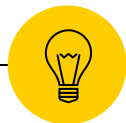


# Leveraging public cloud resources for the processing of CMS open data

CHEP - October 19 - 25, 2024



Kati Lassila-Perini  
Helsinki Institute of Physics - Finland  
Tom Cordruwisch, Subash Jayawardhana  
Lapland University of Applied Sciences - Finland

---

1

# CMS Open Data

---

# Decade of CMS open data - with a small dedicated team

Type something Help About ▾

Search

13 result(s) found Sort by Most recent

Current parameters Clear all

News x CMS x

Availability

Include on-demand datasets

Type

- Dataset (42,583)
  - Collision (342)
  - Derived (253)
  - Simulated (41,988)
- Documentation (54)
  - About (4)
  - Activities (13)
  - Authors (3)
  - Guide (27)
  - Help (2)

- CMS releases 13 TeV proton collision data from 2010 and 2011
- CMS completes Run-1 heavy ion open data collection
- CMS completes the release of its entire Run-1 proton-proton data
- First CMS open data from LHC Run 2 released
- CMS releases heavy-ion data from 2010 and 2011
- CERN Open Data Policy for the LHC Experiments
- CMS completes 2010-2011 proton-proton data release
- CMS releases open data for Machine Learning
- Observing the Higgs with over one petabyte of new data
- The Future of Particle Physics is "Open"
- Improving educational content with high-school teachers
- CMS releases new batch of research data from LHC Run 2
- CMS releases first batch of high-level LHC open data

INSPIRE HEP

literature references.reference.doi:10.7483/OPENDATA.CMS

Literature Authors Jobs Seminars Conferences More...

85 results | cite all Citation Summary  Most Recent ▾

**Bridging Worlds: Achieving Language Interoperability between Julia and Python in Scientific Computing** #1  
Ianna Osborne, Jim Pivarski, Jerry Ling (Apr 28, 2024)  
Contribution to: ACAT2024 • e-Print: 2404.18170 [cs.PL]  
[pdf](#) [cite](#) [claim](#) [reference search](#) [0 citations](#)

**Sparks in the Dark** #2  
Olga Sunneborn Gudnadottir, Axel Gallén, Giulia Ripellino, Jochen Jens Heinrich, Raazesh Sainudin et al. (Apr 5, 2024)  
e-Print: 2404.04138 [hep-ex]  
[pdf](#) [cite](#) [claim](#) [reference search](#) [0 citations](#)

**Finetuning foundation models for joint analysis optimization in High Energy Physics** #3  
Matthias Vigl (Tech. U., Munich (main)), Nicole Hartman (Tech. U., Munich (main)), Lukas Heinrich (Tech. U., Munich (main)) (Jan 24, 2024)  
Published in: *Mach.Learn.Sci.Tech.* 5 (2024) 2, 025075 • e-Print: 2401.13536 [hep-ex]  
[pdf](#) [DOI](#) [cite](#) [claim](#) [reference search](#) [5 citations](#)

**Parametric Matrix Models** #4  
Patrick Cook, Danny Jammooa, Morten Hjorth-Jensen, Daniel D. Lee, Dean Lee (Jan 22, 2024)  
e-Print: 2401.11694 [cs.LG]  
[pdf](#) [cite](#) [claim](#) [reference search](#) [3 citations](#)

Date of paper

Number of authors

- Single author 15
- 10 authors or less 74

Exclude RPP

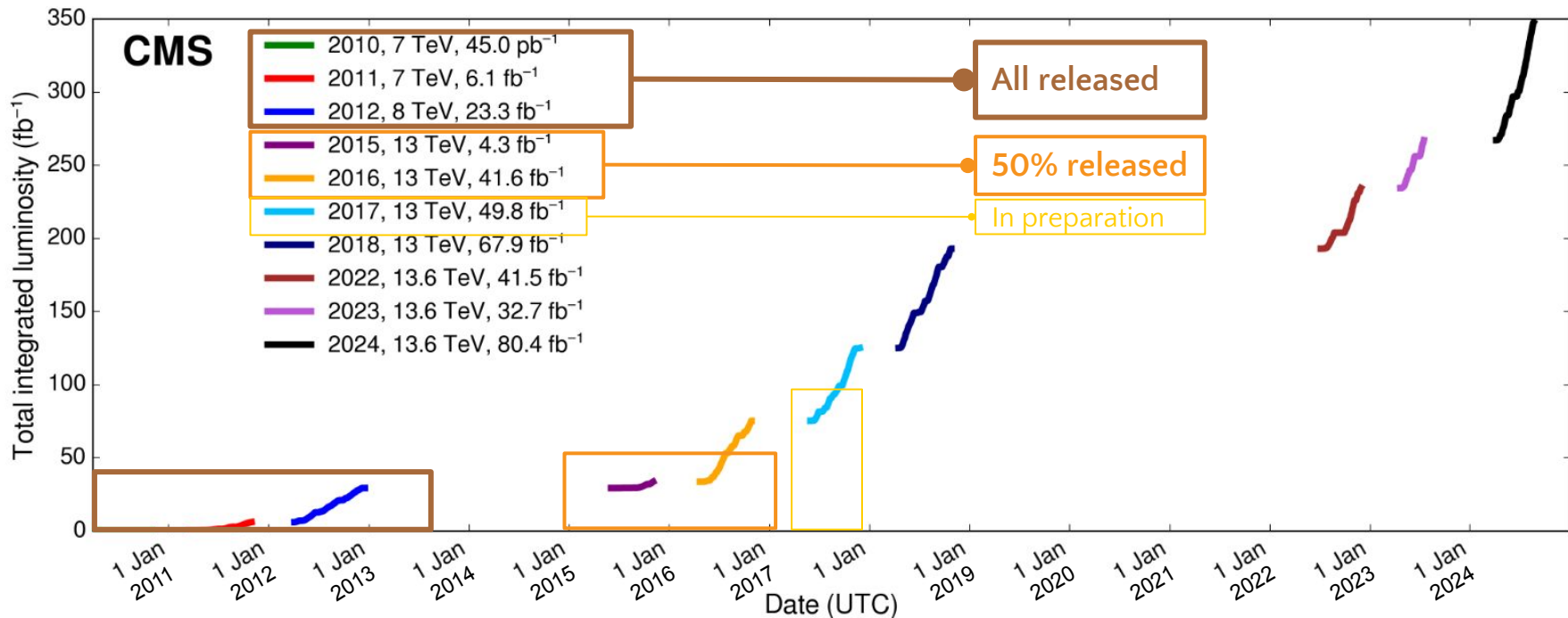
- Exclude Review of Particle Physics 85

Document Type

- article 53
- published 40
- conference paper 26
- thesis 7

## CMS open data in use





Open data quantity & time defined in the CMS open data policy

2

# Public cloud and Kubernetes

Why and how?



