Blink资源配置的三年双11(2016~2018)回顾

资源配置指的是配置作业运行时各个阶段的资源,包括并发度,CPU,内存(分为heap memory, native memory,

directory memory)等资源,让作业能在最少的资源占用下达到最优的性能。资源配置的太多会造成浪费,配置的太少

● 代码中设置,如:topologyBuilder.setBolt("2", new WordSpliter(), 4)

版本加入社区的 setResource API接口(FLINK-5133)。通过这三种方式就能从粗粒度到细粒度来设置作业的资源。

接下来,通过回顾2016~2018这三年双十一来说下Blink资源配置的演变历程。

2016: 文本配置 2016年的双十一,是Blink的第一年双十一,成功支持了主搜索离线计算,以及搜索和推荐全链路实时化。采用的Blink

通过Java和Scala开发的Blink程序,如果把资源都通过API的方式写死在代码里面,显然不能满足快速调优和业务迭代的

版本是 0.8.x, 任务的类型以DataStream和TableAPI为主。

需求。因此业务方很多是自己定义了一套业务相关的配置文件,比如配置 taobao_source1, wireless_source2 的cpu和并 发度,然后在代码里面解析后设置到对应的节点上。这样造成的问题是,每个业务方都得定义自己的一套格式,进行解

析后设置,并不通用,重复工作较多。 因此Blink提供了通过JSON格式资源文件的方式来配置资源的方式。具体做法是:

• 通过 compile 命令,将代码预编译,输出各个transformation节点和默认资源配置到json文件中 • 用户修改json文件,可以设置各个transformation节点的并发度,cpu, memory等 • 用 run 命令提交任务的时候通过 -conf 指定这个资源配置的json文件 JSON文件的格式类似于: "nodes": [

"id": 1,

- "uid": "1",
- "pact": "Data Source", "name": "Source: CustomTypeTT4TableSource-wireless_search_pv", "slotSharingGroup": "default",

"chainingStrategy": "HEAD",

```
"direct_memory": 64,
      "native_memory": 0,
      "state_size": 0
    },
      "id": 2,
      "uid": "2",
      "pact": "Operator",
      "name": "SourceParser:EternalTTSourceCollector",
      "slotSharingGroup": "default",
      "chainingStrategy": "ALWAYS",
      "parallelism": 256,
      "vcore": 0.2,
      "heap_memory": 128,
      "direct_memory": 32,
      "native_memory": 0,
      "state_size": 0
    },
JSON资源配置文件的优点:
 ● 统一了资源配置的格式,业务方不需要在代码里面耦合资源设置的代码了,只需要关注业务逻辑就行了,减轻了很
   多工作量。
但同时也引入了新的问题:
 • 复杂任务的JSON资源文件长达几千行,有几百个transformation,每次改动都是个艰巨的工作
 • transformation节点和真正的物理执行节点的对应关系不是很明确,有时候需要多次反复修改才能确定
还记得2016年双十一,给搜索的珠峰项目做压测和调优的时候,和同事一起通过 vim 编辑json文件,批量替换,反复修
改定位节点,耗时几小时才能完成一个稍复杂任务的调优。
```

较自然的。

SQL任务的资源的方式。因为SQL任务的底层最终也是翻译成transformation,所以提供JSON文件的资源配置方式是比

具体使用流程是:用户在Bayes平台上写完SQL语句,通过『获取资源配置』就能生成一份默认的JSON资源文件。然后

Data Source(ID=1) Data Source(ID=5)

Source:

CustomTypeTT4TableSource-

wireless_search_pv

HEAP MEMORY: 1024 NATIVE MEMORY: 0 STATE SIZE: 0 PARALLELISM: 512

UID:

NAME:

VCORE:

PACK:

STATE SIZE: 0

PARALLELISM: 256

Source:

