

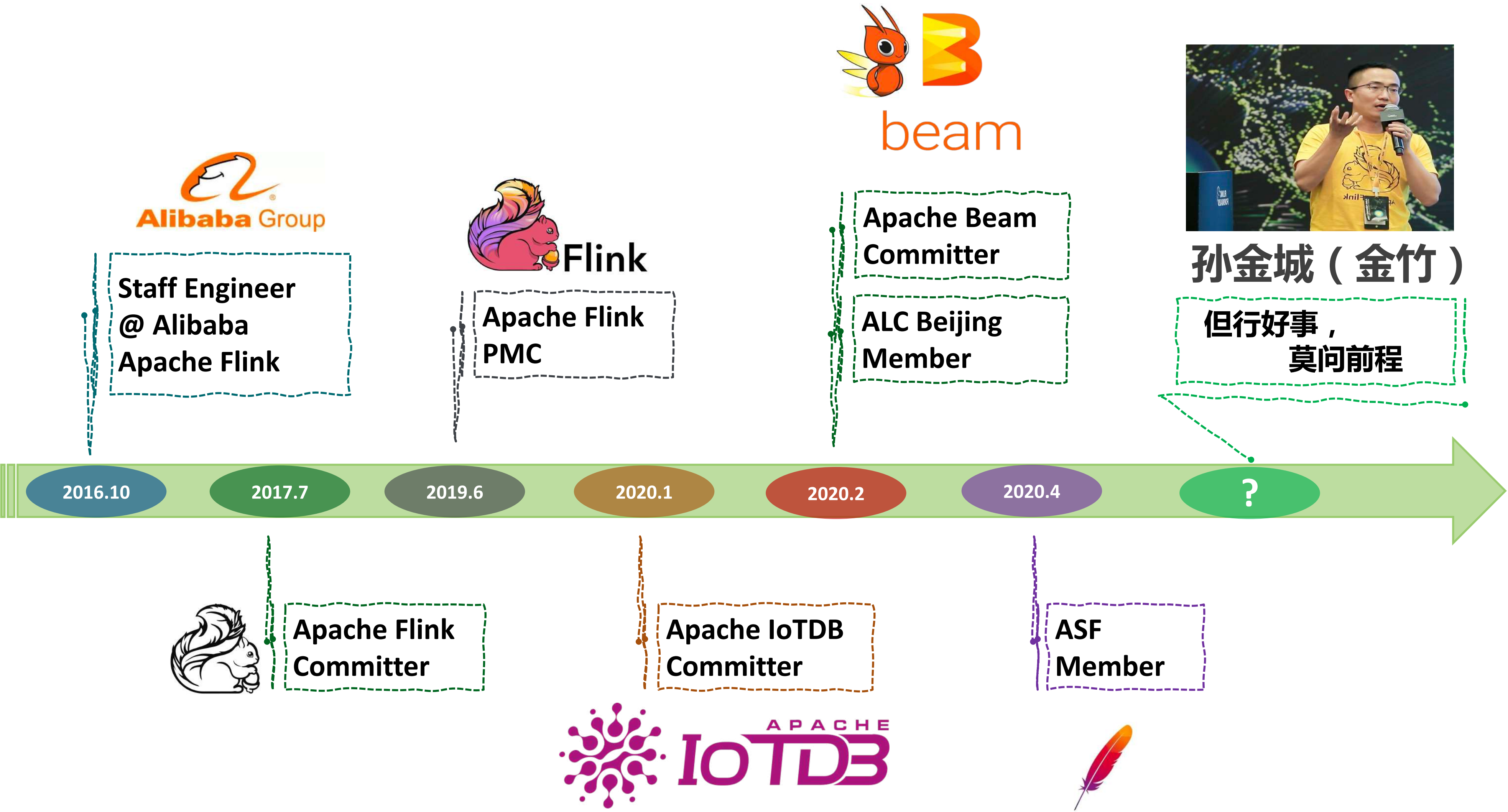
Python on Flink & Flink on Zeppelin

章剑锋(简锋)&孙金城(金竹)

Staff Engineer @Alibaba

04/22/2020

About Me



Agenda

- ① Why PyFlink
- ② Status of PyFlink
- ③ Future of PyFlink
- ④ Flink on Zeppelin



直播
85
分钟



博客
两篇



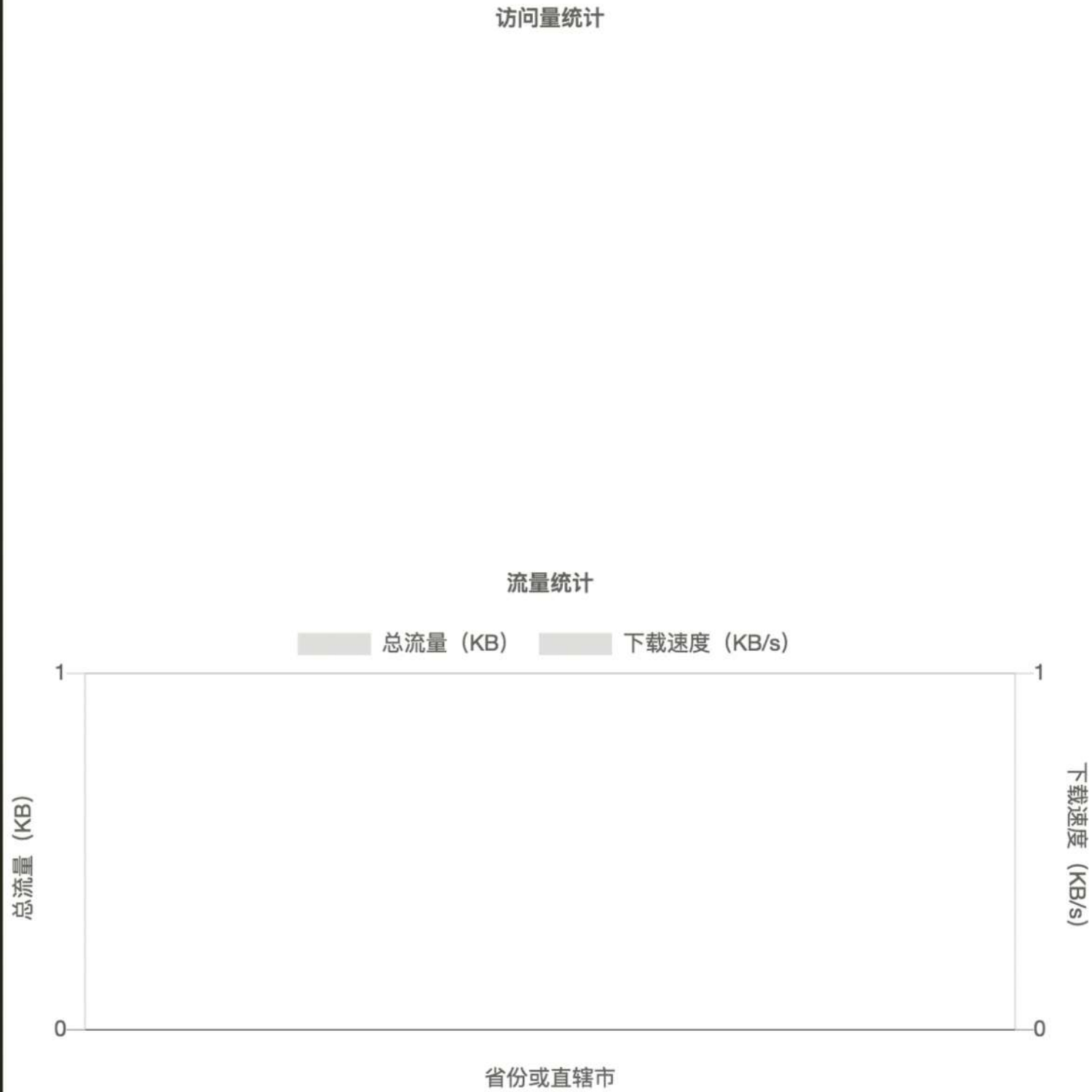
《必修课！一小时吃透PyFlink》

```
1 import os
2
3 from pyflink.datastream import StreamExecutionEnvironment
4 from pyflink.table import StreamTableEnvironment, EnvironmentSettings
5 from enjoyment.cdn.cdn_udf import ip_to_province
6 from enjoyment.cdn.cdn_connector_ddl import kafka_source_ddl, mysql_sink_ddl
7
8 # 创建Table Environment, 并选择使用的Planner
9 env = StreamExecutionEnvironment.get_execution_environment()
10 t_env = StreamTableEnvironment.create(
11     env,
12     environment_settings=EnvironmentSettings.new_instance().use_blink_planner().build()
13 )
14 # 创建Kafka数据源表
15 t_env.sql_update(kafka_source_ddl)
16 # 创建MySQL结果表
17 t_env.sql_update(mysql_sink_ddl)
18
19 # 注册IP转换地区名称的UDF
20 t_env.register_function("ip_to_province", ip_to_province)
21
22 # 添加依赖的Python文件
23 t_env.add_python_file(
24     os.path.dirname(os.path.abspath(__file__)) + "/enjoyment/cdn/cdn_udf.py")
25 t_env.add_python_file(os.path.dirname(
26     os.path.abspath(__file__)) + "/enjoyment/cdn/cdn_connector_ddl.py")
27
28 # 核心的统计逻辑
29 t_env.from_path("cdn_access_log")\
30     .select("uuid, "
31            "ip_to_province(client_ip) as province, " # IP 转换为地区名称
32            "response_size, request_time")\
33     .group_by("province")\
34     .select( # 计算访问量
35            "province, count(uuid) as access_count, "
36            # 计算下载总量
37            "sum(response_size) as total_download, "
38            # 计算下载速度
39            "sum(response_size) * 1.0 / sum(request_time) as download_speed") \
40     .insert_into("cdn_access_statistic")
41
42 # 执行作业
43 t_env.execute("pyflink_parse_cdn_log")
44
45
46
47
```

