

**VIETTEL SOLUTIONS**

**Viettel**

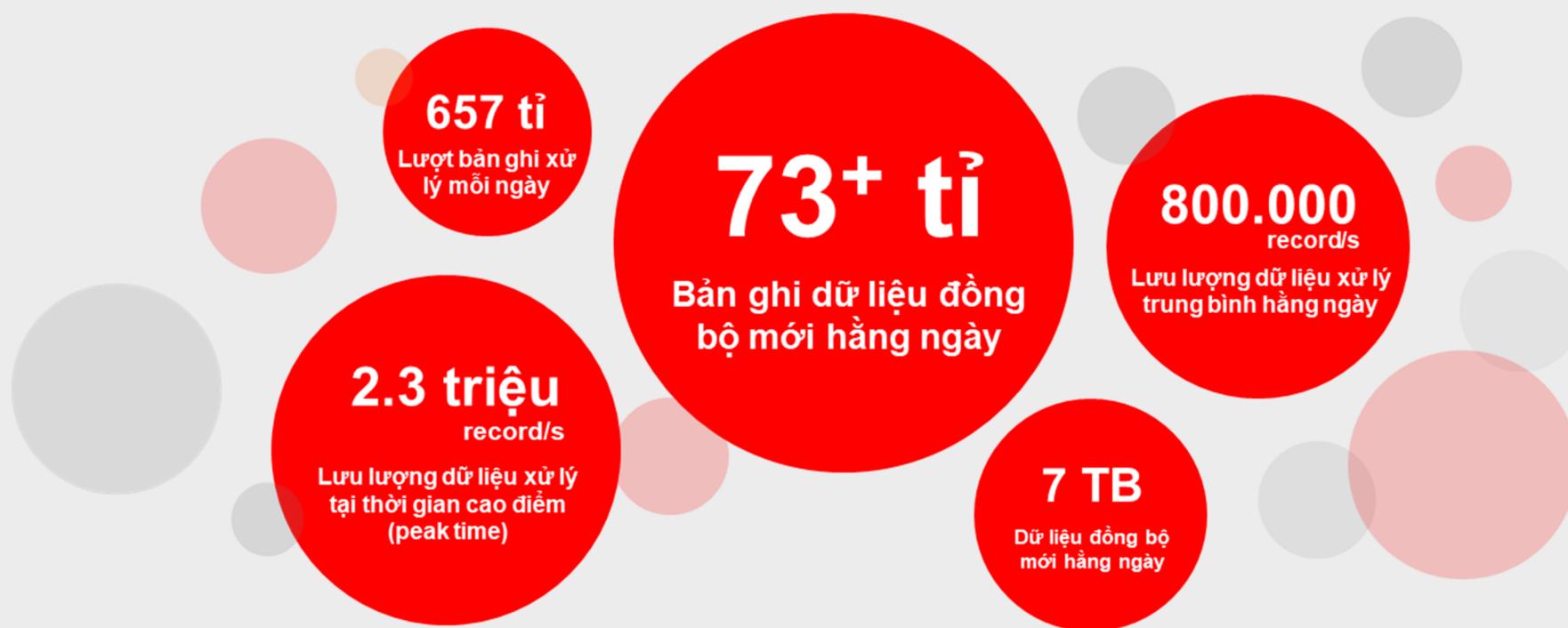
## **VIETTEL DATA PLATFORM CATALOG**

Nền tảng dữ liệu lớn Viettel

## Nền tảng của chúng tôi cung cấp

- ✓ License phần mềm theo năm hoặc trọn đời. Hỗ trợ kỹ thuật 1 năm chính hãng hoặc mua gói hỗ trợ theo năm.
- ✓ Có khả năng triển khai on-premise và môi trường đám mây (cloud)
- ✓ Hỗ trợ các hệ điều hành phổ biến Windows, Linux, Unix
- ✓ Có khả năng cài đặt, kiểm tra hệ thống, môi trường các máy chủ phân tán tự động trên một công cụ toàn trình giao diện web
- ✓ Tự động áp dụng các cấu hình nâng cao hiệu năng hệ thống tự động trong quá trình cài đặt.
- ✓ Khả năng mở rộng máy chủ không giới hạn

### Năng lực xử lý



# “TÍNH NĂNG

## LƯU TRỮ DỮ LIỆU PHÂN TÁN

- Cung cấp khả năng lưu trữ dữ liệu lớn (> 1000TB) và lưu trữ phân tán trên các máy chủ
- Các loại dữ liệu: cấu trúc, bán cấu trúc, phi cấu trúc.
- Tích hợp với các công nghệ xử lý dữ liệu lớn(Apache Spark, Apache MapReduce, Apache Hive, Apache HBase)
- Loại phần cứng lưu trữ: SSD, HDD
- Chịu lỗi và hỗ trợ tính toán song song thông qua nhân bản dữ liệu, đảm bảo tính sẵn sàng cao của dữ liệu
- Hỗ trợ các thuật toán nén dữ liệu phổ biến: Gzip, Snappy, LZO
- Hỗ trợ phân vùng dữ liệu, phân dữ liệu thành các vùng để tăng tốc độ truy cập dữ liệu
- Truy cập, hỗ trợ các giao thức truy cập dữ liệu phổ biến: HDFS, S3
- Sao lưu và khôi phục dữ liệu tự động, định kỳ theo lịch được lập
- Đảm bảo dữ liệu được bảo vệ an toàn trong quá trình lưu trữ và truyền tải.

**HDFS:** hệ thống quản lý tệp phân tán giúp lưu trữ lượng dữ liệu khổng lồ trên các máy chủ khác nhau:

- **Kiến trúc Master-Slave:**
  - **NameNode:** Quản lý metadata của hệ thống file.
  - **DataNode:** Lưu trữ dữ liệu thực tế dưới dạng block.
  - Cơ chế tự động phân phối dữ liệu giữa các **DataNode**
- **Cơ chế HA(High Availability):**
  - **Zookeeper Quorum**
  - Replicate dữ liệu trên nhiều **DataNode**
- **Mã hóa dữ liệu lưu trữ(Encryption at Rest):** HDFS Transparent Encryption
- **Mã hóa dữ liệu trên đường truyền(Encryption in Transit):** TLS/SSL
- **Kiểm soát truy cập bằng POSIX Permissions, Access Control Lists (ACLs), Apache Ranger**
- **Khả năng mở rộng theo chiều dọc; mở rộng theo chiều ngang không giới hạn.**

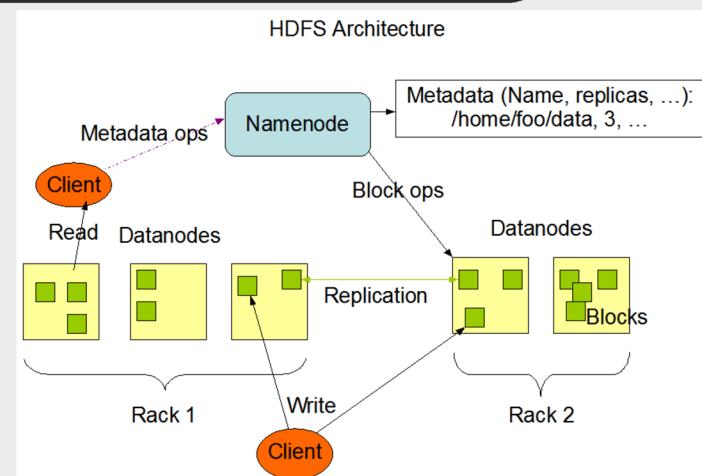
**Apache HBase :** Cơ sở dữ liệu NoSQL phân tán, column-oriented, xây dựng trên HDFS

