

# Fast Columnar Physics Analyses of Terabyte-Scale LHC Data on a Cache-Aware Dask Cluster ([2207.08598](https://arxiv.org/abs/2207.08598))

Niclas Eich, Martin Erdmann, Peter Fackeldey, Benjamin Fischer, [Dennis Noll](#), Yannik Rath

IRIS-HEP Topical Meeting

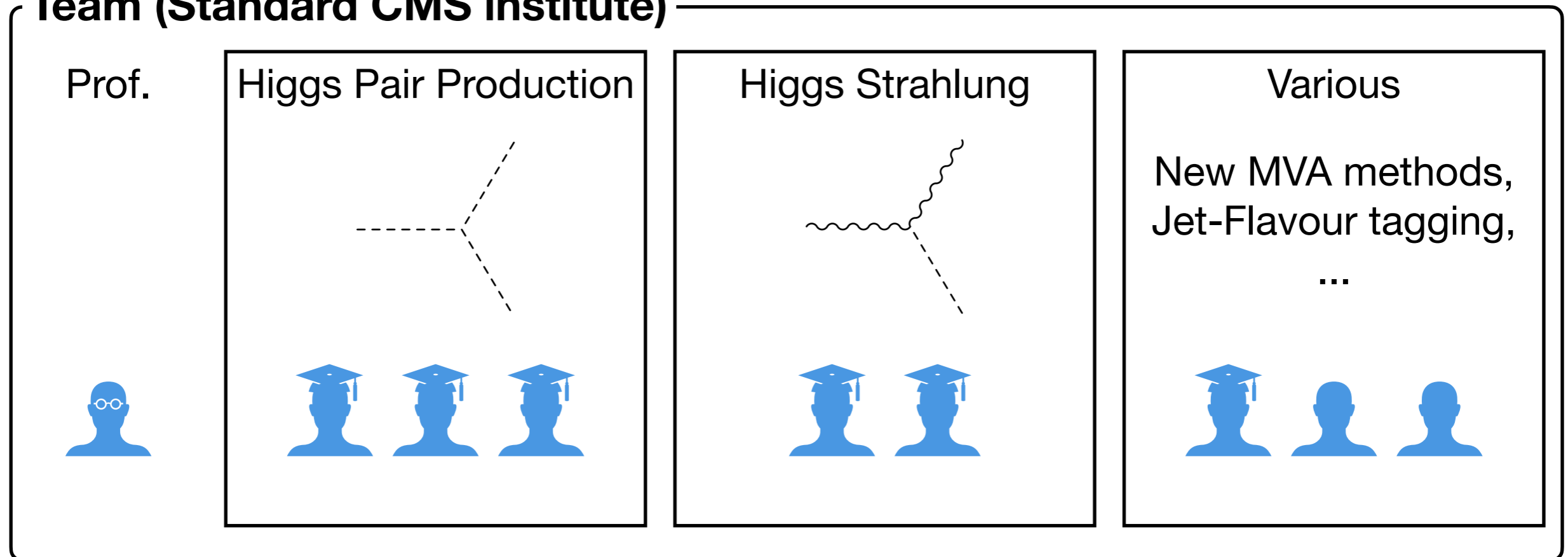
17.10.22





- Small time-to-insight drives high physics potential:
  - Tuning of existing methods
  - Investigations of new methods
- Ultimate goal: Analysis in duration of a coffee break
- Also: Crucial for HL LHC analyses

## Team (Standard CMS institute)

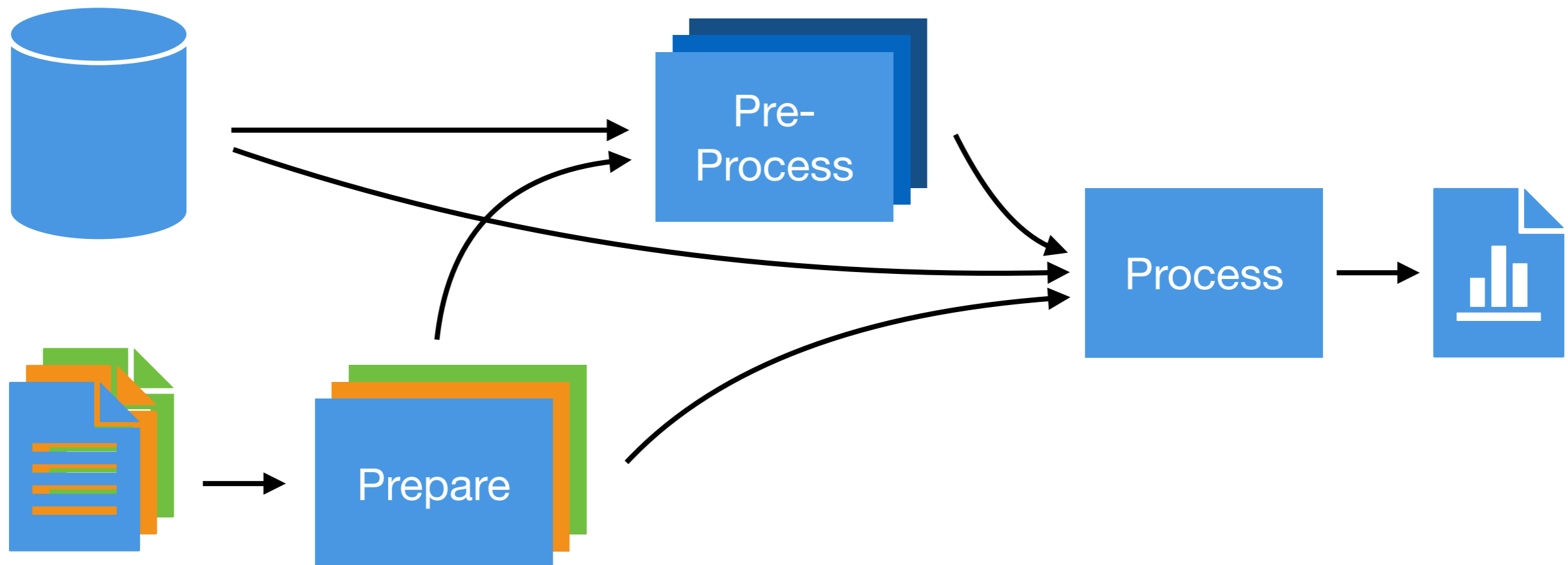


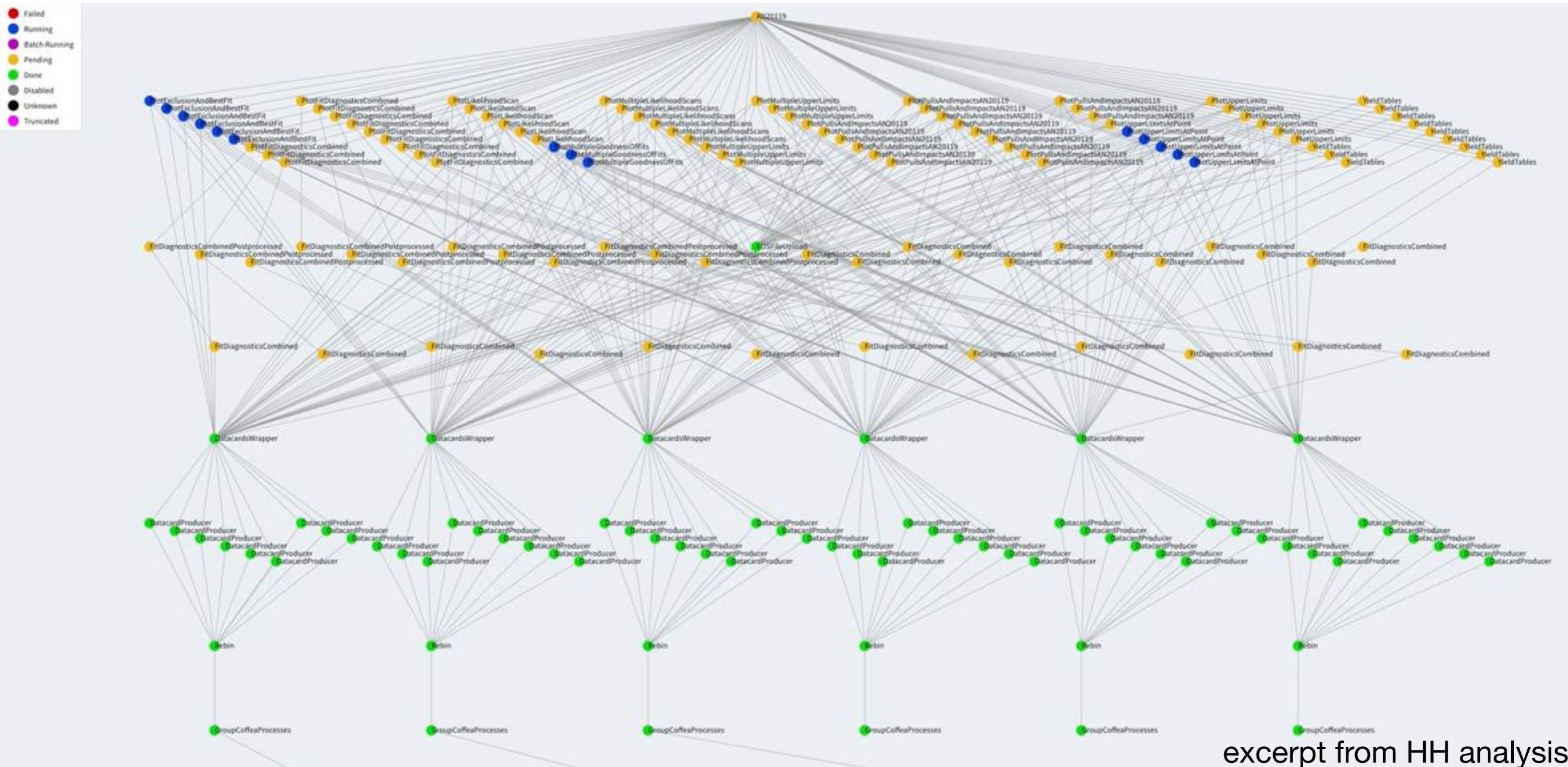
- Analyses have common challenges:
  - Many inputs:
    - Simulations and recorded data
    - Meta data (Efficiency factors, Corrections, ...)
  - Heavy computations (Event reconstruction, MVA algorithms, ...)
  - Bookkeeping
- Shared software repository

- Inputs: Event data, meta data
- O(100) different computing steps:
  - Preparation of meta data
  - Pre-processing of event data
  - Main-processing of event data
- Outputs: Histograms, plots, ...
- Tasks connected using law

## Typical workflows

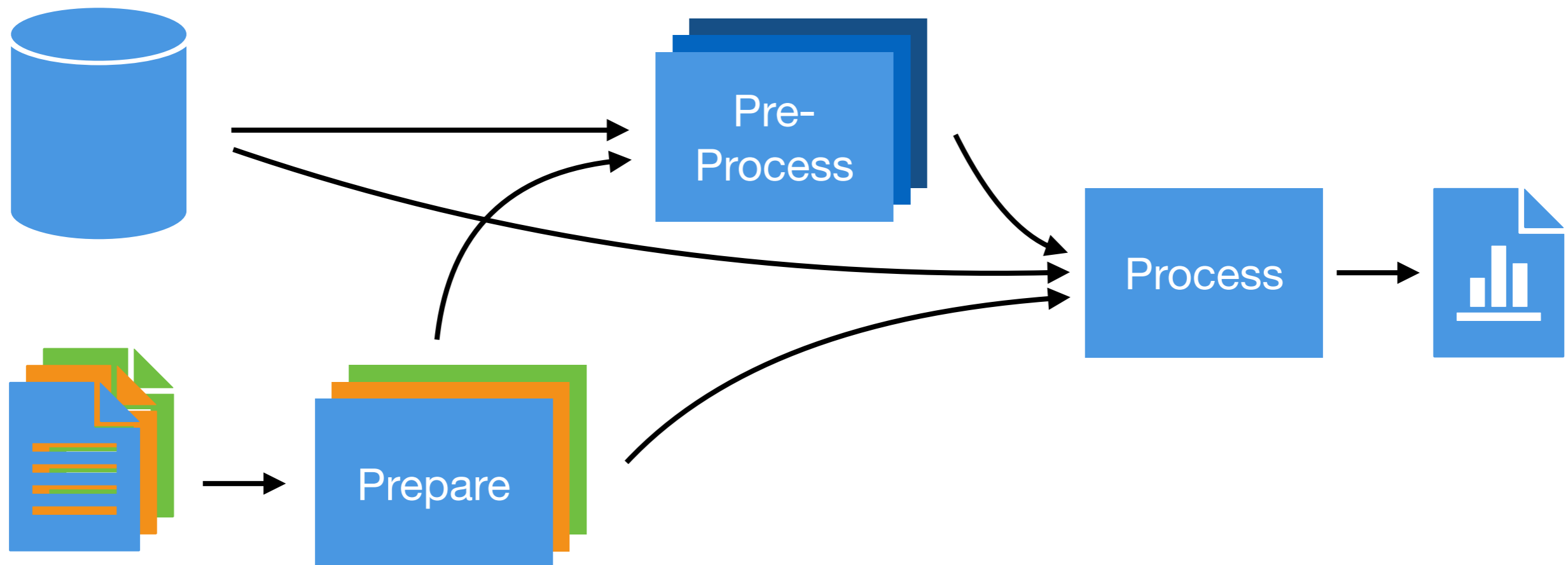
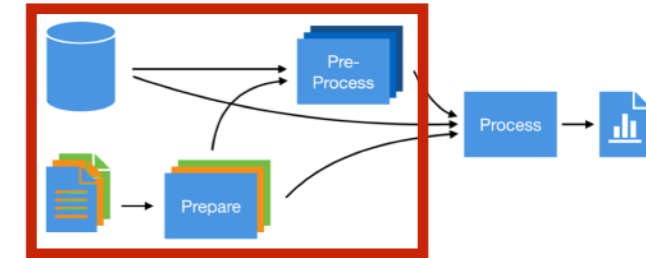
- Full analysis
- Development (repeated)
- Debugging (repeated)
- Plotting features of data



