

背景

前几天知识星球有人提问说：

星主你好，请教个问题，yarn 显示 flink job 还在运行中,但是实际上已经失败了,这个怎么判断 flink job 确实挂了呢?因为我需要监控 flink 程序是否正常在运行，好告警。

莫忘 提问：星主你好，请教个问题，yarn显示flinkjob还在运行中,但是实际上已经失败了,这个怎么判断flinkjob确实挂了呢?因为我需要监控flink程序是否正常在运行，好告警，星主有什么好点的方式么？

Job状态如图

app	Flink	Apache	root:spark	0	Tue May 19	Tue May 19	N/A	RUNNING	UNDEFINED	1
实时	Flink				19	13:41:17				
计算					13:41:17	+0800	2020			
提交					2020					
结束					2020					

Start Time	Duration	End Time	1
2020-05-19 13:41:43	2s	2020-05-19 13:41:45	

依赖yarn的状态会被坑死，我就因为历史包袱导致因为这个背了几个故障

真实建议查到作业的task级别的状态去判断整个作业的运行状态（PS：有时候作业的状态为running，但task也有fail的），可以在控制台去看下这个接口



[查看详情 >](#)

ithan 觉得很赞

zhisheng：等这周末我写篇我们利用 flink UI 上的 rest api 做的一些工作
2020/5/19

如果你没遇到这个问题，那么说明你遇到的场景还比较少，还得多历练；

如果你也遇到这个问题，那么我们握握手，以示相互安慰，都是同命苦的小伙伴。

这里不得不又得提一下我就因为这个问题导致背过两个生产事故的故障。这个问题确实会在 Flink 里面出现，起码我现在碰到的次数就好些次了。

问题剖析

这个问题主要的意思是表示通过 YARN 看到作业 AppID 的状态是 Running 的，点进去查看 Flink UI 的状态也是 Running 的，但是 Flink 作业的 Task 其实是有挂（failed）了的。

The top screenshot shows a YARN application list with a table containing columns: ID, User, Name, Application Type, Queue, StartTime, FinishTime, State, FinalStatus, Progress, and Tracked. One entry is highlighted with a red box around the 'State' column, which shows 'RUNNING'. The bottom screenshot shows the Flink UI for a job named 'kafka_decouple_metrics_cluster'. It displays 'Available Task Slots' as 0 and 'Running Jobs' as 1. Below this, a 'Running Job List' table shows the job's status as 'RUNNING'.

The screenshot shows the Flink UI for a job named 'kafka_decouple_metrics_cluster'. The job is in a 'RUNNING' state. Below the job overview, there is a 'Task Status' table with columns: Name, Status, Bytes Received, Records Received, Bytes Sent, Records Sent, Parallelism, Start Time, Duration, End Time, and Tasks. The table shows several tasks in a 'CANCELED' state, indicating a failure in the job.

