

# Deploying Apache Ranger in the Big Data ecosystem at CERN

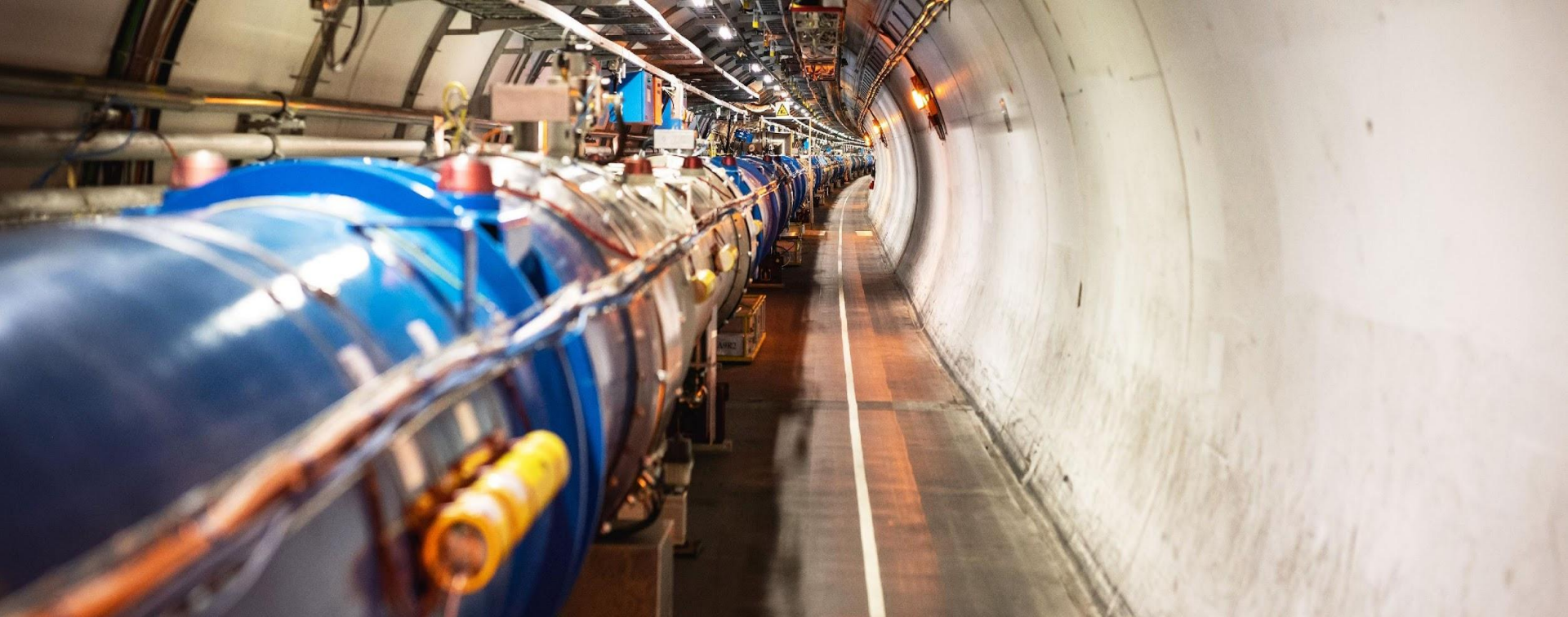
Emil Kleszcz @ CERN

09 Oct 2023

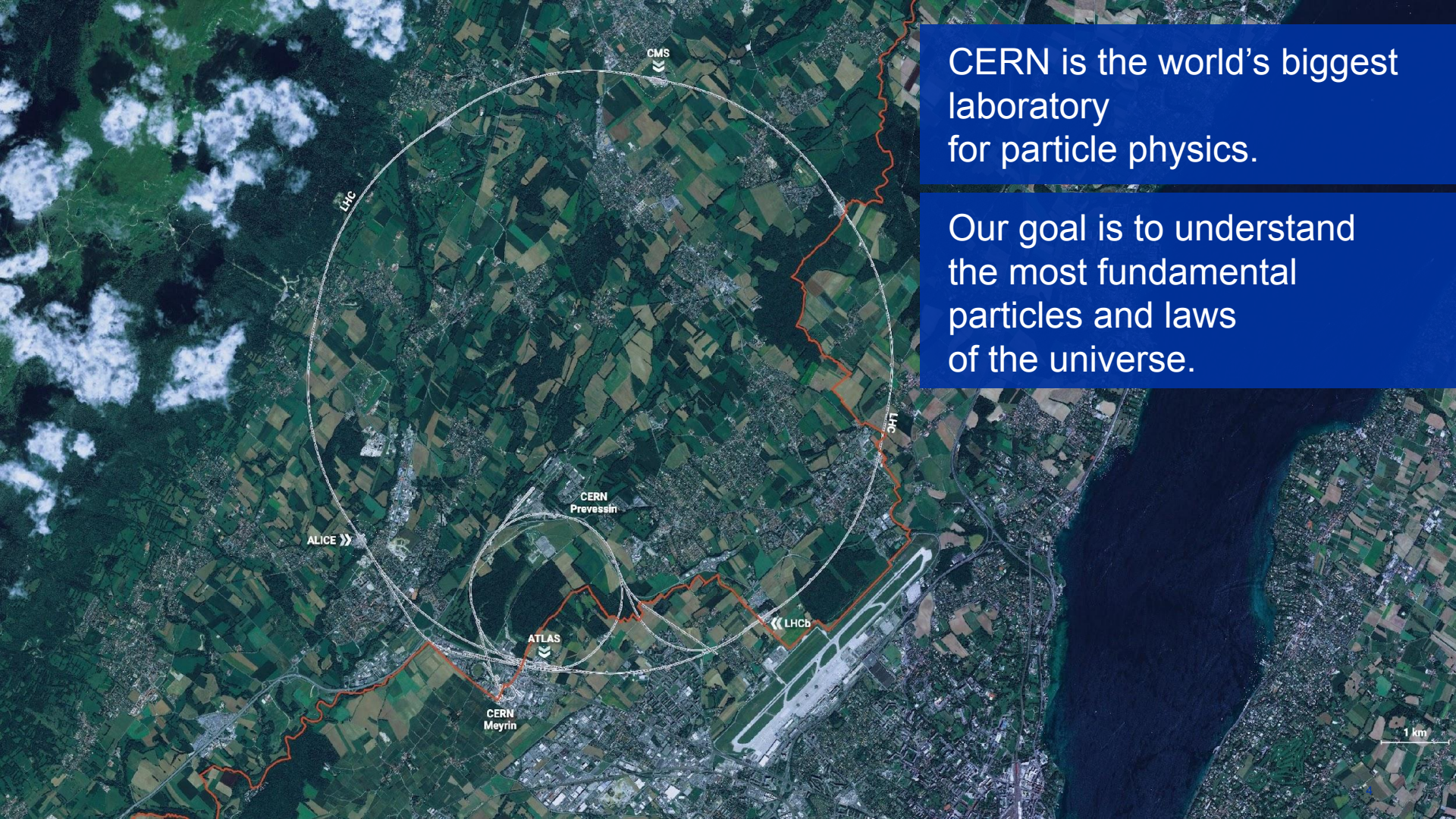
# Presentation content

1. About CERN
2. Big Data ecosystem @ CERN
3. Overview of Apache Ranger project
4. Apache Ranger @ CERN
5. Summary





# About CERN



CERN is the world's biggest laboratory for particle physics.

Our goal is to understand the most fundamental particles and laws of the universe.

A low-angle photograph of several tall flagpoles against a clear blue sky. The sun is visible in the background, creating a lens flare. Various national flags are flying from the poles, including Spain, Greece, Italy, Denmark, Bulgaria, Finland, Germany, and the United Kingdom. A large orange circle is overlaid on the left side of the image, containing the word 'COLLABORATION' in white capital letters.

# COLLABORATION

# CERN was founded in 1954 with 12 European Member States



## 23 Member States

Austria – Belgium – Bulgaria – Czech Republic  
Denmark – Finland – France – Germany – Greece  
Hungary – Israel – Italy – Netherlands – Norway  
Poland – Portugal – Romania – Serbia – Slovakia  
Spain – Sweden – Switzerland – United Kingdom

## 3 Associate Member States in the pre-stage to membership

Cyprus – Estonia – Slovenia

## 7 Associate Member States

Croatia – India – Latvia – Lithuania – Pakistan  
Türkiye – Ukraine

## 6 Observers

Japan – Russia (suspended) – USA  
European Union – JINR (suspended) – UNESCO

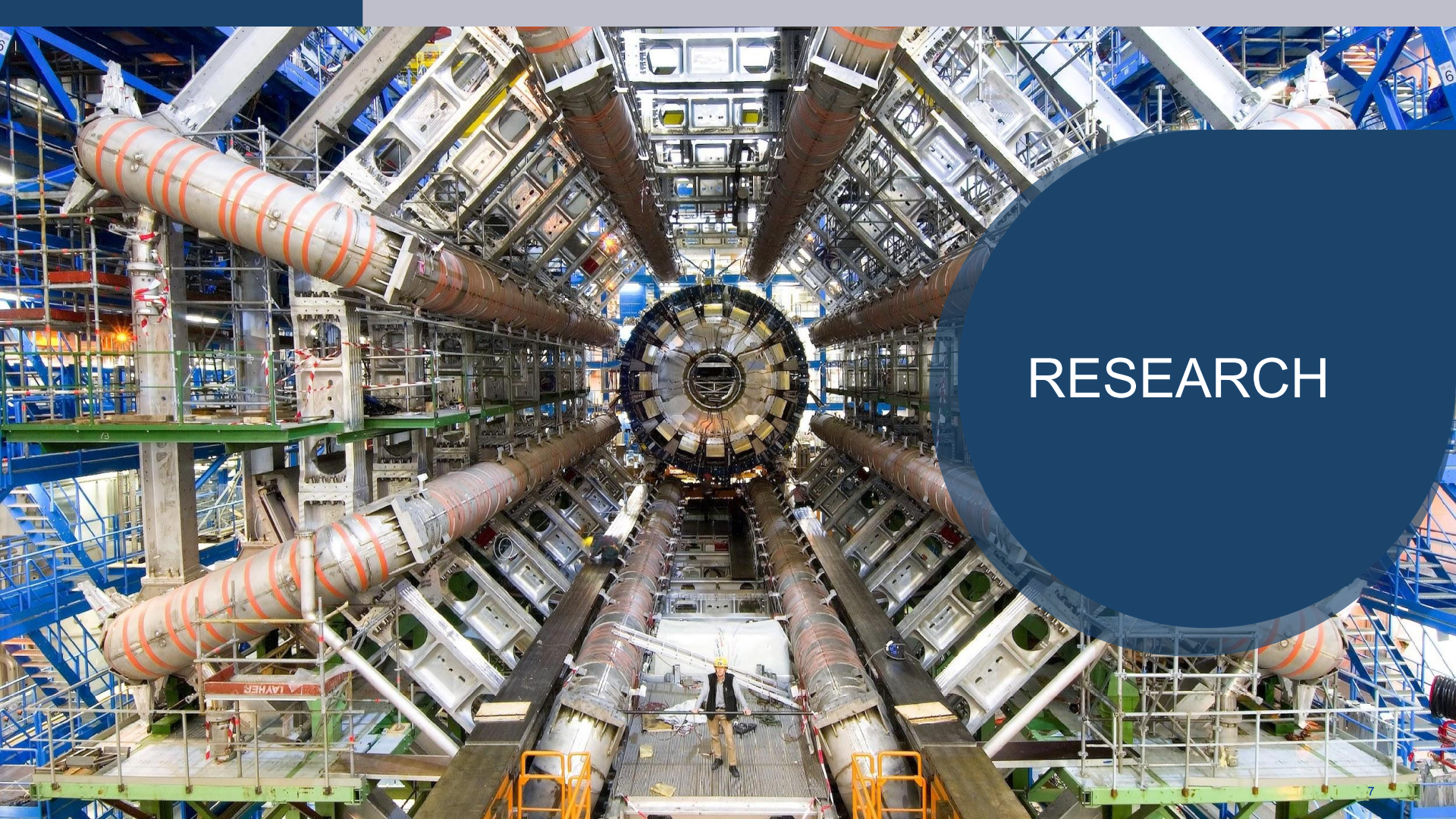
Brazil signed Associate Membership Agreement  
on 3 March 2022. Ratification process in Brazil  
is underway.