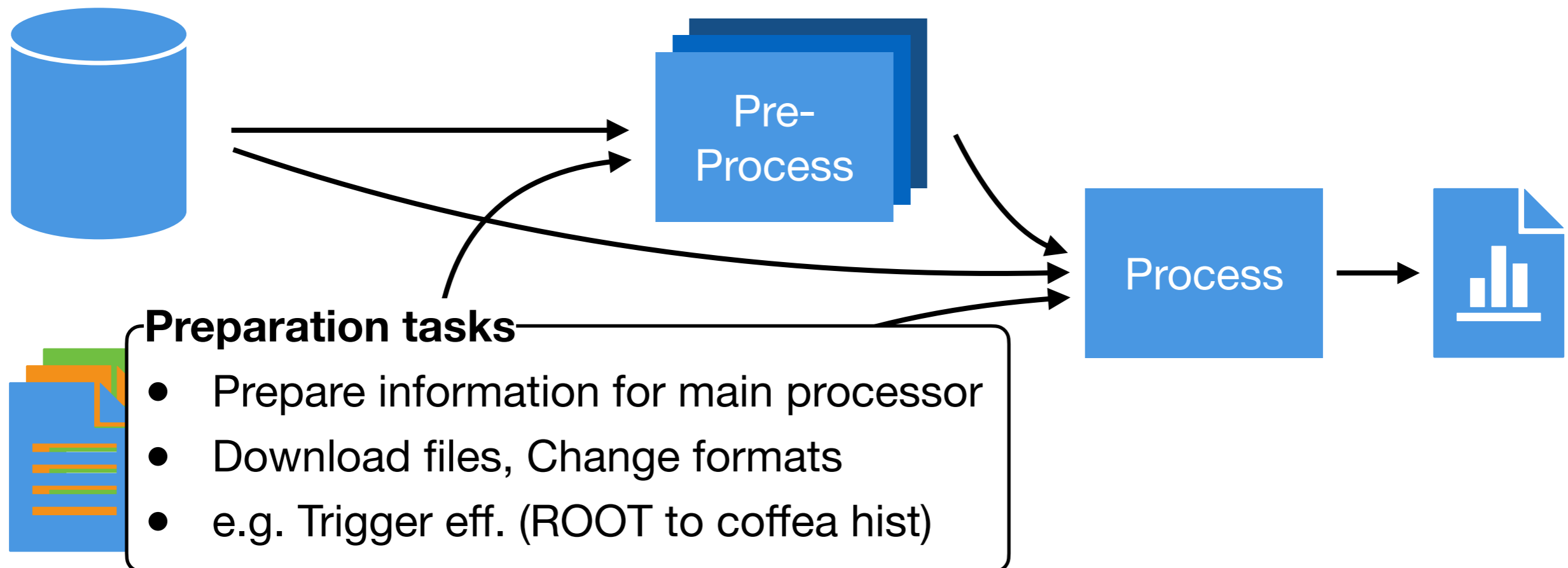
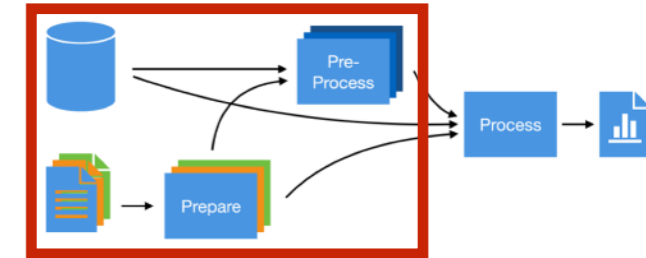


- Make-like execution of whole analysis (`law run FullAnalysis`)
- Visual task graph representation using Luigi Scheduler
  - Used for overview of run status, structural improvements, debugging

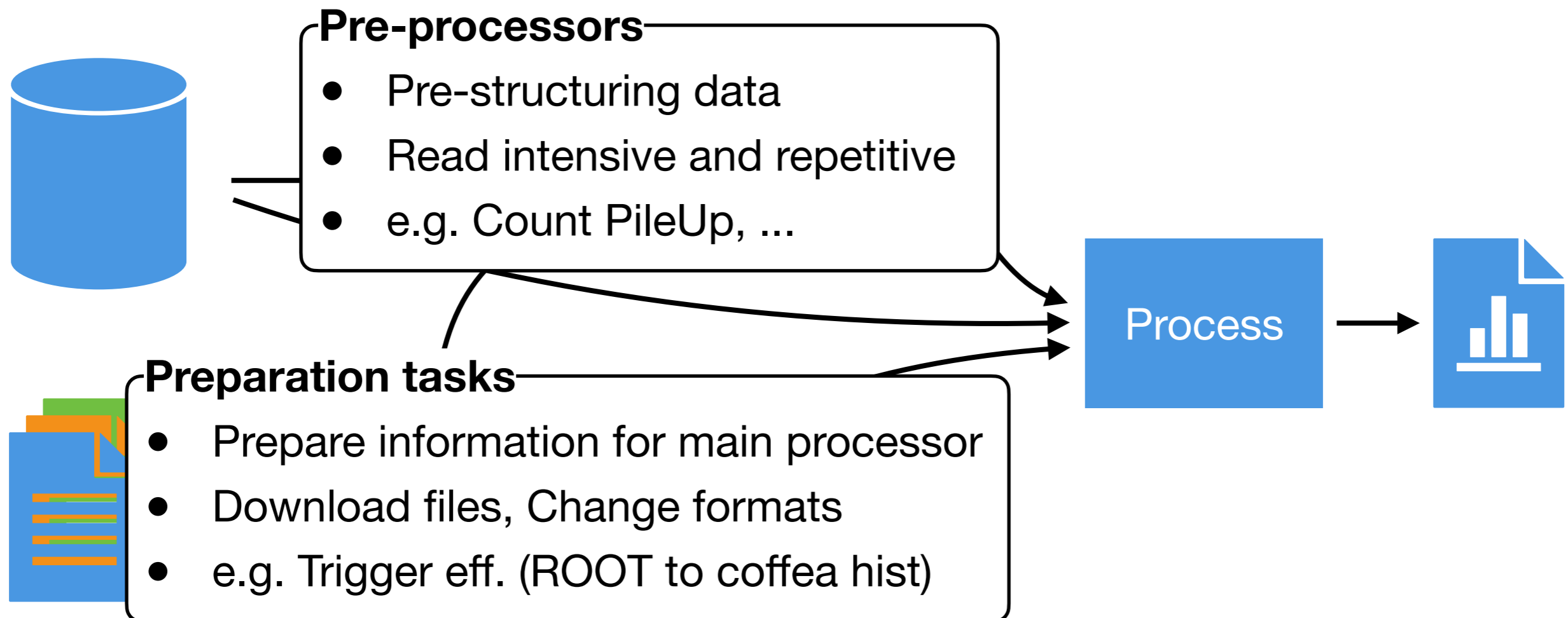
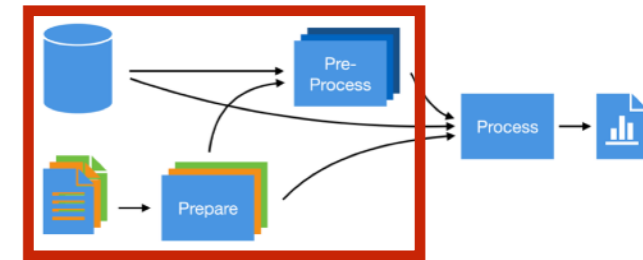
- Different kinds of inputs:
  - Event data  $O(10\text{TB})$ :
    - Recorded data & Simulation
    - NanoAOD format  $\sim 1\text{kB}$  per event
  - Meta data  $O(\text{MB})$ :
    - Efficiency measurements, scaling factors, ...
    - Twiki pages, JSON files, Custom ROOT files, ...



- Different kinds of inputs:
  - Event data  $O(10\text{TB})$ :
    - Recorded data & Simulation
    - NanoAOD format  $\sim 1\text{kB}$  per event
  - Meta data  $O(\text{MB})$ :
    - Efficiency measurements, scaling factors, ...
    - Twiki pages, JSON files, Custom ROOT files, ...

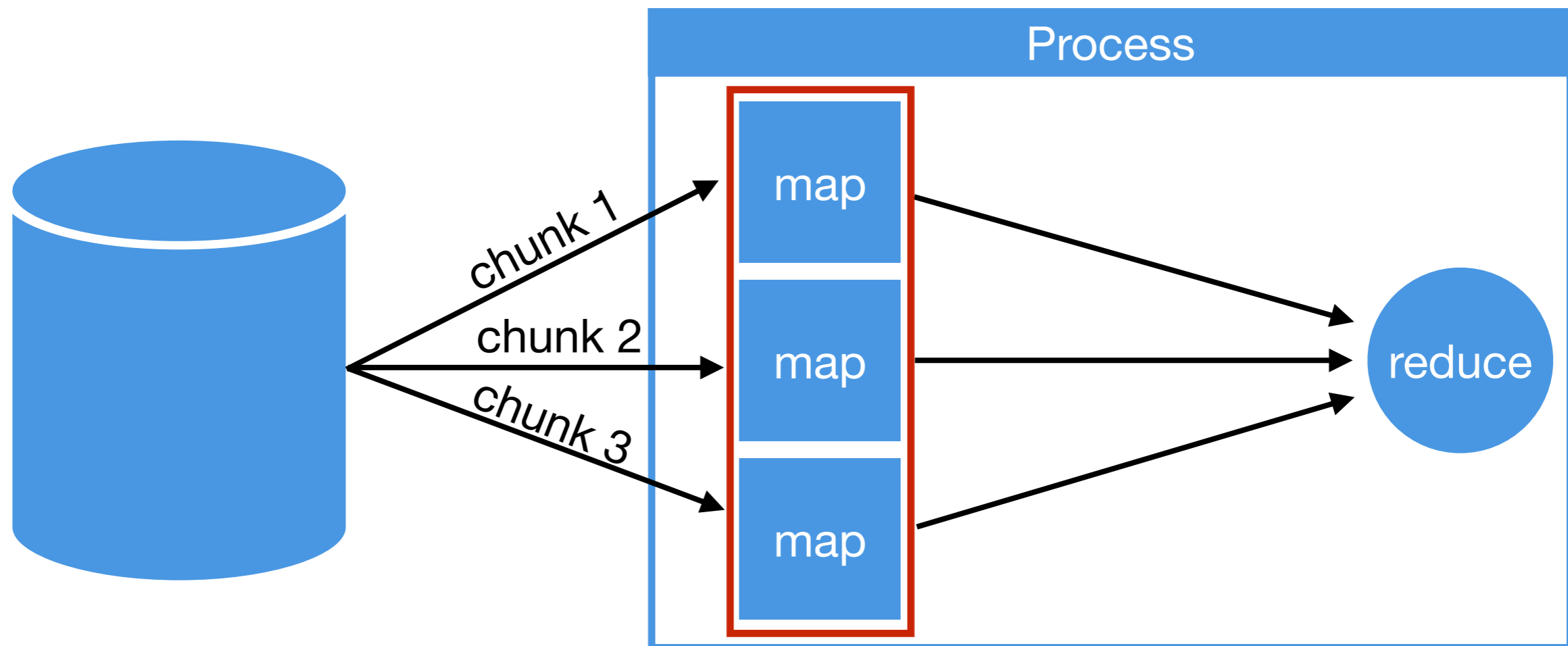
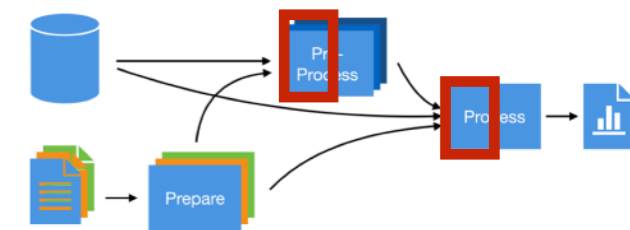


- Different kinds of inputs:
  - Event data  $O(10\text{TB})$ :
    - Recorded data & Simulation
    - NanoAOD format  $\sim 1\text{kB}$  per event
  - Meta data  $O(\text{MB})$ :
    - Efficiency measurements, scaling factors, ...
    - Twiki pages, JSON files, Custom ROOT files, ...

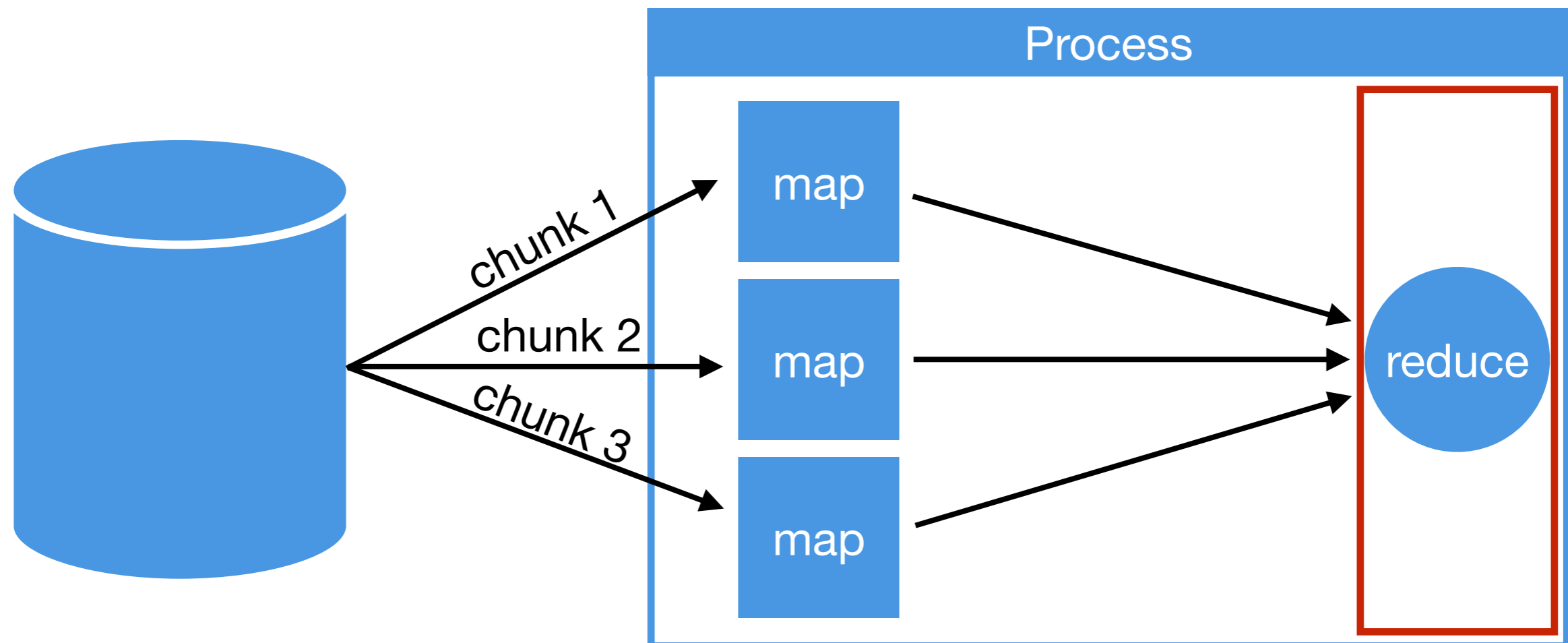
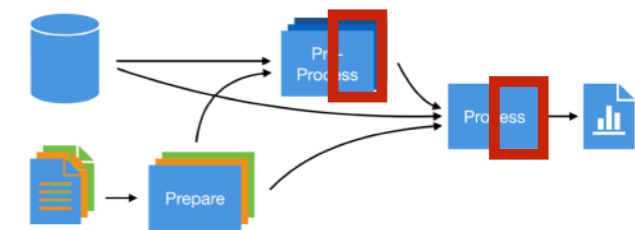