上面的问题可以用上面三个图来描述，总而言之就是三者的状态不同步。

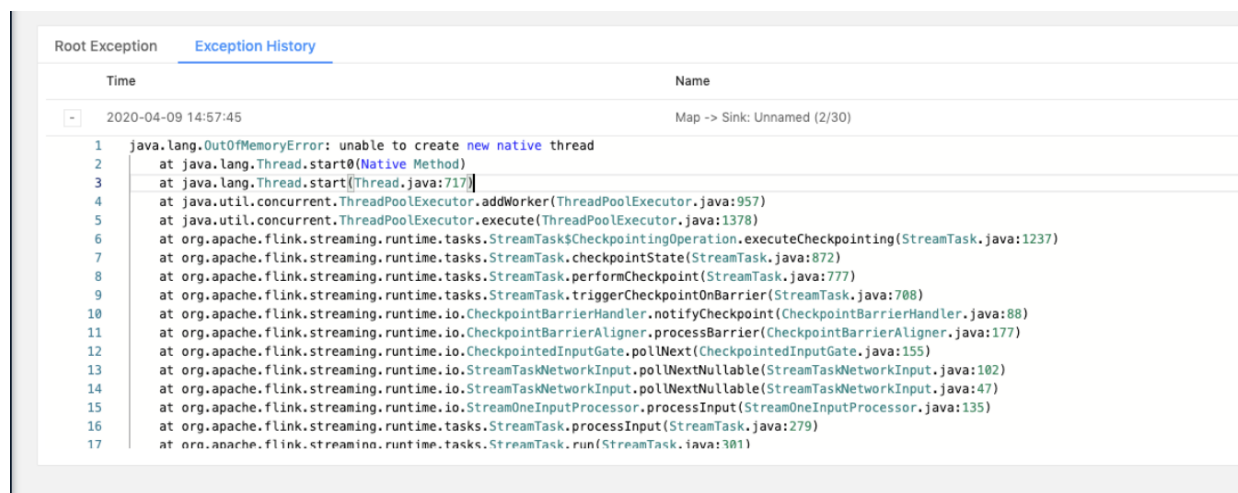
问题影响

如果对作业的状态进行告警只关注 Yarn 或者 Flink 显示的状态，那么自然会错漏掉这种情况，从而导致最终这个作业挂了而一直没被发现。作为作业的 Owner 来说，肯定不能接受公司实时计算平台针对这种情况就不告警了。当遇到大故障的时候，如果对于这种作业不能及时的发现，那么就很可能会有很严重的故障，尤其是这个作业还是非常重要的情况下（数据延迟不能超过公司制定的故障级别）。所以，我就因此问题背过两次该问题的故障，两次相隔没几天，都没来得及修复上线。

复盘一下这两次出现的故障：

- 第一次是因为一个业务方在 IDEA 测试完功能后就直接去生产测试，因为测试环境无数据，他的作业里面引入公司的一个日志中间件，然后他的作业代码里面有个 HBase 的依赖和公司的日志

中间件的依赖有冲突，在消费 Kafka 的数据后写到 HBase 之前会利用公司的日志中间件把一份数据打到另一个 Kafka Topic 里面，因为冲突的问题，导致中间件频繁的连接，且数据消费的模式是 startFromEarliest，最后导致创建非常多的线程，然后这个作业 Task 所在的机器都遭殃了，所有的机器线程都被打满了，导致这些机器上面部署调度到的其他作业 Task 也都受到了影响（这个影响是导致 Checkpoint 失败，如下图所示）。



同时伴随着部分作业会一直做 failover，当找到这个导致故障的根本作业后（这个找的技巧后面可以单独说），直接把那个作业 kill 掉了，后面慢慢部分作业的状态都 OK 了，但是就有上面提到的这种作业（状态是 OK，但是 Task 其实已经 failed 的）。针对这种如果没及时发现，就会酿成大故障，在此也提醒大家遇到故障后要及时梳理全排查受影响的范围和严重性，切勿匆匆打扫战场！！

- 第二次是在上一次故障的后几天，果然验证了那句老话：祸不单行！！！故障也是类似问题，只不过这次是换成写 Redis 了，也是创建了大量的线程导致的。

总结一下两次故障：

因为作业的问题，导致机器线程打满，从而影响了机器上其他的 Flink 作业。暴露的问题主要有两点：

- 监控不够深，对于 Task failed 的未及时告警
- 作业未做到资源相互隔离（导致后面强烈要上 Cgroup）

解决方法

这里 Cgroup 的解决方案就不在这篇文章讲了，其实公众号前几天已经发过一篇文章。

这里讲一下，作业要做到 Task 级别的告警方法（主要利用 Flink Rest API）。

你可以通过该 https://xxx/application_1575453055442_4323/jobs/overview 接口拿到作业的状态和 Task 状态信息，如下两图所示：

The screenshot displays a web application interface with a sidebar on the left containing navigation links: Overview, Jobs, Running Jobs, Completed Jobs, Task Managers, Job Manager, and Submit New Job. The main content area is divided into two sections: 'Available Task Slots' and 'Running Jobs'.

Available Task Slots: Shows 0 slots available. Below this, it indicates 'Total Task Slots 320' and 'Task Managers 320'.