- Lưu trữ và truy vấn ngẫu nhiên (random read/write)
- Truy vấn độ trễ thấp

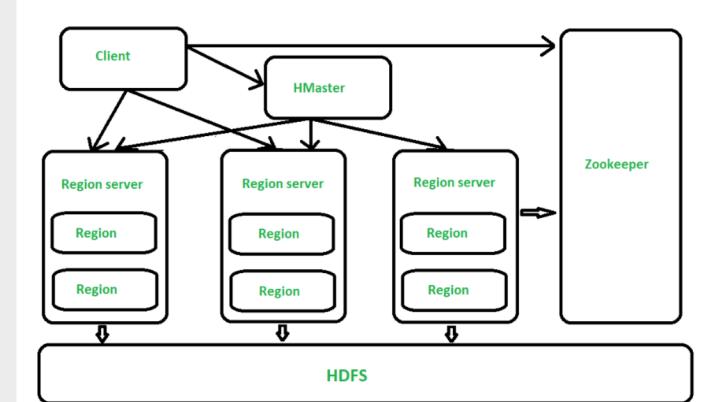
## Stack công nghệ



## Kiến trúc công nghệ HDFS



## Kiến trúc công nghệ HBase



# “TÍNH NĂNG

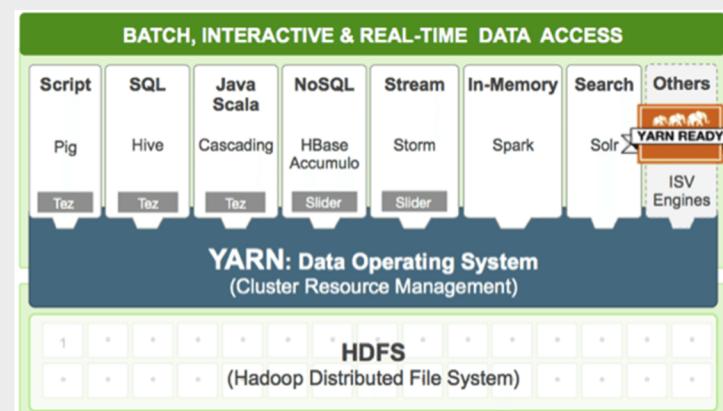
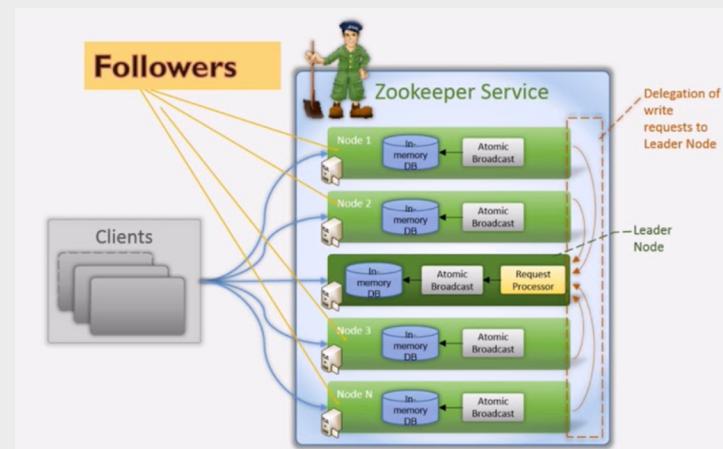
## QUẢN LÝ VÀ ĐIỀU PHỐI TÀI NGUYÊN PHÂN TÁN

Viettel Data Platform sử dụng Framework YARN làm công cụ quản lý tài nguyên của toàn bộ hệ thống. Ngoài ra Viettel Data Platform sử dụng Zookeeper để cung cấp các cơ chế đồng bộ hóa giữa các máy chủ và dịch vụ trong môi trường phân tán, đảm bảo tính nhất quán và phục hồi lỗi của các thành phần trong nền tảng.

**Zookeeper:** ZooKeeper là một hệ thống quản lý và đồng bộ hóa phân tán, được sử dụng rộng rãi trong các hệ thống phân tán để quản lý cấu hình, đồng bộ hóa trạng thái và cung cấp dịch vụ điều phối cho các ứng dụng phân tán.

**YARN:** đảm nhận vai trò quản lý các tài nguyên của hệ thống lưu trữ liệu kết hợp với việc chạy các ứng dụng phân tích.

- **Resource Manager:** Là thành phần master trong YARN, tiếp nhận các yêu cầu thực thi từ người dùng, điều phối tài nguyên của hệ thống cho các tác vụ thực thi và quản lý, giám sát các ứng dụng phân tán, có khả năng cung cấp:
  - **Pre-warmed Containers** – Giữ sẵn tài nguyên để job không mất thời gian khởi tạo.
  - **Long-running Services** – Giữ dịch vụ chạy liên tục để giảm latency.
- **Node Manager:** Là thành phần thực thi ứng dụng.
- **Timeline Server:** Dịch vụ chịu trách nhiệm thu thập và lưu trữ thông tin về trạng thái, hiệu suất, và lịch sử của các ứng dụng đang chạy trong cluster. Dịch vụ hỗ trợ người dùng và quản trị viên trong việc phân tích hiệu suất hoặc xử lý sự cố
- **History Server:** Dịch vụ lưu trữ thông tin về các ứng dụng đã hoàn thành
- **DNS Registry:** Dịch vụ đóng vai trò như một dịch vụ đăng ký và định vị các thành phần trong cụm, giúp các thành phần như ResourceManager, NodeManager, và Timeline Server có thể định vị nhau một cách linh động



# “TÍNH NĂNG

## THU THẬP VÀ CHIA SẺ DỮ LIỆU

- Giao diện quản lý truy cập, xác thực người dùng truy cập vào giao diện điều khiển và thiết kế các luồng thu thập, chia sẻ dữ liệu
- Cung cấp hệ thống dashboard tổng hợp thông tin các luồng thu thập, chia sẻ dữ liệu
- Quản lý các luồng đồng bộ dữ liệu (thêm, sửa, xoá)
- Giám sát theo dõi trạng thái của quá trình thu thập dữ liệu
- Thiết kế và quản lý các mẫu luồng đồng bộ dữ liệu để tái sử dụng
- Khả năng tiền xử lý dữ liệu, biến đổi dữ liệu, chuẩn hóa dữ liệu
- Kết nối đa dạng đến các nguồn dữ liệu bên ngoài thông qua các giao thức: API, FTP, STMP, Kafka,...
- Thu thập dữ liệu thời gian thực với độ trễ thấp
- Thu thập dữ liệu theo lô với dung lượng lớn
- Lập lịch tự động thu thập dữ liệu

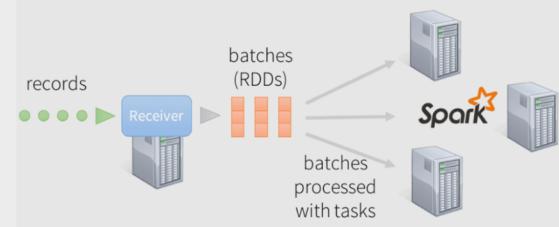
**Kafka:** Là một hệ thống truyền tin (messaging system) mạnh mẽ và có khả năng mở rộng cao, được thiết kế để xử lý và lưu trữ các luồng dữ liệu trong thời gian thực. Kiến trúc gồm 4 thành phần:

- **Broker:** Là server Kafka thực thi việc lưu trữ, nhận và phục vụ dữ liệu. Một cụm Kafka gồm nhiều broker, phối hợp với nhau để chia sẻ và nhân bản dữ liệu.
- **Zookeeper:** Quản lý metadata của cluster (thông tin broker, topic, partition, leader election). Đảm bảo tính nhất quán và chịu lỗi.
- **Producer:** gửi (publish) các bản ghi (records) vào một hoặc nhiều topic trên Kafka
- **Consumer:** Đọc (subscribe) và tiêu thụ (consume) các bản ghi từ topic

**Nifi:** Apache NiFi là một nền tảng tích hợp dữ liệu mạnh mẽ và linh hoạt, được thiết kế để tự động hóa luồng thu thập, chia sẻ dữ liệu giữa các hệ thống khác nhau.

- Web Server
- Flow Controller
- Provenance Repository
- Content Repository
- FlowFile Repository
- NiFi Registry

**Spark Streaming:** được thiết kế để xử lý luồng dữ liệu theo thời gian thực. Spark Streaming hoạt động dựa trên mô hình micro-batching, trong đó dữ liệu thời gian thực được chia thành các lô nhỏ (micro-batches) với khoảng thời gian cấu hình (batch interval).



## CÔNG NGHỆ SỬ DỤNG



# “TÍNH NĂNG

## XỬ LÝ DỮ LIỆU PHÂN TÁN

- Cung cấp môi trường làm việc cho các kỹ sư dữ liệu:
  - **Trino:** SQL query engine, Ad-hoc query
  - **Apache Hive & Spark:**
    - Xử lý, tổng hợp và phân tích dữ liệu
    - Lưu trữ, xuất dữ liệu
- Xây dựng, quản lý luồng công việc:
  - **Apache Airflow:** Quản lý job, kiểm tra và cảnh báo lỗi các luồng job
- Cung cấp môi trường làm việc cho nhà khoa học dữ liệu, chuyên viên phân tích dữ liệu:
  - **MLFlow:** Quản lý vòng đời các mô hình Machine Learning
  - **Apache Zeppelin:** web-based notebook
- Xử lý các luồng dữ liệu có tải lên tới 2 triệu message/s, hoạt động 24/7 đảm bảo tính đúng đắn, đầy đủ dữ liệu lên đến 99.9% với KPI sai số không quá 0,01%

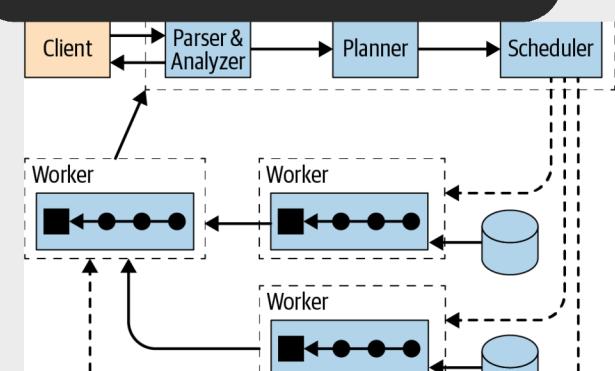
**Trino:** công cụ SQL mạnh mẽ, truy vấn dữ liệu lớn tốc độ near-realtime từ nhiều nguồn khác nhau trong phân tích, BI, và Data Lakehouse.

