



A Big Data Platform for heterogeneous data collection and analysis in large-scale data centers

S. ROSSI TISBENI

A. CARBONE, C. CAVALLARO, D. CESINI, D. C. DUMA, A. FALABELLA, E. FURLAN, M. GALLETTI, J. GASPARETTO, B. MARTELLI, D. MICHELOTTO, F. MINARINI, L. MORGANTI, E. RONCHIERI, D. SALOMONI, G. SERGI.

INFN-CNAF

CNAF is the national center of INFN (Italian Institute for Nuclear Physics) for **Research and Development on Information and Communication Technologies**

- Focus on **Software development and distributed systems**: development of **Grid** middleware for WLCG and of **CLOUD** technologies (coordination of the distributed INFN Cloud infrastructure)
- Hosts the Italian **Tier-1** data center for HEP experiments
 - Provides resources, support and services needed for data storage, distribution and processing 36PB of data on disk storage, 86PB of data on tape storage, 400k HS06
- Hosts the **INFN-TTLab** that through Technology Transfer activities provides industrial partners with resources and know-how

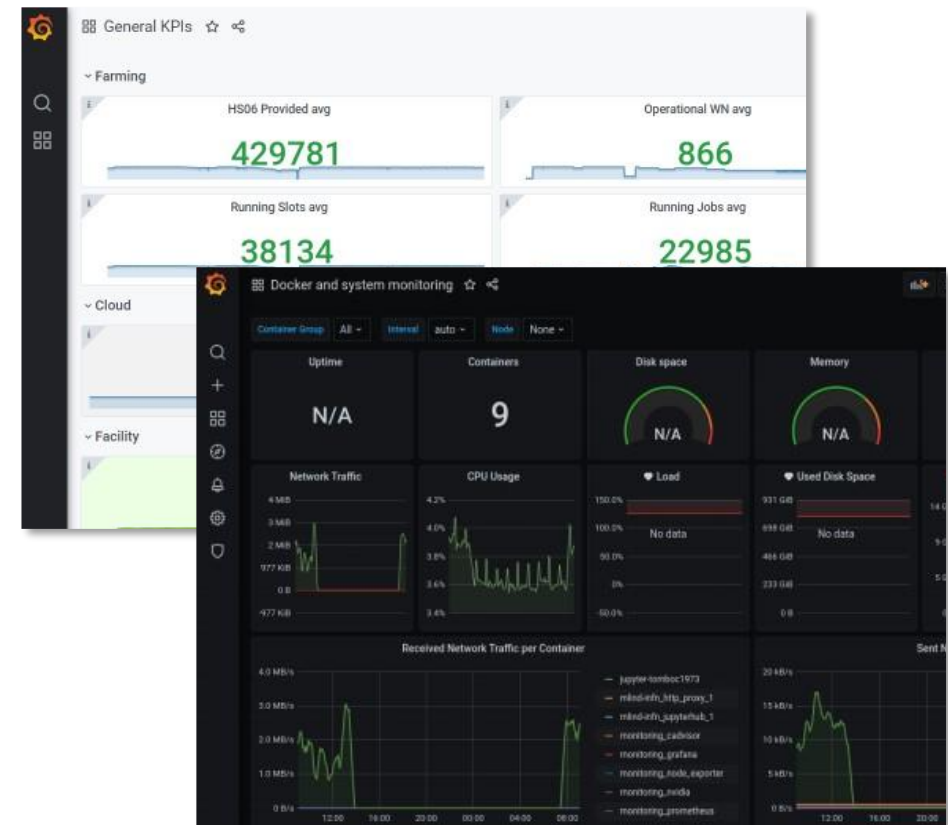
Context (1/2)

INFN Tier-1 datacenter

- Monitoring system based on Sensu, InfluxDB, Grafana
 - General datacenter KPIs
 - Specific KPIs for batch system, storage, network and cloud services
 - Specific dashboards for critical services based on Elastic stack

INFN Cloud infrastructure across several sites

- Monitoring system based on Zabbix, ELK, Grafana
 - Cloud infrastructure monitoring and accounting
 - Tenant VMs and applications monitoring

