# Apache Linkis Data Processing Practice

李孟

2022.10.22

# Background

The basic composition of the original version of the big data platform framework:
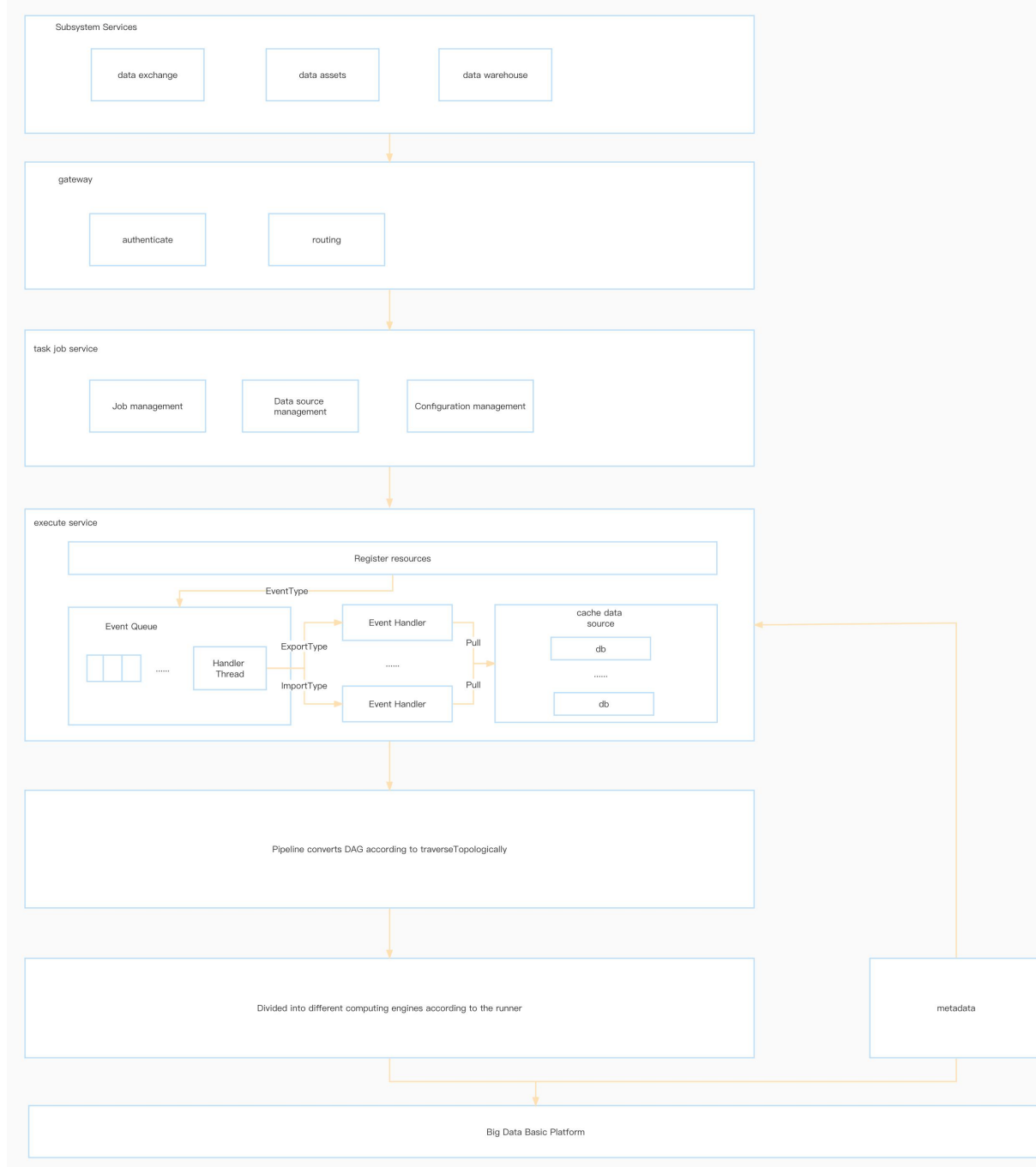
Beam data processing and Hive Hook as metadata

# Background

shortcoming

- external access single

- high development cost

- metadata access is complex

# Selection

Apache Livy

Apache Zeppelin

Netflix Geine

openLooKeng

Apache Linkis

# Selection

**Linkis builds a decoupling computing middleware layer with the ability to connect, expand, manage, orchestrate and reuse**

- Computing Governance Services

- Public Enhancement Services

- Microservice Governance Service

- Complete component base

- Community

# Scenario and Value

DataSphereStudio(DSS)

