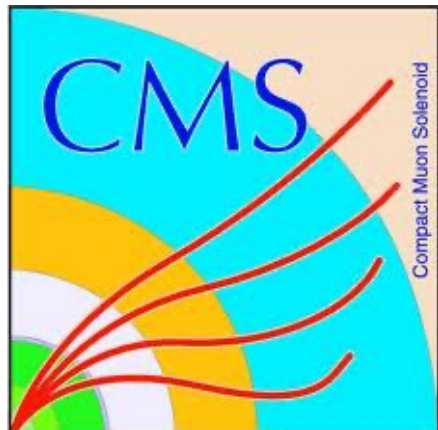# Object Stores for CMS data

**Nick Smith**, Bo Jayatilaka, David Mason, Oliver Gutsche, Alison Peisker, Robert Illingworth, Chris Jones (FNAL)
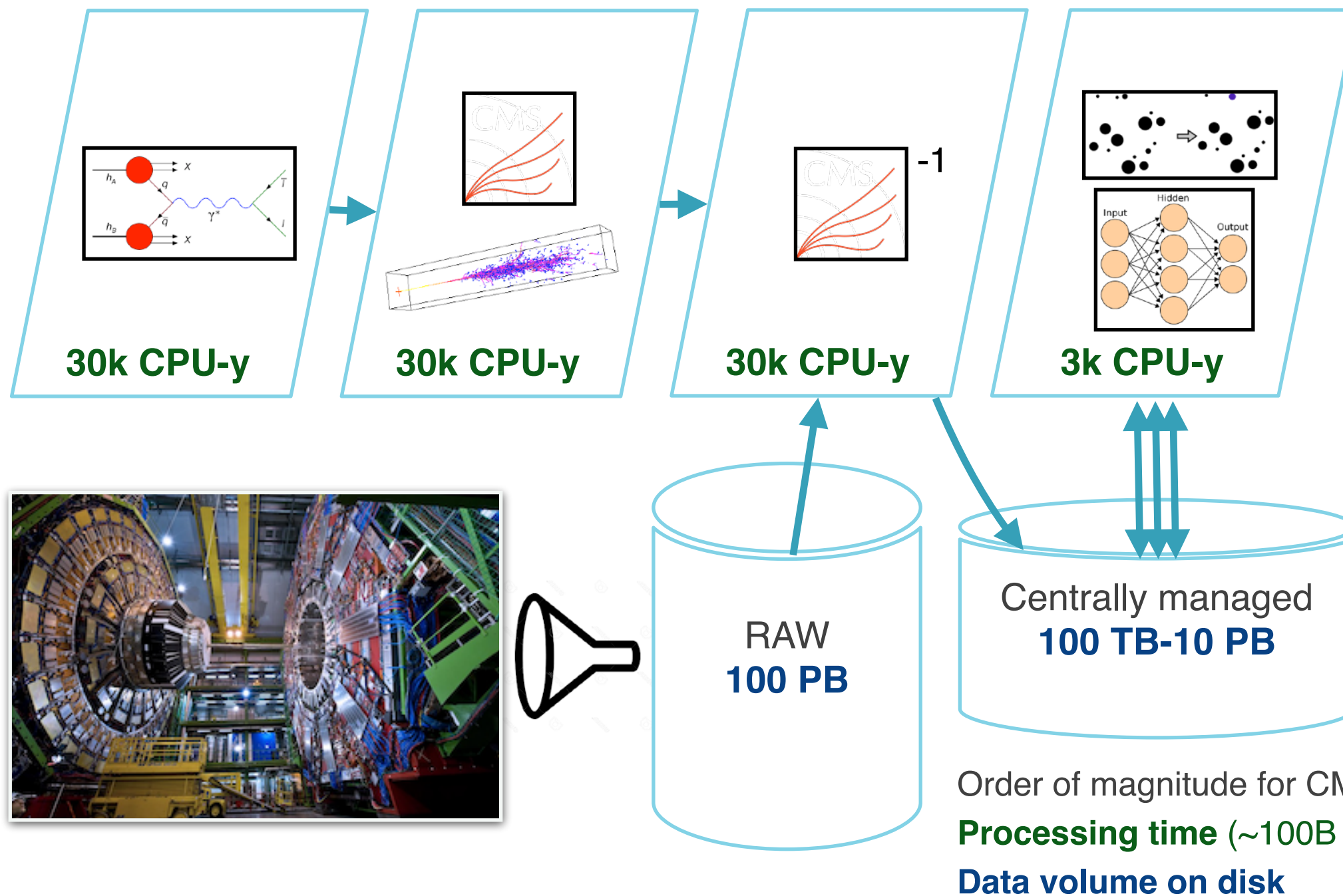
CHEP 2023
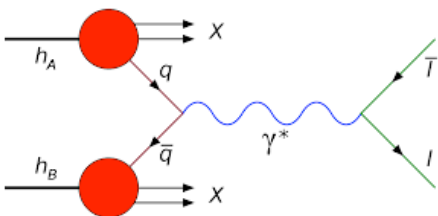
9 May 2023



Rothko, Number 19 (1949)

# Our inference pipeline



**30k CPU-y**  **30k CPU-y**  **30k CPU-y**  **3k CPU-y**

RAW
**100 PB**

Centrally managed
**100 TB-10 PB**

Order of magnitude for CMS Run-II (2016-18)
**Processing time** (~100B events)
**Data volume on disk**

🔷 **Fermilab**

# Centrally managed data



Primary dataset

Abstract, "what kind of events."
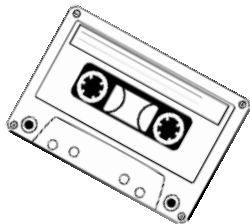e.g. hard scatter process for simulation, trigger filter for data

**Data tiers**

**AOD** — 1e5/event

Data columns pertaining to
low-level reconstruction

.root — 1e9/file

**MiniAOD** — 1e4/event

Calibrated physics objects
Particle-flow candidates

.root — 1e9/file

. . .

**Data volume**
order of magnitude
[bytes]

🔷 **Fermilab**

# Centrally managed data



Primary dataset

Abstract, "what kind of events."
e.g. hard scatter process for simulation, trigger filter for data
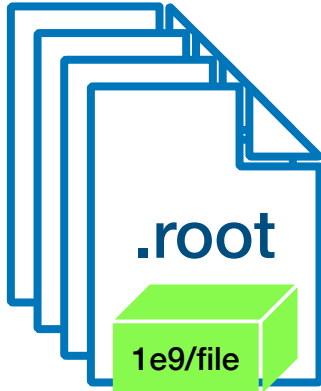
Data tiers

AOD · 1e5/event

Data columns pertaining to low-level reconstruction

.root · 1e9/file

MiniAOD · 1e4/event
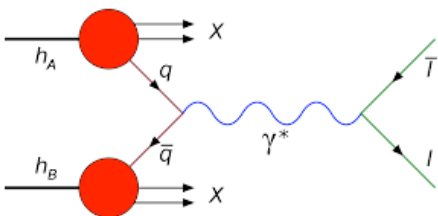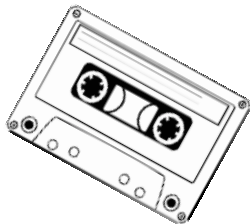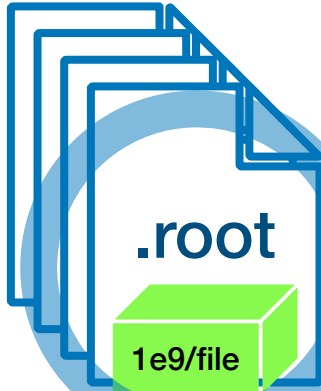
Calibrated physics objects
Particle-flow candidates

.root · 1e9/file

Data volume
order of magnitude
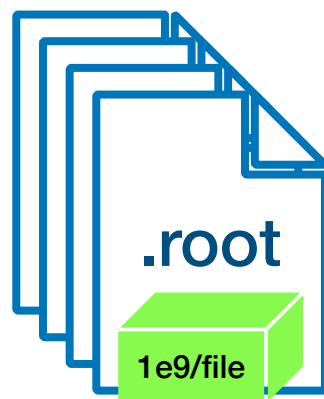[bytes]

Fermilab

# Primary

Abstract, "what kind
e.g. hard scatter proc

## Data tiers

### AOD

Data columns pe
low-level reconst

.root

`1e9/file`

Accessed

Not accessed



lume

gnitude
s]

🔷 **Fermilab**

# Projected usage

**Fermilab**

# Projected usage

**CMS** *Public*
Total Disk
*2022 Estimates*

- No R&D improvements
- Weighted probable scenario
- 10 to 20% annual resource increase



**CMS** *Public*
Total Tape
*2022 Estimates*

- No R&D improvements
- Weighted probable scenario
- 10 to 20% annual resource increase



**CMS** *Public*
Total Disk HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

CACHE: 13%
AODSim: 11%
MINIAOD: 13%
AOD: 12%
ALCARECO: 4%
USER: 4%
SKIM: 7%
RECOSim: 2%
MINIAODSim: 23%
RECO: 5%
RAWSim: 4%
NANOAODSim: 3%
PREMIX: 3%
OPERATIONS: 10%
Other: 5%

**900 PB**



**CMS** *Public*
Total Tape usage HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

HIAOD: 8%
AODSim: 12%
HIRAW: 12%
AOD: 11%
MINIAOD: 2%
ALCARECO: 4%
MINIAODSim: 3%
SKIM: 6%
Other: 4%
RAW: 39%

**2400 PB**

**Fermilab**

# Strawman

- What if we stored batches of events for each data product individually?
    - No more merge jobs!
- Most content does not change with re-processing
    - Even for UltraLegacy, already two MiniAOD versions
    - Keeping only new products would save a lot of disk

## Data-tier scheme

| MiniAOD Data product | KB per event | |
|---|---|---|
| | v1 | v2 |
| packed+pruned genParticles | 5.7 | 5.7 |
| slimmedElectrons | 1.3 | 1.3 |
| Others | 48.7 | 48.7 |
| Total | 55.7 | 55.7 |