Operator

CustomTypeTT4TableSource-

app_taobao_page_sys

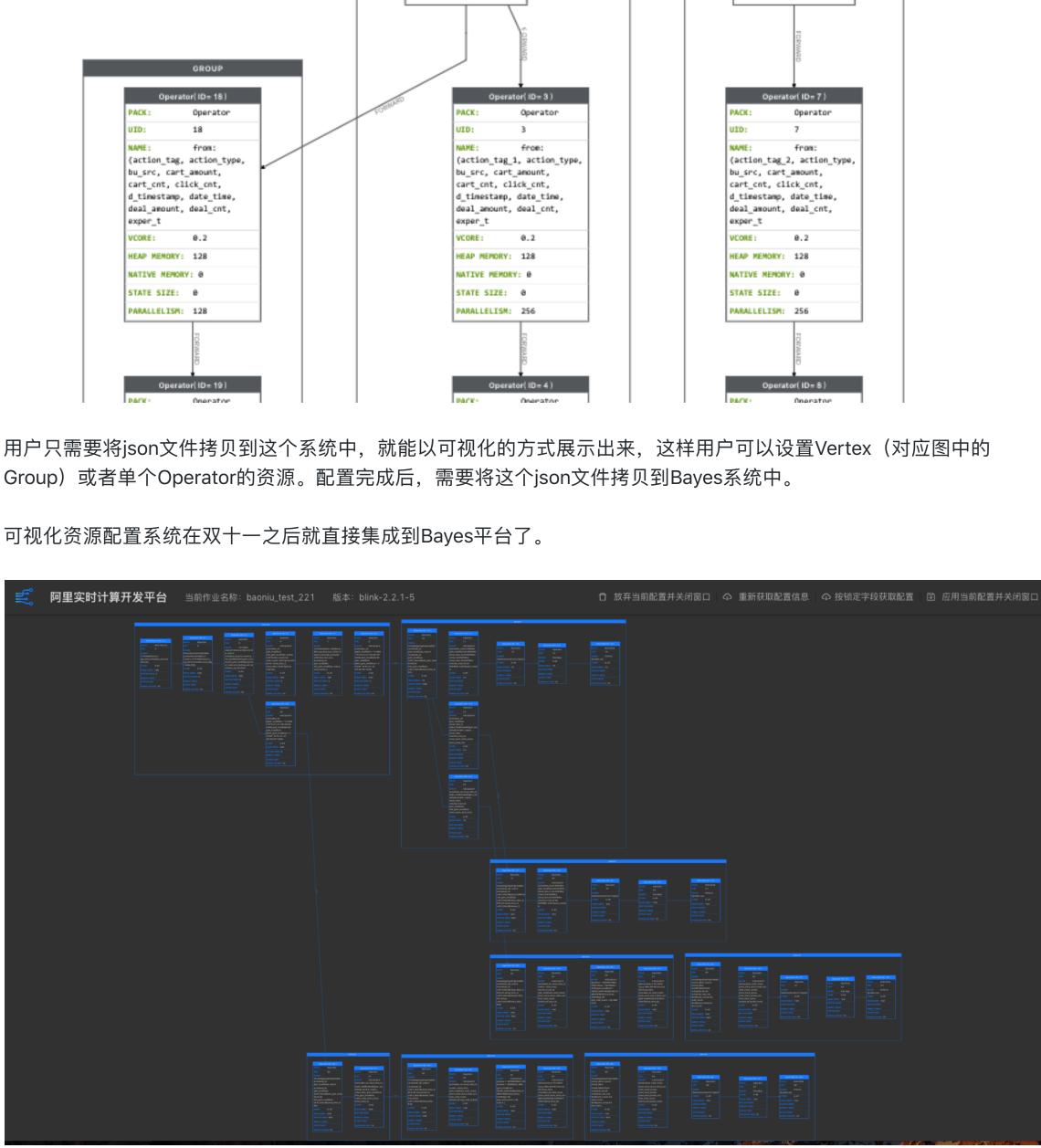
HEAP MEMORY: 256

PARALLELISM: 256

通过Bayes平台进行资源配置的优点是:流程更简单,资源配置生成和修改在一个界面,鼠标点击就能完成。

SourceParser:EternalTTSo SourceParser:EternalTTSo urceCollector urceCollector VCORE: VCORE: HEAP MEMORY: 128 HEAP MEMORY: 128 NATIVE MEMORY: 0

GROUP



经过一段时间的开发测试,Blink团队于2017年12月推出了1.4.5版本,提供自动配置(AutoConf)功能。其基本原理 是:采样收集任务运行时的指标(Metrics)信息,在下次运行前通过算法给出推荐的资源配置。 AutoConf 使用流程

息相关。Blink团队一直在积极探索更加智能的资源配置方式,希望能降低用户的工作量,同时能节省集群资源。

『资源配置』和『性能调优』一直是每年双十一工作的重点。性能调优不仅仅局限于资源配置,但是资源配置和性能息

2018: 半自动化配置

上线新版本

AutoConf 工作原理

同时基于以下假设:

其基本工作模式:

memory.

