

GOTC

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

OPEN SOURCE , OPEN WORLD

「综合技术」专场

新一代大数据调度平台
Apache DolphinScheduler 最新进展 & Roadmap

代立冬 2021年08月01日

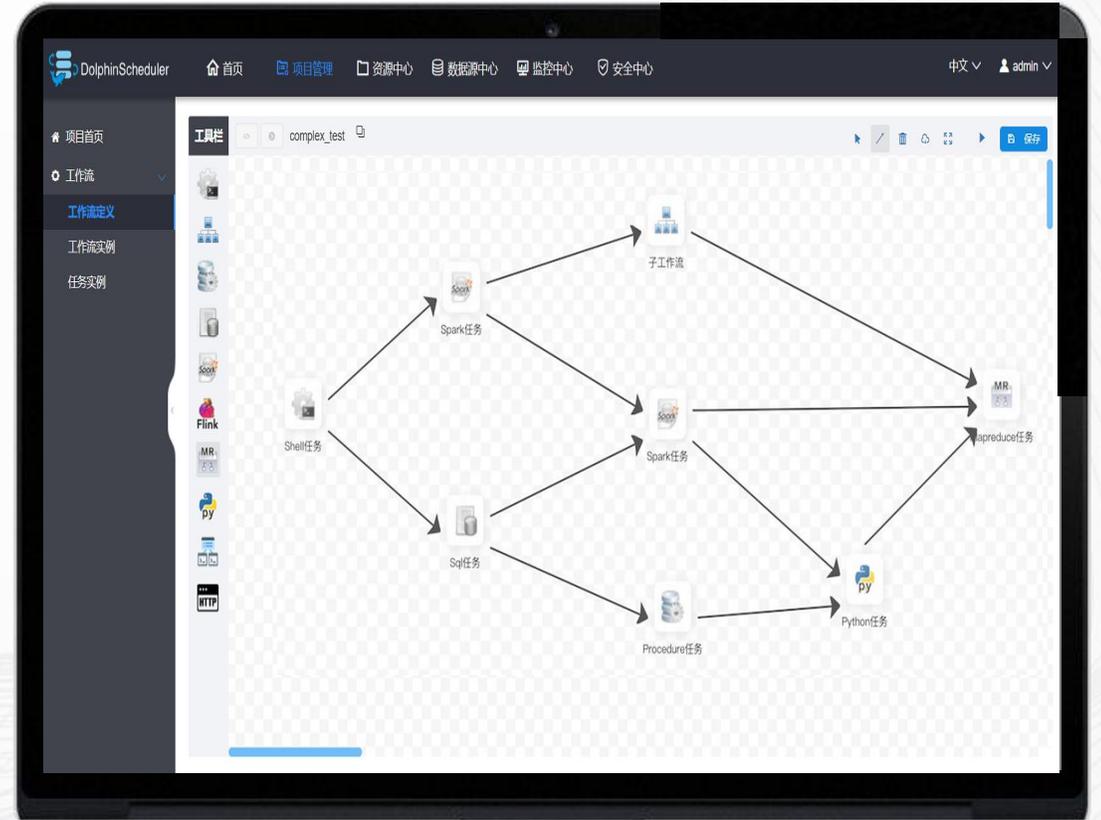
DolphinScheduler 介绍

Apache DolphinScheduler 是一个云原生并带有强大可视化界面的大数据 workflow 调度系统。

2021 年 04 月 09 日正式成为 Apache 顶级项目。

首个由国人主导并贡献到 Apache 基金会的大数据 workflow 领域顶级项目。已累计有 400+ 公司在生产上使用。

DolphinScheduler 致力于在数据 workflow 编排中“解决复杂的大数据任务依赖及触发关系，让各种大数据任务类型开箱即用”。



DolphinScheduler 部分用户(不分先后)

GOTC



社区建设情况

- Apache DolphinScheduler Group 7 (499)
- Apache DolphinScheduler Group 4 (460)
- Apache DolphinScheduler Group 2 (482)
- Apache DolphinScheduler Group 1 (495)
- Apache DolphinScheduler group 10 (229)
- Apache DolphinScheduler Group 5 (435)
- Apache DolphinScheduler Group 6 (449)
- Apache DolphinScheduler Group 9 (499)
- Apache DolphinScheduler Group 8 (459)
- Apache DolphinScheduler Group 3 (474)
- DolphinScheduler Developer Group (304)
- DolphinScheduler 贡献者种子孵化群 (211)



贡献者公司分布(100+公司)

Contributors 219



+ 208 contributors

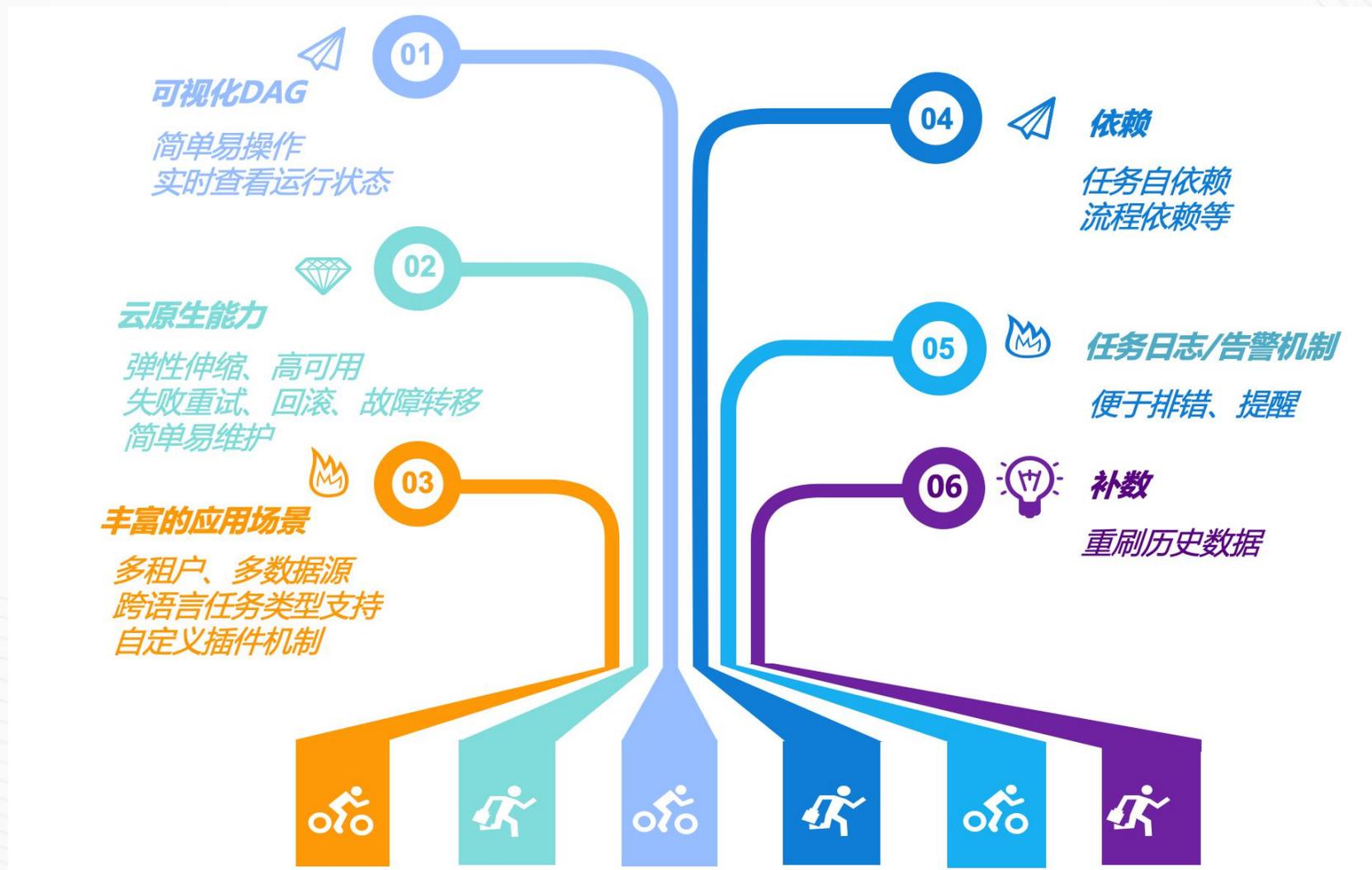
代码贡献者

Contributors 72



+ 61 contributors

文档贡献者



DolphinScheduler 优势



高可靠性

- 去**中心化**的多 Master 和多 Worker, 自身高可用能力
- 采用任务队列来避免过载, 不会造成机器卡死



丰富的使用场景

- 支持暂停恢复操作
- 支持多租户, 权限管理等大数据应用场景
- 支持近 20 种任务类型, 如 Spark, Hive, MR, Python, Sub-Process, Shell 等



简单易用

- 一键部署 – 简化部署, 易维护
- **可视化 DAG** 界面, 所有流程定义都是可视化, 通过拖拽任务形成 workflow 模板
- 支持 Open API 方式与第三方系统对接