Around 50 Cooperation Agreements  
with non-Member States and Territories

CERN's annual budget  
is 1200 MCHF (equivalent  
to a medium-sized European  
university)

As of 31 December 2021  
Employees:  
**2676** staff, **783** fellows

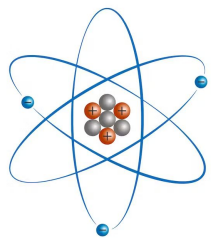
Associates:  
**11175** users, **1556** others



RESEARCH

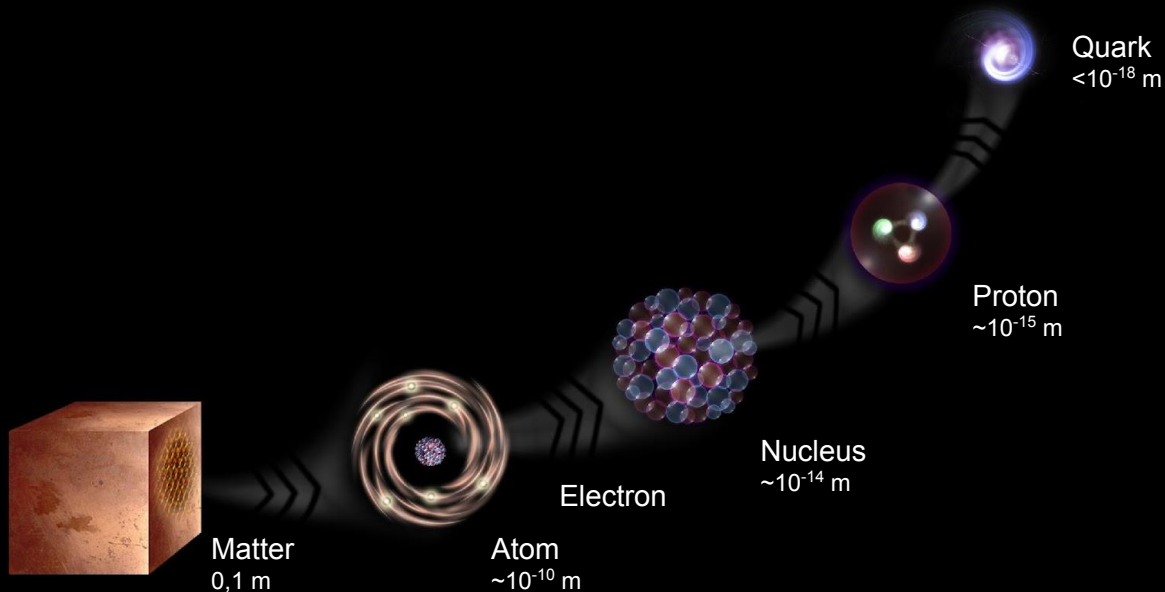
# What is the universe made of?

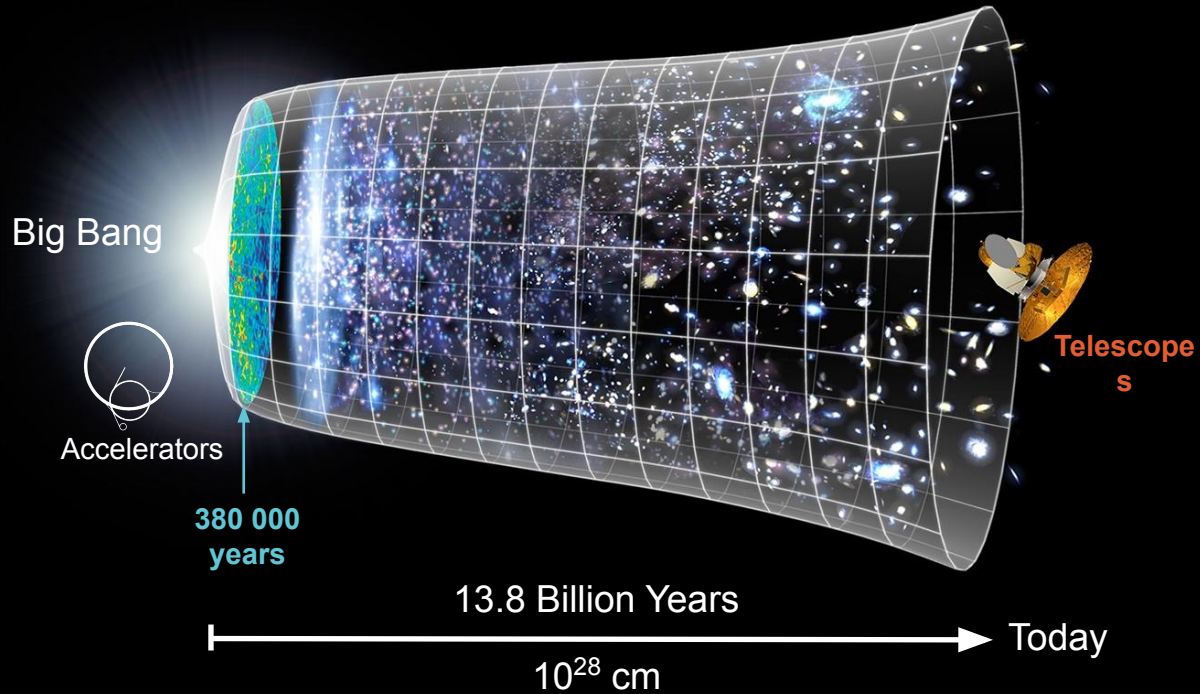
We study the elementary building blocks of matter and the forces that control their behavior



Atom structure

- Proton
- Neutron
- Electron





## How did the universe begin?

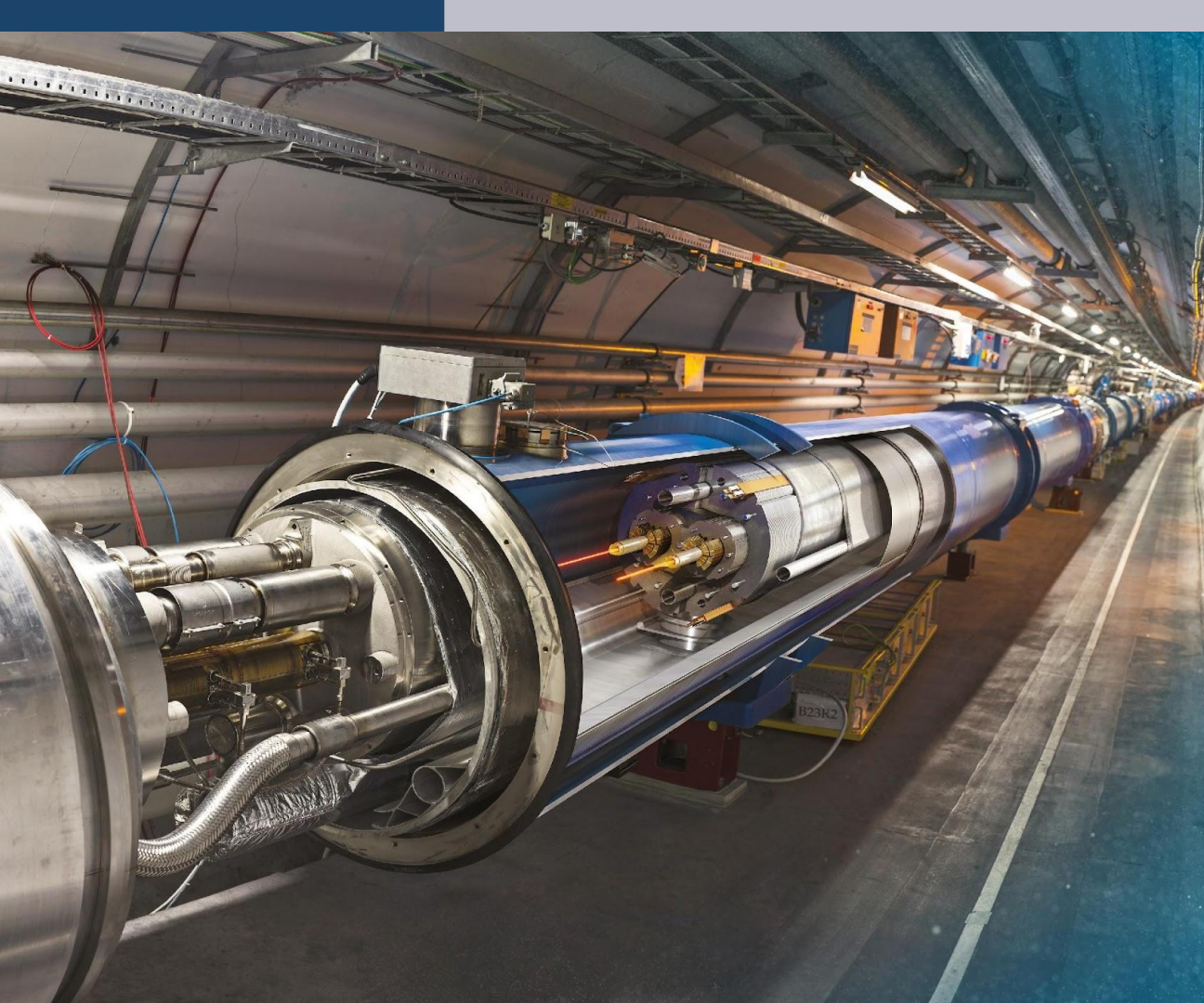
We reproduce the conditions a fraction of a second after the Big Bang, to gain insight into the structure and evolution of the universe.

