



When Things Get Big

Scaling Cloud Native Workloads and Software Distribution

Ricardo Rocha, CERN @ahcorporto

<https://kcdzurich.ch/>

Big Machines



CMS

CERN

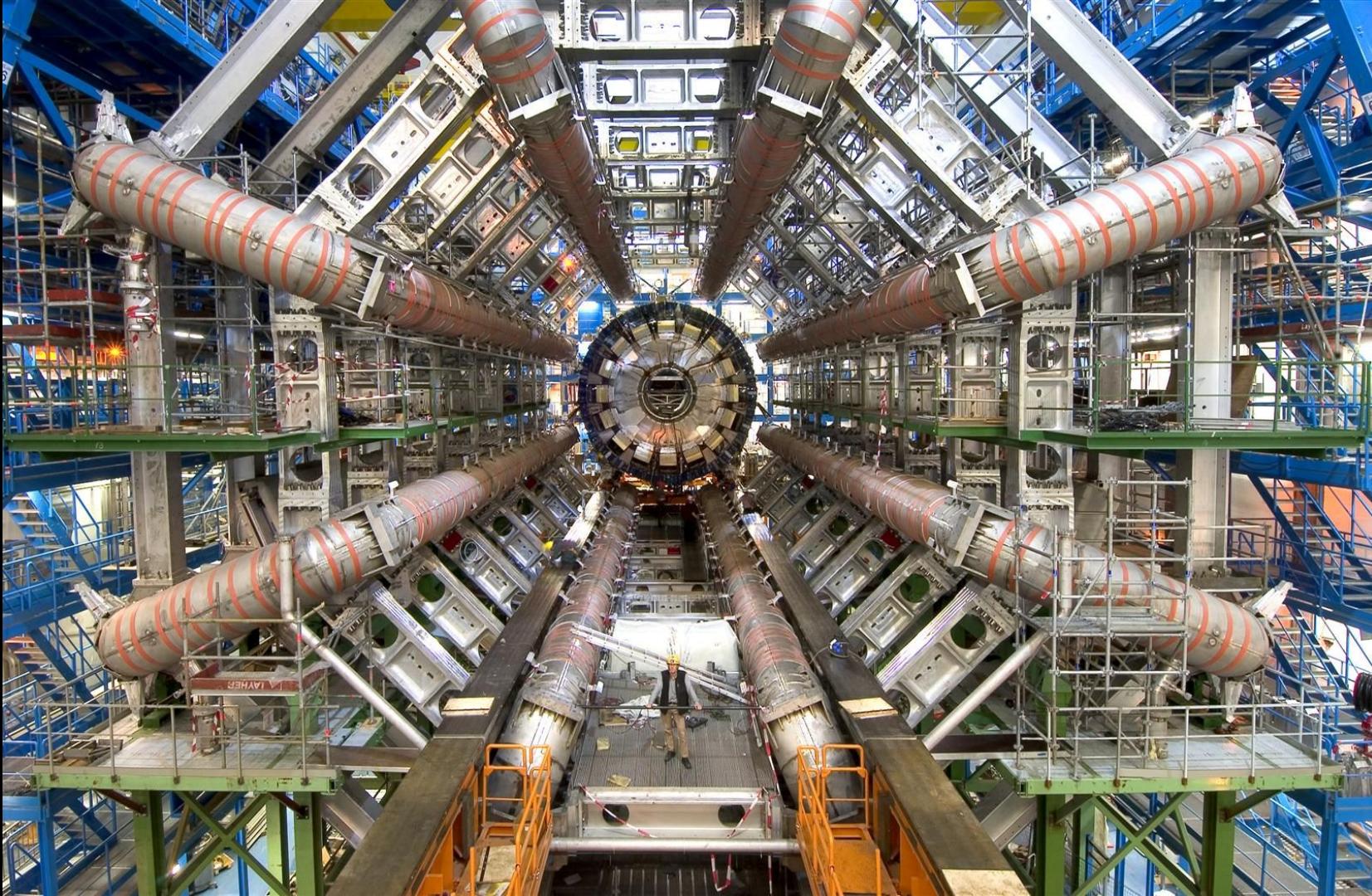
LHC

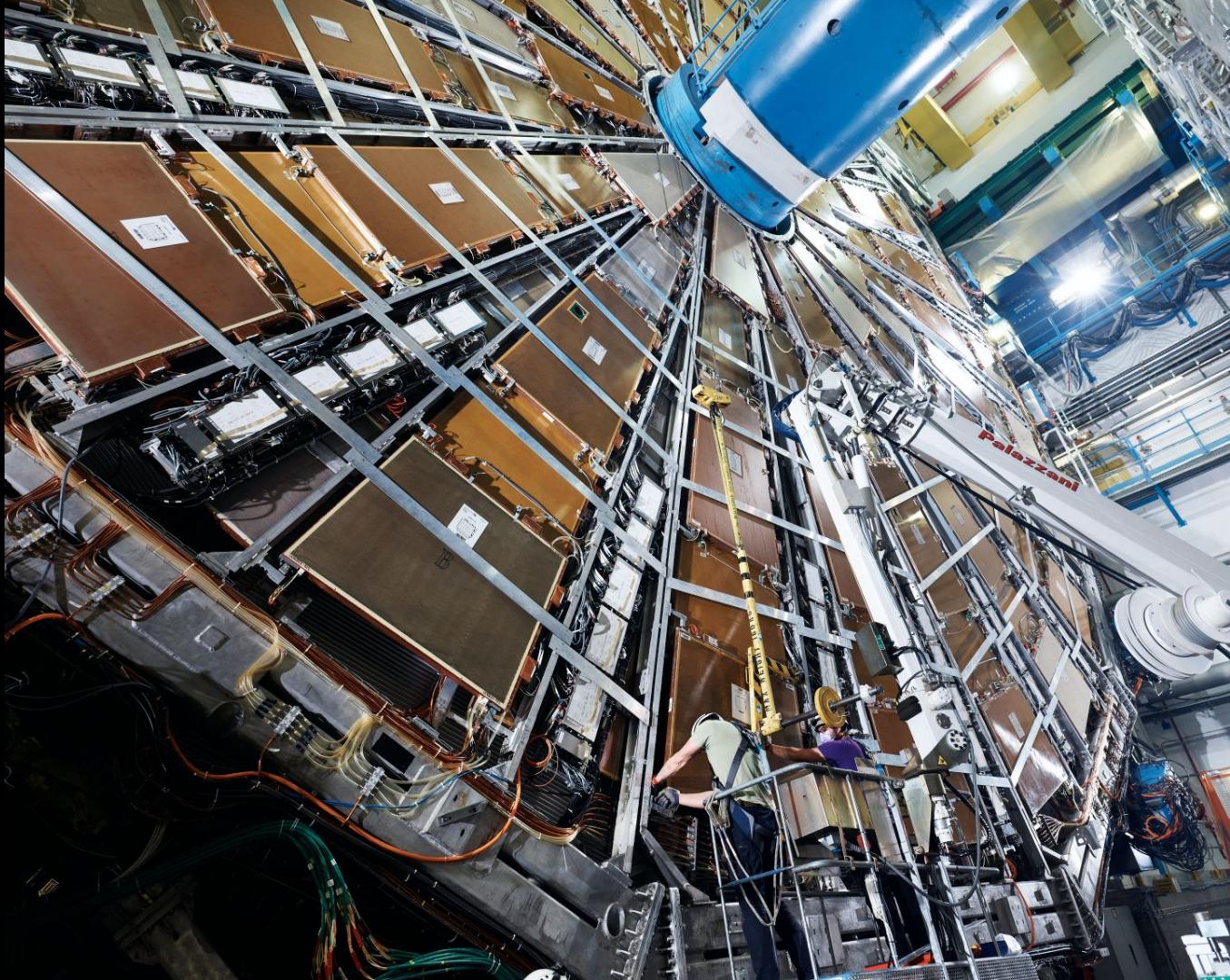
[Large Hadron Collider]