高扩展性、云原生能力

- 支持自定义任务类型
- 调度器使用分布式调度, 调度能力随集群线性增长
- 弹性伸缩, Master 和 Worker 支持动态上下线



- Task 以 DAG 形式连接, 实时监控任务的状态



- 支持 Shell、MR、Spark、SQL、依赖等 10 多种任务类型



- workflow 优先级、任务优先级, 全局参数及局部自定义参数



- workflow 可定时、依赖、手动、暂停/停止/恢复



- 支持多租户、补数、日志在线查看及资源在线管理



- 完善的系统服务监控, 任务超时告警/失败

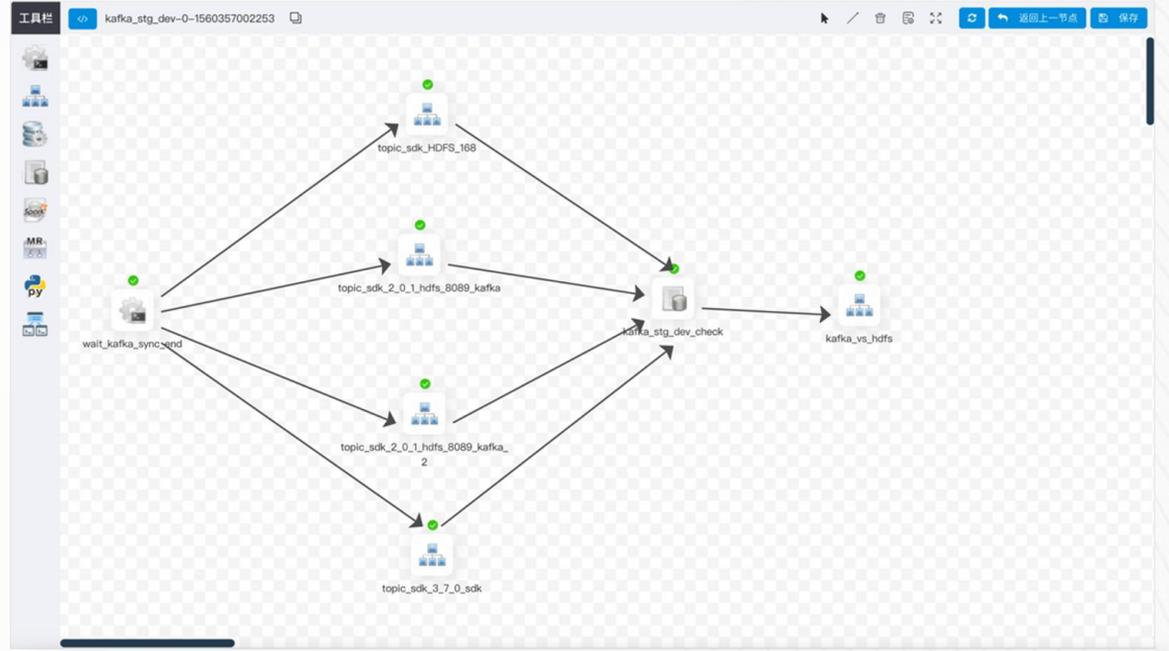
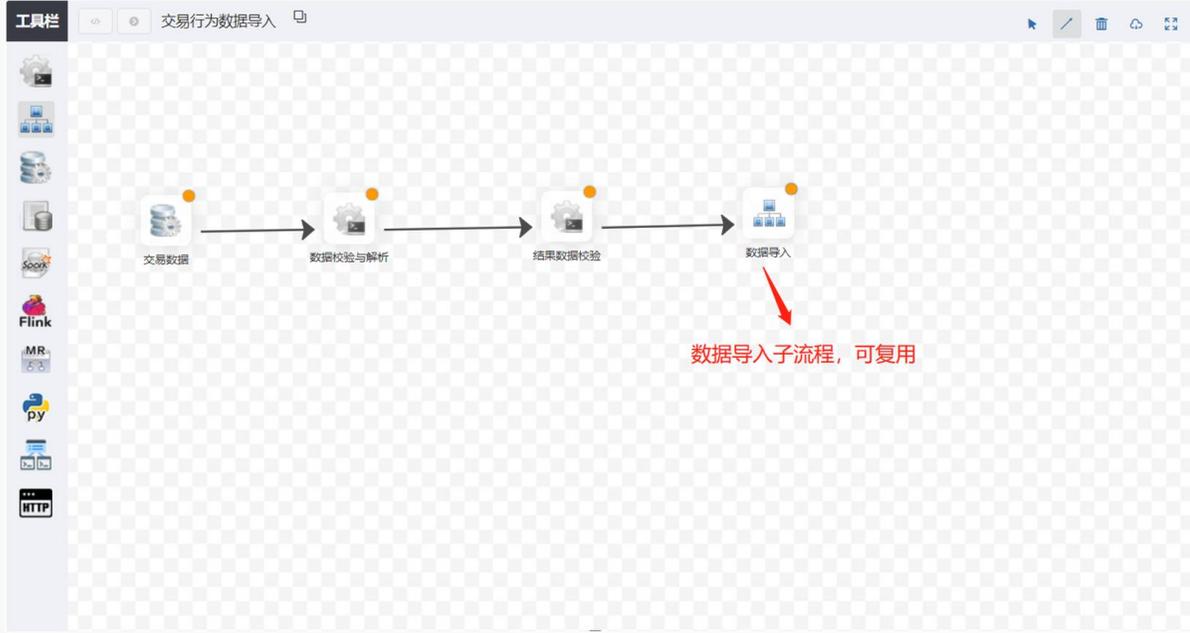


- 去中心化设计确保系统的稳定、高可用



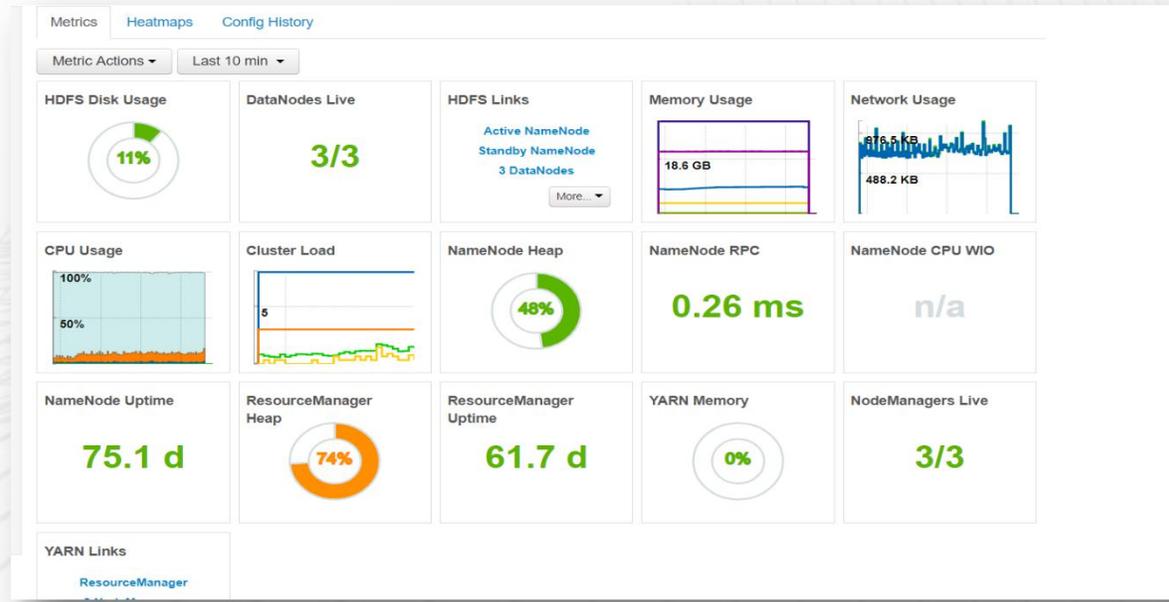
- 支持每日百万量级任务稳定运行
- 弹性伸缩

DolphinScheduler 可视化流程



工作流实例

名称	host	状态									
编号	工作流名称	运行类型	调度时间	开始时间	结束时间	运行时长	运行次数	host	容错标识	状态	操作
1	var_test-0-1589082290636	启动工作流	-	2020-04-28 21:58:11	2020-04-28 21:58:15	4	1	192.168.220.241	NO	成功	刷新
2	var_test-0-1588082277549	启动工作流	-	2020-04-28 21:57:58	2020-04-28 21:58:02	4	1	192.168.220.241	NO	成功	刷新
3	var_test-0-1588082270447	启动工作流	-	2020-04-28 21:57:50	2020-04-28 21:57:54	4	1	192.168.220.241	NO	成功	刷新
4	var_test-0-1587729893626	启动工作流	-	2020-04-24 20:04:54	2020-04-24 20:05:00	6	1	192.168.220.241	NO	成功	刷新
5	var_test-0-1586532775630	启动工作流	-	2020-04-10 23:32:56	2020-04-10 23:33:00	4	1	192.168.220.241	NO	成功	刷新
6	test-0-1586500590633	调度执行	2020-04-09 18:00:00	2020-04-10 14:36:31	2020-04-10 14:36:35	4	1	192.168.220.241	NO	成功	刷新
7	test-0-1586500590648	调度执行	2020-04-09 19:00:00	2020-04-10 14:36:31	2020-04-10 14:36:35	4	1	192.168.220.241	NO	成功	刷新
8	test-0-1586500590662	调度执行	2020-04-09 20:00:00	2020-04-10 14:36:31	2020-04-10 14:36:35	4	1	192.168.220.241	NO	成功	刷新
9	tes										
10	tes										





