

从数据到洞察：利用DataOps加速AI模型开发

代立冬

白鲸开源科技联合创始人 & CTO

关于我



白鲸开源科技联合创始人

Apache 孵化器导师

Apache DolphinScheduler PMC Chair

Apache SeaTunnel PMC

ApacheCon 亚洲大数据湖仓论坛出品人

中国科协 “2023开源创新榜” 优秀人物

目录 CONTENTS

01

AI 大模型开发流程

02

DataOps 核心理念
白鲸开源 DataOps 实践

03

实践案例分析

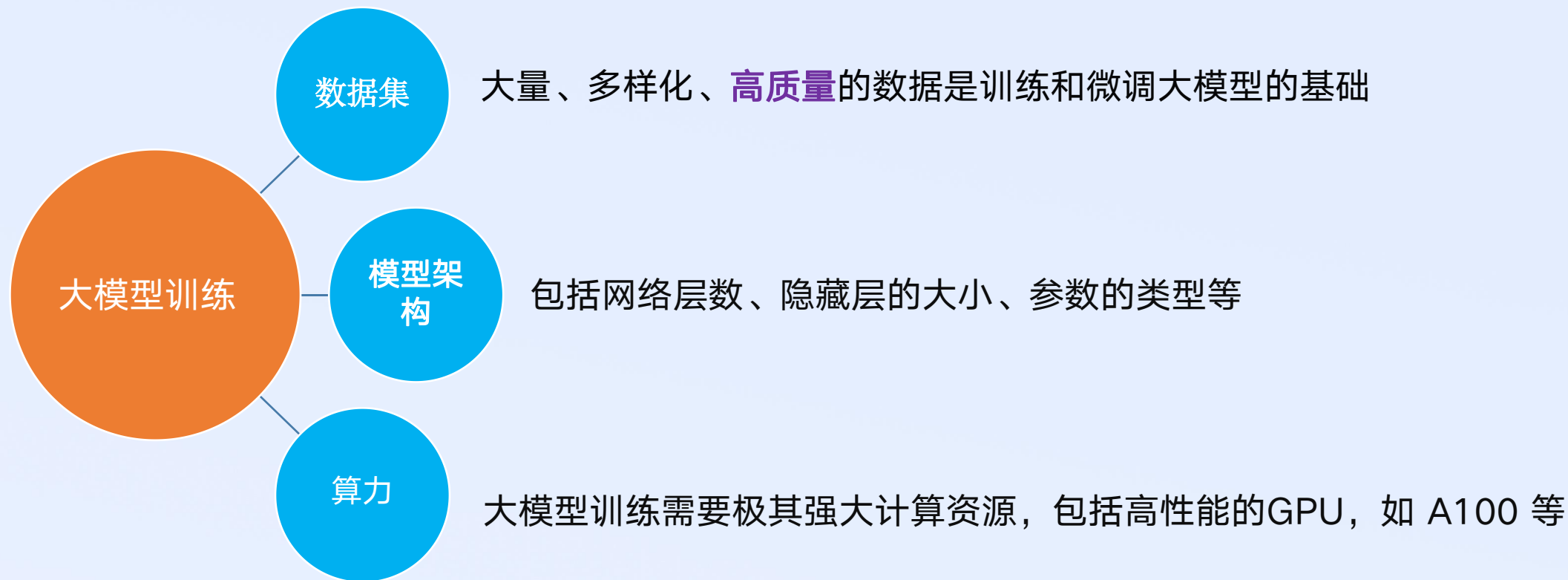
04

DataOps 未来

PART · 01

DataOps 核心理念

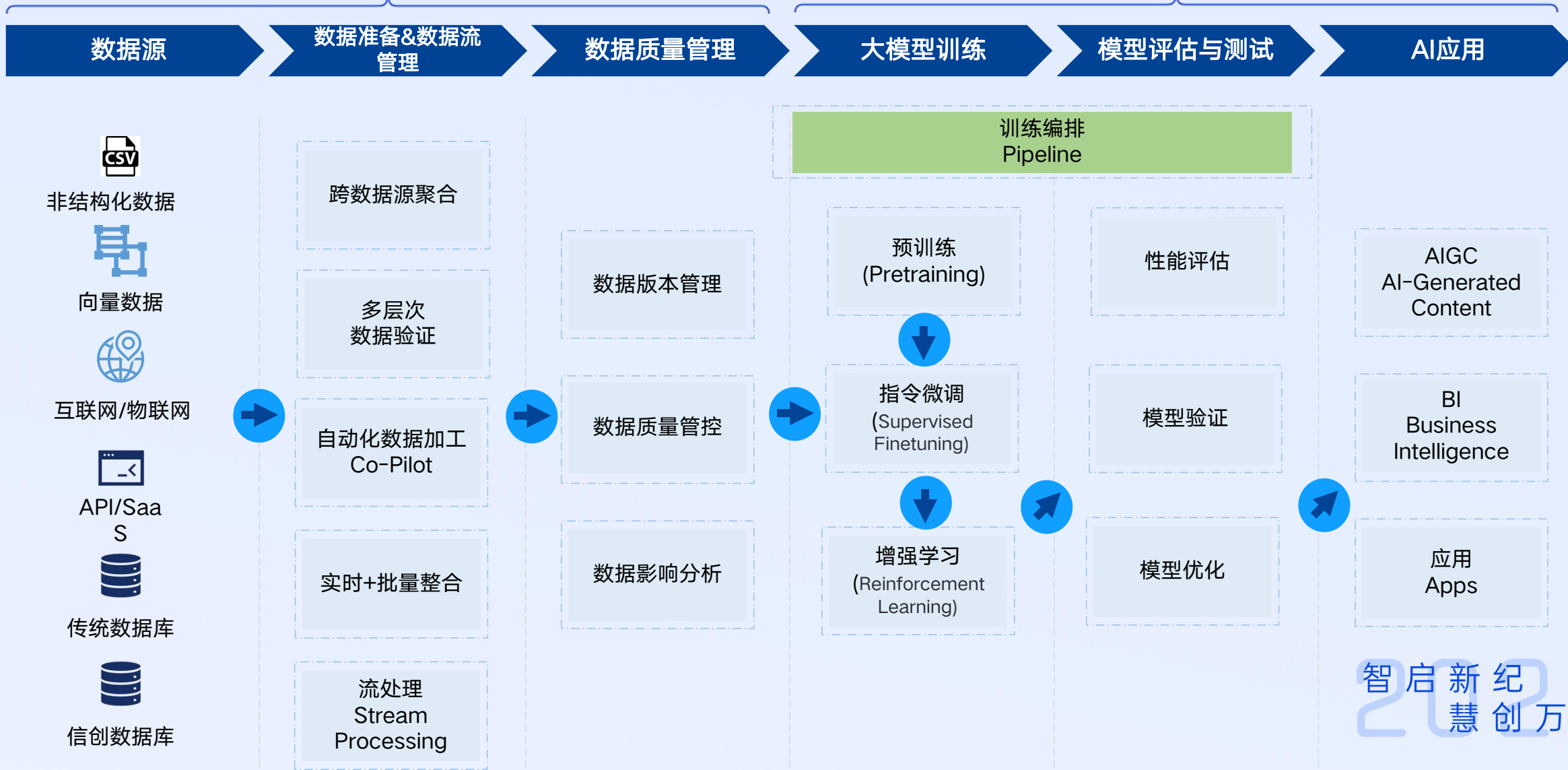




DataOps与AI模型开发的融合，将加速AI模型的开发周期，提升模型的准确性和效率。

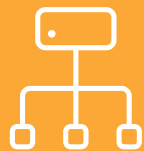
数据工程阶段

算法工程阶段





- 企业内拥有多组“数据平台”，数据资源和流程分散在各部门，难以掌控



- 大数据、流数据、AI数据加工缺乏工具管控形成了企业新的“蜘蛛网”



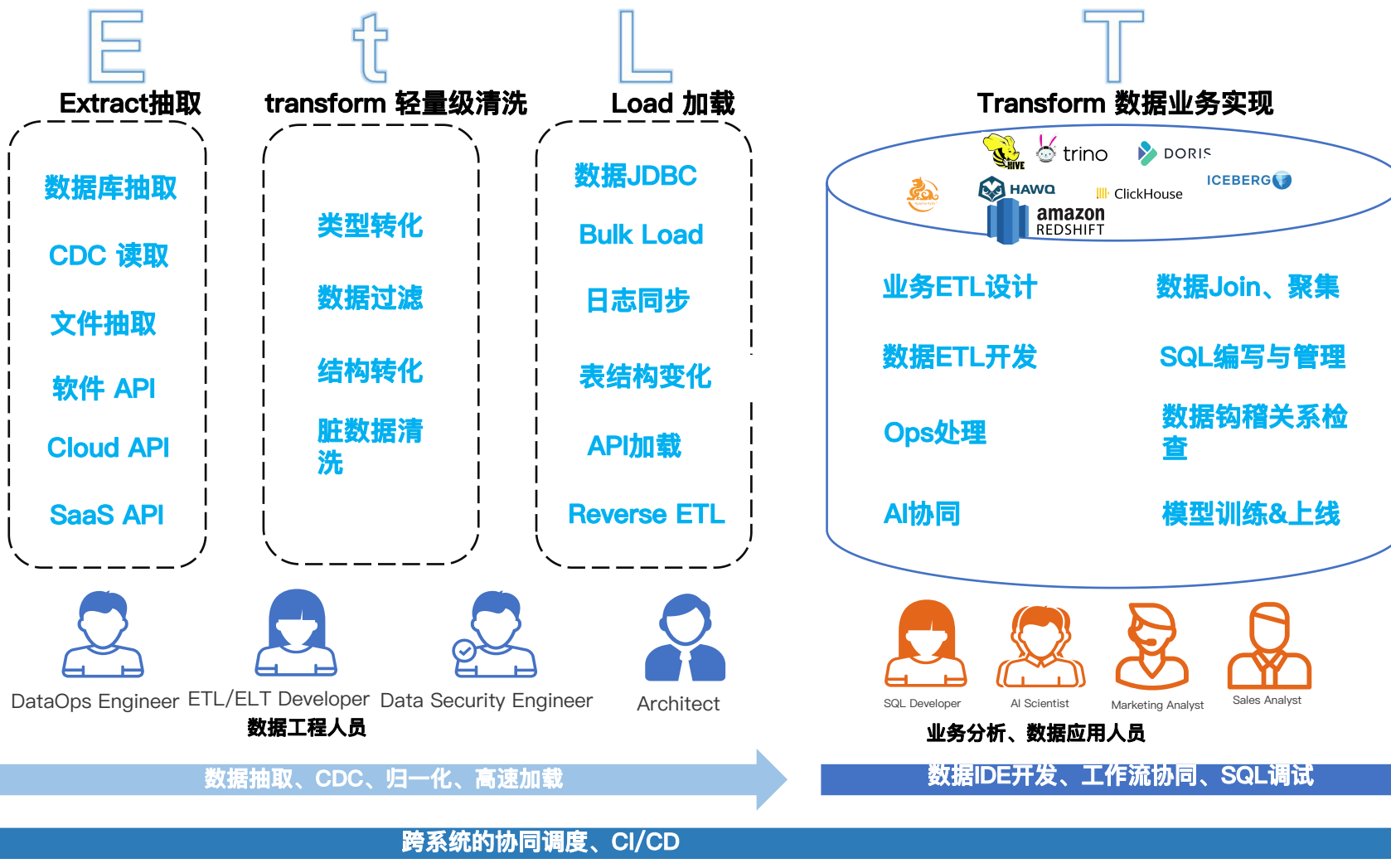
- 企业大数据开发处于“野蛮生长状态”，整体研发管理距离应用开发DevOps流程相差甚远