ALICE

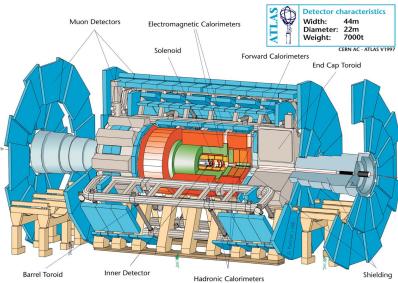
LHCb

ATLAS





Big Applications



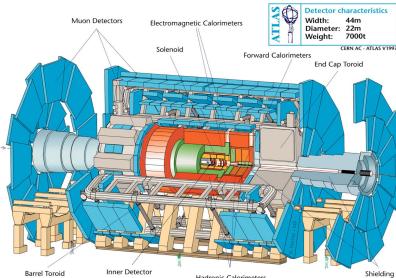
1 PB / sec



ATLAS Event Filter



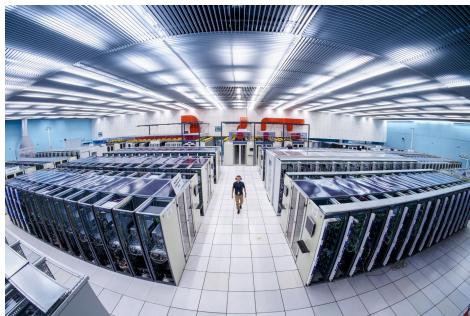
< 10 GB / sec



1 PB / sec



ATLAS Event Filter



< 10 GB / sec

Typically split into Hardware and Software Filters
(this might change too)

40 million particle interactions / second

~3000 multi-core nodes

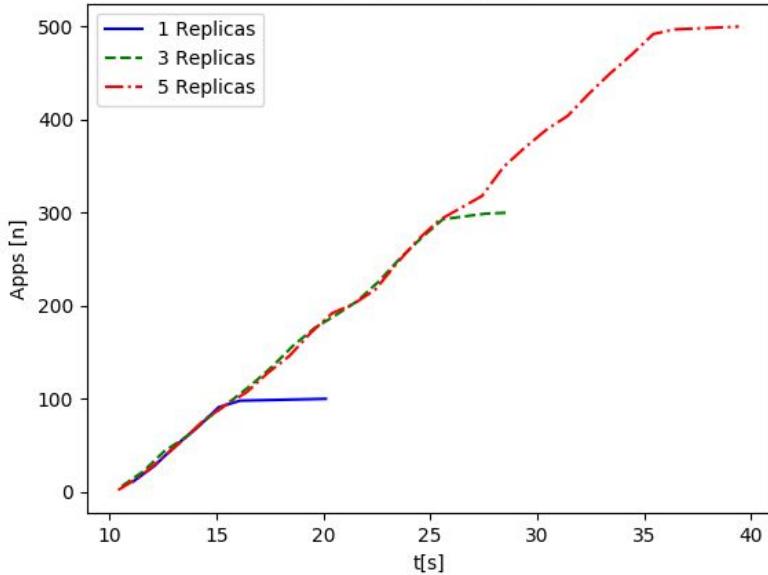
~30.000 applications to supervise

Critical system, sustained failure means data loss

Can it be improved for Run 4?

First study 2017, Mattia Cadeddu, Giuseppe Avolio

Kubernetes 1.5.x



First results... 100 node cluster

40 sec to start **500** applications on 100 nodes

Container launch rate **12Hz**

Not particularly promising ...

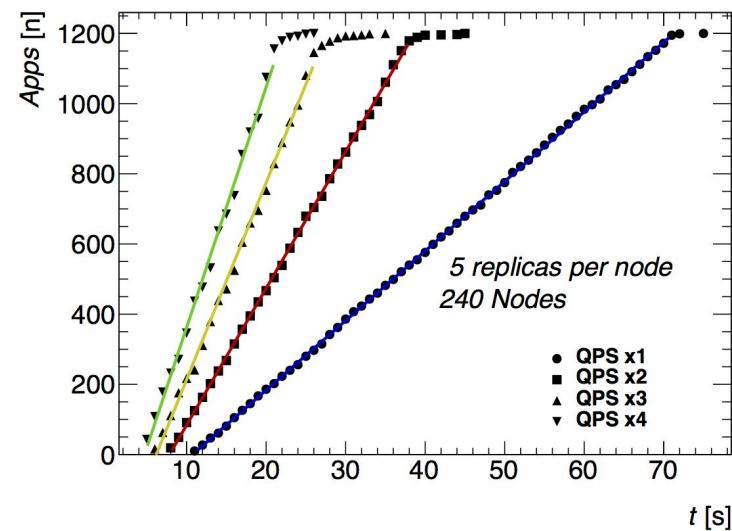
Issues Found

QPS defaults were very conservative

Tested x1, x2, x3, x4 values

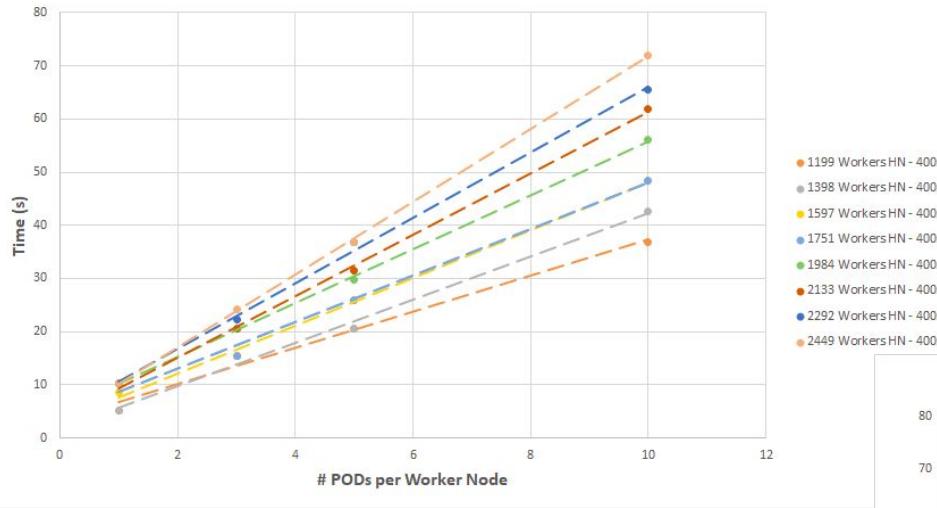
QPS x4 gives 57Hz: > **330 % improvement**

Does Kubernetes prefer larger clusters ?