Running Jobs: Shows 1 job running. Below this, it indicates 'Finished 0', 'Canceled 0', and 'Failed 0'.

Running Job List: A table with columns: Job Name, Start Time, Duration, End Time, Tasks, and Status. It lists one job with the following details:

Job Name	Start Time	Duration	End Time	Tasks	Status
[Redacted]	2020-05-09 11:55:19	13d 7h 51m 44s	-	344 344	RUNNING

Below the job list, there is a network log section. The 'Network' tab is selected, showing a list of requests. The first request is highlighted, showing its details:

- Name: overview
- Request URL: https://[Redacted]application_1575453055442_4323/jobs/overview
- Request Method: GET
- Status Code: 200
- Remote Address: 47.114.8.159:443

The 'Preview' tab is also visible, showing the JSON response of the request. The response is a nested object with 'jobs' and 'tasks' arrays. The 'jobs' array contains one job object, and the 'tasks' array contains one task object. The job object has the following properties:

- jobs: [{jid: "31cfece1f50c6d004d439777e72ea6e3", name: "HDSpaceX-core", state: "RUNNING",...}]

The task object has the following properties:

- tasks: {total: 344, created: 0, scheduled: 0, deploying: 0, running: 344, finished: 0, canceling: 0, canceled: 0, failed: 0, reconciling: 0, finished: 0, last-modification: 1590115403638, start-time: 1588996519067, end-time: -1, duration: 1151459038, jid: "31cfece1f50c6d004d439777e72ea6e3", name: "HDSpaceX-core", state: "RUNNING",...}

Red arrows point from the 'Job' and 'Task' labels to the corresponding objects in the JSON response.