# Large Hadron Collider (LHC)



Protons  
(ions)

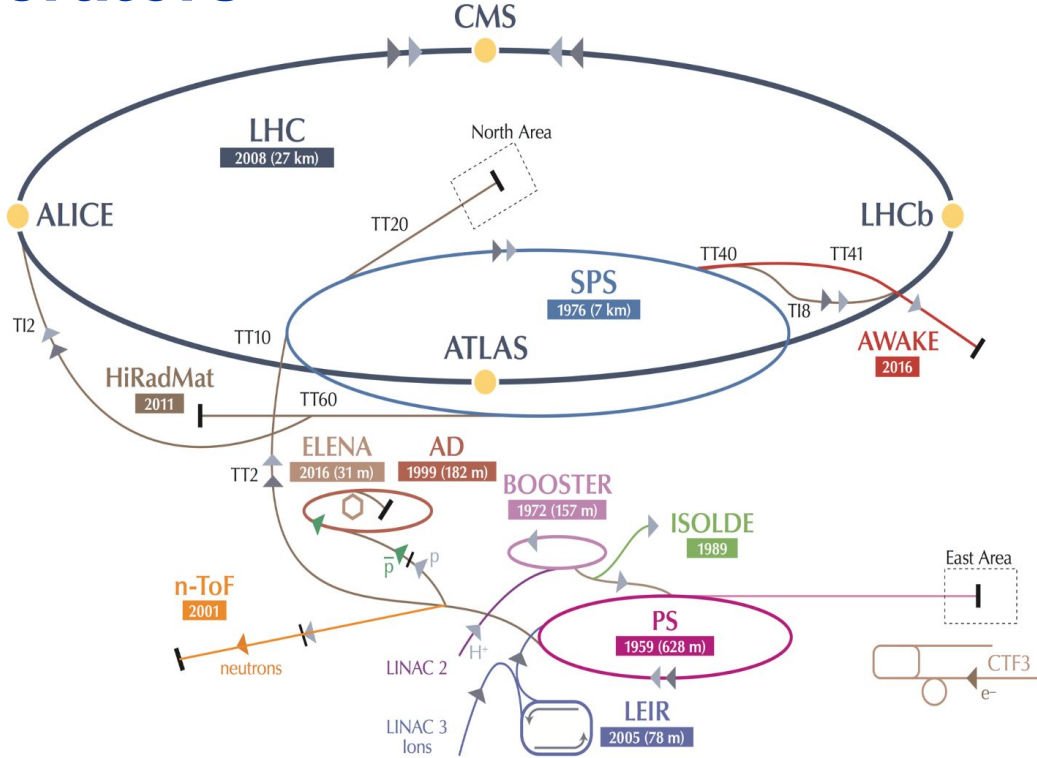
Protons  
(ions)



# Large Hadron Collider (LHC)

- 27 km in circumference
- About 100 m underground
- Superconducting magnets steer the particles around the ring
- Particles are accelerated close to the speed of light

# CERN Accelerators

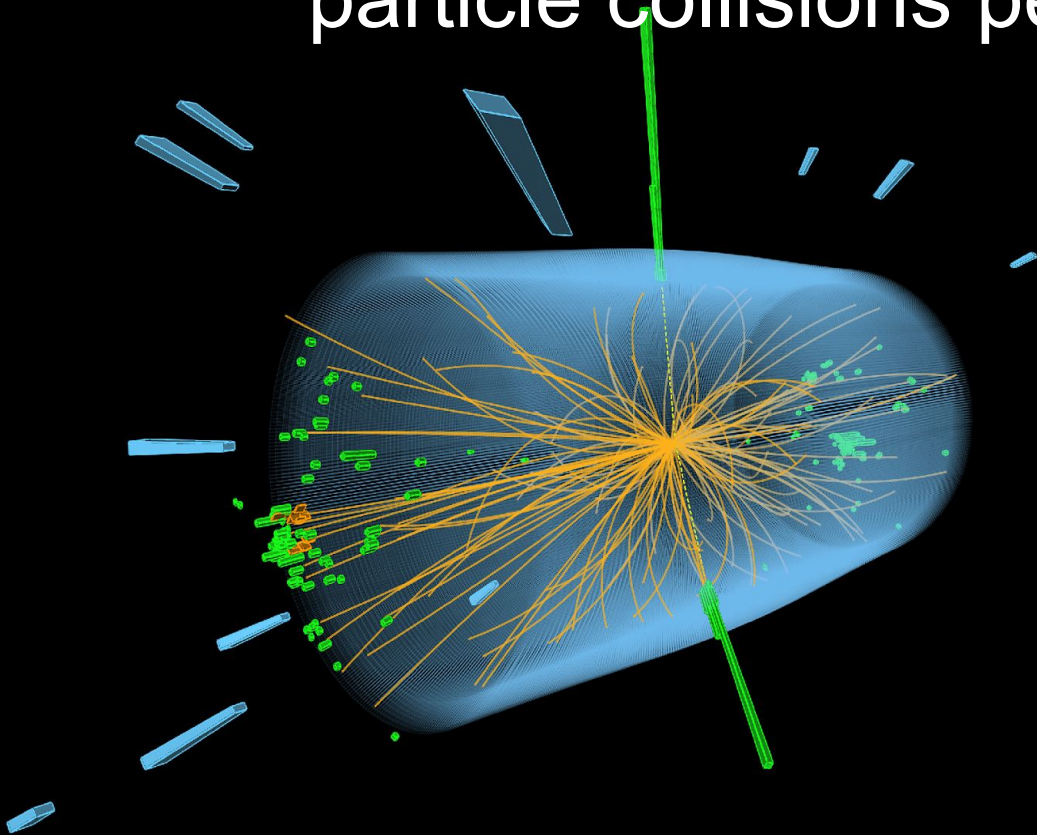


~20 projects other than LHC with > 1200 physicists



DATA  
CENTRE

# The LHC produces more than 1 billion particle collisions per second



The energy of the particles in collision is converted into new particles.

# CERN DC and WLCG (Worldwide LHC Computing Grid)

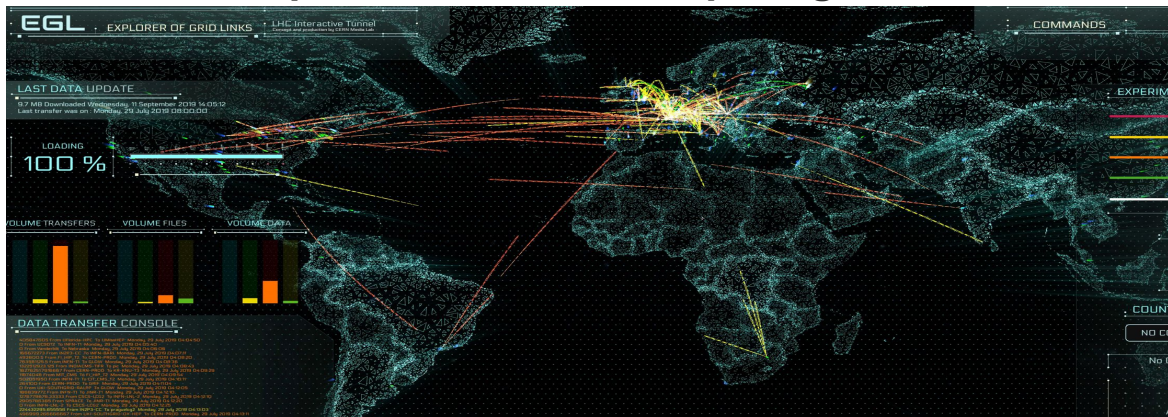
## CERN Science is Data Intensive Science

- Run 1 (2009-2013) we stored 65 PB
- Run 2 (2015-2018) we stored 209 PB
- Run 3 (2022-2026) we expect to store 600 PB

## • CERN Data Centre in Meyrin

- **13k servers, 450k CPU cores, 320 PB** of storage
- New data centre being built in Preveessin site

## Solution is to combine the power of all our computing centres

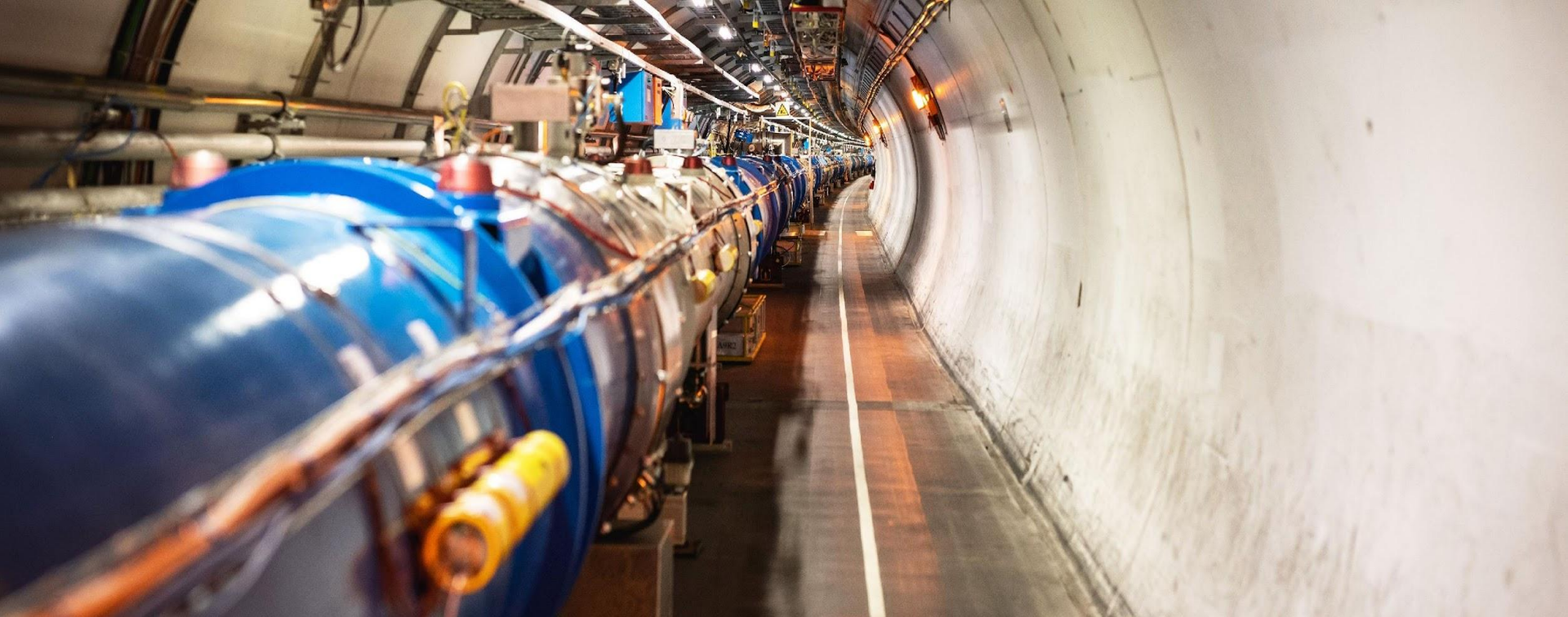


Used to store, distribute, process and analyse data.