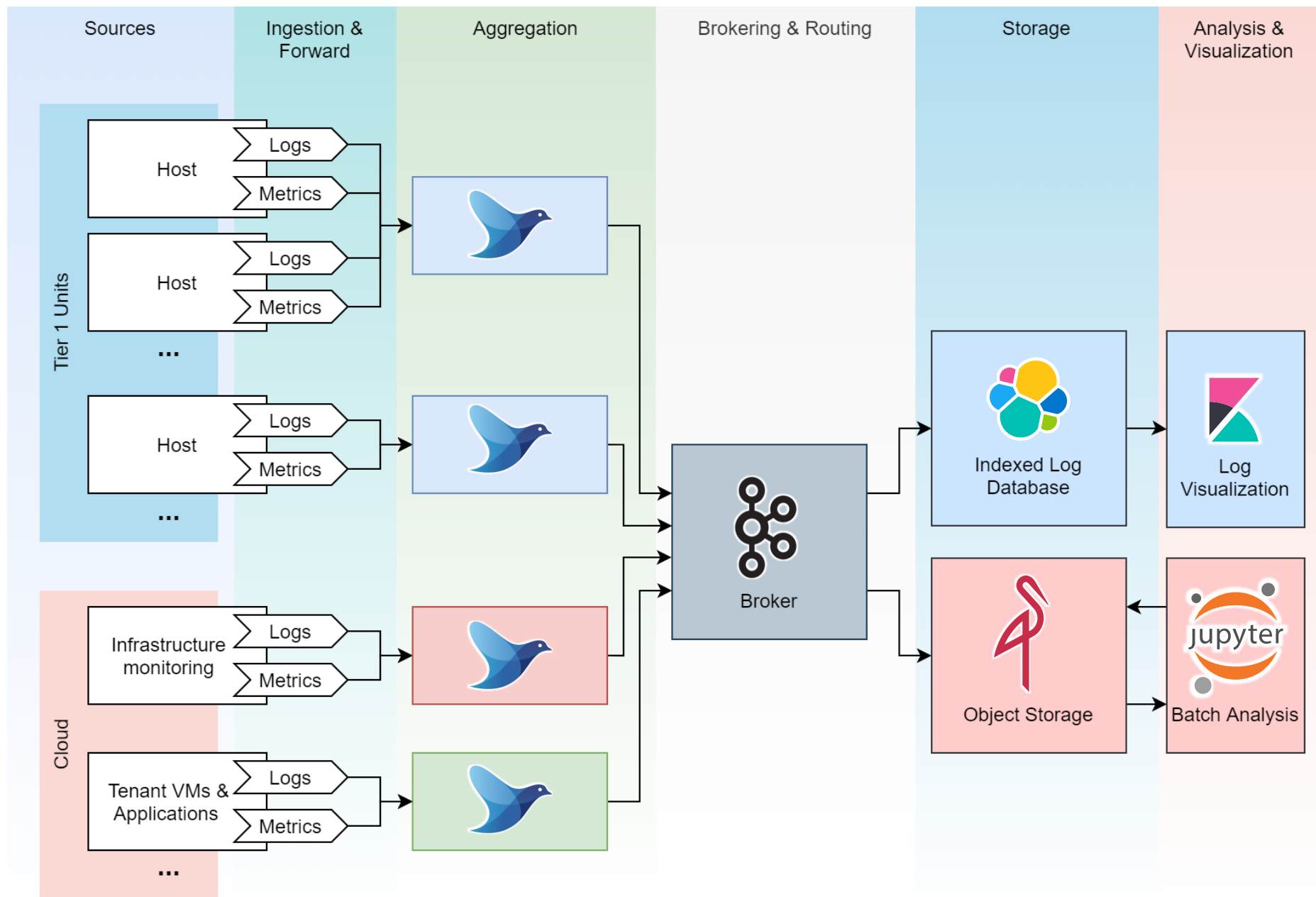


Context (2/2)

- IoTwins EU project to create a big data platform for optimized and replicable industrial and facility management models
 - Methodology for the monitoring infrastructure reuse and deployment in new and different contexts
 - Focus on monitoring of IoT/Edge/Cloud integrated infrastructure
 - Predictive analysis of fault and support on troubleshooting
 - <https://iotwins.eu>
- EPIC Cloud (for applications with high security requirements)
 - Alarm system based on Sensu/Uchiwa/Slack
 - Monitoring system based on Sensu/Grafana/InfluxDB/Prometheus
 - In the future a system for monitoring security events and enrich these data with a remote threat intelligence source will be probably needed

Big Data Platform / Goals

- Centralize **collection** of logs, metrics and data from **heterogeneous** data sources and present them in a **flexible and actionable** format
- Harmonize the current **monitoring** infrastructures in production at the center
- Provide different storage solution with proper retention policies based on security and privacy requirements
- Provide policy-based, **authenticated** access to the data
- Provide sysadmins and users with tools for **information discovery** from system data and monitoring
- Develop a reliable, extensible, scalable and manageable platform for collection and analysis of **Big Data** to be offered as a service



Ingestion & Aggregation



- **Filebeat** for log ingestion
 - Tail on file with variable structure and multiline support
 - Can be deployed as puppet module on production machines
 - Rsyslog already in use for critical service
- **Fluentd** for aggregation, format, and filtering
 - Small footprint, easy to install, open source
 - JSON native
 - Flexible in refactor and filtering
 - Supports various sources and outputs via Plugins

Local persistent cache with buffer

Guarantee At Least Once delivery to the broker



Broker

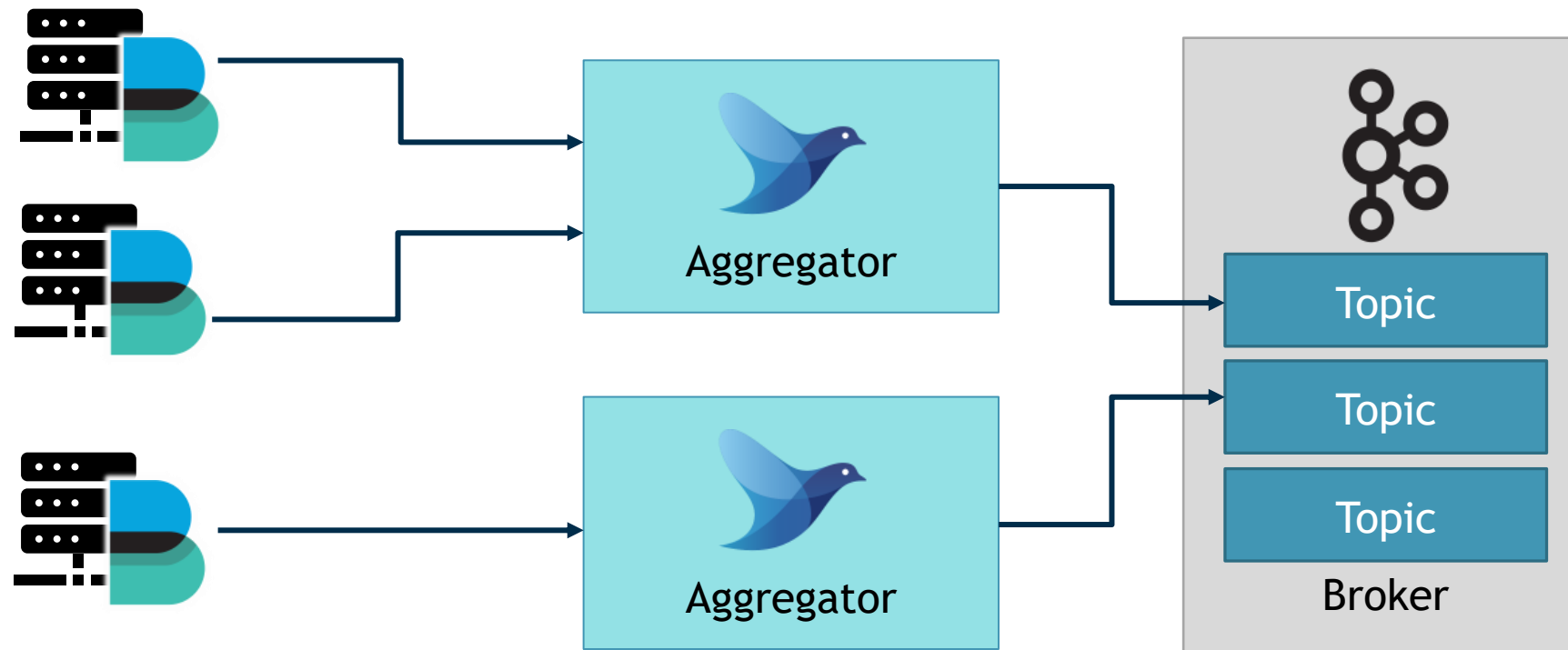
Apache Kafka, topic based publish-subscribe data distribution system

