# 基于k8s的Flink SQL大规模数据相似度计算

汇报人：22121630 汪江豪
22120721 阮金桐

CONTENT

目 录

# 1.Flink多节点部署方式

flink集群

flink多节点部署

k8s集群

## 传统集群方式

设置ssh无密码互联

修改flink/conf/flink-conf.yaml文件：
obmanager.rpc.address: master
taskmanager.numberOfTaskSlots: 2

修改conf目录下master文件为
master:8081

修改conf目录下worker文件为:
slave1
slave2

slave1,slave2节点上设置同理
额外修改flink-conf.yaml文件中的
taskmanager.host:slave1/slave2
为自己主机名

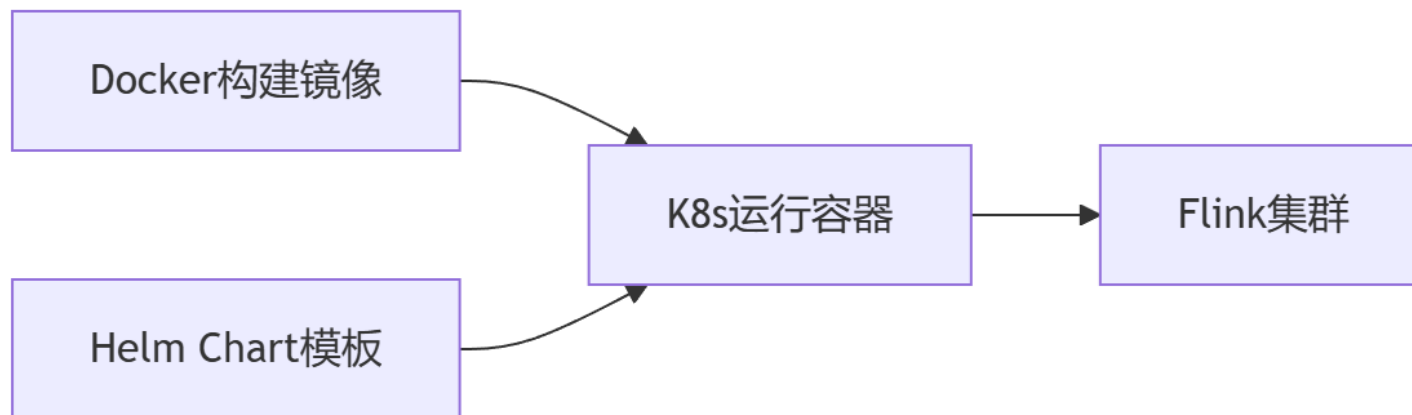在master节点运行start-cluster.sh，提交
jar包

**传统集群部署**

调度节点 — master

工作节点 — slave1
slave2

Docker构建镜像 → K8s运行容器 → Flink集群

Helm Chart模板 → K8s运行容器

K8s是容器编排平台，用于自动化部署、扩展和管理容器化应用

Flink中的jobmanager和taskmanager作为pod运行在k8s集群中

Docker是容器化工具，可将应用及其依赖打包成轻量级、可移植容器镜像

Helm是k8s的包管理工具，类似于linux的apt

# K8s集群方式部署Flink

启动k8s集群，节点加入集群

下载证书管理器
kubectl create -f https://github.com/jetstack/cert-manager/releases/download/v1.8.2/cert-manager.yaml

添加软件源
helm repo add flink-operator-repo https://downloads.apache.org/flink/flink-kubernetes-operator-1.11.0/

安装flink的k8s管理工具
helm install flink-kubernetes-operator flink-operator-repo/flink-kubernetes-operator

安装后operator作为pod保持运行，看到running说明安装成功
kubectl get pods

随后可通过yaml文件提交flink作业
Kubectl create –f **basic.yaml**

若要查看flink的webui，可进行端口转发，并允许外来地址访问
kubectl port-forward svc/largejob-rest 8081 --address 0.0.0.0

