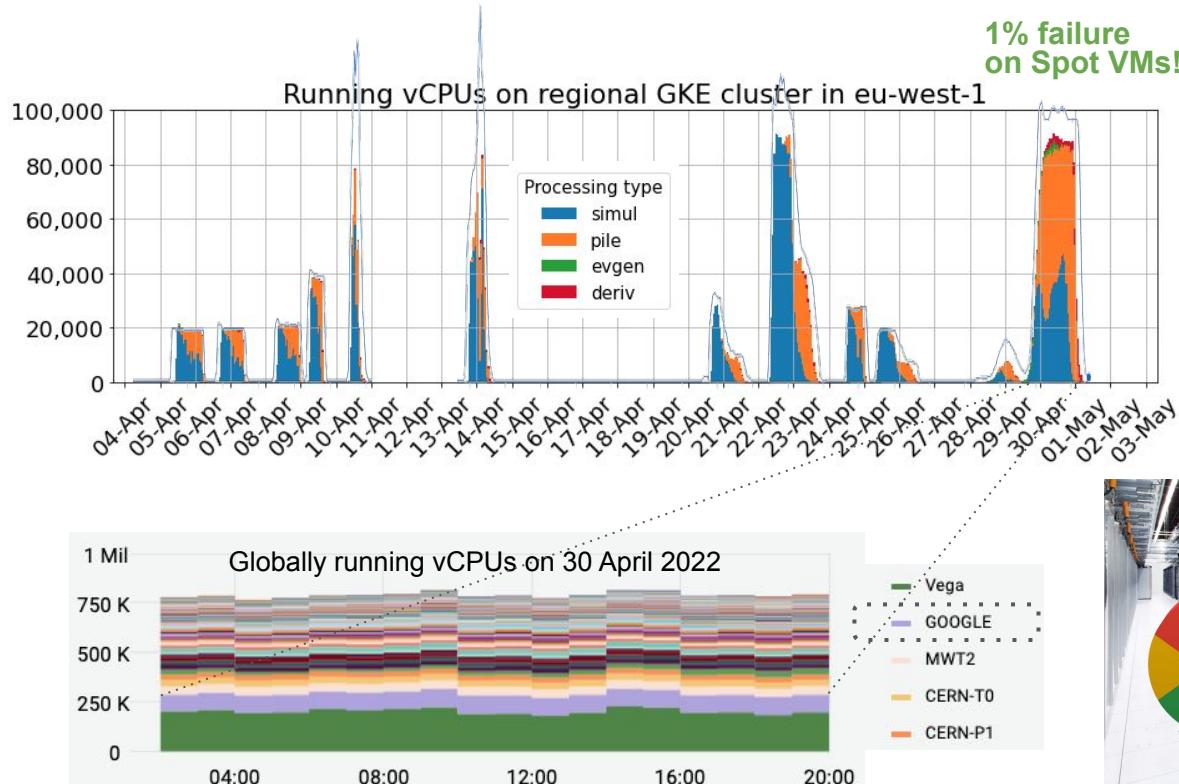


FRESNO STATE



Leibniz-Rechenzentrum  
der Bayerischen Akademie der Wissenschaften

# Elastic cloud scale out



1% failure  
on Spot VMs!

$O(10^8)$  events  
 $O(10^5)$  vCPUs  
 $O(10^4)$  Pods  
 $O(10^3)$  Nodes  
 $O(1)$  day  
 $O(1)$  Harvester instance  
<1 Engineer

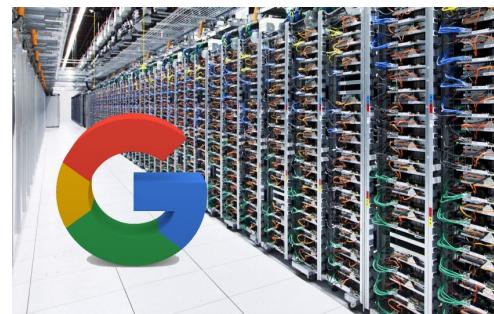
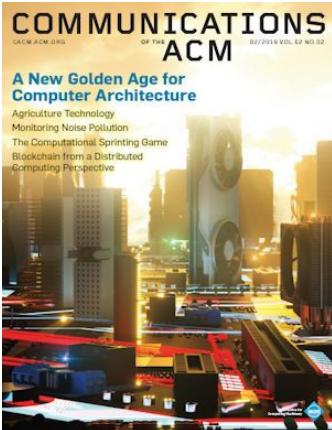