1 million processing cores in about 170 data centres and 42 countries.

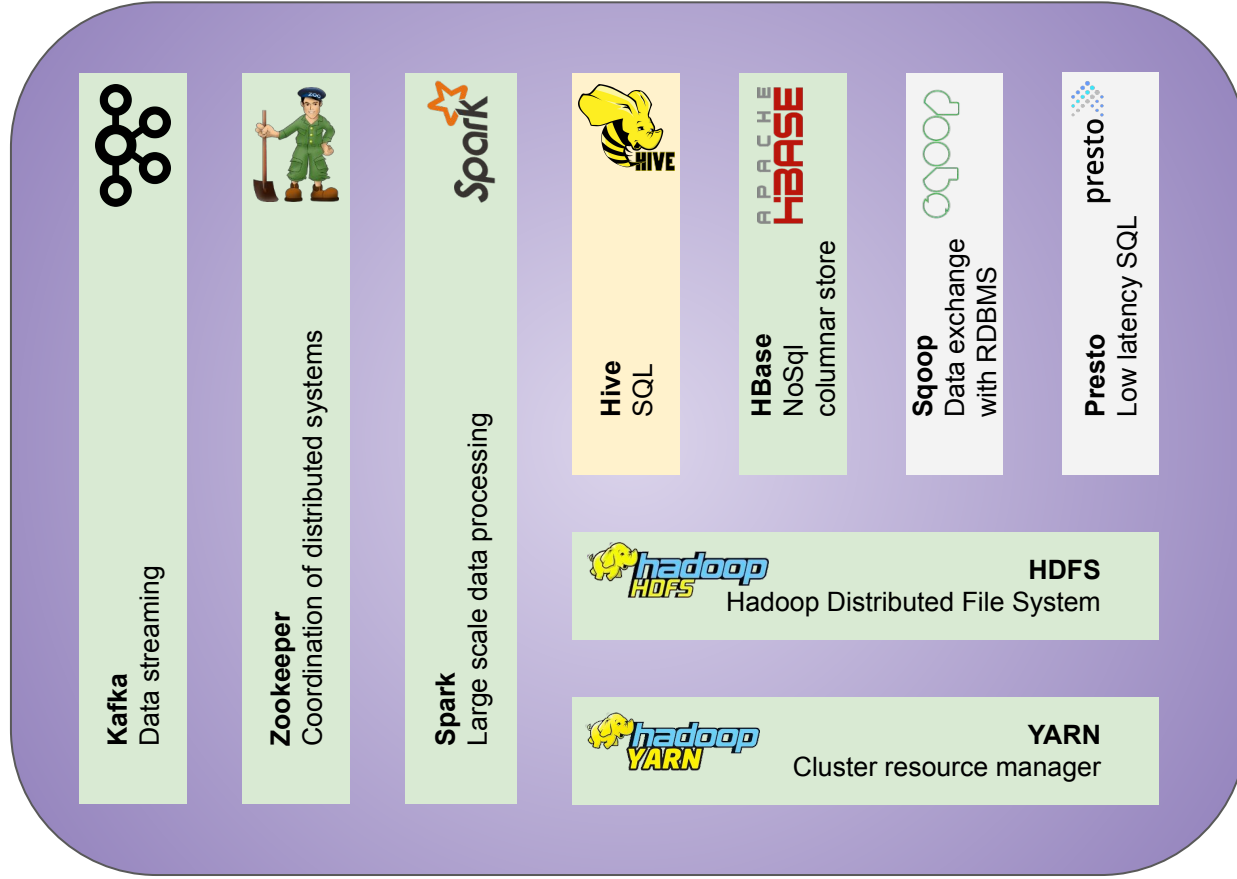
More than 1 Exabyte of CERN data stored world-wide.



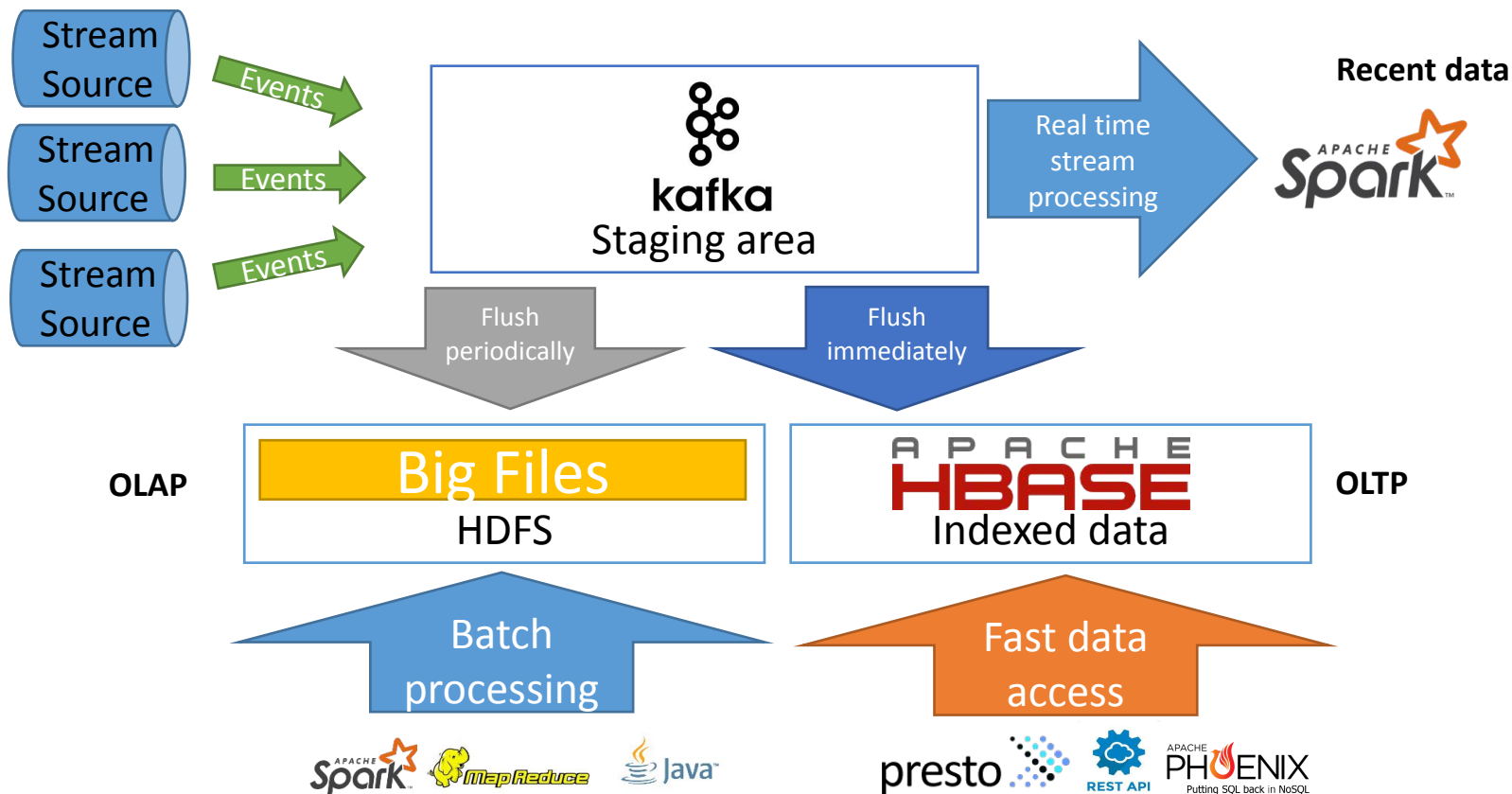


# Big Data ecosystem @ CERN

# Holistic view on the Big Data components



# Hadoop ecosystem @ CERN



# Key users and workload

- Users:
  - IT Security
  - IT Monitoring
  - Experiments
  - Accelerators Monitoring
- HBase and HDFS clusters



# Key statistics

- **5 prod clusters** (+ 3QA & multiple DEV)
- Hardware: Intel-based servers, continuous capacity expansion, hyper-threading enabled, HDDs and SSDs, **180+ bare metal, 60+ VMs**
- **Users: 600+** (general-purpose), **1000+** (accelerators) and growing
- Storage capacity: **40PB+**
- Service started in 2013



# Current security state

- Integration with CERN's LDAP (e-groups)
- Kerberos and SSL
- FW rules defined
- Potential improvements:
  - Policies and ACLs defined manually
  - Lack of central management
  - No user-friendly auditing (raw logs)
  - No CERN SSO integration for Web UIs

```
# changing owner
```

```
hdfs dfs -chown <owner> <hdfs_path>
```

```
# granting new permissions to a user
```

```
hdfs dfs -setfacl -m user:<grantee>:<permission> <hdfs_path>
```