## Motivation: 1 - Why?

**NanoAOD:**  
compact format  
(1 kB/evt)  
good for many analyses  
no need for CMSSW

ease of  
use

CMS Open-Data Workshop 2024  
CERN IdeaSquare  
Jul 29 - Aug 1, 2024  
08:00 - 18:00 CEST  
Instructors: Matt Bellis, Julie Hogan, Kati Lassila-Perini, Tom McCauley, Sezen Sekmen  
Helpers: Xavier Tirtin, Daniela Merizalde, David Mena

completeness

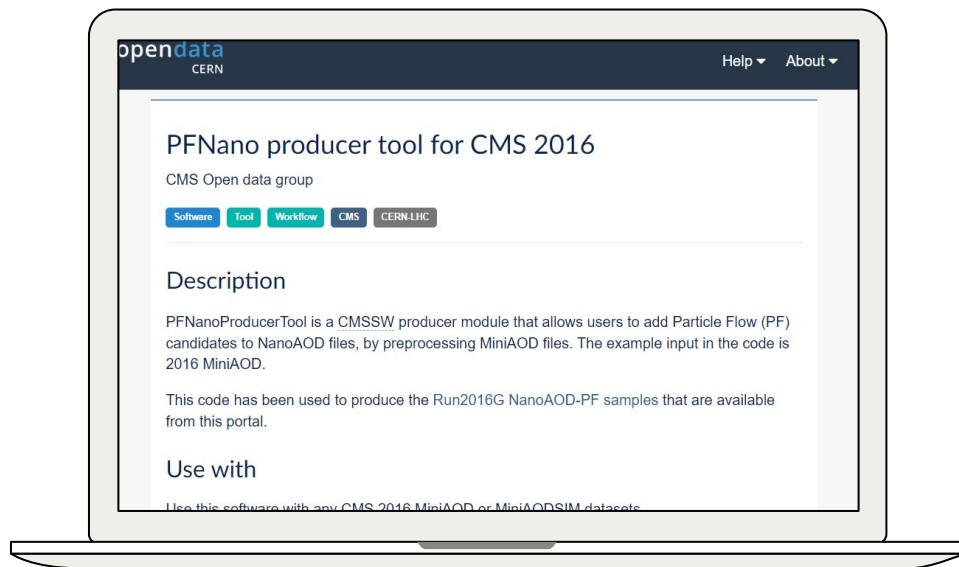
**Custom  
NanoAOD:  
enrich  
NanoAOD  
with what you  
need**

**MiniAOD:**  
richest Run-2 OD format  
(30-40 kB/evt)  
good for ~ all analyses  
requires CMSSW

CMS Open-Data Workshop 2023  
Fermilab  
July 11-14, 2023  
8:00 am - 12:00 am (CDT, UTC-5)  
Instructors: M. Bellis, J. Hogan, K. Lassila-Perini, T. McCauley, S. Sekmen, X. Tirtin, R. Trujillo  
Helpers: V. Morina, J. Nelson, H. Pederson, X. Shen



## Motivation: 2 - How?



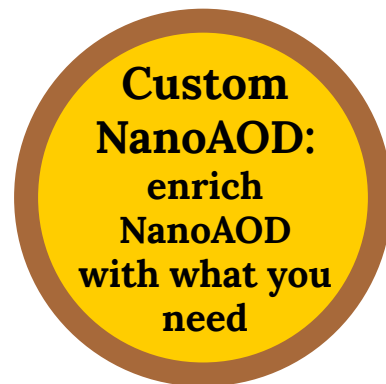
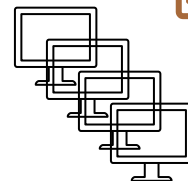
code:



env:



resources:





# Why cloud and Kubernetes?

**CMS Open Data Workshop at WHEPP XVII, 2024**  
WHEPP XVII, IIT Gandhinagar  
Jan 3-10, 2024  
6:45 - 7:15 pm IST

**Instructors:** Julie Hogan, Kati Lassila-Perini, Matt Bellis  
**Student Facilitators:** Aravind Sugunan, Ritik Saxena, Atul Jaiswal, Pruthvi Suryadevara, Mukund Shelake

		Do you have?	Want to get?	
1	Computing resources	✓	✗	<ul style="list-style-type: none"><li>No need for cloud resources</li><li>Ask your IT support</li></ul>
2	Tools	✗	✓	<ul style="list-style-type: none"><li>Read the CMS OD workshop tutorial (HTCondor)</li><li>Ask your IT support</li></ul>
3	Skills	✗	✓	<ul style="list-style-type: none"><li>Read the CMS OD workshop tutorial (HTCondor)</li><li>Ask your IT support</li></ul>
4	Your custom CMS OD	✗	✓	<ul style="list-style-type: none"><li>Process what you need</li><li>Download and store locally</li><li>Analyze on your own resources</li></ul>