- Decouple data ingestion and consumption
- Allow broadcasting data to multiple subscribers simultaneously
- Allow for fault tolerance and high performance

Proposed setup:

- 3 brokers with 3 replicas for each partition
- Retention set to 3 days

Architecture / Ingestion



Storage

The **Elastic Stack** acts as Short-Term Storage database for log files

- **Logstash:**
 - Routing data from Kafka to the Elastic Index
 - Wide use of filters to parse information from logs
 - At Least Once buffered delivery
- **Elasticsearch:**
 - Performs indexing of data
 - Stores logs as JSON objects
 - Expose SQL like interface for queries
- **Kibana:**
 - Discovery and visualization through dashboard

[PoS\(ISGC2019\)027](#)

Storage

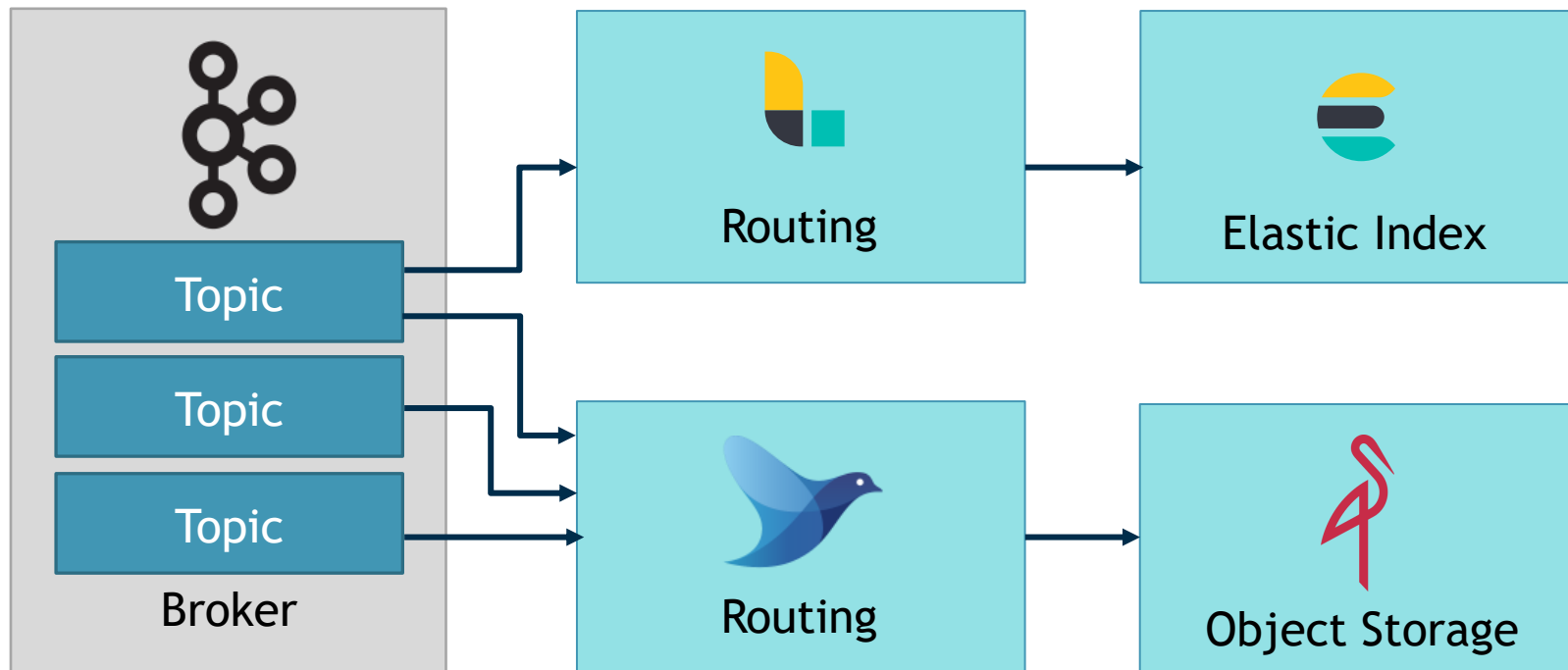


MinIO as Long-Term Object storage solution

- Disaggregate storage and computing resources
- Cloud native architecture, can be federated with other clusters
- High-performance, bucket-based object storage
- Policy based, OIDC Authenticated access through INDIGO IAM
- S3 and REST API for data access

Integrated with storage technologies from INFN Cloud

Architecture / Consumers



Batch Analysis



- **Jupyter**
 - Provides a web-based frontend for users to access the data
 - Open Source
 - Supports multiple programming languages (Python, Scala, Java)
 - Simultaneous access to host from multiple users
 - Integration with OpenID connect through INDIGO IAM
- **Spark:**
 - Supports high-performance batch analysis with map/reduce paradigm
 - Access to diverse data sources, including time-series database and s3 Object storage
 - Run natively on cloud in orchestrated clusters
 - Supports Python, Scala and R
 - Powers dataframe handling and machine learning libraries