2022

Start Time

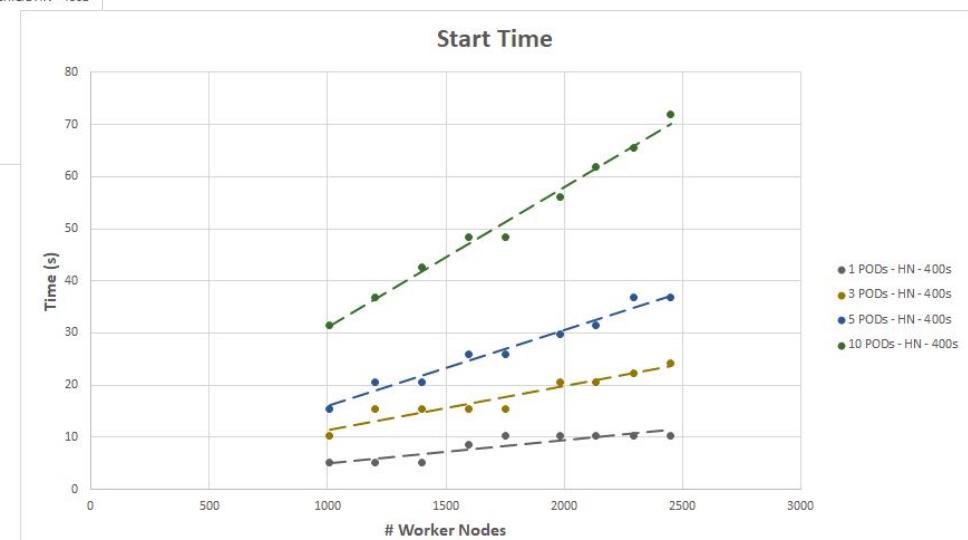


Prediction 5000 Node Cluster

PODs per Host	Start Time (s) [95% CL]	Stop Time (s) [95% CL]
1	[13.8 – 29.1]	[18.7 – 26.2]
3	[39.0 – 48.4]	[44.4 – 51.8]
5	[75.1 – 80.3]	[81.1 – 82.9]

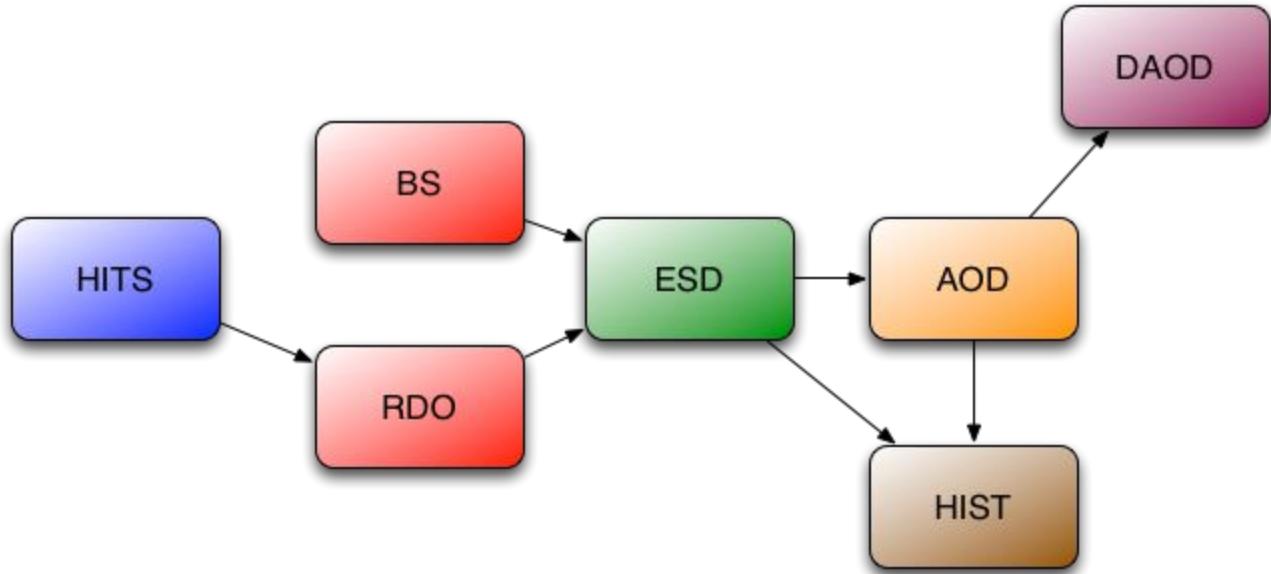
Kubernetes 1.21.2

2500 Nodes



Credit: Giuseppe Avolio, ATLAS

Big Workloads



11/22/2013 5:55:18 p.m.

Running jobs: 244151
Transfer rate: 40.08 GiB/sec



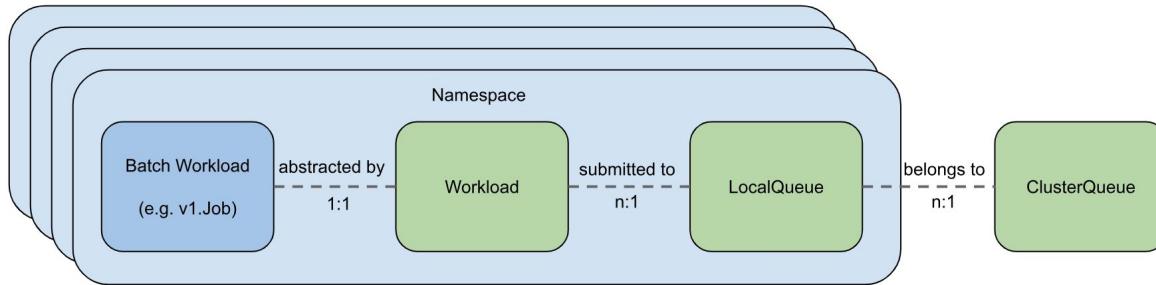
WLCG

US Dept of State Geographer
© 2013 Google
Data SIO, NOAA, U.S. Navy, NGA, GEBCO
Image Landsat

Google earth

Fecha de las imágenes: 4/10/2013 66°43'28,18" N 8°52'37,10" O alt. ojo 16085.50 km