## Why cloud and Kubernetes?

		Do you have?	Want to get?	Why cloud / Kubernetes?
1	Computing resources	✗	✓	<ul style="list-style-type: none"><li>• Short-term, immediate resources</li><li>• Pay what you use</li><li>• Compatible with CMS OD environment</li></ul>
2	Tools	✗	✓	<ul style="list-style-type: none"><li>• <a href="#">Kubernetes</a> as a basic tool</li><li>• <a href="#">Terraform</a> to deploy resources</li><li>• <a href="#">Argo workflows</a> to manage jobs</li><li>• Open source and free</li></ul>
3	Skills	✗	✓	<ul style="list-style-type: none"><li>• Applicable within and beyond research</li><li>• Attractive for early careers and young-minded</li><li>• Example setup provided - we did it for you!</li></ul>
4	Your custom CMS OD	✗	✓	<ul style="list-style-type: none"><li>• Process what you need</li><li>• Download and store locally</li><li>• Analyze on your own resources</li></ul>

3

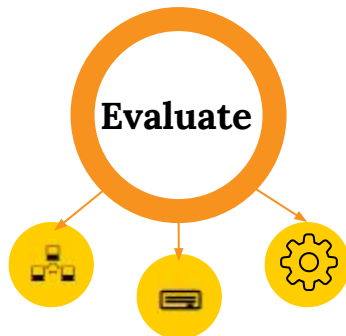
## How much time and money?

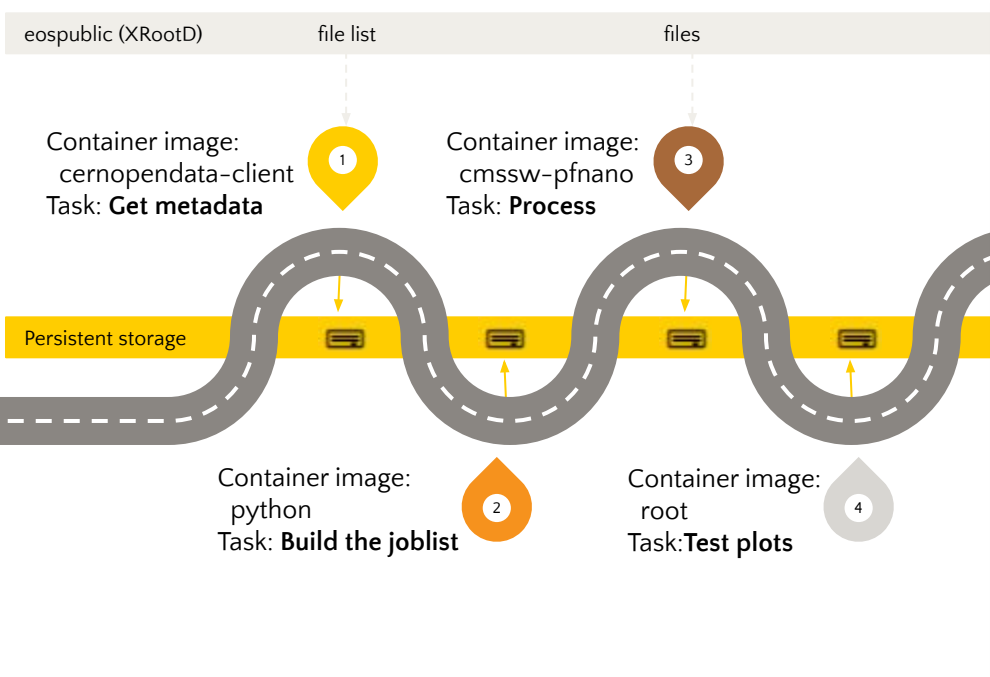
Resources from a Google Cloud Research credit 303424260



## Benchmarking use case

- Processing custom-NanoAOD:
  - input: MiniAOD data
  - full NanoAOD processing + Particle-Flow (PF) candidates
    - same as in the already-provided PF-Nano datasets on the portal
  - image: CMS OD image with PFNano processing code precompiled
  - using Argo workflow to run the processing
    - steps: get dataset metadata → make a joblist → process  $\times N_{\text{jobs}}$  → (test plot)





```

$ argo get -n argo @latest
Name: pfnano-process-6rwwj
Namespace: argo
ServiceAccount: argo-service-account
Status: Running
Conditions:
PodRunning True
Created: Fri Sep 27 23:48:05 +0200 (1 minute ago)
Started: Fri Sep 27 23:48:05 +0200 (1 minute ago)
Duration: 1 minute 28 seconds
Progress: 2/355
ResourcesDuration: 1s*(1 cpu),25s*(100Mi memory)
Parameters:
nEvents: -1
recid: 30511
nJobs: 353
bucket: test-gcs-argo-bucket-regional

STEP TEMPLATE PODNAME DURATION MESSAGE
● pfnano-process-6rwwj cms-od-example
├─✓ get-metadata get-metadata-template pfnano-process-6rwwj-get-metadata-template-4259457820 17s
├─✓ joblist joblist-template pfnano-process-6rwwj-joblist-template-3232337233 11s
├─● runpfnano(0:eventsinjob:-1,firstfile:1,it:1,lastfile:1) runpfnano-template pfnano-process-6rwwj-runpfnano-temp
├─● runpfnano(1:eventsinjob:-1,firstfile:2,it:2,lastfile:2) runpfnano-template pfnano-process-6rwwj-runpfnano-temp
[...]
├─● runpfnano(348:eventsinjob:-1,firstfile:349,it:349,lastfile:349) runpfnano-template pfnano-process-6rwwj-runpfn
├─● runpfnano(349:eventsinjob:-1,firstfile:350,it:350,lastfile:350) runpfnano-template pfnano-process-6rwwj-runpfn
├─● runpfnano(350:eventsinjob:-1,firstfile:351,it:351,lastfile:351) runpfnano-template pfnano-process-6rwwj-runpfn
├─● runpfnano(351:eventsinjob:-1,firstfile:352,it:352,lastfile:352) runpfnano-template pfnano-process-6rwwj-runpfn
├─● runpfnano(352:eventsinjob:-1,firstfile:353,it:353,lastfile:353) runpfnano-template pfnano-process-6rwwj-runpfn

```

## Workflow structure





## Cluster types

- ◉ Vocabulary
  - ◉ **GKE** = Google Kubernetes Engine
  - ◉ **Cluster** consists of nodes (machines with local boot disks)
  - ◉ Jobs run on **pods** (containerized applications) on the nodes
  - ◉ **Persistent storage** is a disk available for all nodes (and pods)
- ◉ Two types of GKE clusters: standard and auto-pilot:

### Using: **GKE Standard clusters**

- allows defining all cluster components (type and number of nodes)
- auto-scales (nodes deleted automatically) if so configured
- cost goes with time and depends on the cluster setup
- can create nodes with the container image from a secondary boot disk

### Not tested: GKE Auto-pilot clusters

- creates the cluster based on the resource request
- did not allow for a NFS disk server (container in the privileged mode)
- auto-scales (nodes deleted automatically)
- creating nodes with the image from the secondary boot disk did not work out of the box



## Disk types and cost

- NFS disk
  - predefined size
  - the cost / disk size, not / usage
  - requires an NFS server on the cluster
- Google Cloud Storage bucket
  - size not predefined
  - the cost / actual usage
- Storage:
  - negligible (for a cluster lifetime 1-2d)
- Download (“egress”):
  - costly

→ Our choice: GCS bucket  
for the ease of use

europa-west4 (Netherlands)	NFS	GCS bucket
Storage /month	0.44\$/10GB (22\$/500GB)  SSD: 1.87\$/10GB	0.20\$/10GB (10\$/500GB)
Download cost	0.12\$/GB ⚠️👉 (60\$/500GB) 👈⚠️	
Download time	Local download speed  Slight dependence on distance	Idem  Independent of distance



## Cluster and job configuration

- Node type
  - number of CPUs and amount of memory (“standard”/”highcpu”/”highmem”)
- Match with the job resource needs
  - use a test workflow running a single job / node to see the resource needs
  - for a quick, provider-independent check:

```
kubectl top node (check node resources in use)
```

```
kubectl top pods -n argo (check single pod usage)
```

```
$ kubectl top pods -n argo | grep runpfnano
pfnano-process-g2m2k-runpfnano-template-1439319107 998m 1394Mi
pfnano-process-g2m2k-runpfnano-template-2136346710 999m 1460Mi
```

- resource requests define how argo distributes the jobs to the nodes
  - goal: close to full occupancy
    - lack of memory kills
    - lack of CPU slows down

```
resources:
  requests:
    cpu: "800m"
    memory: "1.8Gi"
    ephemeral-storage: "5Gi"
```

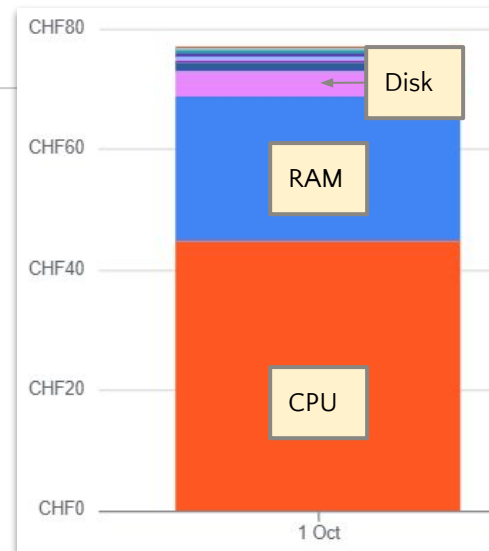
→ 1 job / vCPU



# Cost of full dataset processing

▼	Compute Engine	CHF72.94
▲	Cloud Monitoring	CHF1.69
●	Networking	CHF1.65
■	Kubernetes Engine	CHF0.76
◆	Cloud Storage	CHF0.08

SKU	Service	SKU ID	Usage	Cost
■	E2 Instance Core running in Netherlands	012A-5DBB-1352	2,191.27 hour	CHF44.69
●	E2 Instance Ram running in Netherlands	D9EA-4FF0-E394	8,765.17 gibibyte hour	CHF23.96
▼	Storage PD Capacity in Netherlands	AE8C-46C3-4994	110.43 gibibyte month	CHF4.13
●	Prometheus Samples Ingested	A4E4-DF03-CD86	33,221,907 Samples	CHF1.69
★	Regional Kubernetes Clusters	B561-8FBD-1264	8.91 hour	CHF0.76
■	Network Intelligence Center Network Analyzer Resource Hours	9BF8-CD36-F9B8	617 count	CHF0.58
■	Network Intelligence Center Topology and Google Cloud Performance Resource Hours	D9AD-28F8-05D8	617 count	CHF0.58
+	Network Intelligence Center Internet to Google Cloud Performance Resource Hours	BDBA-22FA-3925	617 count	CHF0.42
■	Network Internet Data Transfer Out from Netherlands to EMEA	EF9D-7A2D-D06F	1.22 gibibyte	CHF0.12
▲	Standard Storage Netherlands	89D8-0CF9-9F2E	4.7 gibibyte month	CHF0.08
▼	Networking Private Service Connect Partner Select End Point	6882-8FBD-87E5	8.82 hour	CHF0.07
◆	Storage Image	DAA2-253C-6680	0.35 gibibyte month	CHF0.01



An example job of 33M events, MuonEG (1.1 TB)  
 90-node e2-standard-4 regional cluster  
 (4 vCPUs, 16GB memory / node) - auto-scale  
 80 CHF - 9 hours

! ➔ Add data download ➔ !  
 ! ➔ ~ 40\$ / 350GB ➔ !