数据加工平台任务监控总览

workflow实例

名称	host	状态	操作
var_test-0-1586532775630	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
test-0-1586500590633	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
test-0-1586500590648	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
test-0-1586500590662	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
test-0-1586500590683	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
test-0-1586500590702	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
var_test-0-1586500590722	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
var_test-0-1586500590741	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
var_test-0-1584434190810	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫
var_test-0-1583067398834	192.168.220.241	NO	🔄 🏠 📄 🗑️ 🚫

流程实例状态查看



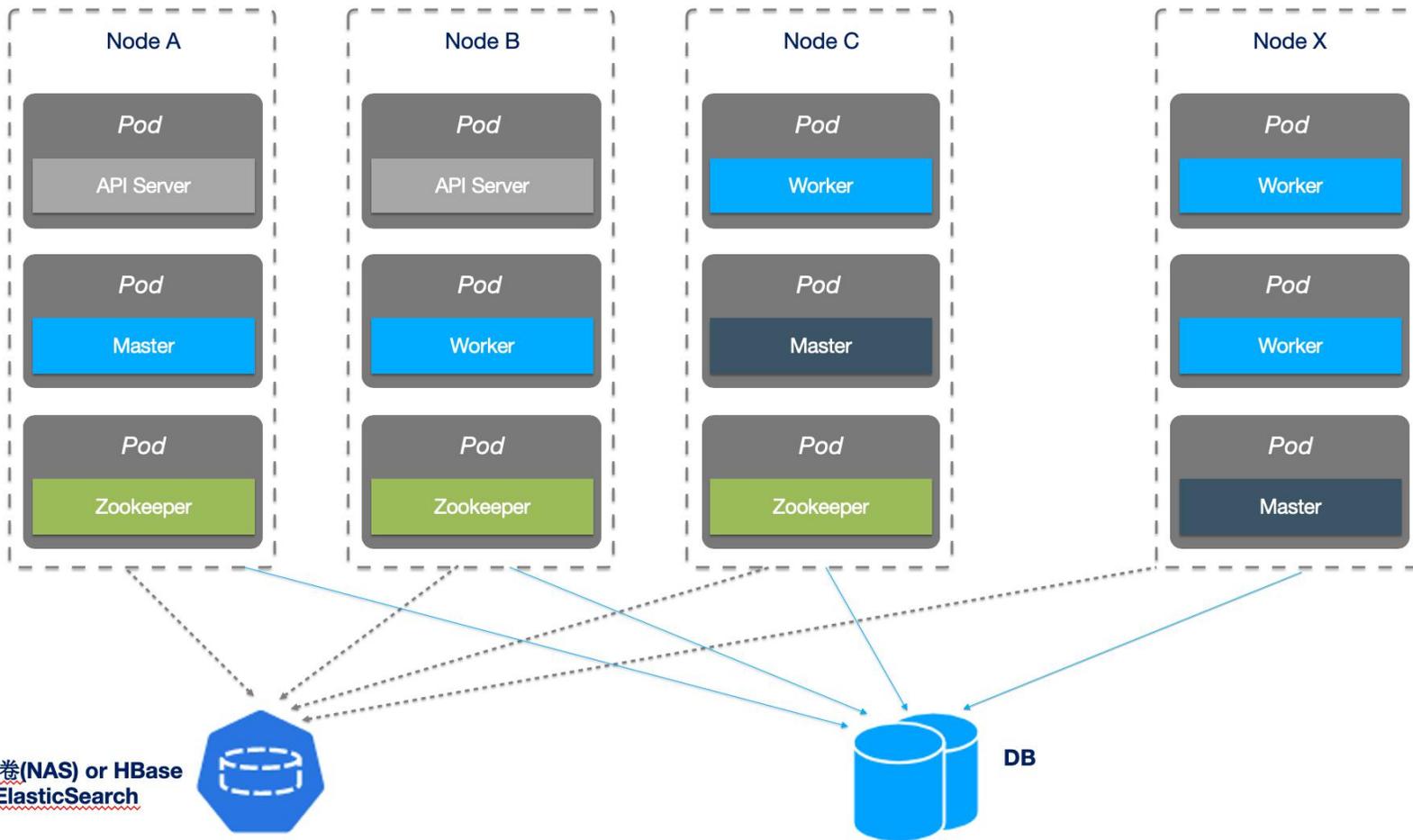
任务执行状态回溯

任务实例

编号	名称	时间	host	运行时长(s)	重试次数	操作
1	shell	04-10 23:32:58	192.168.220.241	2	0	🔄
2	spark	04-10 14:36:33	192.168.220.241	1	0	🔄
3	spark	04-10 14:36:33	192.168.220.241	1	0	🔄
4	spark	04-10 14:36:33	192.168.220.241	1	0	🔄

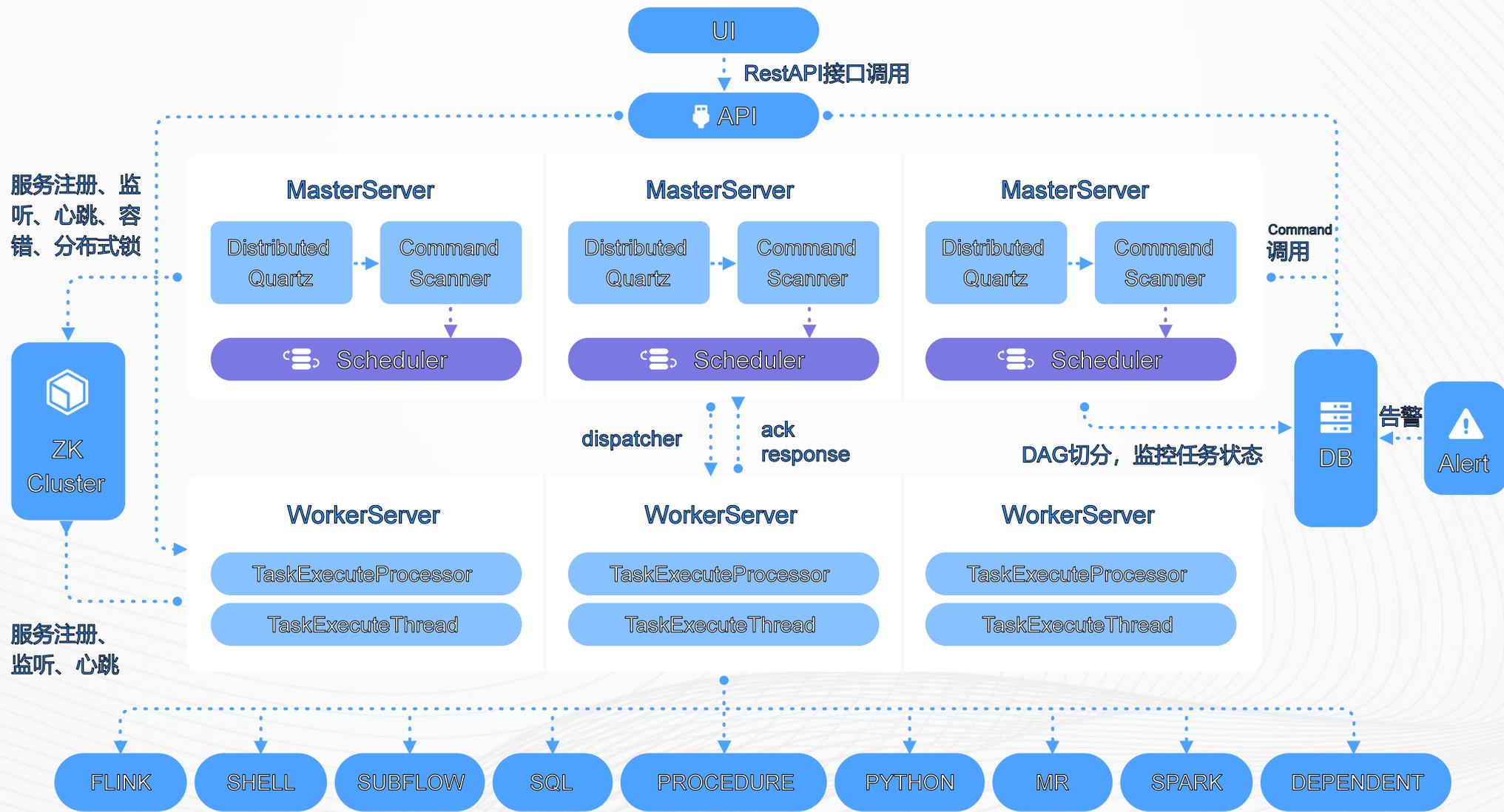
```
[INFO] 2020-04-10 14:36:31.808 - [taskAppId=TASK-1-551-1212]:[63] - spark task params {"mainArg s":"","driverMemory":"512M","executorMemory":"2G","programType":"SCALA","mainClass":"","driver Cores":1,"deployMode":"cluster","executorCores":2,"mainJar":{"res":"shelltest.sh"},"sparkVersion":"SPA RK2","numExecutors":2,"localParams":[],"others":"","resourceList":[]}  
[INFO] 2020-04-10 14:36:31.809 - [taskAppId=TASK-1-551-1212]:[113] - spark task command : ${SPARK HOME2}/bin/spark-submit --master yarn --deploy-mode cluster --class d --driver-cores 1 --driver-memor y 512M --num-executors 2 --executor-cores 2 --executor-memory 2G --queue default shelltest.sh  
[INFO] 2020-04-10 14:36:31.809 - [taskAppId=TASK-1-551-1212]:[99] - tenantCode user:hdfs, task dir:1 551_1212  
[INFO] 2020-04-10 14:36:31.809 - [taskAppId=TASK-1-551-1212]:[103] - create command file:/tmp/dolp
```