- 多种新兴数据引擎、云原生、新数据架构的变化缺乏管控，数据血缘、同步、调度与数据发展严重落后

新技术环境下 EtLT 架构出现

云、SaaS、本地混合数据源让传统的数据处理流程从ETL、ELT变为能更加快速满足业务需求的 EtLT 架构，EtLT 能更加敏捷应对离线/实时数据湖、数据仓库、AI 模型训练当中的复杂多变的数据需求场景。



白鲸开源是一家开源原生的 DataOps 商业公司，由多个 Apache Foundation Member 成立，80% 员工都是 Apache Committer，主导 2 个 Apache 顶级开源项目 (Apache DolphinScheduler, Apache SeaTunnel)。

同时根据全球最佳实践发布商业版版本 -- WhaleStudio，帮助企业解决内部多数据源、多数据系统复杂的数据集成，数据开发、工作流编排运维及部署、数据质量管控、团队敏捷协作等一些列问题，并在 6000 多家企业中得到实践和使用。

DataOps 的优势



增强数据质量

DataOps 强调数据质量的重要性，通过自主数据校验，提高数据的准确性，确保模型微调所需的优质数据

提升开发效率

DataOps 通过自动化和优化数据处理流程，减少人工干预，提高数据处理和 AI 模型的开发效率。

促进团队协作

DataOps 提倡团队间协作和沟通，通过共享工具和平台，有助于构建高效、敏捷的数据 & AI 模型开发团队。

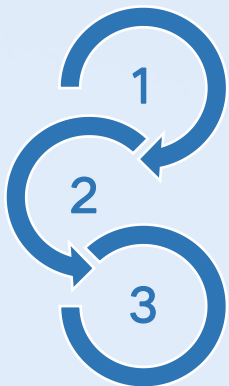
DataOps 关键实践

台

之任务调度平



大数据
工作流
痛点

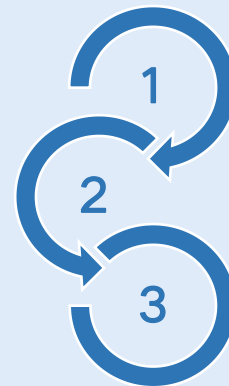


多个任务单元
存在时间顺序
存在依赖关系

+

执行频率高
数据量、任务量大
云原生

企业级
场景
痛点



脚本 + 大数据组件



代码复用性低



集群化部署与拓展
过于复杂



开源数据组件
更新升级频率高

上一代调度引擎



非可视化设计
智能化程度低



扩展能力弱
受限于单点瓶颈



多云异构数据
管理能力弱

新一代调度引擎



Star: 12.1k



可视化任务调度
支持多任务类型



去中心分布式设计
高稳定可用



百万数据量级
任务稳定运行

5 +
年社区运营

520 +
开发者贡献

6000 +
用户使用



联通数科早期使用原商业调度系统支撑着全域数据平台加工与调度，以接口机配合Shell（HiveSQL）为主的开发编排运维模式，处理日均数万的流程实例和日均**上百万的 Job** 作业，对比闭源调度工具、Airflow、Azkaban后，最终选择DS。

- 满足业务需求和调度功能要求
- 满足大数据量要求
- 用户使用成本低



SHEIN

早期使用Airflow调度全球任务，但因为分布式支持问题、无可视化问题导致系统开发效率和稳定性堪忧，同时也无法支持K8S和全球的云原生部署。选择从Airflow迁移至DS。

- 全球云部署、K8S支持
- 分布式去中心化以保证稳定
- 全量替换Airflow
- 解决全球大数据调度易用问题，赋能分析团队快速开发调度任务



荔枝 FM

过去大数据调度使用SQL/Shell/Python脚本和其他大数据组件完成整个AI流程，面临不易用且难复用问题。使用基于DS的AI开发平台后，荔枝FM将获取数据、数据预处理、模型训练、模型预测、模型评估和模型发布过程抽象成组件，用 DAG 串联，使用拖拽和配置的方式实现低代码开发。

- 实现对海量数据存算
- 可以复用ML流程
- DAG 执行引擎



高性能、大批量数据调度



全球云部署、易用数据开发



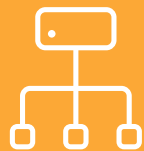
AI/ML Ops调度融合

DataOps关键实践 之数据集成工具





- 数据源多达几百种，版本间不兼容，而且不断有新的出现



- 数据丢失与重复，无法一致性
- 出现问题无法回滚或者断点继续执行
- 同步过程不透明，缺少监控



- 频繁读取 binlog 对数据源端影响大
- 大事务、Schema 变更影响下游
- 低吞吐高时延导致数据无法及时到达



- 离线同步和实时同步常被分开管理，维护困难
- 数据割接人工进行

上百种源数据库/地点

数据同步与集成

目标数据库/地点

MySQL PostgreSQL Kafka

MongoDB Elastic Hive

Druid Redis AWS Aurora

Hudi Kudu HBase

InfluxDB Neo4j Feishu

原有解决方案

Spark Flink

Sea Tunnel

SeaTunnel Universal API

Table API
Source API Sink API
Engine API

SeaTunnel Engine

- 批量数据全量、增量集成
- 实时数据集成
- 批量无主键增量集成等

160+ 3x 500+
连接器数量 版本迭代 接口数增长

其他解决方案

Airbyte DataX

比 Airbyte性能快30倍 单机比DataX性能快2.6倍

MySQL PostgreSQL Kafka

MongoDB Elastic TiDB

Druid Redis Hive

Iceberg Kudu HBase

AWS Redshift SelectDB StarRocks

AWS S3 ClickHouse Snowflake



- 简单易用，开箱即用，不依赖HDFS, Flink, Spark
- 全可视化操作



- 丰富的数据源支持，目前已经支持 100+ 种数据源



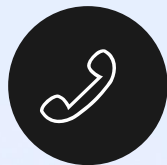
- 批流、实时、CDC一体化配置
- 无主键增量数据集成
- 整库同步、表结构自动变更



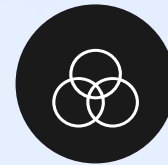
- 存算分离架构设计
- 高性能数据同步
- 支持节点动态伸缩



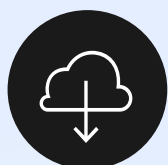
- 全量到增量无锁化自动切换
- 读缓冲(一个源到多个目标数据源, 只用一次读取)
- 动态速率控制, 对源端和目标端压力可控



- 支持 Schema Evolution
- 断点续传
- 实现 Exactly-Once 一次语义, 保证数据一致性



- 云组件支持
 - K8s支持
 - AWS Redshift、S3, RDS, DynamoDB
 - 阿里 OSS File, Max, TableStore等



- 无中心化设计确保系统的高可用, 支持多云
- 支持每日千亿级数据量同步

智启新纪
慧创万物

其它引擎

01

容错粒度大

当进行多表同步时，任何一张表出现问题，整个作业都会失败停止，导致所有表的同步都会延迟

02

资源浪费严重

每个作业只能同步一张表，当我们有多张表需要同步时，需要启动多个Job来运行，这种方式非常的浪费资源

03

JDBC连接数过多

每个task只能处理一张表，所有每张表至少需要一个JDBC连接来读取或写入数据。当进行多表同步和整库同步时，需要大量的JDBC连接

04

不支持数据缓存

CDC场景下，有些企业设置的数据库日志留存时间比较短，作业失败时，数据库日志处理也会停止，会出现源端数据库日志被清除的情况

05

不支持表结构变更

在CDC场景下，DDL变更的检测和下游应用非常重要，目前Datax/Spark/Flink/FlinkCDC都无法支持

SeaTunnel 同步引擎- Zeta

01

易用省资源

WhaleTunnel 利用Zeta引擎Dynamic Thread Sharing技术，提高CPU利用率，不依赖HDFS，Spark等复杂组件，更好的单机处理性能

02

可视化开发与运维

WhaleTunnel 使用智能辅助的可视化UI技术，帮助用户快速建立单表多字段、多表多字段、SaaS和非结构化到数据库等复杂的数据集成任务并监控

03

多方式确保一致性

支持无中心HA和更细粒度的作业回滚机制，结合多阶段提交与CheckPoint机制，确保数据一致的同时避免大量回滚导致性能下降

04

极致CDC&批量性能

WhaleTunnel 实现 zero-copy 技术，不需要序列化开销，列式内存格式增加大吞吐，可以最大化利用网络带宽，支持Oracle/SQL Server/Mysql 等