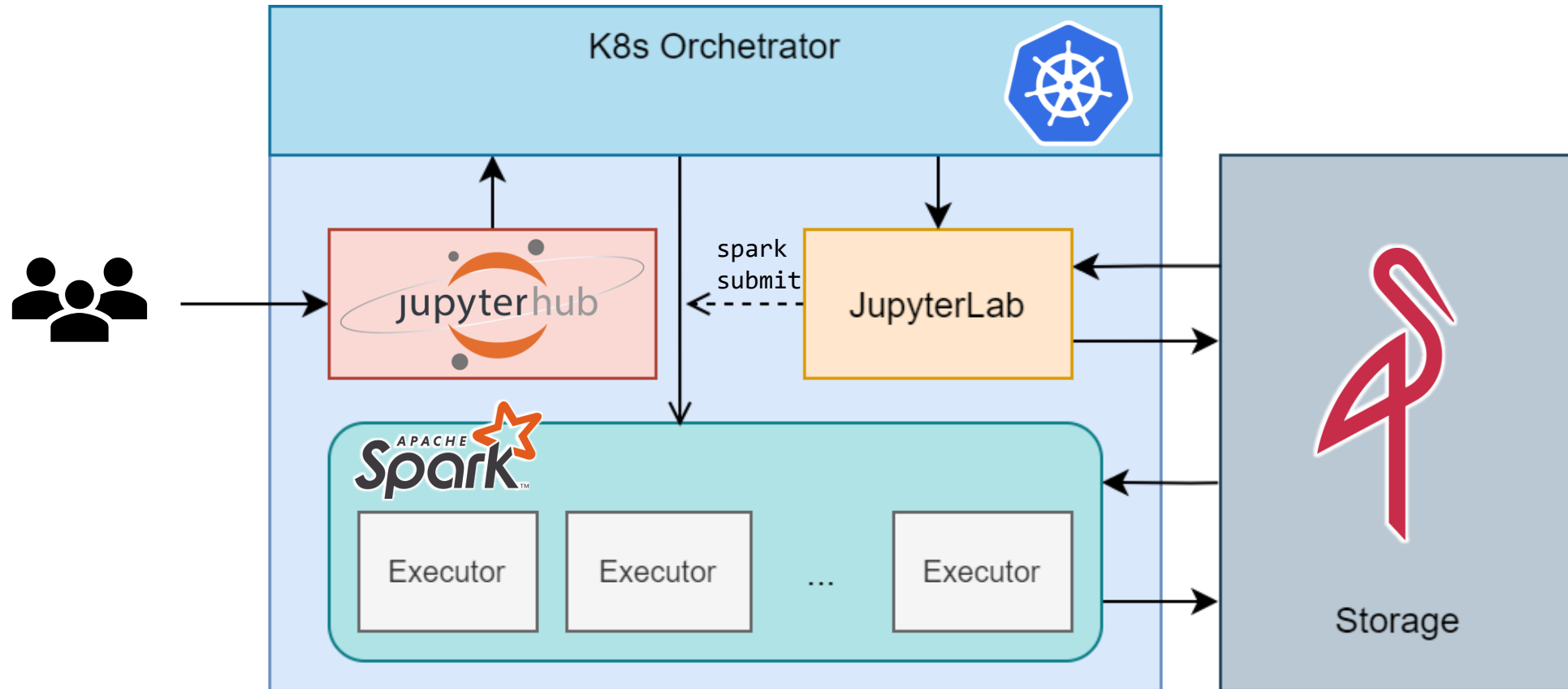
Orchestration

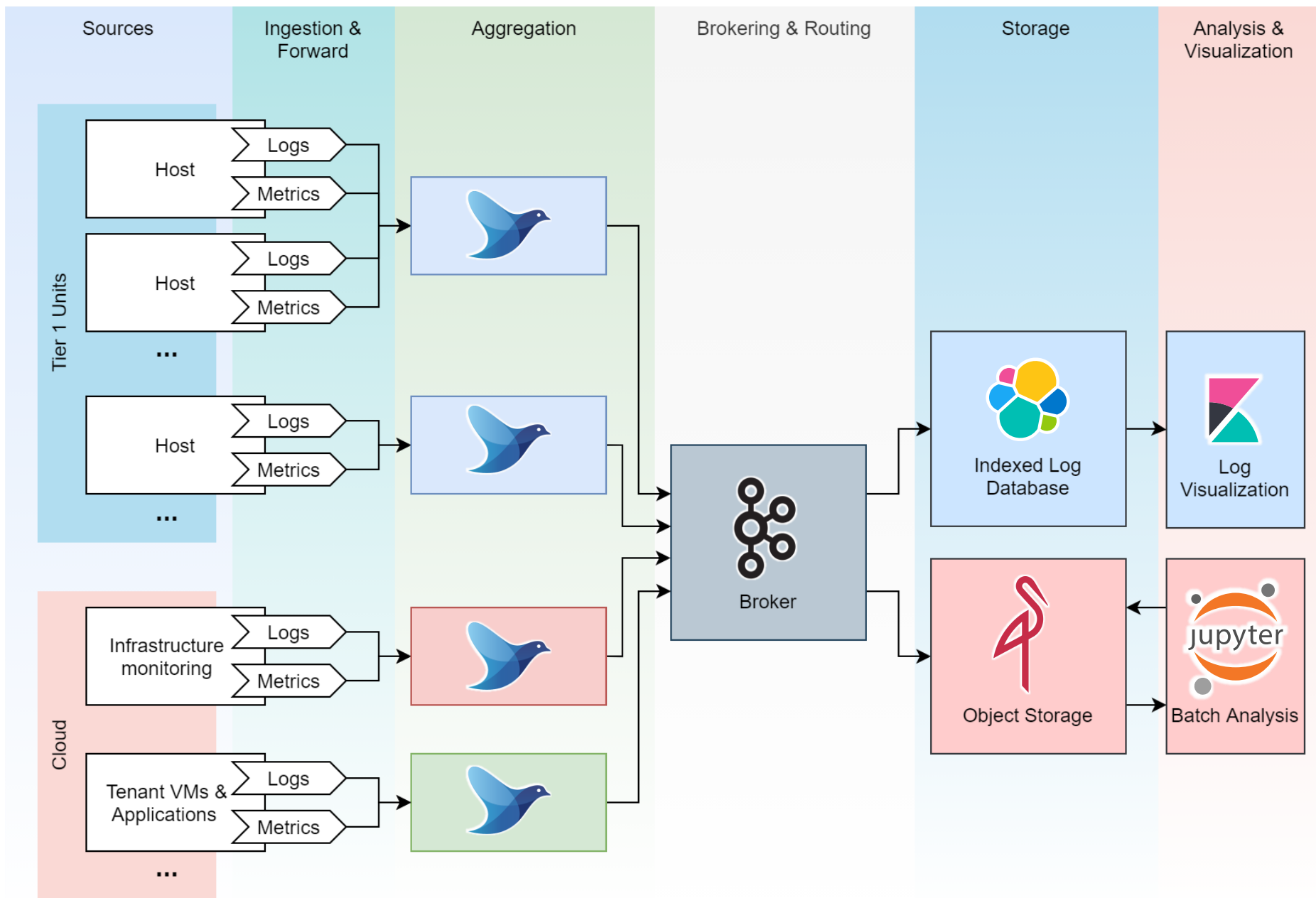


Kubernetes: modern orchestrator for containerized services

- Manages the deployment of JupyterLab single-user containers (spawned by JupyterHub)
- Allows to allocate different virtual resources (CPUs, RAM) to each user
- Allows each user to spawn containerized environment through JupyterHub
- Use sandbox containers as users workspace
- Manages the deployment of Spark executors with fine control on the number of executors and resources for each executor (CPU cores, RAM)
- ReplicaSet guarantee: a new pod is automatically instantiated on the next available node with almost zero downtime if a cluster node goes down