## Column scheme

| MiniAOD Data product | KB per event | |
|---|---|---|
| | v1 | v2 |
| packed+pruned genParticles | 5.7 | - |
| slimmedElectrons | 1.3 | - |
| Others | 48.7 | - |
| Updated slimmedElectrons | - | 1.3 |
| Total | 55.7 | 1.3 |

Numbers sourced from a CMS UL17 TTBar simulation file

🔷 **Fermilab**

# Object store vs. filesystem

- Traditional data storage technology: distributed filesystem
  - e.g. NFS, EOS, dCache, Lustre, HDFS*, …
  - Often with remote access protocol (xrootd)
  - Files are concurrently read/writeable
- Popular new-ish technology: object store
  - Native remote access (http)
  - Objects are immutable (overwrite possible)



attrib

🔷 Fermilab

# Breaking down the ROOT file

- Essentially storing (+ moving) smaller units
  - This is usually a bad thing
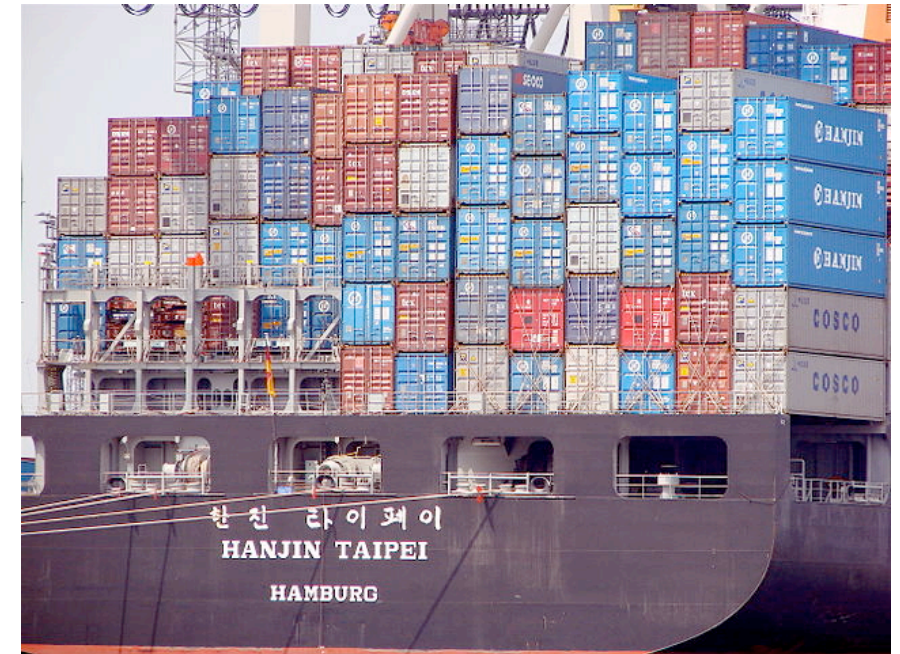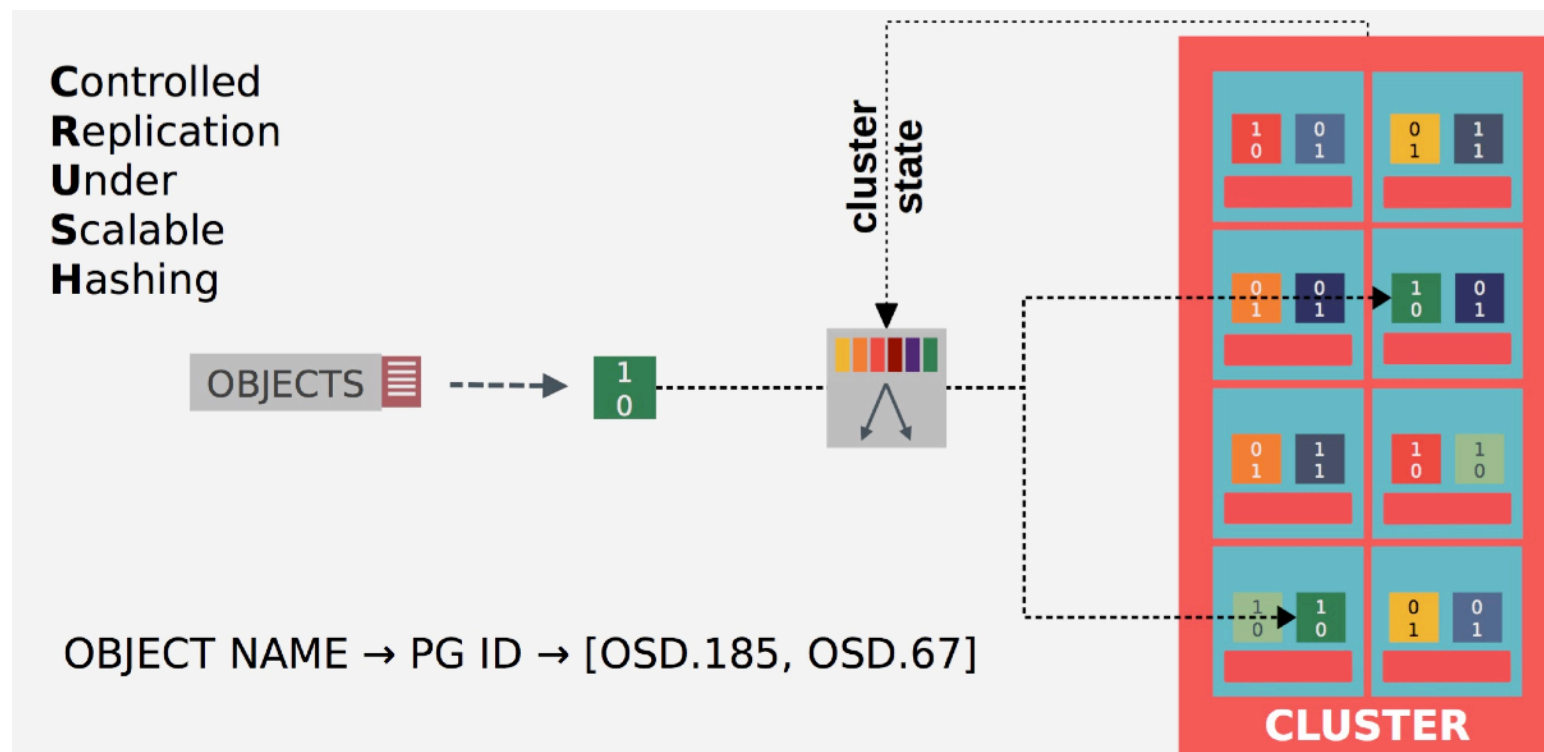


Intermodal container

?