返回的数据如下 JSON 所示：

```
1 {
2   "jobs": [{
3     "jid": "31cfece1f50c6d004d439777e72ea6e3",
4     "name": "xxx",
5     "state": "RUNNING",
6     "start-time": 1588996519067,
7     "end-time": -1,
8     "duration": 1151459038,
9     "last-modification": 1590115403638,
10    "tasks": {
11      "total": 344,
12      "created": 0,
```

```

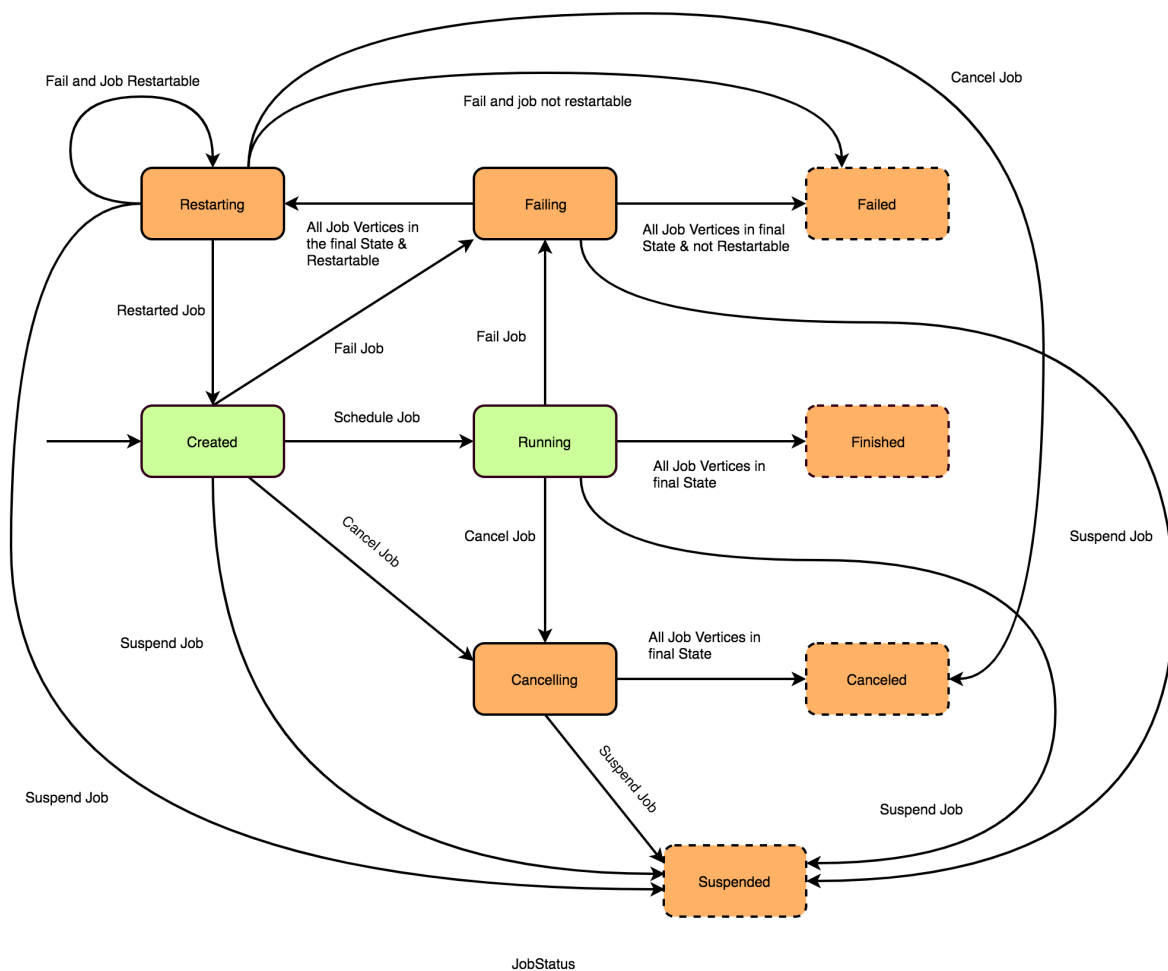
13         "scheduled": 0,
14         "deploying": 0,
15         "running": 344,
16         "finished": 0,
17         "canceling": 0,
18         "canceled": 0,
19         "failed": 0,
20         "reconciling": 0
21     }
22 }
23 }

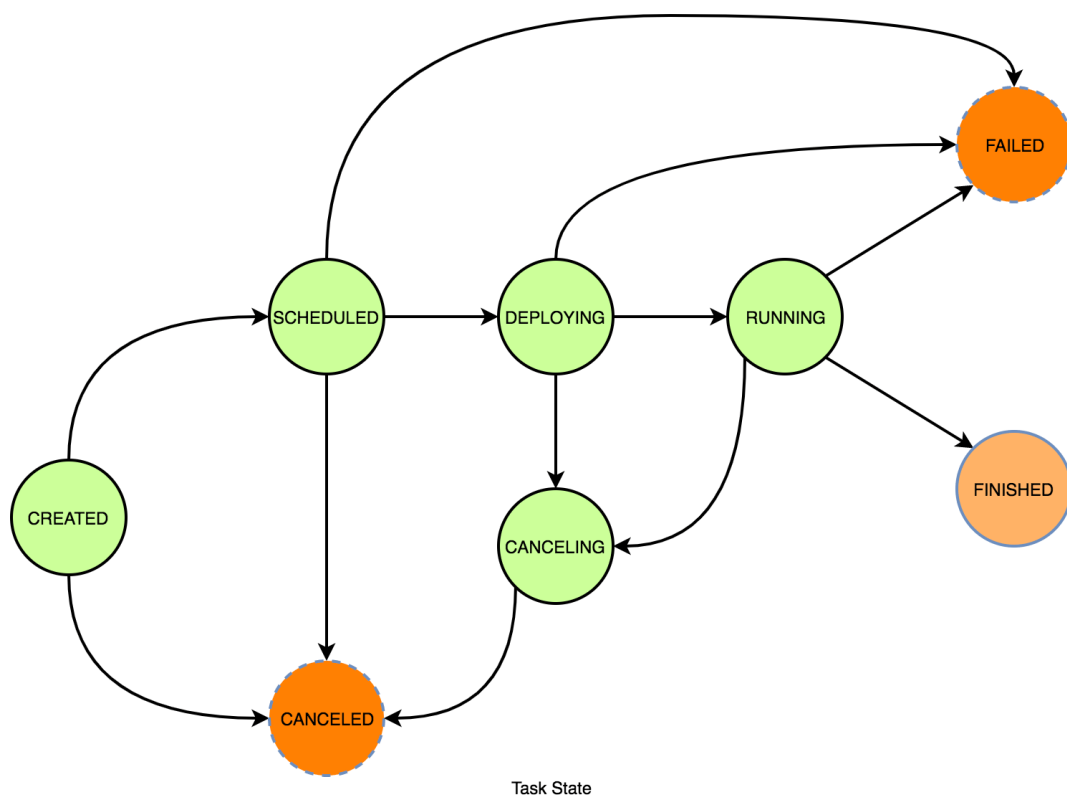
```