Workloads, not just Pods

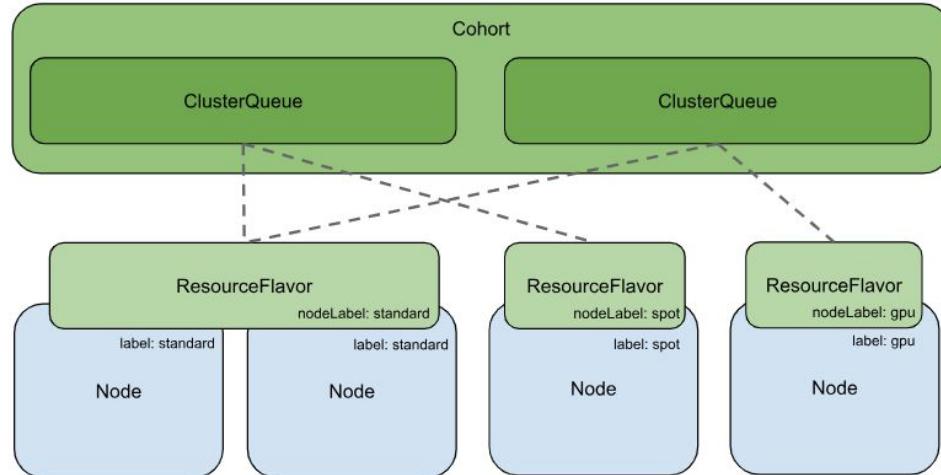


<https://kueue.sigs.k8s.io/>

Workloads, not just Pods

Queues and Priorities, not all workloads are equal

Fair Sharing, resource usage optimization



Workloads, not just Pods

Queues and Priorities, not all workloads are equal

Fair Sharing, resource usage optimization

Gang Scheduling, Array Jobs



Kubernetes
AI + HPC DAY
NORTH AMERICA

NOVEMBER 6, 2023

CHICAGO, IL
#K8SAIHPCDAY

Kubernetes Batch WG

<https://github.com/kubernetes/community/tree/master/wg-batch>

CNCF Batch System Initiative

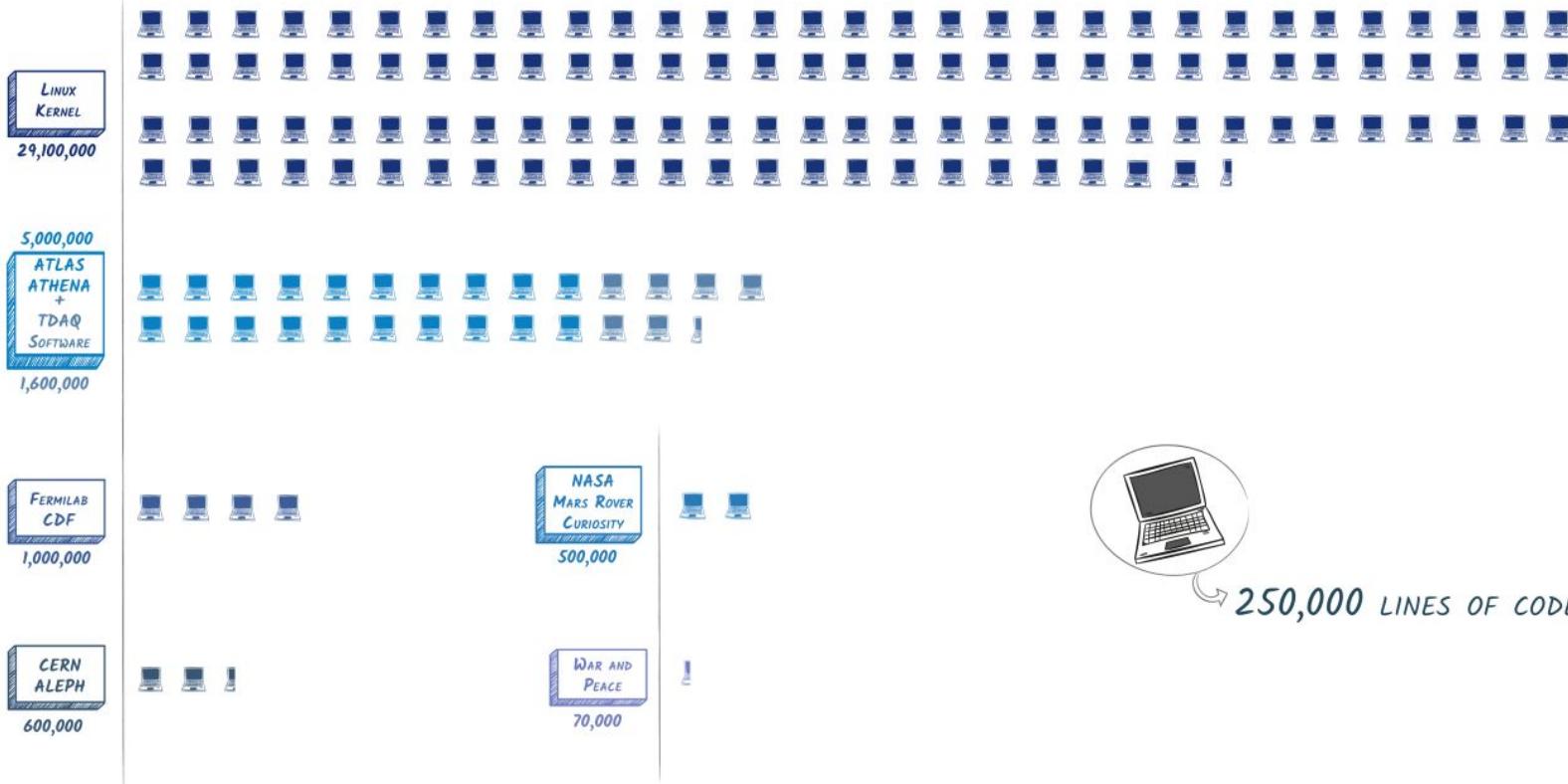
<https://github.com/cncf/tag-runtime/issues/38>

CNCF Research User Group

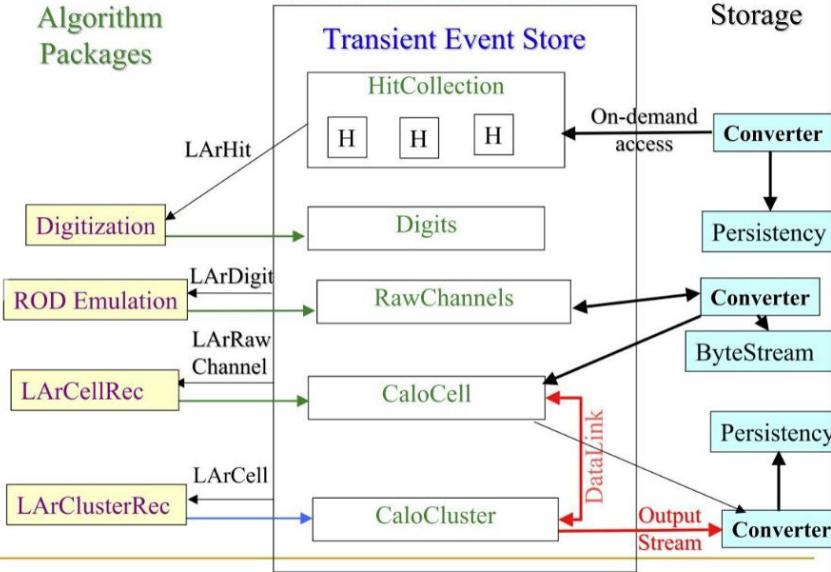
<https://community.cncf.io/research-end-user-group/>

