

RheosProposal

Rheos Proposal

Abstract

Rheos is an open source stream data platform for large-scale cloud computing environments built around Kafka, and other distributed real-time continuous and periodic stream computation systems.

Proposal

While Kafka has given us core capabilities in stream processing, managing a large, highly distributed, real-time data pipelines running on cloud spanning across security zones and data centers is hard without automation, and supporting services. Rheos is intended to address majority of those concerns.

With Rheos, you can:

- create and perform lifecycle management on a data pipeline that consists of Kafka, Zookeeper, and the supported stream processing technologies such as Storm, Flink and etc.
- provide a real-time data pipeline at scale in a highly available manner
- browse existing or discover new streams
- provide shared and curated streams in a multi-tenant distributed cloud environment spanning across data centers
- provide a data pipeline that is easy to use by the developers and easy to manage
- mirror data streams from one datacenter to another for HA reasons
- move data streams across security zones

Overview of Rheos

Rheos includes:

Streaming As A Service (StraaS): A cloud service that provisions and provides full lifecycle management (LCM) for Zookeeper, Kafka, Storm, and [Mirror Maker](#) clusters. StraaS is built on a modular architecture with a pluggable extension and frameworks. This combination allows StraaS to create and perform LCM on a stream pipeline running on any cloud platforms (such as [OpenStack](#), AWS, Google Cloud and etc.).

Stream Metadata Service: A metadata system that provides a system of record for each stream and the associated producer and consumers known to the system. The recorded information includes

- the physical (cluster) location of a topic or a stream processing job/topology
- data durability, retention policy, partition, producer, consumer information
- Source and target data mirroring information
- topic schema information
- and more

Stream Discovery: Prior to accessing a stream or deploying a stream processing job, one must “register” the Kafka topic, stream producer/consumer or the job with the Stream Metadata Service. With this, Kafka topics, broker list along with the associated schemas can easily be discovered or browsed via Rheos Portal or REST API. That is, no hard coding of broker list in client code!

Kafka Proxy Server: To allow any out-of-the-box tool, framework, and programming language to seamlessly produce or consume data in Rheos, a Rheos Kafka proxy server that implements Kafka Protocol is provided to intercept the initial client request so that it can identify which Kafka cluster the topic resides on via the Metadata Service.

Operation Management Service: rheos performs stream, producer, consumer life cycle management operations with a set of pre-defined Standard Operation Procedure (SOP) in the Operation Management Service. Each SOP has a series of steps that need to be performed on the Guest instance via StraaS.

Data Mirroring Service: A service that is built around Kafka's [MirrorMaker](#). It is used to set up Mirror Maker instances and mirror a group of topics from a Kafka cluster to one or more target Kafka clusters via a REST API. Through the API, one can start and stop the mirroring of a topic group.

Data Mirroring is used to mirror Kafka data across regions and within a single region across availability zones for high availability reasons. In addition, Data Mirroring is used to provide data movement from one security zone to another.

Background

Access to near real-time data is no longer an option, but a must for eBay to be more efficient and effective to compete in today's eCommerce landscape. Due to lack of data freshness, feeds that were used to place automatic bids on Google adWords, to surface eBay listings on Google Paid Search were using 24-36 hours old data by daily Teradata extracts. That often led to producing ads with expired, deleted or out-of-stock listings which ultimately results in fewer clicks, conversions and eventual loss of revenue. With Rheos in place, eBay is able to deliver feeds to Google for bidding and ad content every one minute in Paid Search to promote hot listings, trending items and daily deals. The improved frequency and quality of feeds present great opportunity for big GMB growth.

Besides growth opportunity, the stream data platform will also enable better risk and fraud detection for our brand, eBay buyers and sellers provide an effective monitoring, alerting and notification system provide eBay user behavioral data in seconds latency surface eBay site data through the Data Change Stream and other opportunities in analytics and machine learning

eBay started this journey since 2015 with StraaS (Streaming as a Service), Kafka and Storm.

To date, Rheos has been able to process ~ 100 billion messages daily.

Rationale

Rheos allows you to provision and perform life cycle management on the assets on the cloud platform of your choice. Because of StraaS' modular architecture with pluggability, one can easily adopt a new distributed stream and batch data processing technology with little development effort. Along with this, the established core concepts and capacities (such as stream namespace, data mirroring to fulfil both availability and data movement across security zones, stream discovery with a system of record metadata system, offset-based failover mechanism and a Kafka Proxy Server) will undoubtedly benefit many data stream platforms. More importantly, adopt Rheos an Apache project allows faster innovation from the community.

Current Status

Meritocracy

Rheos was initially developed based on ideas from eBay Data Platform developers. As a project under incubation, we are committed to expanding our effort to build an environment which supports a meritocracy. We are focused on engaging the community and other related projects for support and contributions. Moreover, we are committed to ensure contributors and committers to Rheos come from a broad mix of organizations through a merit-based decision process during incubation. We believe strongly in Rheos' initial goals of built for ease of use for developers and ease of management for site operations.