Break-bulk cargo

Fermilab

# Breaking down the ROOT file

- Essentially storing (+ moving) smaller units
  - This is usually a bad thing
- Calculated placement
  - Like a hash, client-side
  - Downside: cluster state change causes reshuffle
    - Consistent hashing to minimize movement
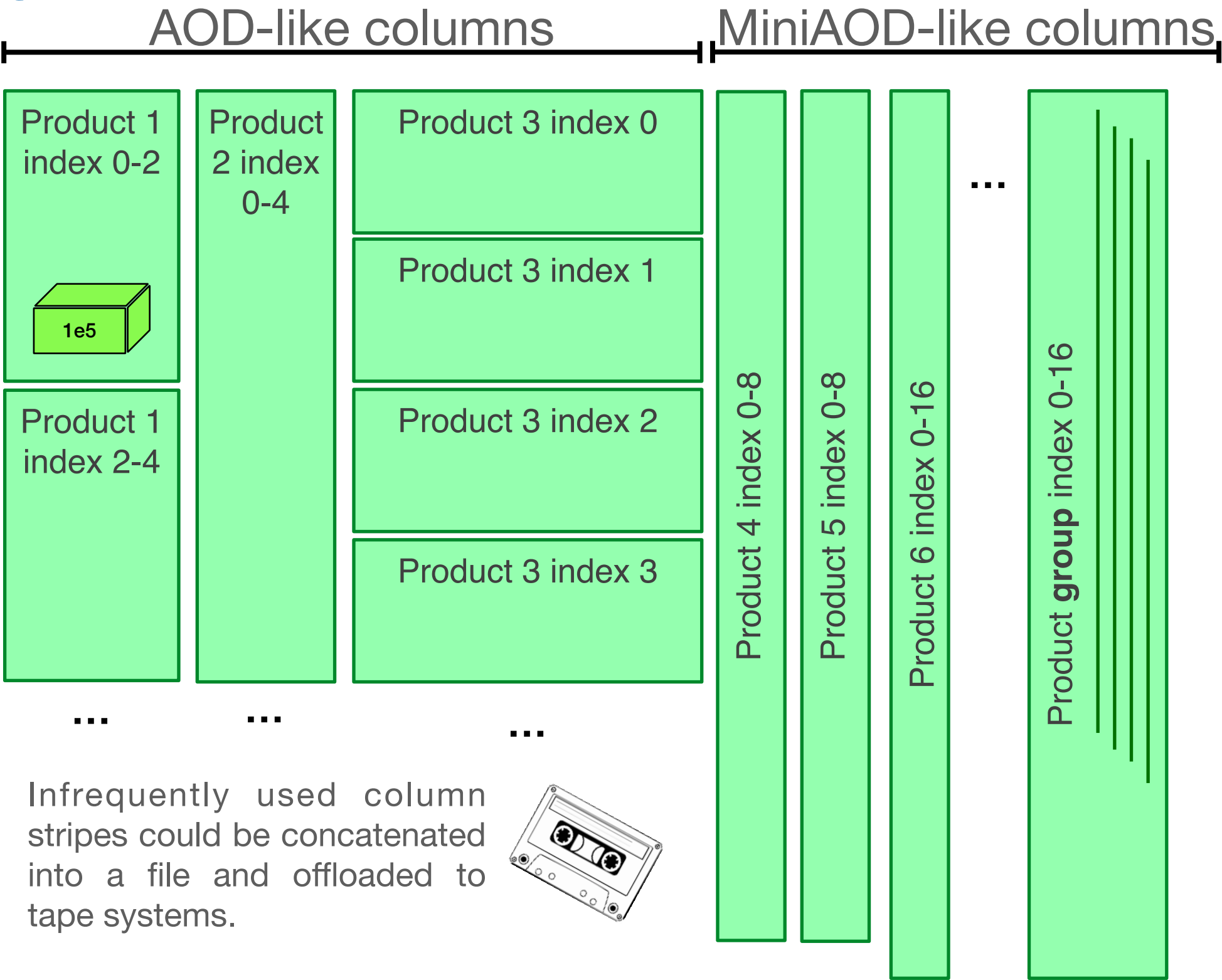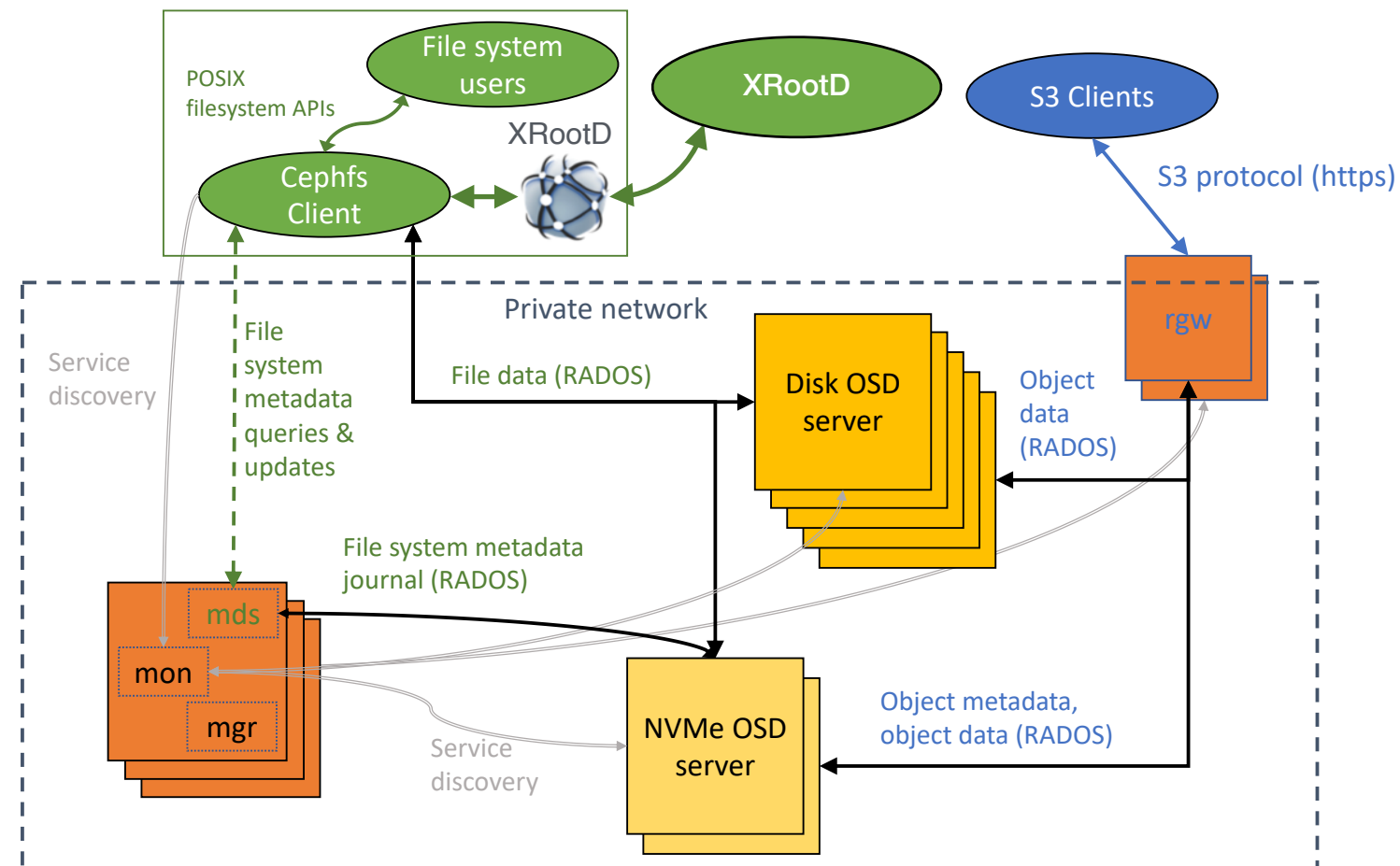


Intermodal container



**C**ontrolled
**R**eplication
**U**nder
**S**calable
**H**ashing

OBJECTS

cluster state

OBJECT NAME → PG ID → [OSD.185, OSD.67]

CLUSTER



amazon S3

ceph

Break-bulk cargo

🔷 Fermilab

# Object data format

AOD-like columns

MiniAOD-like columns

**Index object**

Product 1 metadata

Product 2 metadata

Product 3 metadata

...

Event batch index

EventID

EventID

...

Event batch index

EventID

EventID

...

...

1e5

Product 1 index 0-2

1e5

Product 1 index 2-4

Product 2 index 0-4

Product 3 index 0

Product 3 index 1

Product 3 index 2

Product 3 index 3

Product 4 index 0-8

Product 5 index 0-8

Product 6 index 0-16

Product **group** index 0-16

...

...

...

...

Infrequently used column stripes could be concatenated into a file and offloaded to tape systems.

🟦 **Fermilab**

# Test cluster

- Ceph pilot cluster setup at FNAL
  - 9 retired dCache machines
    - Total 2 PB HDD, circa 2014-2018
    - 288 OSDs
  - Two servers for metadata
    - 20TB NVMe (32 OSDs)
- Edge machines for:
  - xrootd door to CephFS
  - Ceph management daemons
  - RadosGW
    - Implements S3 protocol
    - Auth: pre-shared key or OIDC token
- Obviously not production-grade
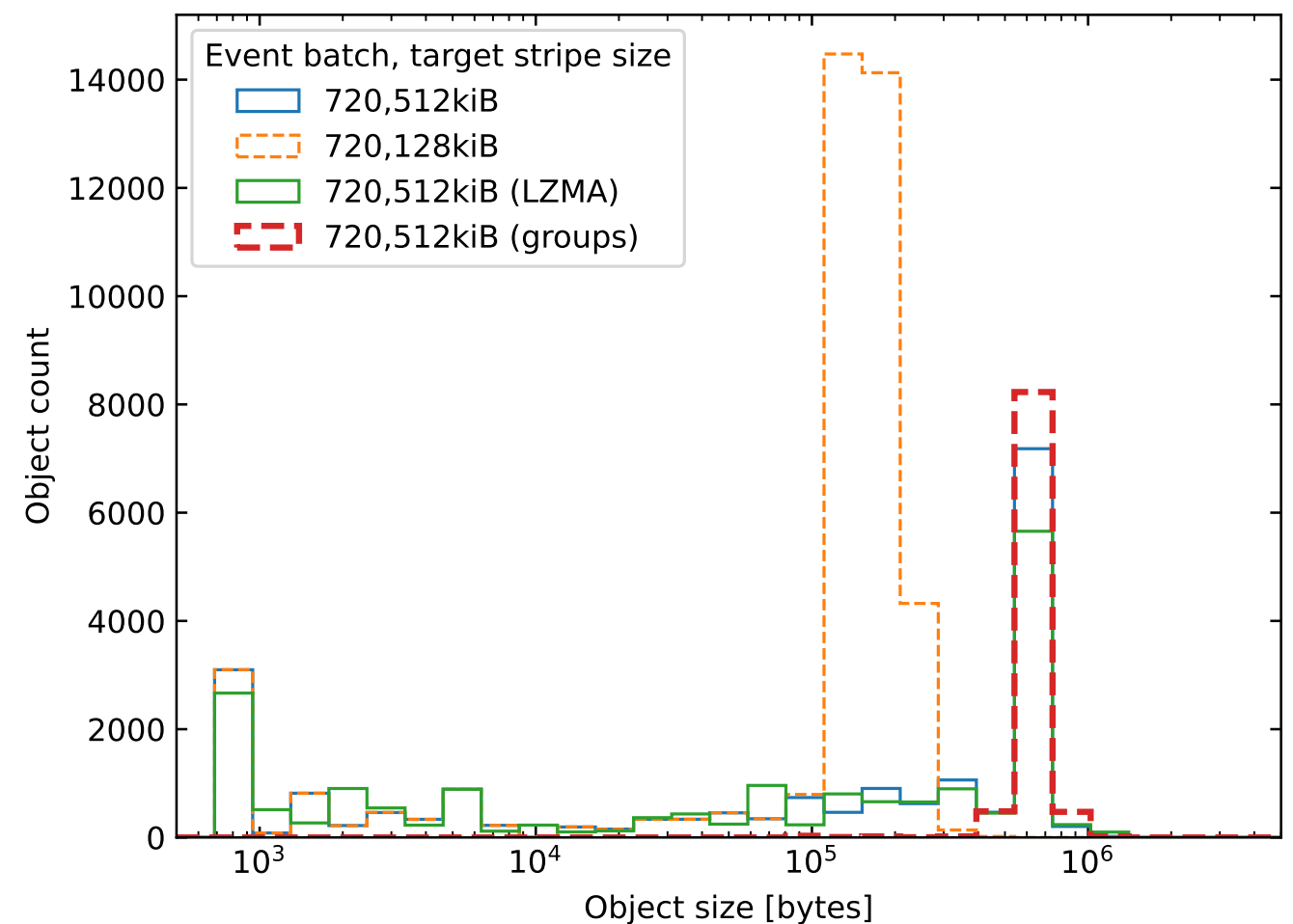  - Good for us: experience with failures!