05

多复用/数据库日志多表解析

WhaleTunnel支持多表或整库同步，解决JDBC连接过多的问题，支持多表或整库数据库日志读取解析，解决CDC多表同步场景下需要重复解析日志的问题。

智启新纪

慧创万物

跨云数据准备

JPMorgan & Chase

美国最大商业银行

解决多云异构环境下，需要异构数据打通，将AWS Aruora, DynamoDB, SFTP数据实时同步到ES, S3, Snowflake下

异构数据实时数据同步

bilibili

超大型客户

解决多数据源数据每日出入数据库以及每日出入仓同步数据问题，数据集群规模几十台，日均记录数上千亿，日均数据量在100T以上。

滴滴

腾讯云

去哪儿旅行
总有你想要的低价

oppo

vip.com 唯品会

多点DMALL

bilibili

Shopee

YOOZOO
游族网络FOXCONN®
富士康科技集团新浪微博
weibo.com

coupang

ByteDance
字节跳动

ZUIYOU

松果出行

BOTON 宝通

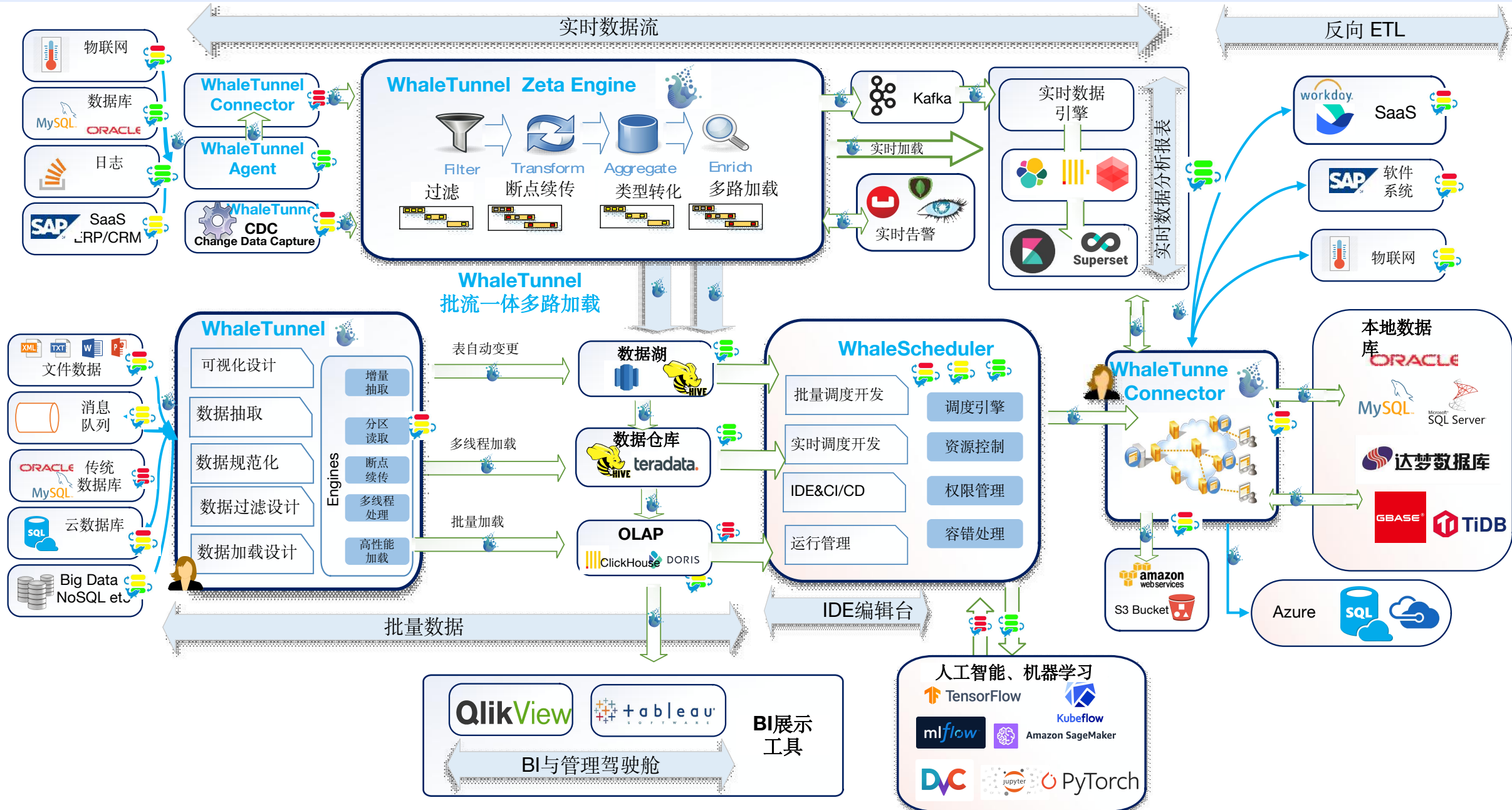
中国移动
China Mobile虎牙直播
huya.com国家电网
STATE GRID
国网重庆市电力公司
STATE GRID CHONGQING ELECTRIC POWER COMPANY震旦集團
AURORA GROUP

GRIDSUM 国双

PART · 02

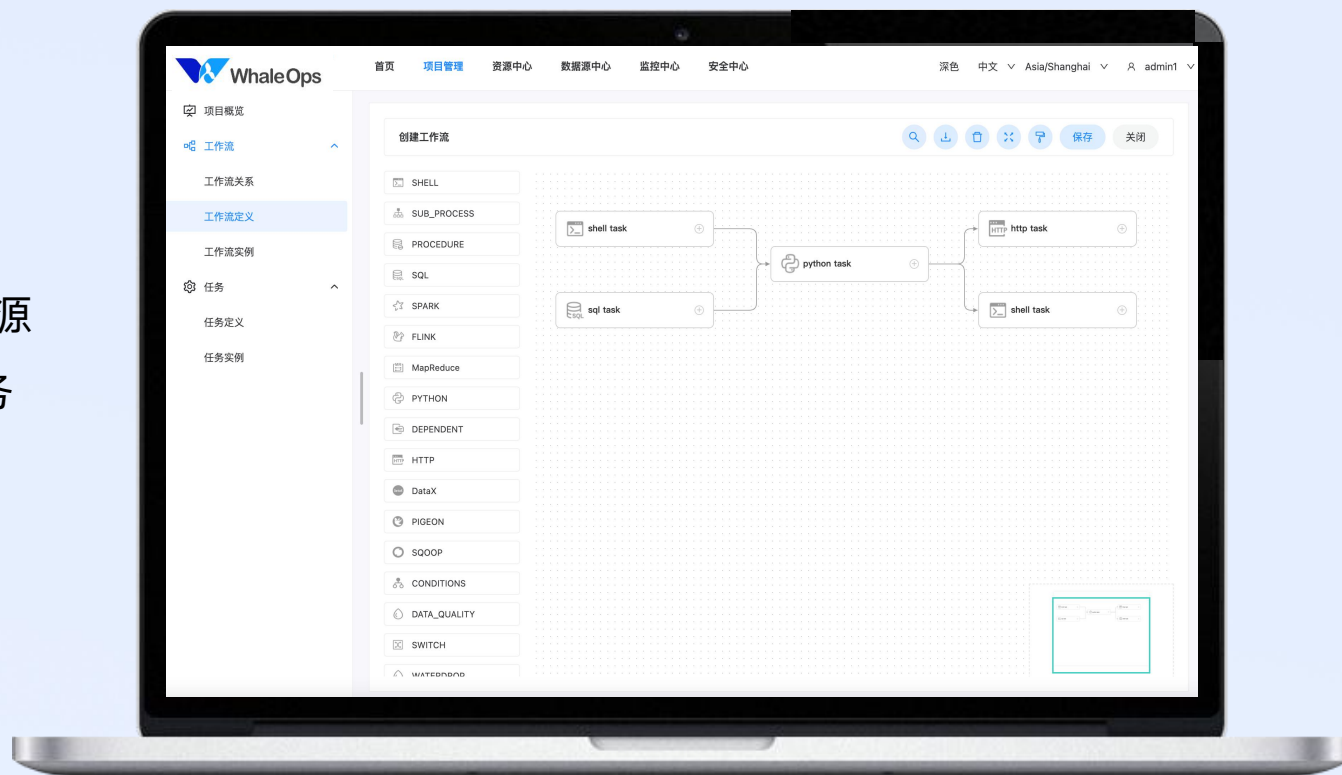
白鲸开源 DataOps 实践





基于 Apache DolphinScheduler 和 Apache SeaTunnel 的产品 WhaleStudio，是分布式、云原生并带有强大可视化界面的 DataOps 系统，增加了商业客户所需的企业级特性：

- 完全自主研发，上下游生态圈广阔，支持 160+ 种数据源
- 全面支持云原生—云、仓、湖 实时/离线批流一体化任务管控
- 低代码实现企业大数据的操作系统和高速公路
- 完善的 DataOps 流程可无缝集成代码工具
- 丰富的数据源对接和传统 ETL 数据组件支持
- 一站式完成从开发-》测试-》上线-》运维闭环



全部项目 10

项目概览

工作流

调度任务

离线任务实例

实时任务实例

加冕列表

隔离列表

离线同步

离线任务定义

离线任务实例

实时同步

实时任务定义

实时任务实例

Test-2430854580 暂停

恢复

同步概览 运行日志 运行记录 告警列表 运行中断时间: -

自动刷新 5s

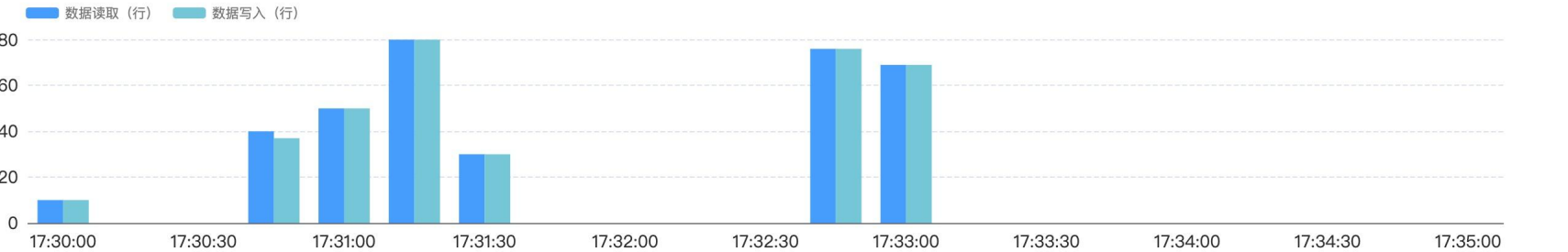
已读取数据量 (行) **3,560,000** Row MB 已写入数据量 (行) **3,560,000** Row MB