Linkis

Scriptis

DolphinScheduler
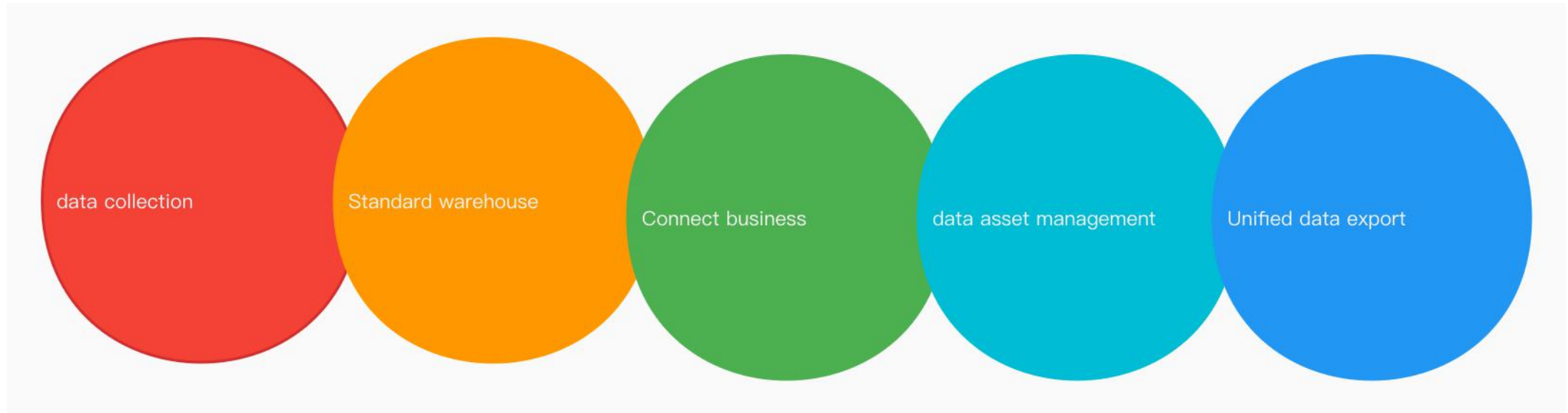
Exchangis

Streamis

......