🔷 Fermilab

# Client design

- Framework to evaluate alternative I/O strategies ([github](github))
  - Mimics CMS event processor design: TBB thread pool + tasks
  - Easy to add new output modules, simulate event processing, and test I/O
  - Serialization of data products: ROOT TBufferFile
- Developed S3 source and output module in framework
  - Using [libs3](libs3) + libcurl for protocol, async event loop separate from thread pool
  - Key features:
    - Parallel stream compression
    - Asynchronous I/O
    - Row-wise to column-wise pivot
- In following slides: stress testing the RadosGW server
  - Using many clients in parallel
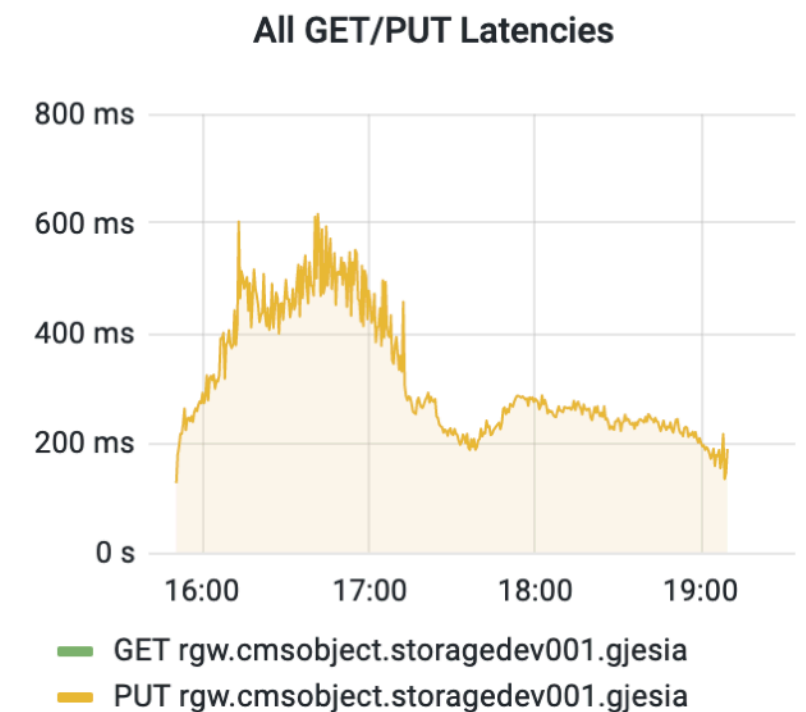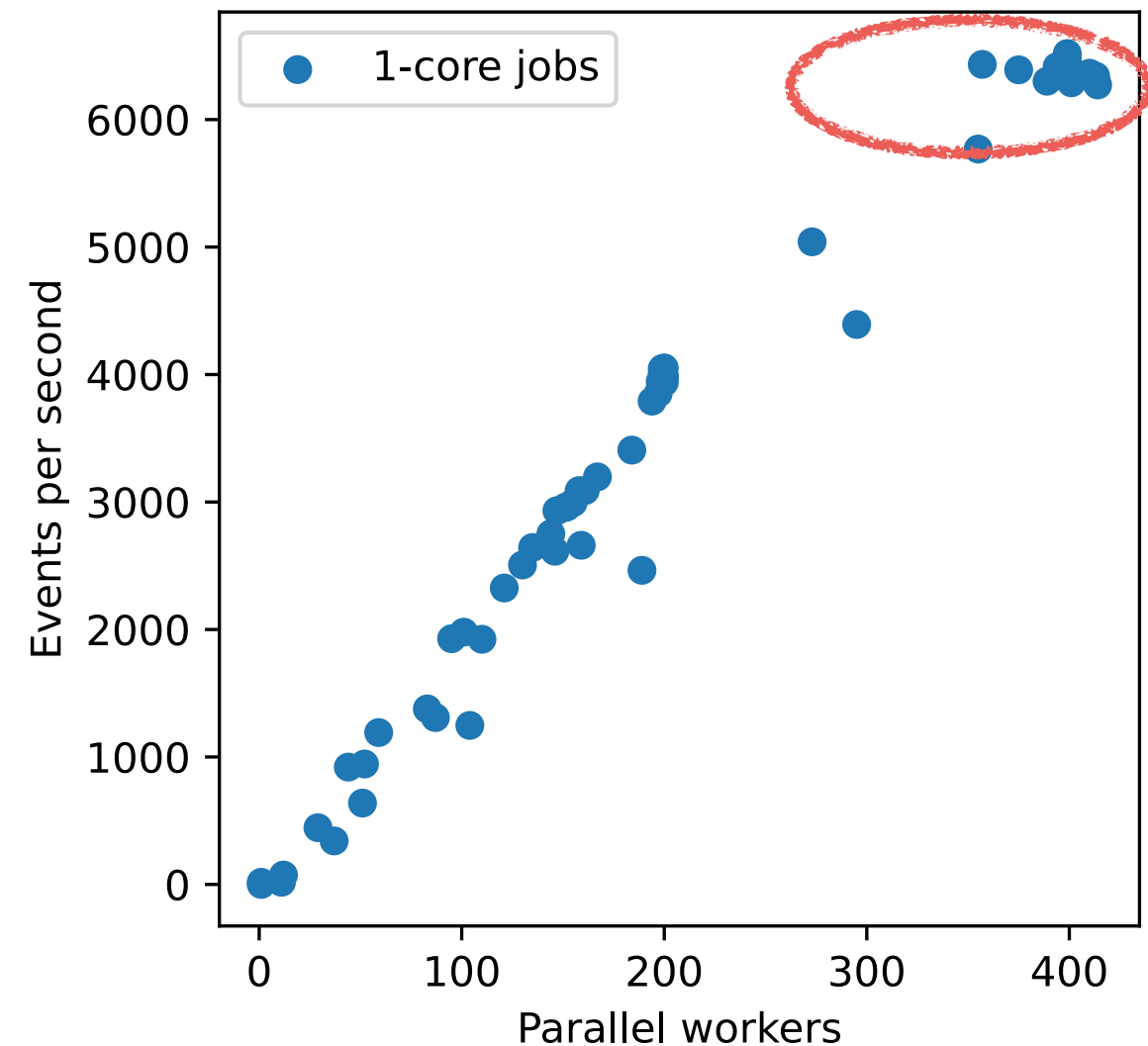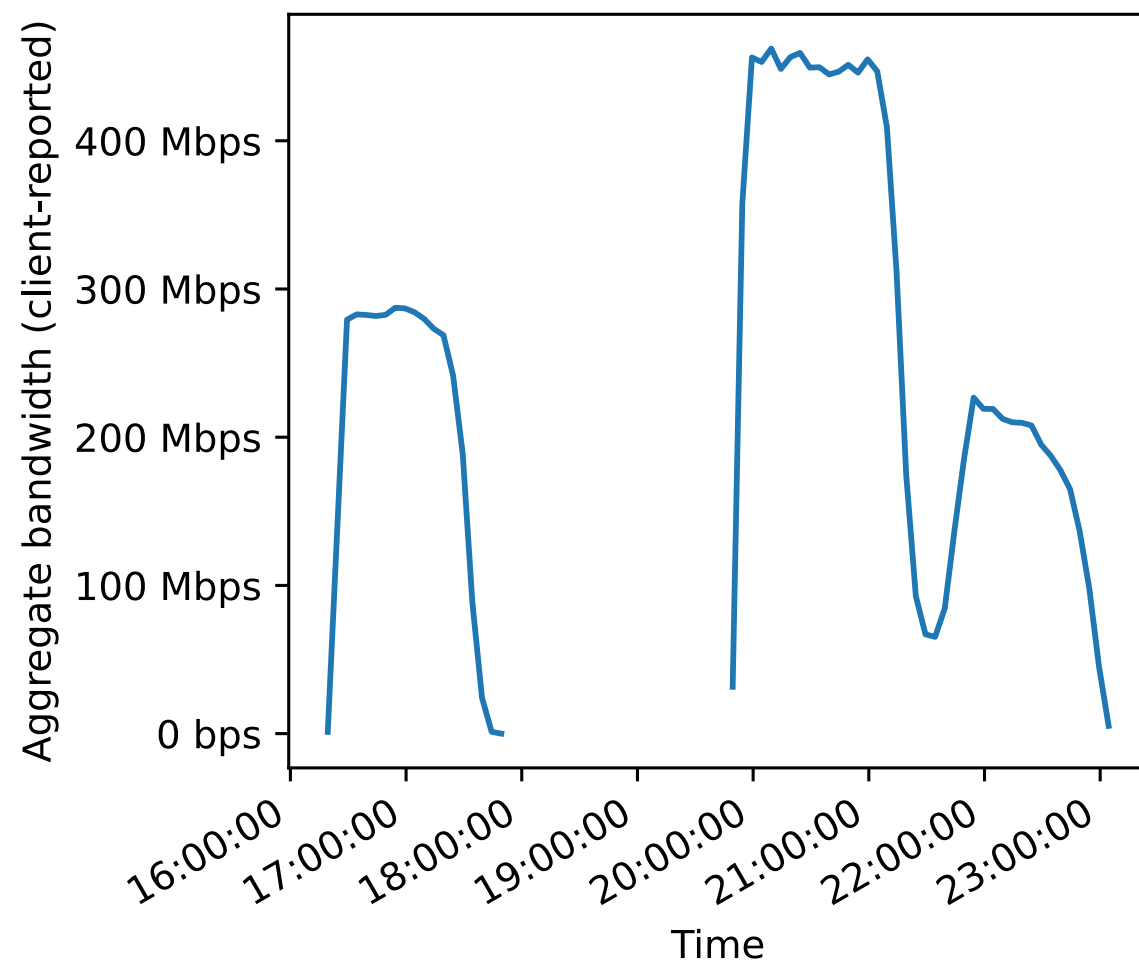
🔷 **Fermilab**

# Storage efficiency

- Input: 80k event MiniAOD file
  - LZMA compression
- Various S3 output configurations tested
- For erasure-coded Ceph pools, minimum object granularity of k*4kiB
  - Implies wasted space (vs. overhead for data resiliency)
  - Wasted space for EC4+2 in % listed below