我这边检测的时候是会既判断作业的状态和 Task 的 running 个数是否和 total 个数一致，因为我们生产还有的作业是批和流一起用的，所以只判断 running 个数是否和 total 个数一致还不能判断作业真正的状态，还得再加一个判断是否 failed 的个数是大于 0。如果大于 0 则表示作业有 Task 失败，然后会给作业的 Owner 告警。

加了这个检测后，作业的状态告警更准，更符合我们的需求，另外我还写成了接口统计所有异常状态的作业，方便故障后的及时梳理工作。

这里再拿出我之前画的关于作业和 Task 的运行状态变换图：





除了利用上面的这个接口去做作业状态的告警，另外我们还利用了下面这些 API 来做一些日常的统计。

统计大状态作业

我们生产中也有不少的大状态作业，刚接手实时平台底层的负责的时候，碰到的就是这种大状态作业，做 Checkpoint 会有很多问题，我们就想每隔一段时间统计一下集群的这种大状态作业，然后我这边自己有个心理准备，以防 Checkpoint 的时候打爆机器网卡带宽（也碰到过这种故障）。

ID: b730c1591a99462abf76497493881655 | Start Time: 2020-01-02 15:44:40 | Duration: 75d 19h 33m 19s

Overview Exceptions TimeLine **Checkpoints** Configuration

ID	Status	Acknowledged	Trigger Time	Latest Acknowledgement	End to End Duration	State Size	Buffered During Align
+ 10917	FAILED	43/60	11:11:11	11:15:05	5m 0s	51.8 GB	8.69 MB
+ 10916	FAILED	43/60	11:01:11	11:04:48	5m 0s	51.8 GB	4.84 MB
+ 10915	FAILED	43/60	10:51:11	10:54:41	5m 0s	51.8 GB	4.69 MB
+ 10914	FAILED	43/60	10:41:11	10:44:42	5m 0s	51.7 GB	3.91 MB
+ 10913	FAILED	43/60	10:31:11	10:34:50	5m 0s	51.7 GB	3.06 MB
+ 10912	FAILED	43/60	10:21:11	10:24:47	5m 0s	51.8 GB	6.47 MB
+ 10911	FAILED	43/60	10:11:11	10:14:47	5m 0s	51.7 GB	3.44 MB
+ 10910	FAILED	43/60	10:01:11	10:04:40	5m 0s	51.7 GB	5.50 MB
+ 10909	FAILED	43/60	09:51:11	09:54:49	5m 0s	51.6 GB	3.13 MB
+ 10908	FAILED	43/60	09:41:11	09:44:49	5m 0s	51.6 GB	4.38 MB

那么我们也是利用 Rest API 来检测的，帮助我们不少问题。

大状态作业

- [application_1575453055442_3698](#) 当前状态大小为：10.3GB，历史平均状态大小为：8487.5MB
- [application_1575453055442_4634](#) 当前状态大小为：2.2GB，历史平均状态大小为：830.3MB
- [application_1575453055442_3693](#) 当前状态大小为：1.6GB，历史平均状态大小为：2332.3MB
- [application_1575453055442_4640](#) 当前状态大小为：1.7GB，历史平均状态大小为：1328.7MB
- [application_1575453055442_4652](#) 当前状态大小为：3.7GB，历史平均状态大小为：2410.7MB
- [application_1575453055442_3603](#) 当前状态大小为：24.4GB，历史平均状态大小为：8000.3MB