Insert: n Update: n Delete: n DDL: n

当前数据延迟 (秒) **1.03** 10% ①
Maximum: 20.33 Average: 2

全量阶段同步速率
读取速率(行/秒): 33214 处理速率(行/秒): 33145
读取速率(MB/秒): 10 处理速率(MB/秒): 10

同步数据 最近 5min 1hour 1day 1week



实时增量 搜索表名称

源表	目标表	Insert	Update	Delete	DML	DDL	上次同步时间
testwyr1	TESTWYR1	100,000	100,000	0	200,000	0	2024-1-1 17:33:04
testwyr2	TESTWYR2	100,000	100,000	18	200,018	5	2024-1-1 17:31:15
testwyr3	TESTWYR3	5,000	0	0	5,000	0	2024-1-1 17:30:44

1 10 / 页 跳至

任务基本信息

同步任务定义 [同步任务定义名称](#)

业务模型 实时同步

状态 暂停

执行用户 admin

开始时间 2024-1-1 15:22:22

结束时间 -

最近操作用户 admin

最近操作时间 -

任务进度

结构迁移 表3/3

全量初始化 表3/3

实时增量 ● 已运行 2h

启动参数

运行类型 直接启动

优先级 medium

Worker分组 default

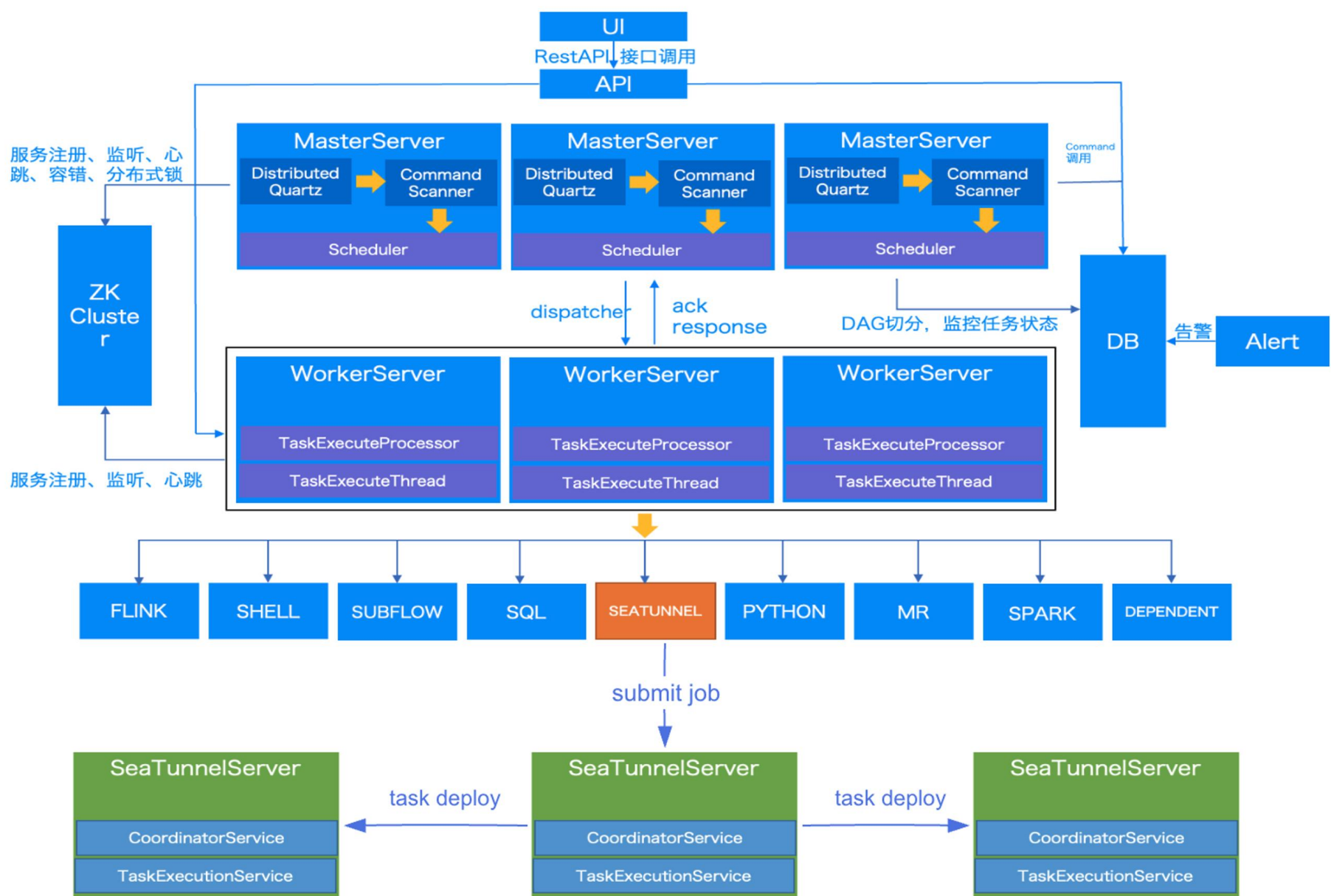
告警通知

失败告警 运维1组

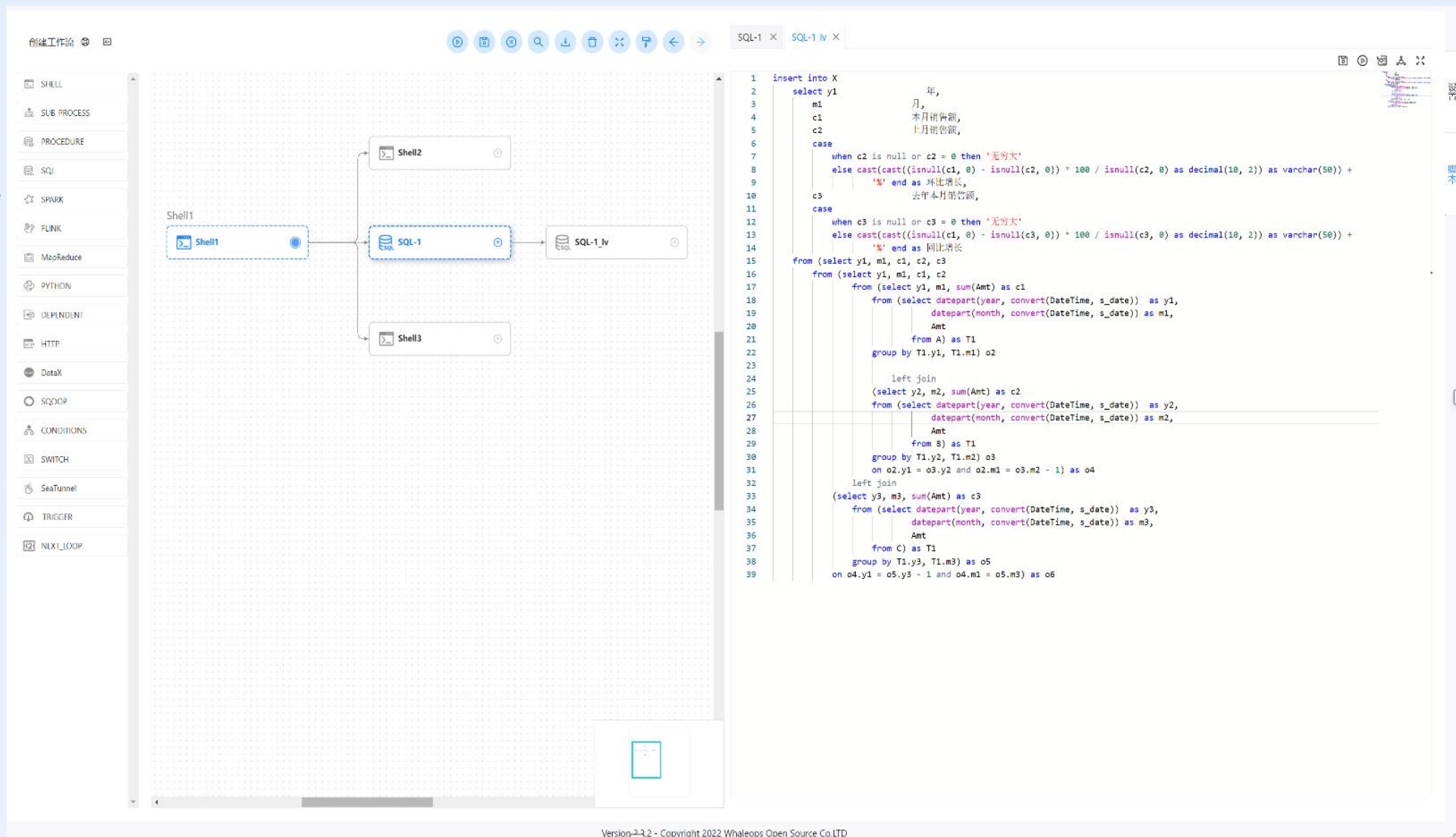
实时增量延迟 运维1组

DDL事件通知 运维1组

全量同步完成通知 运维1组



1. 支持各类计算任务组件:
Amazon DMS、Azure Datafactory,
Amazon Datasync、Apache Linkis, DataX, Sqoop, SeaTunnel等
2. 支持各类云数据库和计算架构, 支持 K8S、MLDB。
3. 平台采用插件式设计, 支持自由扩展数据源支持。
4. 可视化的数据源管理, 数据源统一集中管理, 一次配置, 到处使用, 大大减少配置修改带来的工作量。