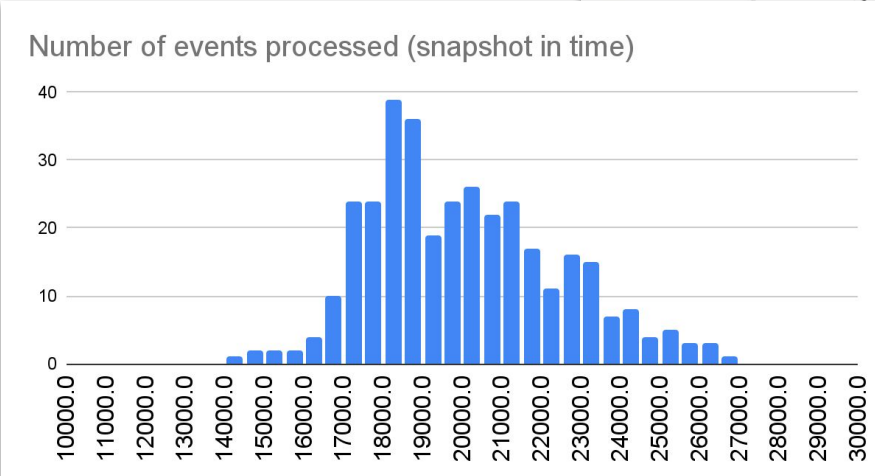


# Auto-scaling - input files and events are not equal!

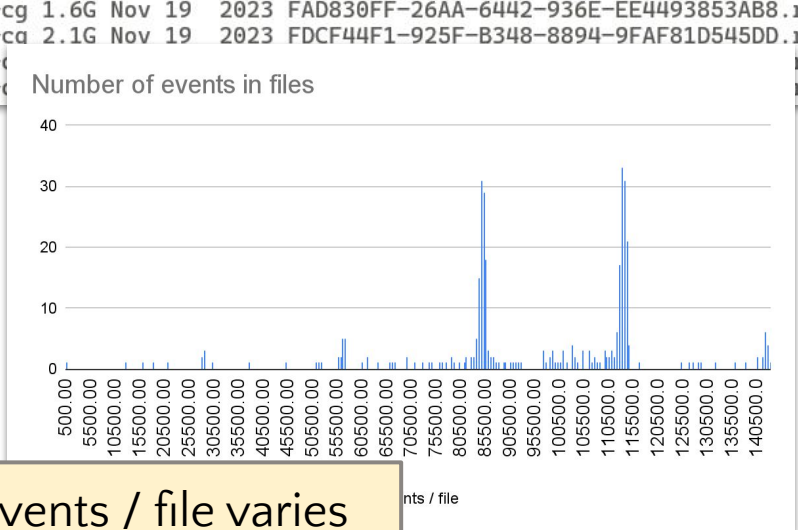
Jobs are not equal in time!

```
W-r--r--. 1 cmsrucio def-cg 3.6G Nov 19 2023 17718268-8398-884B-9088-48B20B18A339.root
W-r--r--. 1 cmsrucio def-cg 3.8G Nov 19 2023 17718268-8398-884B-9088-48B20B18A339.root
W-r--r--. 1 cmsrucio def-cg 3.3G Nov 19 2023 17718268-8398-884B-9088-48B20B18A339.root
W-r--r--. 1 cmsrucio def-cg 2.2G Nov 19 2023 F8EFD4D5-16FB-3B42-9F1B-8EB341817D0C.root
-rw-r--r--. 1 cmsrucio def-cg 3.5G Nov 19 2023 F9EA793C-A6FE-2047-9BB3-DBA4346B6EC1.root
-rw-r--r--. 1 cmsrucio def-cg 3.2G Nov 19 2023 FA08E11D-A7CF-6545-B539-270E5E17EC85.root
def-cg 1.6G Nov 19 2023 FAD830FF-26AA-6442-936E-EE4493853AB8.root
def-ca 2.1G Nov 19 2023 FDCF44F1-925F-B348-8894-9FAF81D545DD.root
def-
def-
```

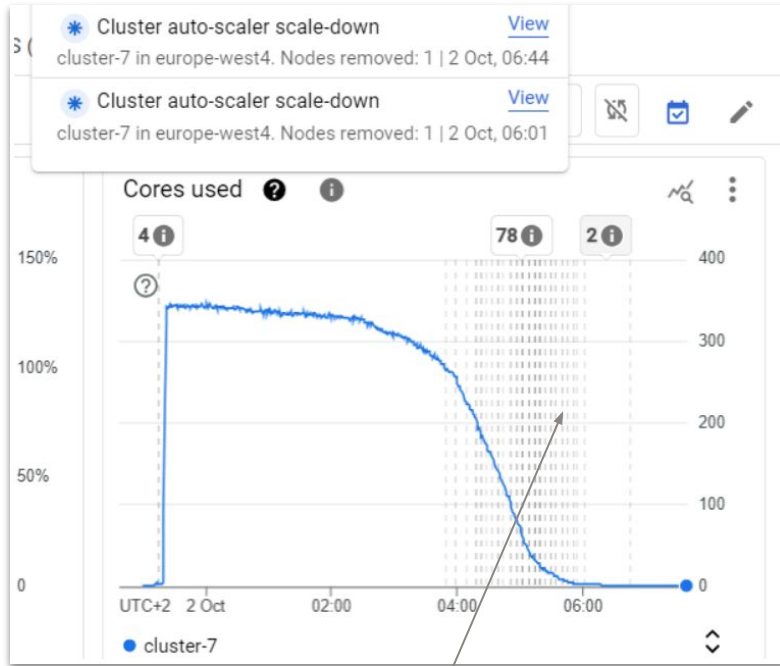
File size varies



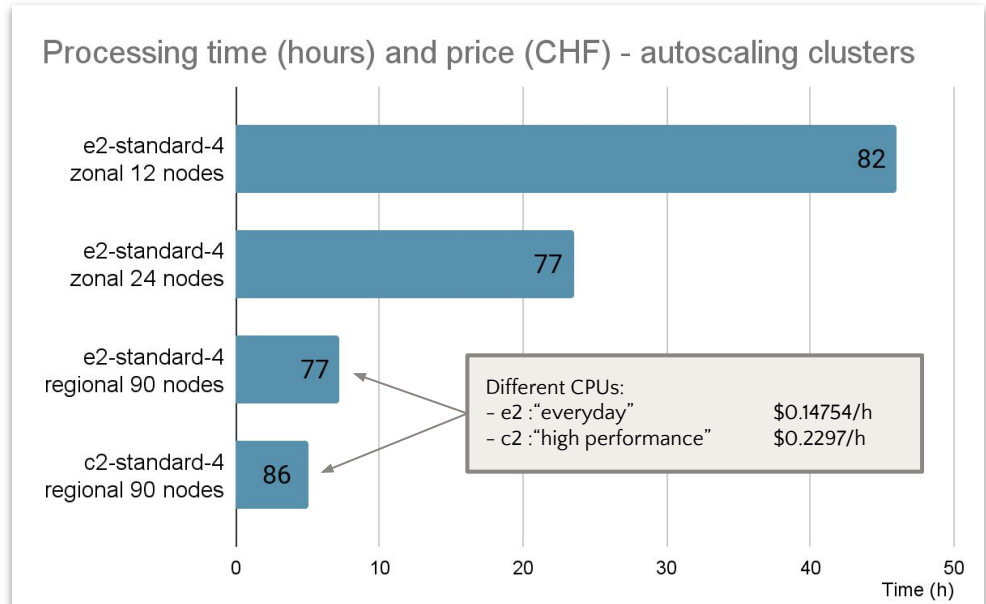
Time / event varies (within a dataset)