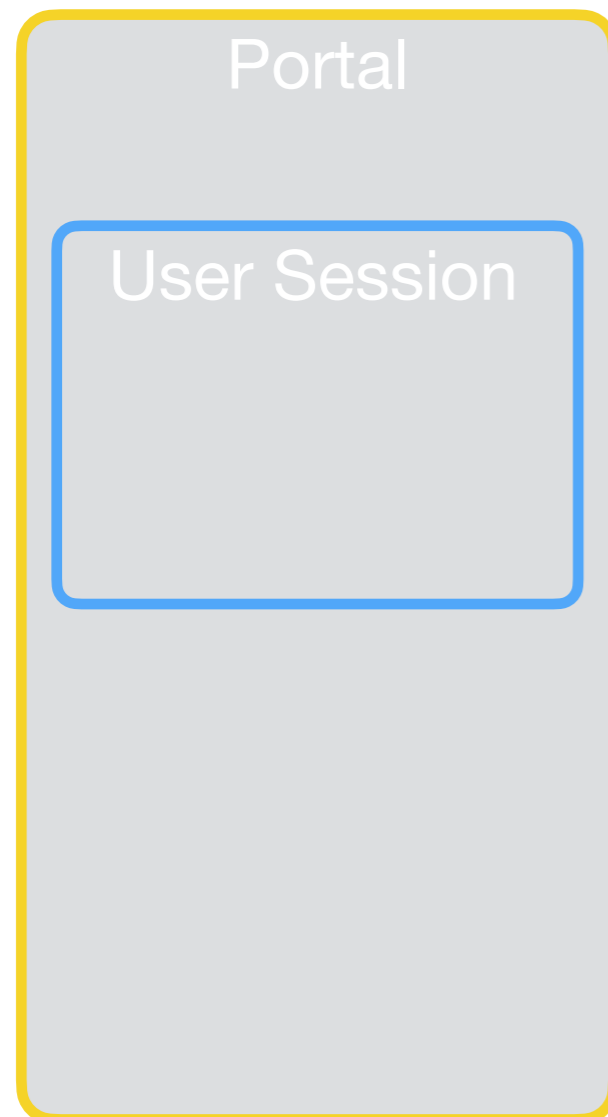
- Processor *maps* perform real computing payload
- Fast:
  - Vectorised processing (awkward, numpy)
  - Central GPU server for MVA evaluation
- Parallel processing in two levels of hierarchy:
  - Datasets (ttbar, ST, DY, ...)
  - Chunks (100k events)



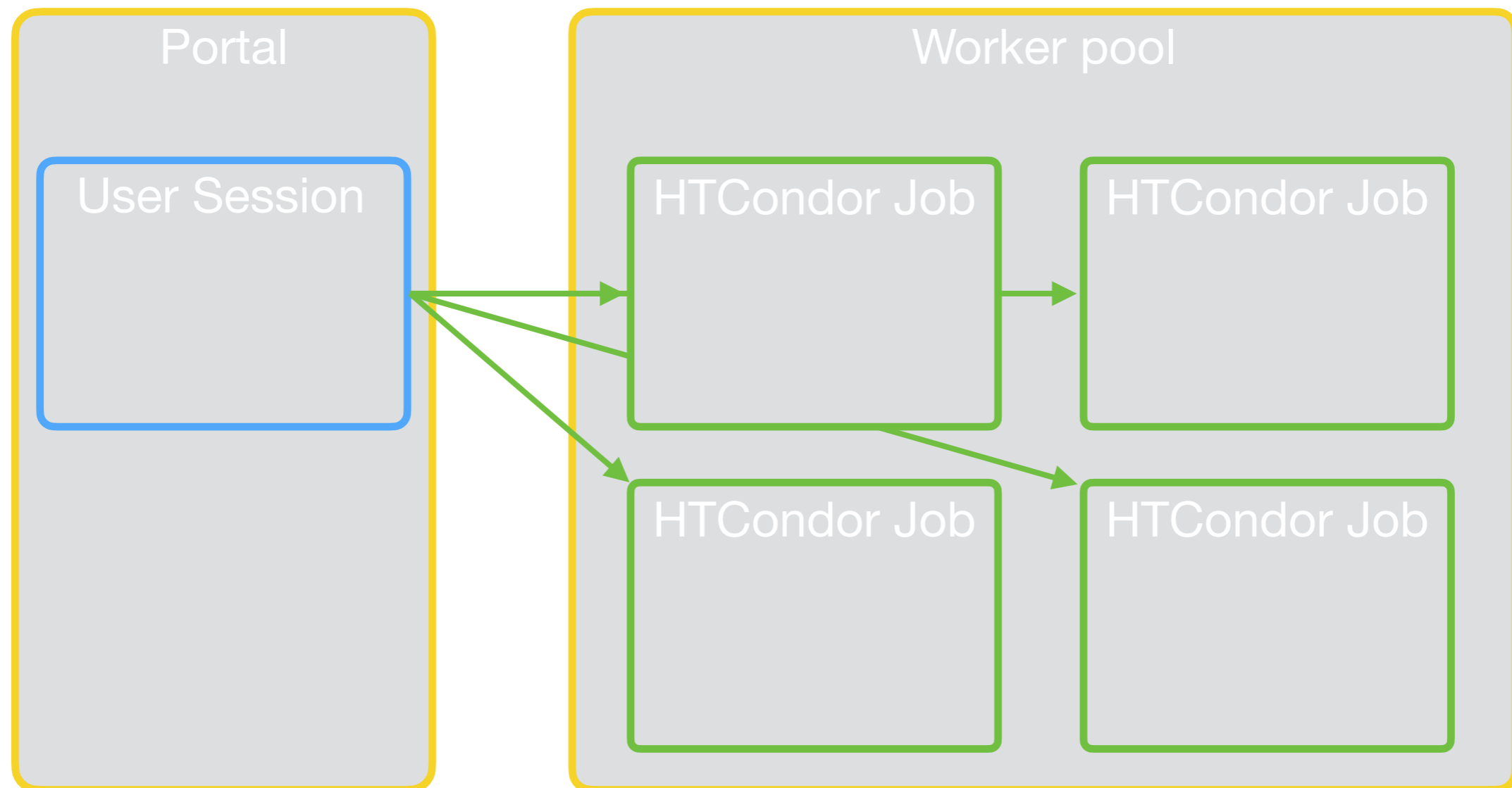
- *Reduces* objects returned by map steps:
  - Histograms, event counts, cutflow statistics, ...
- Histograms are special case:
  - Can be multidimensional and very large O(GB)
  - Implemented fast histogramming methods ([here](#)):
    - hist with h5py backend
    - e.g. Expand, merge, compression, ...



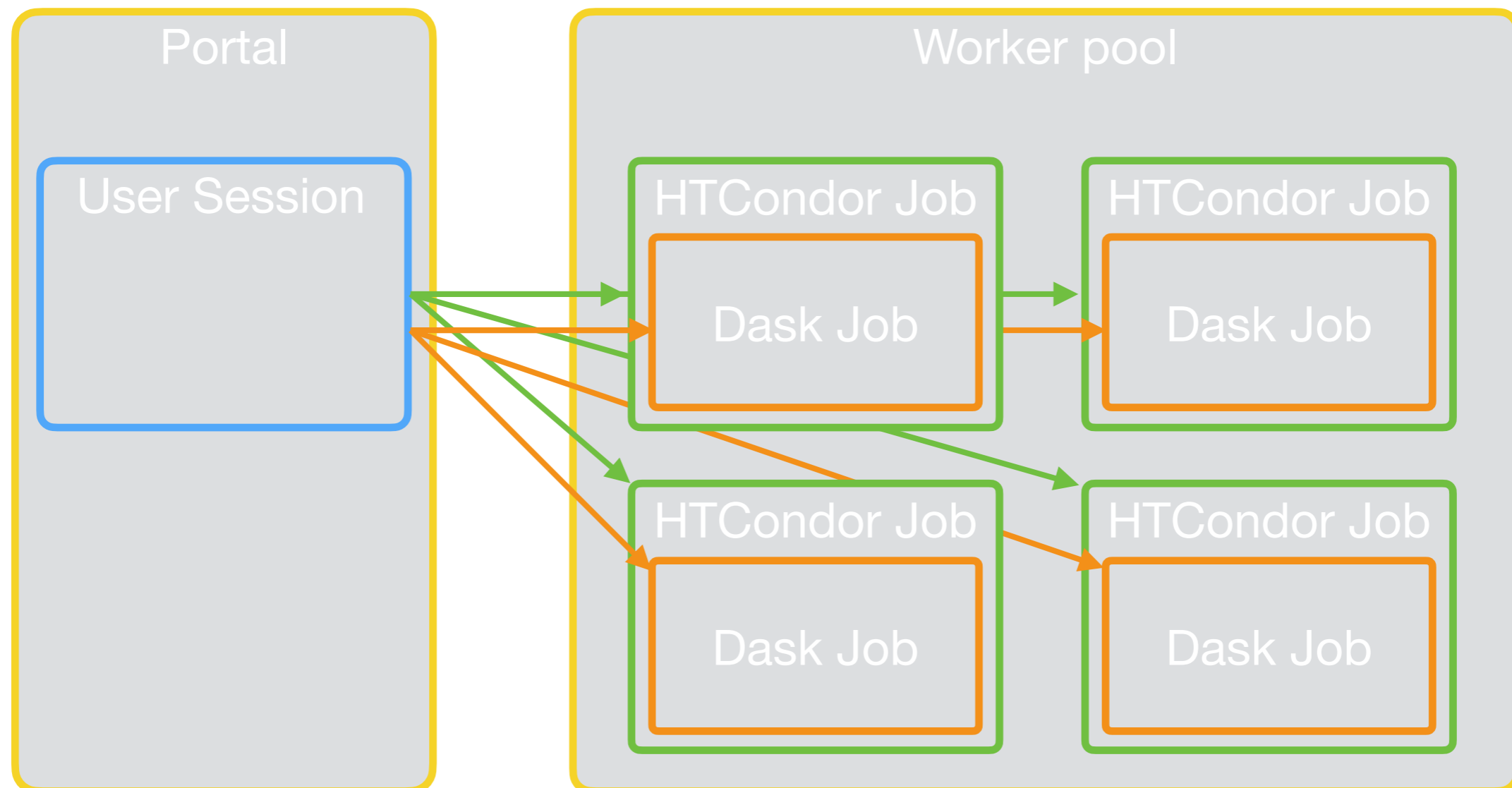
- *Maps* can be parallelised over multiple worker nodes
- Using HPC cluster with one portal node and worker pool
- Two step process:



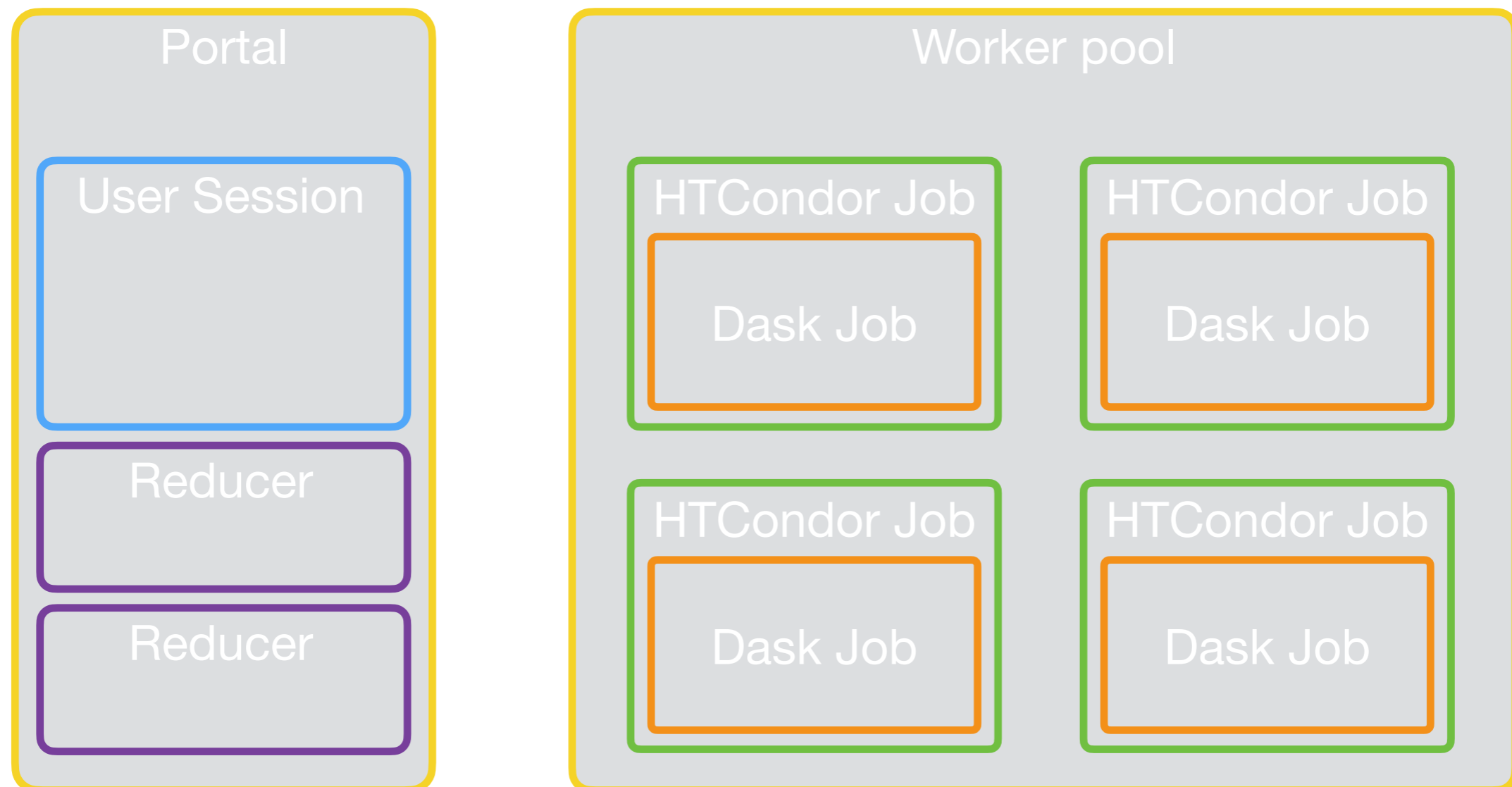
- *Maps* can be parallelised over multiple worker nodes
- Using HPC cluster with one portal node and worker pool
- Two step process:
  1. HTCondor opens jobs in worker pool (1 CPU Thread, 1.5 GB RAM)