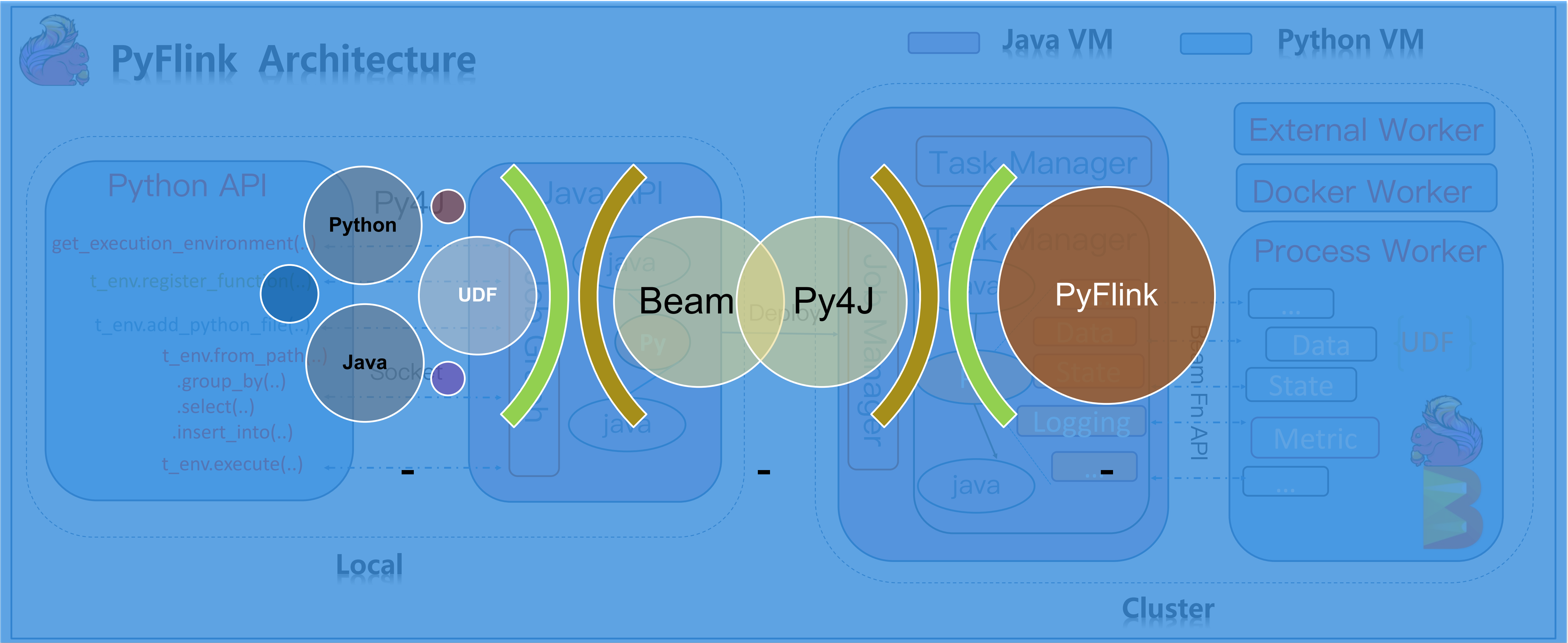
Kafka Message: 5934d8ae-d66f-4095-a01c-68a7ea8701a6,36.24.191.206,260,150280,https://www.aliyun.com/product/bastionhost?spm=5176.224200.h2v3icoap.119.1f956ed6muyb0V&aly_as=96NleqHS