上图已经是我们经过几次联系业务开发进行作业治理后的效果了，大状态作业明显降低了很多，针对这个大状态作业的优化问题也是比较复杂的，后面可以单独讲一篇。

你可以通过 `https://xxx/{application_id}/jobs/{jobid}/checkpoints` 接口查看到 Checkpoint 的信息:

The screenshot shows the Flink Checkpoint UI for a specific job. The job ID is `e47bc5a70b23f4e1576cfe38d02d8ca5`. The UI displays checkpoint counts: Triggered: 401, In Progress: 1, Completed: 400, Failed: 0, Restored: 0. The latest completed checkpoint is ID: 400, with a completion time of 20:12:44 and a state size of 4.86 GB. Below this, there's a table of operators. The browser's network tab is open, showing a GET request to `https://xxx/application_1575453055442_4868/jobs/e47bc5a70b23f4e1576cfe38d02d8ca5/checkpoints` with a status code of 200.

接口返回的数据如下:

```
1 {
2   "counts": {
3     "restored": 0,
4     "total": 401,
5     "in_progress": 1,
6     "completed": 400,
7     "failed": 0
8   },
9   "summary": {
10    "state_size": {
11      "min": 281401159,
12      "max": 7645131875,
13      "avg": 4184641844
14    },
15    "end_to_end_duration": {
16      "min": 5543,
17      "max": 47550,
18      "avg": 17686
19    },
20    "alignment_buffered": {
21      "min": 0,
22      "max": 0,
23      "avg": 0
24    }
25  },
26  "latest": {
27    "completed": {
28      "@class": "completed",
29      "id": 400,
30      "status": "COMPLETED",
31      "is_savepoint": false,
32      "trigger_timestamp": 1590149544874,
33      "latest_ack_timestamp": 1590149564924,
```



```

34         "state_size": 5220376693,
35         "end_to_end_duration": 20050,
36         "alignment_buffered": 0,
37         "num_subtasks": 192,
38         "num_acknowledged_subtasks": 192,
39         "tasks": {},
40         "external_path":
"dfs:/flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-400",
41         "discarded": false
42     },
43     "savepoint": null,
44     "failed": null,
45     "restored": null
46 },
47 "history": [{
48     "@class": "in_progress",
49     "id": 401,
50     "status": "IN_PROGRESS",
51     "is_savepoint": false,
52     "trigger_timestamp": 1590149574874,
53     "latest_ack_timestamp": 1590149579446,
54     "state_size": 3920009423,
55     "end_to_end_duration": 4572,
56     "alignment_buffered": 0,
57     "num_subtasks": 192,
58     "num_acknowledged_subtasks": 127,
59     "tasks": {}
60 }, {
61     "@class": "completed",
62     "id": 400,
63     "status": "COMPLETED",
64     "is_savepoint": false,
65     "trigger_timestamp": 1590149544874,
66     "latest_ack_timestamp": 1590149564924,
67     "state_size": 5220376693,
68     "end_to_end_duration": 20050,
69     "alignment_buffered": 0,
70     "num_subtasks": 192,
71     "num_acknowledged_subtasks": 192,
72     "tasks": {},
73     "external_path":
"dfs:/flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-400",
74     "discarded": false
75 }, {
76     "@class": "completed",
77     "id": 399,
78     "status": "COMPLETED",
79     "is_savepoint": false,
80     "trigger_timestamp": 1590149514874,
81     "latest_ack_timestamp": 1590149520970,
82     "state_size": 4678903399,
83     "end_to_end_duration": 6096,
84     "alignment_buffered": 0,
85     "num_subtasks": 192,
86     "num_acknowledged_subtasks": 192,
87     "tasks": {},
88     "external_path":
"dfs:/flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-399",
89     "discarded": true