保存配置

() 输入表配置

压测作业(文档) 方案一

最大运行时间(分钟):

回追时间:

启用 AutoConf:

压测说明

60

2018-10-10 21:08:57 📋 🛛 💿

AutoConf 的工作原理主要基于以下几点:

用少数实验来推导出可用的配置。

• 并发度和资源的设置可以分开考虑

到的 statistics 越准确,几次后会收敛。

● 在流计算任务中基于稳定状态的 statistics 可以得到足够的信息用来优化配置。

● 每个 task 的平均输出 QPS 和输入 QPS 之间的比例基本不变。

资源配置

数据检查

下一步

上一步

取消

(可用 0.33 CU)

X

完成

Х

上线作业

● (1) 用户指定任务使用多少 CU 资源(1 CU对应 1个CPU核, 4G内存), AutoConf算法会根据历史运行情况计算 生成一份资源配置 如果任务是第一次运行,之前没有历史记录,这时候算法会根据经验值生成一份初始推荐配置。 。 注意:初始值只是根据经验值算出来的,不一定准确,甚至可能和实际情况相差较大(例如用户的source源 partition数很大)。需要用户多进行几次调优才能得到比较合理的配置。 (2) 使用生成的配置运行,在运行期间Blink会采集一些指标信息(比如:输入输出tps, latency,cpu,memory等) • (3) 待任务运行一段时间后, 重复步骤(1) 一般通过几轮迭代能得到一个较优的资源配置。 上线新版本 (2)数据检查 (3)上线作业 资源配置 资源配置方式:
 智能推荐配置(10.00 CU 可用): 指定使用 8 CU ② 使用上次资源配置(手动资源配置) ②





√ 输出表配置

迭代间隔(分钟):

作业最大CU数:

10

200

AutoConf 并不是万能的,算法本身做了很多假设,比如假设每个Task处理的数据是均匀的,每条数据处理的时间是一样 的等。此外,对于SQL本身写法导致的性能问题,单纯通过AutoConf来调大资源使用,收益并不明显,很容易达到性能

新的配置, 也可以选择使用自己已经调好的配置。

配置等都在考虑中。未来可期!

展望

如果用户从SQL和业务代码层面去优化,性能的提升会更明显,因此Blink除了在引擎层面优化SQL执行计划,也提供了 专门的SQL优化文档来指导用户实践。

Blink团队还在对资源配置方式进行不断的优化中。可视化资源配置将更方便地支持具有复杂拓扑结构的任务。AutoConf

算法未来或许可以实现完全无需人工干预的智能配置,将 AutoConf 与 SQL 执行计划生成相结合,流量高低峰期自适应

瓶颈。例如: • SQL里面使用了大量的正则表达式 ● 用户的UDF实现有耗时较大的外部IO导致性能较差,调大并发度和CPU也很难达到预期性能。 ● SQL写法导致数据倾斜,如果不解决好倾斜问题,只是调大并发度效果也不明显。 源表分区数有限制导致无法扩大并发;维表或者结果表性能不足

『自动任务压测』的目的是为了解放用户,让整个调优和压测的过程更加自动化,用户可以选择启用AutoConf每次生成

性能又不够。此外,资源配置的方式如果太复杂,对用户来说成本也比较高。 如何既方便快捷,又合理的配置资源对各种大数据处理系统来说一直都是一个难点。 目前业界的各种大数据处理引擎,比如Hive,JStorm,Spark一般都提供了如下几种方式来配置资源: • 默认配置文件, 如 storm.yaml • 提交作业的命令行参数, 如 ./bin/xxx-submit --executor-memory 1G --num-executors 1 --conf xxx.sql.shuffle.partitions=100 Flink中也提供了类似的方案: (1)flink-conf.yaml配置 (2)./bin/flink run命令行 (3)以及由阿里的同学在Flink 1.3

"parallelism": 512, "vcore": 1, "heap_memory": 1024, 2017年的双11, Blink开始通过阿里实时计算平台Bayes面向全集团用户,任务类型以SQL为主,版本为1.4.x。社区的 Flink SQL本身并没有提供资源设置的方式,Blink当时开发了SQLRunner/JobLauncher,提供通过JSON资源文件配置