The screenshot displays the WhaleStudio workflow editor interface. On the left, a sidebar lists various components such as Shell, SQL, and Python. The main workspace shows a workflow diagram with three components: Shell1, SQL-1, and SQL-1.lv. Shell1 is connected to SQL-1, which is then connected to SQL-1.lv. On the right, a SQL query editor shows a complex query with multiple joins and conditional logic. The query is as follows:

```
1 insert into X
2 select y1
3     m1
4     c1
5     c2
6     case
7         when c2 is null or c2 = 0 then '无较大'
8         else cast(cast((isnull(c1, 0) - isnull(c2, 0)) * 100 / isnull(c2, 0) as decimal(10, 2)) as varchar(50)) +
9             '%' end as 环比增长,
10    c3
11    case
12        when c3 is null or c3 = 0 then '无较大'
13        else cast(cast((isnull(c1, 0) - isnull(c3, 0)) * 100 / isnull(c3, 0) as decimal(10, 2)) as varchar(50)) +
14            '%' end as 同比增长
15 from (select y1, m1, c1, c2, c3
16       from (select y1, m1, c1, c2
17            from (select y1, m1, sum(Amt) as c1
18                 from (select datepart(year, convert(DateTime, s_date)) as y1,
19                      datepart(month, convert(DateTime, s_date)) as m1,
20                      Amt
21                 from A) as T1
22            group by T1.y1, T1.m1) o2
23
24            left join
25            (select y2, m2, sum(Amt) as c2
26             from (select datepart(year, convert(DateTime, s_date)) as y2,
27                  datepart(month, convert(DateTime, s_date)) as m2,
28                  Amt
29             from B) as T1
30            group by T1.y2, T1.m2) o3
31            on o2.y1 = o3.y2 and o2.m1 = o3.m2 - 1) as o4
32        left join
33        (select y3, m3, sum(Amt) as c3
34         from (select datepart(year, convert(DateTime, s_date)) as y3,
35              datepart(month, convert(DateTime, s_date)) as m3,
36              Amt
37         from C) as T1
38        group by T1.y3, T1.m3) as o5
39        on o4.y1 = o5.y3 - 1 and o4.m1 = o5.m3) as o6
```


51CTO WOT 支持160种数据源接口，多种数据集成方式

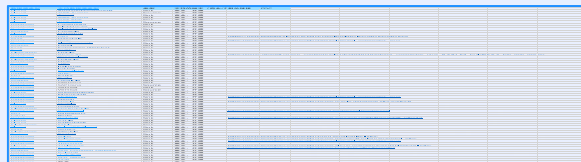
WhaleTunnel支持 160种数据源，例如 MySQL, SAP Hana, Oracle, DB2, SQLServer, Gbase, Kafka, ClickHouse, RedShift、达梦等
平台采用插件式设计，支持自由扩展数据源支持

支持多种

- 批量数据全量、增量集成
- 实时数据集成
- 批量无主键增量集成等

商业版支持商业数据库实时CDC

- Mysql
- PostGreSQL
- SQLServer
- Oracle
- DB2
- AWS Aurora
- 翰高
- StarRocks
- 达梦
- 人大金仓
- PolarDB