Events / file varies



Auto-scaling at work:  
cost = N nodes x time



Optimal configuration:  
a large cluster for a short time

# 80\$ / 1TB

Approximate price of PFNano-type processing / 1 TB of input data

# 7 h / 1 TB

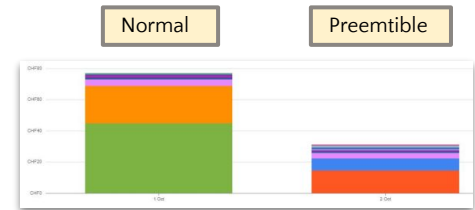
Time to process 1 TB

# 40\$

Example download (PFNano-type content - 35% of the original size)



- Did the jobs fail and why?
  - Some XRootD timeouts (fixed on server-side), rare cluster networking timeouts.
- Can we gain in speed by uploading input files to the cluster?
  - No, not for this workflow. → See [back-up](#)
- How to handle the big CMSSW container image?
  - Use a secondary boot disk. → See [back-up](#)
- Is there an overhead for Kubernetes and Argo services?
  - Nothing significant → See [back-up](#)
- What about spot / preemptible nodes (cheap but deletable)?
  - Cheaper but unreliable, definitely worth a follow-up → See [back-up](#)



- I/O can be expensive
  - Message logger every event increases run time 1-2 %.
- Mounted disk on kubernetes pods is shared to the persistent storage at the end of the step
  - Internal networking in the cluster can be surprisingly slow.
- Multithreading in kubernetes clusters is not obvious
  - CMSSW jobs can be configured to run in parallel threads within a job.
  - But: in a cluster, 4 single-thread jobs go faster than a 4-thread job in a 4 vCPU node.

## FAQ and other observations

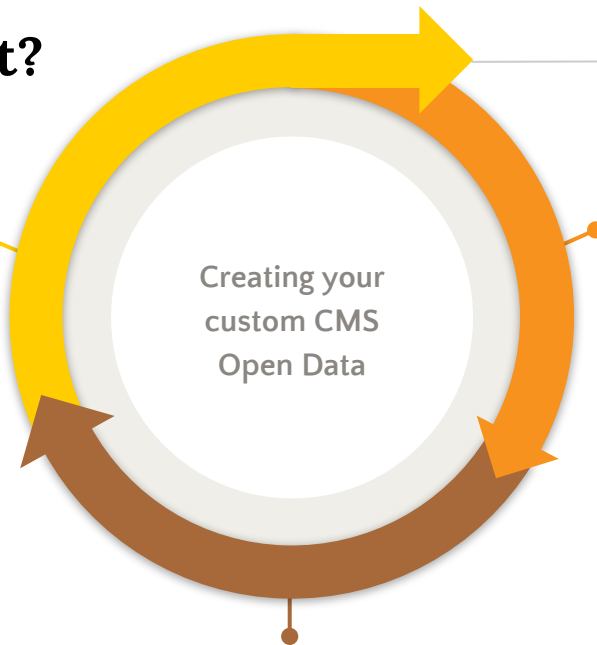




## Is this difficult?

ADAPT FOR YOUR USE and OFF YOU GO!

Run a test job with *your* processing code.  
Evaluate *your* output file size and adapt the  
disk size accordingly.  
Adapt the cluster size and type or use our  
suggested values.



KNOW YOUR PHYSICS and CMS OD

Have your research idea!  
Learn about CMS OD:  
- Workshops / docs / support  
Need more than NanoAOD?  
- No? You could have skipped this talk.  
- Yes? Adapt the example to your needs.

GET STARTED WITH OUR INSTRUCTIONS and EXAMPLES

Create a Google Cloud project. / Install: terraform, gcloud (or use Google Cloud shell), argo CLI, kubectl.  
Deploy resources using example Terraform scripts.  
Run our example job with Argo workflows.



## Conclusion and outlook

---

**CMS Open Data can be used without complications, NanoAOD format is a streamlined and condensed storage format that can be analyzed directly by open data users.**

For analyses requiring detailed event content, we have demonstrated that using public cloud resources for custom NanoAOD processing is feasible both for time and cost.

We have shown how to optimize disposable cloud resources for a typical processing task.

Containerized CMS Open Data workflows can easily be run in a modern kubernetes environment.