Image credits: ICT Network News

# Integration of heterogeneous architectures

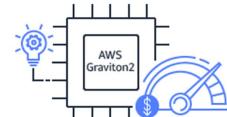
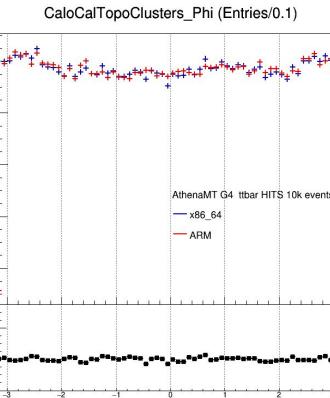
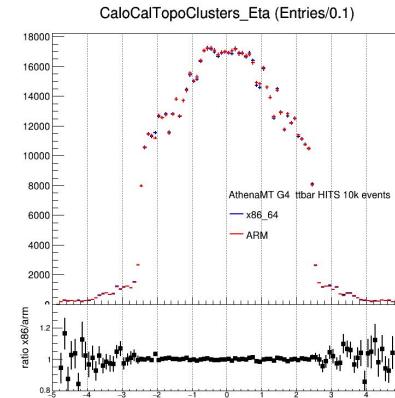


- Cloud queues backed by resources not commonly available on prem
- Straightforward to integrate different architectures, e.g. ARM, GPU
- Multi-arch Docker images doing the heavy lifting