# basic.yaml文件

```yaml
1   apiVersion: flink.apache.org/v1beta1
2   kind: FlinkDeployment
3   metadata:
4     namespace: default
5     name: largejob
6   spec:
7     image: flink:1.16
8     flinkVersion: v1_16
9     flinkConfiguration:
10      taskmanager.numberOfTaskSlots: "12"
11    serviceAccount: flink
12    jobManager:
13      resource:
14        memory: "4096m"
15        cpu: 2
16      podTemplate:
17        spec:
18          # nodeSelector:
19            # flink-role: jobmanager
20          tolerations:
21            - key: "node-role.kubernetes.io/master"
22              operator: "Exists"
23              effect: "NoSchedule"
24          containers:
25            - name: flink-main-container
26              volumeMounts:
27                - name: my-jar
28                  mountPath: /opt/flink/largejob
29          volumes:
30            - name: my-jar
31              hostPath:
32                path: /home/ubuntu/tableapp/upload
                  type: Directory
```

```yaml
      taskManager:
35      resource:
36        memory: "51200m"
37        cpu: 12 # 每个taskmanager分配12个CPU
38      replicas: 3 # 启动2个taskmanager pod, 调度到2个节点
39      podTemplate:
40        spec:
41          nodeSelector:
42            flink-role: taskmanager
43          tolerations:
44            - key: "node-role.kubernetes.io/master"
45              operator: "Exists"
46              effect: "NoSchedule"
47          containers:
48            - name: flink-main-container
49              volumeMounts:
50                - name: my-jar
51                  mountPath: /opt/flink/largejob
52          volumes:
53            - name: my-jar
54              hostPath:
55                path: /home/ubuntu/tableapp/upload # 本地JAR所在目录
56                type: Directory
57    job:
58      jarURI: local:///opt/flink/largejob/inputfromfiledemo-1.0.jar
59      entryClass: cn.edu.shu.large_result # 填写作业的主类
      parallelism: 36
```

将本地数据集和jar包放在所有节点同一目录下，并修改权限为777
挂载到flink镜像中的某个目录

这里本地数据集和jar包位于/home/ubuntu/tableapp/upload

挂载目录为flink镜像中的/opt/flink/largejob

# 2.解法思路

# 解决办法

| 方法 | 配置 | 实际使用资源 | 总计耗时 |
|---|---|---|---|
| Flink集群 | 3台4核8G<br>2个TM节点 | 2*4核8G | 34m50s |
| K8s集群 | 3台16核64G | 3*12核50G | 8m50s |

# 解法思路

| | | |
|---|---|---|
| 1 | 116420525701620762252 | 106501936839371135489 |
| 2 | 104979585499014359063 | 103039370087174812731 |
| 3 | 101130571432010257177 | 111183544898345861357 |
| 4 | 116331515612347682756 | 114390577443742844396 |
| 5 | 110656253137238747097 | 102565924973578852934 |
| 6 | 109800388691366698136 | 102545157386069758716 |
| 7 | 106752695486123789059 | 115120856388820348743 |
| 8 | 117594348100980996964 | 100517144772812557903 |
| 9 | 109405659400238396060 | 105452294703789324242 |
| 10 | 102857824121129353719 | 104147032621576433597 |

约7万个        约10万个

原始数据特点：
6832726行、21位数字字符串、无重复行

存储类型：
INT：10位×
BIGINT：19位×
STRING、DECIMAL性能开销大×

**解决办法——转换数据集：**

排序——两列数据分别映射到0-7万，0-10万

保存映射字典，供后续转换使用

使用INT存储

# 映射字典：

第一列

```
large_referrer_dict.csv
68196    118441596793086402969,68194
68197    118441827513137270649,68195
68198    118441866522954267546,68196
68199    118443887616207357079,68197
68200    118443964499336832769,68198
68201    118444229992664950227,68199
68202    118444248646004441122,68200
68203    118444997653815563102,68201
68204    118445277475148270304,68202
68205    118445438046067769625,68203
68206    118446147598193798543,68204
68207    118446153145043212994,68205
68208
```

第二列

```
large_referee_dict.csv
101210    118444005060455277166,101208
101211    118444229992664950227,101209
101212    118444401187415650231,101210
101213    118444576709377182815,101211
101214    118444881207894739198,101212
101215    118444997653815563102,101213
101216    118445277475148270304,101214
101217    118445438046067769625,101215
101218    118446147598193798543,101216
101219    118446153145043212994,101217
101220    118446297469672800612,101218
101221    118446413811161034753,101219
101222
```

原始数据集：286M

```
1     116420525701620762252    106501936839371135489
2     104979585499014359063    103039370087174812731
3     101130571432010257177    111183544898345861357
4     116331515612347682756    114390577443742844396
5     110656253137238747097    102565924973578852934
6     109800388691366698136    102545157386069758716
7     106752695486123789059    115120856388820348743
8     117594348100980996964    100517144772812557903
9     109405659400238396060    105452294703789324242
10    102857824121129353719    104147032621576433597
```

→

转换后数据集83M

```
large_relation.csv
1     60777,35921
2     18389,16589
3     4084,61626
4     60452,79054
5     39424,13946
6     36298,13838
7     25039,83137
8     65162,2781
9     34885,30107
10    10436,22831
11    50032,38691
12    20105,43662
13    56302,55433
14    1755,60672
15    49336,4561
16    37054,55332
17    406,62598
18    6139,88632
19    54334,86173
20    15570,5331
21    36673,36743
```