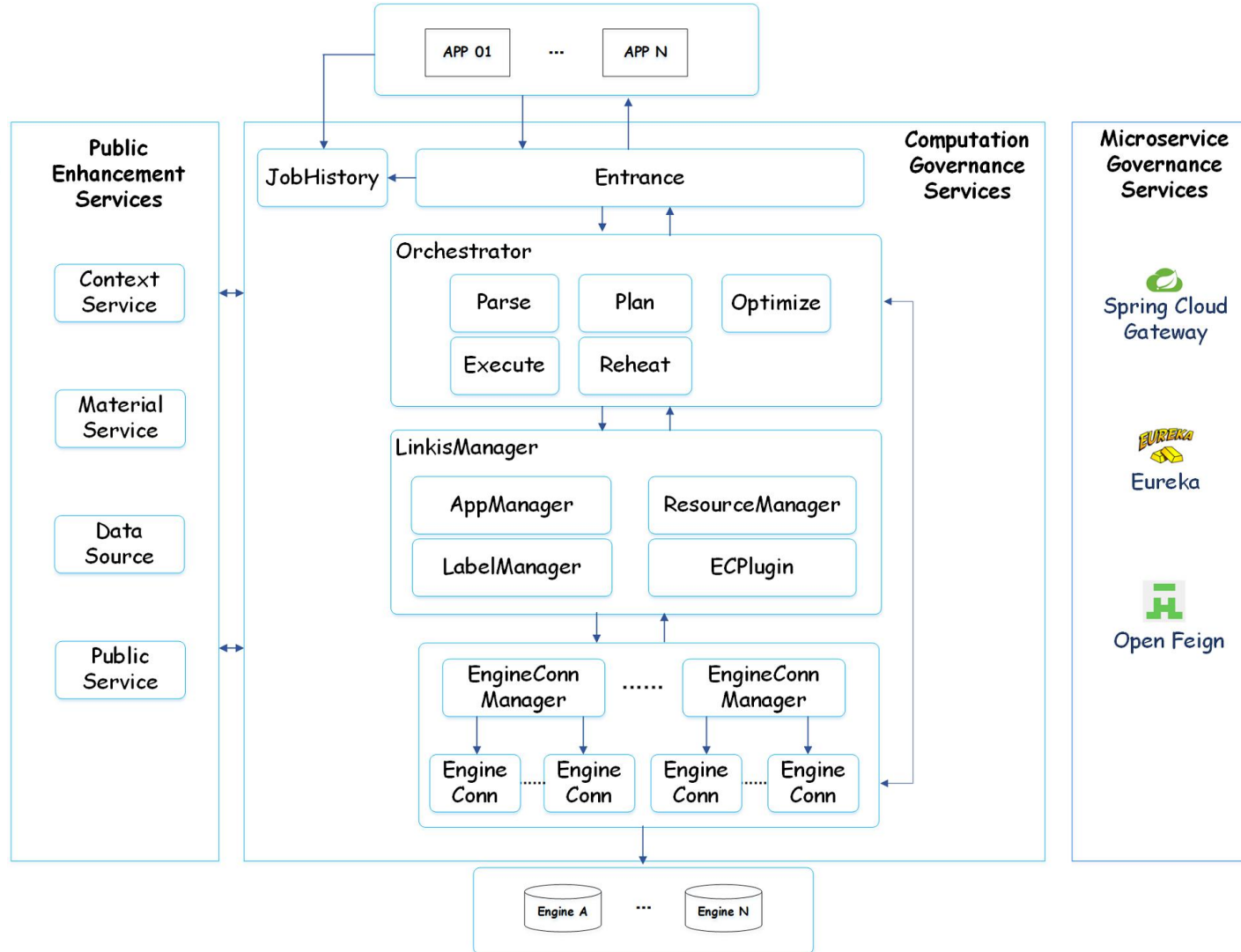# Scenario and Value

**Provide the basis for data processing**

- data assets

- data model

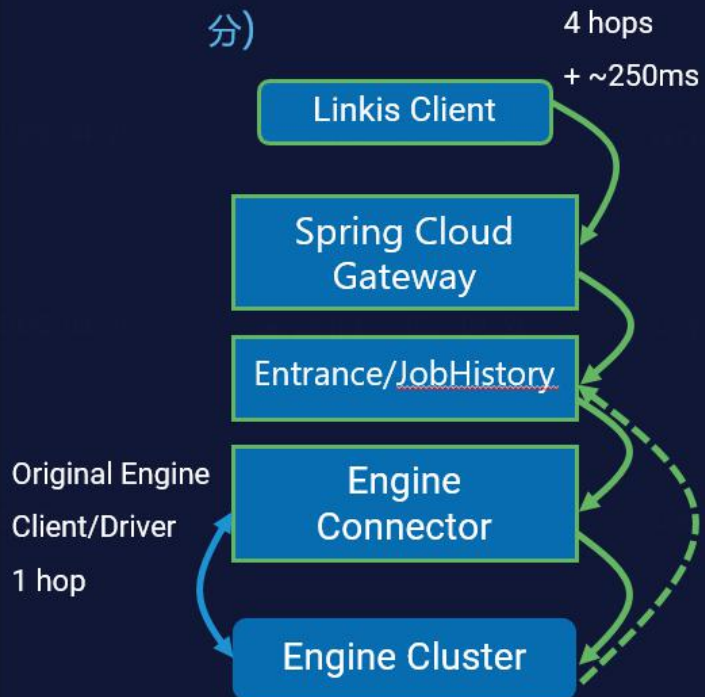- data development

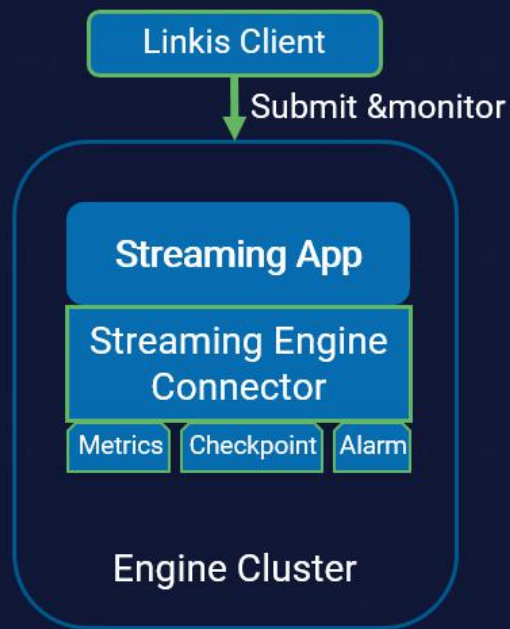- data scheduling

# Scenario and Value

# Architecture-Linkis

# Architecture-Linkis
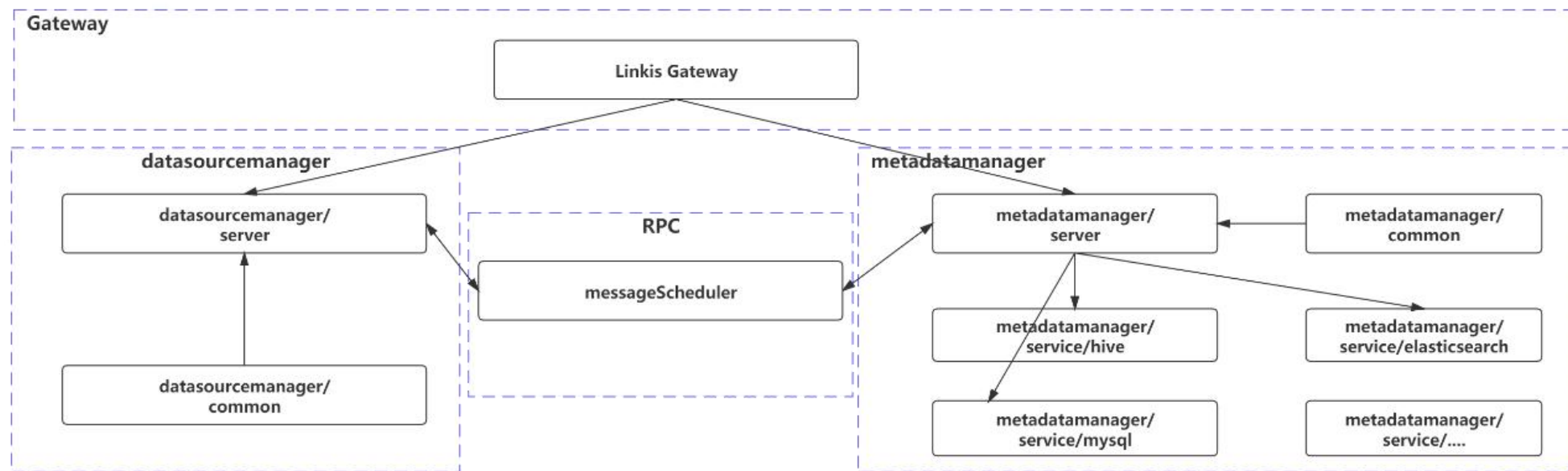
# Architecture-DSS

# Data Governance-DataSource

- Basic functions of a data source service provider

- The data source service provides some basic information about the data source environment

- Supports parameter dynamic lookup library table generation for multiple data source types

- Supports data sources to be divided according to user, creation source, and creation system