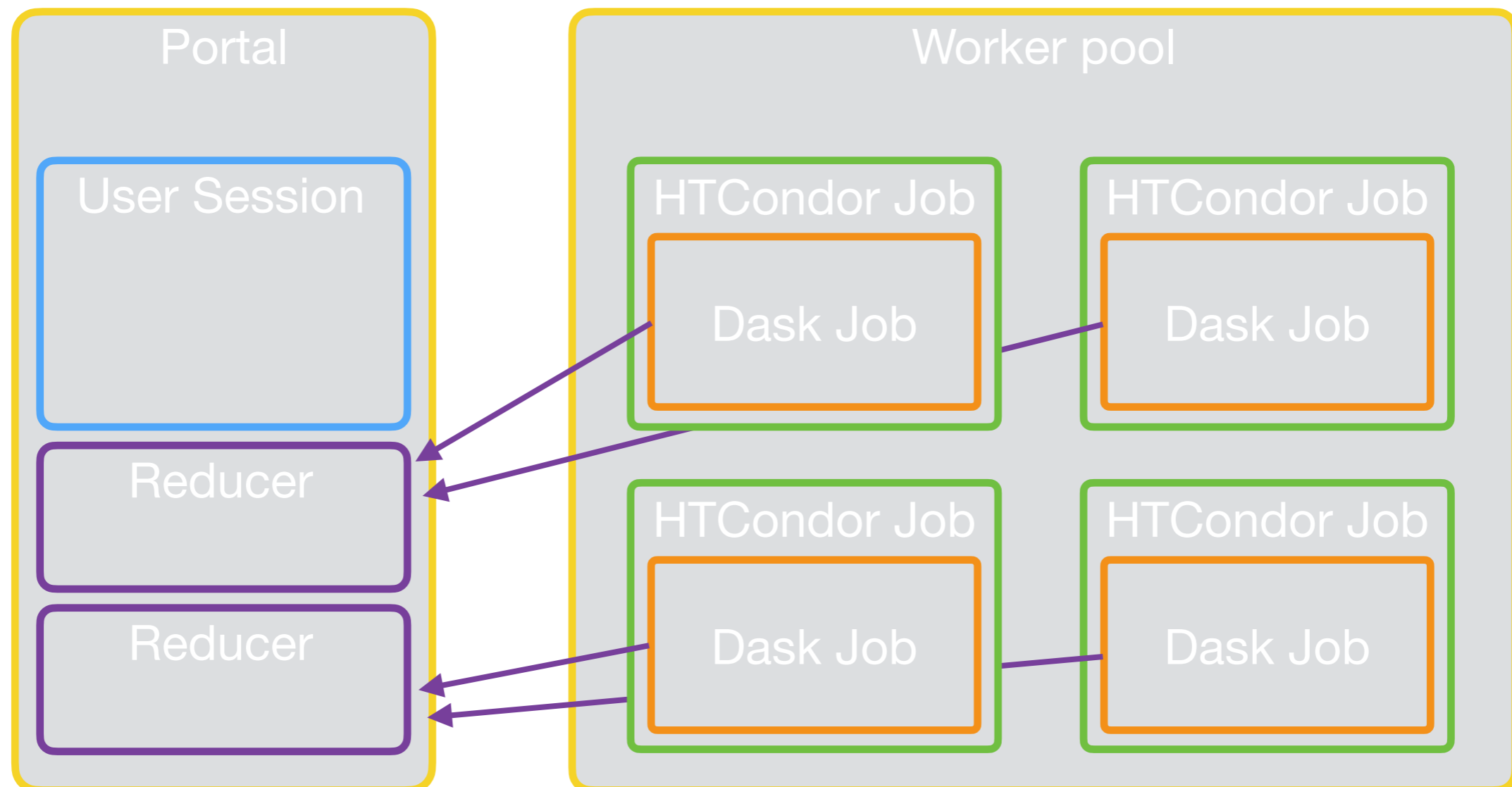
- *Maps* can be parallelised over multiple worker nodes
- Using HPC cluster with one portal node and worker pool
- Two step process:
  1. HTCondor opens jobs in worker pool (1 CPU Thread, 1.5 GB RAM)
  2. Dask uses HTCondor slots



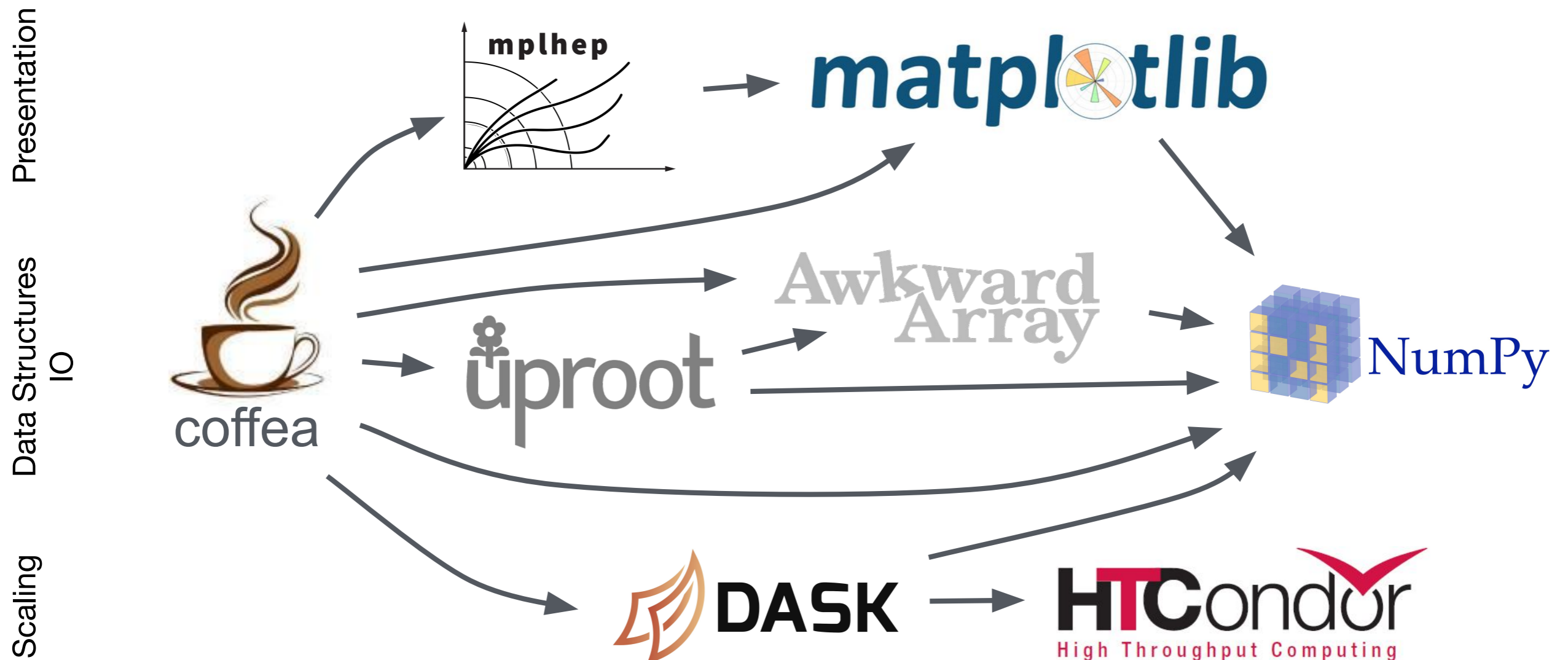
- Outputs of parallel map jobs (e.g. histograms) must be *reduced*
- Ten parallel reducers on portal node
- Each reducer has two reducing instances:
  - Normal reduce: map jobs finished regularly, reduce job pulls
  - Early reduce: map jobs need space, map job pushes



- Outputs of parallel map jobs (e.g. histograms) must be *reduced*
- Ten parallel reducers on portal node
- Each reducer has two reducing instances:
  - Normal reduce: map jobs finished regularly, reduce job pulls
  - Early reduce: map jobs need space, map job pushes

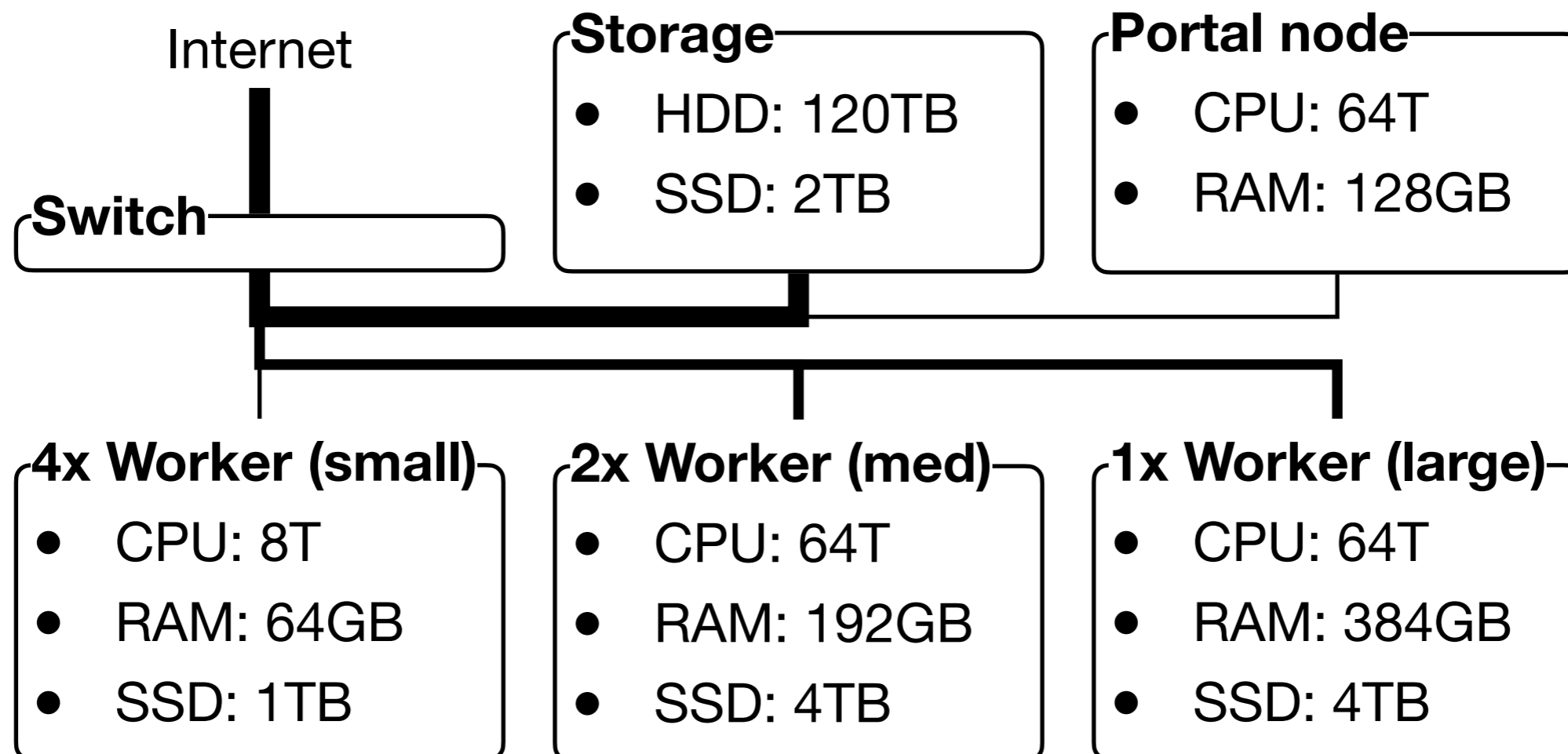


- Everything implemented in **python** (analysis repository)
- Modular software framework:
  - Standard tools: NumPy, matplotlib, Dask, HTCondor, ...
  - HEP specific tools: AwkwardArray, uproot, mplhep, coffea, ...



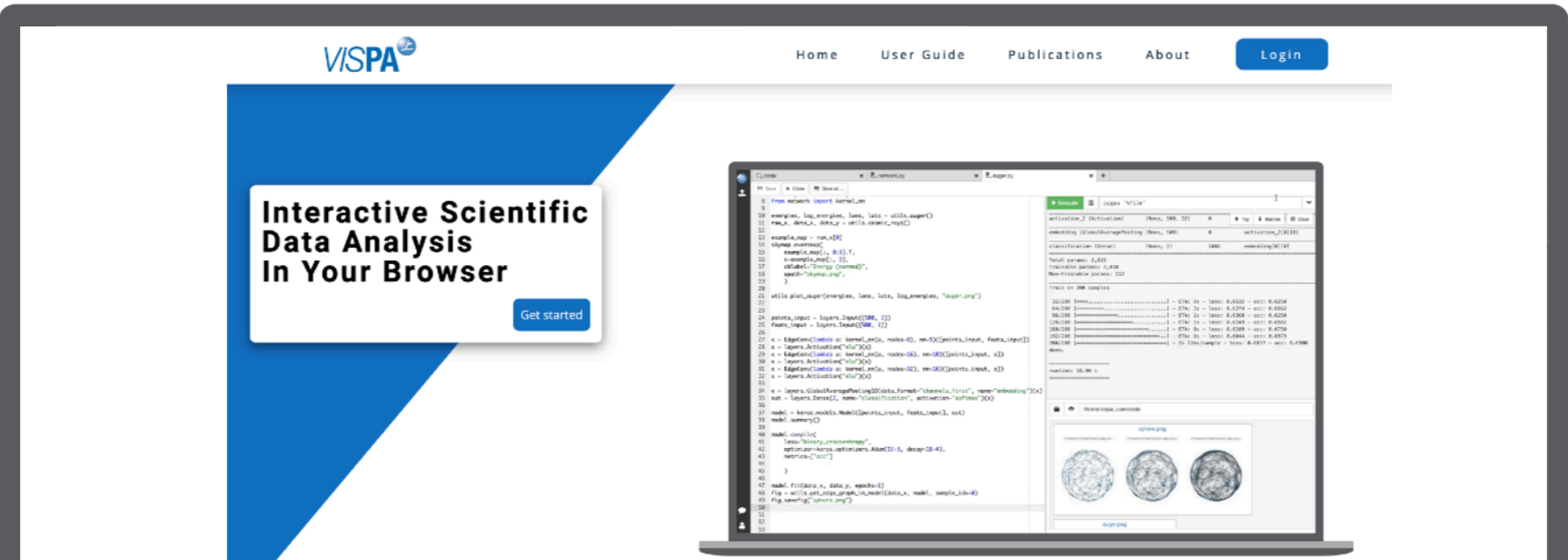


- All analyses run on local institute cluster
- Optimised for scientific data analysis and machine learning
- Hardware: see below
- Software:
  - System (Ubuntu OS) via ansible
  - Analysis via conda, shared over network



Total: CPU: 224T, RAM: 832GB, GPU: 17 (various)


- User base:
  - 10 Researchers on daily basis (e.g. CMS experiment, Auger observatory)
  - Courses with up to 200 participants (e.g. Nuclear physics, ML in Physics)
  - Schools and workshop with up to 50 participants
- Front-end for data analysis in your web browser ([link](#))



The screenshot shows the VISPA website interface. On the left, there's a blue banner with the text "Interactive Scientific Data Analysis In Your Browser" and a "Get started" button. The main content area features a code editor with Python code for data analysis, a terminal window showing execution output, and a visualization of three circular plots. The website header includes navigation links: Home, User Guide, Publications, About, and a Login button.