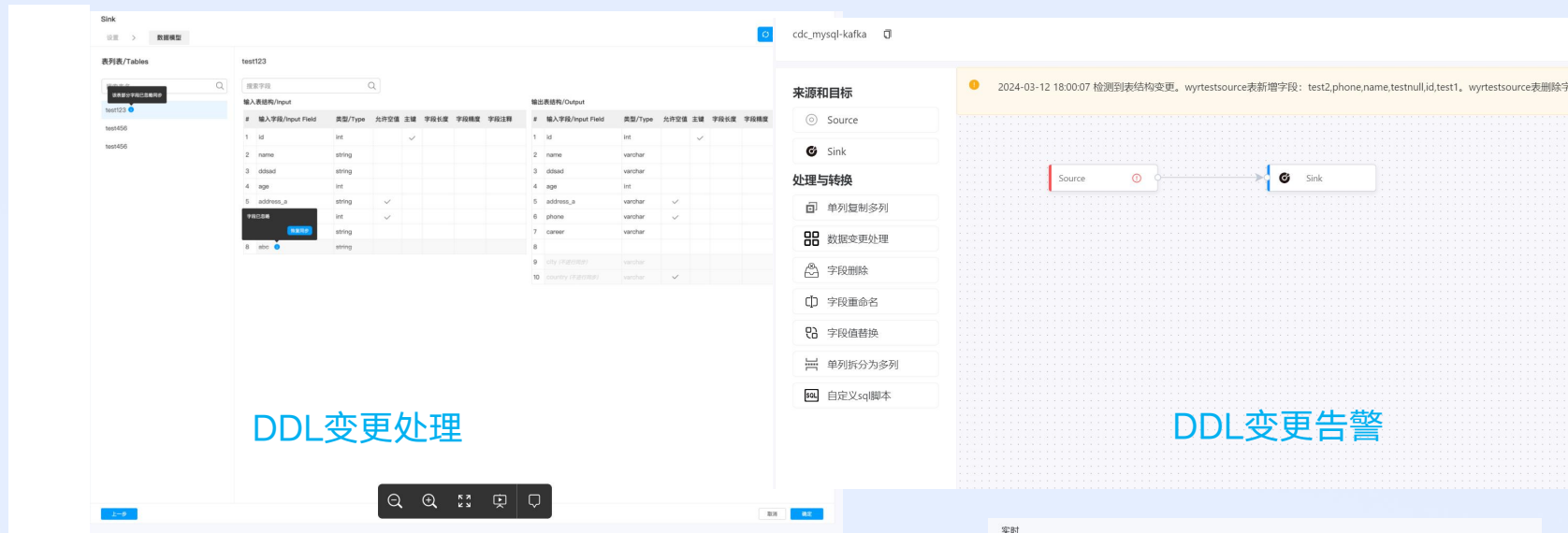


*目前支持160种数据源，详情参见Excel

实时数据处理支持多种实时数据监测

处理：

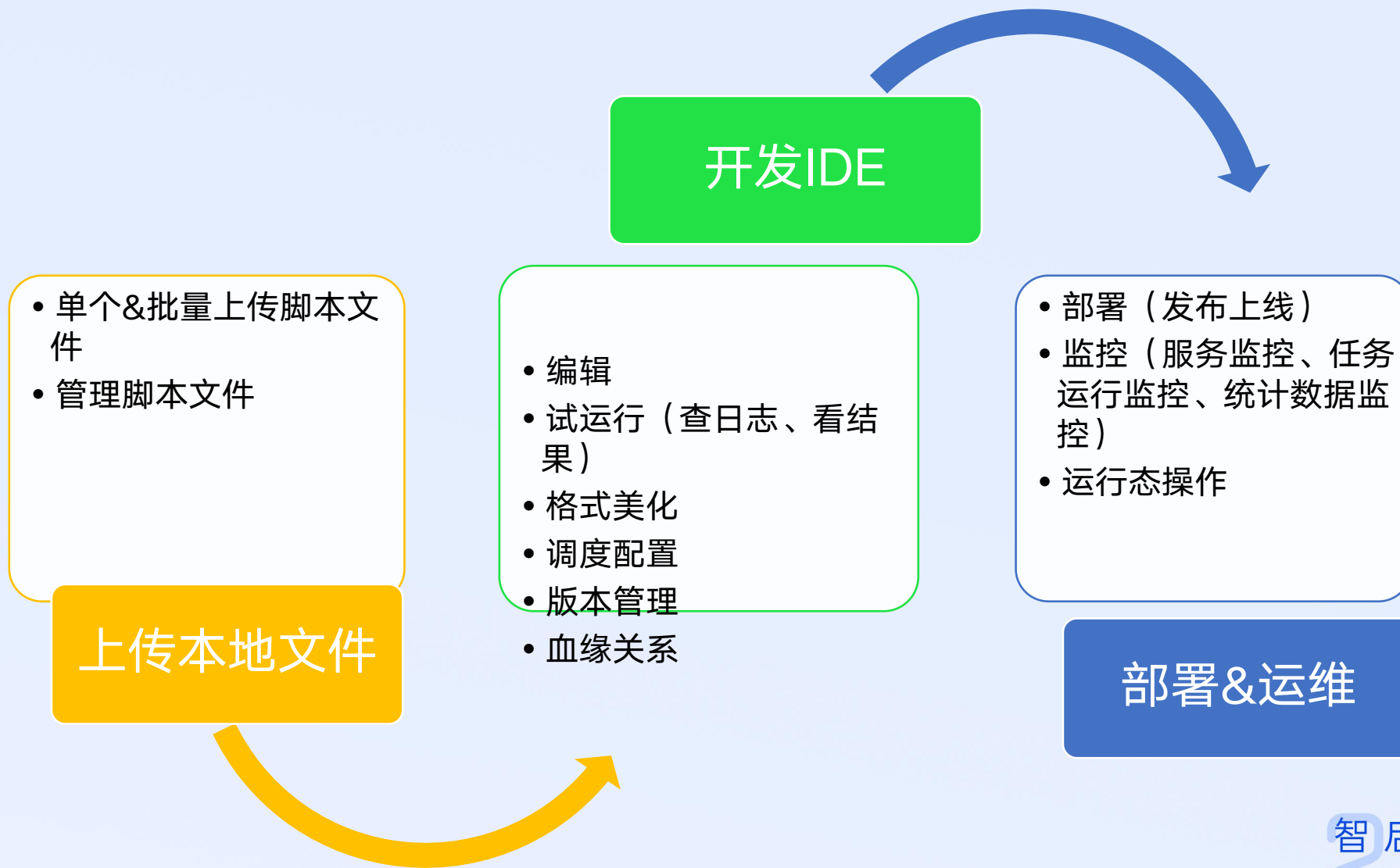
- DDL变更暂停
- DDL变更告警
- DDL暂停加表
- DDL手工处理



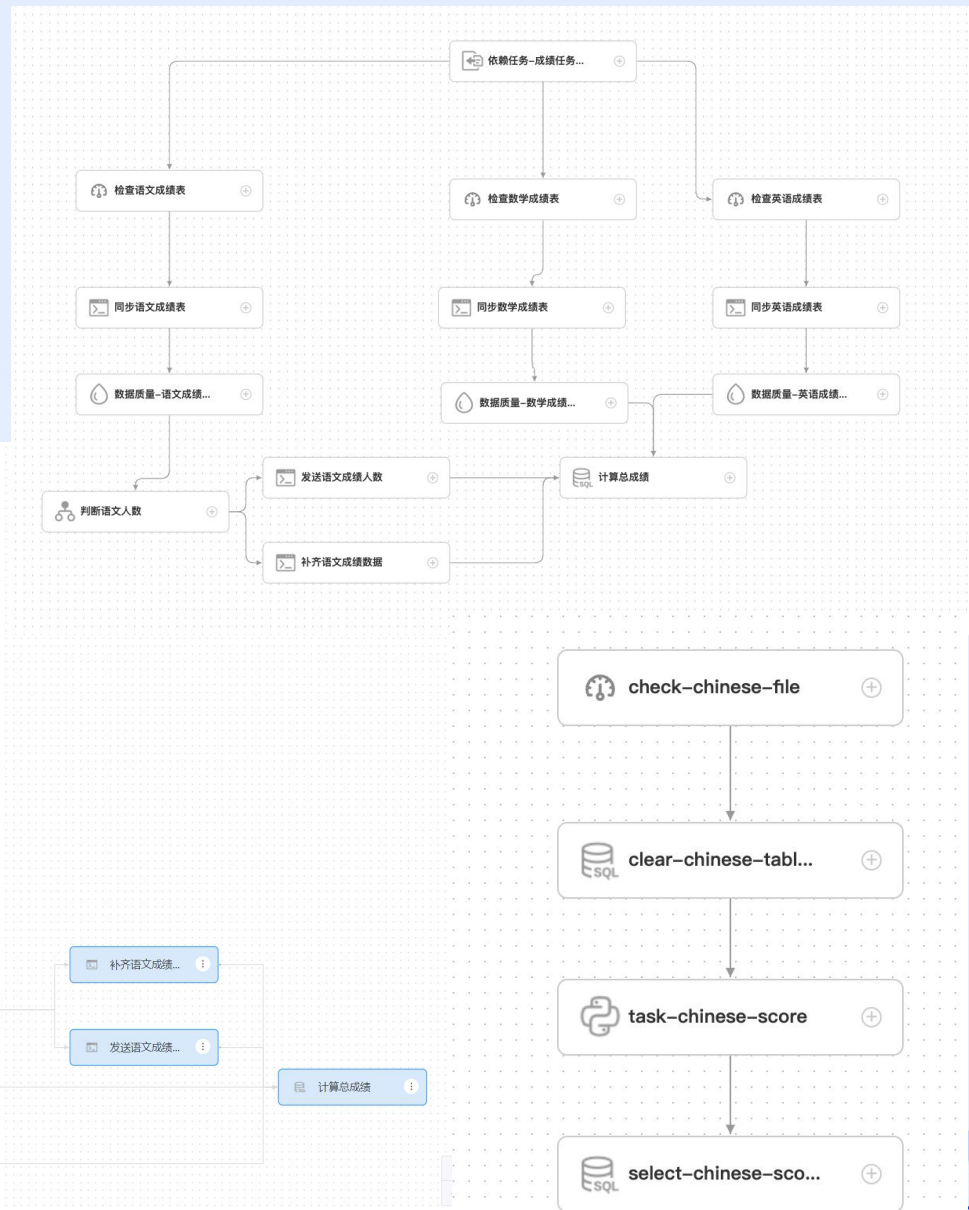
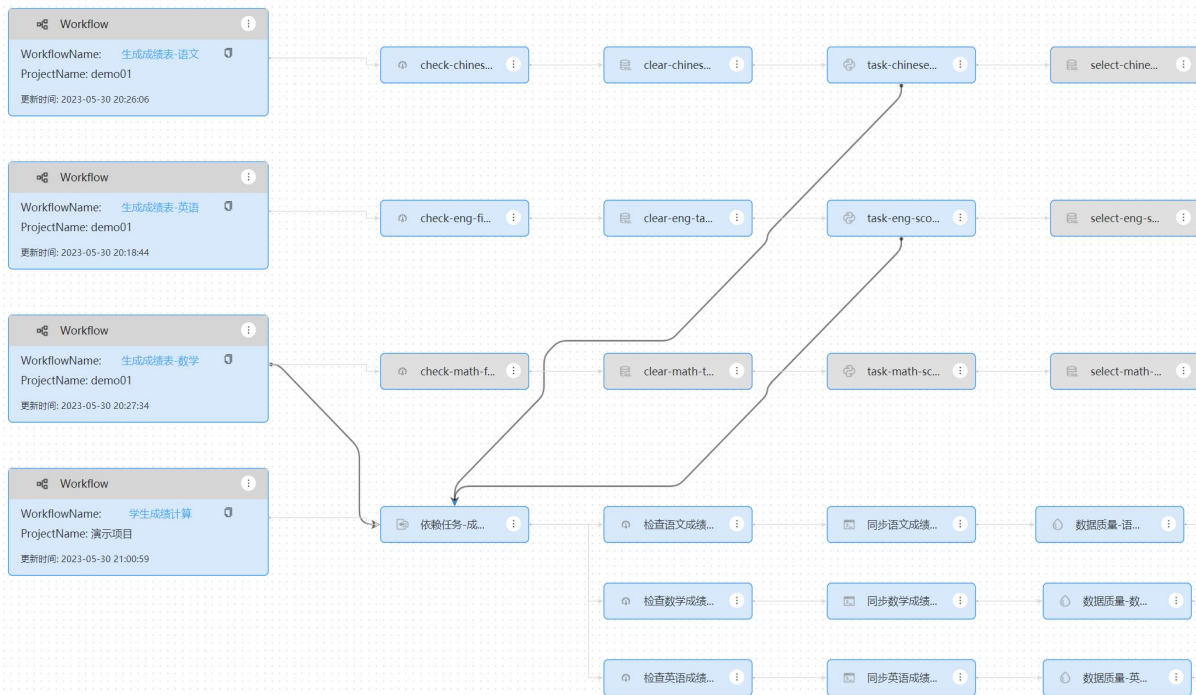
支持多种方式控制&监测速率：

- 数据采集速率控制
- 并发控制
- 数据延迟告警
- 数据全量完成告警
- 数据CDC增量启动告警



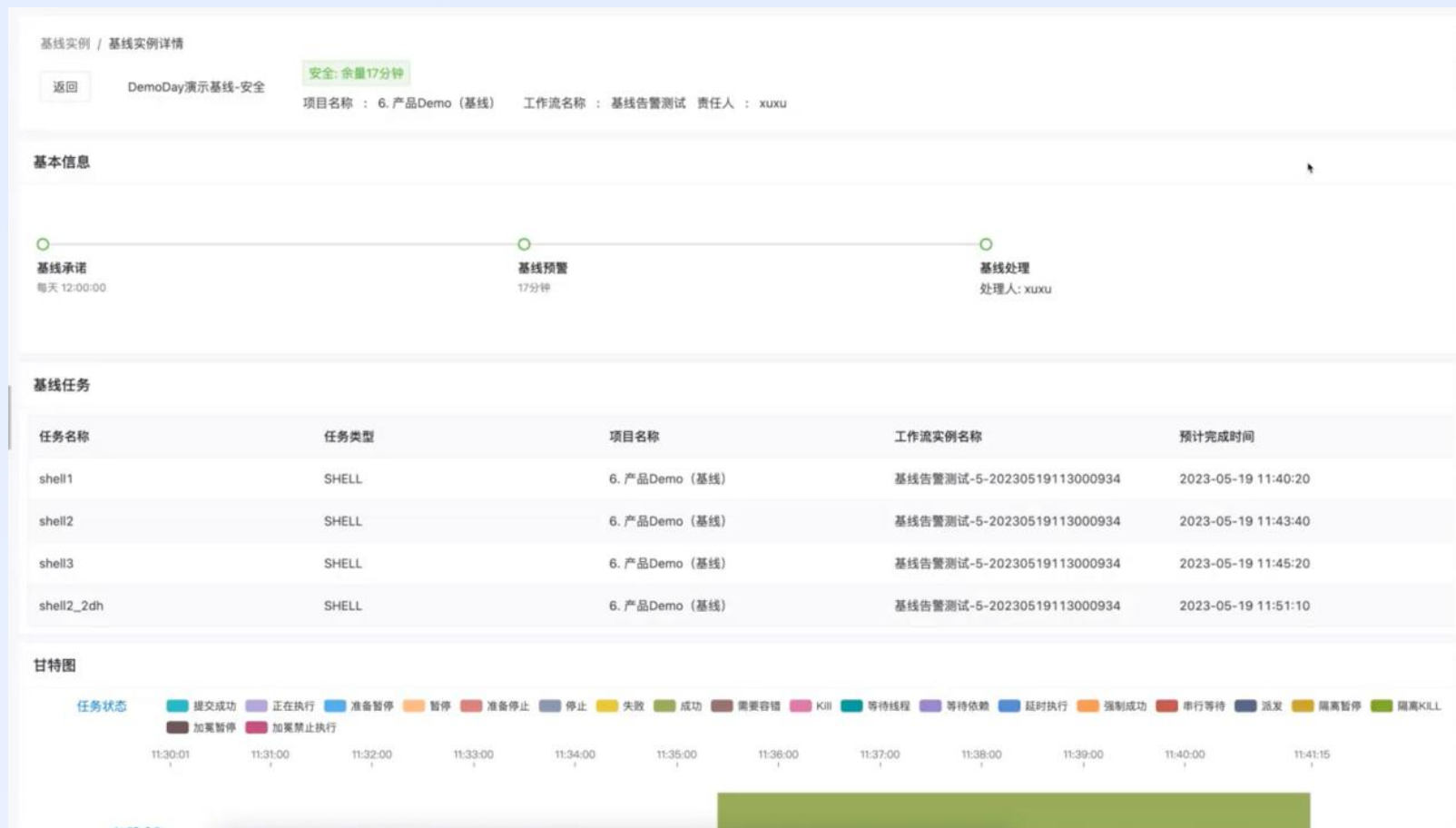


- 全局跨工作流的任务和实例间的依赖关系
- 结合任务与表定义，实现表及血缘分析以及任务操作
- 支持实例级别的依赖链路展示
- 支持全局视图进行停止、暂停、重跑、依赖链重跑等操作