---

# Thank you!

## Questions?

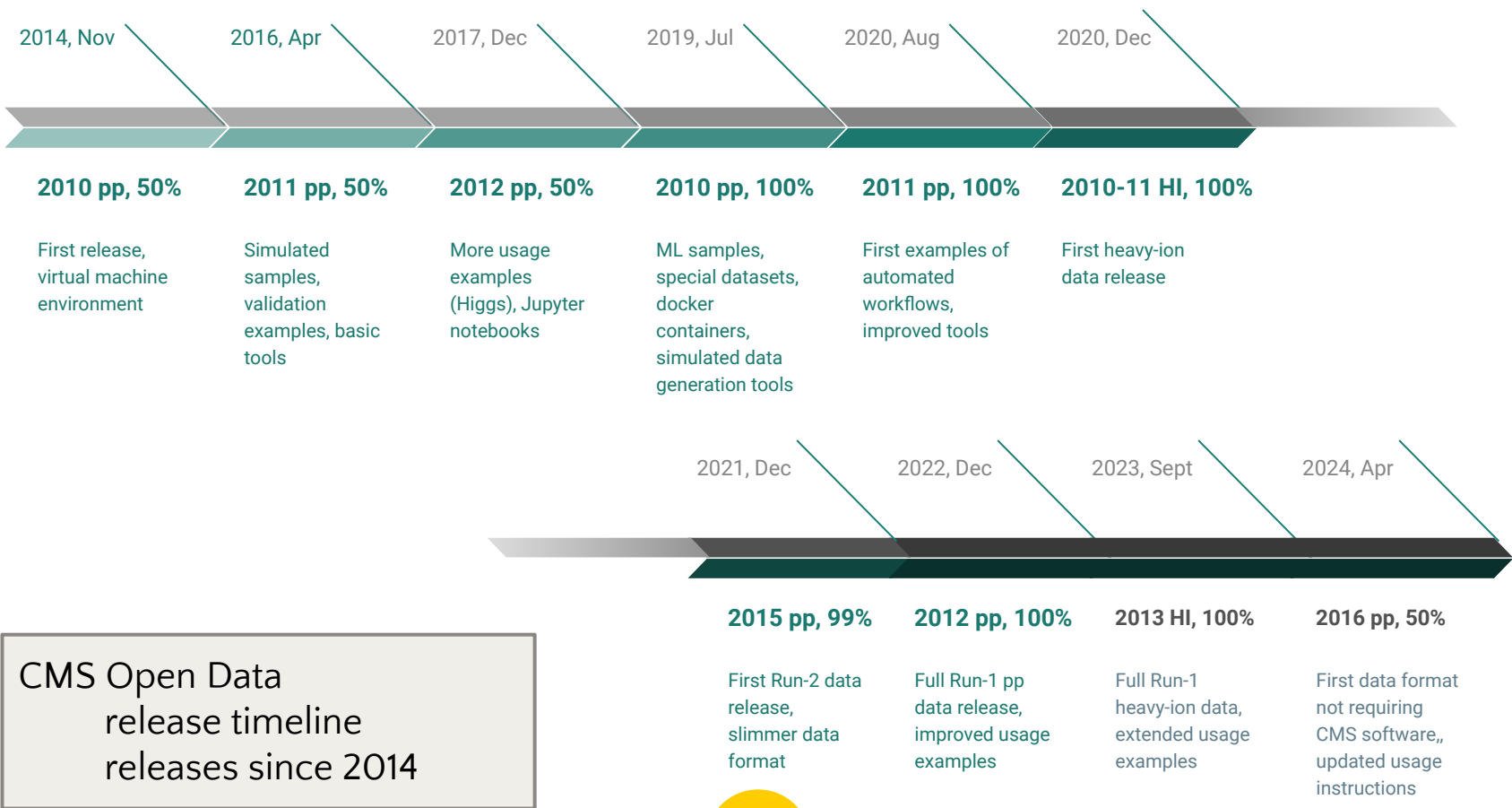
And thanks to [SlidesCarnival](#) for this free presentation template



## **Back-up**

---

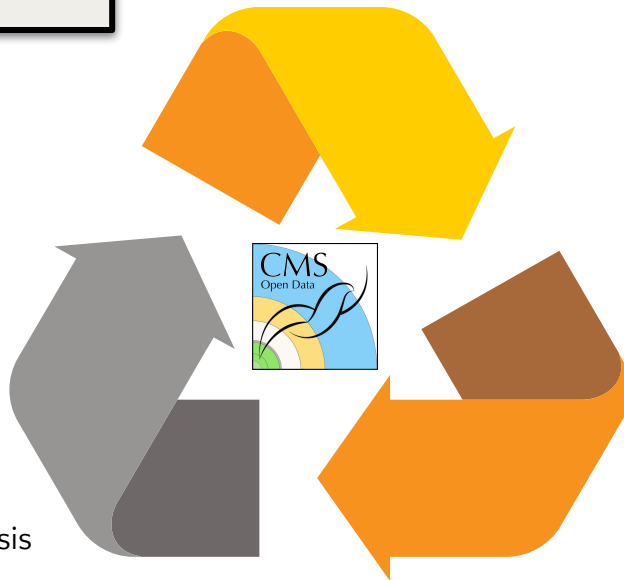




**CMS open data:  
full research-level data  
- not an “open-data” reduction**

Tools:

- software
- environments
- interfaces



Data:

- collision data
- simulations
- additional data for analysis

Knowledge:

- instructions
- actionable examples
- understanding of experimental data





## Input data streaming vs upload

- Data is streamed with xrootd protocol from eospublic at CERN
  - No significant difference between locations close to or far from CERN
  - Processing time dominates over data access time.
- Any faster if input data uploaded?
  - Uploaded files to the container local disk before processing
    - Fairly fast upload with xrscp (upload faster close to CERN)
  - But:
    - No speed-up for the processing time (even slightly slower from local files)
    - Explored differences with file:<filepath> (normal local file access) and root://<xrootdserver>/<filepath> (local xrootd server in the container): no significant difference
    - xrootd server version on eospublic more recent than in the container
- No significant gain.