2017: 可视化配置

用户可以在Bayes页面上直接修改json文件,或者拷贝出来修改,保存后上线和运行。

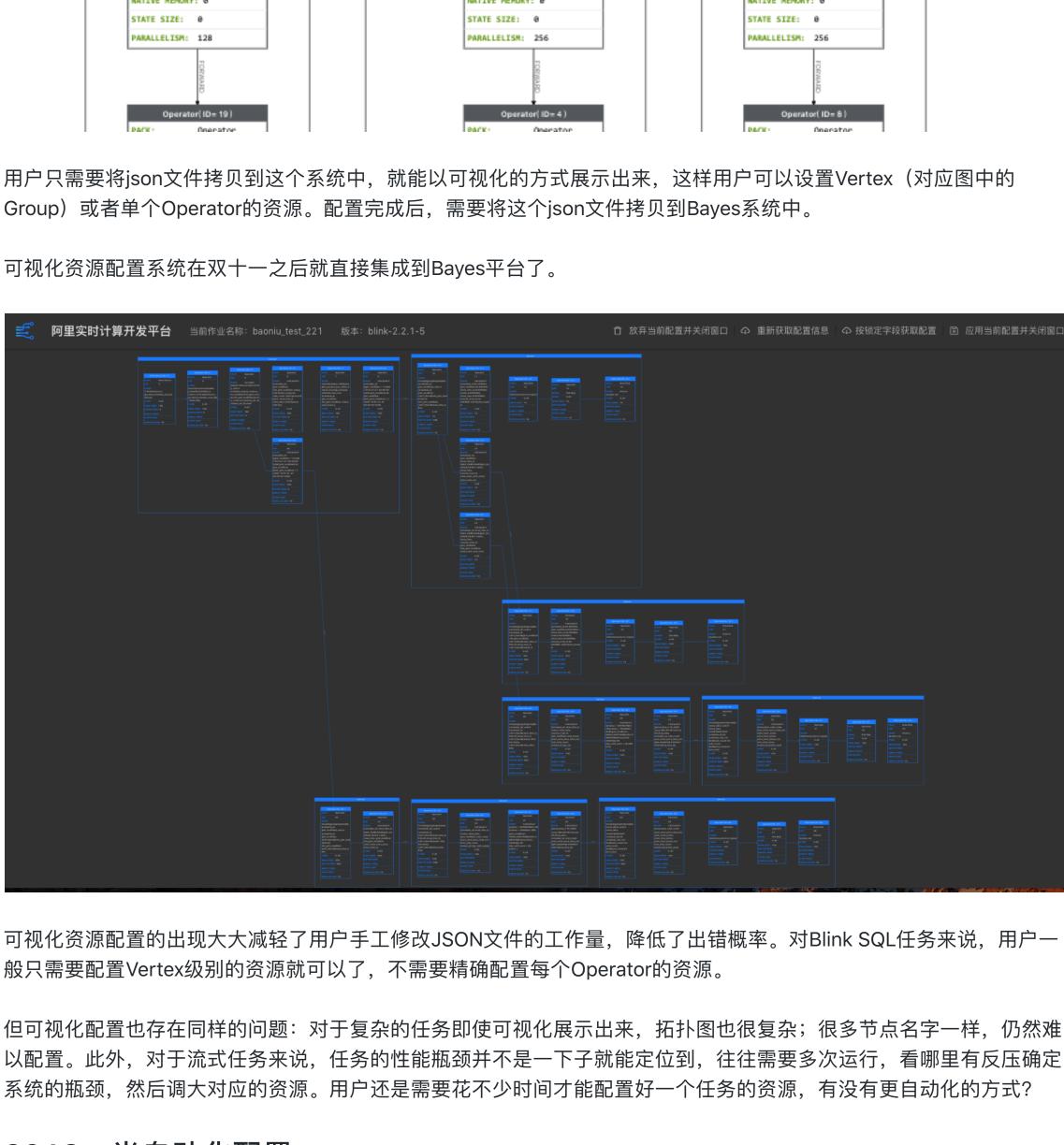
但缺点也一样存在:对于长达几千行的复杂json资源文件,每次修改都很繁琐,费时费力。

NAME:

VCORE:

因此当时开发了一个『可视化资源配置』的功能:

Operator(ID=2) PACK: Operator NATIVE MEMORY: 8 STATE SIZE: 0 PARALLELISM: 256



上线检查通过,预计消耗CU 7.25 个

● Blink 任务一次实验(修改配置+运行+计算 statistics)的代价(时间+资源+用户操作+任务影响)很高,所以尽量

● 在资源配置足够的情况下,一个节点单个并发度(task)的一个 event 平均净处理时间和固定 overhead 基本不变。

● 一开始配置略宽松的资源让任务能够正常运行(不需要达到目标流量),以便拿到 metrics 来计算 statistics。

• 节点并发度的比例关系可以通过workload的比例算出;节点的资源可以通过metrics得到,比如cpu, heap

从稳定状态的 statistics 可以得到比较准确的资源配置。稳定状态越接近目标状态(数据流量,并发度,资源)得

accessid='xxx',

accesskey='xx

topic='xx type='tt'

lengthcheck='pad',