- Truy vấn dữ liệu phân tán từ nhiều nguồn (Có thể tận dụng các máy chủ xử lý dữ liệu để trả về kết quả truy vấn nhanh chóng):
  - Data Lake: HDFS, S3, ADLS, Data Warehouse: Hive, Iceberg, Delta Lake
  - Database: MySQL, PostgreSQL, Oracle, SQL Server
  - Streaming: Kafka
  - NoSQL: MongoDB, Cassandra
- Federated Query
- Tích hợp dễ dàng với: Superset, Tableau, Power BI
- Hỗ trợ khám phá dữ liệu, tìm kiếm dữ liệu dựa trên mô tả, dựa trên siêu dữ liệu của dữ liệu

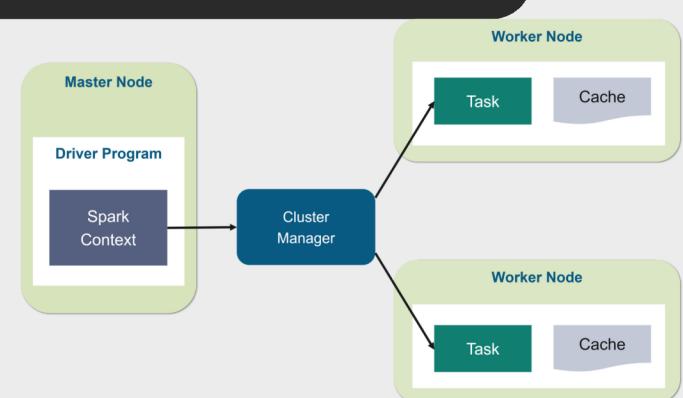
**Apache Spark:** công cụ xử lý dữ liệu phân tán In-Memory:

- Kiến trúc Master-Slave:
  - Master (Driver) chịu trách nhiệm phân tích, lập kế hoạch và phân phối công việc.
  - Workers ( Các Executors) nằm trên nhiều máy chủ, thực hiện công việc nhận task, thực thi và trả kết quả.
  - Cluster Manager (YARN,..) cấp phát tài nguyên cho cả master và workers.
- Sử dụng DAG (Directed Acyclic Graph) để ghi lại toàn bộ các bước xử lý, khôi phục lại kết quả tính toán khi có lỗi xảy ra
- Tích hợp với các thư viện học máy(Spark MLlib)
- Xử lý theo lô (SparkSQL)
- Xử lý near-realtime(Spark Streaming)
- Hỗ trợ đa ngôn ngữ: Python, Java, Scala, R, SQL

Kiến trúc công nghệ Trino



Kiến trúc công nghệ Apache Spark



# “TÍNH NĂNG

## XỬ LÝ DỮ LIỆU PHÂN TÁN

- Cung cấp môi trường làm việc cho các kỹ sư dữ liệu:
  - **Trino:** SQL query engine, Ad-hoc query
  - **Apache Hive & Spark:**
    - Xử lý, tổng hợp và phân tích dữ liệu
    - Lưu trữ, xuất dữ liệu
- Xây dựng, quản lý luồng công việc:
  - **Apache Airflow:** Quản lý job, kiểm tra và cảnh báo lỗi các luồng job
- Cung cấp môi trường làm việc cho nhà khoa học dữ liệu, chuyên viên phân tích dữ liệu:
  - **MLFlow:** Quản lý vòng đời các mô hình Machine Learning
  - **Apache Zeppelin:** web-based notebook
- Xử lý các luồng dữ liệu có tải lên tới 2 triệu message/s, hoạt động 24/7 đảm bảo tính đúng đắn, đầy đủ dữ liệu lên đến 99.9% với KPI sai số không quá 0,01%

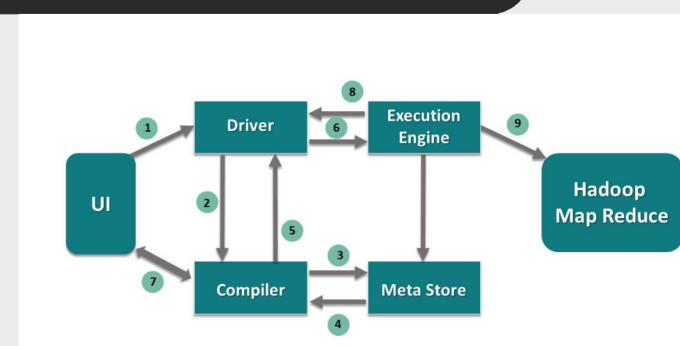
**Apache Hive:** công cụ xử lý dữ liệu lớn theo lô(batch-data):

- Sử dụng ngôn ngữ truy vấn giống SQL – HiveQL
- Phù hợp với các tập dữ liệu có quy mô terabyte đến petabyte.
- Hỗ trợ phân vùng (Partition) & phân mảnh (Bucketing)

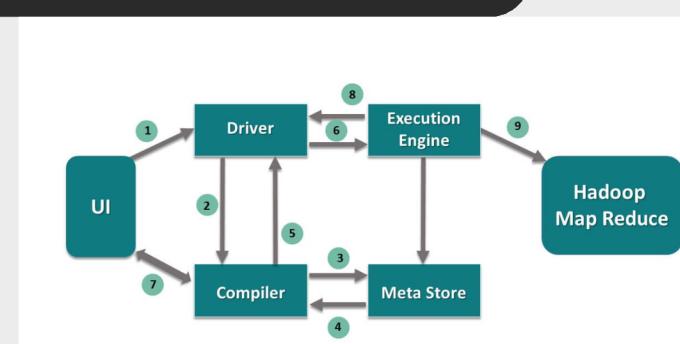
**Apache Airflow:** công cụ lập lịch, quản lý luồng công việc

- Các thành phần chính:
  - **Scheduler:** Đọc DAG, lên lịch các task, đưa các task cần chạy vào Message Broker
  - **Web Server:** Giao diện UI để quản lý, theo dõi DAG, trigger, xem log, ...
  - **Message Broker:** Hàng đợi (RabbitMQ hoặc Redis) lưu các message đại diện cho task chờ thực thi
  - **Celery Workers:** Nhóm worker chạy song song các tác vụ, lắng nghe broker, lấy task về thực thi, cảnh báo bất thường và gửi kết quả về backend
  - **Result Backend:** Lưu trạng thái và kết quả của task (Redis, database như MySQL/Postgres, hoặc S3...)
- Khả năng xây dựng luồng xử lý dữ liệu phức tạp, gồm nhiều tác vụ thực hiện tuần tự hoặc song song thông qua ngôn ngữ lập trình: Python, Java, Scala
- Tính năng retry\_on\_failure, on\_failure\_callback(gửi email, cảnh báo, ghi log,...)
- Có khả năng truyền biến thời gian khi chạy lại DAG

Kiến trúc công nghệ Apache Hive



Kiến trúc công nghệ Apache Airflow



# “TÍNH NĂNG

## XỬ LÝ DỮ LIỆU PHÂN TÁN

- Cung cấp môi trường làm việc cho các kỹ sư dữ liệu:
  - **Trino:** SQL query engine, Ad-hoc query
  - **Apache Hive & Spark:**
    - Xử lý, tổng hợp và phân tích dữ liệu
    - Lưu trữ, xuất dữ liệu
- Xây dựng, quản lý luồng công việc:
  - **Apache Airflow:** Quản lý job, kiểm tra và cảnh báo lỗi các luồng job
- Cung cấp môi trường làm việc cho nhà khoa học dữ liệu, chuyên viên phân tích dữ liệu:
  - **MLFlow:** Quản lý vòng đời các mô hình Machine Learning
  - **Apache Zeppelin:** web-based notebook
- Xử lý các luồng dữ liệu có tải lên tới 2 triệu message/s, hoạt động 24/7 đảm bảo tính đúng đắn, đầy đủ dữ liệu lên đến 99.9% với KPI sai số không quá 0,01%

**MLFlow:** công cụ quản lý toàn bộ vòng đời của các mô hình Machine

Learning

- **MLflow Tracking:** Ghi lại và theo dõi các thí nghiệm (experiments), bao gồm:

- Parameters (tham số hyperparameter)
- Metrics (độ chính xác, độ lỗi...)
- Artifacts (mô hình đã huấn luyện, đồ thị, log...)