## WHY VISPA

The Visual Physics Analysis (VISPA) project offers a flexible desktop-like development environment using only your web browser.




**Access From Anywhere**

VISPA runs in your web browser with no setup required, enabling you to develop and execute your data analysis whenever and wherever you want.



**Desktop-like UI**

VISPA offers a full desktop-like environment with standard tools such as a code editor and file browser, with specialized extensions to improve your analysis workflow.

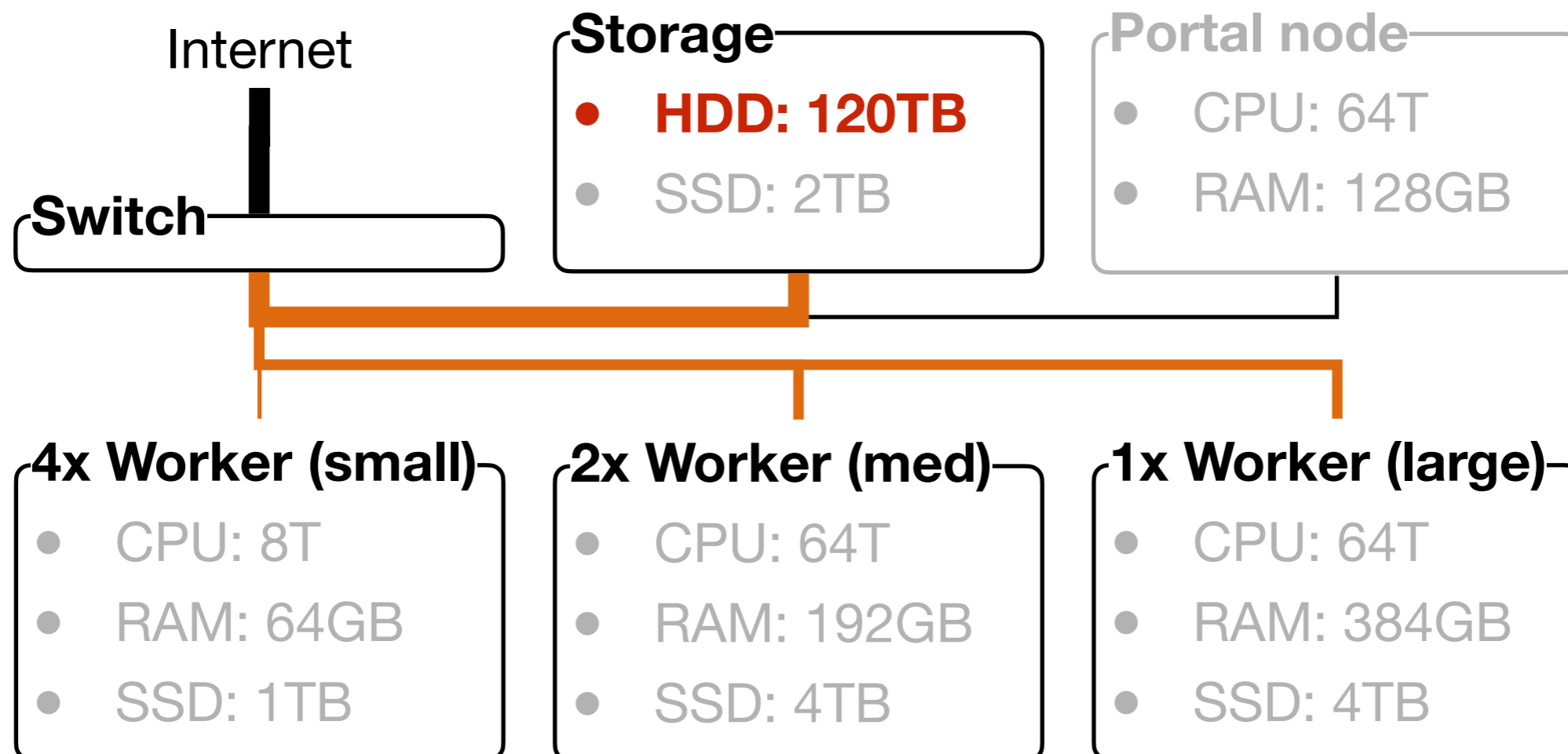


**Execute Anywhere**

Get directly started on our cluster using pre-installed libraries, create your own environment, or use any resource available to you.

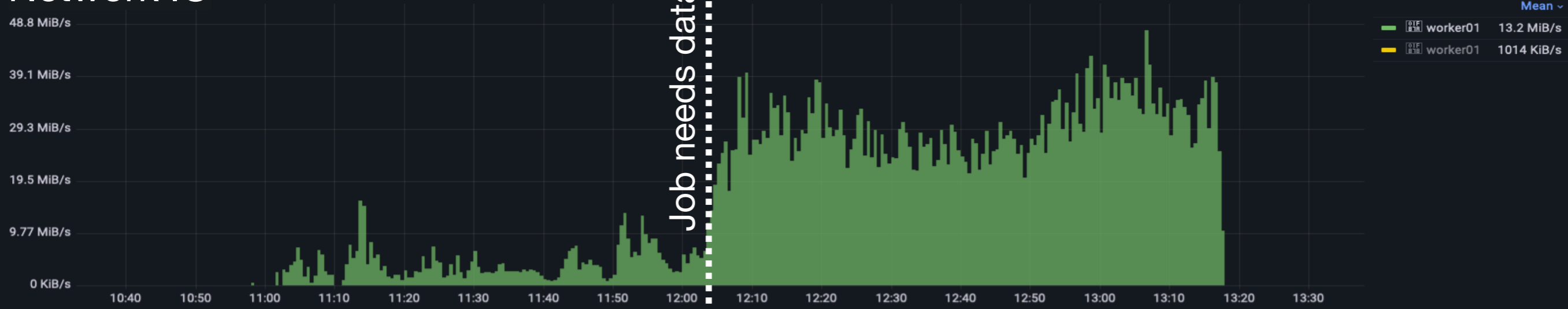
Live Demo

- Data is streamed from central network storage to worker nodes
- Map jobs are very time efficient and fast using vectorised processing
- Central network storage has two limitations:
  - **Read speed of HDDs**
  - **Limited network bandwidth**
- Streaming of data becomes critical bottleneck!



# 16 A New Bottleneck (2/2)

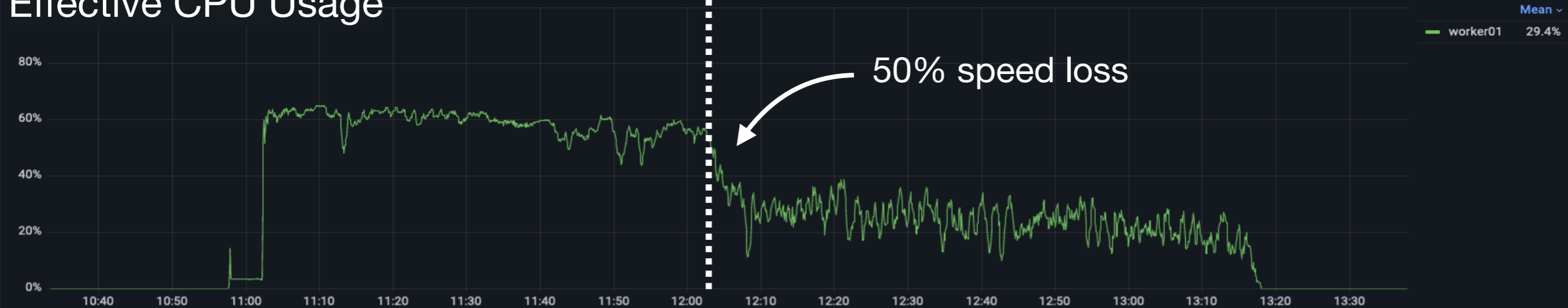
## Network IO



## CPU Wait



## Effective CPU Usage



- Critical bottleneck is streaming of data to the processing elements
- Various storage types and locations available to solve the problem
- Two stage solution using central network storage and on-worker SSD

## WLCG

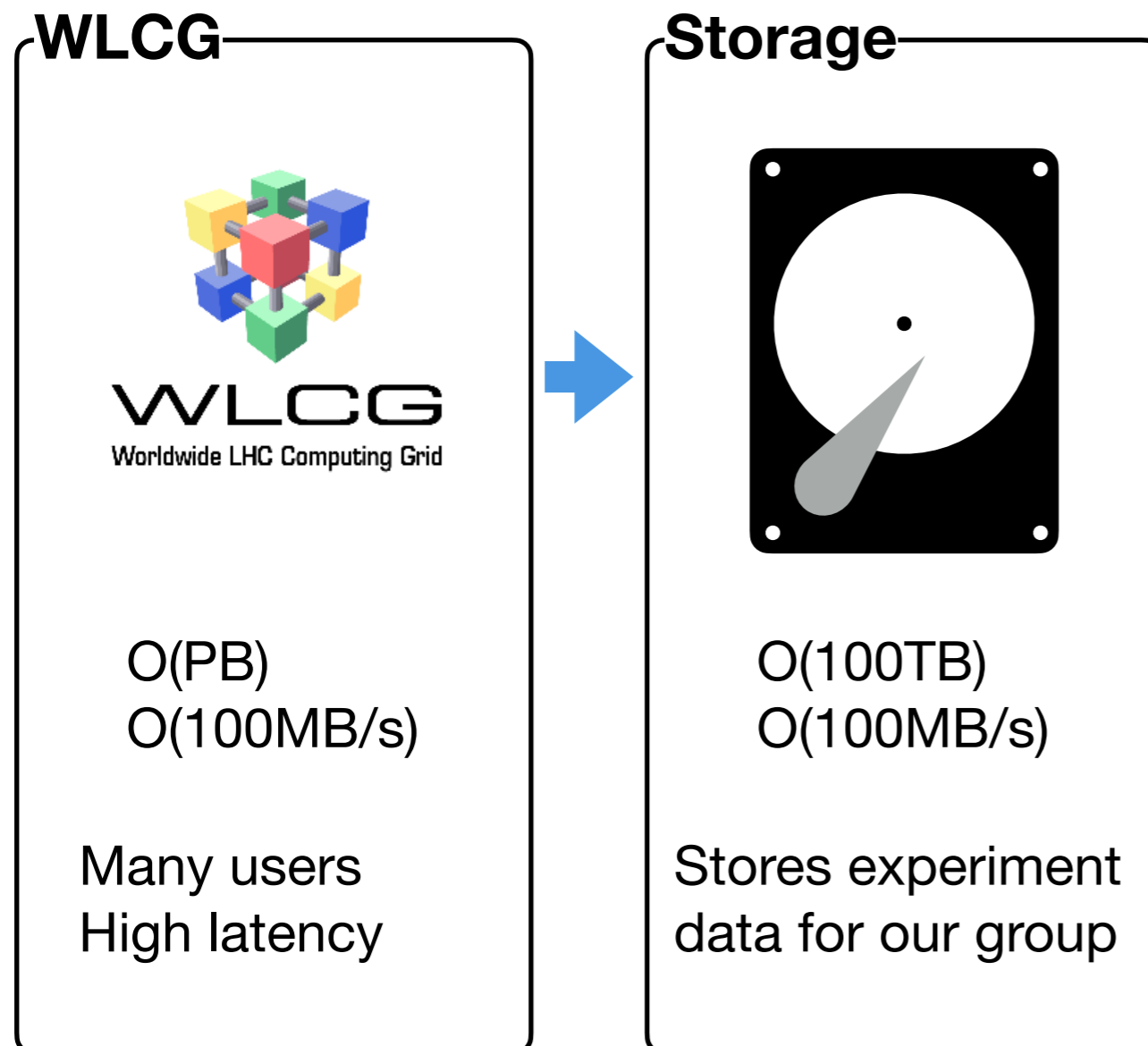


**WLCG**  
Worldwide LHC Computing Grid

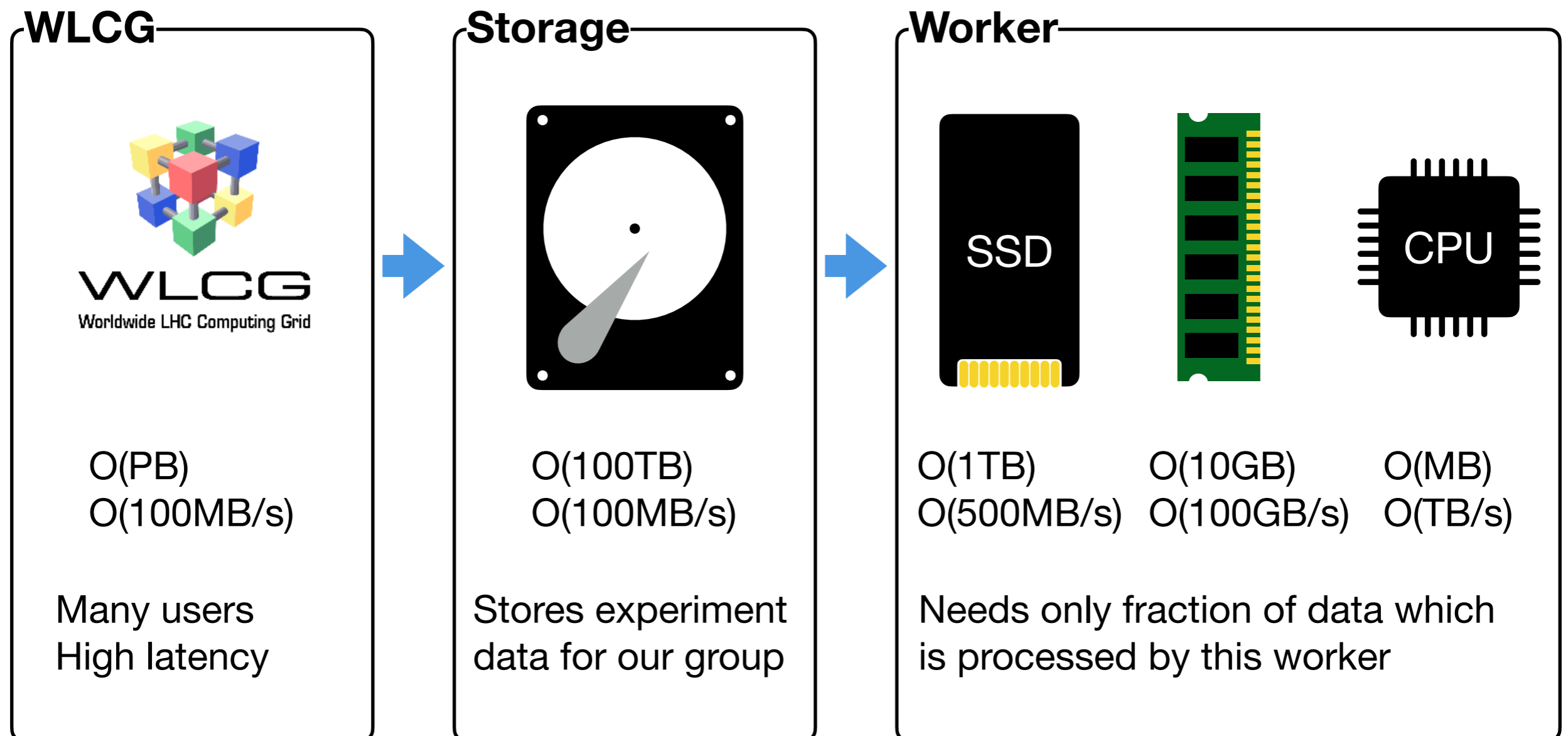
O(PB)  
O(100MB/s)

