

vmware®

EXPLORE

2022

Greenplum Streaming Server 流式数据处理架构与实践

刘晖亮

VMware 研发经理

内容概要

- 流式数据简介
- Greenplum Streaming Server架构
- GPSS Kafka数据源实践
- 常见问题解答

vmware®
EXPLORE
2022

流式数据简介

vmware®
EXPLORE
2022

什么是流式数据?

- 概念
 - 持续生成的动态数据流
- 特征
 - 数据持续到达且到达速度快
 - 单条数据尺寸小 (通常为KB单位, 甚至更小)
 - 数据实效性强, 且可能无法预测数据边界
 - 低延迟处理需求
 - 多种一致性需求

流式数据与批式数据对比

	批数据	流数据
实效性	非实时、高延迟	实时、低延迟
数据大小	大批量数据	单条记录或包含几条记录的微批量数据
分析	用来做复杂分析	简单的响应函数、聚合处理

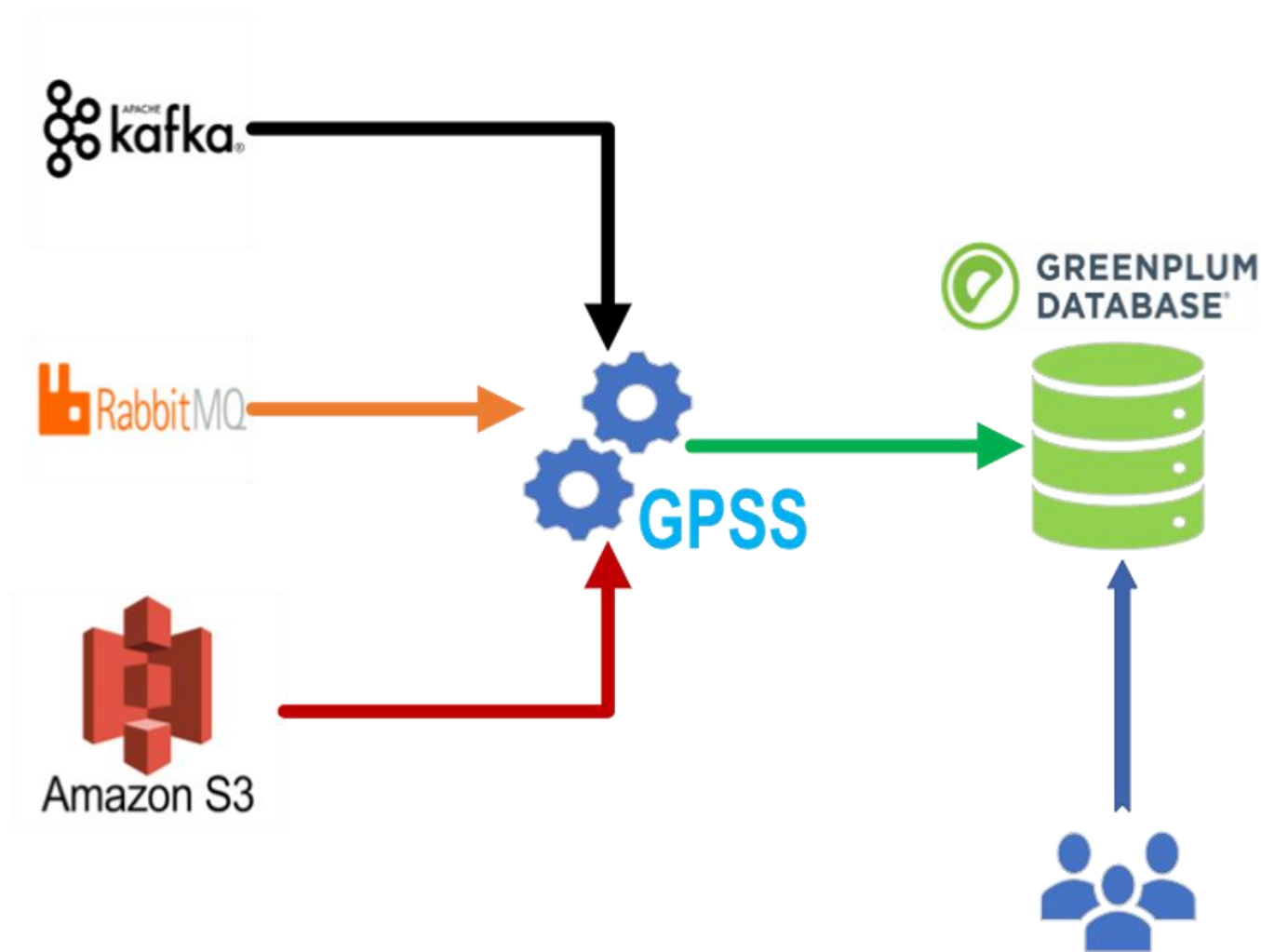
- 流批一体
 - 同一个业务，使用同一个SQL逻辑来实现大数据的流计算和批计算
 - Greenplum Streaming Server