### • MLflow Models

- Quản lý và đóng gói mô hình đã huấn luyện theo nhiều định dạng (Flavors) như:

- Python-function
- TensorFlow
- Scikit-learn
- PyTorch
- ONNX

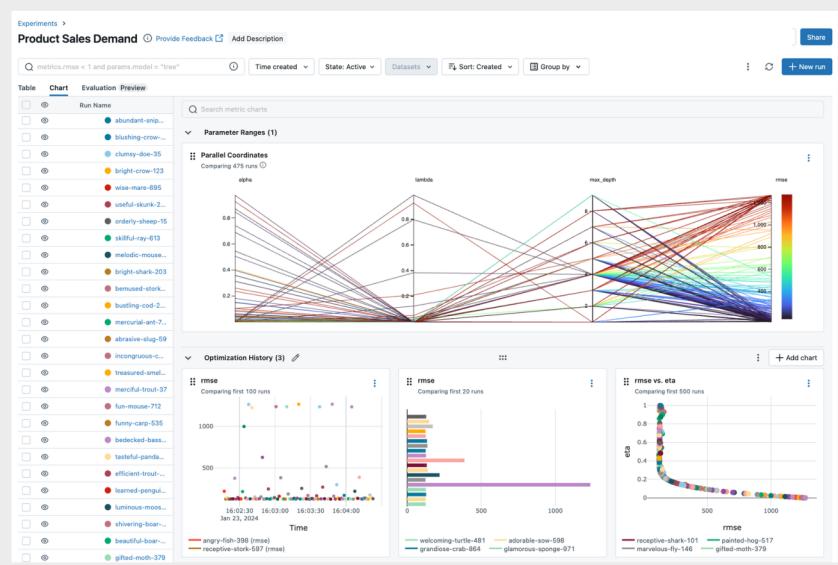
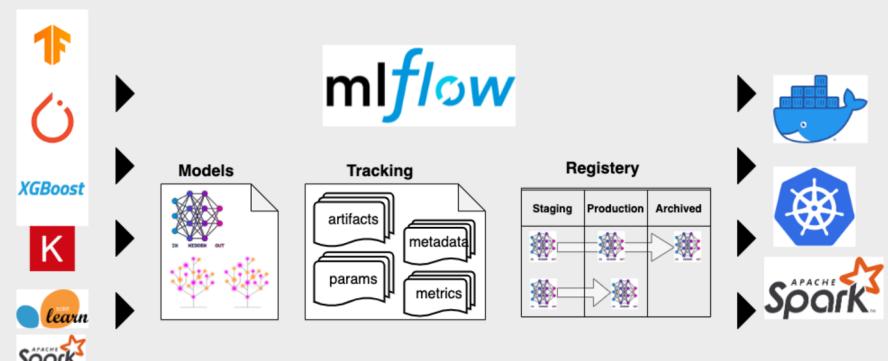
- Triển khai mô hình đã được đào tạo trên một máy chủ hoặc cluster và cung cấp API để gửi dữ liệu đầu vào và nhận kết quả từ mô hình học máy

### • MLflow Projects

- Đóng gói code, môi trường và dependencies thành project (MLproject file).
- Cho phép chạy lại experiments

### • MLflow Model Registry

- Lưu trữ trung tâm cho các mô hình đã đóng gói.
- Hỗ trợ versioning, staging, production, và archiving.



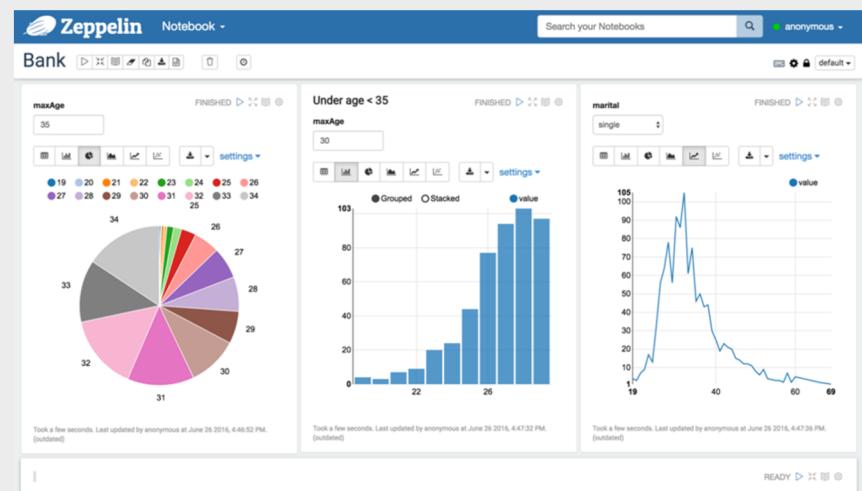
# “TÍNH NĂNG

## XỬ LÝ DỮ LIỆU PHÂN TÁN

- Cung cấp môi trường làm việc cho các kỹ sư dữ liệu:
  - **Trino:** SQL query engine, Ad-hoc query
  - **Apache Hive & Spark:**
    - Xử lý, tổng hợp và phân tích dữ liệu
    - Lưu trữ, xuất dữ liệu
- Xây dựng, quản lý luồng công việc:
  - **Apache Airflow:** Quản lý job, kiểm tra và cảnh báo lỗi các luồng job
- Cung cấp môi trường làm việc cho nhà khoa học dữ liệu, chuyên viên phân tích dữ liệu:
  - **MLFlow:** Quản lý vòng đời các mô hình Machine Learning
  - **Apache Zeppelin:** web-based notebook
- Xử lý các luồng dữ liệu có tải lên tới 2 triệu message/s, hoạt động 24/7 đảm bảo tính đúng đắn, đầy đủ dữ liệu lên đến 99.9% với KPI sai số không quá 0,01%

**Apache Zeppelin:** cung cấp các trình thông dịch đa ngôn ngữ, hỗ trợ nhiều người cùng làm việc trên cùng một báo cáo hay phân tích:

- Khả năng tích hợp với nhiều công cụ học máy thông qua trình thông dịch, cho phép:
  - Lựa chọn thuật toán học máy (MLlib, TensorFlow, Scikit-learn, H2O.ai)
  - Tiền xử lý dữ liệu (Spark, Hive, NiFi)
  - Huấn luyện mô hình (Spark MLlib, trên YARN)
  - Tinh chỉnh mô hình (Hyperparameter tuning, Grid Search)
  - Đánh giá mô hình (Cross-validation, Metrics)
- Hỗ trợ xuất báo cáo, biểu đồ
- Hỗ trợ làm việc nhóm, người dùng có thể:
  - Cùng chỉnh sửa notebook.
  - Gắn comment, markdown.
  - Trình bày kết quả dễ hiểu thông qua biểu đồ.



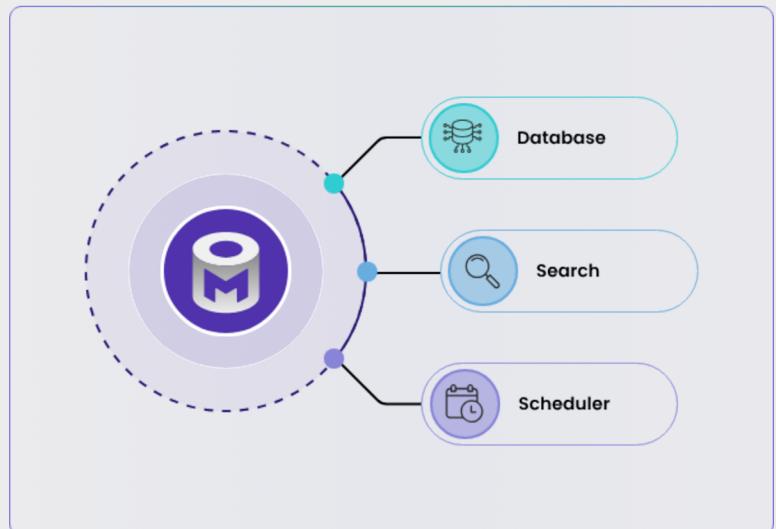
# “TÍNH NĂNG

## QUẢN TRỊ SIÊU DỮ LIỆU VÀ CHẤT LƯỢNG DỮ LIỆU

- Khả năng thu thập thông tin về các siêu dữ liệu gắn với dữ liệu, định dạng dữ liệu, tính đầy đủ của dữ liệu, tần suất cập nhật của dữ liệu
- Phân tích chất lượng dữ liệu của dữ liệu nguồn: tính toàn vẹn dữ liệu, thiếu dữ liệu, trùng lặp dữ liệu
- Nắm bắt siêu dữ liệu về quá trình xử lý dữ liệu để kiểm soát và tra cứu gốc dữ liệu
- Khả năng lưu dấu vết các hoạt động thêm, xoá, sửa dữ liệu do người dùng thực hiện