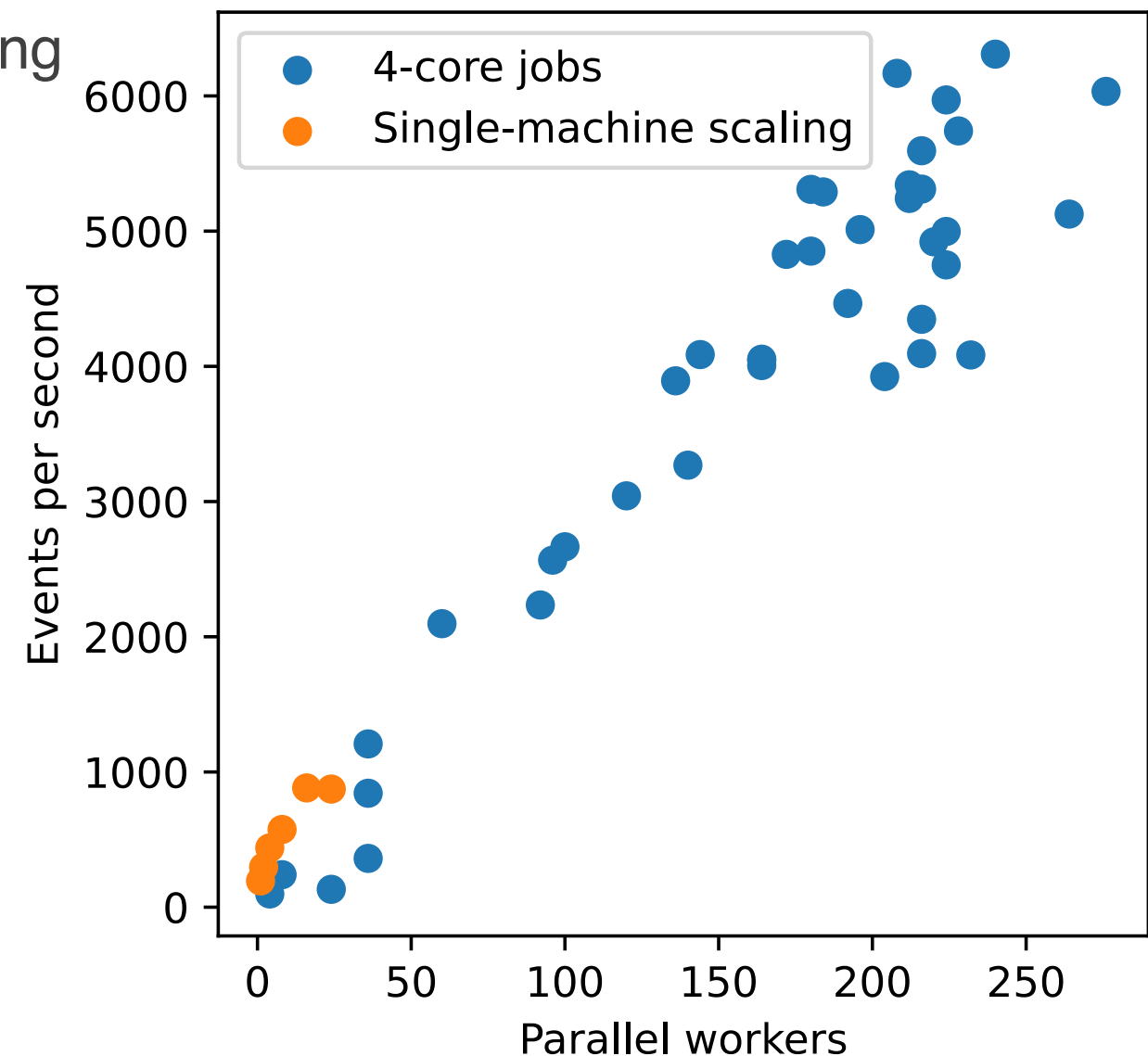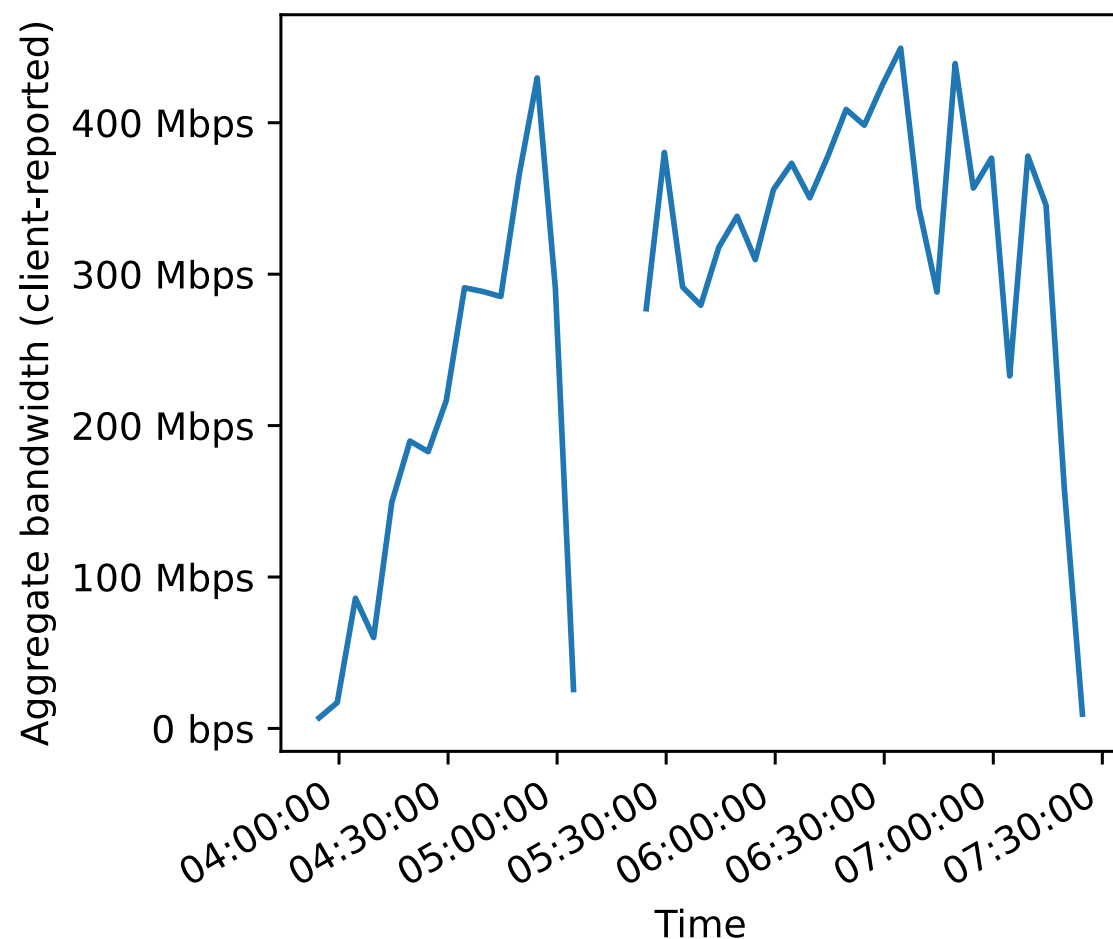| Format | KB per event |
|---|---|
| **MiniAOD input** | **55.7** |
| **Objects:**<br>- **event batch size 720**<br>- **target stripe size 128kiB** | **71.4 + 6.5%** |
| **Objects:**<br>- **event batch size 720**<br>- **target stripe size 512kiB** | **70.6 + 3.5%** |
| **Objects (LZMA):**<br>- **event batch size 720**<br>- **target stripe size 512kiB** | **61.8 + 3.7%** |
| **Objects (+product groups):**<br>- **event batch size 720**<br>- **target stripe size 512kiB** | **70.6 + 1.4%** |

🧭 **Fermilab**

# Write-only stress test

- Submit 1-core condor jobs
  - Read MiniAOD from FNAL dCache, write to S3
  - Handling up to 500 PUT/s
    - Past experience: can do ~1500 for smaller objects
  - Wrote 4.5 TB, 7.4 Mobj total
- Saturation at ~400 clients, ~400 MB/s*