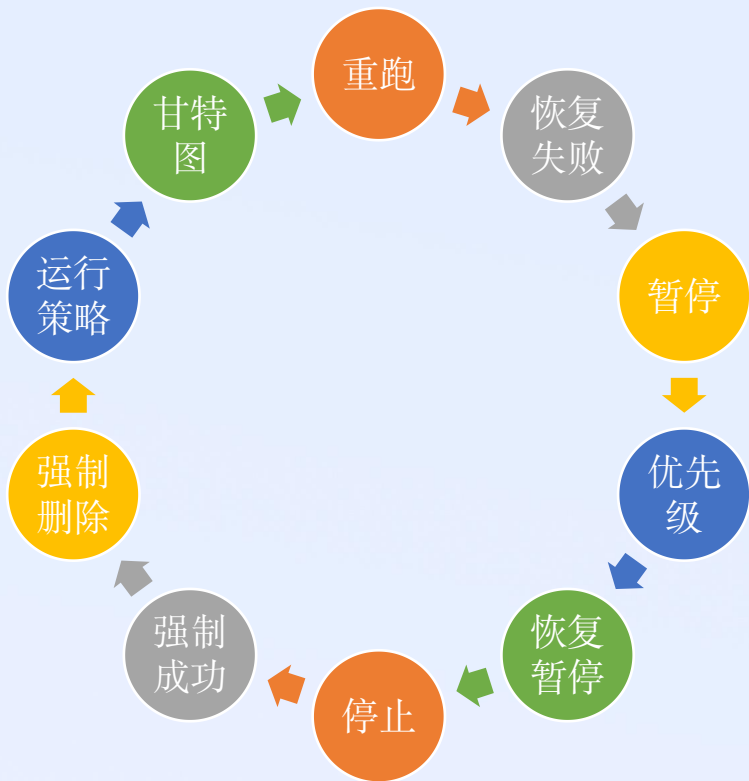
配置智能基线帮助“智能”告警:

- ✓ 定义核心任务基线，多一双“智能”的眼镜
- ✓ 根据任务的执行历史只能推算时长
- ✓ 设置安全预警时间，智能告警

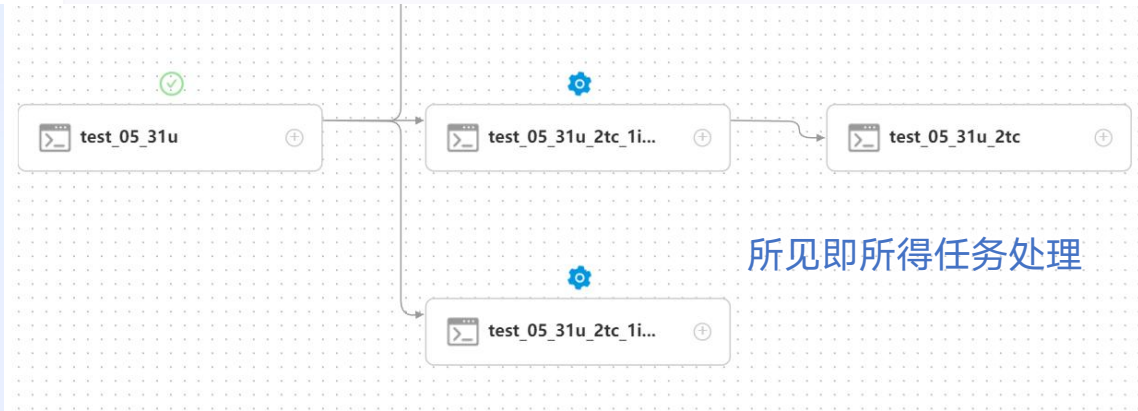
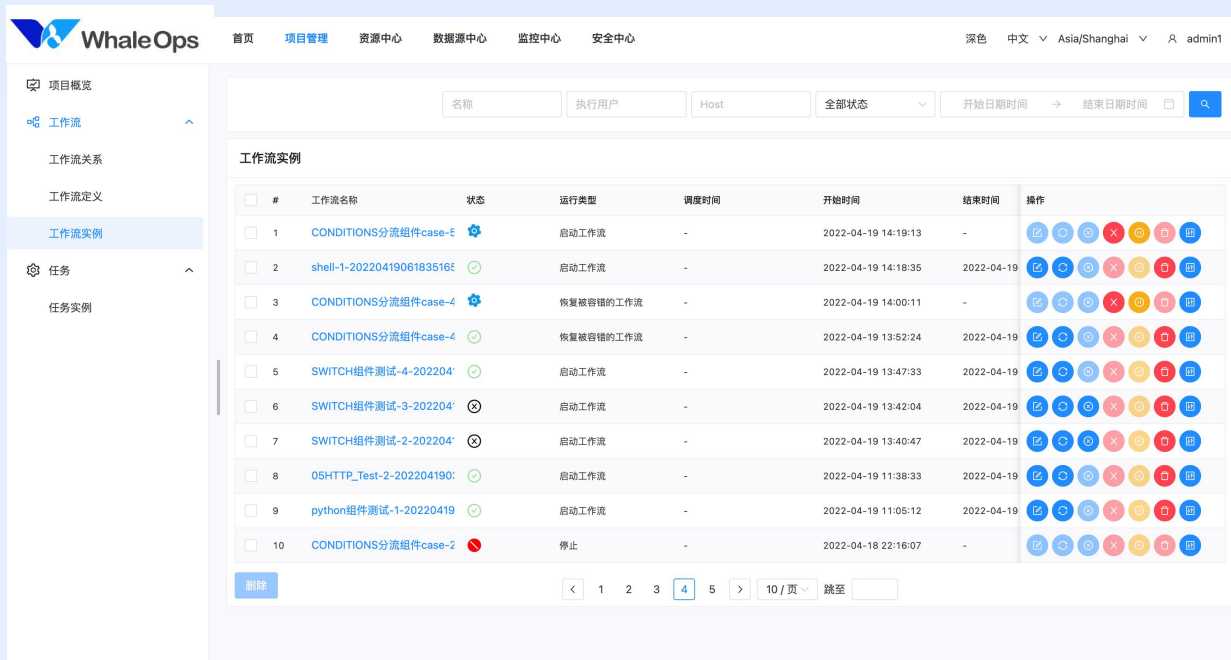