Community

Rheos seeks to develop developer and user communities during incubation.

Core Developers

Rheos is currently being designed and developed by engineers from eBay Inc.

Alignment

The ASF is a natural host for Rheos given that it is already the home of Kafka, Storm, Zookeeper, Spark, Flink and other emerging cloud and streaming software projects. Rheos intends to provide a cloud agnostics with built-in flexibility to adopt new streaming technology easily.

Known Risks

Orphaned Products

The core developers of Rheos team plan to work full time on this project. There is very little risk of Rheos getting orphaned since at least one large company (eBay) is extensively using it as their near real-time data platform in production.

Inexperience with Open Source

eBay believes strongly in open source and the exchange of information to advance new ideas and work. Examples of this commitment are active OSS projects such as Apache Kylin (<http://kylin.apache.org/>) and Apache Eagle (<https://eagle.apache.org/>). Our submission to the Apache Software Foundation is a logical extension of our commitment to open source software.

Homogenous Developers

The initial committers in this proposal belong to eBay. We expect our entry into incubation will allow us to expand the number of individuals and organizations participating in Rheos development.

Reliance on Salaried Developers

Rheos has been developed by salaried developers supporting eBay Near Real-Time (NRT) Data Platform. Though the open source version of Rheos will have default implementation in certain areas, eBay NRT Data Platform running in production will absolutely be based off the open source version. We expect our reliance on salaried developers will decrease over time during incubation.

Relationships with Other Apache Products

Rheos directly interoperates with or utilizes the following Apache projects, libraries, and frameworks

Basic:

- Apache Maven
- Apache Commons
- Apache Logging Services

- Google Guava

Distributed computing:

- Apache Zookeeper
- Apache Curator
- Apache Kafka
- Apache Storm

Data I/O, frameworks and libraries:

- Apache Avro
- Confluent.io Schema Registry
- Jackson
- [OpenStack](#) Common Libraries (Oslo)
- Taskflow

REST Service:

- Spring

Being part of Apache's Incubation community, could help with a closer collaboration among these four projects and as well as others.

An Excessive Fascination with the Apache Brand

Rheos is proposing to enter incubation at Apache in order to help efforts to diversify the committer-base, not so much to capitalize on the Apache brand. The Rheos project is in production use already inside eBay, but is not expected to be an eBay product for external customers. As such, the Rheos project is not seeking to use the Apache brand as a marketing tool.

Documentation

Coming soon

Initial Source

Rheos has been under development since 2015 by a team of engineers at eBay Inc.

External Dependencies

All external dependencies are licensed under an Apache 2.0 or Apache-compatible license. As we grow the Rheos community we will configure our build process to require and validate all contributions and dependencies are licensed under the Apache 2.0 license or are under an Apache-compatible license.

Cryptography

Rheos adheres to eBay Security Control and Guidelines.

Rheos client connections to Kafka must be authenticated and authorized via the broker's security measures with eBay's Identity Service. Sensitive data will only be available in dedicated Kafka cluster(s) with TLS enabled.

Required Resources

Mailing List

rheos-private for private PMC discussions (with moderated subscriptions)

- dev@rheos.incubator.apache.org
- user@rheos.incubator.apache.org
- private@rheos.incubator.apache.org
- commits@rheos.incubator.apache.org

Source Control

Git is the preferred source control system. We request a Git repository for Rheos with mirroring to [GitHub](#) enabled

<https://git-wip-us.apache.org/repos/asf/incubator-rheos.git>

Issue Tracking

We request the creation of an Apache-hosted JIRA. Jira ID: (RHEOS)

Other Resources

The existing code already has unit tests so we will make use of existing Apache continuous testing infrastructure. The resulting load should not be very large.

Initial Committers

- Ankur Bansal [now at Uber]
- Michael Chiocca [mchiocca@ebay.com]
- Lubin Liu [lubliu@ebay.com]
- Subash Ramanathan [sramanathan@ebay.com]
- Viswa Vutharka [vvutharkar@ebay.com]
- Tao Xiao [taxiao@ebay.com]
- Xin Xu [xinxu1@ebay.com]
- Connie Yang [cyang@ebay.com]
- Wayne Zhou [jianbzhou@ebay.com]

Affiliations

Most of the initial committers are employees of eBay Inc and one from Uber. The nominated mentors are employees of Hortonworks, eBay.

Sponsors

Champion

- Henry Saputra [hsaputra@ebay.com]

Nominated Mentors

- Julian Hyde [jhyde@hortonworks.com]
- Taylor Goetz [tgoetz@hortonworks.com]
- Henry Saputra [hsaputra@ebay.com]

Sponsoring Entity

The Apache Incubator