**OpenMetadata:** nền tảng quản lý metadata mạnh mẽ, bao gồm các thành phần

- Metadata Ingestion:** Kết nối & thu thập metadata tự động từ nhiều nguồn: databases (MySQL, Postgres...), data lake (S3, HDFS), data warehouse (Snowflake, BigQuery...), BI tools (Tableau, PowerBI), pipelines (Airflow, dbt)...
- Metadata Store:** Cơ sở dữ liệu trung tâm lưu trữ metadata (schemas, bảng, cột, dashboards, pipelines...)
- Data Catalog:** Giao diện web cho phép người dùng tìm kiếm, duyệt, và xem chi tiết metadata
- Lineage:** Hiển thị đồ thị lineage (dòng chảy dữ liệu) giữa các đối tượng (tables → jobs → dashboards)
- Data Quality & Profiling:** Tích hợp kiểm tra chất lượng dữ liệu, profiling (số lượng null, phân phối giá trị...)
- Glossary & Business Metadata:** Xây dựng business glossary, gắn nhãn (tag), định nghĩa business terms cho metadata
- Access Control & Security:** Phân quyền chi tiết trên metadata, tích hợp với LDAP/SSO/Kerberos
- APIs & SDKs:** Cung cấp REST API và client SDK (Python, Java) để tự động hóa và tích hợp sâu với hệ thống
- Notifications & Lineage Alerts:** Cấu hình cảnh báo khi metadata thay đổi hoặc vi phạm quality rules



# TÍNH NĂNG

## BẢO MẬT

- Các dịch vụ được sử dụng trong Viettel Data Platform bao gồm : Ranger, Kerberos, Knox

**Ranger:** là một dịch vụ quản lý an ninh và kiểm soát truy cập trong các hệ thống dữ liệu lớn. Ranger cung cấp khả năng quản lý quyền chi tiết và chính sách bảo mật tập trung cho nhiều dịch vụ dữ liệu như Hadoop, Hive, Kafka, HBase, và nhiều dịch vụ khác. Các tính năng của Ranger:

### Các tính năng của Ranger:

- Centralized Authorization:** Quản lý các chính sách truy cập (access control policies) từ một giao diện tập trung, thay vì cấu hình phân tán ở từng service.
- Fine-Grained Access Control:** Kiểm soát truy cập chi tiết theo cấp độ: database, table, column, row, topic (Kafka), file path (HDFS)...
- Role-Based Access Control (RBAC):** Phân quyền theo vai trò người dùng (Role), giúp dễ dàng gán nhiều quyền cho nhóm người dùng cùng loại.
- Xác thực người dùng qua LDAP hoặc Active Directory
- Attribute-Based Access Control (ABAC):** Cho phép tạo rule truy cập dựa trên thuộc tính của người dùng (user attributes).
- Policy Conditions:** Chính sách có thể điều kiện theo IP, thời gian, nhóm người dùng, v.v.
- Data masking:** Ẩn dữ liệu nhạy cảm (data masking) và lọc dòng dữ liệu dựa theo người dùng.
- Audit Logging:** Ghi log chi tiết tất cả hành động truy cập.
- Giao diện quản lý trực quan

Policy ID	Policy Version	Event Time	Application	User	Service Name / Type	Resource Name / Type	Access Type	Permission	Result	Access Enforcer	Agent Host Name	Client IP	Cluster Name
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	memory catalog	select	allow			ranger-acl	my-localhost-trino		
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	pxn catalog	select	allow			ranger-acl	my-localhost-trino		
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	tpch catalog	select	allow			ranger-acl	my-localhost-trino		
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	system catalog	select	allow			ranger-acl	my-localhost-trino		
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	tpcds catalog	select	allow			ranger-acl	my-localhost-trino		
--	--	09/26/2021 11:42:29 PM	trino	ranger-admin	tpch catalog	select	denied			ranger-acl	my-localhost-trino		
--	--	09/26/2021 11:42:29 PM	trino	ranger-admin	system catalog	select	denied			ranger-acl	my-localhost-trino		
--	--	09/26/2021 11:42:29 PM	trino	ranger-admin	memory catalog	select	denied			ranger-acl	my-localhost-trino		
--	--	09/26/2021 11:42:29 PM	trino	ranger-admin	jmx catalog	select	denied			ranger-acl	my-localhost-trino		
--	--	09/26/2021 11:42:29 PM	trino	ranger-admin	tpcds catalog	select	denied			ranger-acl	my-localhost-trino		

User Name	Email Address	Role	User Source	Sync Source	Groups	Visibility	Sync Details
admin		Admin	Internal	--	--	Visible	--
rangerenablesync		Admin	Internal	--	--	Visible	--
rangerenabg		Admin	Internal	--	--	Visible	--
hdfs		User	External	hdfs	hdfs, hdfs	Visible	--
yarn-all		User	External	hdfs	hdfs	Visible	--
hive		User	External	hdfs	hdfs	Visible	--
infra-all		User	External	hdfs	hdfs	Visible	--
zookeeper		User	External	hdfs	hdfs	Visible	--
oozie		User	External	hdfs	hdfs, users	Visible	--
ranger		User	External	hdfs	ranger	Visible	--
fec		User	External	hdfs	users	Visible	--
adops		User	External	hdfs	adops	Visible	--
zeppelin		User	External	hdfs	zeppelin	Visible	--

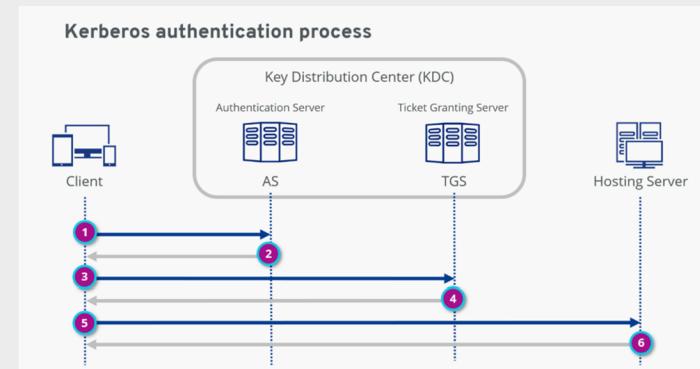


# “TÍNH NĂNG

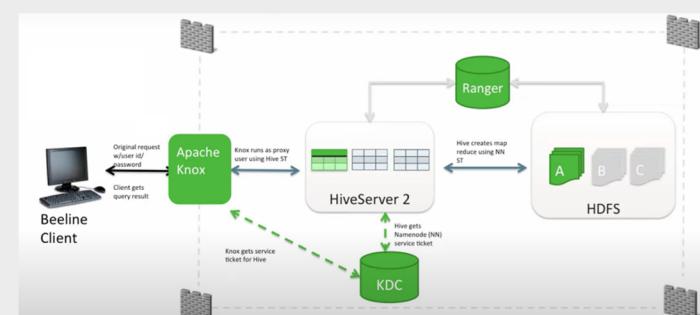
## BẢO MẬT

- Các dịch vụ được sử dụng trong Viettel Data Platform bao gồm : **Ranger, Kerberos, Knox**

**Kerberos:** Kerberos là một giao thức xác thực qua mạng được phát triển tại Viện Công nghệ Massachusetts (MIT) và hiện là một phần quan trọng trong các hệ thống bảo mật mạng. Kerberos sử dụng cơ chế mã hóa đối xứng và một hệ thống dựa trên ticket để đảm bảo rằng danh tính của người dùng hoặc dịch vụ được xác thực trước khi truy cập tài nguyên. Điểm nổi bật của Kerberos là khả năng bảo vệ chống lại các cuộc tấn công nghe lén và giả mạo, nhờ việc hạn chế truyền mật khẩu qua mạng.