51CTO WOT 多种运维管理手段 —— 帮助运维人员快速处理故障

任务上线之后，面对各种突发情况，有多种手段来确保在任务发生异常时可以协助运维人员快速处理异常。



对 workflow 实例和任务实例进行丰富的人工干预功能

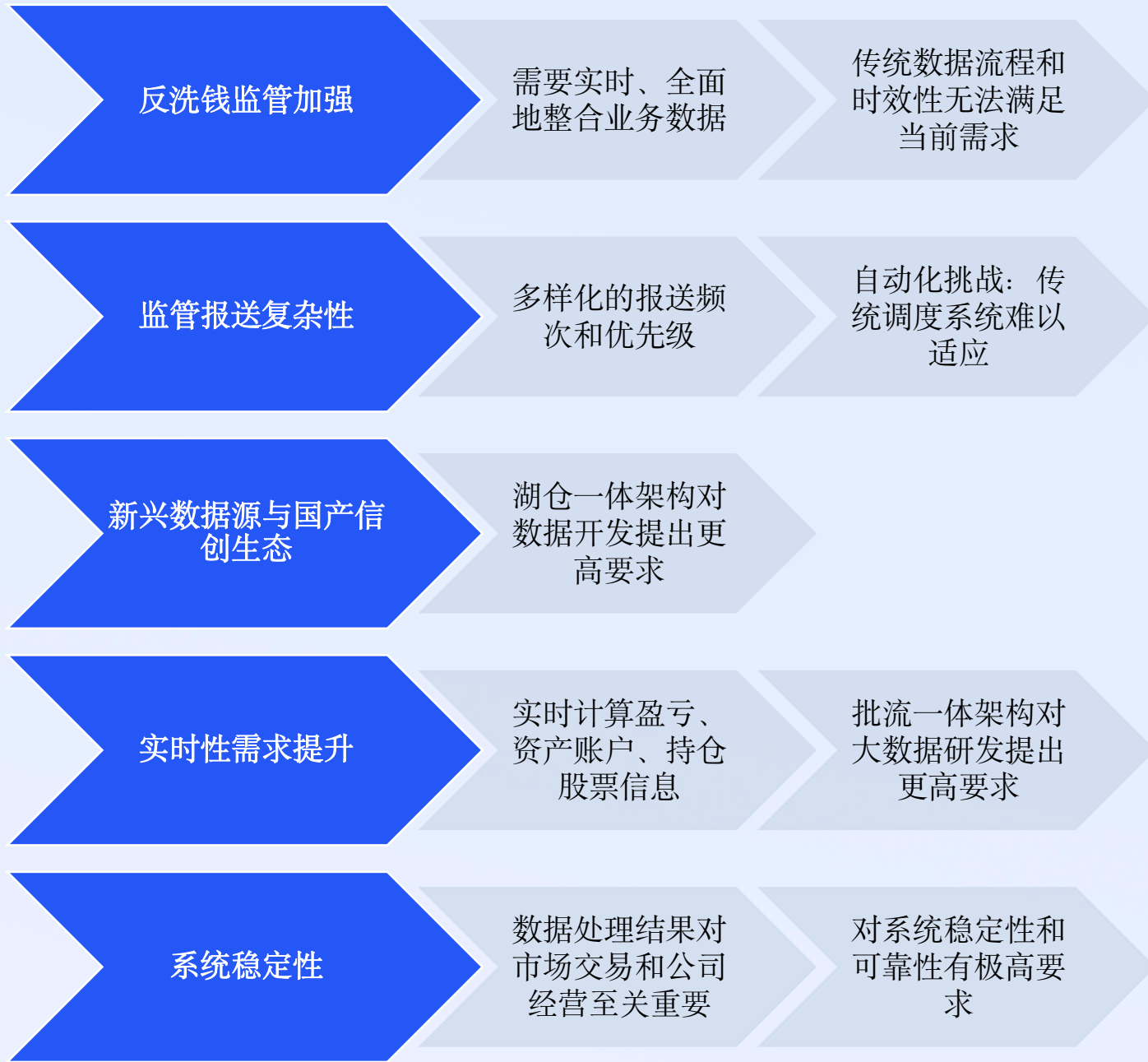


PART · 03

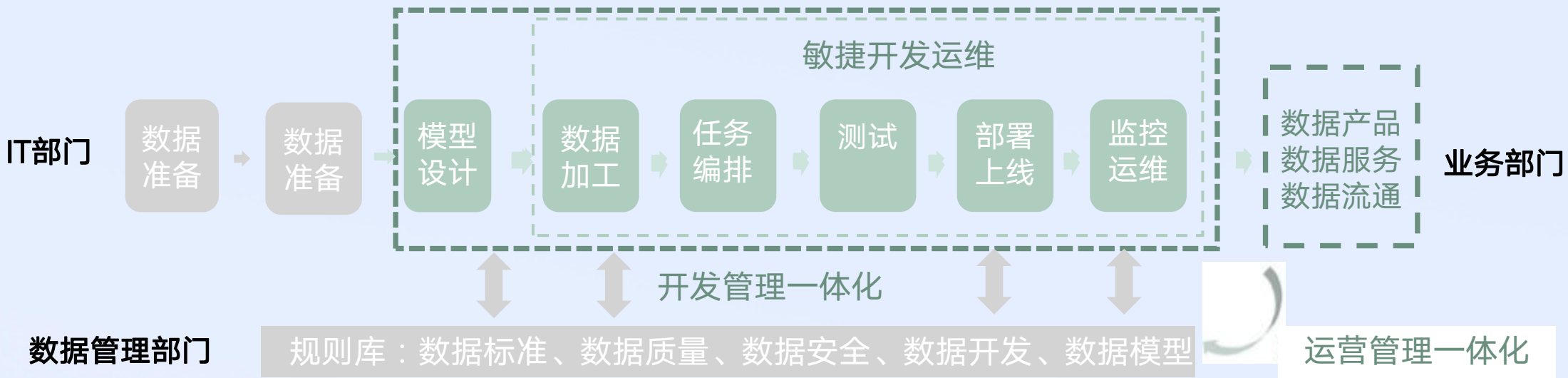
实践案例



中信建投——当前数据开发中存在的挑战



WhaleStudio满足用户数据一体化编辑、上线、数据管控、复杂时间管理等需求，充分提高中信建投数据研发效率。平台应用于公司反洗钱、实时盈亏计算、监管报送、数据精算等多个核心应用，累计编排定义 workflows 超过3000个，上线任务数量接近16000个，交易日平均运行 workflow 实例数量超过5000个，日均任务执行任务数量超过20000个。目前公司各业务线数据处理任务还在持续上线DataOps平台，整个平台规模还在持续增长中。



统一运维管理	白鲸 WhaleStudio 统一运维监控	控制台管理	平台管理
	任务运维 实例监控	集群配置	多租户管理
统一开发	白鲸 WhaleStudio 离线开发	白鲸 WhaleStudio 实时开发	
	代码开发 任务配置	代码开发 任务配置	
资源调度	白鲸 WhaleScheduler 统一调度		
计算引擎	Flink	oushu	
数据底座	OushuDB		
数据采集	白鲸 WhaleTunnel 批量采集	白鲸 WhaleTunnel 实时采集	
数据源	TDSQL	GreatDB	Starrocks Vastbase TDHS_Hive HashData OushuDB GreenplumDB

实际日历		20220321		20220322		20220323		20220324		20220325		20220326		20220327		20220328			
		周一	周二	周三	周四	周五	周六	周日	周一										
切日时间		是否交易日	是否跑批日	是否交易日	是否跑批日	是否交易日	是否跑批日	是否交易日	是否跑批日	是否交易日	是否跑批日	是否交易日	是否跑批日	是否交易日	是否跑批日	是否交易日	是否跑批日		
16点前	交易日	T-1	跑	20220318	跑	20220321	跑	20220322	跑	20220323	跑	20220324	不跑	不跑	不跑	跑	20220328		
		T-2	跑	20220317	跑	20220320	跑	20220321	跑	20220322	跑	20220323	不跑	不跑	不跑	跑	20220327		
	跑批日	T-1	不跑	不跑	20220321	跑	20220322	跑	20220323	跑	20220324	不跑	不跑	不跑	跑	20220328			
		T-2	不跑	不跑	20220320	跑	20220321	跑	20220322	跑	20220323	不跑	不跑	不跑	跑	20220327			
	交易跑批日	T-1	不跑	不跑	20220321	跑	20220322	跑	20220323	跑	20220324	不跑	不跑	不跑	跑	20220328			
		T-2	不跑	不跑	20220320	跑	20220321	跑	20220322	跑	20220323	不跑	不跑	不跑	跑	20220327			
	自然日	T-1	跑	20220320	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	跑	20220326	跑	20220327	
		T-2	跑	20220319	跑	20220320	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	跑	20220326	
	16点后	16点后翻牌为当日日期, 16点前为昨日日期																	
		交易日	T-1	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	不跑	不跑	不跑	跑	20220328	
			T-2	跑	20220320	跑	20220321	跑	20220322	跑	20220323	跑	20220324	不跑	不跑	不跑	跑	20220327	
		跑批日	T-1	跑	20220321	跑	20220322	跑	20220323	跑	20220324	不跑	不跑	不跑	不跑	跑	20220328		
T-2			跑	20220320	跑	20220321	跑	20220322	跑	20220323	不跑	不跑	不跑	不跑	跑	20220327			
交易跑批日		T-1	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	不跑	不跑	不跑	跑	20220328		
		T-2	跑	20220320	跑	20220321	跑	20220322	跑	20220323	跑	20220324	不跑	不跑	不跑	跑	20220327		
自然日		T-1	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	跑	20220326	跑	20220327	跑	20220328	
		T-2	跑	20220320	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	跑	20220326	跑	20220327	
24点切日		自然日	T-1	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	跑	20220326	跑	20220327	跑	20220328
			T-2	跑	20220320	跑	20220321	跑	20220322	跑	20220323	跑	20220324	跑	20220325	跑	20220326	跑	20220327

PART · 04

DataOps 未来



大模型在数据处理流程中可以扮演多种角色，提高整个数据处理流程的效率和智能化水平。大模型将应用于以下方面：

智能调度策略

数据处理涉及复杂的任务调度，大模型可以分析历史作业执行情况、资源使用状况，从而预测未来的工作流需求，智能地调度任务和分配资源。减少延迟，提高整体处理速度，并优化资源利用率。

数据质量检测与清洗

在数据同步过程中，大模型可以辅助自动检测数据质量问题，比如识别异常值、缺失数据或不一致性。通过机器学习算法，模型可以学习数据特征，自动清洗和修正数据，确保数据同步后的质量。

智能数据分类与标签

对于需要分类或标签化的数据，大模型可以自动分析数据内容，对其进行分类或附加有意义的标签，特别是在多模态数据处理场景下，这对于后续的数据分析和应用至关重要。

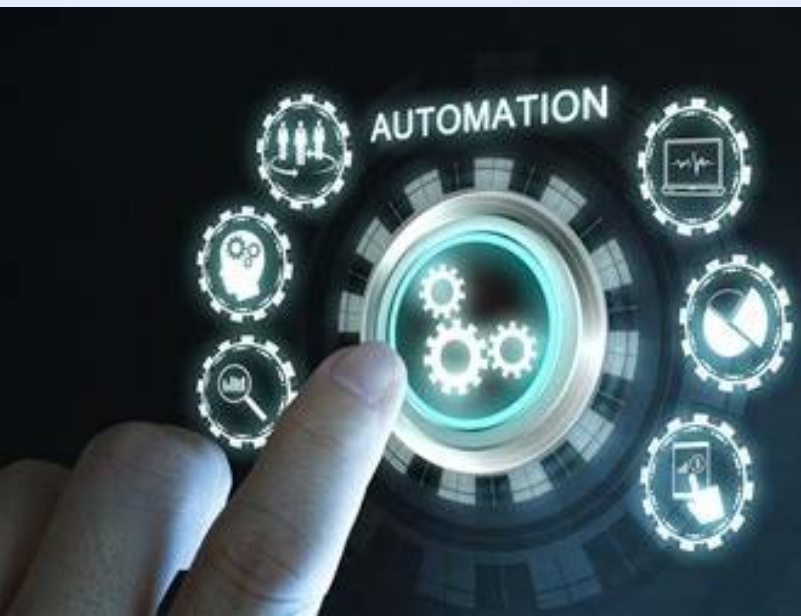
自适应数据同步策略

根据网络状况、数据变化频率和业务需求，大模型可以帮助动态调整数据同步策略，比如选择最合适的同步频率、确定优先级高的数据流，以优化同步效率和减少带宽消耗。

自动化异常处理

在数据传输或处理过程中遇到异常时，大模型可以基于历史数据和模式识别，自动识别异常原因并触发相应的处理机制，减少人工干预，提高处理效率。

Copilot



智能化数据处理

随着AI技术的发展，DataOps将越来越智能化，实现更高效、自动化的数据处理。



数据安全保护

数据安全保护将成为DataOps未来发展的重要方向，如何在确保数据安全的前提下实现数据共享和利用将成为关键。



跨平台/云数据治理

随着企业业务的多平台化和全球化，跨平台的数据治理将成为DataOps的重要发展趋势。



<https://seatunnel.apache.org>



<https://dolphinscheduler.apache.org>



白鲸开源官网: <https://www.whaleops.com>



谢谢观看

THANKS