All GET/PUT Latencies

GET rgw.cmsobject.storagedev001.gjesia
PUT rgw.cmsobject.storagedev001.gjesia
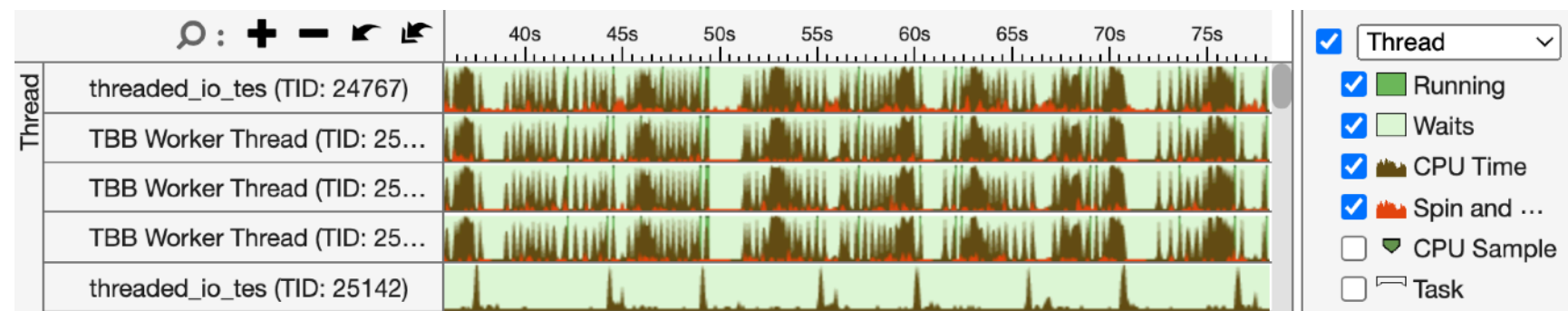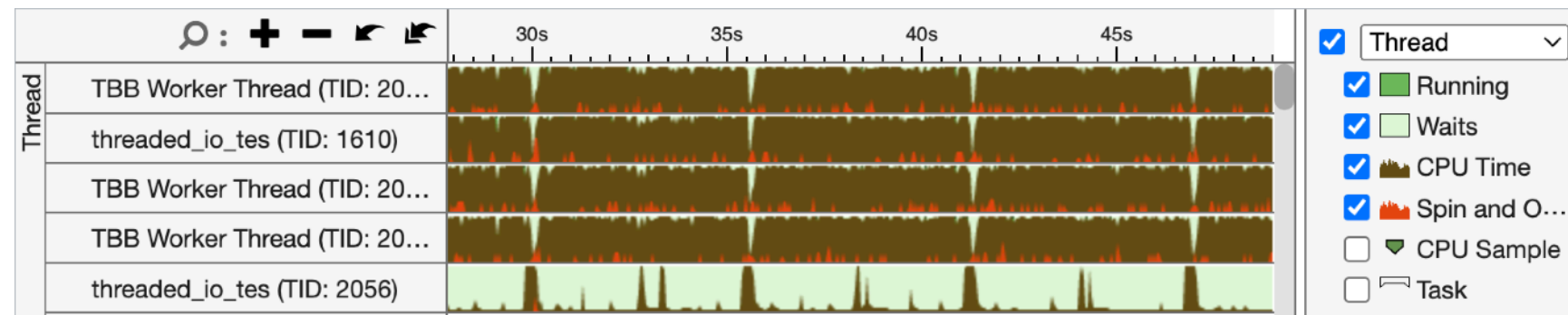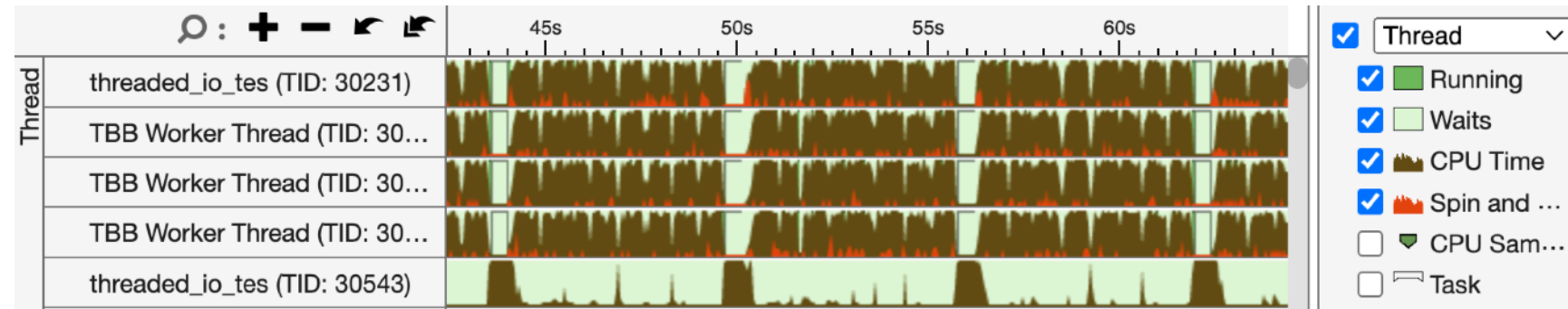
🔷 Fermilab

# Read-only stress test

- Submit 4-core condor jobs
  - Read from S3, decompress, deserialize
- Unable to reach saturation
  - Poor condor queue priority
  - Performance in line with single-machine scaling
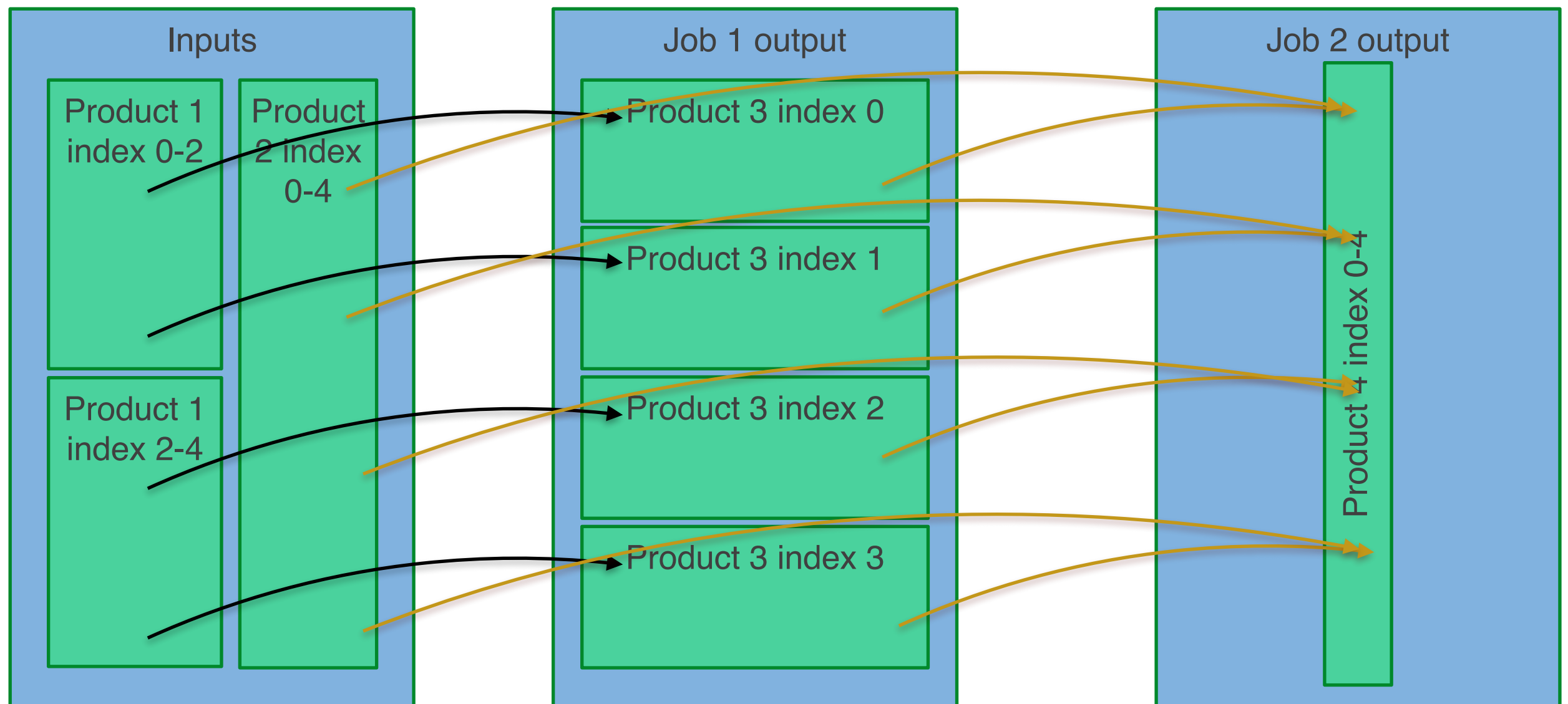    - (As shown at ACAT22)

🔷 Fermilab

# Client performance considerations

- Client application CPU inefficiency driven by I/O latency: either waiting for inputs or to flush output



- By pre-fetching input stripes and using a "fire-and-forget" output technique, CPU efficiency improves substantially



- When server is saturated, client CPU efficiency degrades significantly

🧵 **Fermilab**

# Next steps

- Demonstrate use case: job 2 reads job 1 and input products concurrently
  - Best example of advantage for column-level storage?

🎔 Fermilab

# Summary

- Object data formats provide new data management capabilities
  - Compared to current tier-based EDM file model
  - Reduce disk storage requirements for re-processing
  - Obviate the need to define data tiers
- In a prototype framework accessing a Ceph S3 service
  - On-disk data and metadata volume is as expected
  - Service scaling is promising: one RadosGW can serve ~400 client threads
- To fully utilize, more software development will be needed