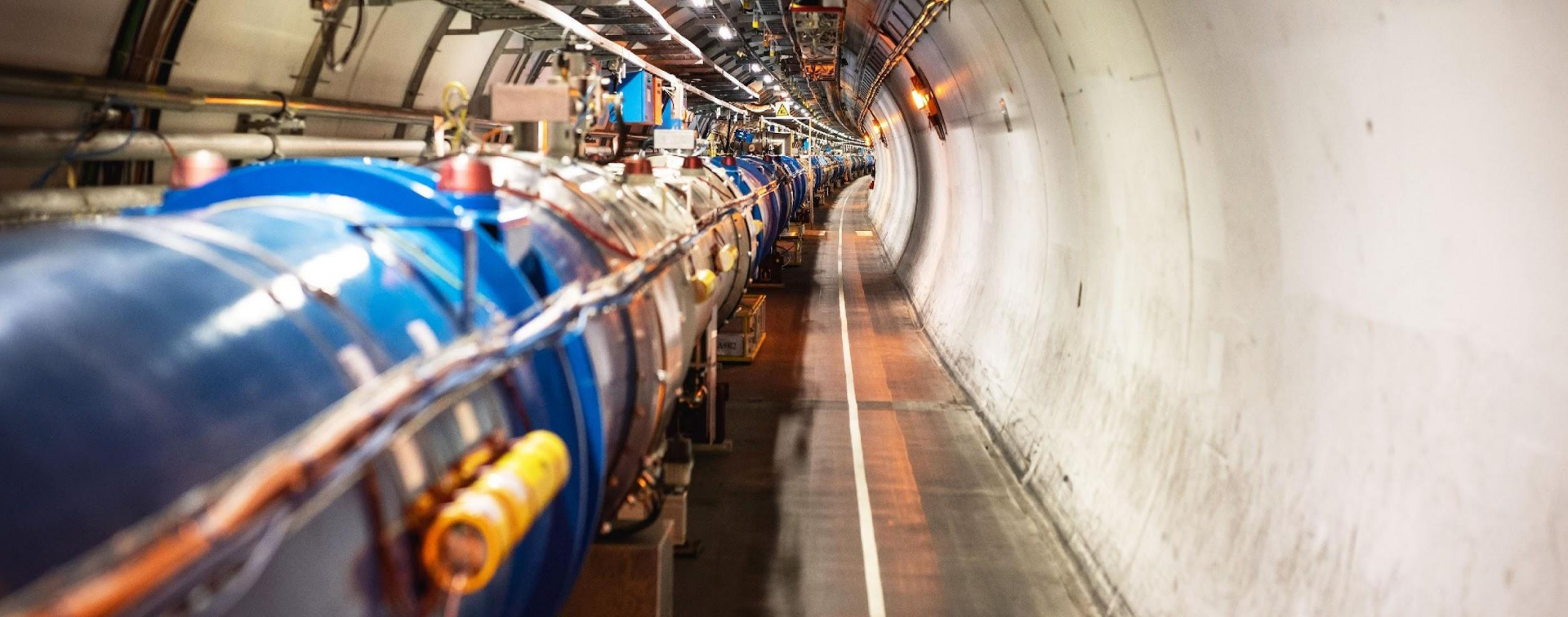
```
# granting new permissions to a group or an e-group
```

```
hdfs dfs -setfacl -m group:<group_name>:<permission> <hdfs_path>
```

**Apache Ranger to improve security and ease management?**



***Apache Ranger***



# Overview of Apache Ranger project

# About the tool



## *Apache Ranger*

- Apache Ranger™ is a framework to enable, monitor and manage comprehensive data security across the Hadoop platform and beyond
- It supports multiple projects such as: Apache Hadoop, Apache HBase, Apache Kafka, YARN, Apache Hive and some more...
- The plugin agent pulls the policy-changes using REST API
- Replaced Apache Sentry



THE  
**APACHE**®  
SOFTWARE FOUNDATION

<https://github.com/apache/ranger>

<https://ranger.apache.org/>

# What Apache Ranger offers

- Fine-Grained Access Control
- Centralized Policy Management
- Dynamic Policy Enforcement
- Centralized auditing
- Delegate administration of policies to group owners
- Integration with LDAP
- Resource and tag-based Policies
- Extensible Lightweight Plugin Architecture
- No single point of failure (if configured)

# WebUI and REST APIs

Ranger REST API 3.0.0-SNAPSHOT

<http://rhdgdev-4iles2101.com.ch.0802api/docs/swagger.json>

Apache Ranger is a framework to enable, monitor and manage comprehensive data security across the Hadoop platform. Apache Ranger currently provides a centralized security administration, fine grain access control and detailed auditing for user access within Apache Hadoop, Apache Hive, Apache HBase and other Apache components

Apache 2.0 License

## AssetREST

GET	/assets/accessAudit
GET	/assets/assets
POST	/assets/assets
GET	/assets/assets/count
POST	/assets/assets/testConfig
DELETE	/assets/assets/{id}
GET	/assets/assets/{id}
PUT	/assets/assets/{id}
GET	/assets/credstores
POST	/assets/credstores
PUT	/assets/credstores
GET	/assets/credstores/count
DELETE	/assets/credstores/{id}
GET	/assets/credstores/{id}

Ranger Access Manager Audit Security Zone Settings

Service Manager cm\_hdfs Policies Create Policy

## Create Policy

### Policy Details :

Policy Type **Access**

Policy Name \*

**enabled**

normal

Policy Label

Resource Path \*

**recursive**

Description

Audit Logging

**YES**

### Allow Conditions :

Select Role	Select Group	Select User	Delegate Admin
<input type="text" value="Select Roles"/>	<input type="text" value="Select Groups"/>	<input type="text" value="Select Users"/>	<input type="checkbox"/>
<a href="#">Add Permissions</a> +			<input type="checkbox"/>

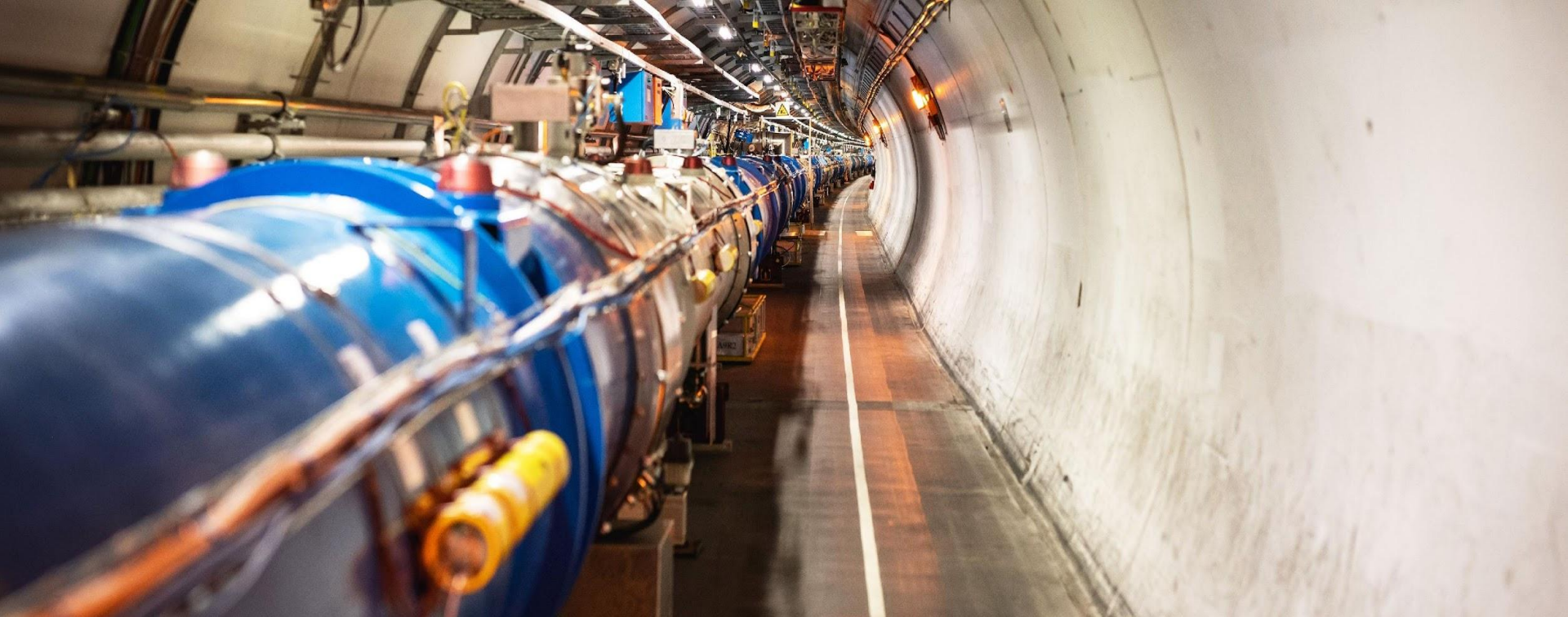
### add/edit permissions

- Read
- Write
- Execute
- Select/Deselect All

# Web interface experience

The screenshot displays the Apache Ranger web interface. The top navigation bar includes 'Resource Policies', 'Tag Policies', 'Audit', 'Security Zone', 'Settings', and 'Reports'. The user is logged in as 'admin'. The left sidebar shows a tree view of services, with 'hbase\_test' selected under the 'Hbase' category. The main content area shows the 'List of Policies : hbase\_test' page. It features a search bar, an 'Add New Policy' button, and a table of policies. The table has columns for Policy ID, Policy Name, Policy Labels, Status, Audit Logging, Roles, Groups, Users, and Action. Two policies are listed: Policy ID 12 with name 'all - table, column-family, column' and Policy ID 20 with name 'UC1: deny cilucas and allow ekleszcz on de...'. Both policies are 'Enabled' and have 'Enabled' audit logging. The first policy is associated with 'ranger' and 'cilucas' users, while the second is associated with 'ekleszcz' and 'cilucas' users. The bottom of the page includes a license notice: 'Licensed under the Apache License, Version 2.0'.