Greenplum Streaming **vmware** Server 架构

EXPLORE
2022

GPSS的功能和特点

- 加载流式数据到Greenplum
- 保证数据的强一致性
- 自动聚合和转换数据
- 支持丰富的数据格式

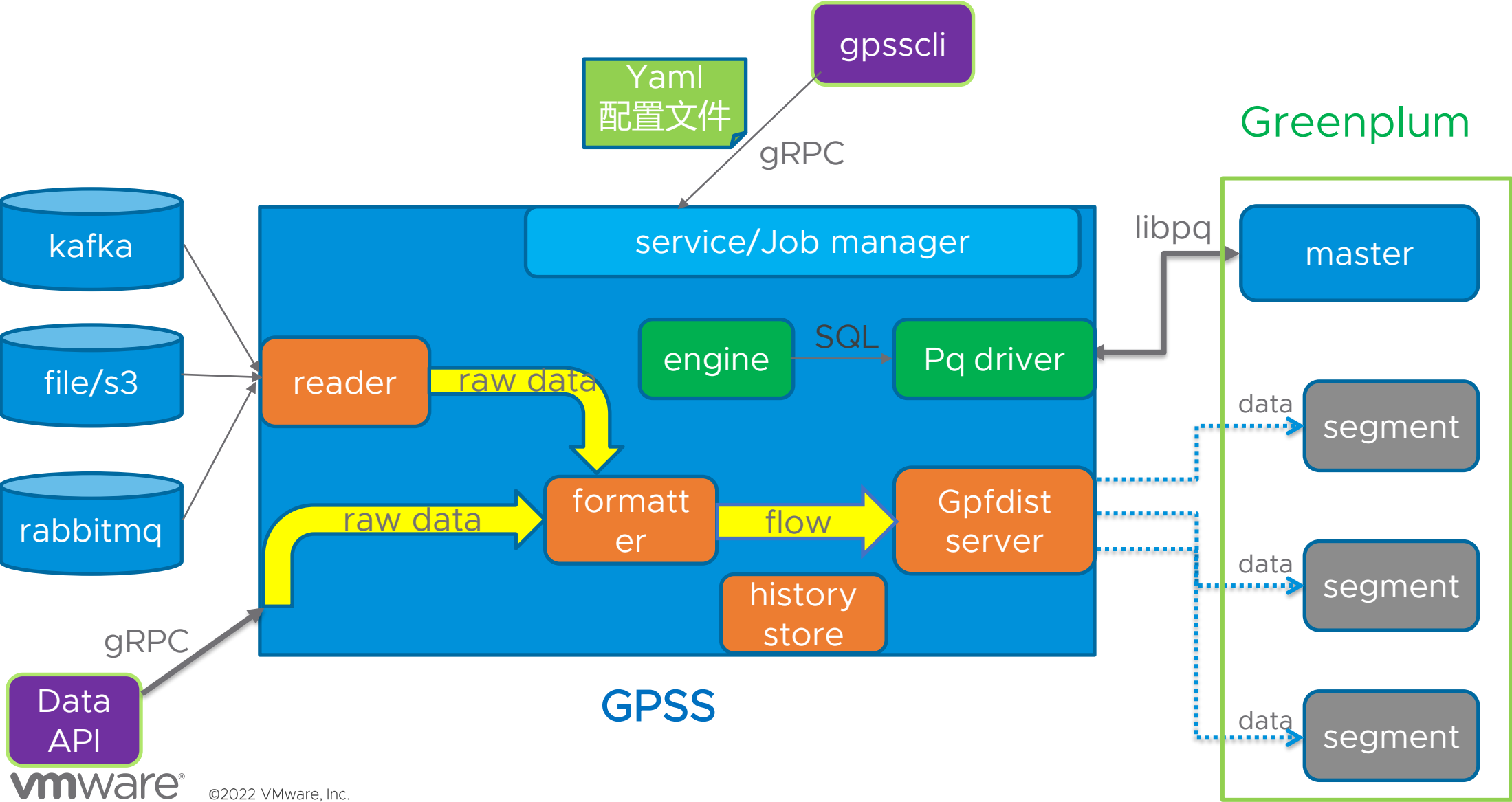


GPSS的组成模块

- GPSS
- GPSSCLI
- Extension



GPSS的架构



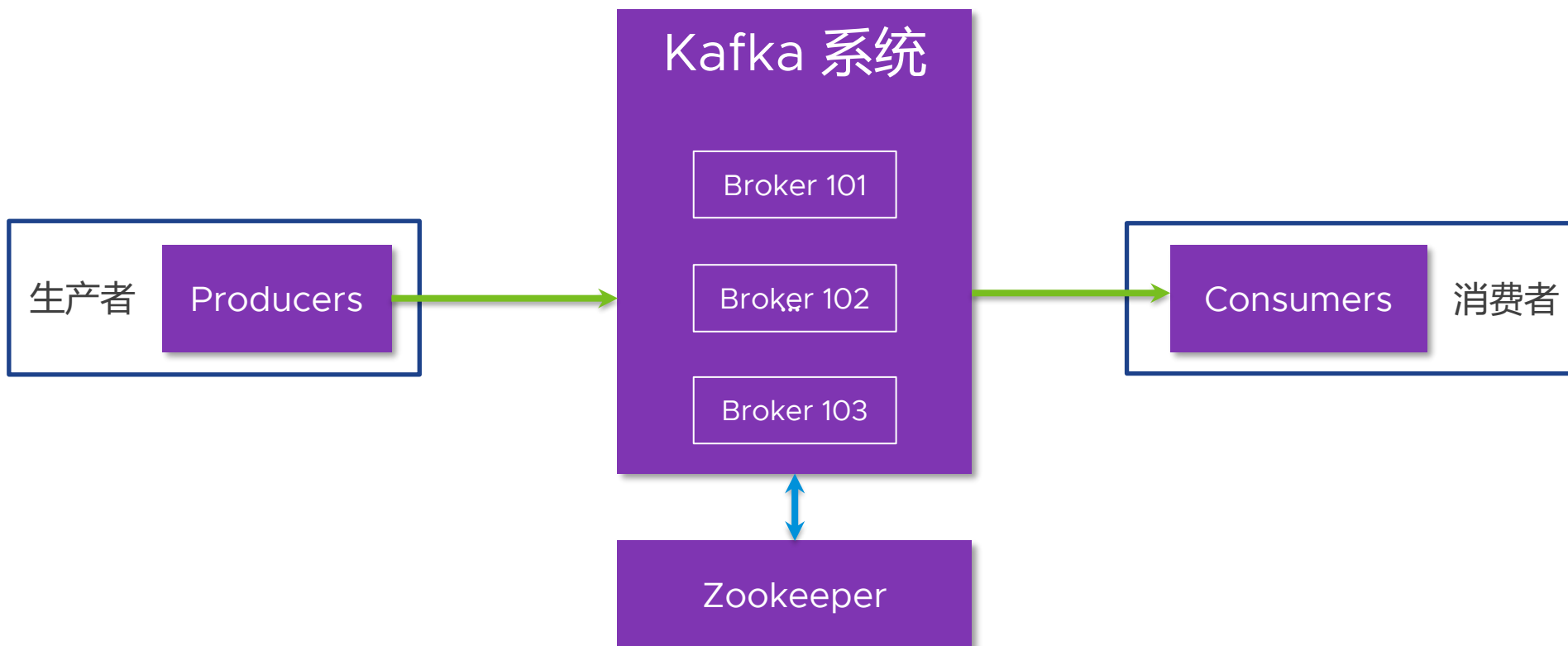
GPSS Kafka 数据源 实践

vmware®
EXPLORE
2022

Kafka系统的组成

Broker: Kafka实例，分布式部署

Zookeeper: 保存集群的元信息并管理broker



Kafka基本概念

Topic

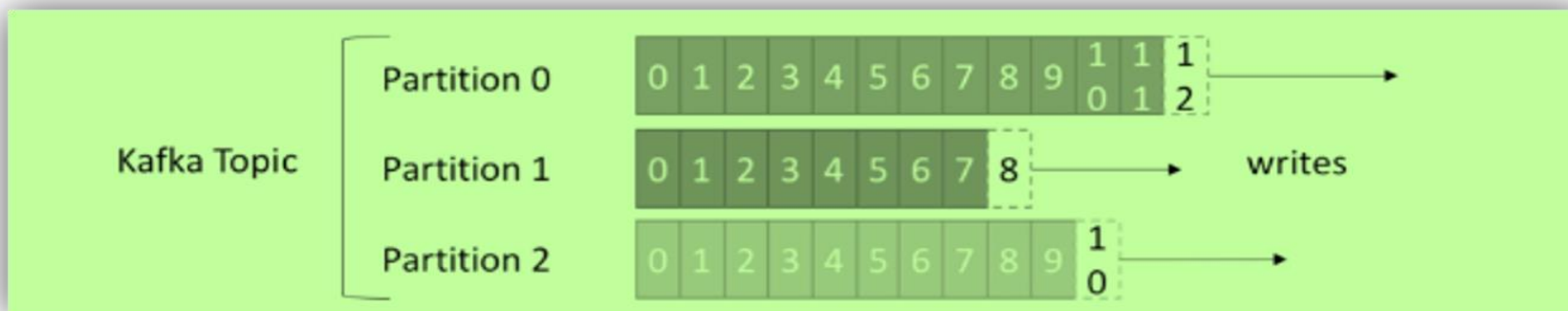