- Support for different data source connection tests

- A metadata query service that supports a certain data source

# Data Governance-DataSource

# Data Governance-Apache Atlas

Atlas is a data governance and metadata framework tightly coupled to the Hadoop ecosystem.
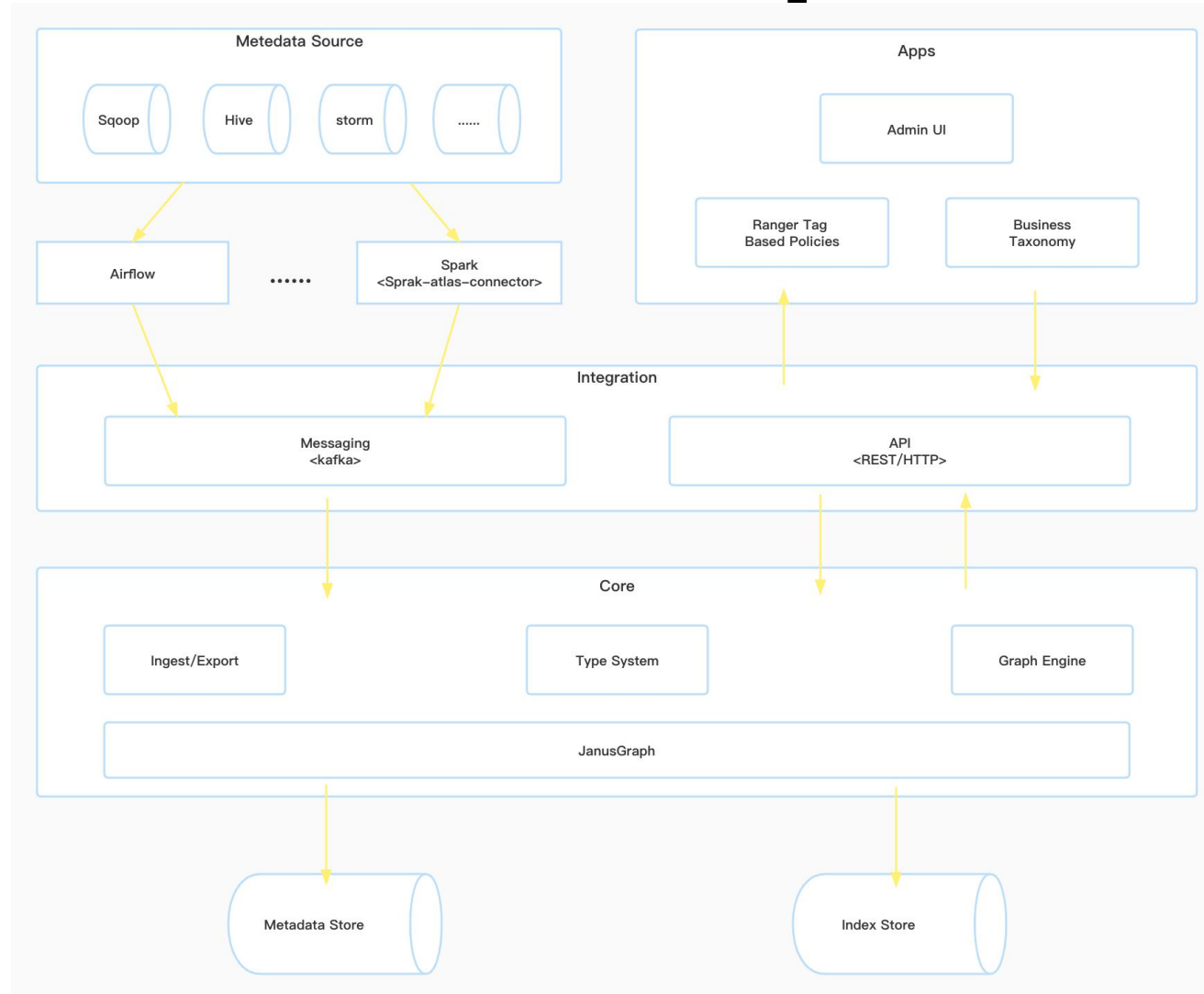
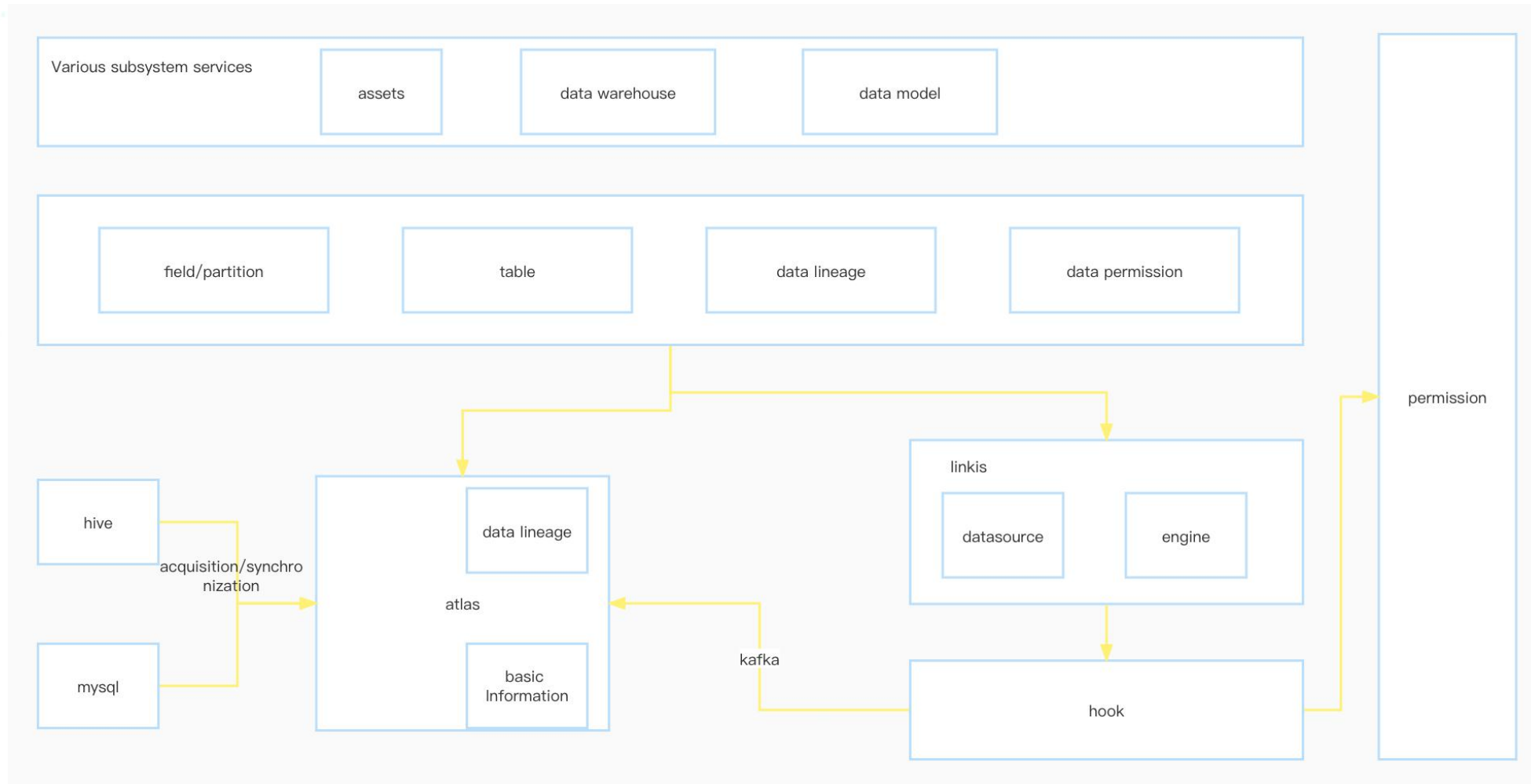# Data Governance-Apache Atlas

- type system
- graphics engine
- capture/export
- api
- messaging
- metadata sources
- atlas admin ui
- tag based policies

# Data Governance-Apache Atlas

# Data Governance-Apache Atlas

# Data Governance-Assets

Data asset management to make data traceable, usable, and trusted

- asset overview

- asset catalog

# Data Governance-Assets

# Data Governance-DataWarehouse

**Data warehouse standard specification management.**

- Subject Domain Management

- Data warehouse management

- Modifier management

- Statistical cycle management

# Data Governance-DataWarehouse

# **Data Governance-DataModel**

The content described by the data model has three parts, namely data operations, data constraints and data structures.