Policy ID	Policy Name	Policy Labels	Status	Audit Logging	Roles	Groups	Users	Action
12	all - table, column-family, column	--	Enabled	Enabled	--	--	ranger cilucas	
20	UC1: deny cilucas and allow ekleszcz on de...	--	Enabled	Enabled	--	--	ekleszcz cilucas	



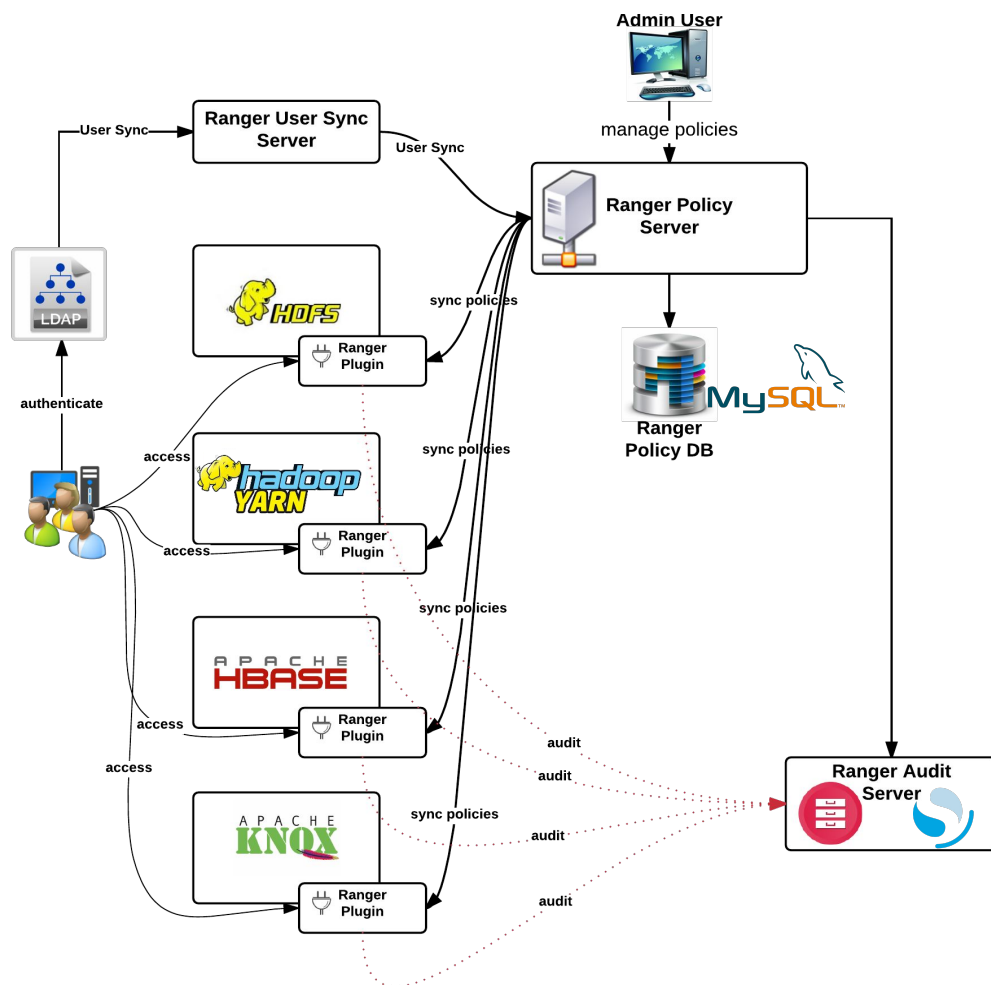
# Apache Ranger @ CERN

# Building the tool regular way

- Included only plugins for our components in use
- Tested against the master branch and latest release 2.4.0
- Build in CentOS 7 with Java 11

```
mvn clean compile package install -DskipTests=true -Dspotbugs.skip=true  
-Dchckstyle.skip=true -Dassembly.plugin.version=3.1.0 -Dhadoop.version=3.3.0  
-Dhbase.version=2.3.4 -Dzookeeper.version=3.6.1 -Dhive.version=3.0.0  
-Dmysql-connector-java.version=8.0.28  
-pl '!plugin-ozone, !plugin-solr, !plugin-nifi, !plugin-nifi-registry, !plugin-kudu,  
!plugin-kms, !ranger-ozone-plugin-shim, !storm-agent, !ranger-storm-plugin-shim,  
!ranger-solr-plugin-shim, !ranger-atlas-plugin-shim, !plugin-atlas, !plugin-kylin,  
!ranger-kylin-plugin-shim'
```

# Our setup



# Usersync plugin

- No performance overhead
- Integrates well with thousands of users and groups
- Responsive/robust
- Supports internal and external sync source
- Great for building Roles
- Integrated with the policies' definitions

# Usersync plugin

The screenshot displays the Apache Ranger web interface with the Usersync plugin active. The top navigation bar includes 'Resource Policies', 'Tag Policies', 'Audit', 'Security Zone', 'Settings', and 'Reports'. The user 'admin' is logged in. The left sidebar shows navigation options: 'User', 'Group', 'Role', and 'Permissions'. The main content area is titled 'Users/Groups/Roles' and shows a 'User List' table. A search bar is present above the table. The table columns are: User Name, Email Address, Role, User Source, Sync Source, Groups, Visibility, and Sync Details. The table lists several users, including 'admin', 'rangerusersync', 'ranger', 'hbase', 'ekleszcz', 'ivikin', 'avani', and 'padley'. Each user row has a checkbox, a role button, a user source button, a sync source button, a groups button, a visibility button, and a sync details button.

<input type="checkbox"/>	User Name	Email Address	Role	User Source	Sync Source	Groups	Visibility	Sync Details
<input type="checkbox"/>	admin		Admin	Internal	--		Visible	--
<input type="checkbox"/>	rangerusersync		Admin	Internal	--		Visible	--
<input type="checkbox"/>	rangertagsync		Admin	Internal	--		Visible	--
<input type="checkbox"/>	ranger		User	External	--	h-hadoop-analytix-users h-hadoop-development	Visible	--
<input type="checkbox"/>	hbase		User	External	--	h-hadoop-analytix-users users by home cernhome users by letter h	Visible	--
<input type="checkbox"/>	ekleszcz		User	External	LDAP/AD	3f-coffeeroom-rola acc-all-extra afe-resource-users ai-admins + More..	Visible	
<input type="checkbox"/>	ivikin		User	External	LDAP/AD	cms-mephi-tpi-mjtl russian-nationals users by home cernhome users by letter l	Visible	
<input type="checkbox"/>	avani		User	External	LDAP/AD	atlas-canada-ibx-production atlas-upgrade-ibx-stip-sensors atlas-women ibx-uoftl + More..	Visible	
<input type="checkbox"/>	padley		User	External	LDAP/AD	904-project bdg-32-4a british1-at-cern cms-b2g-17-010 + More..	Visible	



# *Apache Ranger*



## Plugin deployment

# Demo

The screenshot displays the Apache Ranger web interface. The top navigation bar includes 'Resource Policies', 'Tag Policies', 'Audit', 'Security Zone', 'Settings', and 'Reports'. The user 'admin' is logged in. The main content area shows the 'List of Policies : hdfs\_test' page. A yellow warning banner states: 'By default, fallback to HDFS ACLs are enabled. If access cannot be determined by Ranger policies, authorization will fall back to HDFS ACLs. If this behavior needs to be changed, modify HDFS plugin config - xasecure.add-hadoop-authorization.' Below this is a search bar and an 'Add New Policy' button. The main table lists 17 policies with columns for Policy ID, Policy Name, Policy Labels, Status, Audit Logging, Roles, Groups, Users, and Action.

Policy ID	Policy Name	Policy Labels	Status	Audit Logging	Roles	Groups	Users	Action
1	all - path	--	Enabled	Enabled	--	--	ranger	👁️ 📄 🗑️
2	kms-audit-path	--	Enabled	Enabled	--	--	keyadmin	👁️ 📄 🗑️
3	hbase-archive	--	Enabled	Enabled	--	--	hbase	👁️ 📄 🗑️
5	Personnal folder override	--	Enabled	Enabled	--	--	{USER}	👁️ 📄 🗑️
6	test against cilucas	--	Enabled	Enabled	--	--	clucas	👁️ 📄 🗑️
7	Personnal folder block	--	Enabled	Enabled	--	--	{USER}	👁️ 📄 🗑️
8	project access for cilucas	--	Enabled	Enabled	--	--	clucas	👁️ 📄 🗑️
9	Project test01 access	--	Enabled	Enabled	--	summ2023-supervisors-round1-assigned summer-students	--	👁️ 📄 🗑️
17	permissions for hdfs audit logs	--	Enabled	Enabled	--	--	yarn hbase	👁️ 📄 🗑️

# Setting resource policies

**Policy Details :**

Policy ID **14**

Policy Name \* finance directory **enabled**

Resource Path \*  **recursive** ← **Resource**

Description authorization for /finance directory contents

Audit Logging **YES**

**Allow Conditions :**

Select Group	Select User	Permissions	Delegate Admin
<input type="text" value="/finance"/>	<input type="text" value="Select User"/>	<input checked="" type="checkbox"/> Read <input checked="" type="checkbox"/> Write <input checked="" type="checkbox"/> Execute	<input checked="" type="checkbox"/>

Exclude from Allow Conditions :

**Deny Conditions :**

Select Group	Select User	Permissions	Delegate Admin
<input type="text" value="/interns"/>	<input type="text" value="Select User"/>	<input checked="" type="checkbox"/> Read <input checked="" type="checkbox"/> Write <input checked="" type="checkbox"/> Execute	<input checked="" type="checkbox"/>

Exclude from Deny Conditions :

**Save** **Cancel** **Delete**

**Ranger** Access Manager Audit Security Zone Settings admin

Service Manager > cm\_hdfs Policies > Create Policy

**Create Policy**

**Policy Details :**

Policy Type **Access** **enabled** **normal** **Add Validity Period**

Policy Name \*  **enabled** **normal**

Policy Label

Resource Path \*  **recursive**

Description

Audit Logging **YES**

**Allow Conditions :**

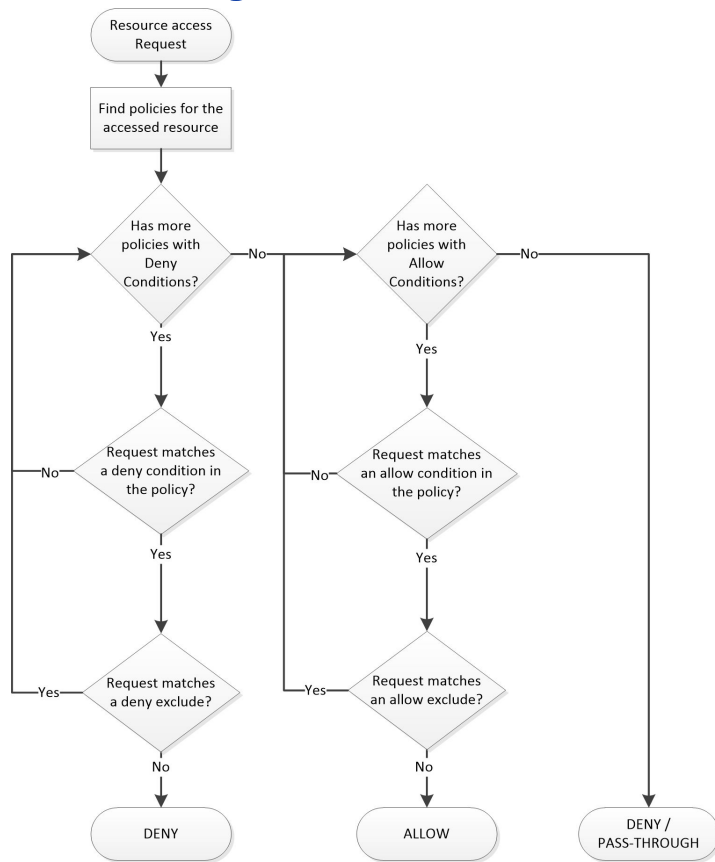
Select Role	Select Group	Select User	Delegate Admin
<input type="text" value="Select Roles"/>	<input type="text" value="Select Groups"/>	<input type="text" value="Select Users"/>	<input type="checkbox"/>