Image credits: Communications of the ACM, February 2019

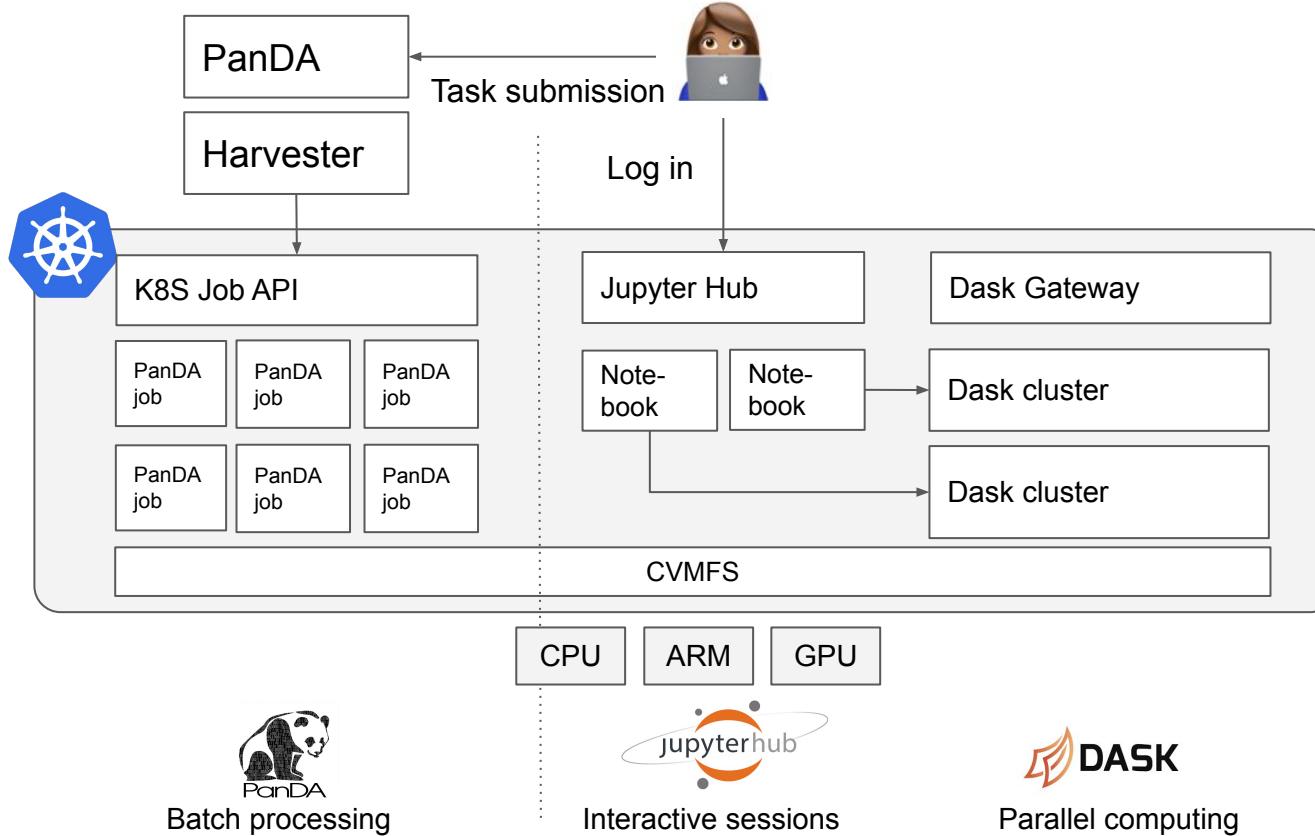


Image credits: NVIDIA GTC May 2020 Keynote



First ATLAS simulation task on ARM processors. Currently under physics validation. Generated on Amazon EKS cluster backed by Graviton 2 nodes.

# Adding interactive analysis facilities



# Distributed, interactive analysis

- Data science community technologies
- Helm charts available: DaskHub
- Configured for scalability and cost effectiveness on GKE

The screenshot shows a JupyterHub interface with several tabs open:

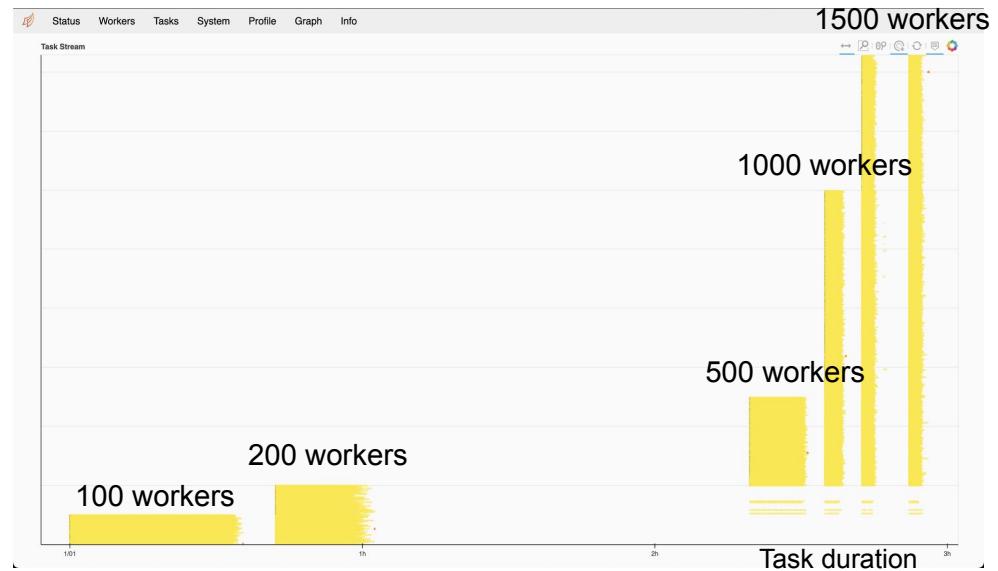
- File Edit View Run Kernel Tabs Settings Help**
- m2mu\_parquet.ipynb**: A code editor containing Python code for m2mu parquet processing.
- utils.py**: A code editor containing utility functions.
- Python 3 (pykernels)**: A kernel tab.

The code in **m2mu\_parquet.ipynb** includes:

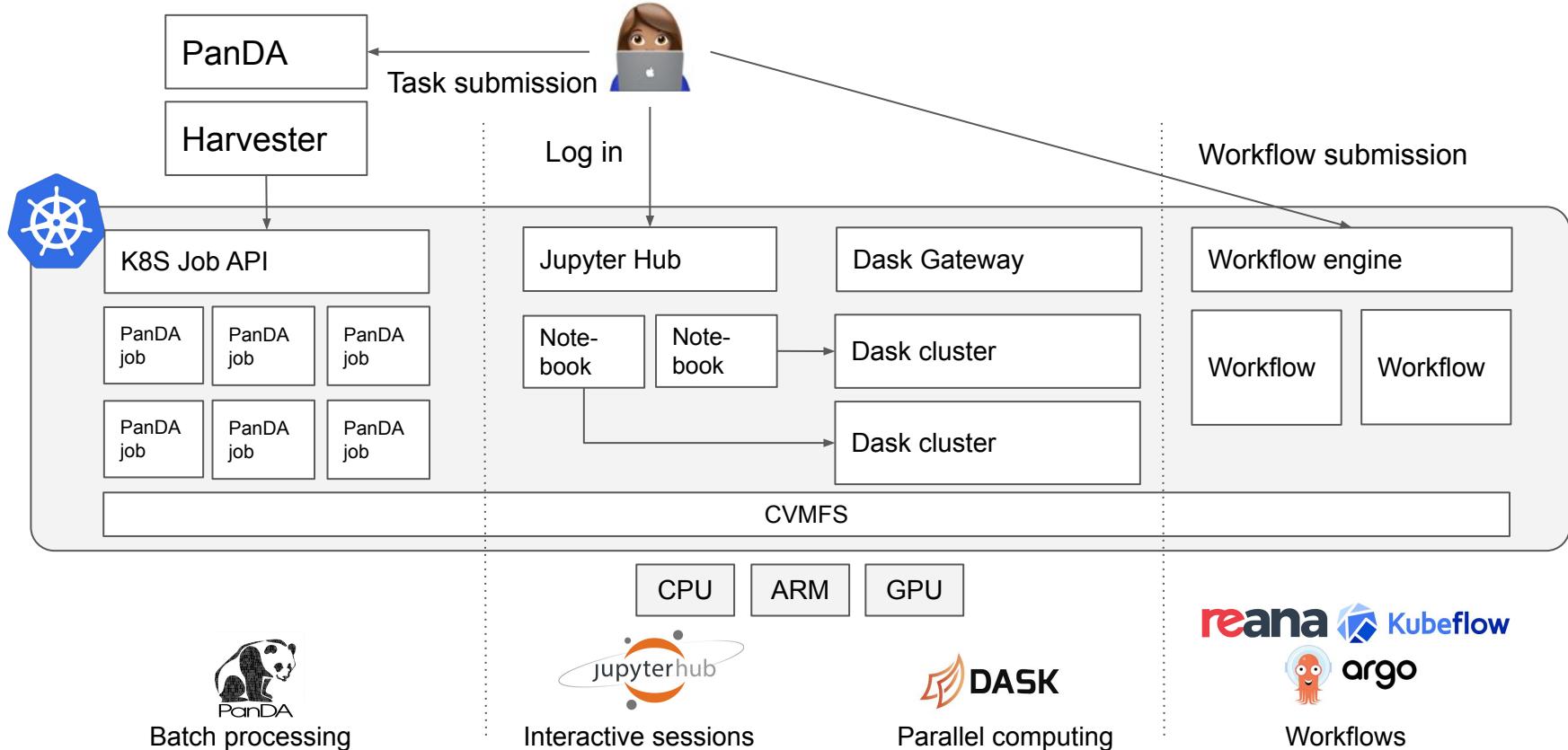
```
def client():
    Name: default.a35e420423d40e0f9bcc3bbd2f0e4ea
    Dashboard: /services/dask-gateway/clusters/default/a35e420423d40e0f9bcc3bbd2f0e4ea/sta
    Scale up cluster
    client = cluster.get_client()
    cluster.scale(465)
    Define distributed methods
    ...
    def get_signed_ur...
    ...
    def get_m2mu_ur...
    ...
    def get_m2mu_ur...
    ...
    tasks = []
    for file in files:
        d_url = dask.get_stores().url_worker[0].upload(data=1, file=file['scope'], file['name'])
        m2mu_list = dask.persist(get_m2mu_ur...
        m2mu = np.concatenate(dask.compute(m2mu_list))
        m2mu = ak.zip(m2mu, with_name='PtEtaPhiM LorentzVector')
        m2mu = ak.zip(ak.combinations(m2mu, 2))
        return ak.to_numpy(ak.flatten(m2mu + m2mu.mass))
    Run the m2mu computation in parallel, then aggregate the results
    ...
    Plot the results
    ...
    plt.hist(m2mu, bins=np.linspace(0, 150000, 500))
    plt.xlabel("m2mu")
    plt.title("m2mu distribution")
    plt.show()
```

**DASK** logo is visible at the bottom right.

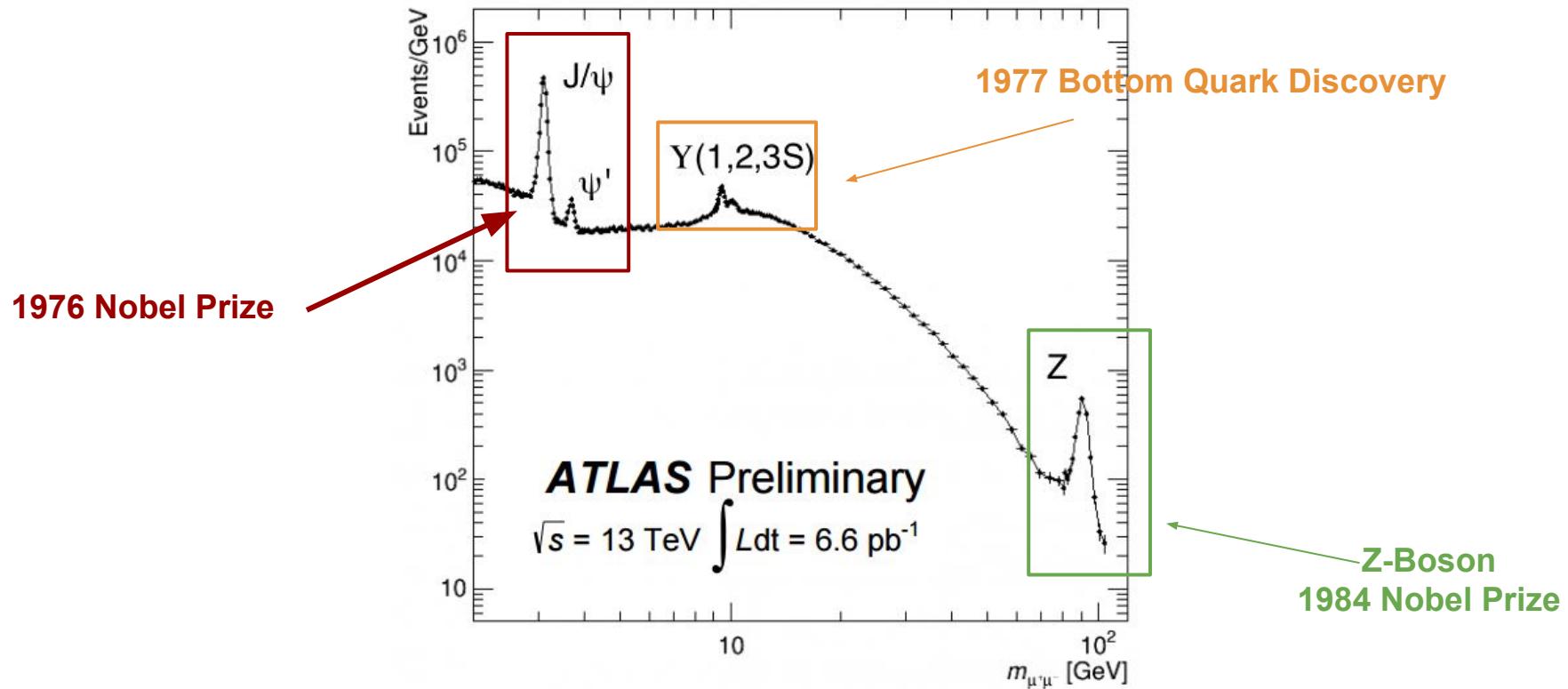
N workers



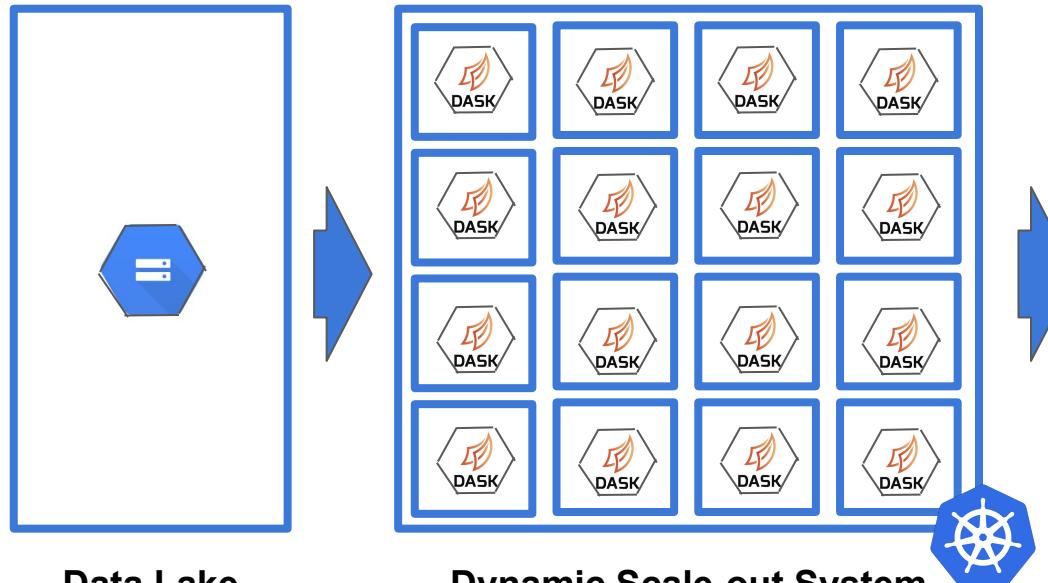
# K8S substrate for ATLAS compute: idealistic view



# Demo: History of Particle Physics in a single Graph



# Demo: History of Particle Physics in a single Graph



Data Lake

100 TB of Data

Dynamic Scale-out System

1-10k vCPU/User

A screenshot of a Jupyter Notebook environment. The left pane shows a file tree for a "kubernetes\_physics\_demo" directory containing various Python scripts and configuration files. The right pane contains two code cells and a terminal window. The top code cell runs a script to analyze particle physics events, and its output is a histogram showing a peak around 10^3 GeV. The bottom code cell shows command-line logs related to Kubernetes and the analysis process.

User-facing UI



# Conclusions

- Kubernetes goes beyond service management
  - Demonstrated K8S native ATLAS batch processing at 100k cores
    - Compatibility to existing infrastructure
    - Reaching scale further than originally planned
  - Next generation, interactive services with high elasticity