Big Software

COUNTING LINES OF CODE



LAr Event Data Model



Aug 27, 2003

Hong Ma, Athena Tutorial

6

Releases

- **Release:** A complete set of code compiled, built and frozen for use
 - Defined by a set of tags for each package,
 - [Tag collector](#)
 - Nightly builds
 - Used by developers to work towards a release.
 - Not guaranteed to work
 - Developer Release: approximately every 3 weeks (currently 6.6.0)
 - Mostly working
 - Major release: approximately every 6 months (currently)
 - Major milestones, such as Data Challenges
- [Release plan](#) and [Release Status](#) can be found on the “Software Development” page
- USATLAS builds same releases as CERN
 - [/afs/usatlas.bnl.gov/software/dist/](http://afs/usatlas.bnl.gov/software/dist/)

	compressed	uncompressed
alpine	3MB	6MB
ubuntu:22.04	30MB	78MB

	compressed	uncompressed
alpine	3MB	6MB
ubuntu:22.04	30MB	78MB
admutils	870MB	1.2GB

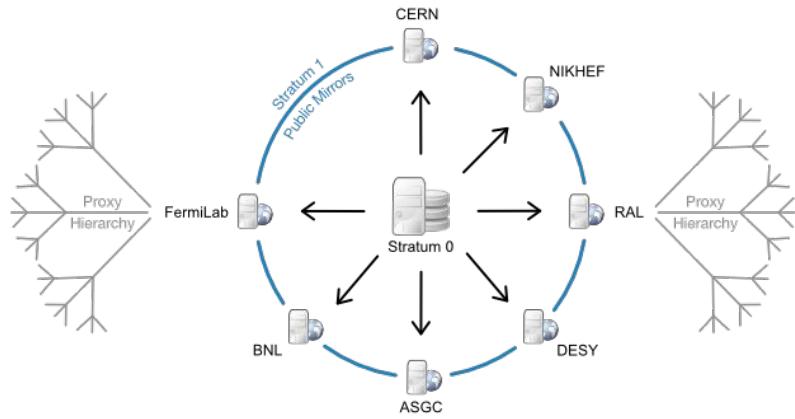
	compressed	uncompressed
alpine	3MB	6MB
ubuntu:22.04	30MB	78MB
admutils	870MB	1.2GB
atlas/athena	18.7GB	61GB

Why the heck is this thing 24 GB?!?!

```
docker images us-docker.pkg.dev/colab-images/public/runtime:latest
REPOSITORY                                TAG      IMAGE ID      CREATED        SIZE
us-docker.pkg.dev/colab-images/public/runtime    latest    01b91234d65c   2 weeks ago   24.1GB
```

@Ricardo Rocha ^ here's another one for your "large images running in production" list. 😊

	compressed	uncompressed
alpine	3MB	6MB
ubuntu:22.04	30MB	78MB
admutils	870MB	1.2GB
atlas/athena	18.7GB	61GB
<??>	76GB	124GB



11:00 → 11:30 Containerd Remote Snapshotter

11:30 → 12:15 Discussion: Containerd Remote Snapshotter

Container SW Distribution

16 May 2019, 12:00 → 17 May 2019, 23:10 Europe/Zurich
513/R-070 - Openlab Space (CERN)
Ricardo Brito Da Rocha (CERN)

Registration You are registered for this event. Check details

THURSDAY, 16 MAY

14:00 ~ 15:00 SM18 Visit

FRIDAY, 17 MAY

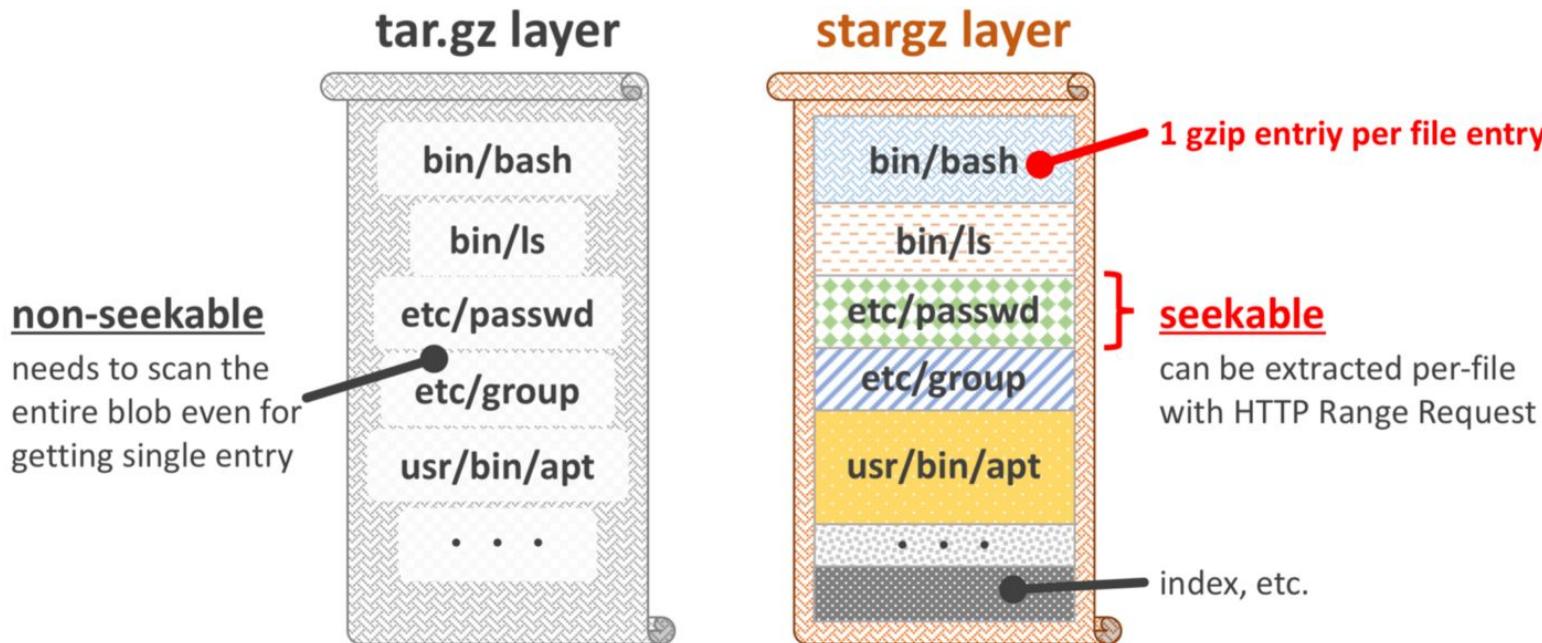
09:00 ~ 09:25 CVMFS and Software Distribution on the Grid
Speaker: Jakob Blomer (CERN)
[cvmfs-containers.pdf](#)