Random Send Clear and Truncate Mysql

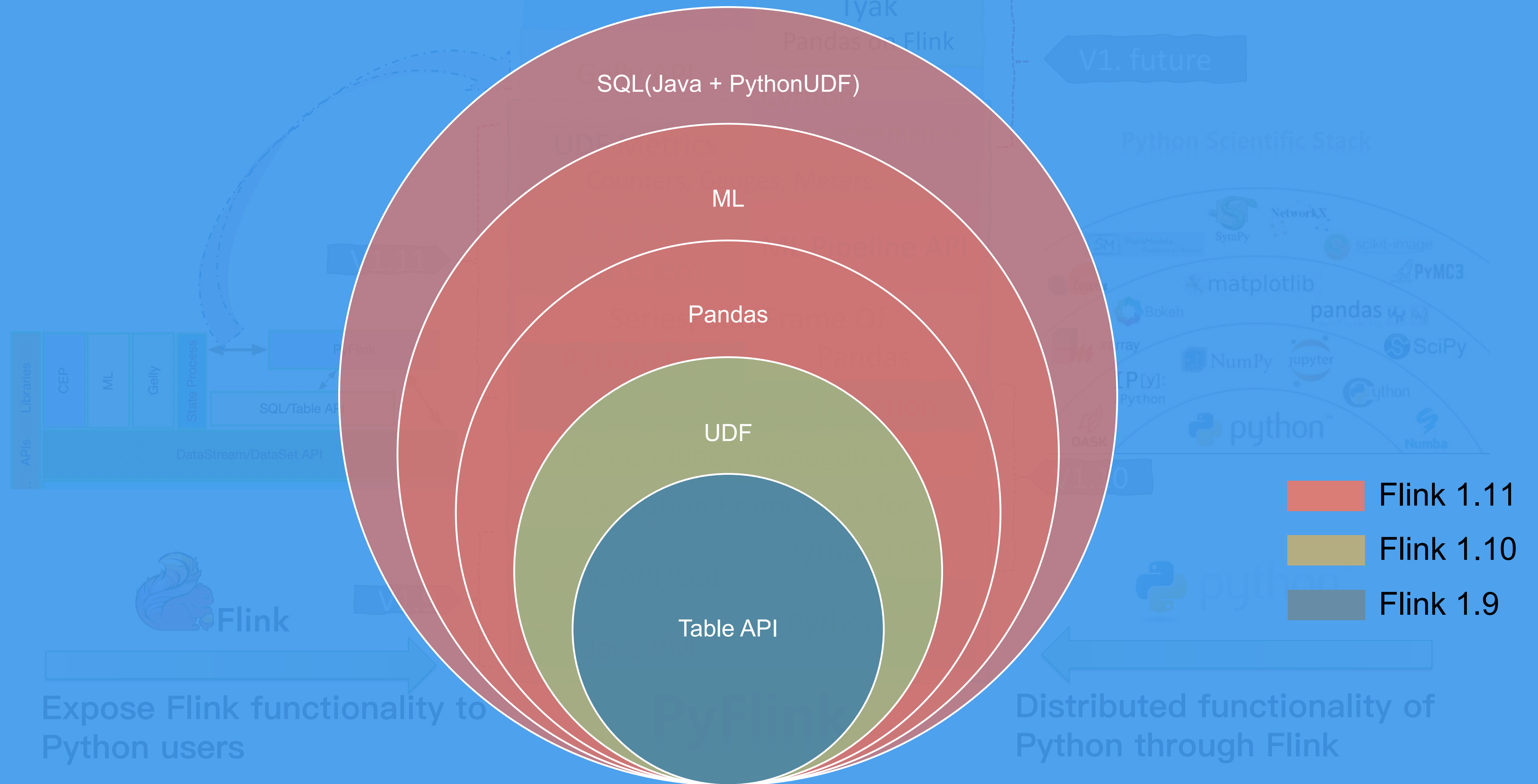
《PyFlink实现CDN日志实时分析》



PyFlink 架构



PyFlink 规划



PyFlink支持PandasUDF

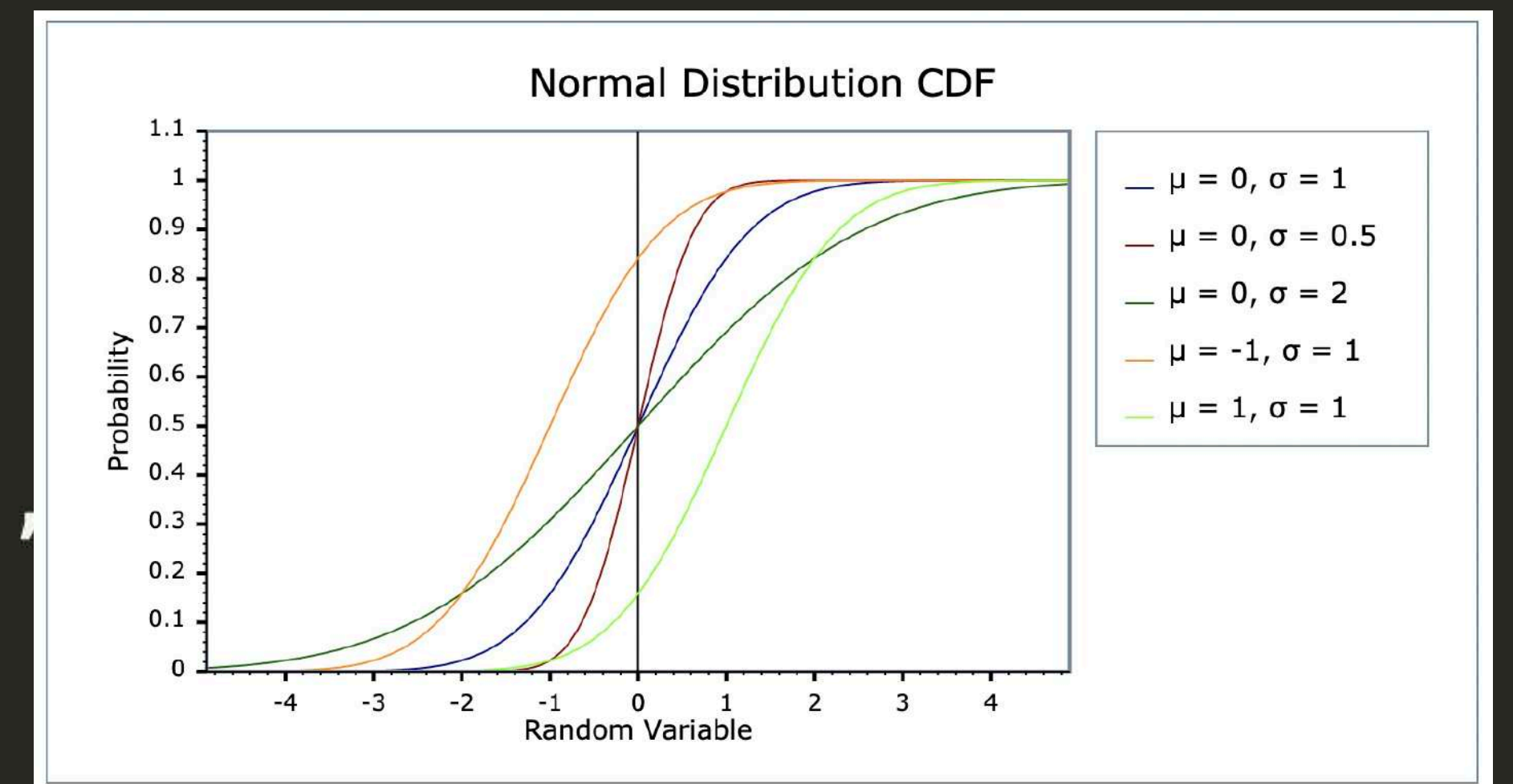
Functionality

Performance

Usability

Pandas UDF

```
1 import pandas as pd
2 from scipy import stats
3 from pyflink.table import DataTypes
4 from pyflink.table.udf import udf
5
6
7 @udf(input_types=[DataTypes.DOUBLE()],
8      result_type=DataTypes.DOUBLE(),
9      udf_type="pandas")
10 def cdf(v):
11     return pd.Series(stats.norm.cdf(v))
12
```