Architecture / Batch Analysis





Sample Workflow

```
Mar 19 07:56:30 ds-816 storm-webdav-server-access: odC.physics_Main.PhysCont_DAOD.t0pro22_v01_a21296_20201123.log.23644713.017554.log.tgz HTTP/1.1" 207 1075>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 2001:1458:201:e3:0:0:100:211 - - - 2021-03-19T06:56:27.135Z "0b071177-5e46-4fd2-a0a9-58bdf1cc7ae0" "DELET>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 188.184.75.208 - - - 2021-03-19T06:56:27.145Z "868485c5-2990-4a20-b571-c962166ac3f2" "DELETE /atlas/atlas>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 188.185.82.95 - - - 2021-03-19T06:56:27.158Z "bcd65495-c459-4f0c-9ddc-162e74a5e3ab" "DELETE /atlas/atlass>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 188.185.82.95 - - - 2021-03-19T06:56:27.159Z "c4cad37b-3160-49d9-af33-22bd6b88ae63" "DELETE /atlas/atlass>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 188.184.75.208 - - - 2021-03-19T06:56:27.176Z "9bb60125-229e-4bb5-a9cf-1ca82b1b9943" "PROPFIND /atlas/atla>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 2001:1458:201:e3:0:0:100:211 - - - 2021-03-19T06:56:27.176Z "aec1fcac-0605-4b28-916e-6a15e9c6f12e" "PROPF>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 188.185.82.95 - - - 2021-03-19T06:56:27.200Z "f6dd121e-7349-4e63-a022-06f96354b382" "PROPFIND /atlas/atla>
Mar 19 07:56:30 ds-816 storm-webdav-server-access: 188.184.75.208 - - - 2021-03-19T06:56:27.202Z "67bac451-aa5d-4dc4-8249-3098624aeedd" "DELETE /atlas/atlas>
:|
```



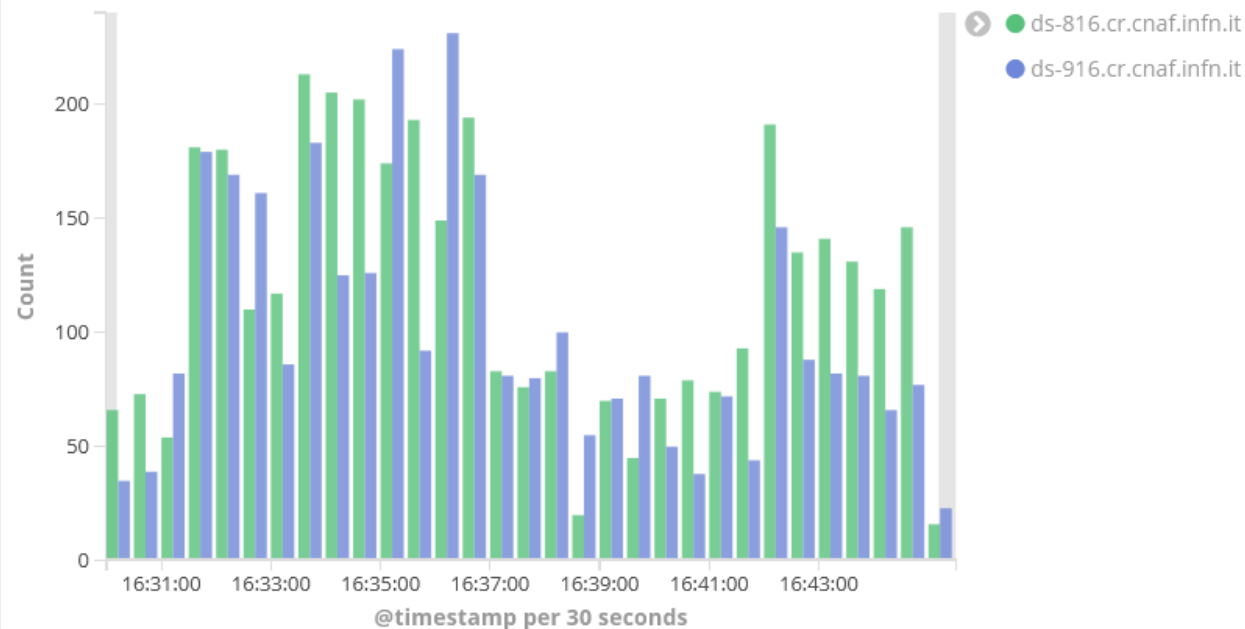
Time ▾	hostname	name	raw	http_code	path	wd_method
> Mar 19, 2021 @ 07:56:30.000	ds-816	storm-webdav-server-access	odC.physics_Main.PhysCont_DAOD.t0pro22_v01_a21296_20201123.log.23644713.017554.log.tgz HTTP/1.1" 207 1075 4	207	-	-
> Mar 19, 2021 @ 07:56:29.761	ds-816	storm-webdav-server-access	2001:1458:301:cd:0:0:100:32f - - - 2021-03-19T06:56:29.761Z "f362cee8-b9f4-418a-90d1-34ec35900820" "GET /.well-known/oauth-authorization-server HTTP/1.1" 200 248 3	200	/.well-known/oauth-authorization-server	GET
> Mar 19, 2021 @ 07:56:29.747	ds-816	storm-webdav-server-access	188.185.82.95 - - - 2021-03-19T06:56:29.747Z "1940a4ba-6e54-4892-9ba8-8d731bf35d6e" "DELETE /atlas/atlasscratchdisk/rucio/user/sgurdasa/da/f0/user.sgurdasa.23649449._013534.output.root HTTP/1.1" 204 0 10	204	/atlas/atlasscratchdisk/	DELETE