- 消息的主题，可以理解为消息的分类，类似于数据库的表

Partition

- Topic的分区，每个topic可以有多个分区，分区中包含有序的消息

Offset

- 每个分区中的消息都有一个递增的唯一ID，称为offset



如何启动GPSS

➤ 准备工作

- 安装GPSS扩展
 - create extension gpss;
 - create extension dataflow;
- GPSS 配置文件 
 - JSON格式: config.json

➤ 启动GPSS

gpss -c config.json

```
1  {
2      "ListenAddress": {
3          "Host": "",
4          "Port": 5000
5      },
6      "Gpfdist": {
7          "Host": "localhost",
8          "Port": 8080,
9          "BindAddress": "0.0.0.0"
10     },
11     "Authentication" : {
12         "Username": "cliAuth",
13         "Password": "SHADOW:CEEWB34BPTXZRVSKIXI7XAMUA"
14     }
15 }
```

GPSSCLI 命令

➤ 格式

- gpsscli [子命令] --gpss-host <HOST> -gpss-port <PORT> [配置文件]
- 常用子命令
 - submit: 提交一个job
 - start/load: 启动一个job
 - stop: 停止一个job
 - remove: 删除一个job
 - list/status: 列出现有job及状态
 - progress: 打印job加载进度信息
 - wait: 等待一个job结束