🔷 Fermilab

# Backup

🎄 Fermilab

# S3Outputer design

Each box is a TBB task
Color = task group

Product 1 has stripes
written every 4 events

Product 2 has stripes
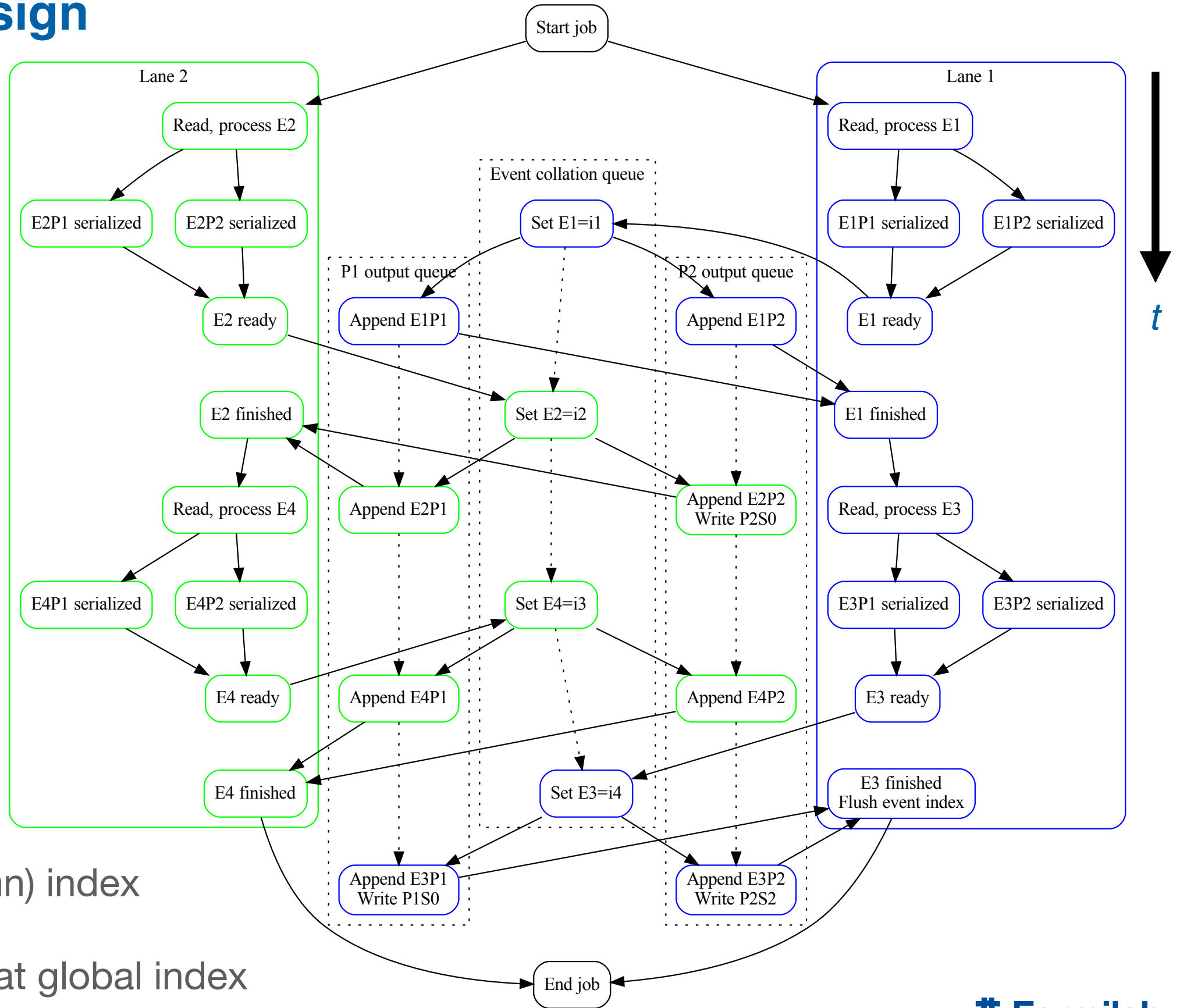written every 2 events

Label convention:
E* = Event number
P* = Product (column) index
i* = Global index
S* = Stripe starting at global index

🔬 Fermilab

# S3Source design

Each box is a TBB task
Color = task group

Product 1 has stripes
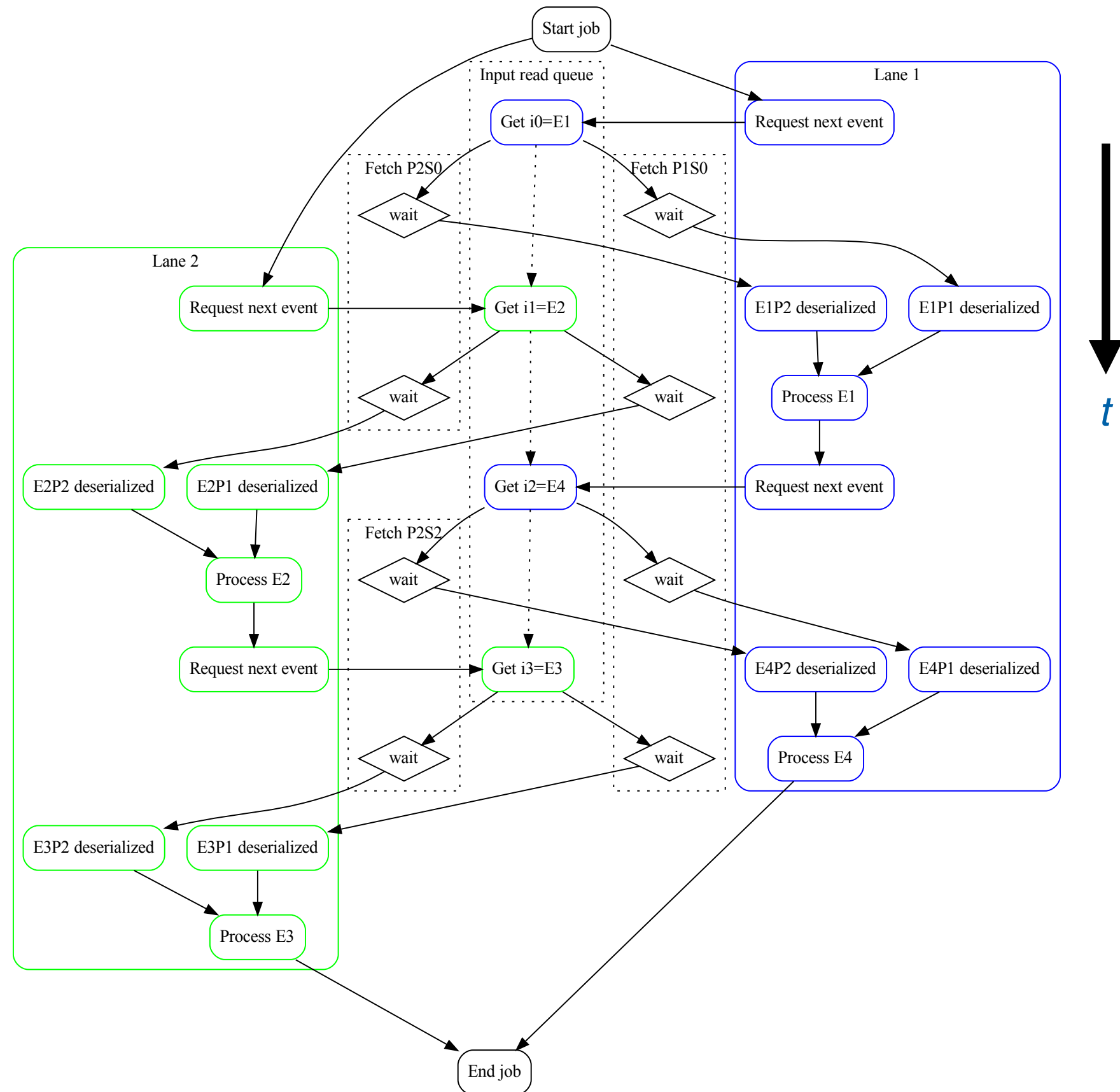read every 4 events