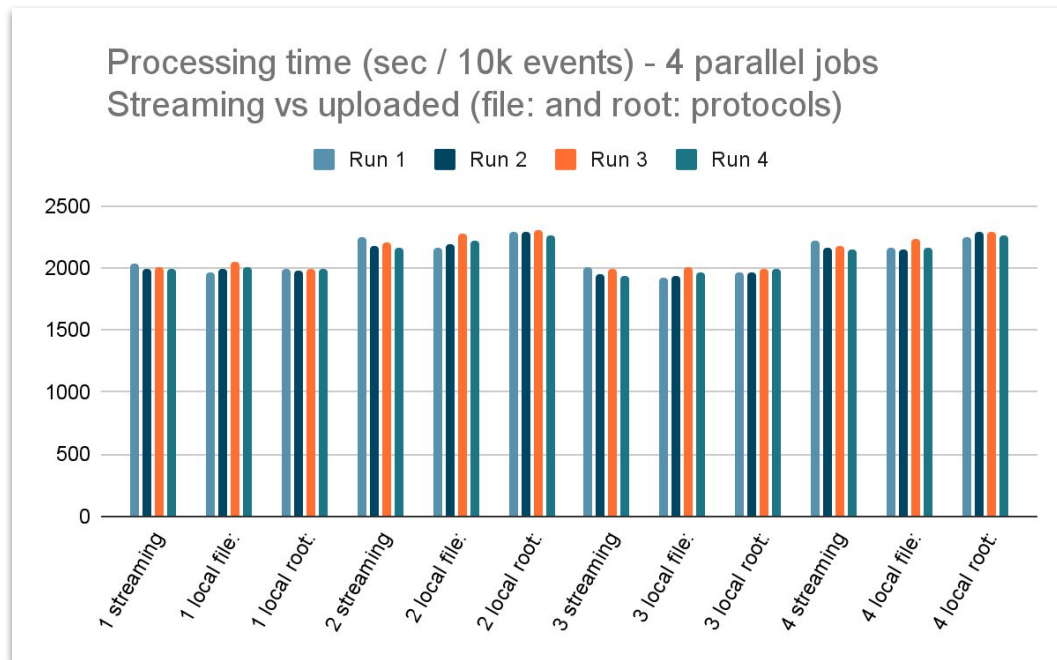


## Processing time: streaming vs local

Compare:

- same 4 jobs
- several times (Run 1-4)
- streaming
- local (read with file;)
- local (read with root:)

No gain observed.





## Container image access

---

- ◉ CMSSW container image is big
  - Initial pull can take 30 mins.
  - Once pulled on the node, it is available to all pods.
  - Run a start job to pull the image to each node.
- ◉ Is there a better solution?
  - Uploading image to Google artifact registry and accessing from there is not significantly faster.
  - Use a new GCP feature: a secondary boot disk with container images preloaded:
    - Build a disk (tools exist), enable image streaming and define the disk in the node pool configuration so that it uses this secondary disk.
    - Immediate start of the jobs 🚀!

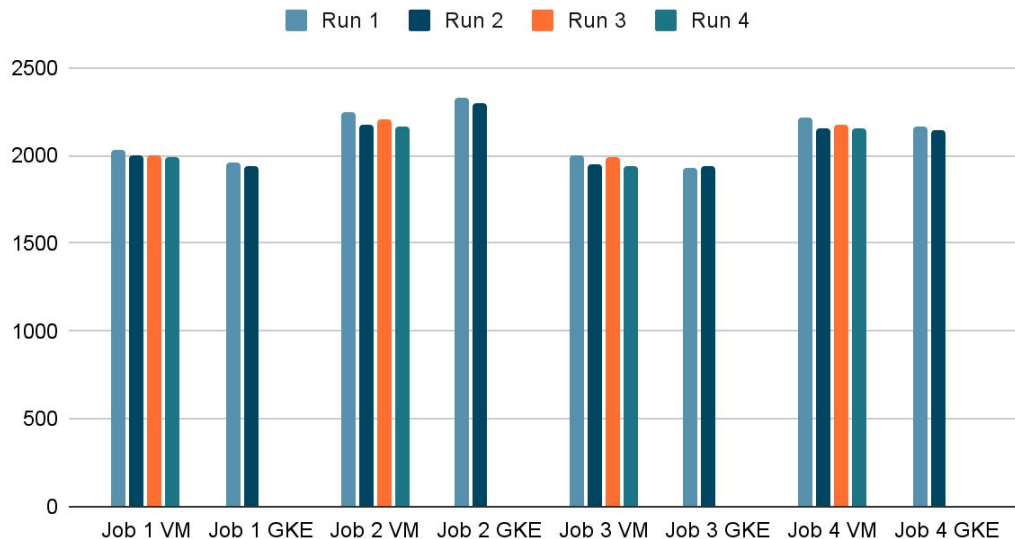


# Kubernetes / argo overhead

NAMESPACE	NAME	CPU(cores)	MEMORY(bytes)
argo	argo-server-5f7b589d6f-6jqcm	1m	18Mi
argo	workflow-controller-864c88655d-gwjcr	4m	26Mi
gke-managed-cim	kube-state-metrics-0	1m	72Mi
gmp-system	collector-2j5dw	5m	139Mi
[...]			
gmp-system	gmp-operator-57874fcf58-8h8w4	2m	62Mi
kube-system	event-exporter-gke-78fb679b7b-fdcra	1m	83Mi
kube-system	fluentbit-gke-28sjh	12m	109Mi
[...]			
kube-system	gke-metrics-agent-x7rft	7m	113Mi
kube-system	connectivity-agent-5f967456fc-2mc7x	1m	40Mi
[...]			
kube-system	connectivity-agent-autoscaler-897d4f648-t97gb	1m	39Mi
kube-system	kube-dns-5fc99b87cb-dgzqj	3m	166Mi
kube-system	kube-dns-5fc99b87cb-nb25q	2m	138Mi
kube-system	kube-dns-autoscaler-6f896b6968-crnrz	1m	39Mi
kube-system	kube-proxy-gke-cluster-4-cluster-4-9248283b-00xm	1m	15Mi
[...]			
kube-system	17-default-backend-6697bb6dfd-swb9h	1m	10Mi
kube-system	metrics-server-v1.30.3-7c8f6576cd-bdc85	4m	78Mi
kube-system	pdcsi-node-2fv5q	5m	47Mi

Compare:  
same jobs  
several times  
VM: no argo, no k8s  
vs GKE cluster  
No significant overhead.

Processing time (sec / 10k events) - 4 parallel jobs VM s GKE





## Preemptible / spot nodes

- Considerably cheaper
  - 1/4 - 1/2 price
- Nodes can be deleted any time
- A trial with a 90-node e2-standard-4 cluster:
  - 13 / 90 nodes terminated
  - 52 / 353 jobs failed
- Requires rerunning of the failed jobs.
  - The price advantage is worth some scripting for automated reruns.