Functionality

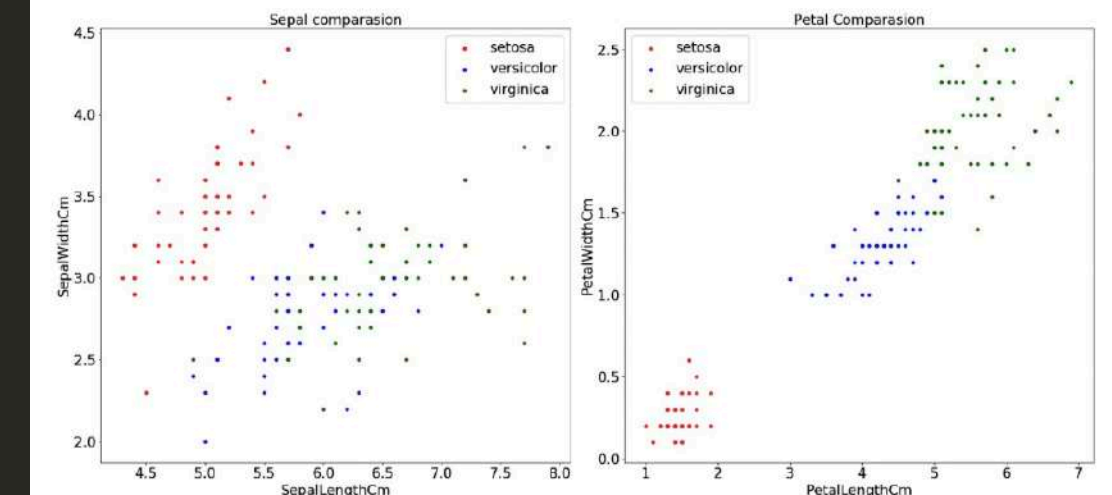
Performance

Usability

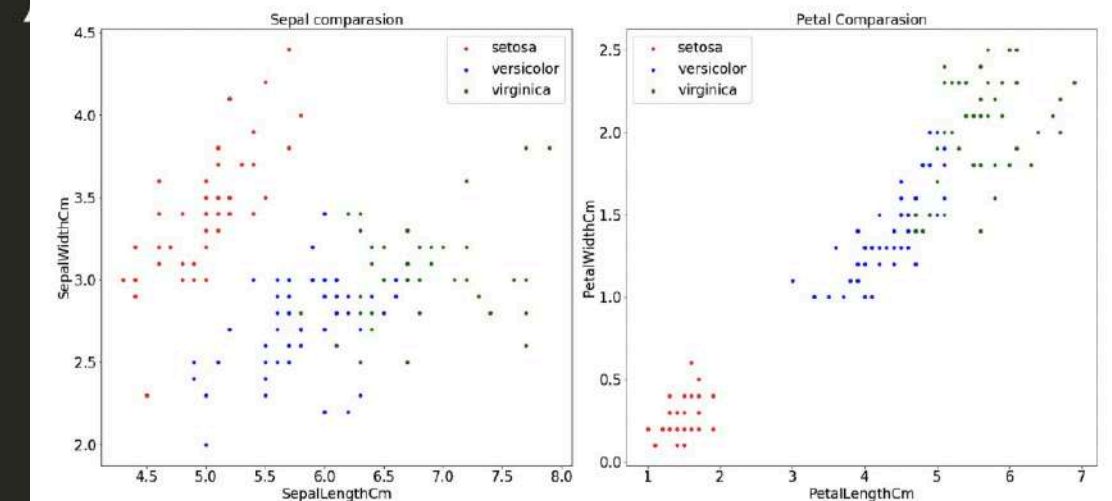
ML Pipeline API

```
25 # transformer
26 va = VectorAssembler()\
27     .set_selected_cols(["sepal_length",
28         "sepal_width", "petal_length", "petal_width"])\
29     .set_output_col("features")
30
31 # estimator
32 kmeans = KMeans() \
33     .set_vector_col("features") \
34     .set_k(3) \
35     .set_reserved_cols(["sepal_length",
36         "sepal_width", "petal_length", "petal_width", "category"])\
37     .set_prediction_col("prediction_result")
38
39 # pipeline
40 pipeline = Pipeline().append_stage(va).append_stage(kmeans)
41 pipeline \
42     .fit(t_env, sourceTable) \
43     .transform(t_env, sourceTable) \
44     .insert_into('kmeansResults')
45
46 t_env.execute('KmeansTest')
```