任务执行日志查看



弹性伸缩
充分利用服务器资源
环境隔离

DolphinScheduler 1.3 架构



高性能 - Master 重构

减少数据库轮询

去分布式锁

减少线程使用

DolphinScheduler 2.0 架构规划与 Roadmap

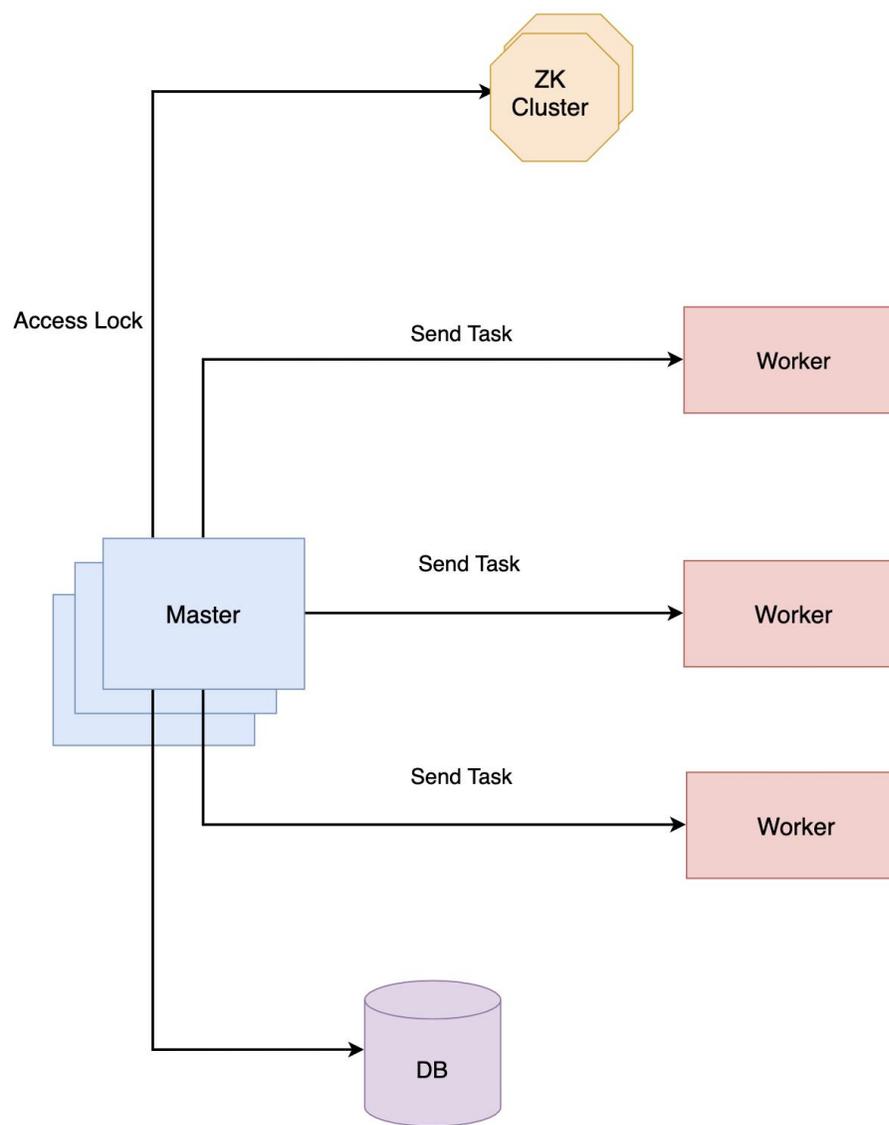
1.3 版本现状 - 分布式锁

Master现状

MasterSchedulerService
获取分布式锁同时轮询
command
生成ProcessInstance

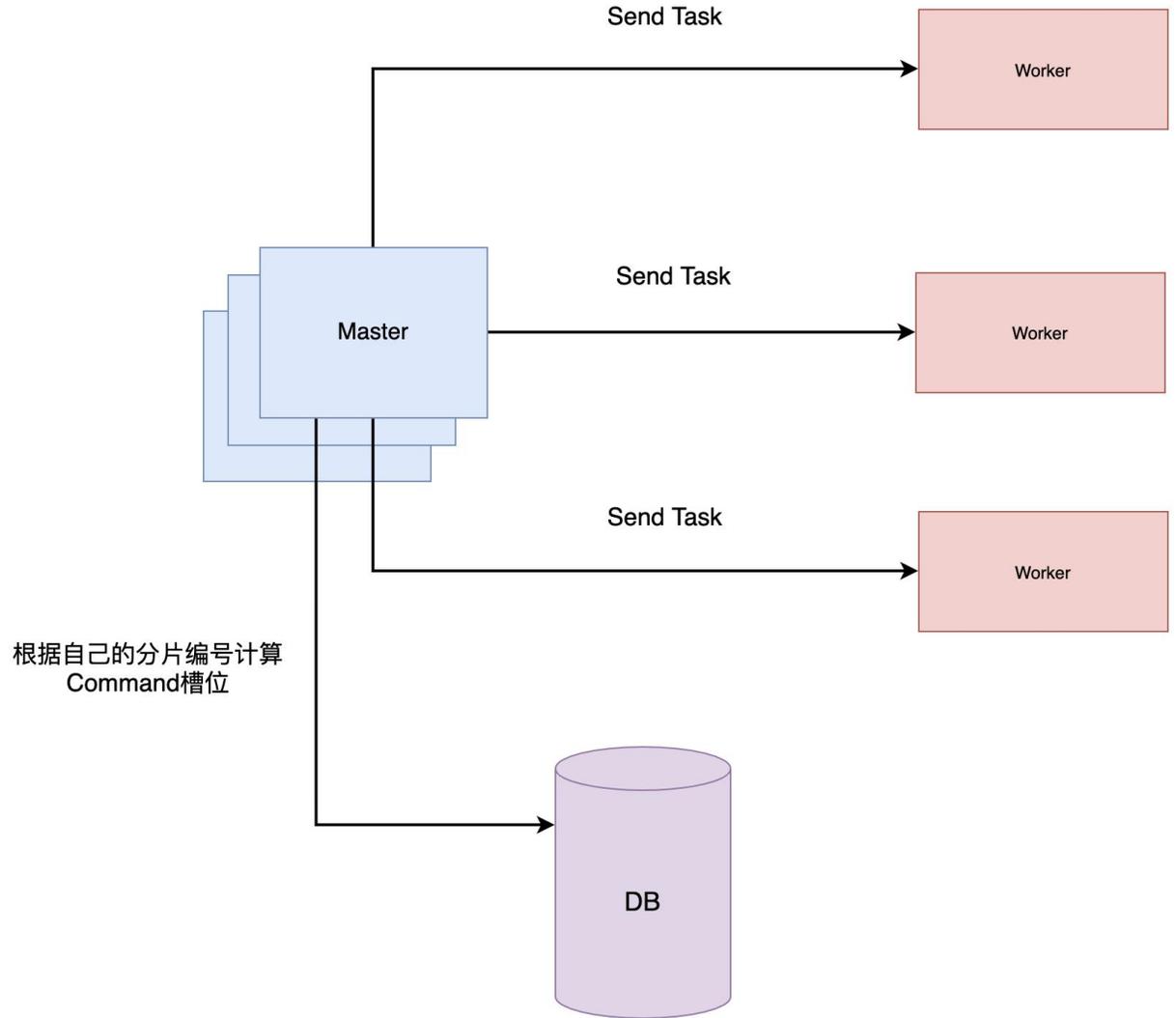
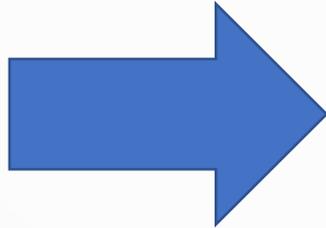
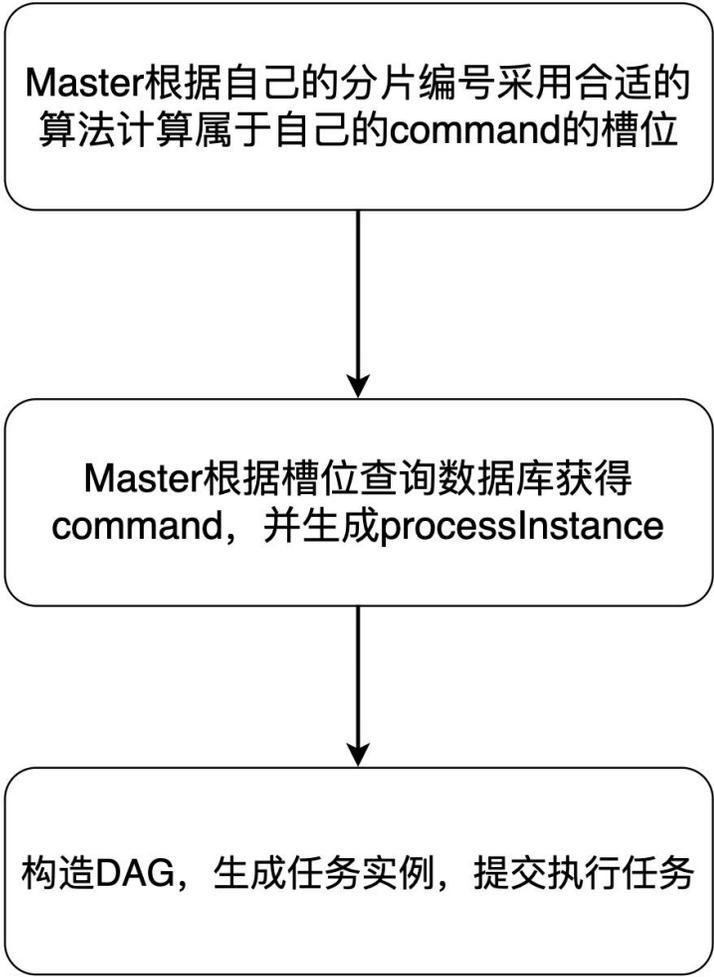
MasterExecThread
构建DAG
轮询 workflow实例
生成任务实例