Many users  
High latency

- Critical bottleneck is streaming of data to the processing elements
- Various storage types and locations available to solve the problem
- Two stage solution using central network storage and on-worker SSD



- Critical bottleneck is streaming of data to the processing elements
- Various storage types and locations available to solve the problem
- Two stage solution using central network storage and on-worker SSD

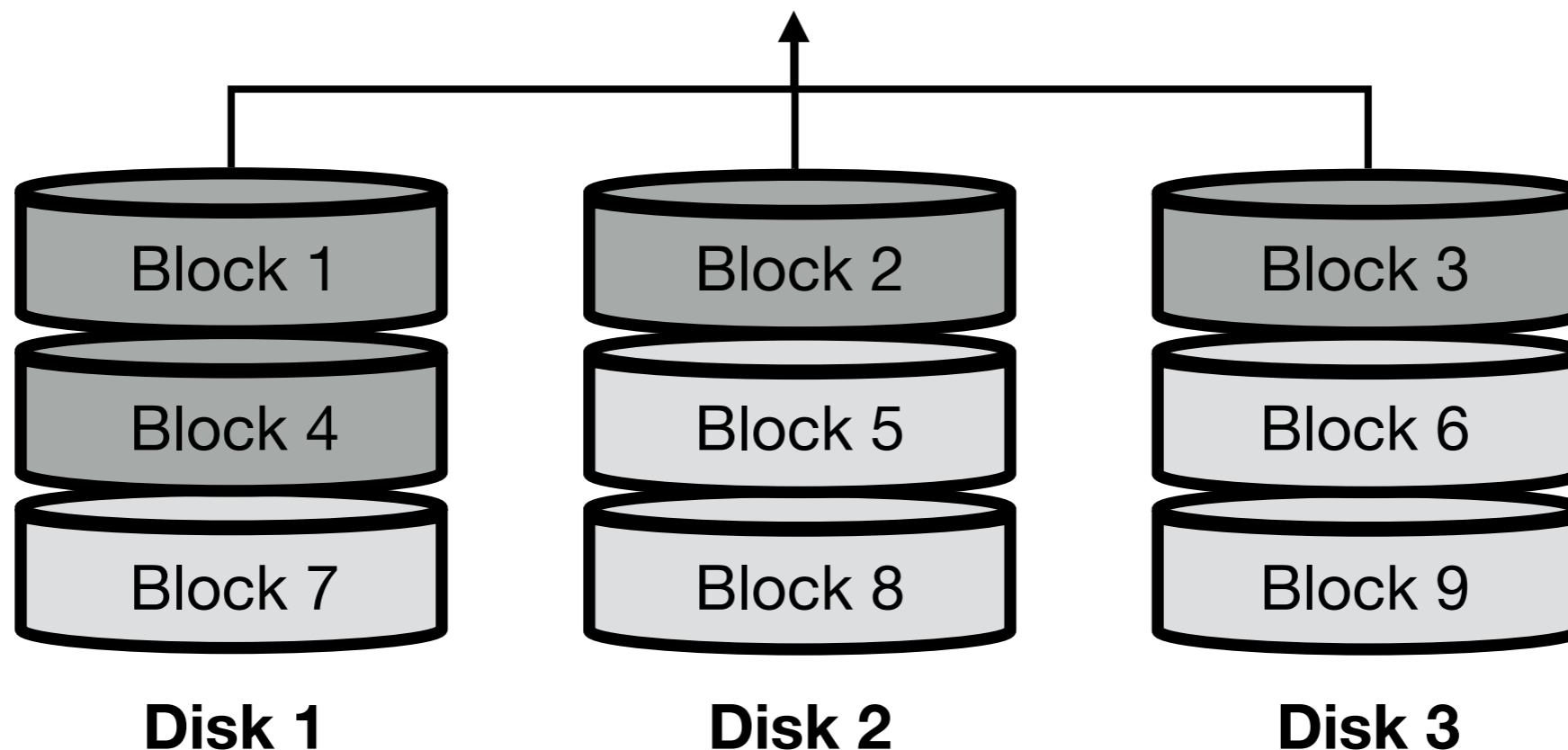
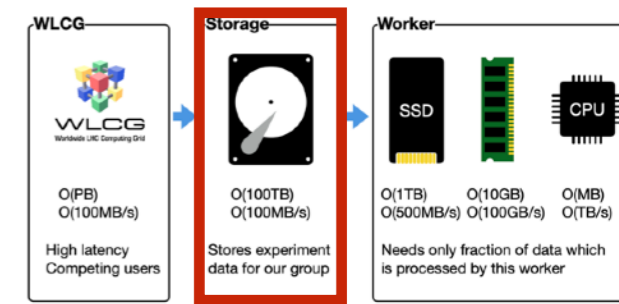


Slide 18

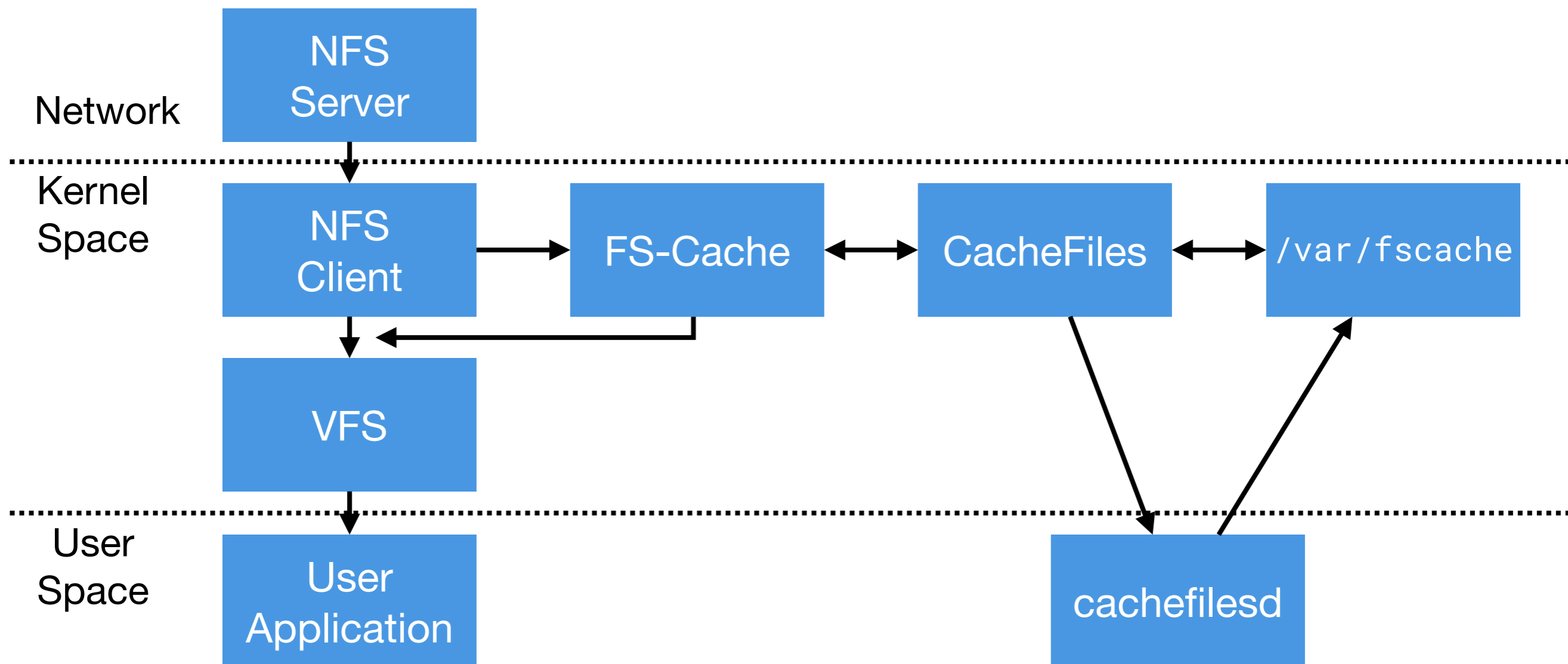
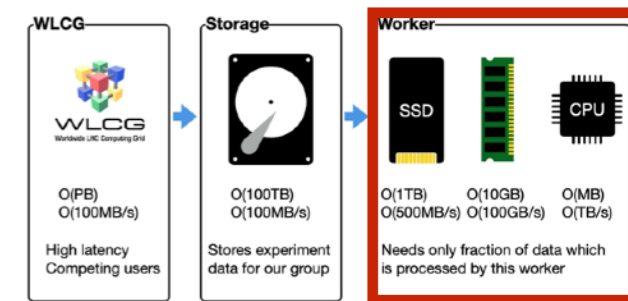
Slide 19-21



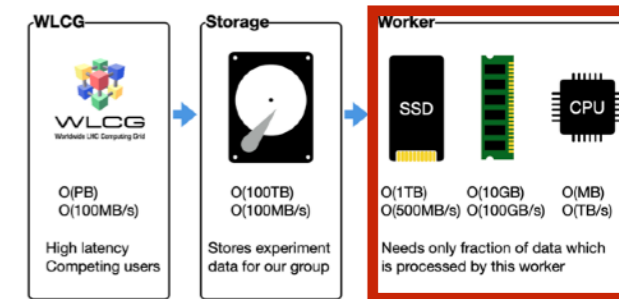
- Save experiment data on-site:
  - Easy and reliable access (no timeouts, credentials, ...)
  - Direct connection to worker nodes (low latency, 10GBit)
- Three storage tiers:
  - /home (2TB): User homes
  - /scratch (24TB): Mirrored experiment data
  - /store (96TB):
    - Un-mirrored experiment data
    - RAID0 striped across 6 x 16TB HDDs enables fast reading



- Using on-worker SSDs to minimise network traffic
- Used software implementation: FS-Cache & cachefilesd
  - Transparent caching system, available in Linux kernel
  - Granularity: Cache on block level (4kB)
  - Strategy: Least recently used (LRU)



- Using on-worker SSDs to minimise network traffic
- Used software implementation: FS-Cache & cachefilesd
  - Transparent caching system, available in Linux kernel
  - Granularity: Cache on block level (4kB)
  - Strategy: Least recently used (LRU)



## Challenge: Cache trashing!

- Needed data gets evicted from cache
- Case 1: By somebody else
  - Everybody uses same files/chunks (where possible)
- Case 2: By yourself
  - Consistent assignment job ↔ worker via affine caching

Network

Kernel  
Space

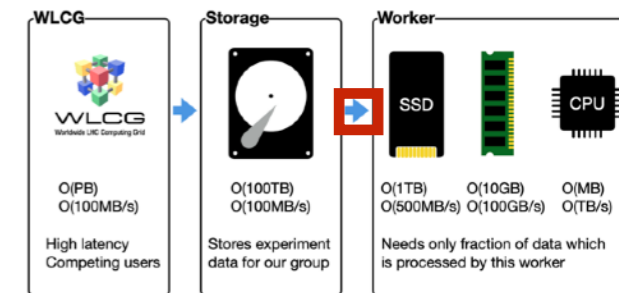
User  
Space

Application

cachefilesd

scache

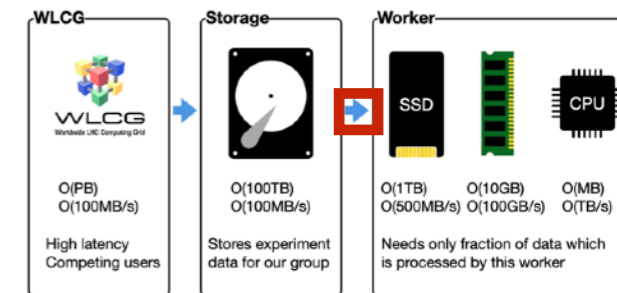
- Always process same data on same worker
- Use 64-dim embedding in hash space:
  - Embedding:



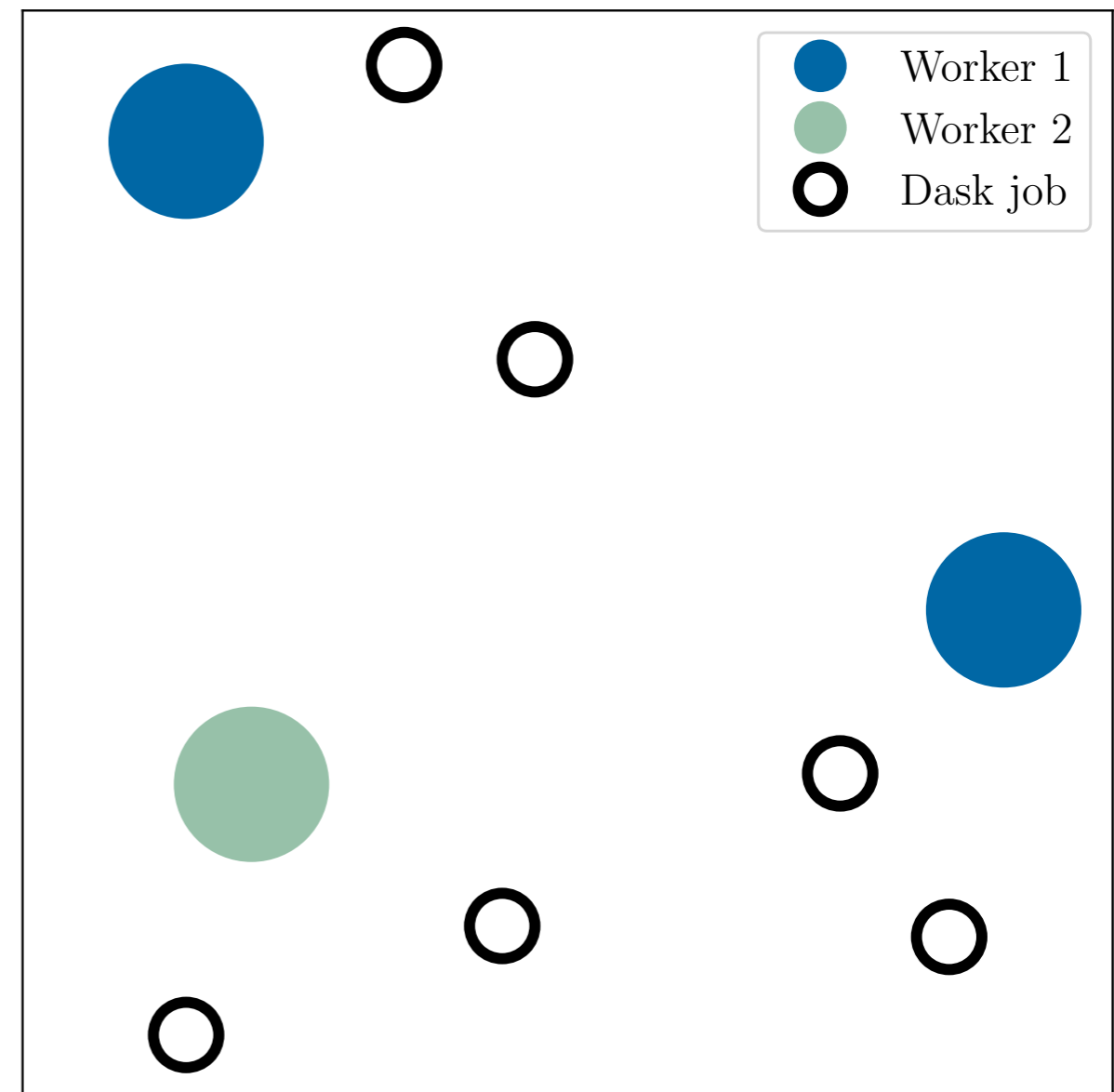
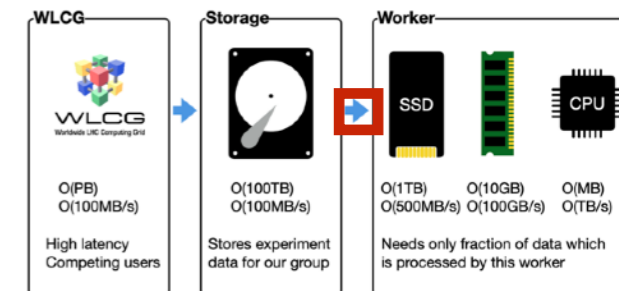
"worker-01"  $\xrightarrow{\text{sha512}}$  10110100 11001001 ...

int ↓  $x_1=180$       int ↓  $x_2=201$

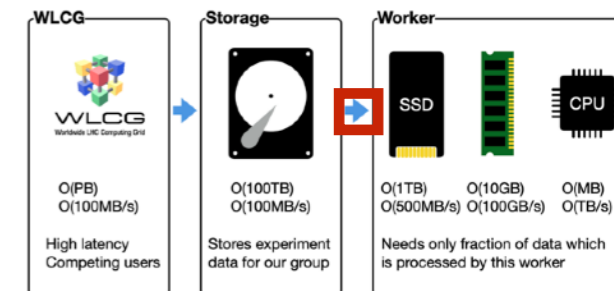
- Always process same data on same worker
- Use 64-dim embedding in hash space:
  - Embedding:  $\text{embed}(\text{"worker-01"}) = (180, 201, \dots)$



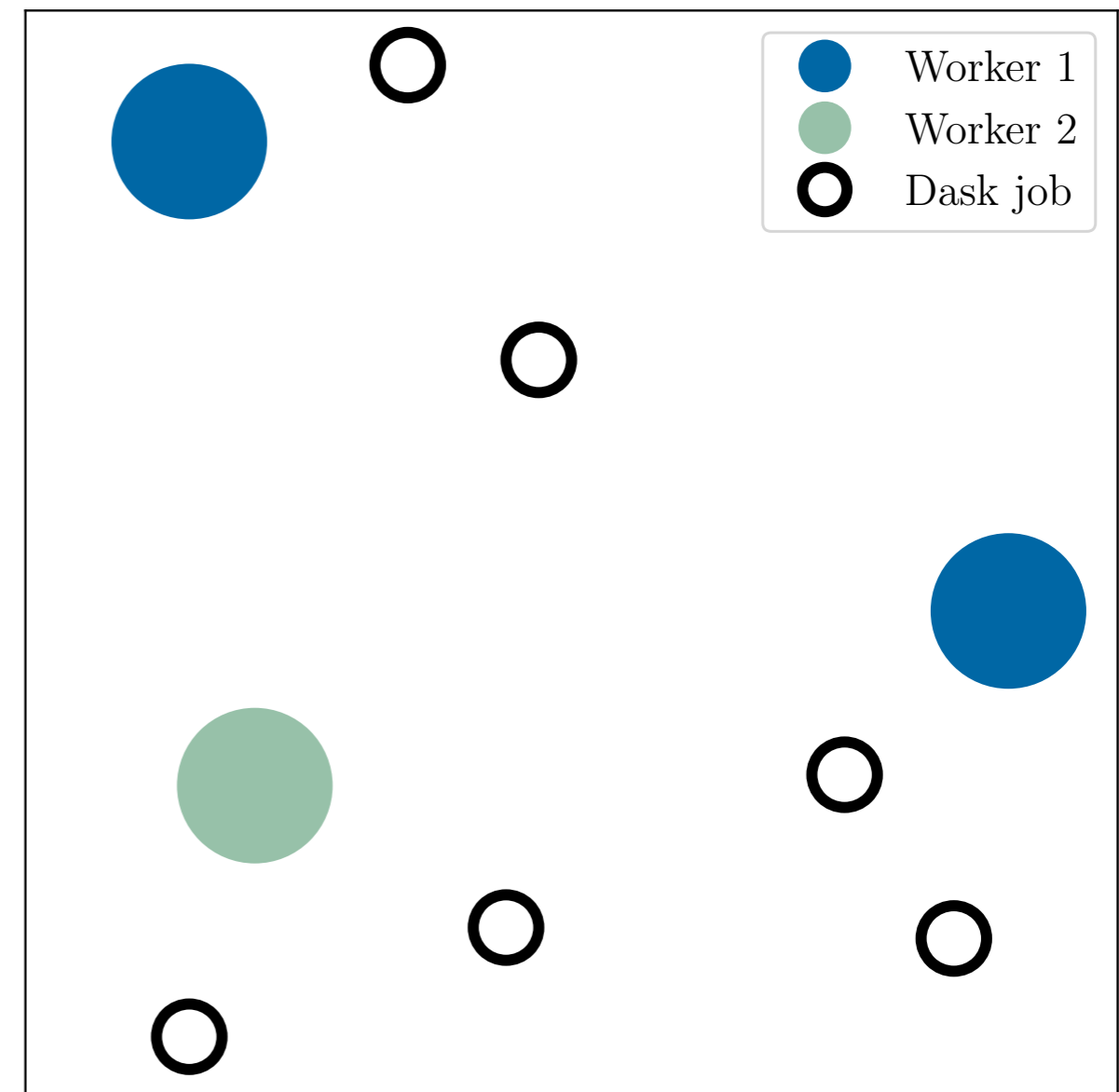
- Always process same data on same worker
- Use 64-dim embedding in hash space:
  - Embedding:  $\text{embed}(\text{"worker-01"}) = (180, 201, \dots)$ 
    - Worker:  $\text{embed}(\text{worker\_name} + \text{id})$
    - Job/data:  $\text{embed}(\text{filename} + \text{event range})$



- Always process same data on same worker
- Use 64-dim embedding in hash space:
  - Embedding:  $\text{embed}(\text{"worker-01"}) = (180, 201, \dots)$ 
    - Worker:  $\text{embed}(\text{worker\_name} + \text{id})$
    - Job/data:  $\text{embed}(\text{filename} + \text{event range})$

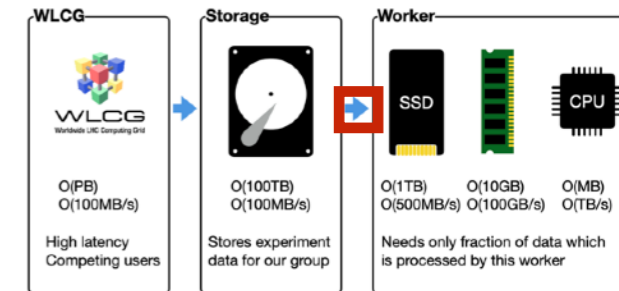


- Two tricks:
  - Worker have weight acc. to size
  - Place worker multiple times (id)

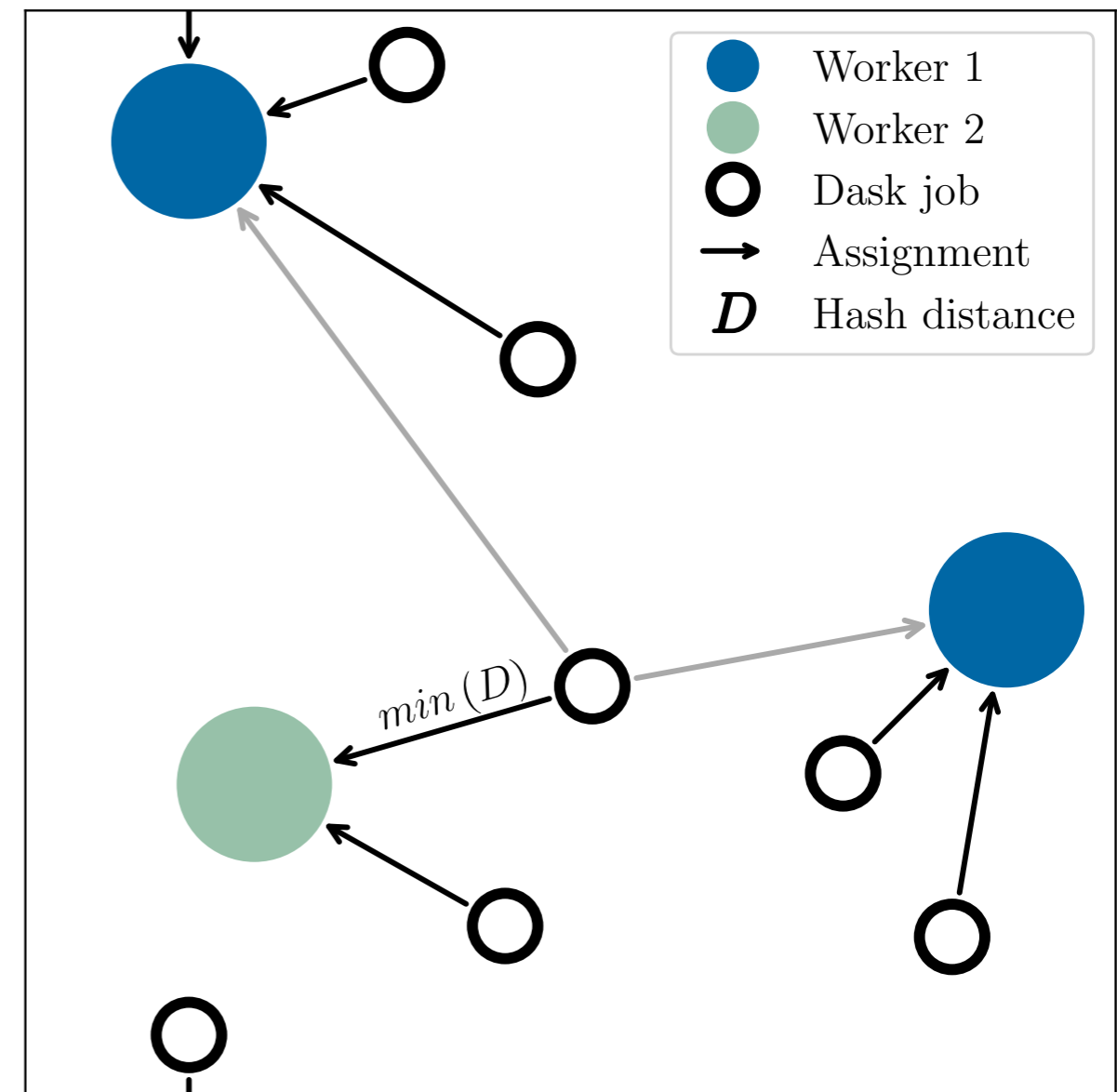


# 21 Assignment Job ↔ Worker

- Always process same data on same worker
- Use 64-dim embedding in hash space:
  - Embedding:  $\text{embed}(\text{"worker-01"}) = (180, 201, \dots)$ 
    - Worker:  $\text{embed}(\text{worker\_name} + \text{id})$
    - Job/data:  $\text{embed}(\text{filename} + \text{event range})$



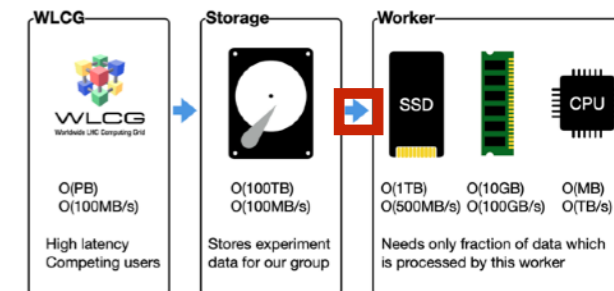
- Two tricks:
  - Worker have weight acc. to size
  - Place worker multiple times (id)
- $D = \text{Weighted euclidean distance}$
- Assignment via minimum  $D$


 Logic similar to ceph ([link](#))



