# Reperforming a Nobel Prize discovery on Kubernetes

**Lukas Heinrich, Ricardo Rocha**

ATLAS and CMS discovered the Higgs boson in 2012.

Since then the LHC has released some of their data publicly.

In this presentation we'll try to reproduce one of the main results from CMS Open Data using modern cloud tools.

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

WLCG

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

Google earth

# Kubernetes

Spun out of Google as an open source

   container orchestration project

Built on lessons from Borg and Omega



Borg, Omega, and Kubernetes

LESSONS LEARNED FROM THREE CONTAINER-MANAGEMENT SYSTEMS OVER A DECADE

BRENDAN BURNS, BRIAN GRANT, DAVID OPPENHEIMER, ERIC BREWER, AND JOHN WILKES, GOOGLE INC.

Though widespread interest in software containers is a relatively recent phenomenon, at Google we have been managing Linux containers at scale for more than ten years and built three different container-management systems in that time. Each system was heavily

Loosely coupled collection of components to deploy, maintain and scale workloads

**Declarative, Load Balancing, Self Healing, Auto Scaling**

Service and Batch Workloads

# Kubernetes

Largest open source project after kernel

**35.000** contributors, **148.000** code commits

**83.000** pull requests, **1.1M** contributions

**2000+** contributing companies

Google, RedHat, VMware, Huawei, Microsoft, IBM, Fujitsu, …

Open community welcome to contributions

Special Interest Groups (SIGs) : Auto-Scaling, Multi-Cluster, Scheduling, ...

Borg, Omega, and Kubernetes

LESSONS LEARNED FROM THREE CONTAINER-MANAGEMENT SYSTEMS OVER A DECADE

BRENDAN BURNS, BRIAN GRANT, DAVID OPPENHEIMER, ERIC BREWER, AND JOHN WILKES, GOOGLE INC.

Though widespread interest in software containers is a relatively recent phenomenon, at Google we have been managing Linux containers at scale for more than ten years and built three different container-management systems in that time. Each system was heavily

# Kubernetes

Lingua franca of the cloud

Managed services offered by all major public clouds

Multiple options for on-premise or self-managed deployments

Common declarative API for basic infrastructure : compute, storage, networking

Healthy ecosystem of tools offering extended functionality

# Cloud Native Landscape
*v0.9.9*

## App Definition & Development

**Database & Data Analytics** | **Streaming** | **SCM** | **Application Definition** | **CI/CD**

## Orchestration & Management

**Scheduling & Orchestration** | **Coordination & Service Discovery** | **Service Management**

- kubernetes — CNCF Project
- CoreDNS — CNCF Project
- linkerd — CNCF Project
- GRPC — CNCF Project
- envoy — CNCF Project

## Runtime

**Cloud-Native Storage** | **Container Runtime** | **Cloud-Native Network**

- containerd — CNCF Project
- rkt — CNCF Project
- CNI — CNCF Project

## Provisioning

**Host Management / Tooling** | **Infrastructure Automation** | **Container Registries** | **Secure Images** | **Key Management**

- notary — CNCF Project
- TUF — CNCF Project

## Cloud

**Public** | **Private**

## Platforms

### PaaS / Container Service

### Serverless/Event-based

## Observability & Analysis

### Monitoring
- Prometheus — CNCF Project

### Logging
- fluentd — CNCF Project

### Tracing
- OPENTRACING — CNCF Project
- JAEGER — CNCF Project

---

This landscape is intended as a map through the previously uncharted terrain of cloud native technologies. There are many routes to deploying a cloud native application, with CNCF Projects representing a particularly well-traveled path.

github.com/cncf/landscape

**CLOUD NATIVE COMPUTING FOUNDATION**

Redpoint    //Amplify PARTNERS

Greyed logos are not open source

# Kubernetes and containers at CERN

Started offering Mesos/DCOS, Swarm and Kubernetes

    Now only Swarm and Kubernetes

Kubernetes by far the most popular solution

    Spark as a Service, WebLogic, JIRA

    INSPIRE-HEP, REANA/RECAST, Jupyter

    OpenStack, Batch / Condor

    And many others

| Clusters | Nodes |
|----------|-------|
| 468 | 2988 |

| Kubernetes | Swarm |
|------------|-------|
| 412 | 47 |

| Mesos | DCOS |
|-------|------|
| 2 | 7 |

# Containers & Open Data in Science

## CONTAINERS IN THE CLOUD

Standardized platforms allow researchers to run each other's software – no installation required. By Jeffrey M. Perkel

"Researchers can be confident that their code will remain valid, whichever platform they choose."

### Easier evaluation

### Languages and clouds

PERSPECTIVE

https://doi.org/10.1038/s41567-018-0342-2

Corrected: Publisher Correction

OPEN

## Open is not enough

Xiaoli Chen[1,2], Sünje Dallmeier-Tiessen[1]*, Robin Dasler[1,11], Sebastian Feger[1,3], Pamfilos Fokianos[1], Jose Benito Gonzalez[1], Harri Hirvonsalo[1,4,12], Dinos Kousidis[1], Artemis Lavasa[1], Salvatore Mele[1], Diego Rodriguez Rodriguez[1], Tibor Šimko[1]*, Tim Smith[1], Ana Trisovic[1,5]*, Anna Trzcinska[1], Ioannis Tsanaktsidis[1], Markus Zimmermann[1], Kyle Cranmer[6], Lukas Heinrich[6], Gordon Watts[7], Michael Hildreth[8], Lara Lloret Iglesias[9], Kati Lassila-Perini[4] and Sebastian Neubert[10]

The solutions adopted by the high-energy physics community to foster reproducible research are examples of best practices that could be embraced more widely. This first experience suggests that reproducibility requires going beyond openness.

### Approaching reproducibility and reuse in HEP

nature

# Containers & (Open) Science

**CONTAINERS IN THE CLOUD**

Standardized platforms allow researchers to run each other's software – no installation required. By Jeffrey M. Perkel

somebody else's ridiculously difficult – neuroscientist London.

...gle command. ...ced friction for ...oducing some- ...u have to build combine it with Lorena Barba, a ...ineer at George ...shington DC. "It ...ss error-prone, ...rcher time."

security, researchers can build private 'BinderHubs' instead. The Alan Turing Institute has two, including one called Hub23 (a reference to Hut 23 at the Second World War code-breaking facility at Bletchley Park, UK), that provides

> **"Researchers can be confident that their code will remain usable, whichever platform they choose."**

greater computational resources and the ability to work with data sets that cannot be

Pyth... lang... Ocea... puta... from... In 20... *Intel*... pilot... revi... *Bioi*... More... the t... direc... and t... thos...

**PERSPECTIVE**

**OPEN**

## Open is not enough

Xiaoli Chen[1,2], Sünje Dallmeier-Tiessen[1]*, Robin Dasler[1], Sebastian Feger[1,3], Pamfilos Fokianos[1], Jose Benito Gonzalez[1], Harri Hirvonsalo[1,4,12], Dinos Kousidis[1], Artemis Lavasa[1], Salvatore Mele[1], Diego Rodriguez Rodriguez[1], Tibor Šimko[1], Tim Smith[1], Ana Trisovic[1,3,4], Anna Trzcinska[1], Ioannis Tsanaktsidis[1], Markus Zimmermann[1], Kyle Cranmer[5], Lukas Heinrich[6], Gordon Watts[7], Michael Hildreth[8], Lara Lloret Iglesias[9], Kati Lassila-Perini[1] and Sebastian Neubert[10]

The solutions adopted by the high-energy physics community to foster reproducible research are examples of best practices that could be embraced more widely. This first experience suggests that reproducibility requires going beyond openness.

Our own experience from opening up vast volumes of data is that openness cannot simply be tacked on as an afterthought at the end of the scientific endeavour. In addition, openness alone does not guarantee reproducibility or reusability, so it should not be pursued as a goal in itself. Focusing on data is also not enough: it needs to be accompanied by software, workflows and explanations, all of which need to be captured throughout the usual iterative and closed research lifecycle, ready for a timely open release with the results.

**nature**

# Challenge: H→4l re-discovery on CMS Open Data

Benchmark analysis based on Open LHC Data.

**Goal**: Fit it within a live demo for 20-minute [Keynote at KubeCon EU 2019](#)
Learn something about cloud-native analysis, reproducibility, Open Data.
Have some Fun.

# Demo

# Challenge: H→4l re-discovery on CMS Open Data

what would this look like in a cloud-native approach?



**Sim Higgs**

**Background**

**Background**

**Real Data**

Event Data

20k+ Core K8s Clusters

Summary Data

Make Plot!

**70 TB** of Physics Data        **~25000** Files

70 TB Dataset

OpenStack Magnum

25000 Kubernetes Jobs

Job Results

Interactive Visualization

Aggregation

70 TB Dataset → Cluster on GKE → Job Results → Interactive Visualization

Max 25000 Cores

Single Region, 3 Zones

25000 Kubernetes Jobs

Aggregation

# Data Upload

Initial dataset (opendata) available on /eos



~1 day for full dataset transfer, done first to Zurich then to NL

**Ingress is free, Ingress is free**...

# Data Upload

Initial dataset (opendata) available on /eos



**Ingress is free, Ingress is free**…

# GCP Analysis Run

Kubernetes clusters on GKE ( Managed Kubernetes service on GCP )

Today's run included

    660 nodes: n1-highmem-16, 104 GB RAM

    10560 cores, 69 TB RAM

| Cluster Creation | → | Image Pre-Pull | → | Data Stage-In | → | Process |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| 5 min | | 4 min | | 4 min | | 90 sec |

# GCP Analysis Run

Network guarantees 2Gb/core up to 16 core nodes ( **32 Gbit per VM** ! )

GCS can handle these rates somehow, and we end up bound by local i/o

Ended up using in-memory filesystems to go around this

| | Zonal standard persistent disks | Regional persistent disks | Zonal SSD persistent disks | Regional SSD persistent disks | Local SSD (SCSI) | Local SSD (NVMe) |
|---|---|---|---|---|---|---|
| **Maximum sustained IOPS** | | | | | | |
| Read IOPS per GB | 0.75 | 0.75 | 30 | 30 | 266.7 | 453.3 |
| Write IOPS per GB | 1.5 | 1.5 | 30 | 30 | 186.7 | 240 |
| Read IOPS per instance | 3,000 | 3,000 | 15,000 - 60,000* | 15,000 - 60,000* | 400,000 | 680,000 |
| Write IOPS per instance | 15,000 | 15,000 | 15,000 - 30,000* | 15,000 - 30,000* | 280,000 | 360,000 |

# GCP Analysis Run

Network guarantees 2Gb/core up to 16 core nodes ( **32 Gbit per VM** ! )

# GCP Pricing

Billing is updated daily, though there are APIs to query for details

Considering a ~10 minutes run it implies (compute table prices, NL region)

$1.043 * 1530 / 6 = **$260** (**~5x cheaper if using pre-emptibles**)

Parking storage cost for the dataset (monthly cost, lots of room for creativity)

$0.020 * 70000 = $1400

Total under $300 usd

Running on credits, **no Committed Use or Sustained Compute discounts**

# Open Questions

A stunt… or could we come up with a usable model?

Technically feasible. What do these technologies imply for LHC computing?

Analysis Models,
Infrastructure,
Funding Models,

….

# Opportunities for Infrastructure

Simplified deployments (Federation), common APIs

Bursting Scale-out to near-arbitrary scales
Auto-scaling

Access to special hardware
(FPGAs, TPUs, …) easily
Integrated into LHC computing

Federation

K8s API

K8s API

K8s
Cluster
1(Microsoft)

K8s API

K8s
Cluster
1(Amazon)

K8s API

K8s
Cluster 1(T-
Systems)

K8s API

K8s
Cluster
1(CERN)

K8s
Cluster
1(Google)

K8s API

K8s
Cluster
1(CERN)

K8s API

K8s
Cluster
1(CERN)

# Opportunities for Analysis Models

Rich gateway to scale-out, adaptive, on-demand computing

out-of-core dataframes

- Many systems natively integrate w/ k8s
- Rich real-time monitoring

How do you move smoothly
Between real-time analysis
and batch/scheduled work?

- k8s supports both well

# Open Data accessible to everyone at scale

LHC experiments part of growing list of experiments with complex
open data problems: data complexity, data volume. large collaborations.



Open Data only useful, if it is feasible for external researchers to analyze it .
Demo goal: Show that  public cloud can provide required scale on-demand.

```
[16:01:21] cmsusr@e6f7bea2253e /Users/lukasheinrich/Code/awesomedemo/higgs-demo/CMSSW_5_3_32/src $ \root -b

  *********************************************
  *                                           *
  *        W E L C O M E  to  R O O T         *
  *                                           *
  *   Version   5.32/00    2 December 2011    *
  *                                           *
  *  You are welcome to visit our Web site    *
  *           http://root.cern.ch             *
  *                                           *
  *********************************************

ROOT 5.32/00 (branches/v5-32-00-patches@42372, Jun 10 2014, 18:26:00 on linuxx8664gcc)

CINT/ROOT C/C++ Interpreter version 5.18.00, July 2, 2010
```

Moving to Cloud-based technologies:
analysis preservation as a by-product

Beyond a VM:  Containerized CMSSW
 ~decade old software to reproduce results

### cmsopendata/cmssw_5_3_32 ☆

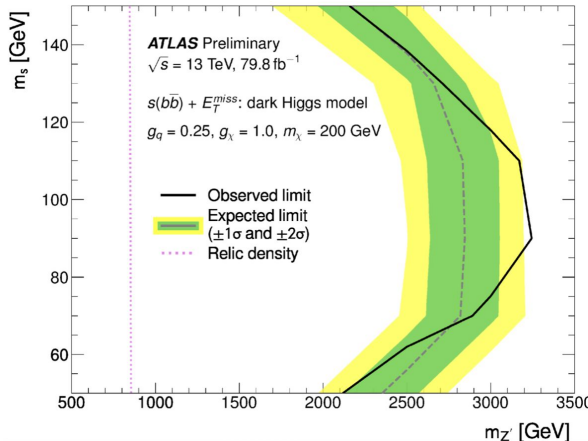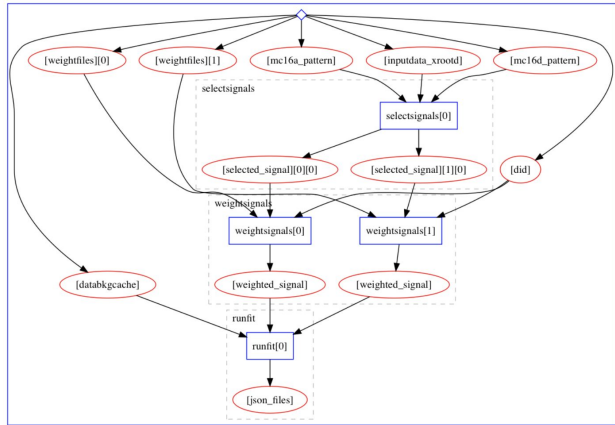By **cmsopendata** • Updated 4 months ago

Container

Effective **re-use** of HEP analysis to generate new science results based on archived software.

Only possible through container-based workflows exposed to the user

# Conclusions

Demonstrated Tbps analysis of CMS Open Data.

Modern cloud computing paradigms can give individuals scale to realistically analyze LHC data, foster reproducibility & reusability of LHC analyses.

Opens up exciting opportunities to evolve the LHC computing landscape as we look towards Run-4 / the HL-LHC etra. Cross-team collaborations are crucial for R&D[*].

Thought expt: if we started today, what would our infrastructure look like?

We did learn a lot & had some fun.

[*] ATLAS  (LH) IT (RR) CMS (Clemens Lange) for R&D towards Run-4