**Knox:** Apache Knox là một cổng dịch vụ (gateway) được thiết kế để bảo mật và đơn giản hóa truy cập vào các dịch vụ dữ liệu lớn trong Viettel Data Platform. Knox cung cấp một điểm truy cập duy nhất, giúp bảo vệ các API của hệ sinh thái Viettel Data Platform thông qua các cơ chế xác thực, ủy quyền, và ghi nhật ký tập trung. Việc tích hợp với Kerberos giúp Knox đảm bảo rằng chỉ các yêu cầu hợp lệ mới được phép truy cập vào tài nguyên.



# TÍNH NĂNG

## QUẢN TRỊ HỆ THỐNG

- Cung cấp khả năng tương tác start, stop với các thành phần ứng dụng.

- Quản trị log ứng dụng trong quá trình vận hành hệ thống, trực quan hóa dữ liệu giám sát qua Ambari Log Search

### Ambari Log Search

Log Search

Service Logs

Fatal: 0, Error: 2897, Warn: 8485, Info: 0, Debug: 6, Trace: 0

Histogram: 12/07/2015 2:04:05,166 - 12/08/2015 18:41:53,944

Event History (1 of 1)

Search: Q Search, Q Search String, Q Include Search, Q Exclude Search, Include Components

Log Data: Components, Hosts, Service Logs

© Hortonworks Inc. 2011 – 2016. All Rights Reserved

### Log Search

Log Time (PDT)

FILE

BSL-0-7035014d1

2017-10-09 15:45:00,753 WARN - [BSL-0-7035014d1]: FAILED: Execution Error, return code 1 from org.apache.hadoop.hive.exec.DDLTask. MetaException(message:java.security.AccessControlException: Permission denied: user:bigsq1, access=WRITE, inode="/tmp/private":hdfs:hdfs:drwx-----) at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:319) at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:219) at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:198) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1955) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1939) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1913) at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.checkAccess(FSNamesystem.java:8750) at org.apache.hadoop.hdfs.server.namenode.NameNodeRpcServer.checkAccess(NameNodeRpcServer.java:2089) at org.apache.hadoop.hdfs.protocolPB.ClientNamenodeProtocolServerSideTranslatorPB.checkAccess(ClientNamenodeProtocolServerSideTranslatorPB.java:1466) at org.apache.hadoop.ipc.ProtobufRpcEngine\$Server\$Protocol\$ClientNamenodeProtocol\$2.callBlockingMethod(ClientNamenodeProtocol\$Protos.java) at org.apache.hadoop.ipc.ProtobufRpcEngine\$Server\$Protocol\$ClientNamenodeProtocol\$2.callBlockingMethod(ClientNamenodeProtocol\$Protos.java) at org.apache.hadoop.ipc.RPC\$Server.call(RPC.java:982) at org.apache.hadoop.ipc.Server\$Handler\$1.run(Server.java:2351) at org.apache.hadoop.ipc.Server\$Handler\$1.run(Server.java:2347) at java.security.AccessController.doPrivileged(Native Method) at javax.security.auth.Subject.doAs(Subject.java:422) at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1866) at org.apache.hadoop.ipc.Server\$Handler\$1.run(Server.java:2345)

2017-10-09 15:45:00,753 WARN - [BSL-0-7035014d1]: Error output : FAILED: Execution Error, return code 1 from org.apache.hadoop.hive.exec.DDLTask. MetaException(message:java.security.AccessControlException: Permission denied: user:bigsq1, access=WRITE, inode="/tmp/private":hdfs:hdfs:drwx-----) at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:319) at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:219) at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:198) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1955) at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkPermission(FSDirectory.java:1939)

# “TÍNH NĂNG

## QUẢN TRỊ HỆ THỐNG

- Khả năng tùy biến cảnh báo hệ thống

The screenshot shows the 'Alerts' section of the Viettel Data Platform. On the left is a sidebar with various service icons: Ranger, NiFi, Spark3, Zeppelin Note..., Airflow, Elastic search, Hue, Kerberos, Open metadata, Redis, Thrift Server 1, Trino, and Hosts. Below these is an 'Alerts' icon. The main area has a header 'Alerts' with sub-links for 'Status', 'Alert Definition Name', 'Service', 'Last Status Changed', and 'State'. It lists ten alerts:

Status	Alert Definition Name	Service	Last Status Changed	State
CRIT	Ambari Server Alerts	Ambari	17 hours ago	Enabled
OK (9)	Ambari Agent Heartbeat	Ambari	15 days ago	Enabled
OK (9)	Spark3 Livy2 Server	Spark3	5 days ago	Enabled
OK (9)	Infra Solr Web UI	Infra Solr	about a day ago	Enabled
OK (9)	NodeManager Web UI	YARN	15 hours ago	Enabled
OK (9)	ZooKeeper Server Process	ZooKeeper	5 days ago	Enabled
OK (9)	DataNode Heap Usage	HDFS	15 hours ago	Enabled
OK (9)	DataNode Process	HDFS	5 days ago	Enabled
OK (9)	DataNode Web UI	HDFS	15 hours ago	Enabled

- Hỗ trợ các kênh thông báo thông tin cảnh báo: SMS, Email

The dialog box is titled 'Create Alert Notification'. It contains fields for 'Method' (set to 'EMAIL'), 'Email To' (containing 'huynd275@viettel.com.vn'), 'SMTP Server' (containing '10.0.311.22'), 'SMTP Port' (containing '9000'), and 'Email From' (containing 'ambariserver@gmail.com.vn'). Below these are fields for 'Use authentication' (checkbox), 'Username' (text input), 'Password' (text input), and 'Password Confirmation' (text input).

- Khả năng tích hợp với các phần mềm trong nền tảng kho dữ liệu
- Khả năng thu thập và lưu trữ thông tin hoạt động của hệ thống

# “TÍNH NĂNG

## QUẢN TRỊ HỆ THỐNG

- Cung cấp khả năng phân quyền cho người vận hành, nhóm người vận hành theo vai trò, mức độ xử lý công việc.

Permissions	Cluster User	Service Operator	Service Administrator	Cluster Operator	Cluster Administrator	Ambari Administrator
View metrics	✓	✓	✓	✓	✓	✓
View status information	✓	✓	✓	✓	✓	✓
View configurations	✓	✓	✓	✓	✓	✓
Compare configurations	✓	✓	✓	✓	✓	✓
View service alerts	✓	✓	✓	✓	✓	✓
Start, stop, or restart service		✓	✓	✓	✓	✓
Decommission or recommission		✓	✓	✓	✓	✓
Run service checks		✓	✓	✓	✓	✓
Turn maintenance mode on or off		✓	✓	✓	✓	✓
Perform service-specific tasks		✓	✓	✓	✓	✓
Modify configurations			✓	✓	✓	✓
Manage configuration groups			✓	✓	✓	✓
Move to another host			✓	✓	✓	✓
Enable HA			✓	✓	✓	✓
Enable or disable service alerts			✓	✓	✓	✓
Add service to cluster				✓	✓	

Danh sách quyền tương tác với service theo vai trò

Permissions	Cluster User	Service Operator	Service Administrator	Cluster Operator	Cluster Administrator	Ambari Administrator
View metrics	✓	✓	✓	✓	✓	✓
View status information	✓	✓	✓	✓	✓	✓
View configuration	✓	✓	✓	✓	✓	✓
Turn maintenance mode on or off				✓	✓	✓
Install components				✓	✓	✓
Add or delete hosts				✓	✓	✓

Danh sách quyền tương tác với các host trong cụm theo vai trò

Permissions	Cluster User	Service Operator	Service Administrator	Cluster Operator	Cluster Administrator	Ambari Administrator
View metrics	✓	✓	✓	✓	✓	✓
View status information	✓	✓	✓	✓	✓	✓
View configuration	✓	✓	✓	✓	✓	✓
View stack version details	✓	✓	✓	✓	✓	✓
View alerts	✓	✓	✓	✓	✓	✓
Enable or disable alerts					✓	✓
Enable or disable Kerberos					✓	✓
Upgrade or downgrade stack					✓	✓

Danh sách quyền tương tác với service theo vai trò

Permissions	Cluster User	Service Operator	Service Administrator	Cluster Operator	Cluster Administrator	Ambari Administrator
Create new clusters						✓
Set service users and groups						✓
Rename clusters						✓
Manage users						✓
Manage groups						✓
Manage Ambari Views						✓
Assign permission and roles						✓
Manage stack versions						✓
Edit stack repository URLs						✓

Danh sách quyền tương tác với giao diện quản trị theo vai trò

# “TÍNH NĂNG

## QUẢN TRỊ HỆ THỐNG

- Cung cấp khả năng thay đổi cấu hình ứng dụng, áp dụng cho toàn bộ ứng dụng liên quan trên tập máy chủ lớn bằng cách thay đổi cấu hình trên giao diện quản lý Ambari và cập nhật cho toàn bộ ứng dụng bằng tính năng Restart All Affected Components

Viettel Data Platform

Services / NiFi / Summary

SUMMARY CONFIGS METRICS

Restart Required: 1 Component on 1 Host

Components: Started (NiFi)

0/0 Live NIFI CERTIFICATE AUTHORITY

RESTART

Restart All Affected

Summary

Components: Started (NiFi)

0/0 Live NIFI CERTIFICATE AUTHORITY

Activate Windows Go to Settings to activate Windows.