真实数据分布



预测结果



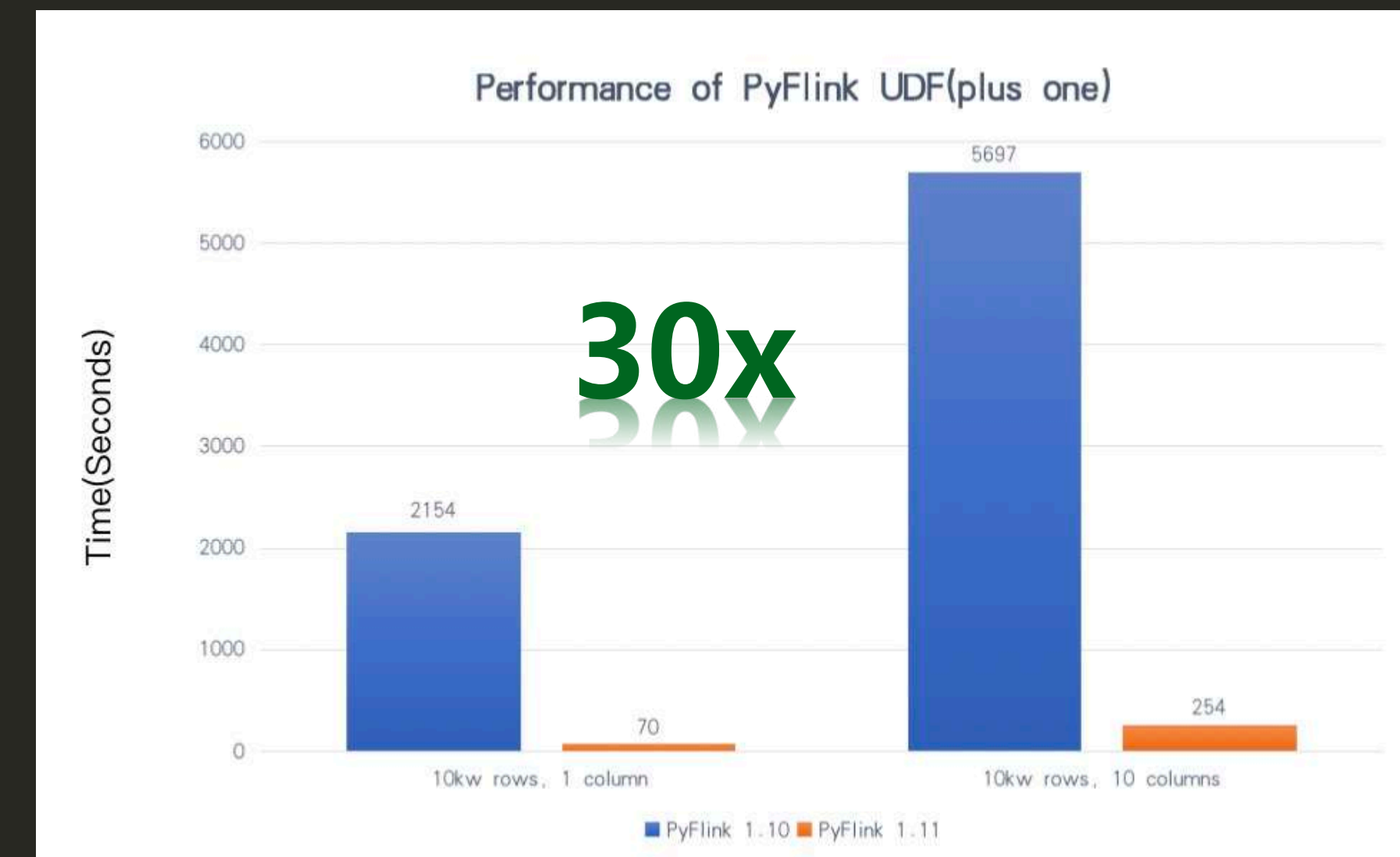
Functionality

Performance

Usability

High Performance

1. Code gen Python UDF
2. Add support for Cython
3. Optimizing Beam WindowValuedCoder
4. Replacing Flink row with Python list
5. ...



PyFlink支持Java SQL + Python UDF


Functionality

Performance

Usability

PyUDF in SQL Client

```
1 Syntax:
2
3 CREATE
4 [TEMPORARY | TEMPORARY SYSTEM] FUNCTION
5 [IF NOT EXISTS]
6 [catalog_name.db_name.]function_name AS identifier
7 LANGUAGE PYTHON;
8
9 Example:
10 CREATE TEMPORARY SYSTEM FUNCTION func1 AS 'pymodule.udf.func1' LANGUAGE PYTHON;
11
```



Welcome! Enter 'HELP;' to list all available commands. 'QUIT;' to exit.

```
Flink SQL> CREATE TEMPORARY SYSTEM FUNCTION func1 AS 'pymodule.udf.func1' LANGUAGE PYTHON;
> SELECT func1(f0) FROM tab1;
```


起点的选择？环境 vs 功能

环境

走那么远，

我们在寻找

一个起点

听说它在... 山那边，海那边...



Py Flink

起点的选择？环境

← → ↻ pypi.org/project/apache-flink/



apache-flink



pip install apache-flink

```
pip install apache-flink
```



jincheng:FFSF jincheng.sunjc\$



起点的选择？环境

在本地尝试以pre-job模式部署作业时，发现会提示如下报错，导致任务提交失败

```
RuntimeError: Python versions prior to 3.5 are not supported for PyFlink [sys.version_info(major=2, minor=7, micro=16, releaselevel='final', serial=0)].
```

```
[libprotobuf FATAL google/protobuf/compiler/cpp/cpp.cc:1370] CHECK failed: GeneratedDatabase()->Add(encoded_file_descriptor, size)
libc++abi.dylib: terminating with uncaught exception of type std::__glibcpp_exception: CHECK failed: GeneratedDatabase()
CHECK failed: GeneratedDatabase()
```

这些问题都可以解决！



但没必要作为学者的起点！

Python 版本问题

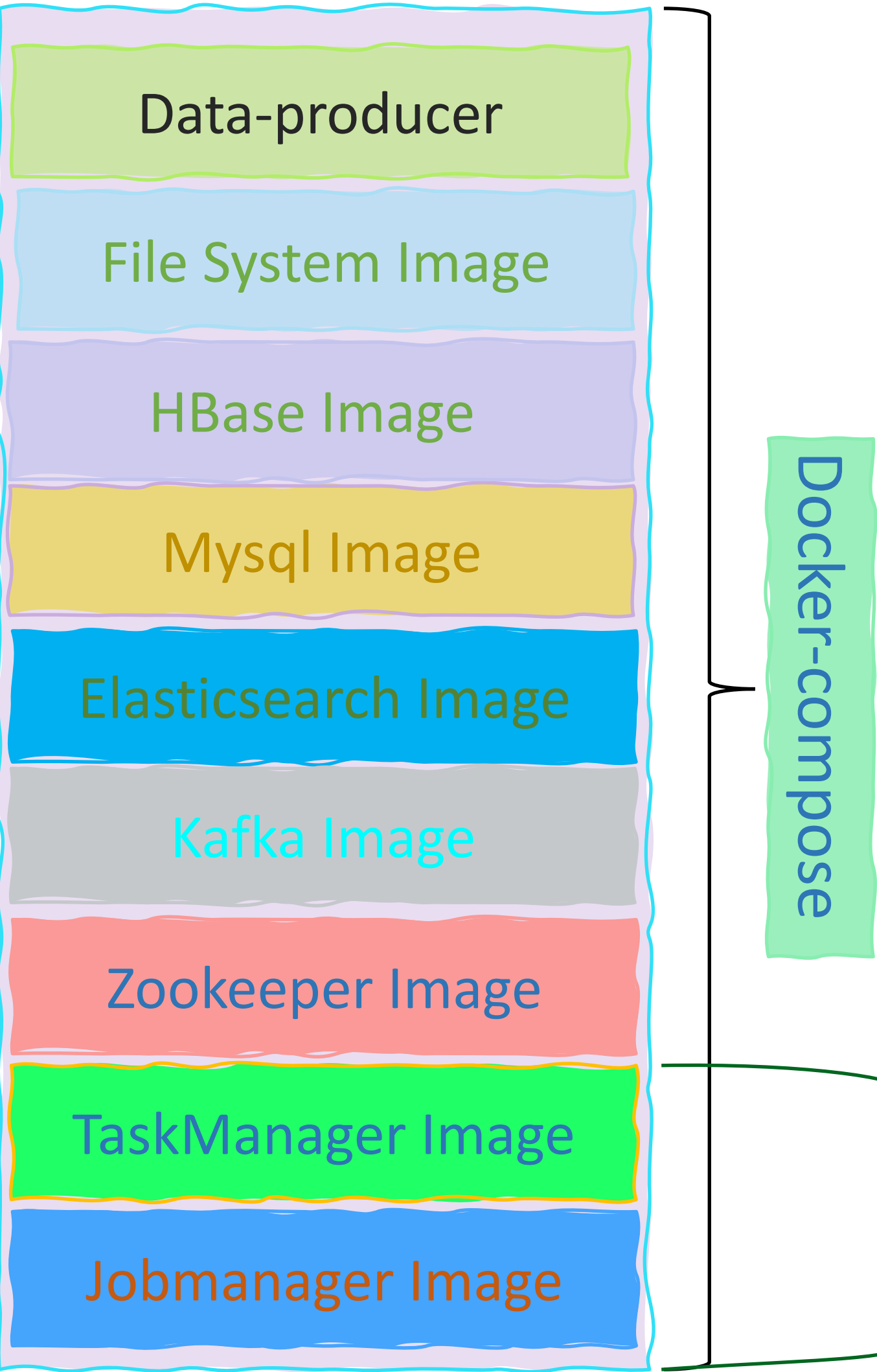
macos 10.15
(Catalina) 引起的

多Connector依赖问题

Apache Flink 扫雷系列 – Flink如何解决多JAR包依赖问题

Posted on 2020-03-31 | Edited on 2020-04-01 | In Apache Flink 扫雷系列 | Comments: | Views: 107