Sample Workflow

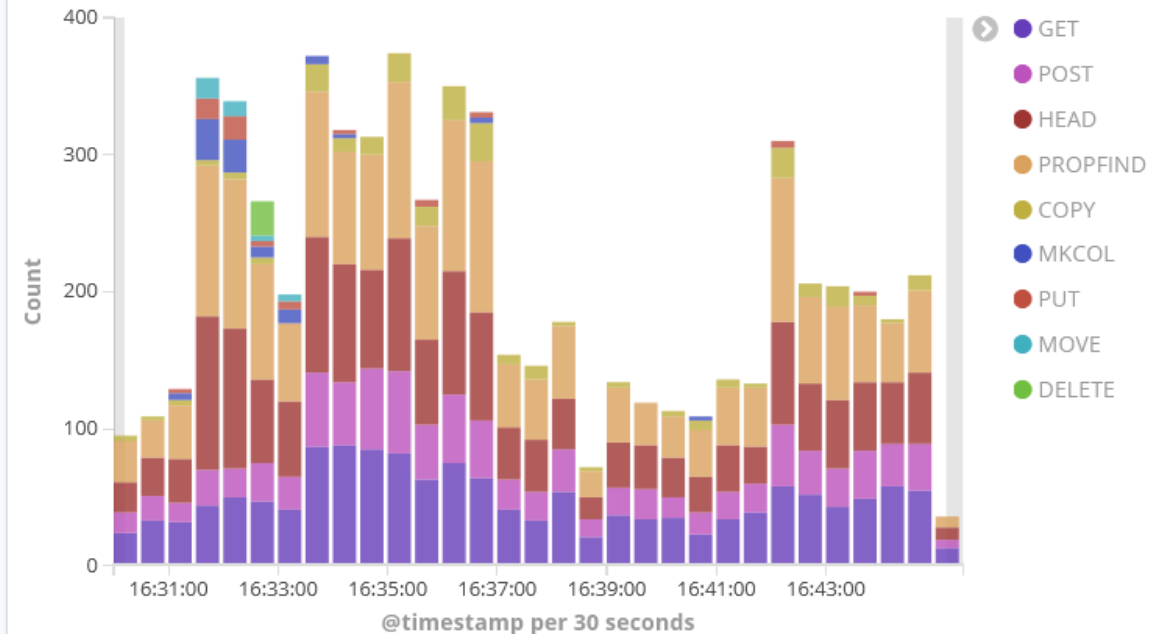
tags.keyword: "atlas"

Add a filter +

Request per WebDAV endpoints



Methods requested to WebDAV



Sample Workflow



Connect to MinIO

Read Logs

```
[3]: conf = SparkConf().set("spark.executor.cores", "2")
spark = SparkSession.builder.getOrCreate()
sc = SparkContext.getOrCreate(conf)
sc
```

[3]: SparkContext

Spark UI

Version

v3.1.1

Master

k8s://https://kubernetes:443

AppName

pyspark-shell

```
[5]: df = spark.read.json("s3a://storage-enr/ds-916/2021/03/20/storm-webdav-server-*.json")
df.toPandas()
```

	@timestamp	@version	hostname	http_code	name	path	raw	throughput	topic	wd_method
0	2021-03-20T07:13:39.685Z	1	ds-916	200	storm-webdav-server-access	/oauth/token	2001:1458:301:cd:0:0:100:32f - - 2021-03-20T...	NaN	storage-test-enr	POST
1	2021-03-20T07:14:12.146Z	1	ds-916	207	storm-webdav-server-access	/atlas/atlasdatadisk/	2001:1458:301:cd:0:0:100:1e3 - - 2021-03-20T...	NaN	storage-test-enr	PROPFIND
2	2021-03-20T07:13:08.402Z	1	ds-916	200	storm-webdav-server-access	/status/metrics?pretty=true	131.154.130.132 - - 2021-03-20T07:13:08.402Z...	2.978875e+06	storage-test-enr	GET
3	2021-03-20T07:13:19.107Z	1	ds-916	200	storm-webdav-server-access	/.well-known/oauth-authorization-server	2001:1458:301:cd:0:0:100:a0 - - 2021-03-20T0...	8.266667e+04	storage-test-enr	GET
4	2021-03-20T07:13:27.496Z	1	ds-916	200	storm-webdav-server-access	/.well-known/oauth-authorization-server	2001:1458:301:cd:0:0:100:2e - - 2021-03-20T0...	6.200000e+04	storage-test-enr	GET
...
101141	2021-03-20T03:07:08.674Z	1	ds-916	200	storm-webdav-server-access	/atlas/atlasdatadisk/	2001:1458:301:cd:0:0:100:313 - - 2021-03-20T...	NaN	storage-test-enr	HEAD
101142	2021-03-20T03:07:10.332Z	1	ds-916	404	storm-webdav-server-access	/atlas/atlasdatadisk/	2001:1458:301:cd:0:0:100:a0 - - 2021-03-20T0...	NaN	storage-test-enr	HEAD
101143	2021-03-20T03:06:58.252Z	1	ds-916	200	storm-webdav-server-access	/oauth/token	2001:1458:301:cd:0:0:100:305 - - 2021-03-20T...	NaN	storage-test-enr	POST
101144	2021-03-20T03:06:58.455Z	1	ds-916	200	storm-webdav-server-access	/oauth/token	2001:1458:301:cd:0:0:100:305 - - 2021-03-20T...	NaN	storage-test-enr	POST
101145	2021-03-20T03:06:58.538Z	1	ds-916	404	storm-webdav-server-access	/atlas/atlasdatadisk/	2001:1458:301:cd:0:0:100:305 - - 2021-03-20T...	NaN	storage-test-enr	PROPFIND