**add/edit permissions**

- Read
- Write
- Execute
- Select/Deselect All

**Add Permissions** +

# Resource-based policy access flow



# More details

- Support for autocompletion
- User can be set with one of the following roles: **KEYADMIN, ADMIN, USER**
- **ROLE** is assigned for a specific user by the Administrator
- **Delegated-Admin** permission allows other resource administrators to manage permissions for their managed-resources
- **USER** role gives ability to manage only resources for which the user has been granted with delegated-admin privilege
- Ranger provides **tag-based policies** too

# Example with 2 policies

## Policies :

1. Block access for a {USER} to /user/\*
2. Override to allow {USER} access to /user/{USER}



disabling an HDFS  
policy doesn't reset  
the previous state

## > kinit cllucas

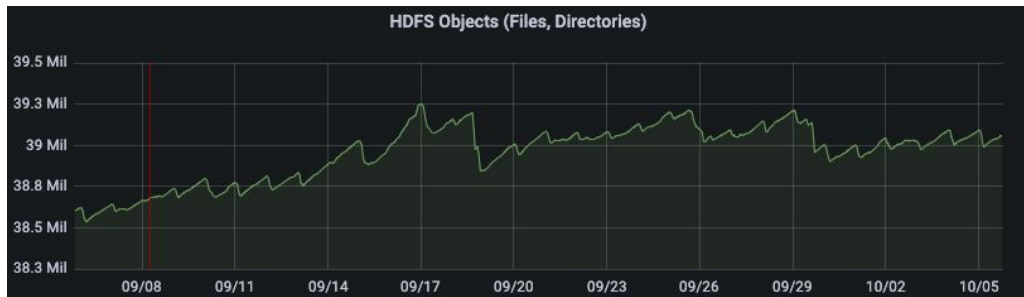
```
[root@ithdpdev-mmartinm01 ranger-2.4.0-hdfs-plugin]# hdfs dfs -ls /user/cllucas
Found 1 items
drwxr-xr-x  - cllucas supergroup          0 2023-09-11 17:08 /user/cllucas/.sparkStaging
```

```
[root@ithdpdev-mmartinm01 ranger-2.4.0-hdfs-plugin]# hdfs dfs -ls /user/ekleszcz
ls: Permission denied: user=cllucas, access=EXECUTE, inode="/user/ekleszcz"
```

# Our experience

- Works with HA clusters (setting namespace)
- No performance overhead
- Extra properties to be set for clusters with Kerb auth in web UI
- Lack of documentation is very painful
- *dfs.namenode.rpc-address.<cluster\_namespace>/<hostname>* must be updated to be done when aliases change (hosts upgraded)

```
[root@ekleszcz]# doAs hdfs hdfs dfs -ls /
Found 13 items
drwxr-xr-x - hbase hdfs 0 2017-06-18 19:57 /#
drwxr-xr-x - hdfs hdfs 0 2018-04-17 15:59 /Training
drwxr-xr-x - hdfs hdfs 0 2019-10-17 11:16 /Trash
drwxr-xr-x - hdfs zp 0 2023-04-05 09:33 /atlas
drwxrwxr-x+ - hdfs zh 0 2023-01-13 18:00 /cms
drwxrwx--x+ - hbase hdfs 0 2023-05-17 11:51 /hbase
drwxr-xr-x - ekleszcz hdfs 0 2023-08-18 15:09 /hbase-snapshots-perf-test
drwxr-xr-x - hdfs hdfs 0 2016-11-04 14:42 /lost+found
drwxr-xr-x - hdfs hdfs 0 2022-07-26 11:42 /project
drwxr-xr-x - hdfs hdfs 0 2023-10-05 19:00 /system
drwxrwxrwt - yarn hadoop 0 2023-10-05 12:02 /tmp
drwxr-xr-x+ - hdfs hdfs 0 2023-10-02 18:00 /user
drwxrwxr-x+ - hbase hadoop 0 2023-10-02 17:11 /var
```





# *Apache Ranger*



## Plugin deployment

# Yarn resources at CERN

## Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Me
194399	31	35	194333	400	5.38 TB	30.07 TB	7 GB

## Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy N
58	0	0	0	0

## Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation
Fair Scheduler	[memory-mb (unit=Mi), vcores]	<memory:2048, vCores:1>	<memory:81920, vCores:32>

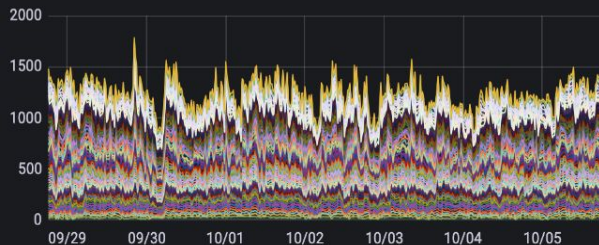
## Application Queues

Legend:  Steady Fair Share  Instantaneous Fair Share  Used  Used (over fair share)  Max Capacity

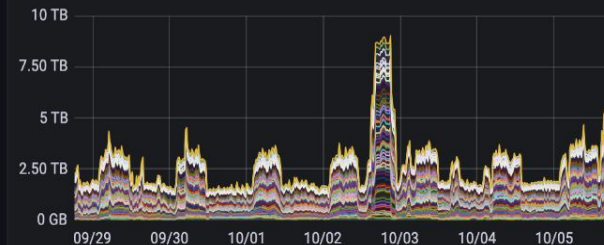
- root
- + root.certlogcompress
- + root.hconfig
- + root.monitops
- + root.certaenr
- + root.cperezde
- + root.itmonops
- + root.tapeops
- + root.atlevind
- + root.\_backup
- + root.hmonitor

## Yarn Fair scheduler

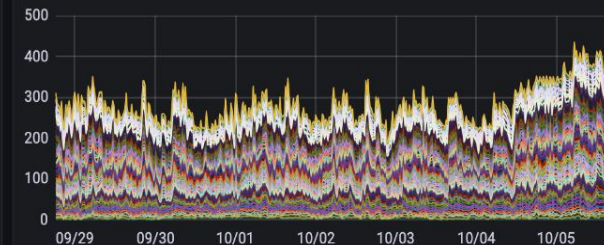
Allocated Cores



Allocated Memory



Containers Running



# Setting policies for a fair scheduler

queue name is: **root.{group}.{user}**

Resources :

Queue \*   Recursive

- **submit-app**: submit apps to the queue
- **admin-queue**: administer a queue (be able to kill an app)

Deny Conditions: hide

Select Role	Select Group	Select User	Policy Conditions	Permissions	Delegate Admin	
<input type="text" value="Select Roles"/>	<input type="text" value="x summer-students"/>	<input type="text" value="x clucas"/>	<i>Add Conditions</i> <input type="button" value="+"/>	<input type="button" value="submit-app"/> <input type="button" value="admin-queue"/>	<input type="checkbox"/>	<input type="button" value="x"/>



*Apache Ranger*

A P A C H E  
**HBASE**

**Plugin deployment**

# HBase ACLs management

Around 10k tables in multiple namespaces

```
hbase(main):006:0> user_permission '@ekleszcz'  
User                               Namespace,Table,Family,Qualifier:Permission  
[REDACTED]                          ekleszcz,,,: [Permission: actions=READ,WRITE,EXEC,CREATE,ADMIN]  
Took 0.0462 seconds
```

```
hbase(main):002:0> user_permission  
User                               Namespace,Table,Family,Qualifier:Permission  
[REDACTED]                          ,,,: [Permission: actions=READ,WRITE,EXEC,CREATE,ADMIN]  
[REDACTED]                          ,,,: [Permission: actions=CREATE]  
[REDACTED]                          ,,,: [Permission: actions=READ,WRITE,EXEC,CREATE]  
[REDACTED]                          ,,,: [Permission: actions=READ,WRITE,EXEC,CREATE,ADMIN]  
[REDACTED]                          ,,,: [Permission: actions=READ,WRITE,EXEC,CREATE,ADMIN]  
Took 0.4352 seconds
```

# Setting policies



disabling an HBase policy  
doesn't reset the previous state

{namespace}:{table}

Resources :

HBase Table *	d	<input type="checkbox"/>
HBase Column-family *	another_test test test2	<input type="checkbox"/>
HBase Column *	x *	<input type="checkbox"/>

Resources :

HBase Table *	x test	<input type="checkbox"/>
HBase Column-family *	c	<input type="checkbox"/>
HBase Column *	cf1 colFam1 colFam2 colFam3	<input type="checkbox"/>

## Permissions :

- Read
- Write
- Create
- Admin
- Execute

Allow Conditions:

Select Role	Select Group	Select User	Policy Conditions	Permissions	Delegate Admin
Select Roles	Select Groups	x ekleszcz	Add Conditions +	Read Write Create Admin Execute	add/edit permissions <input checked="" type="checkbox"/> Read <input checked="" type="checkbox"/> Write <input checked="" type="checkbox"/> Create <input checked="" type="checkbox"/> Admin <input checked="" type="checkbox"/> Execute <input checked="" type="checkbox"/> Select/Deselect All

+  
⚠ Exclude from Allow Conditions:

# Auditing

**Ranger** | Resource Policies | Tag Policies | Audit | Security Zone | Settings | Reports | admin

AUDIT <<

Access

**Admin**

Login Sessions

Plugins

Plugin Status

User Sync

Last Response Time : 09/12/2023 02:21:42 PM

Search for your access logs...