```

```

90     }, {
91         "@class": "completed",
92         "id": 398,
93         "status": "COMPLETED",
94         "is_savepoint": false,
95         "trigger_timestamp": 1590149474586,
96         "latest_ack_timestamp": 1590149512772,
97         "state_size": 4269715948,
98         "end_to_end_duration": 38186,
99         "alignment_buffered": 0,
100        "num_subtasks": 192,
101        "num_acknowledged_subtasks": 192,
102        "tasks": {},
103        "external_path":
104        "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-398",
105        "discarded": true
106    }, {
107        "@class": "completed",
108        "id": 397,
109        "status": "COMPLETED",
110        "is_savepoint": false,
111        "trigger_timestamp": 1590149444586,
112        "latest_ack_timestamp": 1590149462786,
113        "state_size": 4688939894,
114        "end_to_end_duration": 18200,
115        "alignment_buffered": 0,
116        "num_subtasks": 192,
117        "num_acknowledged_subtasks": 192,
118        "tasks": {},
119        "external_path":
120        "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-397",
121        "discarded": true
122    }, {
123        "@class": "completed",
124        "id": 396,
125        "status": "COMPLETED",
126        "is_savepoint": false,
127        "trigger_timestamp": 1590149414586,
128        "latest_ack_timestamp": 1590149421011,
129        "state_size": 4872715569,
130        "end_to_end_duration": 6425,
131        "alignment_buffered": 0,
132        "num_subtasks": 192,
133        "num_acknowledged_subtasks": 192,
134        "tasks": {},
135        "external_path":
136        "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-396",
137        "discarded": true
138    }, {
139        "@class": "completed",
140        "id": 395,
141        "status": "COMPLETED",
142        "is_savepoint": false,
143        "trigger_timestamp": 1590149376379,
144        "latest_ack_timestamp": 1590149412150,
145        "state_size": 4344812828,
146        "end_to_end_duration": 35771,
147        "alignment_buffered": 0,
148        "num_subtasks": 192,

```

```

146         "num_acknowledged_subtasks": 192,
147         "tasks": {},
148         "external_path":
149     "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-395",
150     "discarded": true
151 }, {
152     "@class": "completed",
153     "id": 394,
154     "status": "COMPLETED",
155     "is_savepoint": false,
156     "trigger_timestamp": 1590149346379,
157     "latest_ack_timestamp": 1590149364845,
158     "state_size": 5322745287,
159     "end_to_end_duration": 18466,
160     "alignment_buffered": 0,
161     "num_subtasks": 192,
162     "num_acknowledged_subtasks": 192,
163     "tasks": {},
164     "external_path":
165     "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-394",
166     "discarded": true
167 }, {
168     "@class": "completed",
169     "id": 393,
170     "status": "COMPLETED",
171     "is_savepoint": false,
172     "trigger_timestamp": 1590149316379,
173     "latest_ack_timestamp": 1590149322623,
174     "state_size": 3382116678,
175     "end_to_end_duration": 6244,
176     "alignment_buffered": 0,
177     "num_subtasks": 192,
178     "num_acknowledged_subtasks": 192,
179     "tasks": {},
180     "external_path":
181     "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-393",
182     "discarded": true
183 }, {
184     "@class": "completed",
185     "id": 392,
186     "status": "COMPLETED",
187     "is_savepoint": false,
188     "trigger_timestamp": 1590149279721,
189     "latest_ack_timestamp": 1590149314164,
190     "state_size": 4695354056,
191     "end_to_end_duration": 34443,
192     "alignment_buffered": 0,
193     "num_subtasks": 192,
194     "num_acknowledged_subtasks": 192,
195     "tasks": {},
196     "external_path":
197     "hdfs://flink/checkpoints/e47bc5a70b23f4e1576cfe38d02d8ca5/chk-392",
198     "discarded": true
199 }
200 ]

```

你其实就可以发现上面可以拿到很多信息，比如每次 Checkpoint 的大小、路径、时间、最新的一次、summary 的信息。通过这些你可以拿到作业的状态信息，对于检测大状态（这个判定为大状态可以根据你们公司的场景来定，我目前是根据多个指标关联起来才判断是大状态）特别有效。

上面的 Checkpoint 状态的路径也是一个非常重要的数据，我们也会存储作业最新的一个路径，方便用户作业挂的时候能够从上一个新的 Checkpoint 路径恢复。

Checkpoint 时间配置

针对这个是因为想梳理一下公司目前所有作业的 Checkpoint 的配置，来看到底是否合理，也是有接口的，你可以通过

https://xxx/{application_id}/jobs/{jobid}/checkpoints/config 来查看，返回的结果如下：

```
1  {
2      "mode": "exactly_once",
3      "interval": 120000,
4      "timeout": 600000,
5      "min_pause": 0,
6      "max_concurrent": 1,
7      "externalization": {
8          "enabled": false,
9          "delete_on_cancellation": true
10     },
11     "state_backend": "RocksDBStateBackend"
12 }
```

这几个字段的含义：

- mode 模式：有 EXACTLY_ONCE、AT_LEAST_ONCE、NONE
- interval：单位是毫秒，如果是 9223372036854776000 则表示没开启 Checkpoint
- timeout：单位是毫秒
- state_backend：有 RocksDBStateBackend、MemoryStateBackend、FsStateBackend

另外还可以看到这些作业是否开启了 Checkpoint。

作业 Flink 版本

因为现在我们集群有两个 Flink 版本，所以我们也有接口去统计集群上所有作业 Flink 版本的统计，你可以看看 http://xxx/{application_id}/config 该接口，返回的数据如下：

```

1 {
2   "refresh-interval": 3000,
3   "timezone-name": "China Time",
4   "timezone-offset": 28800000,
5   "flink-version": "1.9.1",
6   "flink-revision": "4d56de8 @ 30.09.2019 @ 11:32:19 CST"
7 }

```

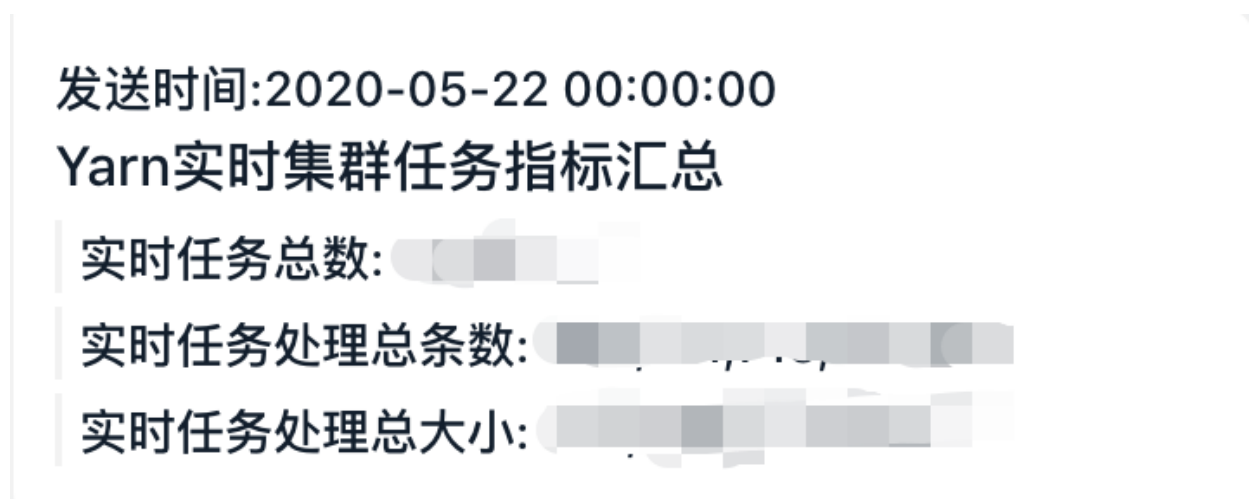
这里可以拿到 `flink-version` 字段的值，但是我们发现有的版本竟然是 `<unknown>` 。

作业处理数据量

还有一个比较重要的就是统计每天我们集群所有的 Flink 作业处理的流数据的数据总量，供我们领导查看，这个就很关键了。

目前也是通过 UI 上的 Rest API 信息查看到的信息，然后统计每个作业的，最后汇总成一个总数。

效果如下：



你可以通过单个作业的查看，比如下图：