起点的选择？功能





```
docker-compose.yml
1 version: '2.1'
2 services:
3   jobmanager:
4     image: pyflink/playgrounds:1.10.0 ①
5     volumes:
6       - ./examples:/opt/examples
7     hostname: "jobmanager"
8     expose: <1 item>
10    ports: <1 item>
12    command: jobmanager
13    environment: <1 item>
15    taskmanager: <7 keys>
29    zookeeper:
30      image: wurstmeister/zookeeper:3.4.6
31      ports: <1 item>
33    kafka:
34      image: wurstmeister/kafka:2.12-2.2.1 ②
35      ports: <1 item>
37      depends_on: <1 item>
39      environment: <3 keys>
43      volumes: <1 item>
45    elasticsearch:
46      image: docker.elastic.co/elasticsearch/elasticsearch-oss:6.3.1 ③
47      environment: <4 items>
52      ports: <2 items>
55      ulimits: <2 keys>
62    data-producer: ⑤
63      image: fhueske/flink-sql-training:1-FLINK-1.9-scala_2.11
64      command: "java -classpath /opt/data/data-producer.jar com.ververica.sql_training.data_producer
65      depends_on: <3 items>
69      environment: <4 keys>
74    db:
75      image: mysql
76      command: --default-authentication-plugin=mysql_native_password ④
77      restart: always
78      ports:
79        - 3306:3306
80      environment: <2 keys>
83      volumes: <1 item>
85    adminer:
86      image: adminer
87      restart: always
88      ports:
89        - 8080:8080
```

FileSystem, HBase, Hive ...

[Docker-compose.yml](#)

Tags

This repository contains 2 tag(s).

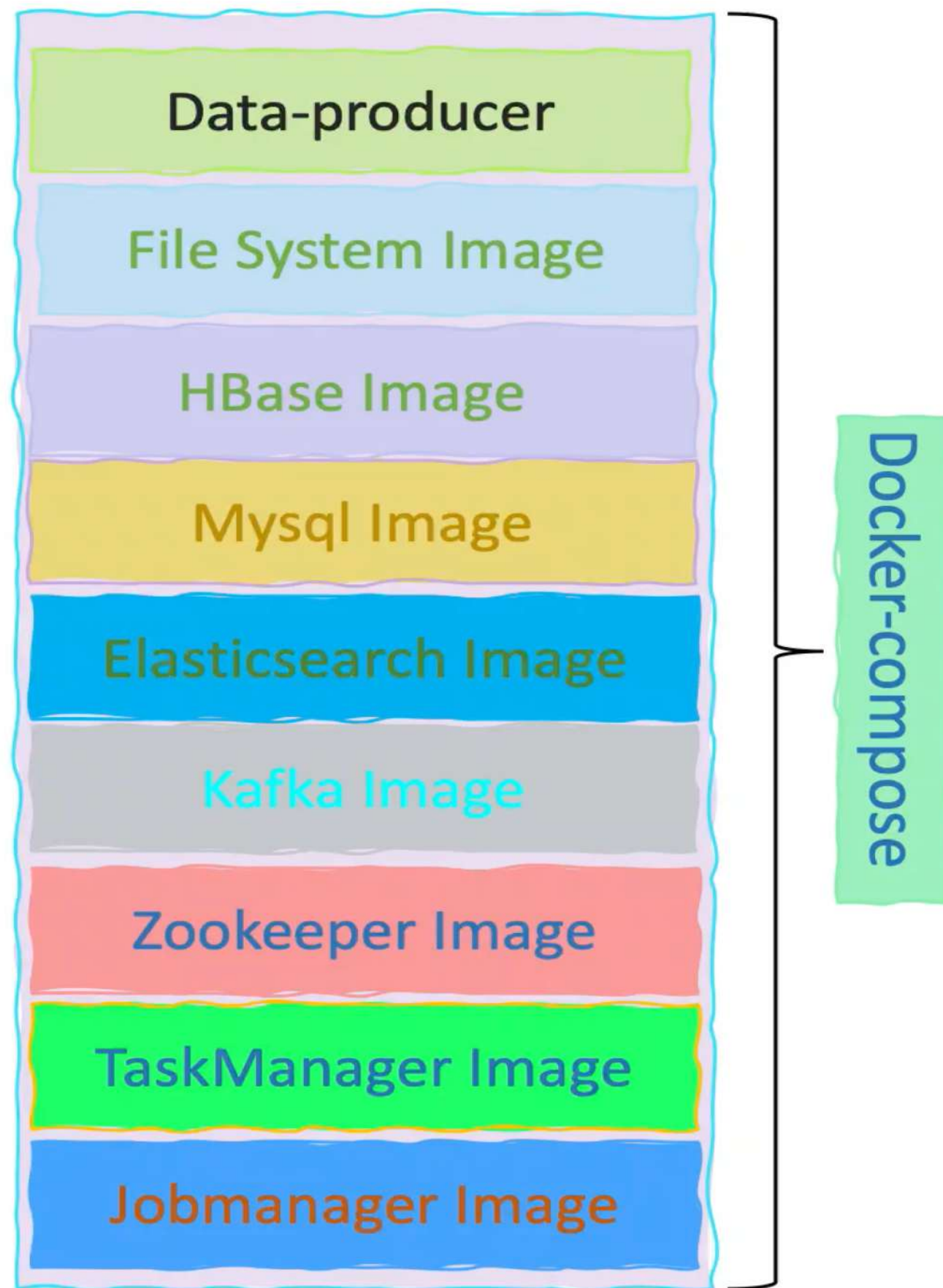
1.10.0		an hour ago
1.11-snapshot		7 days ago

[DockerHub](#)

[See all](#)

起点的选择？功能

Starting point of PyFlink? Start with Functionality



1. Install Docker : <https://www.docker.com/>
2. Config the [image of PyFlink](#) in Docker-compose.yml
3. Download and start-up: `docker-compose up -d`
4. Config volumes: `./examples:/opt/examples`
5. Write and Submit PyFlink Job ...
6. Check the Result
and multiple modifications and attempts
7. Read [PyFlink API Doc](#) for creative development

起点的选择？功能

```
▼ Table(object)
  m __init__(self, j_table)
  m __str__(self)
  f j_table
  m add_columns(self, fields)
  m add_or_replace_columns(self, fields)
  m alias(self, fields)
  m distinct(self)
  m drop_columns(self, fields)
  m fetch(self, fetch)
  m filter(self, predicate)
  m full_outer_join(self, right, join_predicate)
  m get_schema(self)
  m group_by(self, fields)
  m insert_into(self, table_path)
  m intersect(self, right)
  m intersect_all(self, right)
  m join(self, right, join_predicate=None)
  m join_lateral(self, table_function_call, join_predicate=None)
  m left_outer_join(self, right, join_predicate=None)
  m left_outer_join_lateral(self, table_function_call, join_predicate=None)
  m minus(self, right)
  m minus_all(self, right)
  m offset(self, offset)
  m order_by(self, fields)
  m over_window(self, *over_windows)
  m print_schema(self)
  m rename_columns(self, fields)
  m right_outer_join(self, right, join_predicate)
  m select(self, fields)
  m union(self, right)
  m union_all(self, right)
  m where(self, predicate)
  m window(self, window)
```

- select
- as
- filter
- where
- group_by
- distinct
- join(inner, left, right, full)
- join_lateral(inner, left)
- minus
- minus_all
- union
- union_all
- window(window)
- add_or_replace_columns(..)
- drop_columns(..)
- add_columns(..)
- rename_columns(..)
- ...

- OverWindow
- SessionWindow
- TumblingWindow
- SlidingWindow

Powerful
than SQL

博客



<https://enjoyment.cool/>

Stateful computations

over streams

**Unified programming
model**

Python & Java & Scala

推荐阅读

005Three Min Series - How to create UDF in PyFlink 1.10

004Three Min Series - How to using PyFlink Shell

Deep dive how to support Python UDF in Apache Flink 1.10

003Three Min Series - Three Min Series - How PyFlink does ETL

002Three Min Series - Run the Example of WordCount in PyFlink 1.9

001Three Min Series - Setting up the dev environment for Pyflink 1.9

Apache Flink 说道系列 - PyFlink 作业的多种部署模式

Apache Flink 说道系列 - PyFlink 1.11 功能规划讨论, 有你的声音吗?

Apache Flink 说道系列 - 如何在PyFlink 1.10中自定义Python UDF

Apache Flink 说道系列- Python API 中如何使用 Java 自定义 Connector

Apache Flink 说道系列- Python API 中如何使用 Kafka

Apache Flink 说道系列- 如何在IDE中运行使用Java UDFs 的Python 作业

Apache Flink 视频系列 - 2019.08.06直播(德国柏林)

Apache Flink 说道系列- Python Table API 开发环境搭建

Apache Flink 说道系列- Python API 时代的产物

github.com/pyflink/playgrounds

or jump to...

 **pyflink / playgrounds**

推荐阅读

Apache Flink 漫谈系列 - 序

Apache Flink 漫谈系列 - 概述

Apache Flink 漫谈系列 - Watermark

Apache Flink 漫谈系列 - State

Apache Flink 漫谈系列 - Fault Tolerance

Apache Flink 漫谈系列 - 流表对偶(duality)性

Apache Flink 漫谈系列 - 持续查询(Continuous Queries)

Apache Flink 漫谈系列 - SQL概览

Apache Flink 漫谈系列 - JOIN 算子

Apache Flink 漫谈系列 - JOIN LATERAL

Apache Flink 漫谈系列- Temporal Table JOIN

Apache Flink 漫谈系列 - Time Interval(Time-windowed) JOIN

Apache Flink 漫谈系列 - Table API 概述

Apache Flink 漫谈系列 - DataStream Connectors之Kafka



Apache Flink

Agenda

- ① Why PyFlink
- ② Status of PyFlink
- ③ Future of PyFlink
- ④ Flink on Zeppelin