MasterTaskExecThread
提交执行任务
轮询任务状态
取消/暂停/停止任务



DolphinScheduler 2.0 架构 Roadmap

2.0 优化 – 去分布式锁

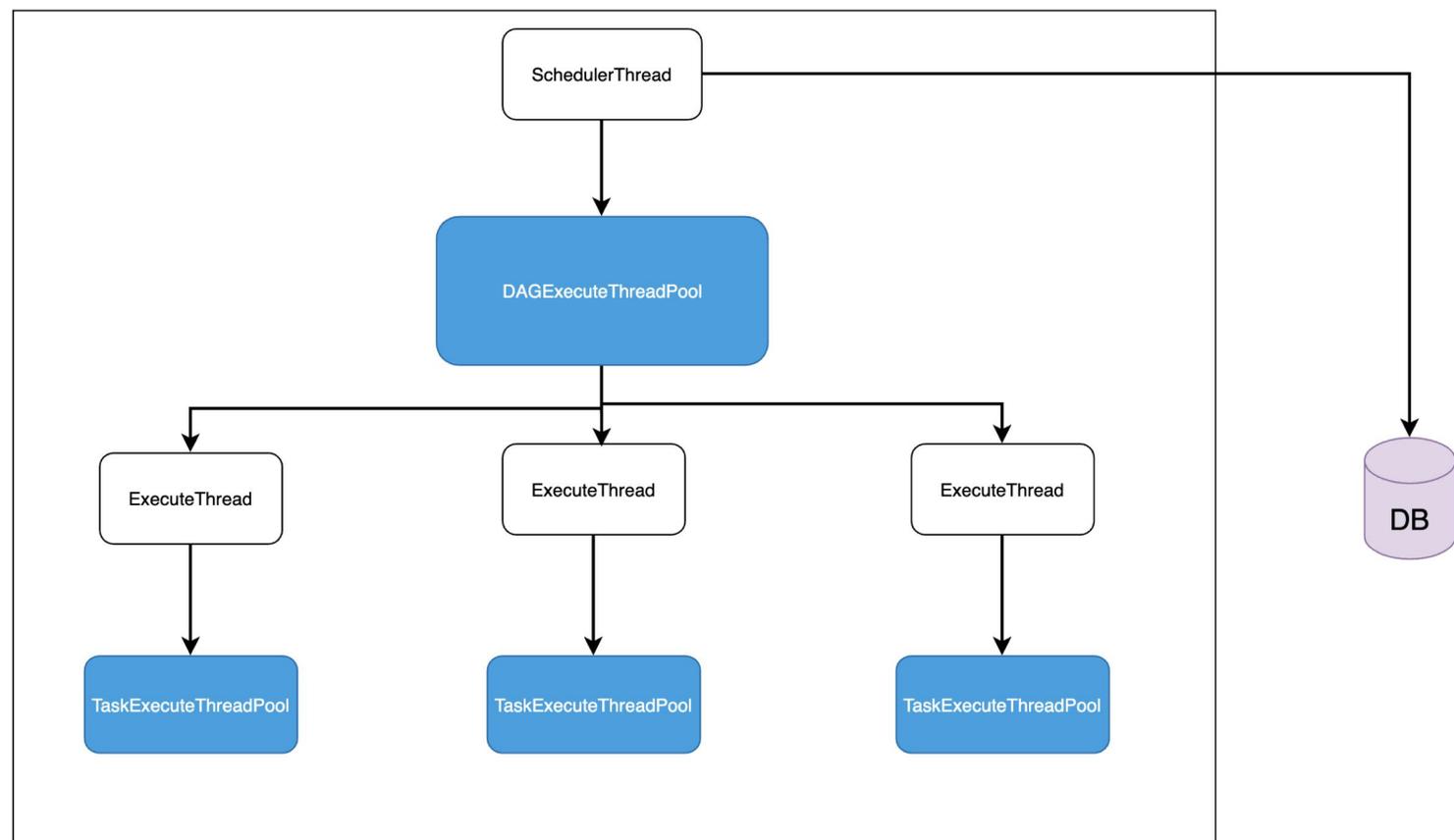


DolphinScheduler 2.0 架构 Roadmap

重构 Master 中的线程模型

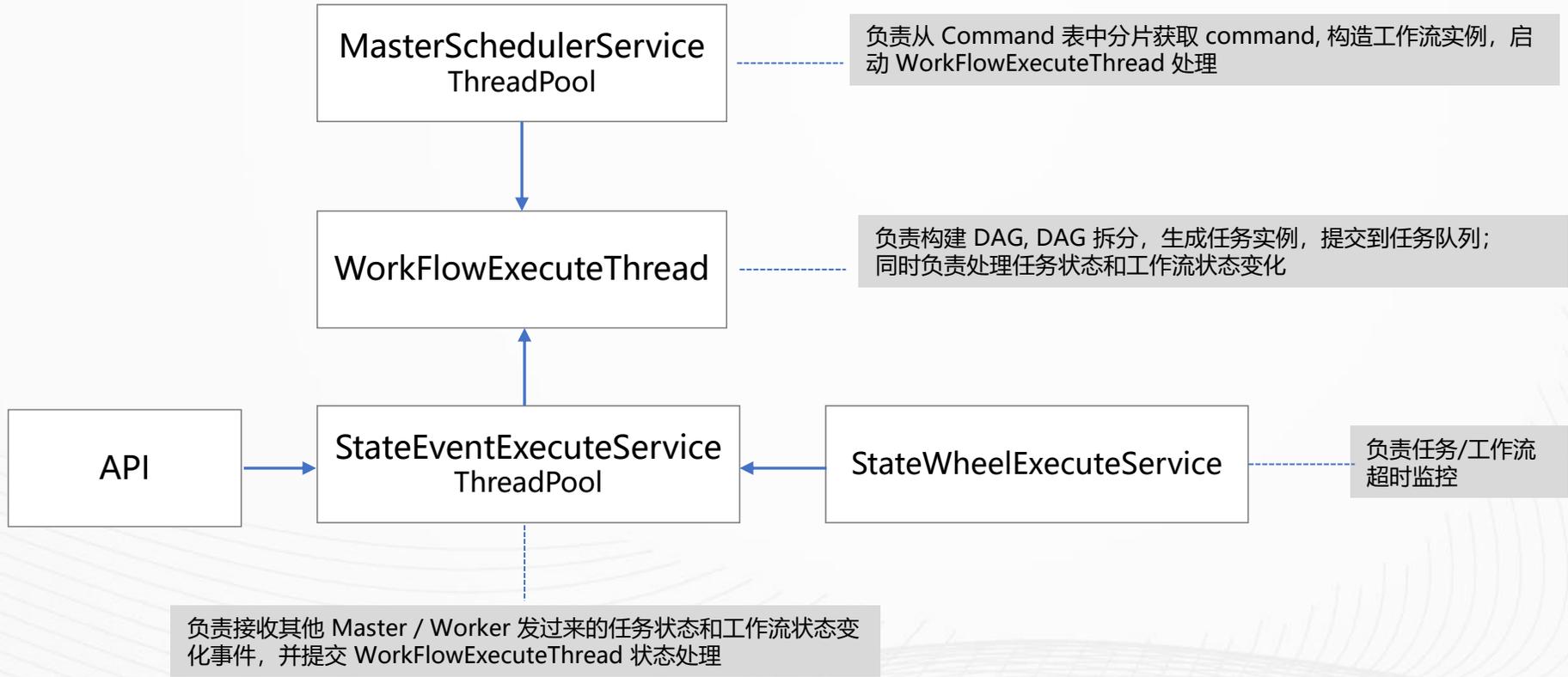
1.3 Master 现状

- 使用一个主线程池执行所有 DAG workflow
- 主线程池执行过程中为每个 workflow 创建一个任务线程池
- Master 能并发处理的工作流为 m , 并发处理的任务为 n , 会产生 $m * n$ 个线程