101146 rows x 10 columns

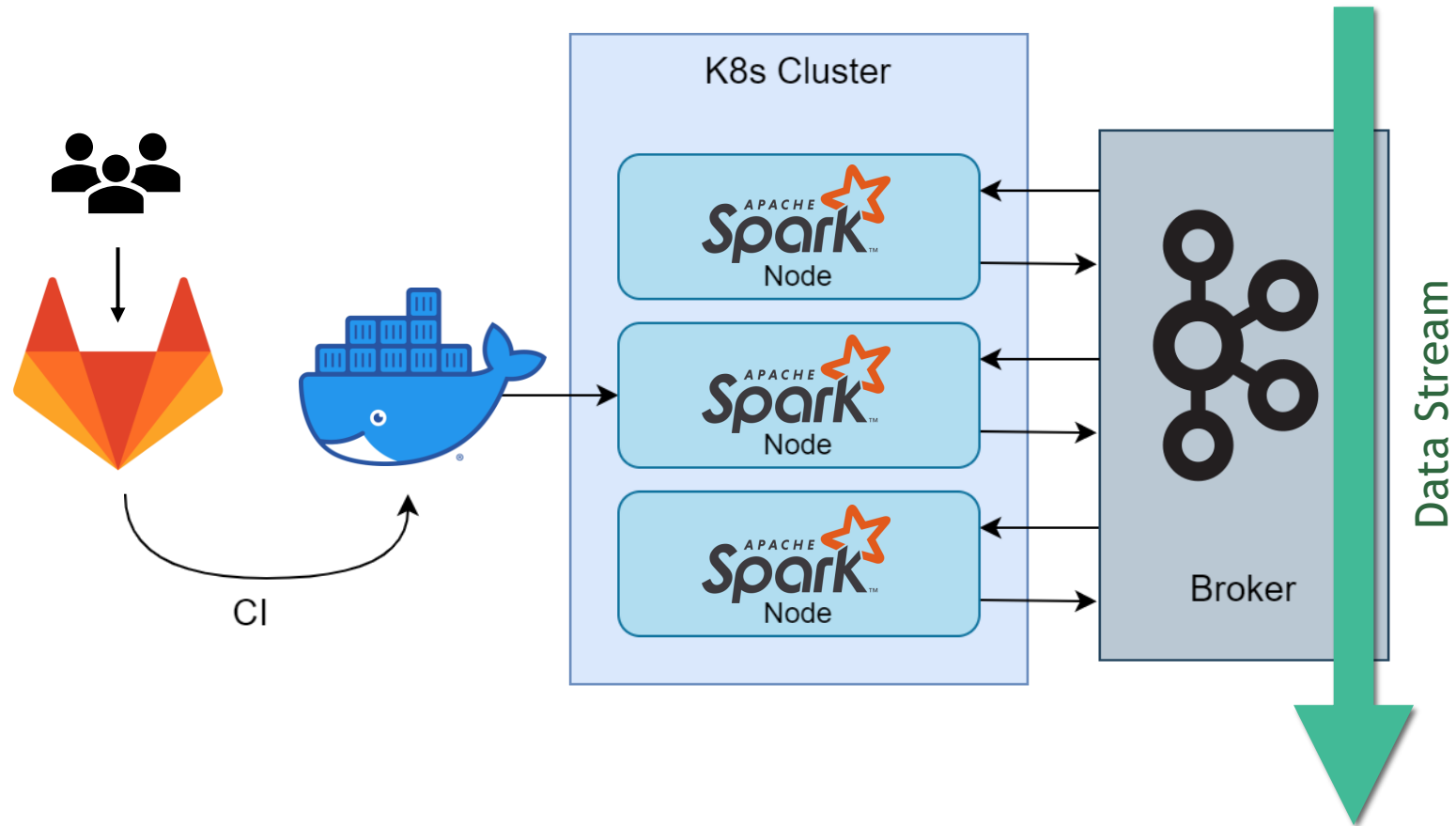
Next Steps / Anomaly Detection use case

In the context of the IoTwins project, replicate Anomaly Detection use case from CINECA datacenter on CNAF data

- Parse unstructured log data to find possible correlations with anomalies
- Perform clusterization to identify clusters on the base of keywords
- Use the result of the clusterization to perform semi-supervised anomaly detection
- Perform the analysis in streaming to label the new entries as they are included in the dataset

This analysis would greatly benefit from the Big Data Platform

Next Steps / Streaming Analysis



Next Steps / Metrics Integration

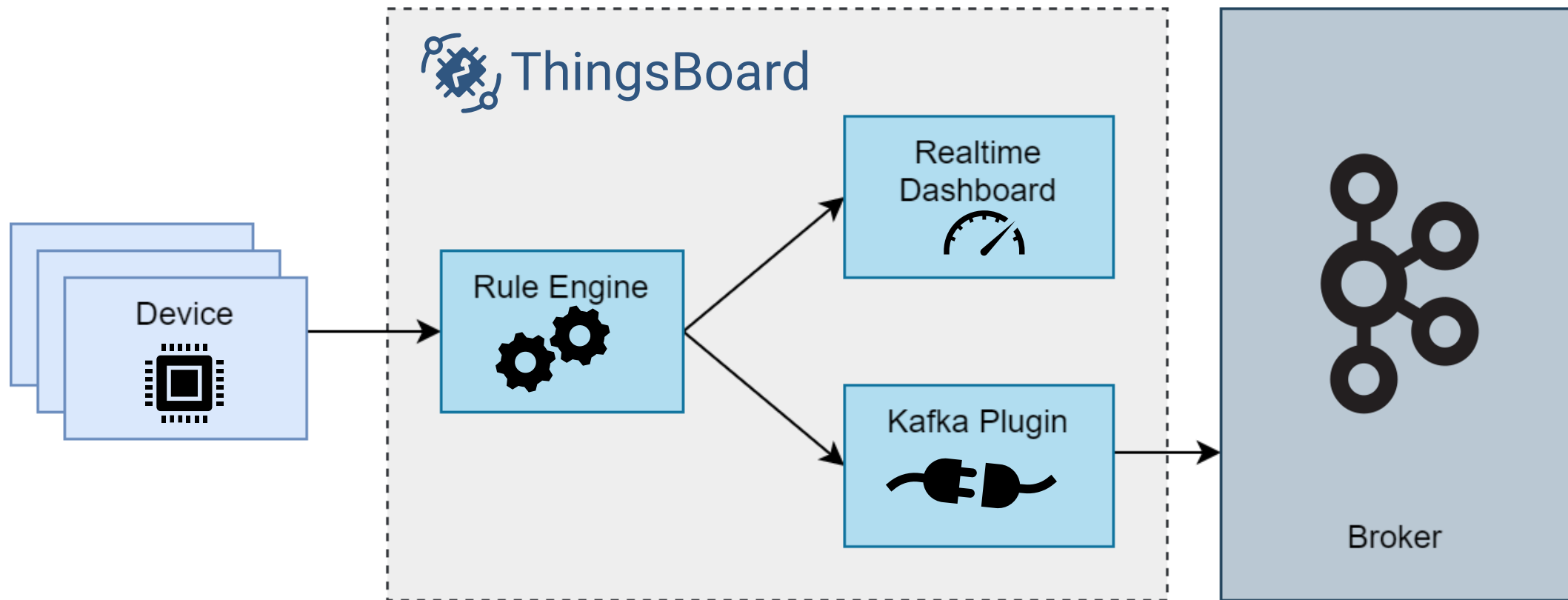
- Current:

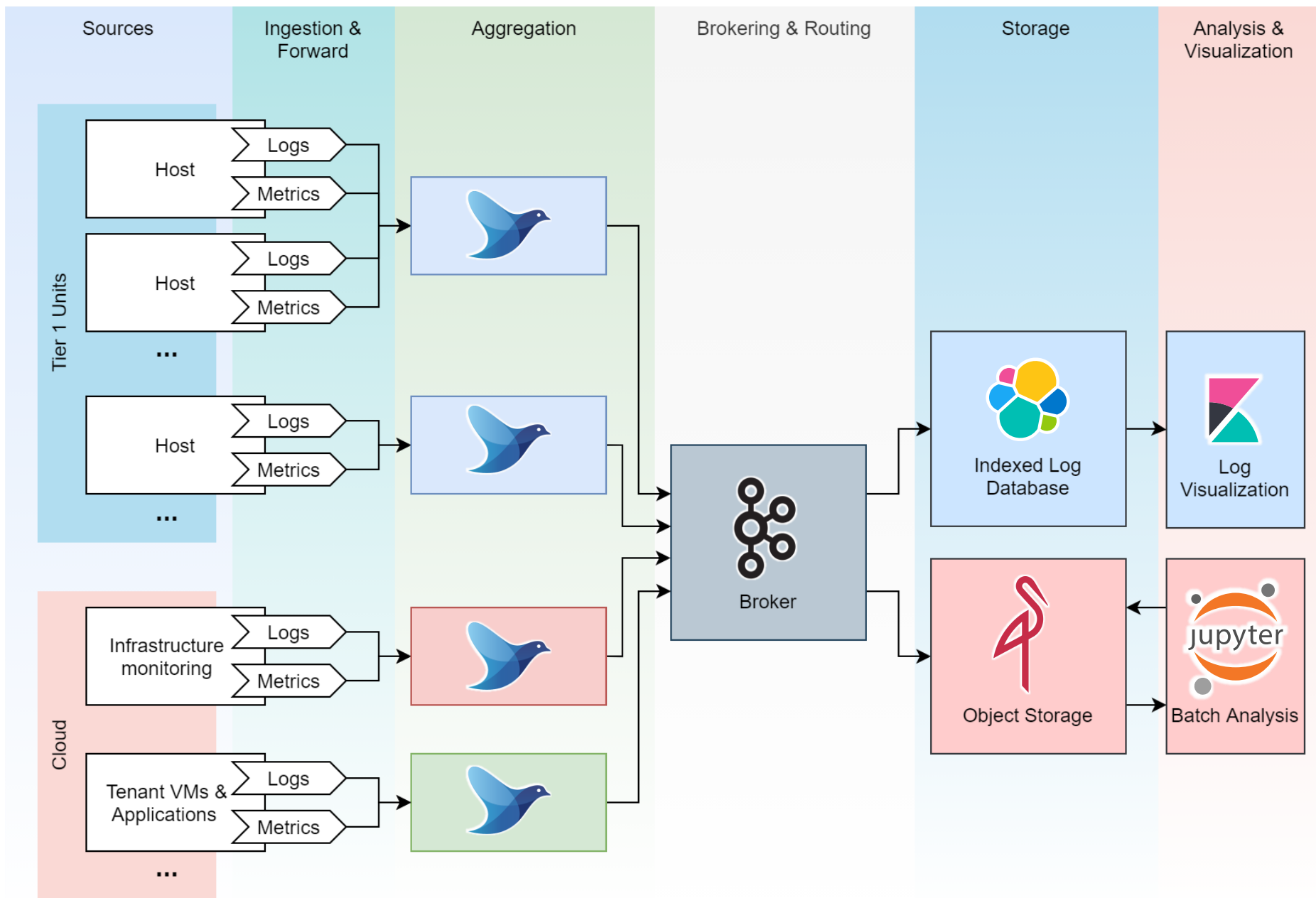


- Future:



Next Steps / Thingsboard





Log file's structure

- Log files range from heavily structured to free flowing and descriptive

0	12/01 03:48:42 : [# 1105 lifetime=18:25:00] S [OK:604700,F:74281,E:0,m:0.000,M:612.382,Avg:0.136] A [OK:61070,F:0,E:0,m:0.00
1	12/01 03:49:42 : [# 1106 lifetime=18:26:00] S [OK:605056,F:74301,E:0,m:0.000,M:612.382,Avg:0.137] A [OK:61111,F:0,E:0,m:0.00
2	12/01 03:50:42 : [# 1107 lifetime=18:27:00] S [OK:605398,F:74324,E:0,m:0.000,M:612.382,Avg:0.137] A [OK:61152,F:0,E:0,m:0.00
3	12/01 03:51:42 : [# 1108 lifetime=18:28:00] S [OK:605896,F:74346,E:0,m:0.000,M:612.382,Avg:0.137] A [OK:61194,F:0,E:0,m:0.00
...	
1410	12/02 03:17:42 : [# 2514 lifetime=41:54:00] S [OK:1341442,F:157428,E:0,m:0.000,M:612.382,Avg:0.126] A [OK:133474,F:0,E:0,m:0
0	00:00:00.140 - ERROR [xmlrpc-5916] - srmRm: File does not exist
1	00:00:00.144 - INFO [xmlrpc-5916] - srmRm: user </DC=ch/DC=cern/OU=Organic Units/OU=Users/CN=atlpilo1/CN=614260/CN=Robot:
2	00:00:00.144 - INFO [Timer-4] - ADVANCED PICKER: dispatching 15 requests.
3	00:00:00.447 - INFO [pool-4-thread-33] - srmPrepareToPut: user</DC=ch/DC=cern/OU=Organic Units/OU=Users/CN=atlpilo1/CN=61
...	
4136257	23:59:56.076 - INFO [GPFSQuotaSubmitter-1] - Submitting GPFS quota info computation for vfs rooted at /storage/gpfs_atlas

Data Format

JSON object:

- Mandatory fields for source and name of the data
 - hostname of the client that sent the event
 - name of the data produced (i.e. name of the program that generated the log, or ident)
- Optional payload based on the data type and source
 - Type of data generated (i.e. log, metric...) for routing
 - Log message or value for metrics and sensor
 - Topic for routing the message to the correct output
 - Timestamp of the event generation
- Supports nested field for complex data structure

Functional Diagram