- table management

- dimension management

- metric management

- indicator management

- tag management

# Data Processing-Exchangis
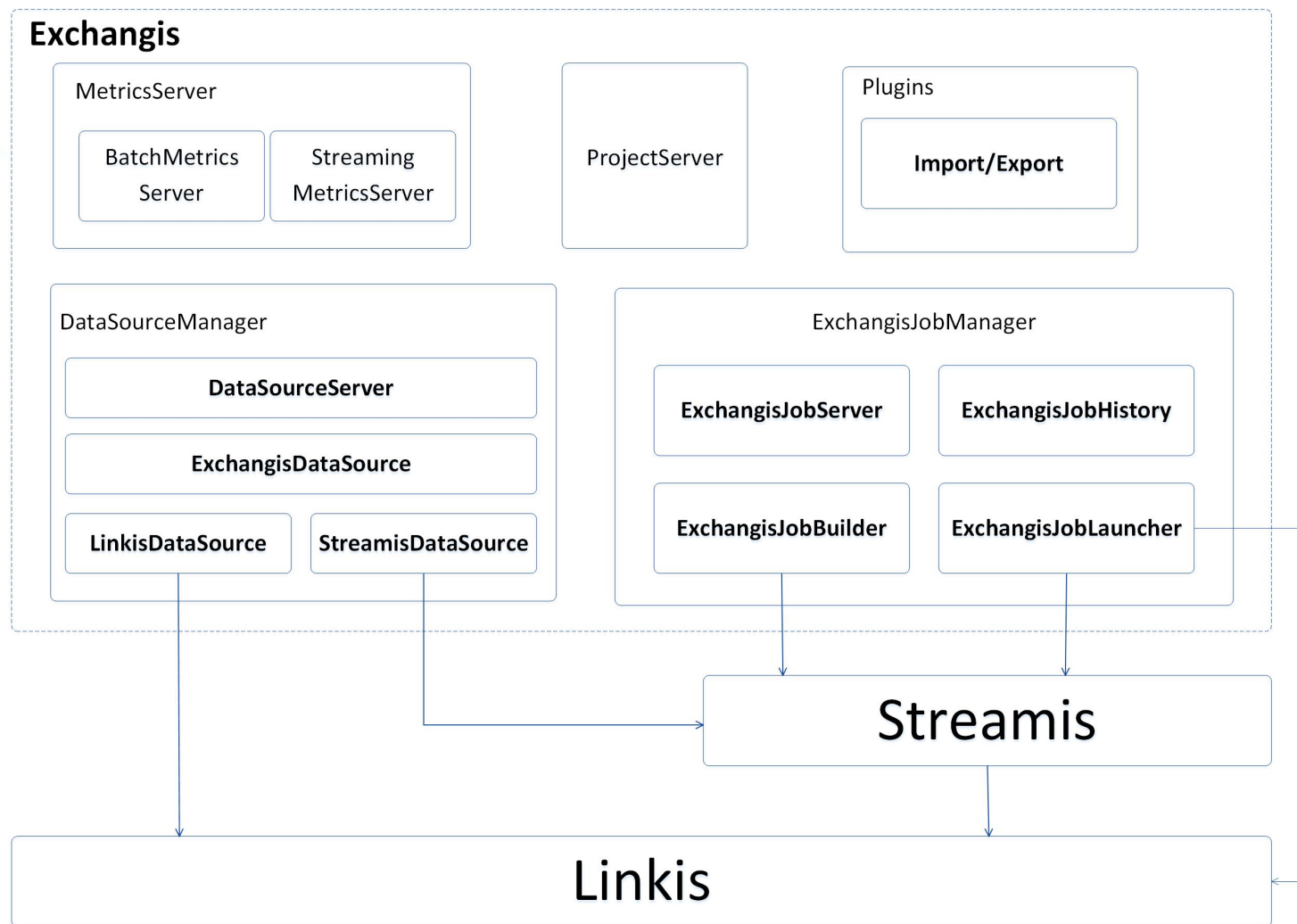
A data exchange tool that supports the synchronization of structured and unstructured data transfers between heterogeneous data sources.
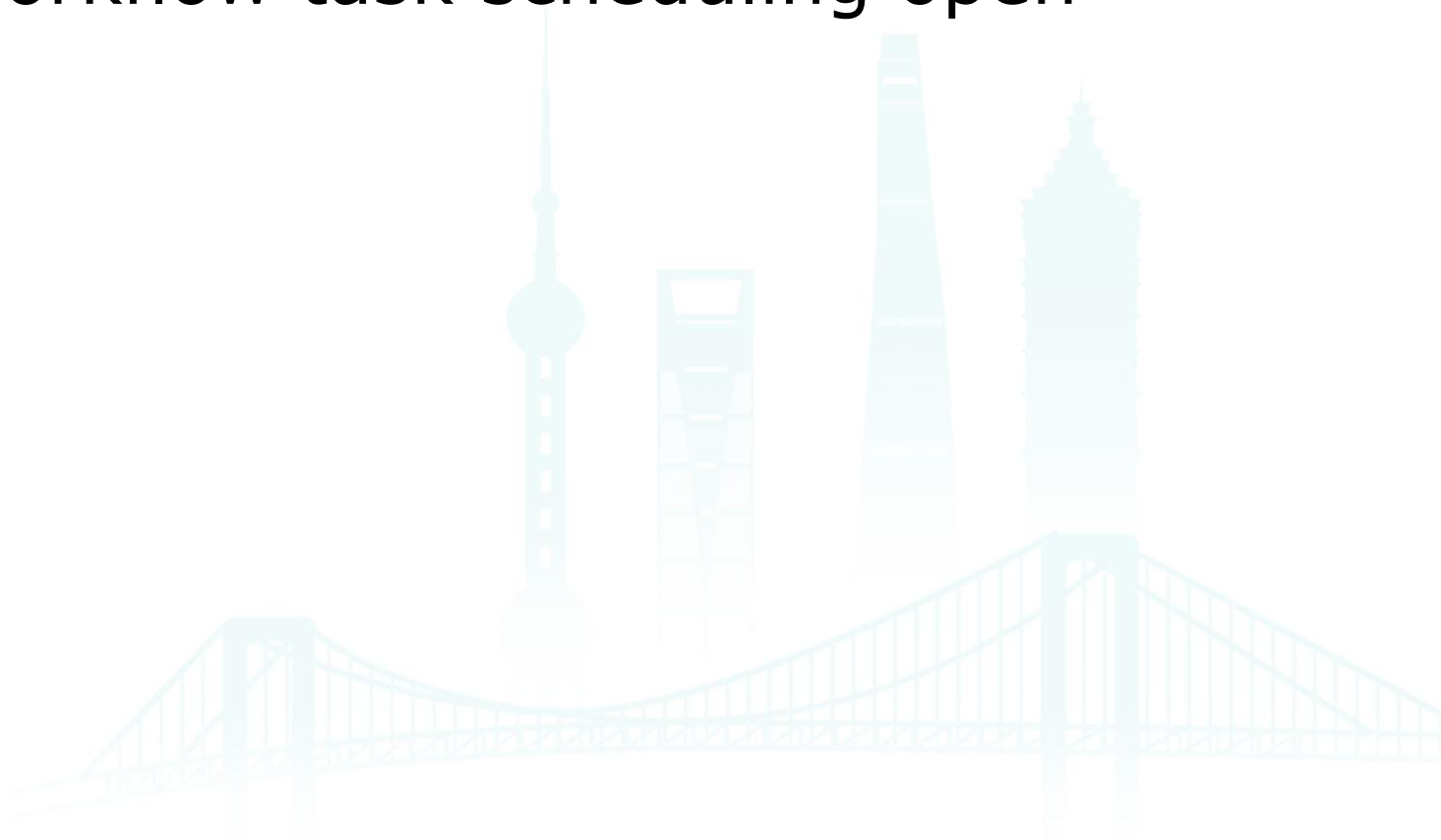
• Lightweight data source management

• High stability, fast response data synchronization task execution

• Open up with DSS workflow, one-stop big data development portal

# Data Processing-Exchangis

**Exchangis**

MetricsServer

| BatchMetrics Server | Streaming MetricsServer |

ProjectServer

Plugins

**Import/Export**

DataSourceManager

**DataSourceServer**

**ExchangisDataSource**

**LinkisDataSource**   **StreamisDataSource**

ExchangisJobManager

**ExchangisJobServer**   **ExchangisJobHistory**

**ExchangisJobBuilder**   **ExchangisJobLauncher**

Streamis

Linkis

# Data Processing-Dolphinscheduler

Apache DolphinScheduler is a distributed and easily extensible visual DAG workflow task scheduling open source system.
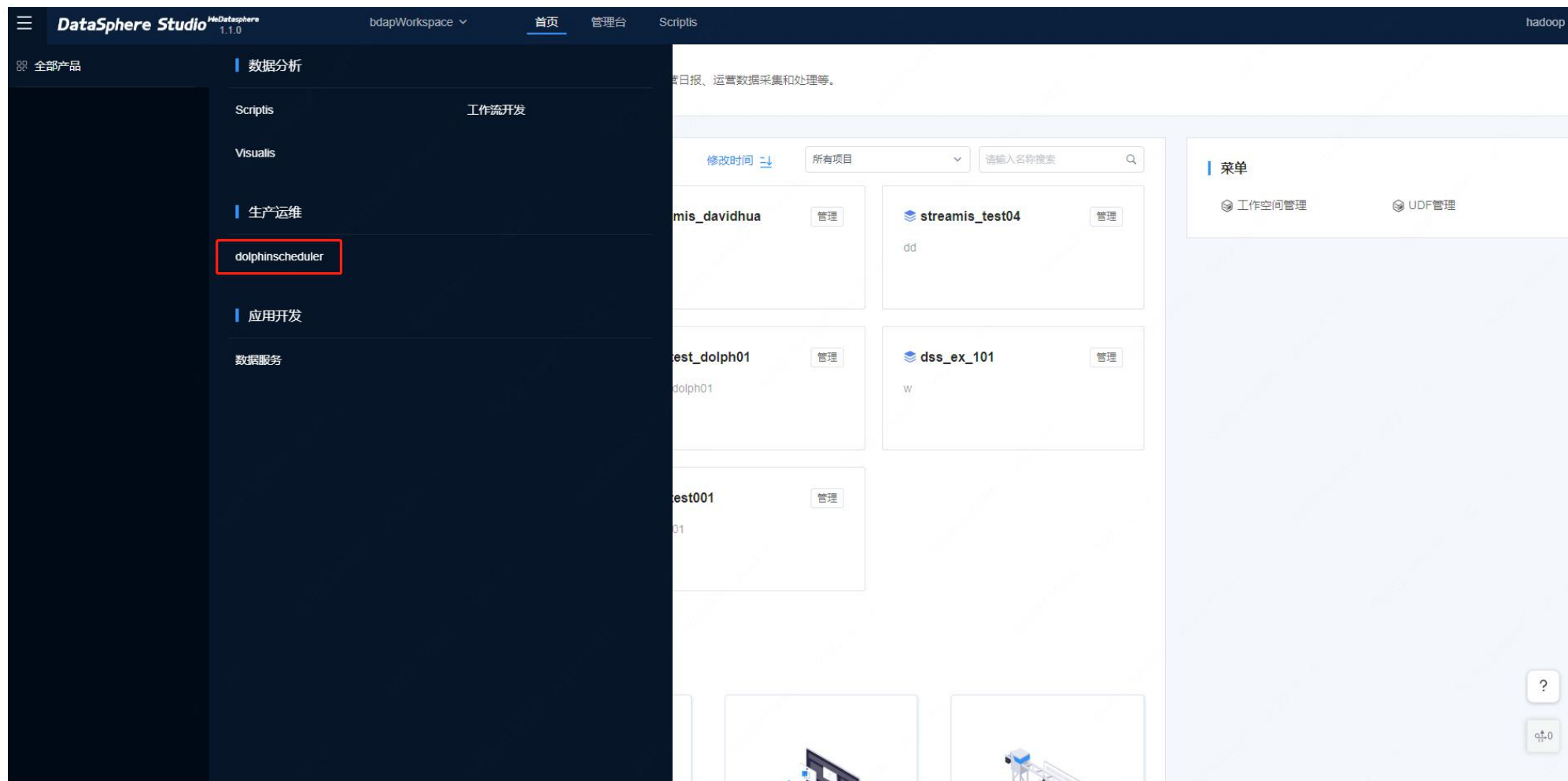
# Data Processing-Dolphinscheduler

- AppConn is the core concept of DSS that can easily and quickly integrate various upper-layer web systems.

- Dolphinscheduler can publish integration components to the scheduler with DSS publishing function
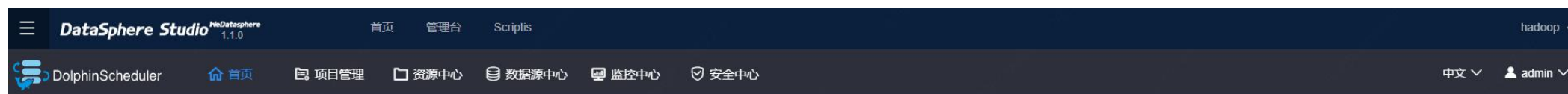
# Data Processing-Dolphinscheduler

DolphinScheduler can normally schedule the workflow node jobs of DataSphere Studio. You also need to install the dss-dolphinscheduler-client plugin, which is used to execute DSS workflow node jobs.

# Data Processing-Dolphinscheduler

# Data Processing-Dolphinscheduler

# Data Link

# Future Outlook

- Improvement of data integration capabilities

- Data governance related service release

- WeDataSphere community-related components continue to improve and strive for Apache

# Open Source

All the components shared this time have been open sourced in the community. I look forward to more friends to join and communicate with each other so that better products can be iterated. I always believe that the open source community is the soul of the product, and it is connected with its pulse to make Products have better development.

Over the years, I have met many friends in the WeDataSphere community, learned a lot, and gave back the results to the community. I will continue to work hard and continue to walk with you.