DolphinScheduler 2.0 架构 Roadmap

重构 Master 中的线程模型



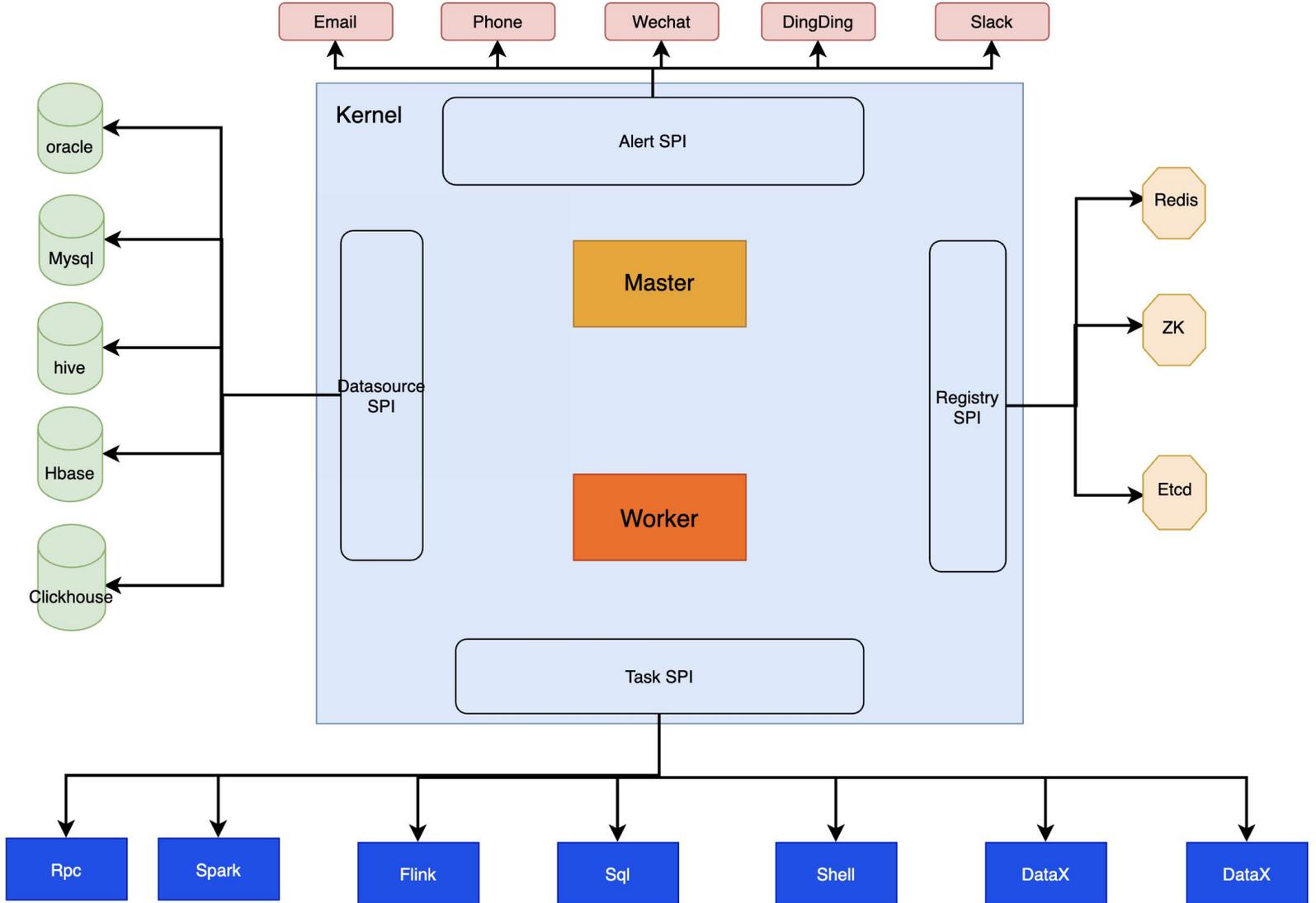
其他重要变化

新增 API 与 Master、Master 和 Master 之间直接通信功能, 负责同步任务及工作流状态的变化

DolphinScheduler 2.0 架构 Roadmap

所有扩展点都采用插件化实现

- 告警 SPI
- 注册中心 SPI
- 资源存储 SPI
- 任务插件 SPI
- 数据源 SPI
-



DolphinScheduler 2.0 架构规划与 Roadmap

支持更多数据组件集成

目前的 10 多种任务类型可能不能满足需求

解决办法:

任务 SPI 化
后续实现热插拔



创建插件

* 插件名称

资源名称

文件名称

描述

插件配置

* 上传资源



有个任务队列设置的页面，
这个页面可以设置多个队列，
每个队列里可以设置任务并发数，
每个任务可以指定他所属的队列

页面爆炸
只关注自己相关的任务

Roadmap

- 2.0 架构改造完成
- 任务血缘分析智能生成 DAG 依赖
- 提供 Python SDK
- 数据质量任务
- 容器调度任务
- CI、CD 等其他系统集成

长期 Roadmap

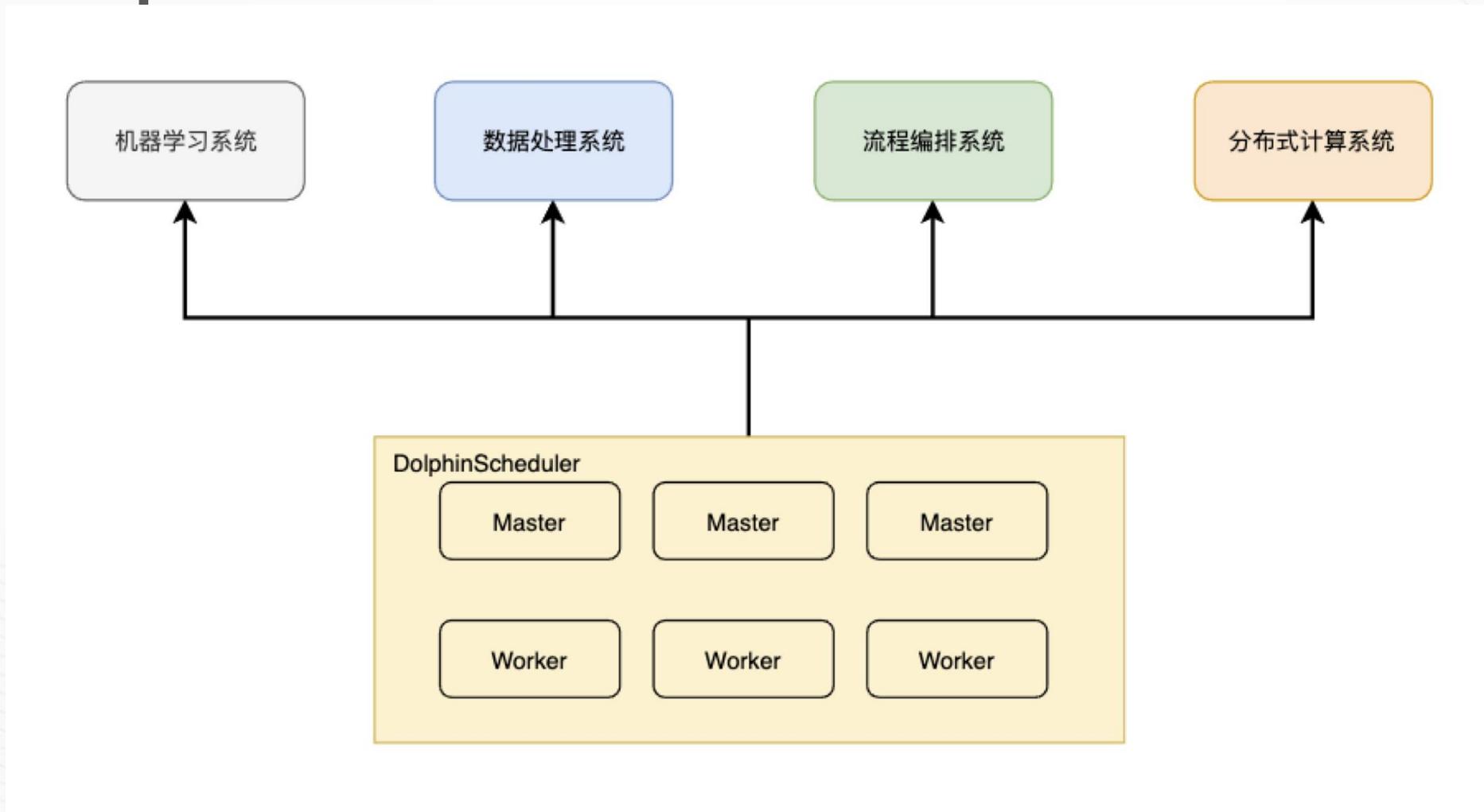
系统更稳、速度更快

GOTC

- 核心依赖组件不可用时，系统仍然能够完成关键调度工作
- 毫秒级超低延迟
- 更高的调度吞吐量
- 智能化运维

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE



使用场景 – 360 的应用

任务失败重试/告警
你想重试几次 每次间隔多久 失败要不要发邮件告诉你?