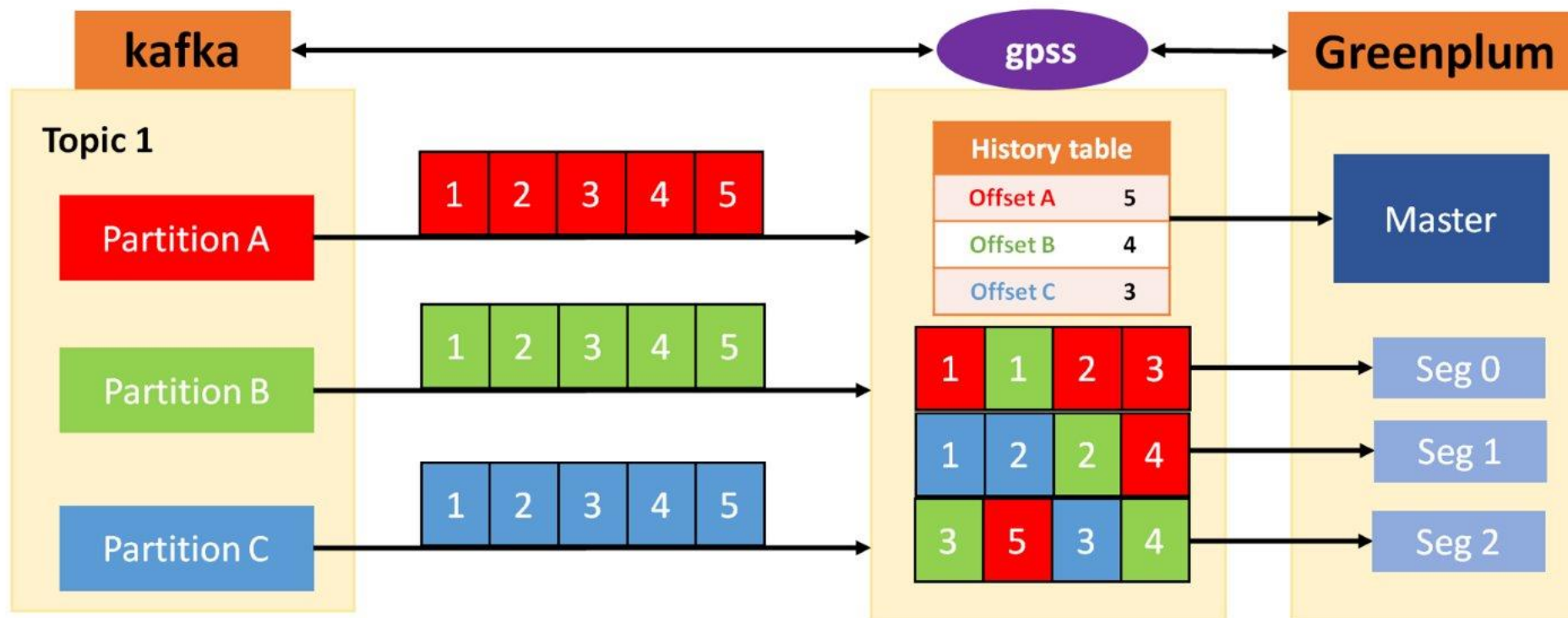
GPSSCLI 配置文件

```
1 DATABASE: ops
2 USER: gpadmin
3 PASSWORD: changeme
4 HOST: mdw-1
5 PORT: 5432
6 VERSION: 2
7 KAFKA:
8     INPUT:
9         SOURCE:
10             BROKERS: kbrokerhost1:9092
11             TOPIC: customer_expenses2
12             PARTITIONS: (1, 2...4, 7)
13         VALUE:
14             COLUMNS:
15                 - NAME: c1
16                 TYPE: json
17             FORMAT: avro
18             AVRO_OPTION:
19                 SCHEMA_REGISTRY_ADDR: http://localhost:8081
20             FILTER: (c1->>'month')::int = 11
21             ERROR_LIMIT: 25
```

```
22 OUTPUT:
23     SCHEMA: payables
24     TABLE: expenses2
25     MAPPING:
26         - NAME: customer_id
27           EXPRESSION: (c1->>'cust_id')::int
28         - NAME: newcust
29           EXPRESSION: ((c1->>'cust_id')::int > 500)::boolean
30         - NAME: expenses
31           EXPRESSION: (c1->>'expenses')::decimal
32         - NAME: tax_due
33           EXPRESSION: (c1->>'tax')::decimal
34     METADATA:
35         SCHEMA: gp kafka_internal
36     COMMIT:
37         MINIMAL_INTERVAL: 2000
```