- VDP cung cấp khả năng quản lý, giám sát các tài nguyên hệ thống

Ambari

Dashboard

Services: HDFS, YARN, MapReduce2, Tez, Hive, HBase, Oozie, ZooKeeper, Ambari Metrics, Kerberos, Spark3

Hosts

Alerts

Cluster Admin

Stack and Versions

Service Accounts

Kerberos

Service Auto Start

METRICS HEATMAPS CONFIG HISTORY

LAST 1 HOUR

Metrics

LAST 1 HOUR

1

2

NameNode Heap: 2%

HDFS Disk Usage: 13%

NameNode CPU WIO: 0.0%

DataNodes Live: 3/3

NameNode RPC: 0.08 ms

Memory Usage: 37.2 GB, 18.6 GB

Network Usage: 48.8 KB

CPU Usage: 100%, 50%

Cluster Load: 5

NameNode Uptime: 8d 4h 27m

ResourceManager Heap: 9%

NodeManagers Live: 3/3

YARN Containers: n/a

HBase Master Heap: 15%

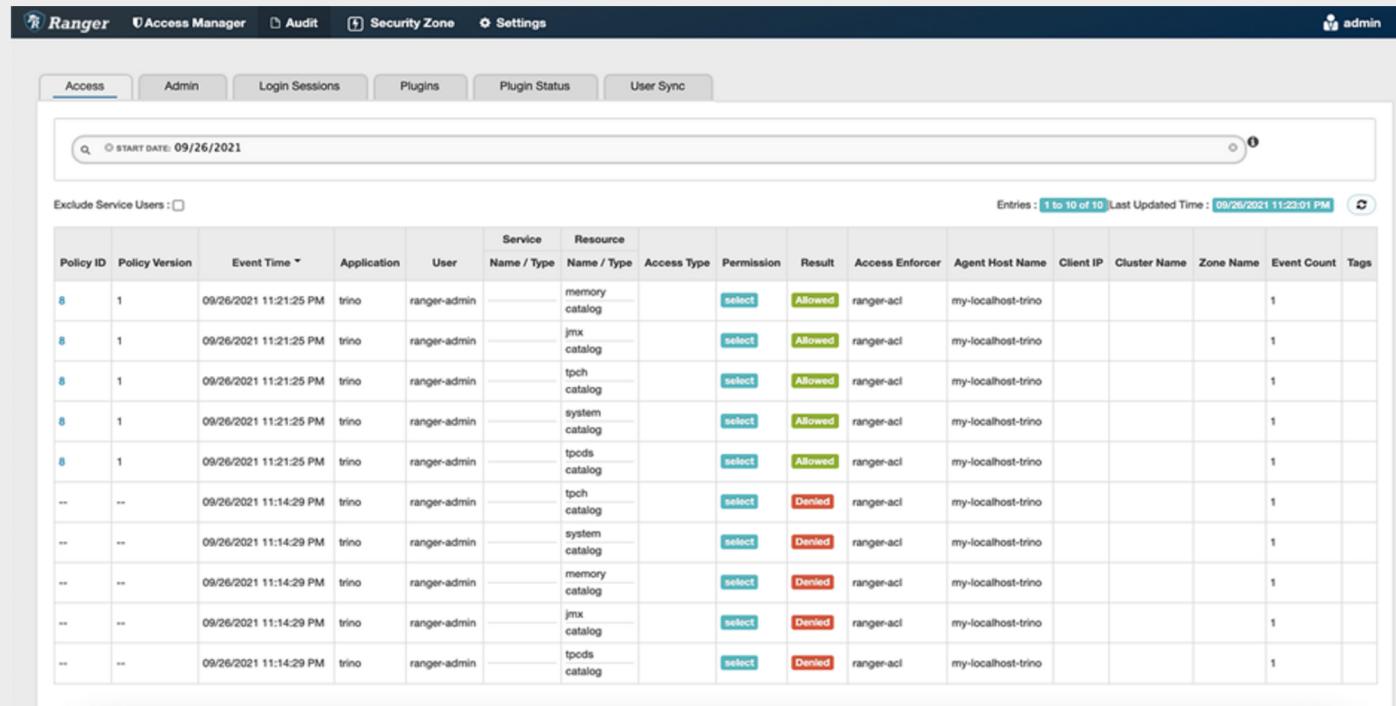
HBase Ave Load: 3

Region In Transition: 0

# “TÍNH NĂNG

## QUẢN TRỊ HỆ THỐNG

- Cung cấp khả năng chia sẻ, phân vùng tài nguyên cho các nhóm người dùng khác nhau qua ứng dụng YARN trên giao diện quản lý Ambari. YARN hỗ trợ nhiều phương thức phân vùng tài nguyên, bao gồm:
  - Queue (Hàng đợi):** Phân chia tài nguyên theo tổ chức, nhóm người dùng hoặc loại workload.
  - Node Labels:** Chỉ định tài nguyên cụ thể (CPU, RAM, GPU) cho một nhóm người dùng nhất định.
  - Resource Pools:** Hạn chế tài nguyên mà mỗi nhóm có thể sử dụng.



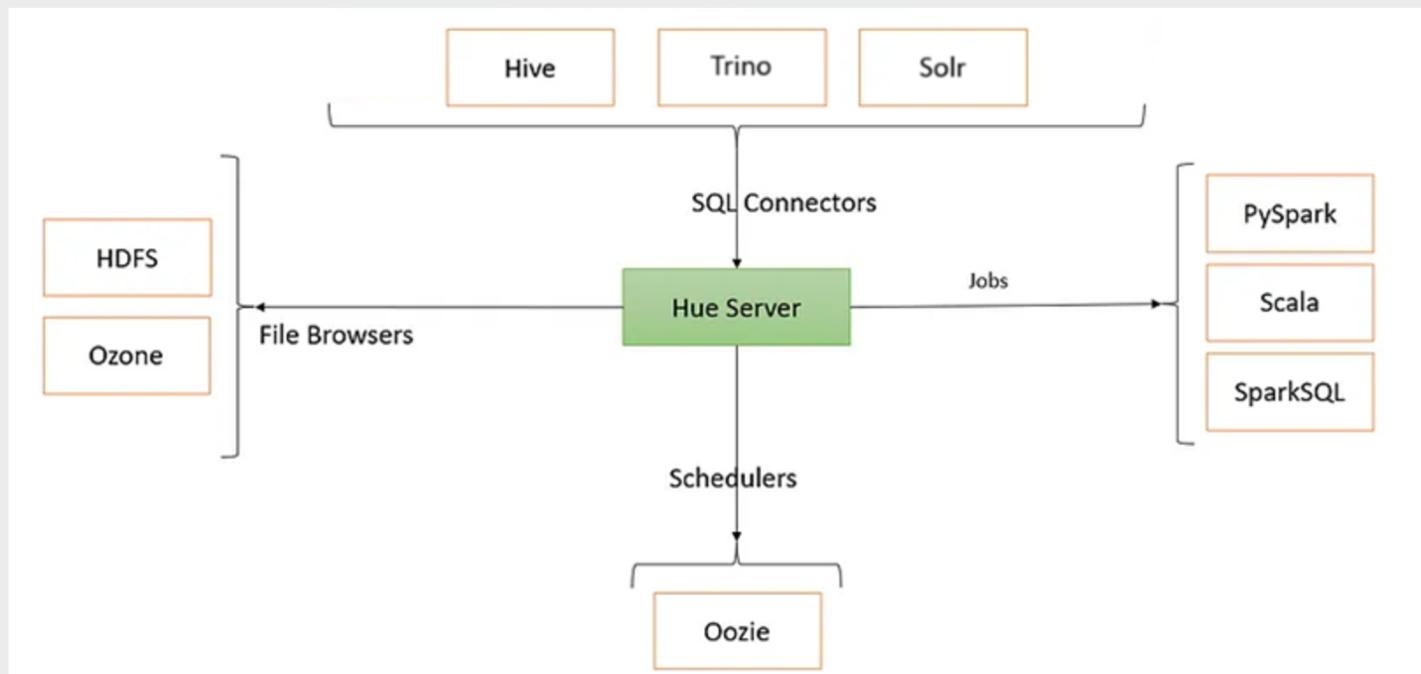
The screenshot shows the Ranger UI interface. At the top, there are tabs: Ranger, Access Manager, Audit, Security Zone, and Settings. The Settings tab is selected, showing the user 'admin'. Below the tabs is a navigation bar with buttons for Access, Admin, Login Sessions, Plugins, Plugin Status, and User Sync. The main area displays a table of access logs. The table has columns: Policy ID, Policy Version, Event Time, Application, User, Service Name / Type, Resource Name / Type, Access Type, Permission, Result, Access Enforcer, Agent Host Name, Client IP, Cluster Name, Zone Name, Event Count, and Tags. The logs show various events for the 'trino' application by 'ranger-admin' user, primarily involving 'memory catalog', 'jmx catalog', 'tpch catalog', and 'system catalog'. Most permissions are 'Allowed' (green), except for some 'Denied' (red) entries. The table includes a search bar at the top and pagination information at the bottom: 'Entries : 1 to 10 of 10 Last Updated Time : 09/26/2021 11:23:01 PM'.