Last Updated Time: 09/12/2023 02:21:36 PM | Entries: 1 to 25 of 56805

Operation	Audit Type	User	Date ( heure d'été d'Europe centrale )	Actions	Session ID
User updated <b>clucas</b>	Ranger User	admin	09/12/2023 10:08:55 AM	<a href="#">Update</a>	674
User updated <b>clucas</b>	Ranger User	admin	09/12/2023 09:23:53 AM	<a href="#">Update</a>	664
User updated <b>clucas</b>	Ranger User	admin	09/12/2023 09:19:49 AM	<a href="#">Update</a>	658
User profile updated <b>clucas</b>	User Profile	admin	09/12/2023 09:19:49 AM	<a href="#">Update</a>	658
Policy updated <b>UC1: deny user clucas on root.it-hadoop-a...</b>	Ranger Policy	admin	09/11/2023 04:55:59 PM	<a href="#">Update</a>	626
Policy updated <b>UC1: deny user clucas on root.it-hadoop-a...</b>	Ranger Policy	admin	09/11/2023 04:55:53 PM	<a href="#">Update</a>	626
Policy updated <b>UC2: deny summer-students in root.it-had...</b>	Ranger Policy	admin	09/11/2023 04:55:33 PM	<a href="#">Update</a>	626
Policy updated <b>UC2: deny summer-students in root.it-had...</b>	Ranger Policy	admin	09/11/2023 04:00:12 PM	<a href="#">Update</a>	626
Policy updated <b>UC2: deny summer-students in root.it-had...</b>	Ranger Policy	admin	09/11/2023 03:42:42 PM	<a href="#">Update</a>	626
Policy updated <b>all - table, column-family, column</b>	Ranger Policy	admin	09/11/2023 02:51:44 PM	<a href="#">Update</a>	623
Policy updated <b>UC1: deny clucas and allow ekleszcz on d...</b>	Ranger Policy	admin	09/11/2023 02:51:32 PM	<a href="#">Update</a>	623
Policy updated <b>all - table, column-family, column</b>	Ranger Policy	admin	09/11/2023 02:51:23 PM	<a href="#">Update</a>	623
Policy updated <b>UC1: deny clucas and allow ekleszcz on d...</b>	Ranger Policy	admin	09/11/2023 02:51:13 PM	<a href="#">Update</a>	623
Policy updated <b>all - table, column-family, column</b>	Ranger Policy	admin	09/11/2023 02:51:01 PM	<a href="#">Update</a>	623

# Auditing: logs kept in HDFS

Stored in /ranger/audit/{plugin}/{date:YYYYMMDD}/{filename}

Example : `hdfs dfs -cat`

```
{/ranger/audit/hdfs/20230912/hdfs ranger audit ithdpdev-mmartinm02.cern.ch.log
```

```
"repoType":1,  
"repo":"hdfs_test",  
"reqUser":"ekleszcz",  
"evtTime":"2023-09-12 15:11:43.818",  
"access":"READ_EXECUTE",  
"resource":"/ranger/audit",  
"resType":"path",  
"action":"read",  
"result":1,  
"agent":"hdfs",  
"policy":-1,  
"reason":"/ranger/audit",  
"enforcer":"hadoop-acl",  
"cliIP":"188.184.11.111",  
"agentHost":"ithdpdev-ekleszcz12.cern.ch",  
"logType":"RangerAudit",  
"id":"bdd67338-b9c2-490f-81e2-66304d5c2673-0",  
"seq_num":1,  
"event_count":1,  
"event_dur_ms":0,  
"tags":[],  
"additional_info":{"forwarded-ip-addresses":"","  
"remote-ip-address":"","188.184.11.111","accessTypes":["execute, read]}",  
"cluster_name":""  
}
```

Result of the request

1 is allowed, 0 is denied

Policy ID used to process the request

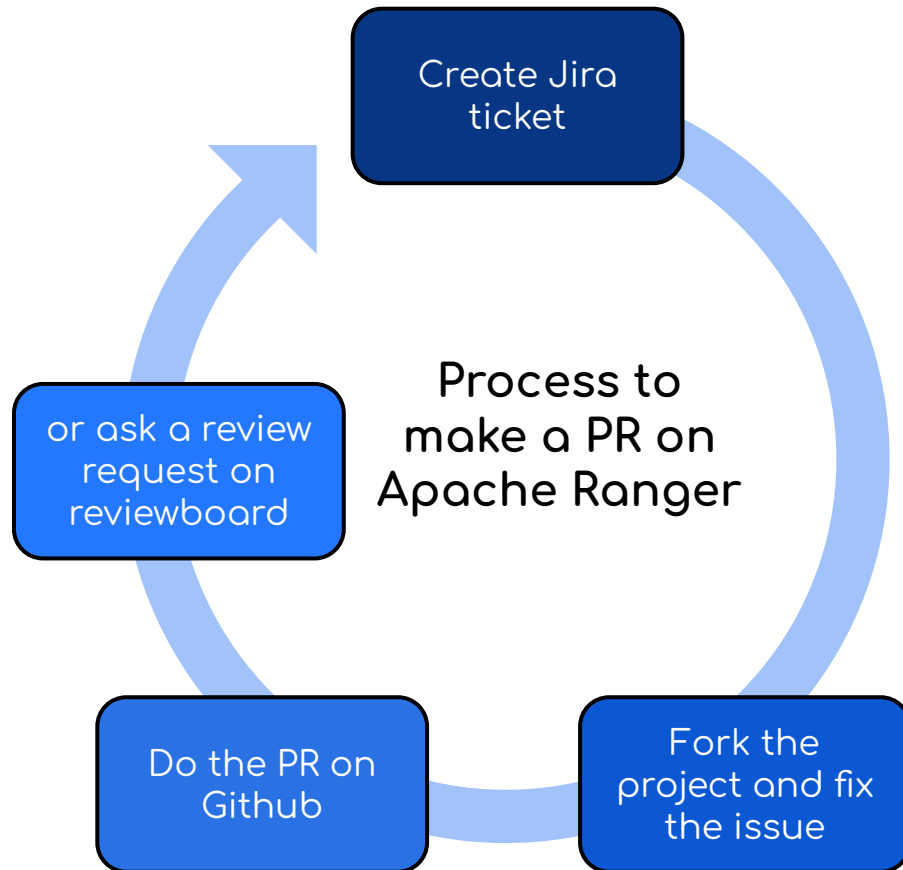
-1 means that no Ranger policy cover the request

# Our experience and concerns

- Documentation
  - Legacy or missing or linked to Ambari or Cloudera but not standalone
  - Little activity for reported issues in GitHub/ mailing list
- Compatibility
  - Certain component versions must be respected - no docs
  - Last official release 2.4.0 from 2017
- Integrations
  - Auditing not supported for OpenSearch, only ES 7.10.2
  - Some missing dependencies in yarn plugin (jars missing, eg. apache-commons-lang3)
- Resource intensive - 1G+ of heap needed with certain plugins enabled
- Pull mechanism only - no enforcement of the current state
- Yarn autocompletion didn't work with fair-scheduler (?)

# Upstream contributions

- Clear way to contribute
- **A contribution submitted** to upstream but not followed up
- A few issues reported to the mailing list
  - ([dev@ranger.apache.org](mailto:dev@ranger.apache.org))

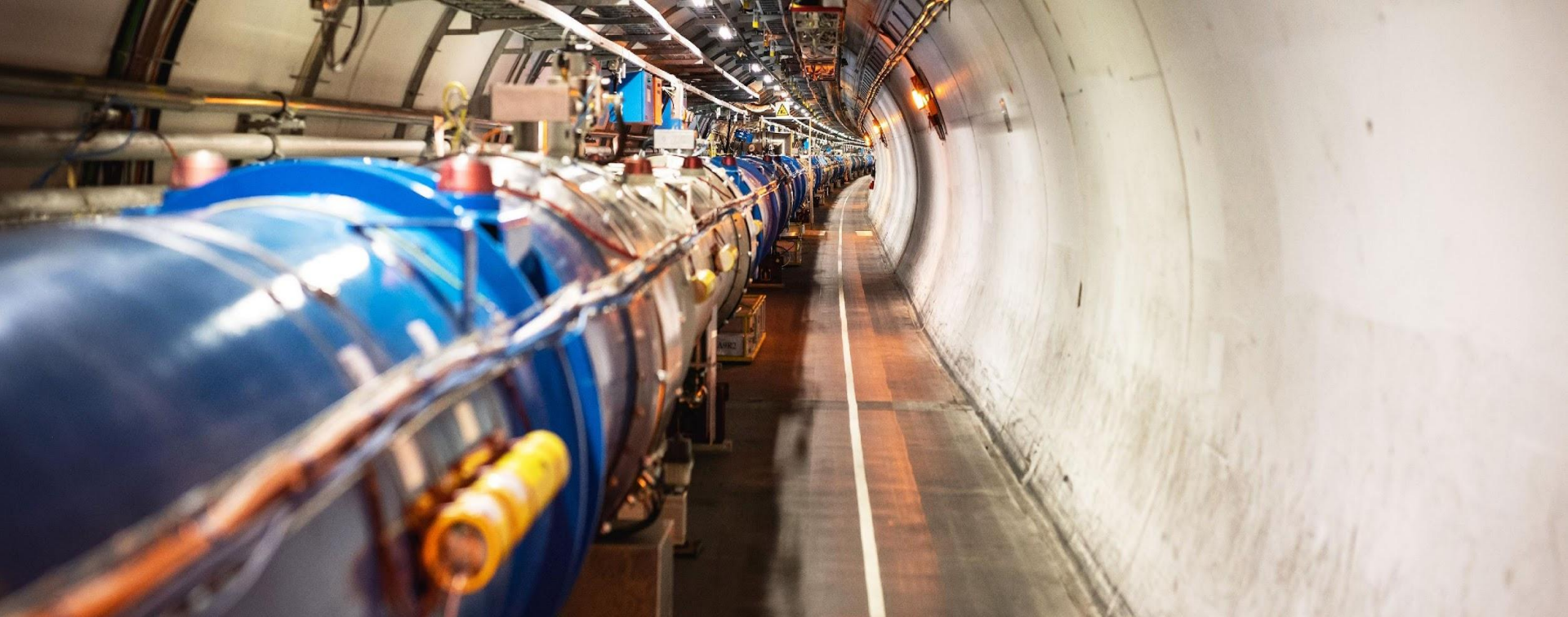


# Future plans



**Apache Ranger**

- Puppetize and deploy to PROD
- SSO integration (Apache KNOX)
- Auditing with OpenSearch backend
  - We moved away from ElasticSearch/OpenDistro
- Outsource mgmt. to project managers via Roles
- Explore HA deployment
- Fix autocompletion for Apache YARN
- Dynamic management of the configs



# Summary

# Summary

- At CERN, **we collect 100s of PBs of data** from various physics experiments
- We have a lot of users who access the data from our Big Data services such as HDFS
- **We need to protect this data and optimize the access rights management**
- Also, we lack a **centralized auditing system**
- **Apache Ranger is a proven-to-work mature project that integrates well** with our main components and provides the solution for our needs
- With Ranger, we faced **a few issues that still need to be tackled**, such as
  - **integration with OpenSearch** for auditing
  - **autocompletion** for the **YARN** plugin
- The tool will be deployed into the prod clusters

# Thank you !



Contact: [emil.kleszcz@cern.ch](mailto:emil.kleszcz@cern.ch)

[home.cern](http://home.cern)