Product 2 has stripes
read every 2 events

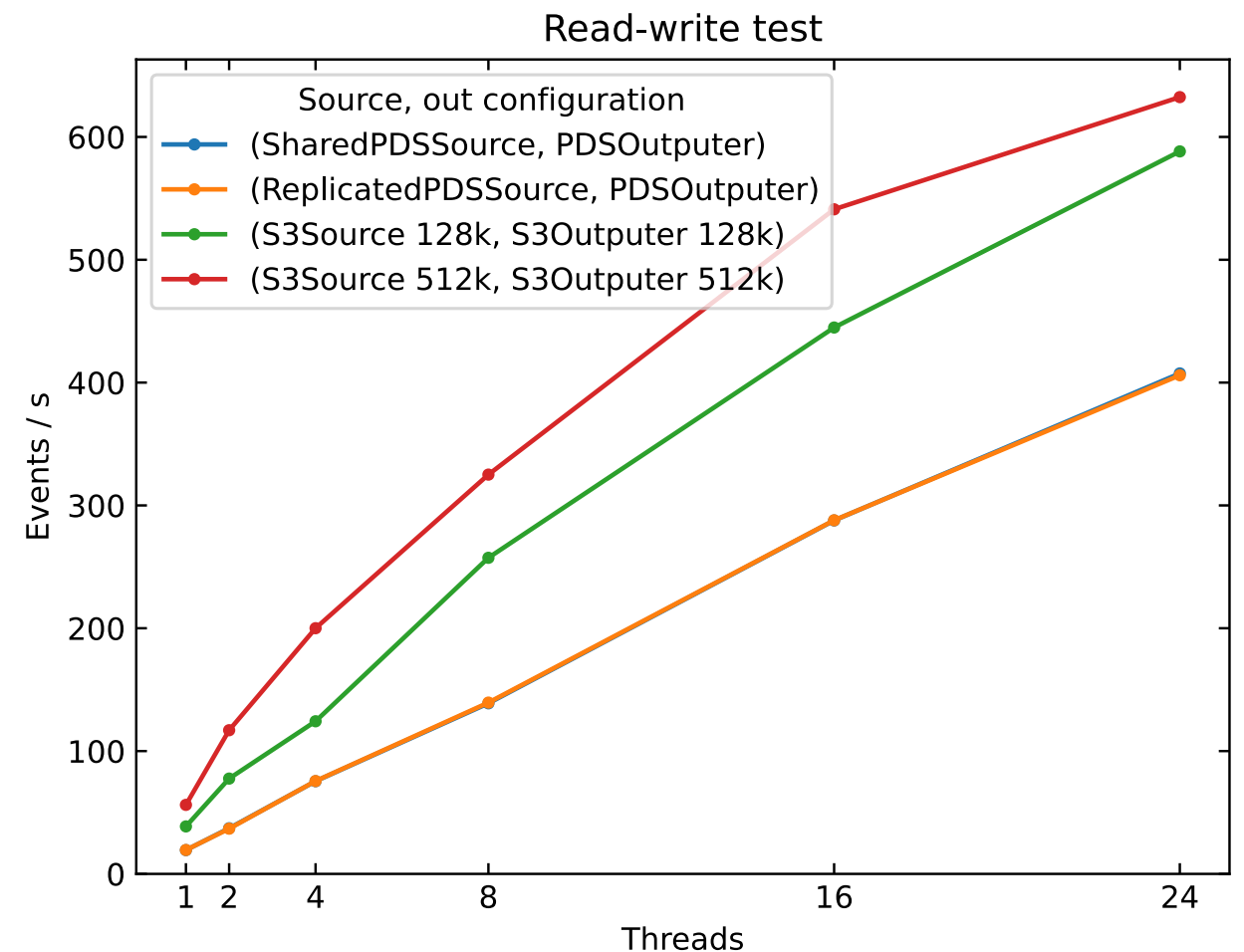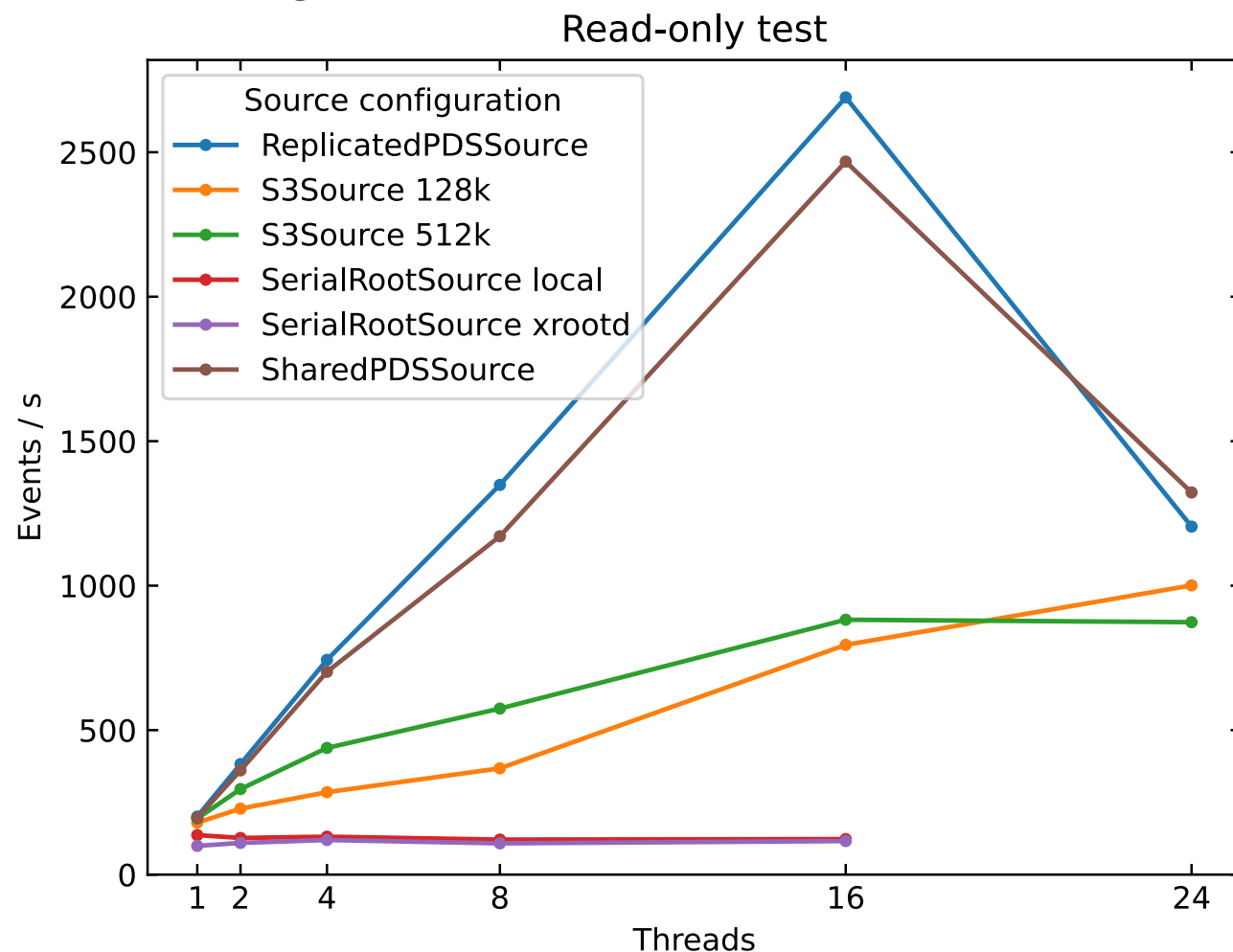Label convention:
E* = Event number
P* = Product (column) index
i* = Global index
S* = Stripe starting at global index

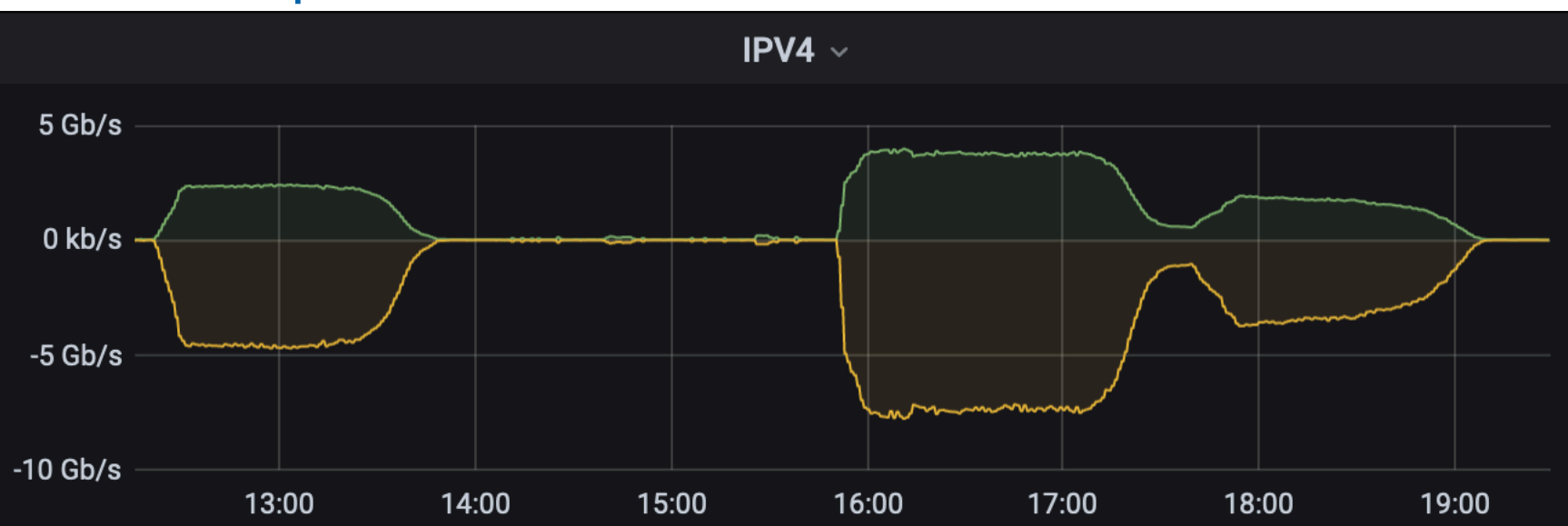🔷 Fermilab

# S3 vs. other Source/Outputers

- ROOT source similar to a CMSSW grid job
  - Read file via xrootd, server has CephFS (same cluster) mounted
  - No ROOT outputer due to bug
- PDS source: write whole events sequentially
  - Very good thread scaling (last data point = all cores on machine)
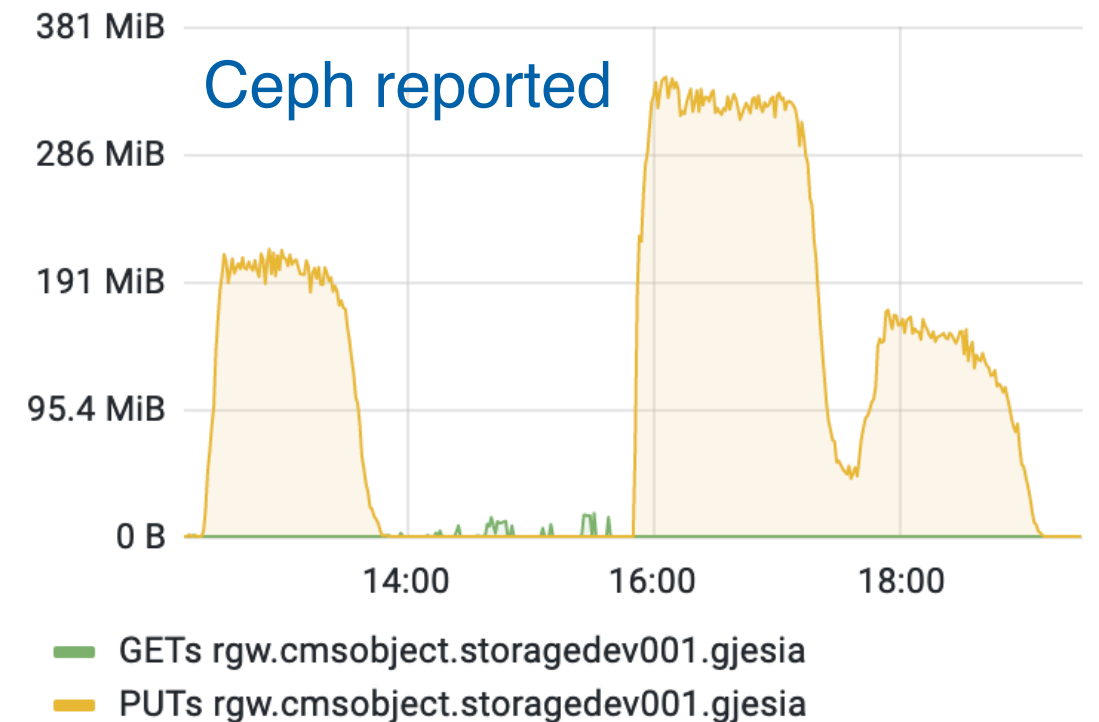  - Writing to local file rather than remote server

# Bandwidth inconsistencies

Ceph cluster prometheus metrics under-report the RadosGW bandwidth compared to server IP traffic & client aggregate bandwidth (recorded by application)



Ceph reported

**Bandwidth by HTTP Operation**

— GETs rgw.cmsobject.storagedev001.gjesia
— PUTs rgw.cmsobject.storagedev001.gjesia

Server reported



Clients reported

🧡 Fermilab