Who is Jeff Zhang



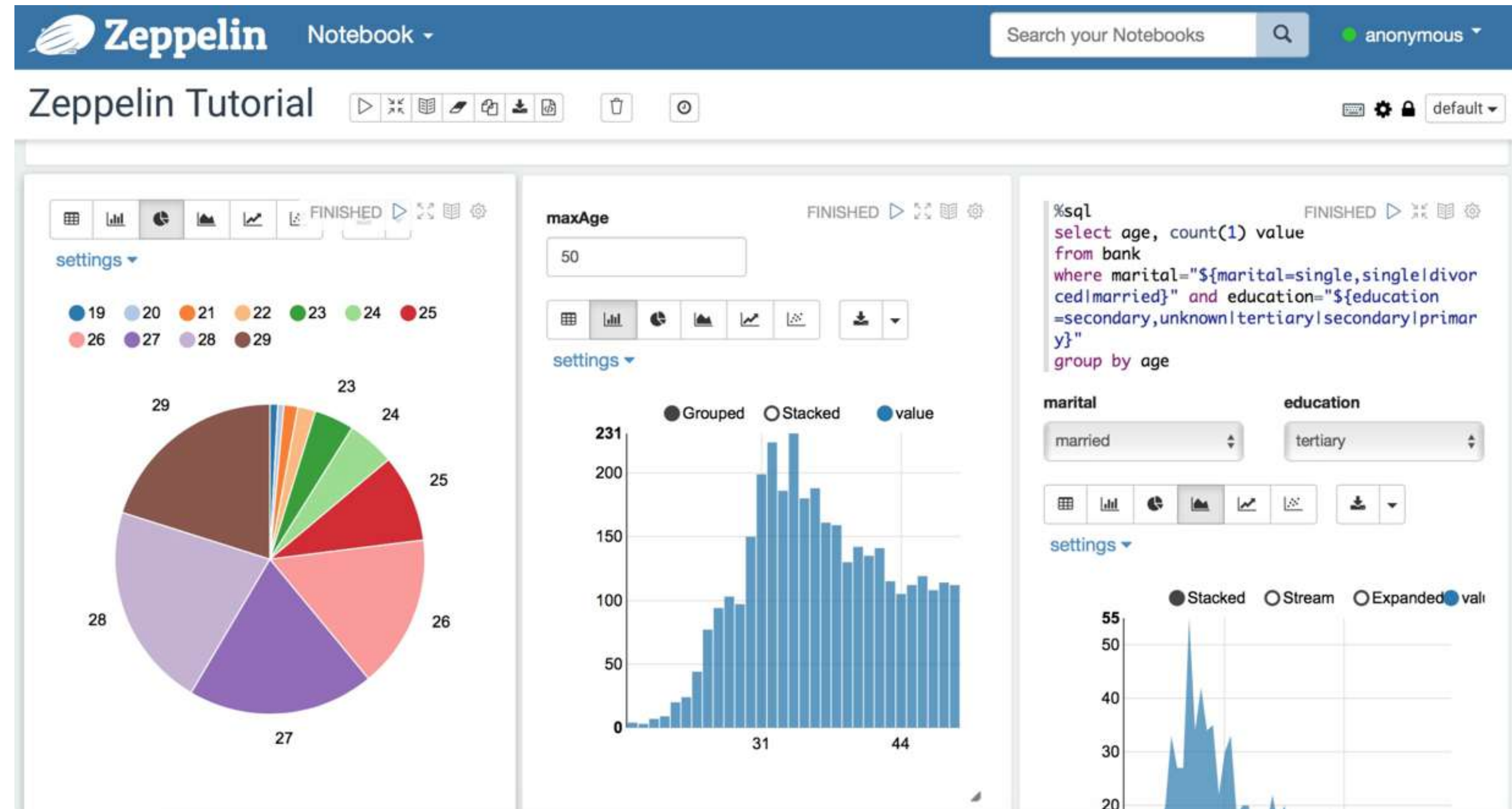
Data



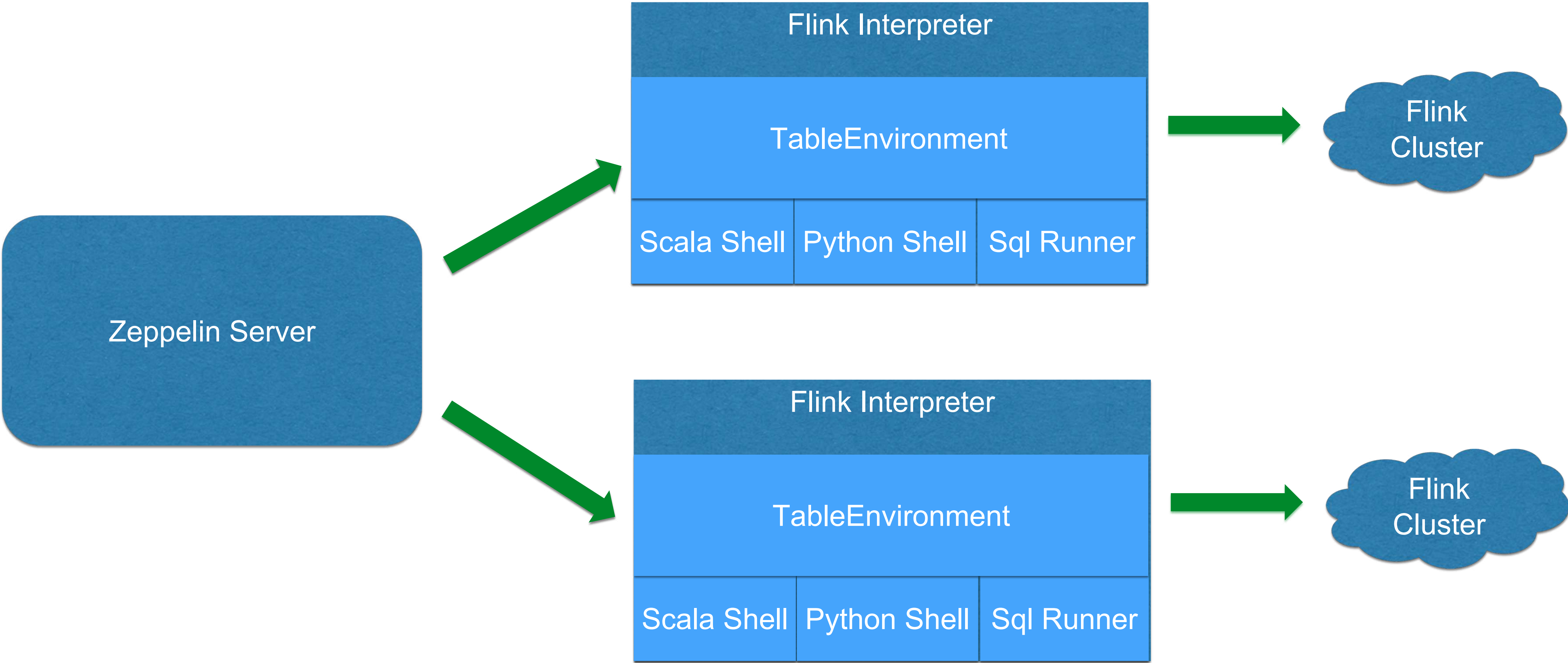
APACHE

What is Apache Zeppelin

Web-based notebook that enables data-driven, interactive data analytics and collaborative documents with SQL, Scala, Python and more.



Flink on Zeppelin (Architecture)



Flink on Zeppelin (Features)

Multiple Language Support
Scala/Python/Sql

Hive Integration

Dependency Management

SQL Result Visualization/Preview

Job Level Configuration

Interactive job control

Multiple thread job submission

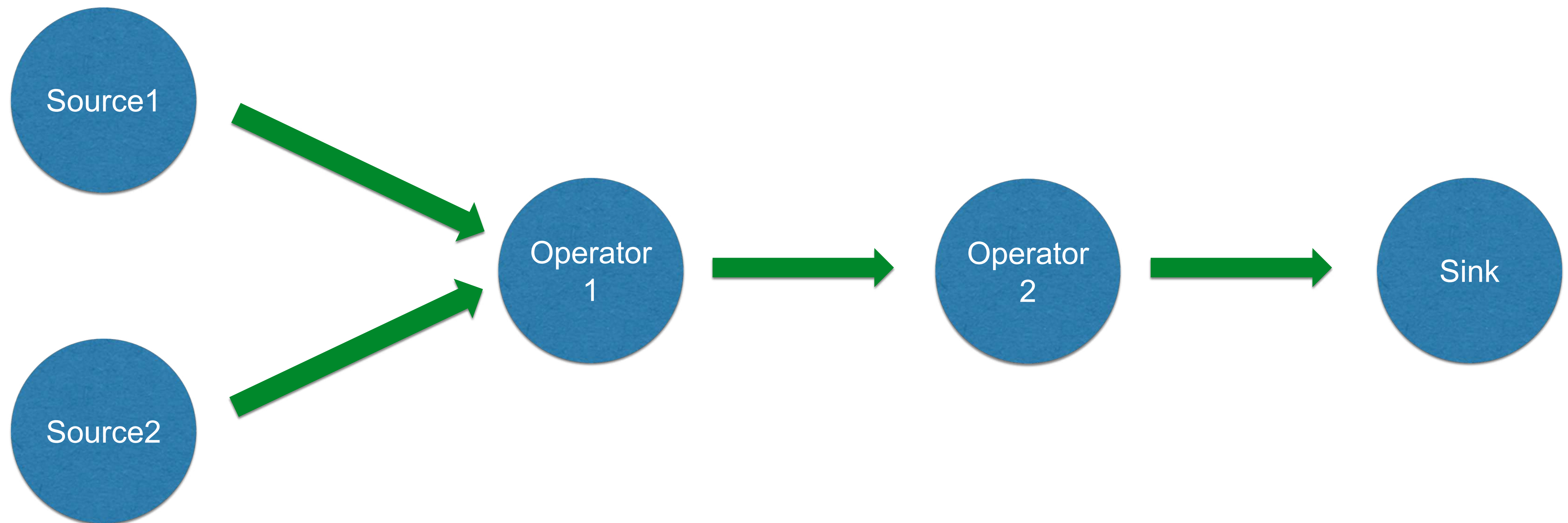
Machine Learning Support

3 essential elements of Flink job

Step 1. Define Source/Sink

Step 2. Build DAG

Step 3. Business Logic



Demo

Materials

- <http://zeppelin.apache.org/>
- Flink on Zeppelin tutorials
 1. Get started <https://link.medium.com/oppqD6dlg5>
 2. Batch <https://link.medium.com/3qumbwRIg5>
 3. Streaming <https://link.medium.com/RBHa2ITlg5>
 4. Advanced usage <https://link.medium.com/CAekyoXIg5>
- 钉钉群