09:25 ~ 09:50 WLCG and the case for userspace containers
Speakers: Alessandra Forti (University of Manchester (GB)), Lukas Alexander Heinrich (CERN)
[ContainersWLDistr.pdf](#)

09:50 ~ 10:15 Overview of Containerd
Speaker: Phil Estes (IBM)
[FOSDEM 2019_con...](#)

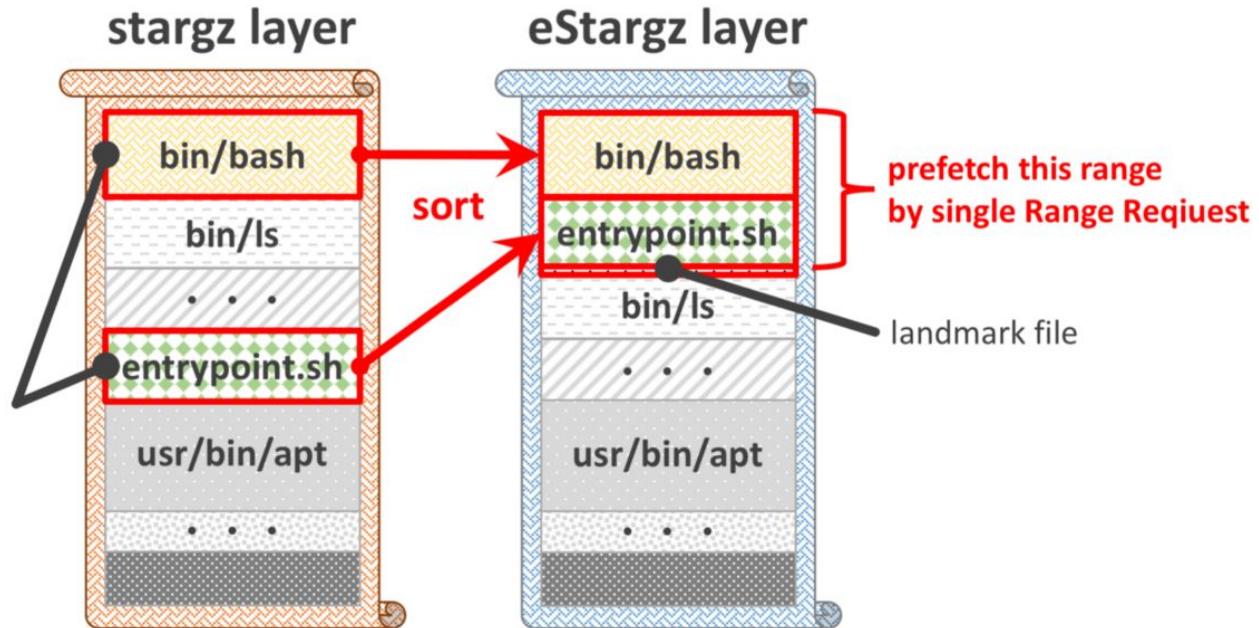
⌚ 30m ⚽ 513/R-070 - Openlab Space

⌚ 45m ⚽ 513/R-070 - Openlab Space

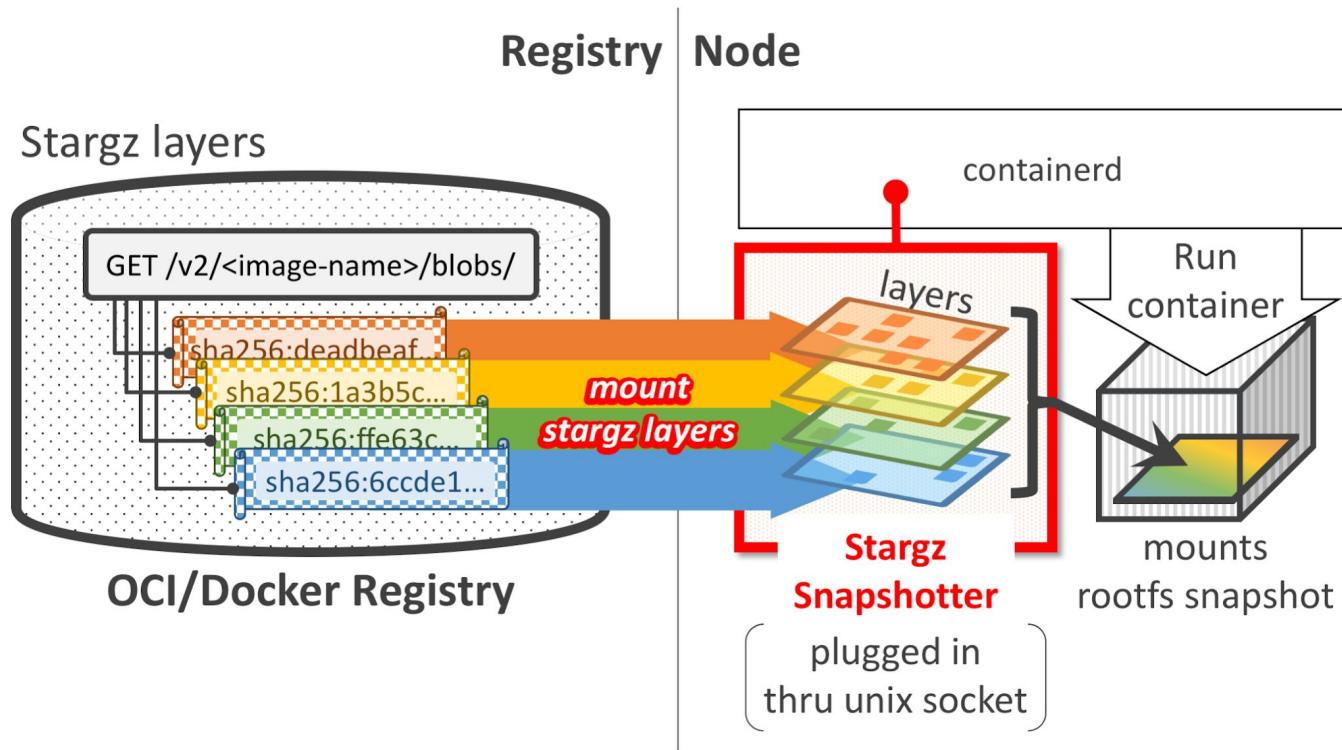


<https://github.com/containerd/stargz-snapshotter>

Prioritised files
most likely accessed
files during runtime



<https://github.com/containerd/stargz-snapshotter>



<https://github.com/containerd/stargz-snapshotter>

Introducing GKE image streaming for fast application startup and autoscaling

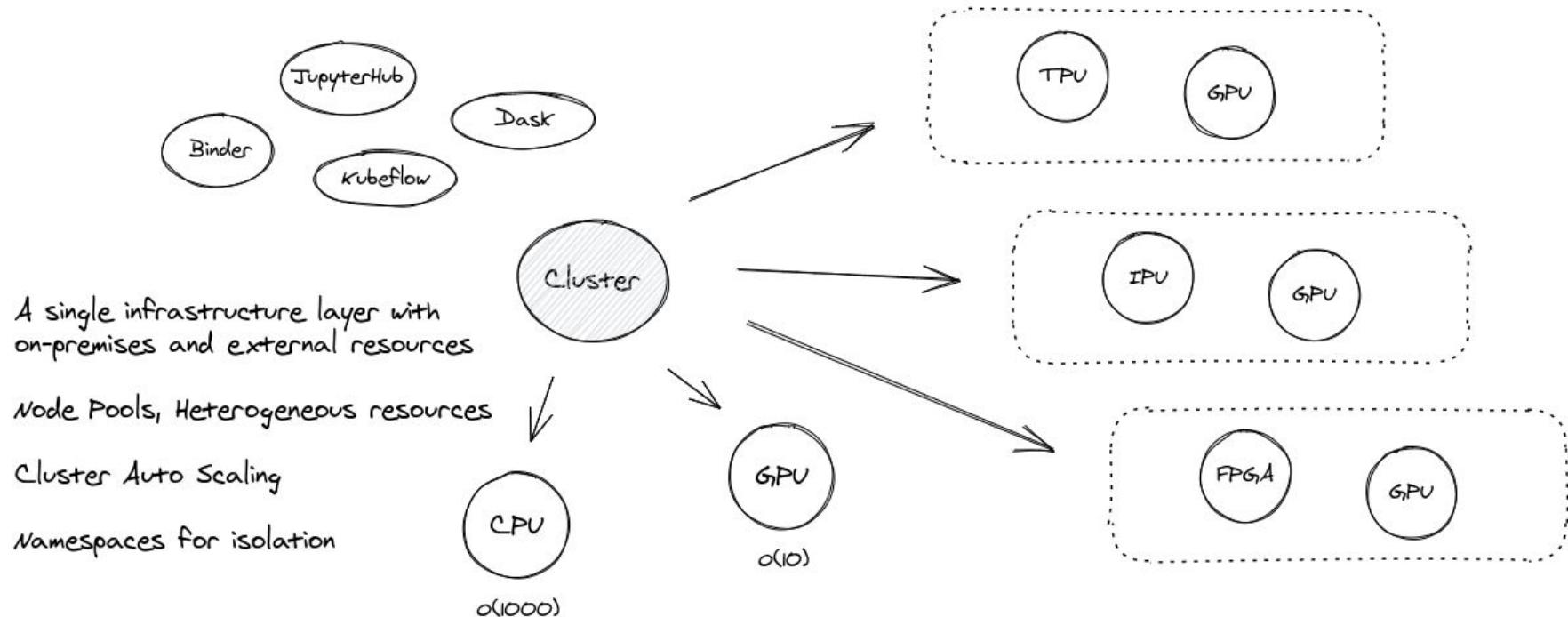
November 4, 2021

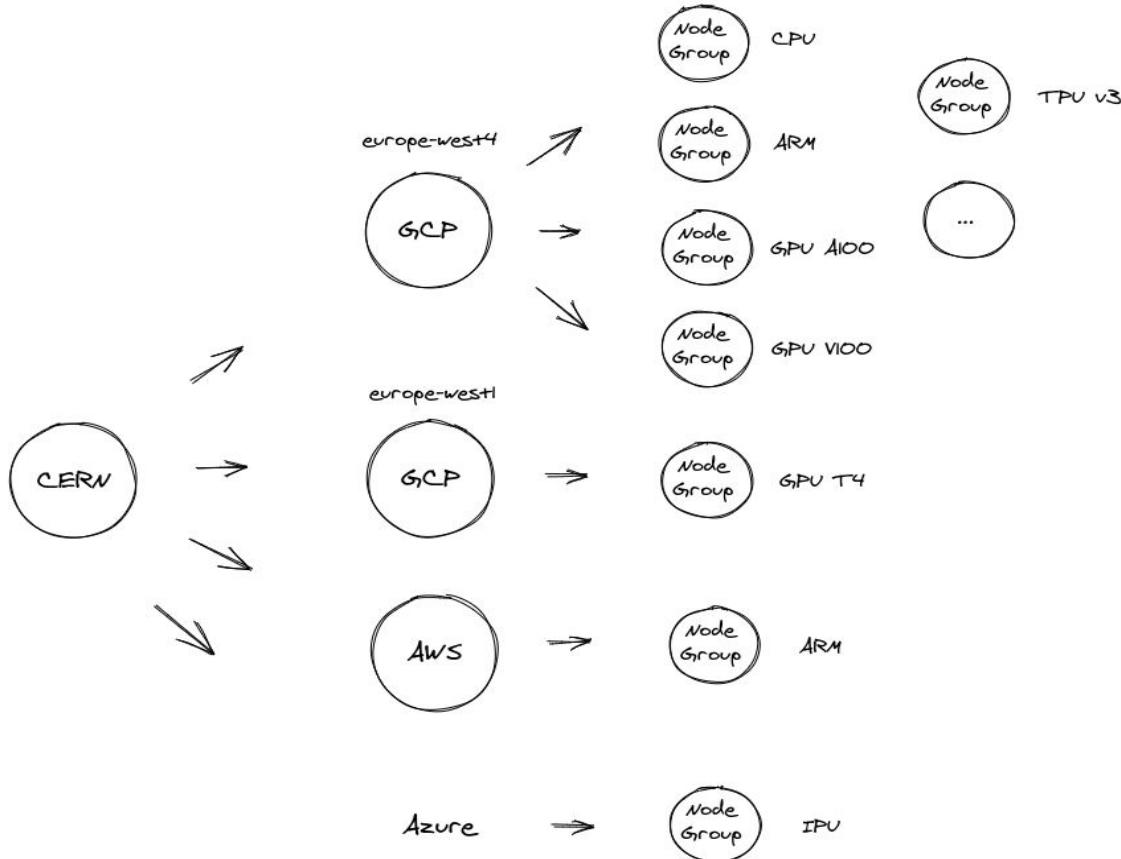
Image streaming is a method of pulling container images in which GKE streams data from eligible images as requested by your applications. You can use Image streaming to allow your workloads to initialize without waiting for the entire image to download, which leads to significant improvements in initialization times. The shortened pull time provides you with benefits including the following:

- Faster autoscaling
 - Reduced latency when pulling large images
 - Faster Pod startup
- You might not notice the benefits of Image streaming during the first pull of an eligible image. However, after Image streaming caches the image, future image pulls on any cluster benefit from Image streaming.

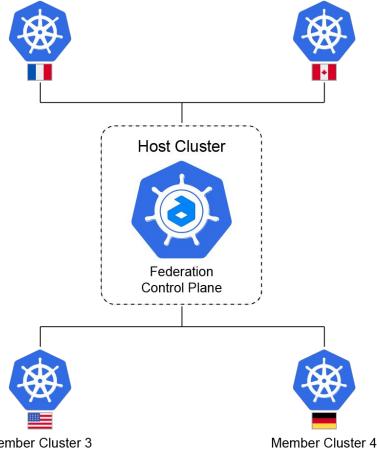
Quick Demo

Big Needs

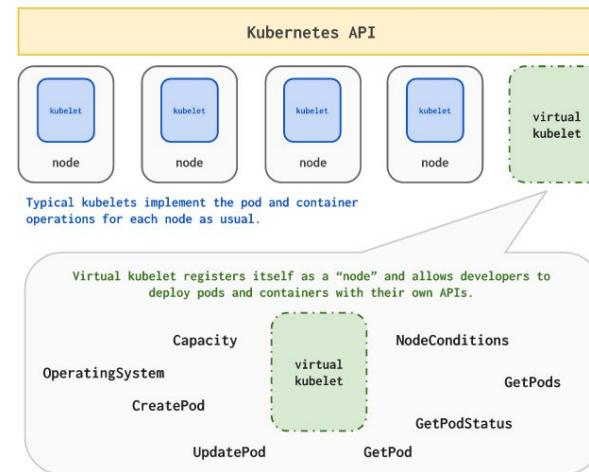




Member Cluster 1

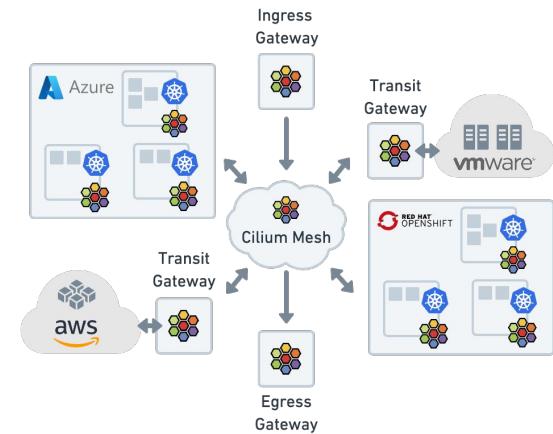
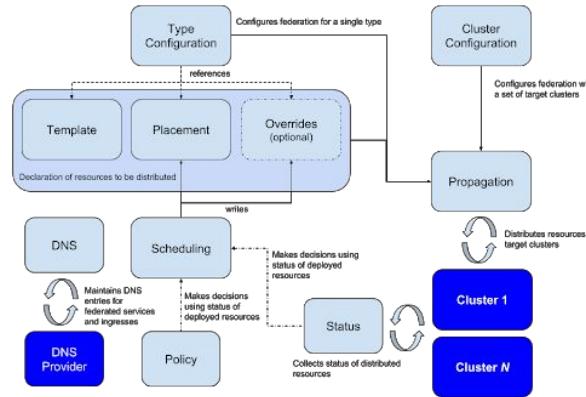


Member Cluster 2



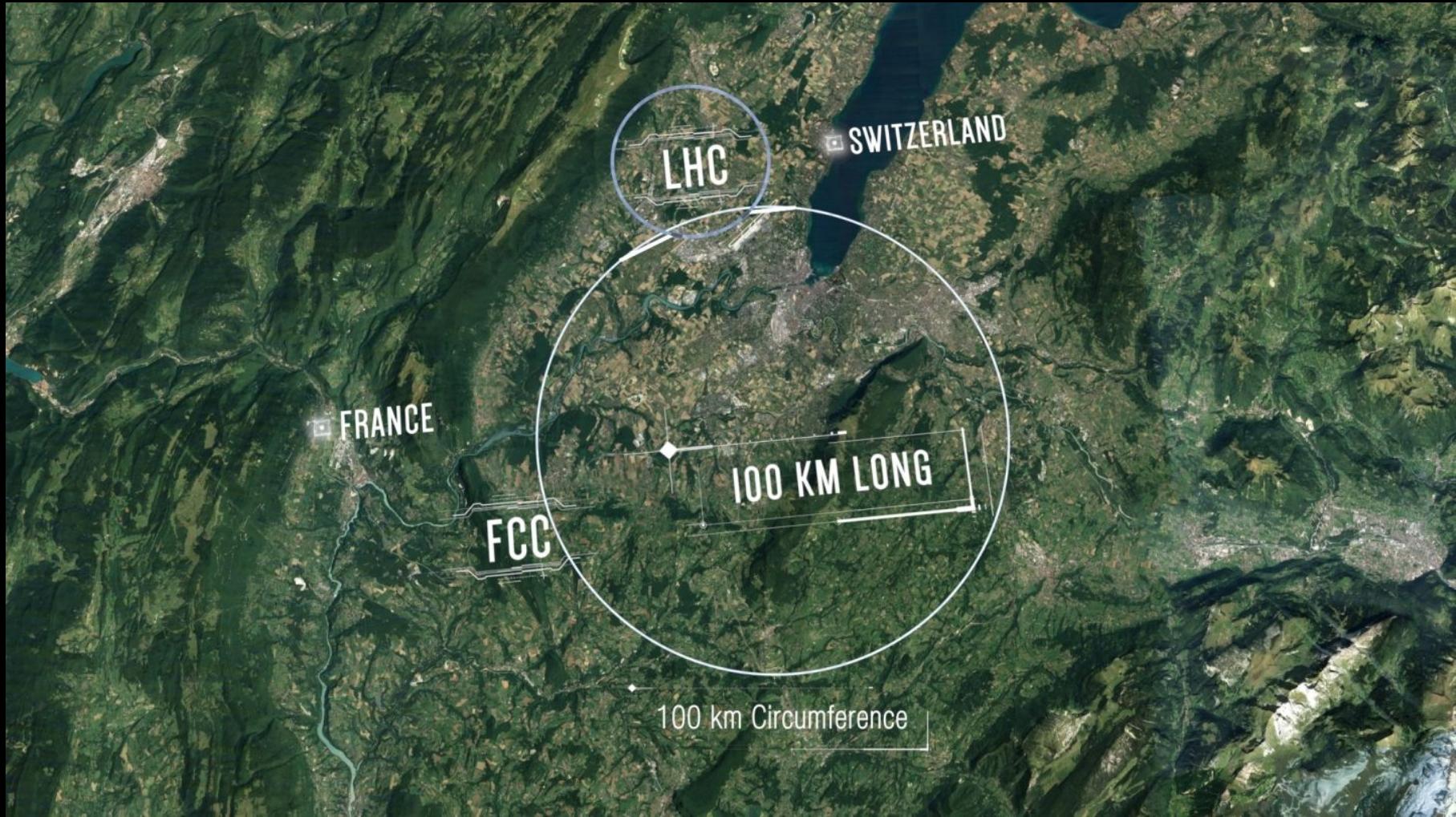
Member Cluster 3

Member Cluster 4



2017

2023



Questions?