各种复杂调度

定时调度、依赖调度、手动调度

丰富的任务类型

spark shell MR HIVE python...

可视化

拖拽生成 DAG

分布式易扩展

无单点问题
资源不够了要扩容

资源文件的在线上传, 管理
jar 包不怕丢

实现集群高可用
集群去中心化

支持多租户
咱俩不能用一个账号

权限管理
我只能访问授权的项目和资源

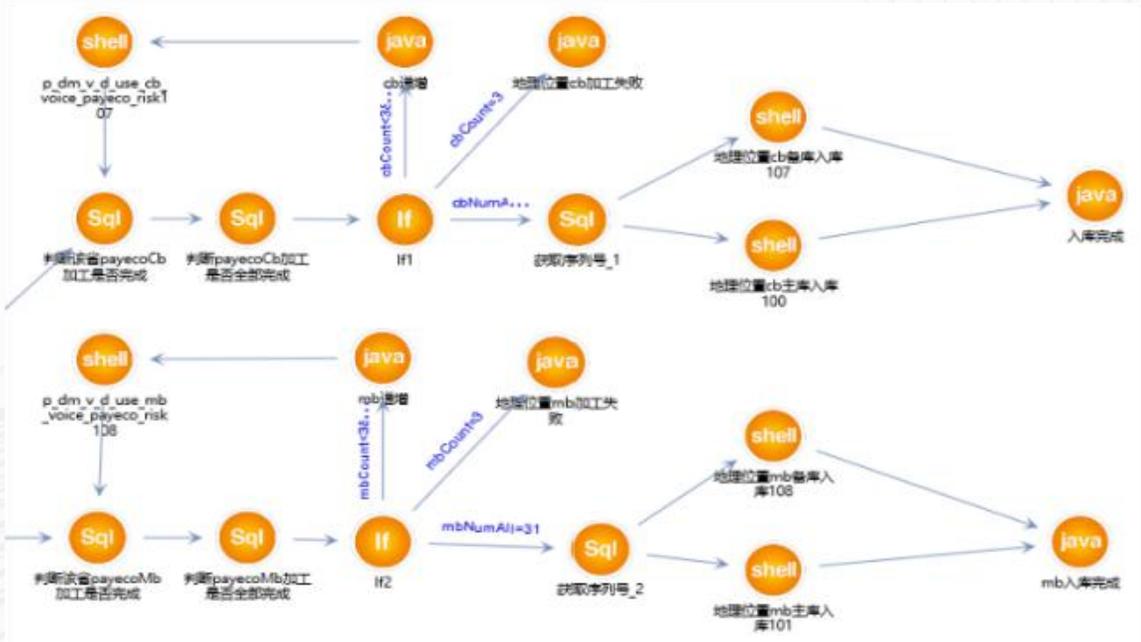
workflows

使用场景 - 联通的应用



目的：集中调度Spark、Hive、Pig、MR job、存储过程、shell等资源，支撑跨系统、跨应用、跨语言的作业任务协同执行，实现调度全局化、透明化

- 复杂多样的作业定义与调度策略：包括父子流程、并行、串行、依赖、和干涉等方式
- 资源负载均衡，实现高效调度：按资源对任务调度分组、并发控制、优先级动态调整

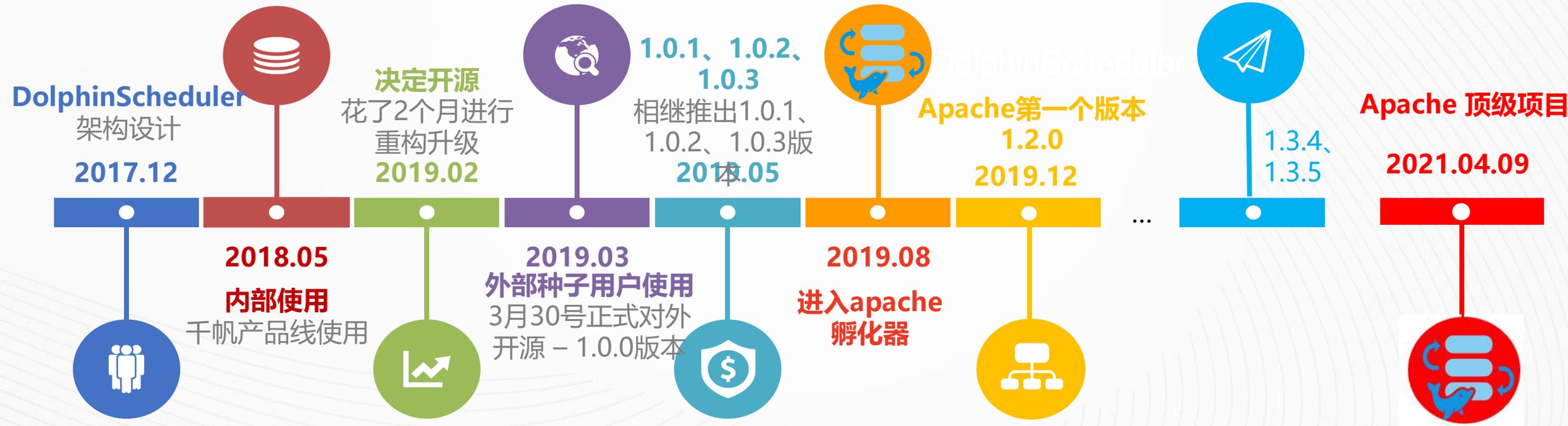


	Airflow	DolphinScheduler
二开成本	Python开发成本较高	Java, 流程清晰, 代码简洁, 成本相对较低
社区活跃度	比较活跃	比较活跃
集群扩展	需重启worker节点	可随时在zk中注册和删除 master、worker, 操作简单, 稳定向高, 可扩展性强
单点故障	主从模式, 存在一定的单点风险	多 master 多 worker, 注册 zk, 一个节点挂了不影响其他节点
过载处理	任务量大, 调度的性能会急剧下降, 甚至卡死服务器	任务过多会依靠队列排队
作业编排	流程定义需要编写Python脚本	支持任务拖拽配置, 本土化亮点突出, 用户交互友好
节点类型	BashOperator、PythonOperator、SSHOperator、HiveOperator, DockerOperator, OracleOperator, MysqlOperator, DummyOperator, SimpleHttpOperator、自定义节点	Shell 节点、子流程节点、存储过程、SQL 节点、SPARK节点、MR、Python 节点、依赖节点、HTTP节点、自定义节点
暂停与重跑	否	支持暂停流程、支持任意节点重跑
全局变量	支持全局变量, 项目内所有流程可见, 缺乏变量权限控制	暂不支持

使用场景 – 联通的应用

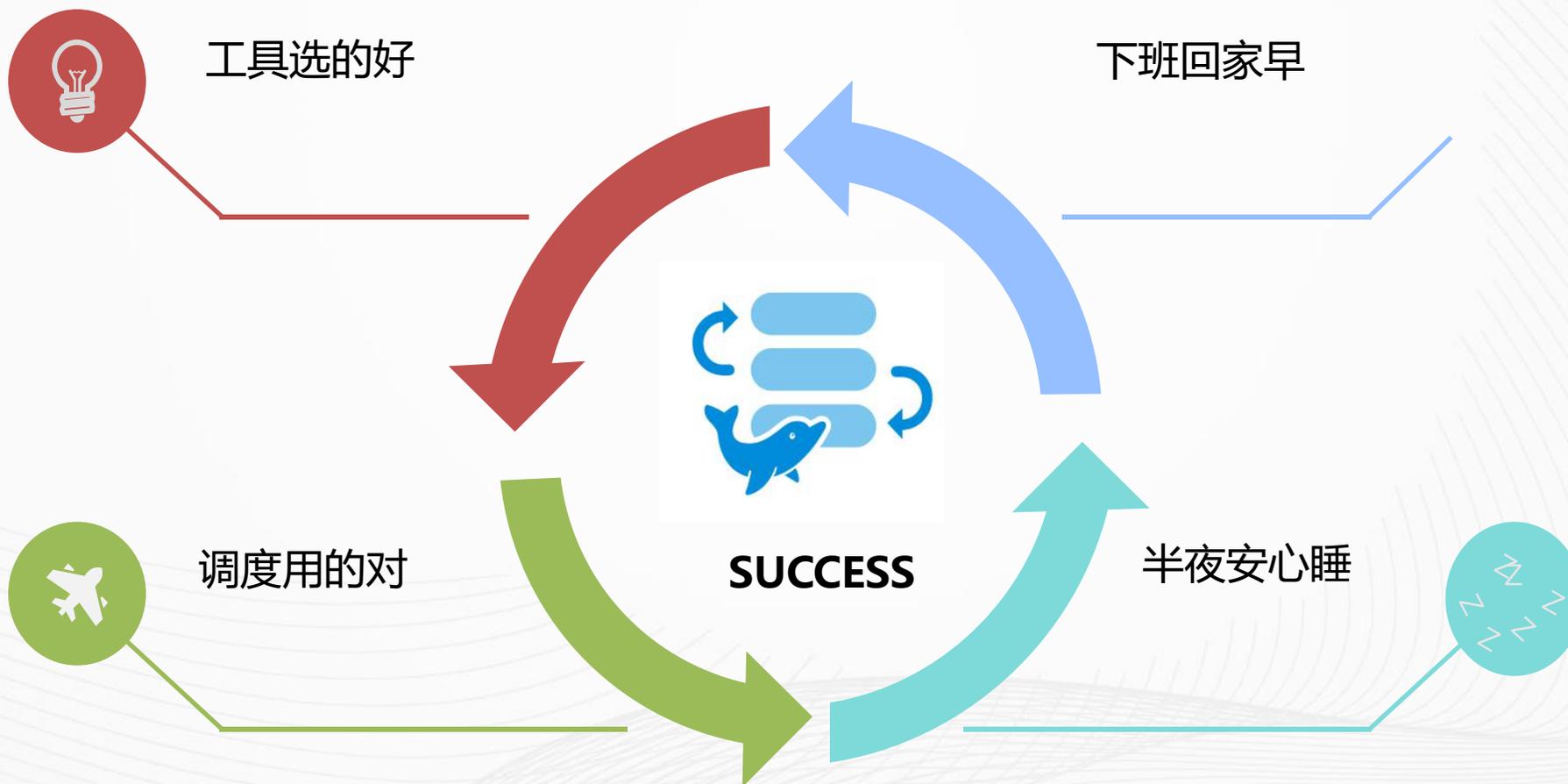
	节点类型	控件功能
	sql节点	从mysql获取参数 得到某个流程子任务的执行结果，会搭配条件判断和延迟节点做循环使用
	Java节点	参数的处理（List转换成map、参数拼接） 遍历list 调用子流程 对代理机的负载均衡
	条件判断节点	根据变量值调用不同节点
	延时节点	延时一段时间