加载Kafka数据到Greenplum

1. 启动一个kafka集群并创建相关topic
2. 通过gpsscli提交gpss kafka job

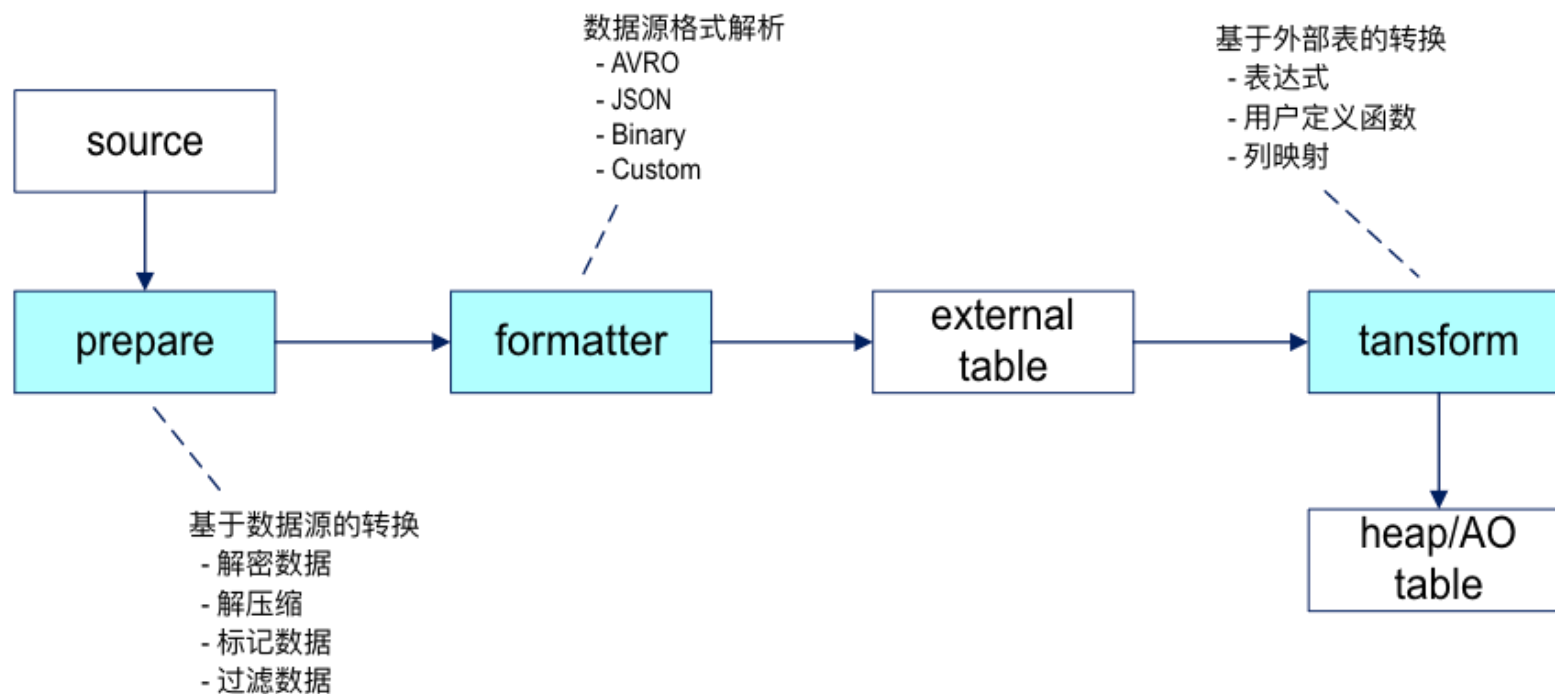


GPSS支持的数据格式

格式	描述
avro	Avro是一种与语言无关的二进制数据序列化格式 GPSS 支持 libz, lzma 或snappy 压缩的avro数据, 并根据schema registry 解析成单个JSON列
binary	GPSS可以从Kafka读取二进制数据并转换成单个bytea 类型的列
csv	逗号分隔值, 一种常见的数据格式.
custom	自定义格式, 可以被指定的formatter解析
delimited	一种可配置分隔符的文本格式
json	一种轻量级的数据交换格式, 采用key:value的方式记录数据