Policy ID	Policy Version	Event Time	Application	User	Service Name / Type	Resource Name / Type	Access Type	Permission	Result	Access Enforcer	Agent Host Name	Client IP	Cluster Name	Zone Name	Event Count	Tags
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	memory catalog		select	Allowed	ranger-acl	my-localhost-trino					1	
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	jmx catalog		select	Allowed	ranger-acl	my-localhost-trino					1	
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	tpch catalog		select	Allowed	ranger-acl	my-localhost-trino					1	
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	system catalog		select	Allowed	ranger-acl	my-localhost-trino					1	
8	1	09/26/2021 11:21:25 PM	trino	ranger-admin	tpcds catalog		select	Allowed	ranger-acl	my-localhost-trino					1	
...	...	09/26/2021 11:14:29 PM	trino	ranger-admin	tpch catalog		select	Denied	ranger-acl	my-localhost-trino					1	
...	...	09/26/2021 11:14:29 PM	trino	ranger-admin	system catalog		select	Denied	ranger-acl	my-localhost-trino					1	
...	...	09/26/2021 11:14:29 PM	trino	ranger-admin	memory catalog		select	Denied	ranger-acl	my-localhost-trino					1	
...	...	09/26/2021 11:14:29 PM	trino	ranger-admin	jmx catalog		select	Denied	ranger-acl	my-localhost-trino					1	
...	...	09/26/2021 11:14:29 PM	trino	ranger-admin	tpcds catalog		select	Denied	ranger-acl	my-localhost-trino					1	

# “TÍNH NĂNG

## CÔNG CỤ TƯƠNG TÁC NGƯỜI DÙNG

**Hue (Hadoop User Experience):** Hue cung cấp một nền tảng trực quan giúp người dùng chạy truy vấn SQL, quản lý dữ liệu trên HDFS, và theo dõi các tác vụ Hadoop mà không cần sử dụng dòng lệnh phức tạp. Hue hỗ trợ tích hợp với nhiều dịch vụ như Hive, Solr, giúp người dùng thực hiện các công việc như phân tích dữ liệu, hoặc quản lý luồng công việc một cách hiệu quả.



The screenshot shows the Hue Query Editor interface. The top navigation bar includes links for Home, Query Editors, Data Browsers, Workflows, and Search. Below the navigation is a sub-navigation bar with Hive Editor and Query Editor selected. The main area has tabs for Navigator, Settings, and DATABASE. The DATABASE tab shows a list of tables under the default database, such as page\_view, tweets, business, and others. To the right of the Navigator is a code editor titled "Sample: Salary growth" with the following SQL query:

```

1 SELECT s07.description, s07.salary, s08.salary,
2 s08.salary - s07.salary
3 FROM
4 sample_07 s07 JOIN sample_08 s08
5 ON s07.code = s08.code
6 WHERE
7 s07.salary < s08.salary
8 ORDER BY s08.salary-s07.salary DESC
9 LIMIT 20
  
```

Below the code editor are buttons for Execute, Save, Save as..., Explain, or create a New query. The bottom section displays the results of the query as a bar chart. The chart has "salary" on the Y-axis and various job descriptions on the X-axis. The chart type is set to "Bar". The data shows salary differences for various professions, with the highest difference being around 190,000 for "Dentists, all Surgeons" and the lowest being around 50,000 for "Rotary drill/Pediatrician".

# “TÍNH NĂNG CHUNG

## SAO LƯU VÀ KHÔI PHỤC

- Hỗ trợ đa dạng các tần suất sao lưu dữ liệu (hàng ngày, hàng tuần, hàng tháng) và tự động sao lưu theo lịch được lập qua cronjob:

### HDFS Snapshot:

- Tạo ảnh chụp nhanh (snapshot) của thư mục trong HDFS mà không ảnh hưởng đến hiệu suất
- Dễ dàng khôi phục dữ liệu về trạng thái trước đó khi gặp sự cố

### HDFS DistCp (Distributed Copy):

- Dùng để sao chép dữ liệu giữa các cluster HDP hoặc giữa các thư mục trong HDFS
- Hỗ trợ sao chép dữ liệu song song, tăng tốc độ truyền tải

### Hive Export/Import:

- Dùng để sao lưu dữ liệu bảng Hive và khôi phục lại khi cần
- Xuất dữ liệu sang HDFS hoặc hệ thống lưu trữ khác

- VDP hỗ trợ các loại sao lưu dữ liệu (tổng thể, tăng tiến):

### Sao lưu toàn bộ (Full Backup):

- Dùng HDFS DistCp để sao chép dữ liệu toàn bộ: hadoop distcp hdfs://namenode/data hdfs://backup-cluster/data
- Dùng HDFS Snapshot để lưu trữ trạng thái của dữ liệu: hdfs dfs -createSnapshot /data backup\_2024\_03\_23

### Sao lưu tăng tiến (Incremental Backup):

- HDFS Snapshot Diff – Xác định & sao lưu chỉ các thay đổi so với snapshot trước đó: hdfs snapshotDiff /data backup\_2024\_03\_22 backup\_2024\_03\_23
- Hive ACID Incremental Backup – Chỉ sao lưu dữ liệu cập nhật trong bảng Hive: EXPORT TABLE sales\_data TO 'hdfs://backup/hive/sales\_incremental\_2024\_03\_23'

# “TÍNH NĂNG CHUNG

## SAO LƯU VÀ KHÔI PHỤC

- VDP có chiến lược phòng chống thảm họa bằng các phương pháp:
  - Sử dụng tính năng sao lưu và khôi phục để tạo bản sao dữ liệu.
  - Sử dụng cấu hình cao khả dụng để đảm bảo dịch vụ luôn sẵn có.
  - Sử dụng các công cụ giám sát để phát hiện sự cố.
  - Sử dụng các công cụ bảo mật để bảo vệ dữ liệu khỏi các truy cập trái phép.
- VDP có khả năng kiểm tra tính năng phòng chống thảm họa gồm việc sao lưu dữ liệu và cấu hình, đồng bộ hóa dữ liệu giữa các trung tâm dữ liệu và thiết lập hệ thống dự phòng.
- VDP có khả năng kiểm tra tính năng sao lưu và khôi phục dữ liệu bằng cách sử dụng kịch bản kiểm thử tính năng hoặc sử dụng tính năng so sánh giữa 2 phiên bản dữ liệu để so sánh nội dung snapshot với dữ liệu gốc: hdfs dfs -diff /data/.snapshot/snapshot\_2024\_03\_23 /data

## KHẢ NĂNG TRUY CẬP ĐỒNG THỜI

- Khả năng truy cập đồng thời của đội vận hành hệ thống với các cơ chế bảo vệ của Ambari:
  - **Web UI & API:** Cho phép nhiều người dùng giám sát và quản lý cụm cùng lúc.
  - **Phân quyền RBAC:** Giúp nhiều người có thể làm việc đồng thời mà không xung đột quyền.
  - **Cluster State Synchronization:** Đảm bảo mọi thay đổi cấu hình được đồng bộ theo thời gian thực.
  - **Giám sát lịch sử thao tác:** Giúp kiểm soát ai đã thực hiện thay đổi nào để đảm bảo bảo mật

**VIETTEL SOLUTIONS**



Nền tảng dữ liệu lớn Viettel

- Địa chỉ: Số 1, Trần Hữu Dực, P. Mỹ Đình 2, Q. Nam Từ Liêm, Tp. Hà Nội
- Điện thoại: [096.133.6133](tel:096.133.6133)
- Email: [cskh@viettel-solution.vn](mailto:cskh@viettel-solution.vn)