DolphinScheduler 部分用户(不分先后)



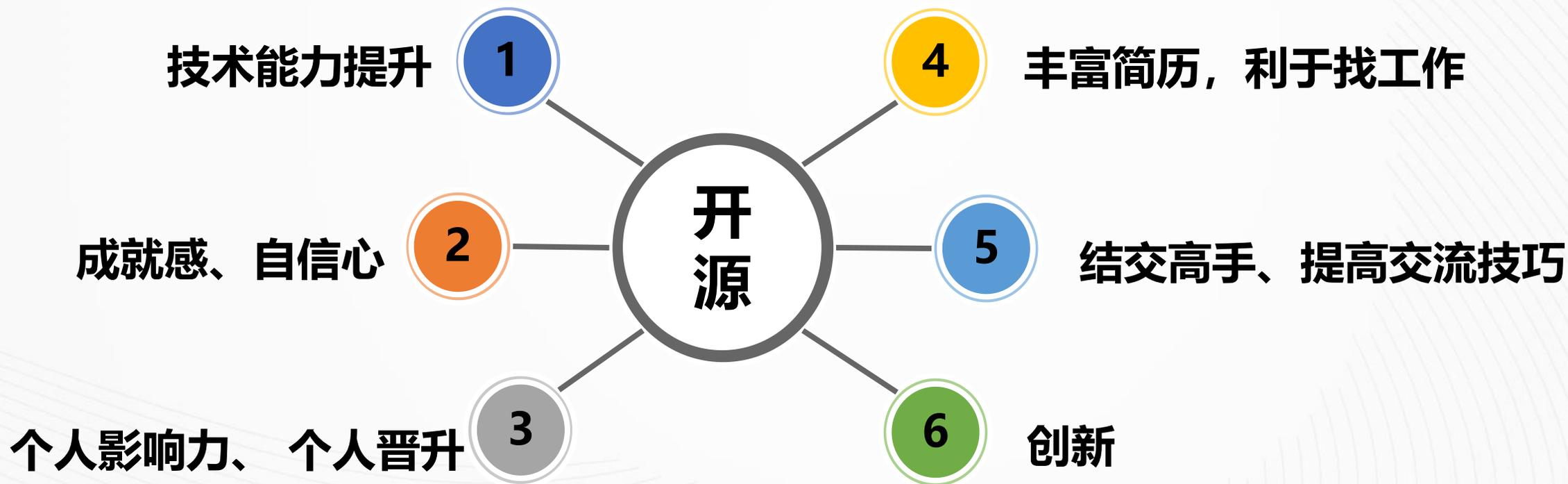
DolphinScheduler Slogan

GOTC



全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE



如何参与开源

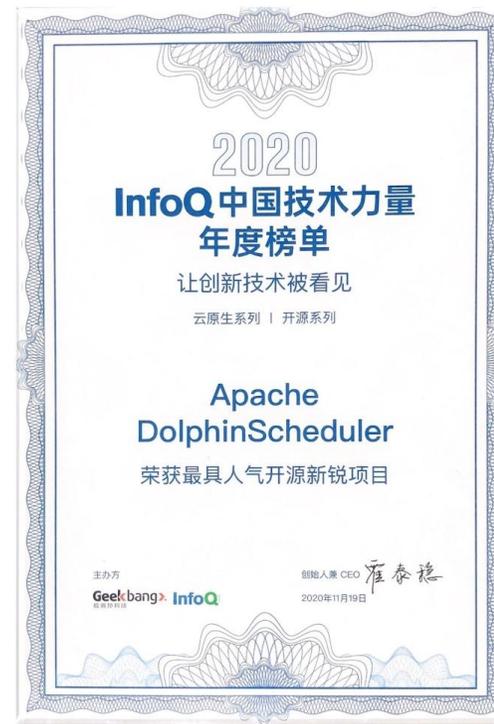
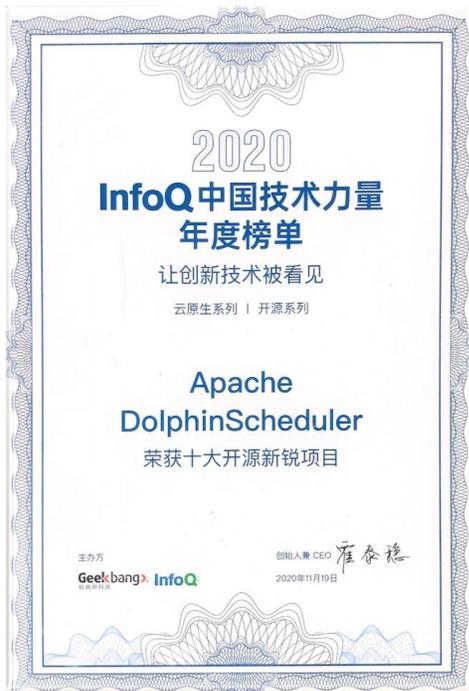
DolphinScheduler 社区参与贡献的方式，包括：

✓ 文档 ✓ 翻译 ✓ 答疑 ✓ 测试 ✓ 代码 ✓ 实践文章、原理文章等 ✓ 会议分享等

欢迎加入贡献的队伍，加入开源从提交第一个 PR 开始，

- 找到带有 easy to fix 标记或者一些非常简单的issue，修复后 [提交PR](#)

DolphinScheduler 2020年荣誉



2020 中国开源年度报告.pdf

我们计算了所有源自中国的共计 21 个 ASF 项目仓库的活动情况，数据如下。

#	name	language	activity	developer_count	issue_comment	open_issue	open_pull	pull_review_comment	merge_pull
1	apache/shardingsphere	Java	2858.72	786	9332	1581	3234	1851	2990
2	apache/incubator-echarts	JavaScript	2307.44	1183	7622	962	301	328	230
3	apache/skywalking	Java	1958.64	525	6904	892	838	3528	725
4	apache/incubator-dolphinscheduler	Java	1929.41	458	9439	1138	1440	722	1135
5	apache/dubbo	Java	1573.76	668	2537	584	558	165	329
6	apache/apisix	Lua	1539.91	300	5874	1030	1000	3363	837
7	apache/incubator-doris	C++	1221.35	149	2327	1055	1410	3543	1227
8	apache/hadoop-ozone	Java	1156.23	70	4270	0	1343	4547	1066
9	apache/incubator-iotdb	Java	996.1	86	4029	116	1485	3467	1268
10	apache/carbondata	Scala	853.36	50	8906	0	518	4767	1
11	apache/rocketmq	Java	825.41	280	1965	396	254	205	106
12	apache/servicecomb-java-chassis	Java	491.5	126	1393	272	326	330	308
13	apache/incubator-brpc	C++	451	170	847	176	70	109	37
14	apache/kylin	Java	416.18	60	1158	0	495	252	374
15	apache/incubator-weex	C++	291.44	146	449	100	46	8	32
16	XiaoMi/pegasus	C++	222.25	18	161	54	152	1065	137
17	apache/incubator-tubemq	Java	214.9	26	251	0	340	82	309
18	apache/incubator-teaclave	Rust	144.51	24	290	55	180	38	177
19	apache/griffin	Java	48.62	15	114	0	21	20	0
20	apache/hawq	C	11.02	5	13	0	3	1	2
21	apache/eagle	Java	2.65	1	1	0	2	0	0

DolphinScheduler 资源

- 官网: <https://dolphinscheduler.apache.org>
- Slack 群: <https://s.apache.org/dolphinscheduler-slack>
- B 站: <https://space.bilibili.com/515596012>
- 公众号: 海豚调度
- Twitter: @dolphinschedule
- 微信群助手: easyworkflow



海豚调度

微信扫描二维码, 关注我的公众号

GOTC

THANKS

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE