# Des péta-octets de données dans Kubernetes, c'est possible!
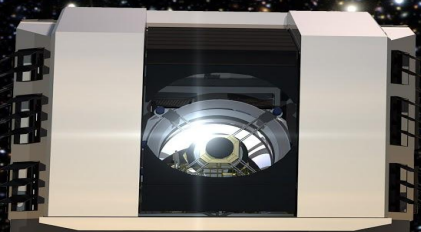
**Fabrice Jammes**

**Scalable Data Systems Expert**
**IN2P3/LSST Corporation**

**Credits:**
**Sabine Elles**
Expert en développement
d'applications
**LAPP**

**Bastien Gounon**
Expert infrastructure Kubernetes
**CC-IN2P3**

# Agenda

1. Large Synoptic Survey Telescope

2. Qserv: LSST Petascale database

3. Benefits of Cloud-Native
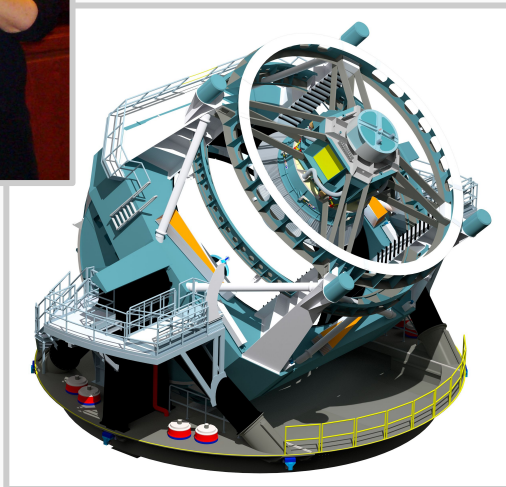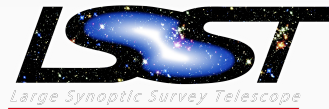
4. On-premise vs Public Cloud

# LSST in short

Large Synoptic Survey Telescope

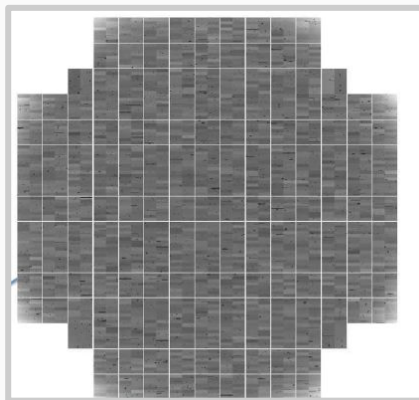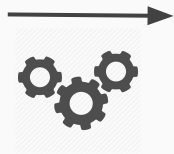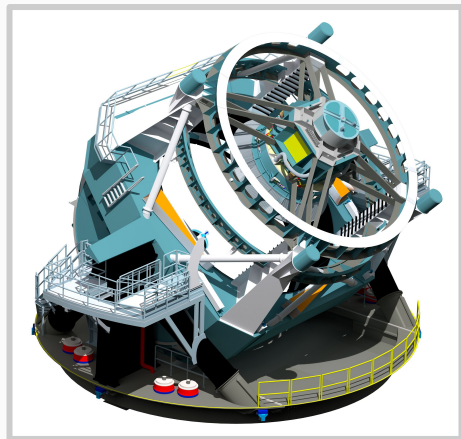Large aperture, wide-field, ground-based survey telescope
**The largest imager ever built for astronomy**

Characteristics
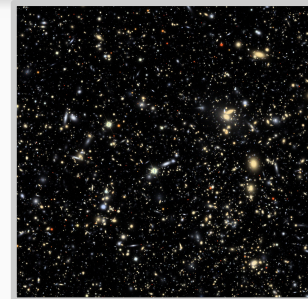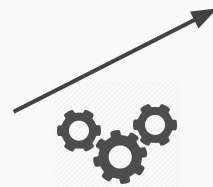
★    All visible sky in 6 bands
★    ~20000□
★    15 seconds exposures, 1 visit/3 days
★    During 10 years!
★    **60 PB of raw data**

# 80+ PB of astronomical catalog


Processed image


Raw data


**Catalog** (stars, galaxies, objects, sources,transients, exposures, etc.)

LSST will build a catalog of 20 billion galaxies and 17 billion stars and their associated physical properties

4

# Data

**Images**
Persisted: **~38 PB**
Temporary: **~½ EB**



★ ~3 million "visits"
★ ~47 billion "objects"
★ ~9 trillion "detections"

★ Largest table: **~5 PB**
★ Tallest table: ~50 trillion rows
★ Total (all data releases, compressed): **~83 PB**

Ad-hoc user-generated data
Rich provenance

# Qserv

The LSST Petascale database

# Who we are

**Database and Data access team**
- ★ 10 engineers at Stanford University + 1 IN2P3
  - ○ *Software development*

**Operations teams**
- ★ 5 sysadmins at NCSA/IN2P3
  - ○ *Large Scale development platforms*
  - ○ *Cloud Native / Kubernetes*
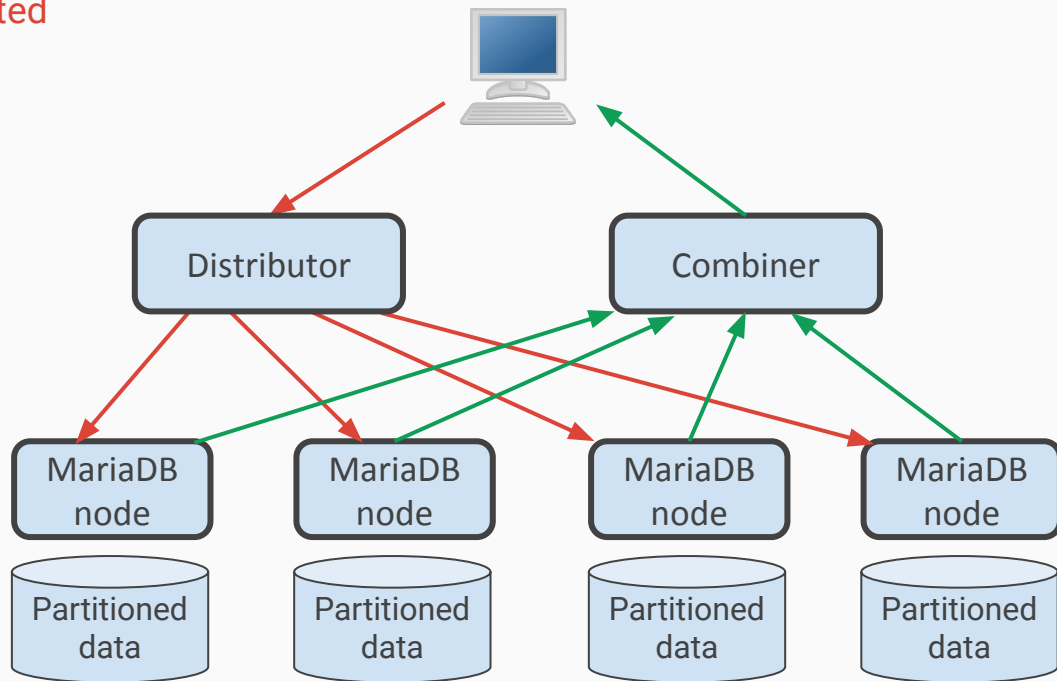  - ○ *System administration, Monitoring*

# Qserv design

Relational database, 100% open source

Spatially-sharded with overlaps
Map/reduce-like processing, highly distributed

# From Cloud-Native to Bare-Metal

## Target for production
~1000 nodes cluster in 2 international Academic data-centers

## Running now
**Development platform (CC-IN2P3)**
*1000 cores, 15 TB memory*
*15 PB storage*
*=> Large scale test: **300 TB synthetized data***
**=> Ingestion of DESC-DC2 data (1 TB)**

**Prototype Data Access Center (NCSA)**
*500 cores, 4 TB memory*
*700 TB storage,*
**=> WISE catalog ("real" dataset)**

dedicated hardware:

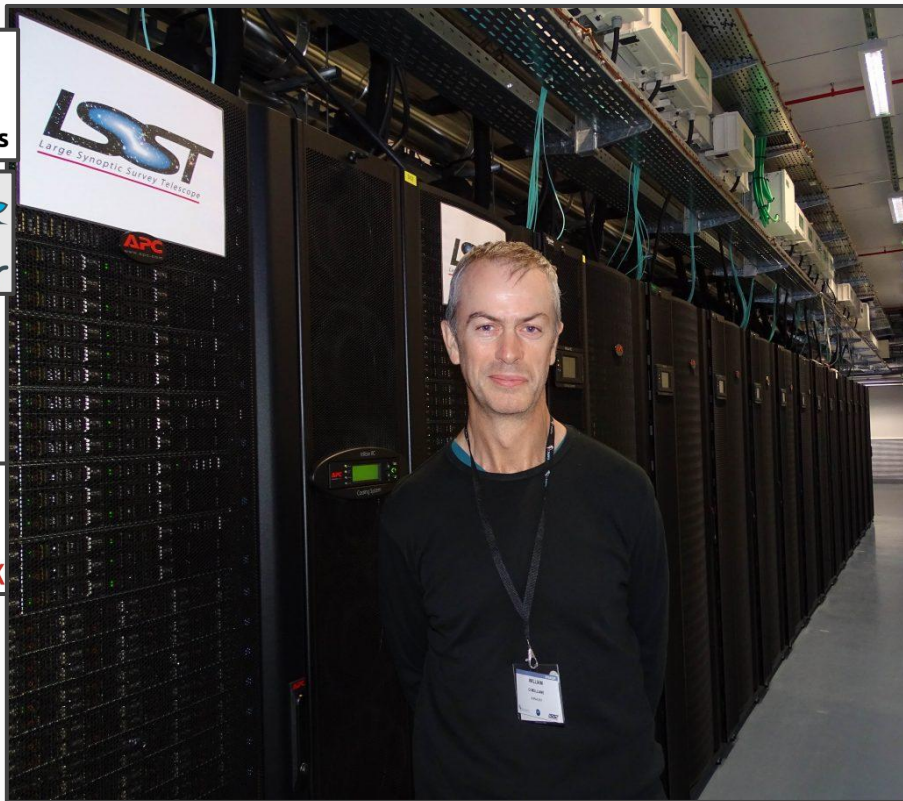3 x Kubernetes masters (40x2.2GHz, 64GB RAM)

     control-plane

2 x Qserv masters (40x2.2GHz, 256GB RAM, 8TB SSD RAID1)

     user interaction, result aggregation

20 x Qserv workers (40x2.2GHz, 256GB RAM, 48TB HDD RAID5)

     database workload and storage

=> 25 nodes Kubernetes cluster (v1.15.3)

deployed via Puppet using puppetlabs/kubernetes plugin

     CRI : containerd

     CNI : weave

     token-based authentication for cluster administration

Qserv Platform @ CC-IN2P3

## ElasticSearch/Grafana activity dashboard

# Benefits of Cloud-Native

# Automated Qserv deployment



Repositories

Workstation

# Automated Qserv deployment



Travis CI

Qserv images

Docker Registry

Repositories

Workstation

# Automated deployment: Cloud Native



Travis CI

Qserv images

Docker Registry

Coming soon:
**Continous Delivery**

GitHub
Repositories

Workstation

gcloud

Terraform

Docker Registry

**Cloud Infrastructure:**
**Google Kubernetes Engine**
Openstack

Google Cloud Platform

openstack

kubernetes

Master VM

Worker 1 VM

Worker 2 VM

Worker 25 VM

**Storage:**
**~ 35TB Catalog**
**Google Persistent Disk**
Ceph

ceph

# Automated deployment: bare-metal

# Automated deployment: CI



Workstation

For each commit
- Build and tag container
- Start kind
- Start Qserv pods
- Launch integration tests
- Push container to Docker hub

Docker Registry

**Kind**

**(embedded in travis-ci)**

Travis CI

Build node

kubernetes

Master

Worker 1

Worker 2

Worker 3

# K8s + Microservice features

- ★ Automated scaling
- ★ Container scheduling
- ★ Auto-healing
- ★ Continuous deployment

- ★ Volume management (storage)
- ★ Easy monitoring
- ★ Healthcheck
- ★ Security

Key

**Node Class Type**

Pod

Container

Component — Contains & Uses Library Component

Uses Storage Component

**Master Class Node (1..n)**

Master Pod

Only runs on 1, eventually 2, of the Master Class Nodes

MySql Proxy
MySql Proxy

MariaDb
MariaDb

CSS Pod (1..2)

MariaDb
QMeta    CSS

CZAR

**Results Local Volume**
Results Database Storage

**CSS Persistent Local Volume**
QMeta Database Storage    CSS Database Storage

includes Secondary Index

**Worker Class Node (1..n)**

Worker Pod

CMSD
CMSD

XRootD
XRootD

Replication Worker
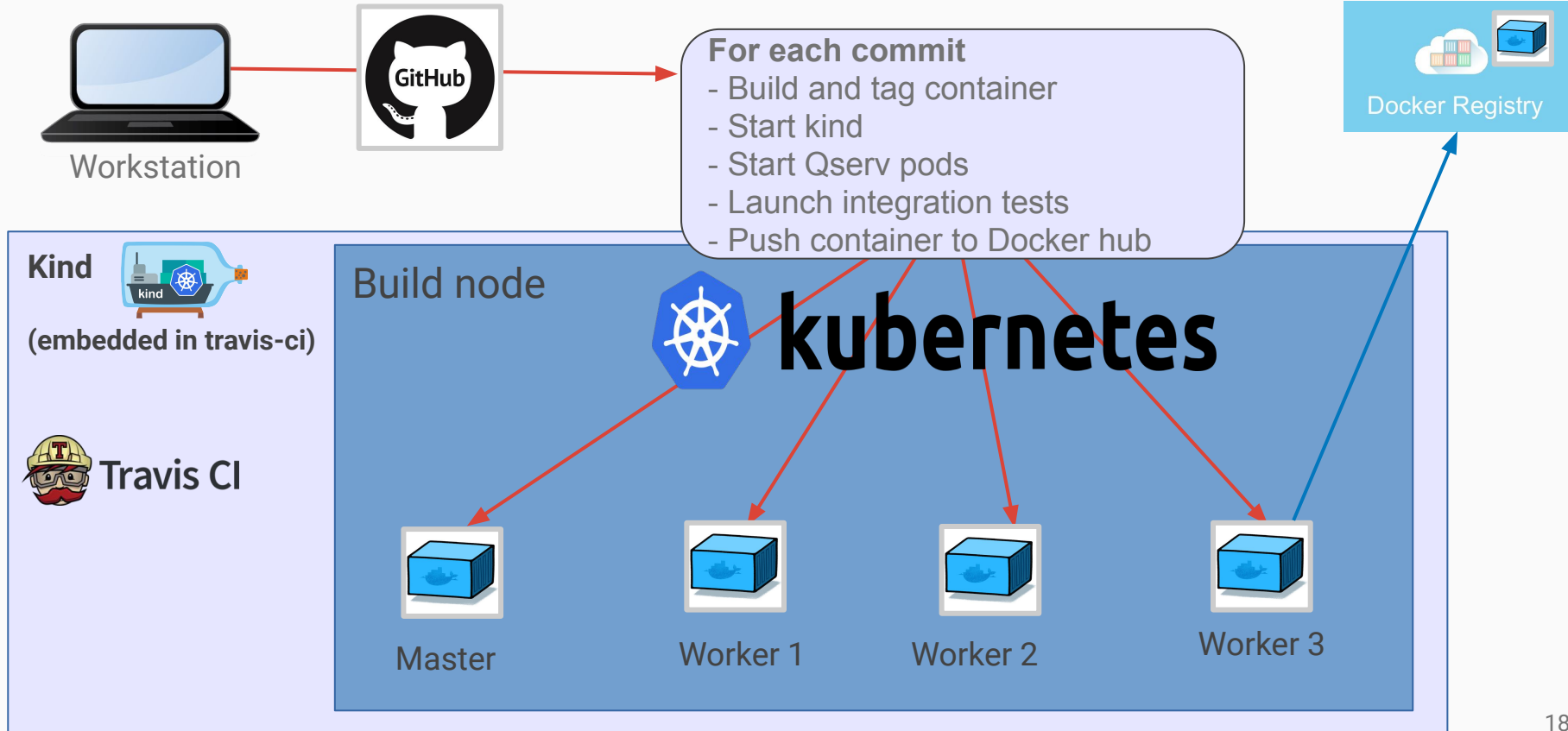Replication Worker

MariaDb
MariaDb

Worker

**Worker Persistent Local Volume**
Catalog Shards

**Unspecified Node(s)**

Redirector Pod (1..2)

CMSD
CMSD

XRootD
XRootD

Replicator Pod (1)

MariaDb
MariaDb

Repl–Master
Repl–Master

**Repl. Persistent Local Volume**
Repl Database Storage

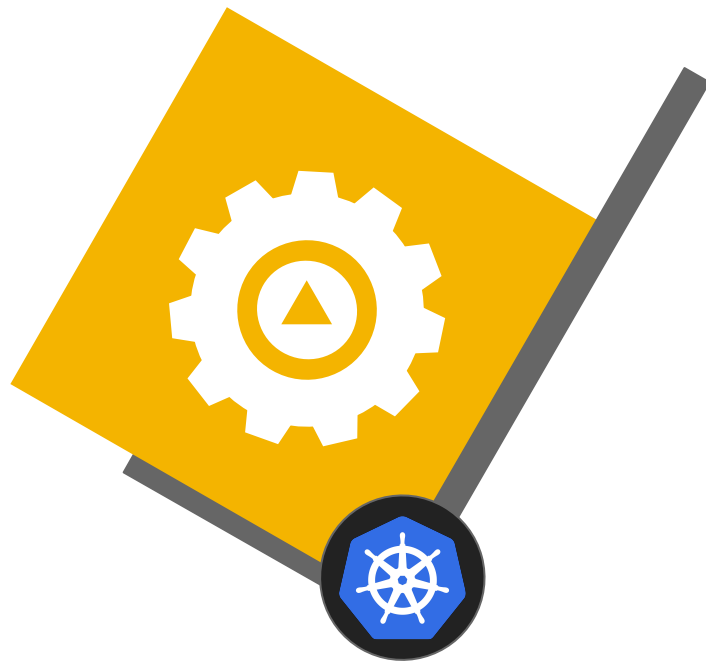# The killer feature: workload portability

**Result: Portability**

Put your app on wheels and move it whenever and wherever you need

Easily move your distributed application anywhere Kubernetes is supported, in seconds.

# Operators: adding sysadmin knowledge inside k8s

⬡ **Resize/Upgrade**

⬡ **Reconfigure**

⬡ **Backup**

⬡ **Healing**



The Sysadmin

# Operators embed ops knowledge from the experts



ops knowledge from the experts

operator
implementation
i.e. k8s controller

Deployments
StatefulSets
Autoscalers
Secrets
Config maps

**See**

- https://kubernetes.io/docs/concepts/extend-kubernetes/operator/

- https://cloud.google.com/blog/products/containers-kubernetes/best-practices-for-building-kubernetes-operators-and-stateful-apps

| Phase I | Phase II | Phase III | Phase IV | Phase V |
|---------|----------|-----------|----------|---------|
| **Basic Install** | **Seamless Upgrades** | **Full Lifecycle** | **Deep Insights** | **Auto Pilot** |
| Automated application provisioning and configuration management | Patch and minor version upgrades supported | App lifecycle, storage lifecycle (backup, failure recovery) | Metrics, alerts, log processing and workload analysis | Horizontal/vertical scaling, auto config tuning, abnormal detection, scheduling tuning |

**HELM** ←——————→

**ANSIBLE** ←——————————————————————→

**GO** ←——————————————————————→

Red Hat

# On-premise
# vs
# Public Cloud

# Containers at Google

Each week, Google launches more than four billion containers across its data centers around the world. These containers house the full range of applications Google runs, including user-facing applications such as Search, Gmail, and YouTube.

Kubernetes was directly inspired by Google's cluster manager, internally known as Borg. Borg allows Google to direct hundreds of thousands of software tasks across vast clusters of machines numbering in the tens of thousands — supporting seven businesses with over one billion users each. Borg and Kubernetes are the culmination of Google's experience deploying resilient applications at scale.
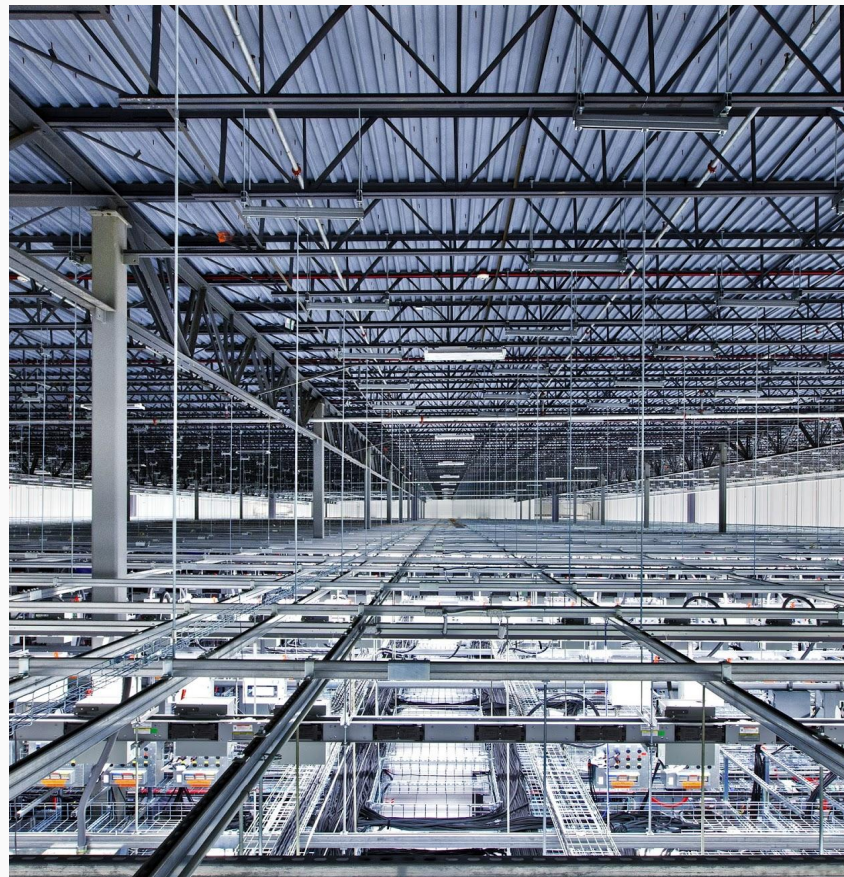
# Public cloud: pros and cons

Pros
* ★ Flexibility for infrastructure provisionning:
  * ○ setup a 40 nodes Qserv cluster in 0.5 days
  * ○ extend it to 50 nodes in 10 seconds
* ★ Excellent support from Google engineers
* ★ Easy to setup development clusters with few maintenance
* ★ Cool proprietary features

Cons,
* ★ Expensive for production platform
  * ○ 100K in 3 months for LSST
* ★ Easy to get stuck with proprietary features
* ★ Hide Kubernetes internals so may be difficult to setup
* ★ Run slower than bare-metal (~25%)

# On-premise: pros and cons

Pros

- ★ Flexibility on cluster setup
  - ○ DIY Kubernetes
  - ○ Fine-tune your components (local HDD)
- ★ Require skilled engineers
- ★ Ease to guarantee your workload portability
- ★ Run faster than public cloud

Cons'

- ★ Difficult to retrieve the global cost
- ★ Require manpower for setup and maintenance
- ★ Hardware upgrade are cost-effective and slow
- ★ Difficult to rebuild the cluster from scratch



© PHOTOTHEQUE IN2P3 / CNRS

# Thanks!

Contact:

Fabrice JAMMES
Formation et conseil Kubernetes
https://k8s-school.fr

fabrice.jammes@k8s-school.fr