# 21 Assignment Job ↔ Worker

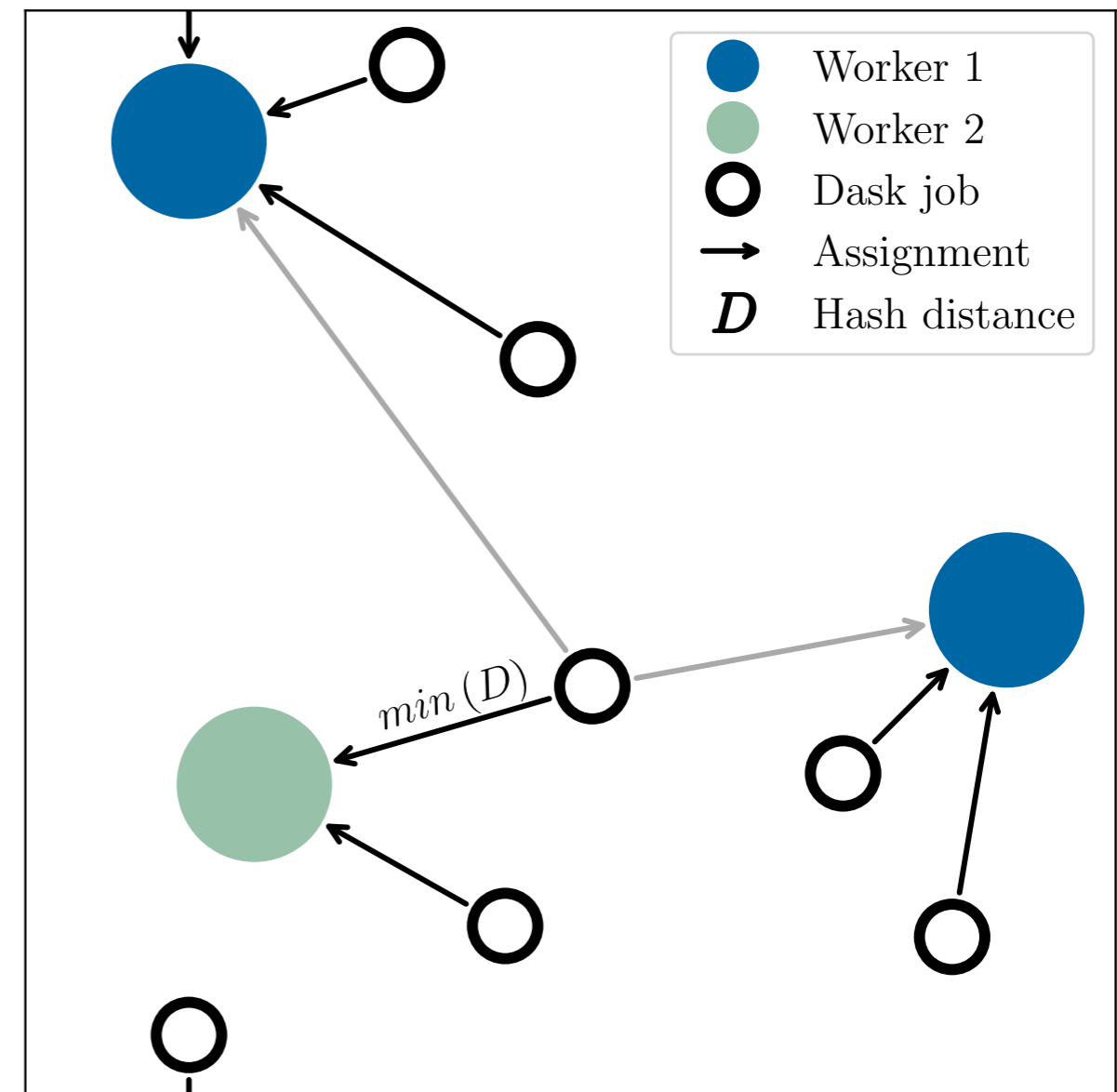
- Always process same data on same worker
- Use 64-dim embedding in hash space:
  - Embedding:  $\text{embed}(\text{"worker-01"}) = (180, 201, \dots)$ 
    - Worker:  $\text{embed}(\text{worker\_name} + \text{id})$
    - Job/data:  $\text{embed}(\text{filename} + \text{event range})$



- Two tricks:
  - Worker have weight acc. to size
  - Place worker multiple times (id)

- $D =$  Weighted euclidean distance
- Assignment via minimum  $D$

- Results in:
  - Affine caching
  - Graceful on failures



## XCache



### XRootD

- No POSIX access
- Not easy with many users

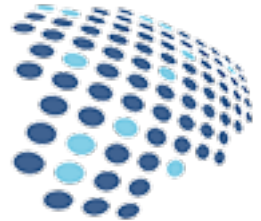
## Distributed File Systems



### ceph

- Does not handle changing files so well

### XTREEMFS



- No cache tiering



### GlusterFS

- Cache tiering removed

### l.u.s.t.r.e. File System

- No graceful failure of parts of storage

### MooseFS

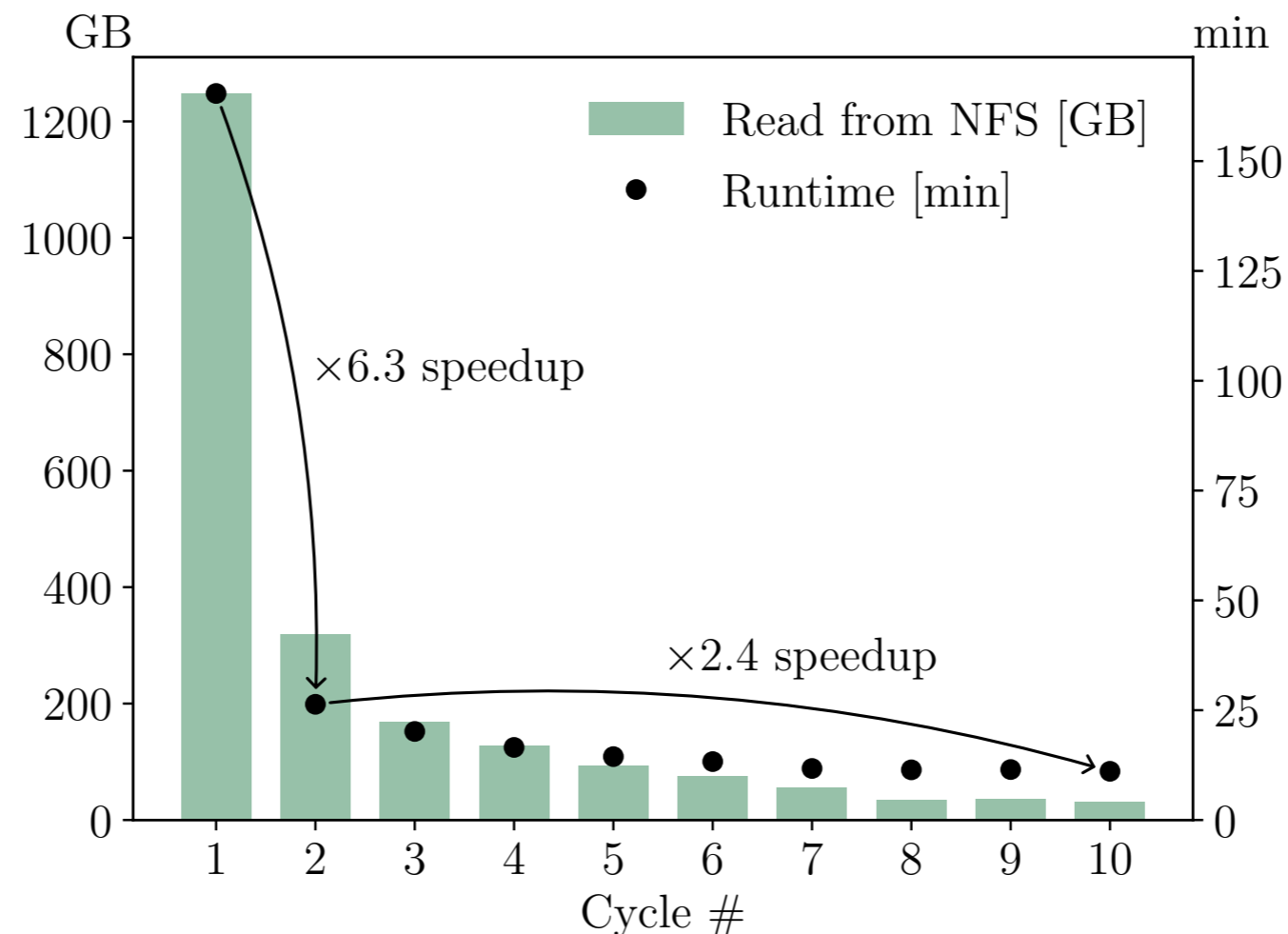


- Caching only based on modification time

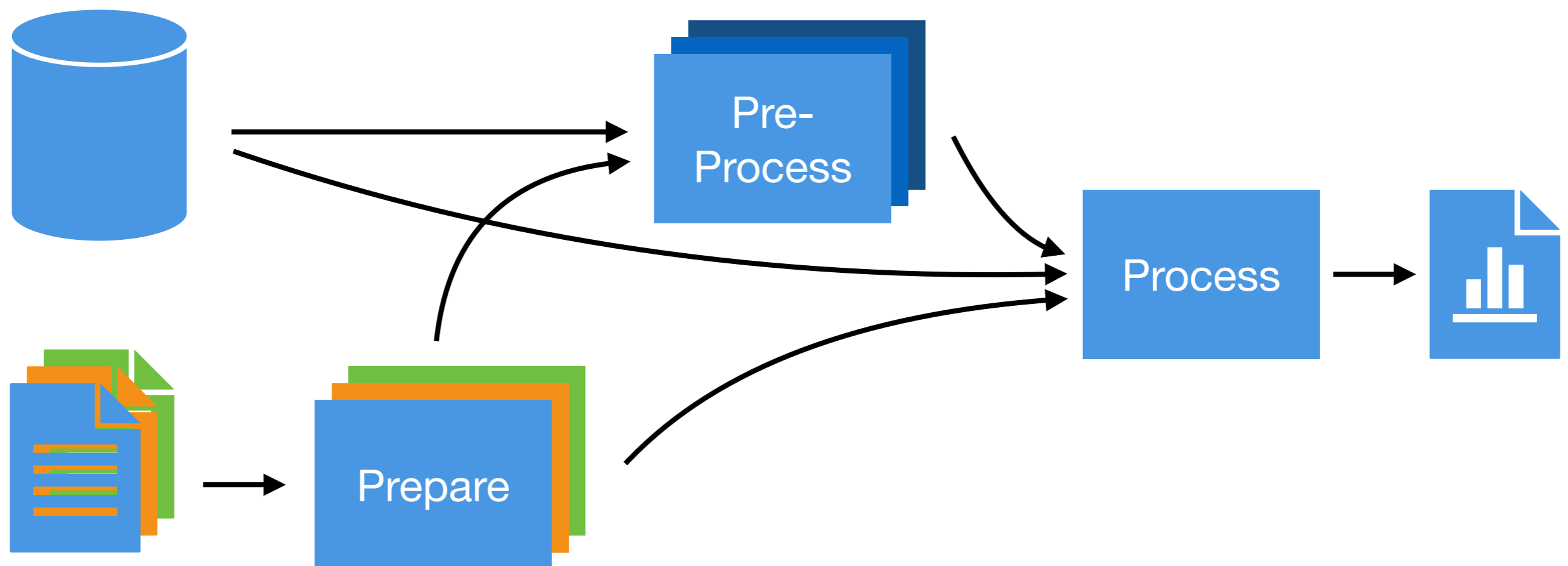


- Seems promising but still very new
- Documentation not complete yet

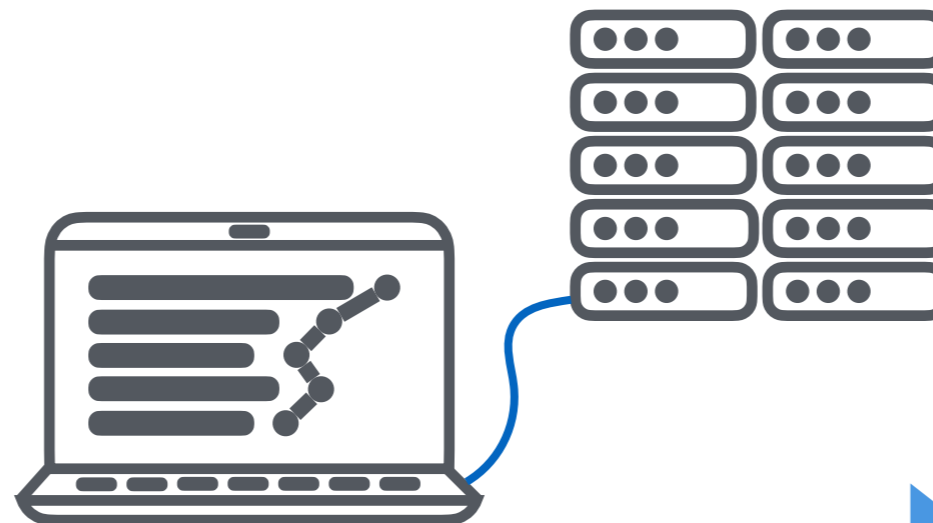
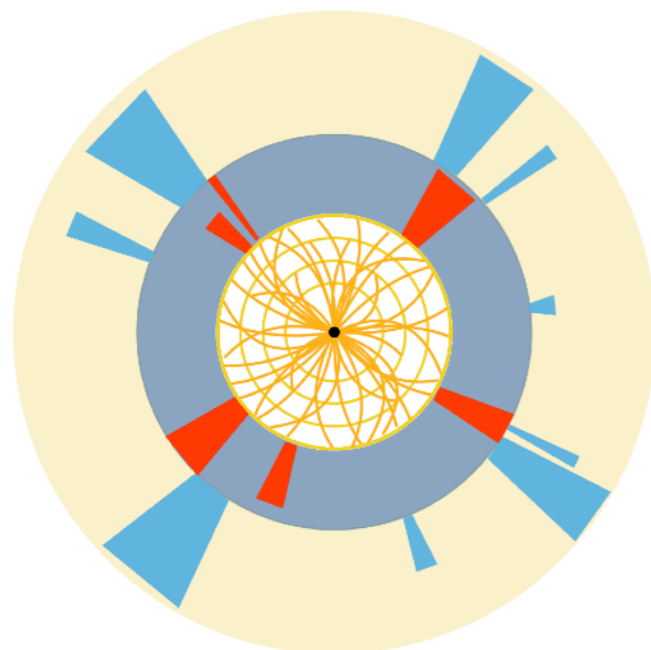
- Read-only task run 10 times (cycles) with 220 workers
- Data:
  - Higgs pair production analysis (1440GB,  $10^9$  events, 120 columns)
  - Using read-optimised compression algorithm (Z-std, 10)
- Results:
  - **Gradual performance: work stealing**
  - **Runtime lower bound: CPU → IO Bottleneck solved**



- Particular speedup for read heavy pre-processors (e.g. PU counting)
- Enables different analysis run-modes for full analysis (e.g. 2016, 3TB):
  - **Explorative** (w/o correction & systematic shifts, only simulations): **20min**
  - **Quick** (w/o computation heavy systematic variations (JEC)): **6h**
  - **Full** publication-ready analysis run: **20h**



- Fast O(TB) Physics Analysis on Small Institute Cluster
- Columnar processing via NumPy, Awkward
- Job distribution via map & reduce (Dask)
- Solved IO bottleneck:
  - Optimized storage distribution
  - Affine caching concept
- Analysis runtimes:
  - Explorative: 20 min
  - Full: 20h



20min - 20h