GPSS支持多阶段数据转换

- 基于数据源的转换
- 数据格式解析
- 基于外部表的转换



常见问题解答

vmware®

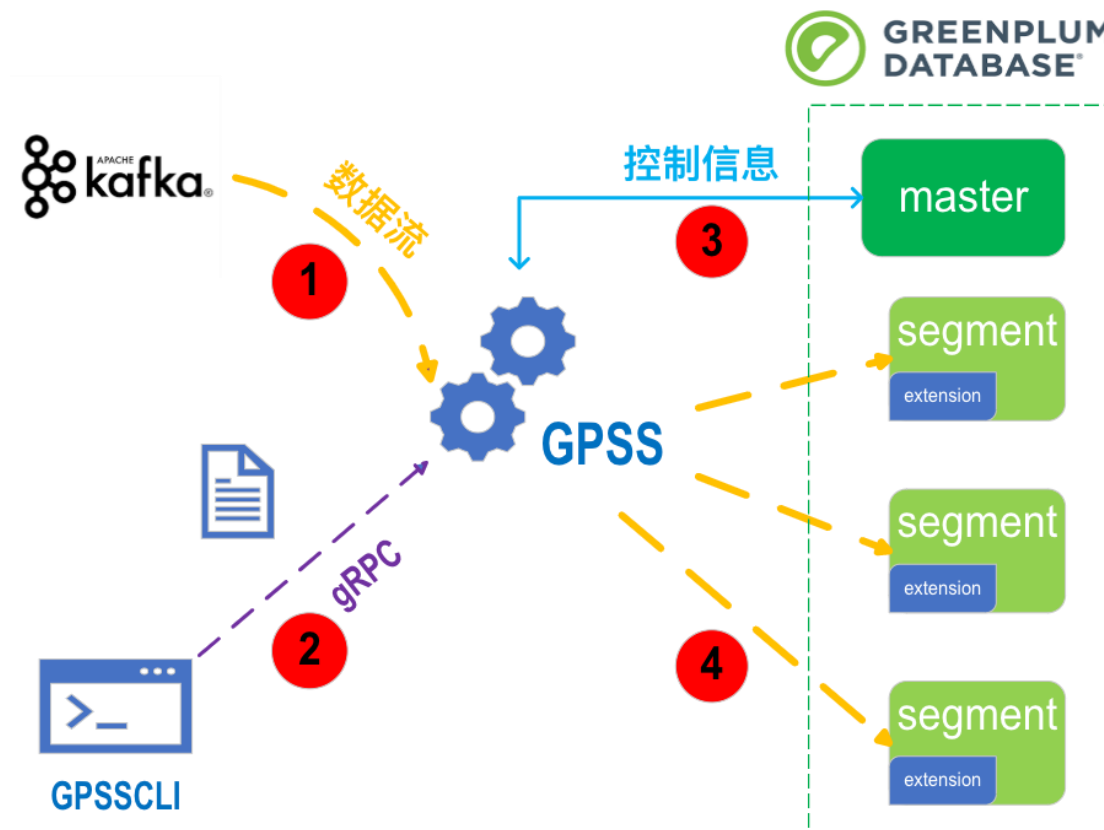
EXPLORE

2022

GPSS支持SSL/TLS加密传输吗？

1. Kafka <- -> GPSS
2. GPSSCLI <- -> GPSS
3. GPSS <- -> Greenplum master
4. GPSS <- -> Greenplum segment

以上都支持SSL/TLS加密传输



GPSS有哪些最佳实践方法？

- GPSS 有两个参数用来控制何时提交数据： MINIMAL_INTERVAL 和 MAX_ROW
- 如果数据流量不太高，可以适当调大这两个参数来增加每次提交的数量，提升效率
- 如果数据流量特别大，每次积攒数据对内存有一定压力，可以减小MAX_ROW来降低对资源的使用
- GPSS Kafka数据源的partitions 参数可以用来指定只加载某个topic的部分分区。这样可以通过起多个GPSS来降低由于分区过多带来的单点瓶颈问题。
- GPSS支持将同一个数据源的数据加载到多个目标表。
- 了解Greenplum和GPSS更多功能，欢迎关注公众号：



谢谢!

vmware®
EXPLORE
2022