  - Other functionalities that can be added
- Next years will determine K8S integration
  - Zero-to-grid-site via GitOps & Helm?
  - Elasticity of resources and availability of “exotic” resources
  - Managed K8S clusters simply work
  - Deploy K8s inside of HPC & university settings
- A lot of work still to be done, but it is a very promising environment

# Acknowledgements

Karan Bhatia [4], Misha Borodin [6], Kaushik De [7], Johannes Elmsheuser [2], Miles Euell [4], Nikolai Hartmann [5], Alexei Klimentov [2], FaHui Lin [7], Tadashi Maeno [2], Usman Qureshi [4], Ricardo Rocha [3], Ming-Jyuan Yang [1]

1. Academia Sinica
2. Brookhaven National Laboratory
3. CERN
4. Google
5. Ludwig Maximilian University
6. University of Iowa
7. University of Texas at Arlington

# Links

CVMFS plugin: <https://github.com/sfiligoi/prp-osg-cvmfs>,

<https://github.com/PanDAWMS/prp-osg-cvmfs>

DaskHub: <https://github.com/dask/helm-chart/tree/main/daskhub>

Dask Gateway: <https://gateway.dask.org/>

Harvester: <https://github.com/HSF/harvester/>

PanDA: <https://panda-wms.readthedocs.io/>

Rucio: <https://rucio.cern.ch/>