# JAVA代码



```java
package cn.edu.shu;    large_result.java is not on the classpath of project wordcount, only sy

import org.apache.flink.api.common.RuntimeExecutionMode;
import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment;
import org.apache.flink.table.api.bridge.java.StreamTableEnvironment;

public class large_result {
    Run | Debug
    public static void main(String[] args) throws Exception {
        StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();
        env.setRuntimeMode(RuntimeExecutionMode.BATCH);
        // env.setParallelism(32);
        StreamTableEnvironment tableEnv = StreamTableEnvironment.create(env);

        String sourceDDL = "CREATE TABLE large_relation ("
            + "referrer INT,"
            + "referee INT"
            + ") WITH ("
            + " 'connector' = 'filesystem', "
            + " 'path' = 'file:///opt/flink/largejob/large_relation.csv', "
            + " 'format' = 'csv', "
            + " 'csv.field-delimiter' = ',', "
            + " 'csv.ignore-parse-errors' = 'true'"
            + ")";
        tableEnv.executeSql(sourceDDL);
```

在flink1.12之后，批处理模式通常比流处理模式在处理静态大数据集时效率更高：

批处理模式会自动做全局优化如(全局排序、资源复用等)。

流处理模式主要为低延迟，适合实时场景。

在传统集群方式中，
java代码设置并行数以提高效率

K8s集群方式，
java代码不用设置并行数，
在yaml文件中指定

# JAVA代码续

对于传统flink集群方式，path为本地数据集路径。

对于k8s集群方式，path为flink镜像中的挂载目录

```java
// 结果输出到文件系统
tableEnv.executeSql("CREATE TABLE SinkTable ("
    + "web1 INT, "
    + "web2 INT, "
    + "similarity DOUBLE"
    + ") WITH ("
    + " 'connector' = 'filesystem',"
    + " 'path' = 'file:///opt/flink/largejob/large_result',"
    + " 'format' = 'csv',"
    + " 'csv.field-delimiter' = ',',"
    + " 'sink.rolling-policy.file-size' = '256MB',"
    + " 'sink.rolling-policy.rollover-interval' = '30 min'"
    + ")");
```

```java
41      // 创建临时视图 ref_count
        tableEnv.executeSql("CREATE TEMPORARY VIEW tmp_ref_count AS "
42          + "SELECT referrer, COUNT(DISTINCT referree) AS web_count "
43          + "FROM large_relation "
44          + "GROUP BY referrer");
45
46      // 创建临时视图 common
47      tableEnv.executeSql("CREATE TEMPORARY VIEW tmp_common AS "
48          + "SELECT a.referrer AS web1, b.referrer AS web2, COUNT(*) AS com_cnt "
49          + "FROM large_relation a "
50          + "JOIN large_relation b ON a.referree = b.referree "
51          + "WHERE a.referrer < b.referrer "
52          + "GROUP BY a.referrer, b.referrer");
53
54      // 计算相似度并插入结果表
55      String insertSQL = "INSERT INTO SinkTable "
56          + "SELECT "
57          + "    common.web1 AS web1, "
58          + "    common.web2 AS web2, "
59          + "    CASE "
60          + "        WHEN r1.web_count + r2.web_count - common.com_cnt > 0 "
61          + "        THEN common.com_cnt * 1.0 / (r1.web_count + r2.web_count - common.com_cnt) "
62          + "        ELSE 0 "
63          + "    END AS similarity "
64          + "FROM tmp_common common "
65          + "JOIN tmp_ref_count r1 ON common.web1 = r1.referrer "
66          + "JOIN tmp_ref_count r2 ON common.web2 = r2.referrer";
67      tableEnv.executeSql(insertSQL);
        }
    }
```

在basic.yaml文件中，将数据集和jar的目录挂载好后，通过命令：

kubectl create –f basic.yaml创建任务，会自动创建镜像，执行查询

# 运行截图

# 运行截图

# 运行截图

# 结果处理——生成36个子文件

# 结果处理——显示结果文件信息

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ ls
task0   task10  task12  task14  task16  task18  task2   task21  task23  task25  task27  task29  task30  task32  task34  task4   task6   task8
task1   task11  task13  task15  task17  task19  task20  task22  task24  task26  task28  task3   task31  task33  task35  task5   task7   task9
```

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ wc -l task
383157187 task
ubuntu@VM-0-7-ubuntu:~/all_tasks$ ls -l
total 11637768
-rw-rw-r-- 1 ubuntu ubuntu 11917070025 May 19 16:12 task
```

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ head -n 10 task
25556,49046,0.02527075812274368
2922,35245,0.01211556383970177
60482,64263,0.06116207951070336
49686,60622,0.02987551867219917
4769,60665,0.0083682008368208
4769,36941,0.016666666666666667
4769,57212,0.0101010101010101
57749,67620,0.030303030303030303
2072,5687,0.03486238532110092
4769,65581,0.0108695652173913
```

上海大学
SHANGHAI UNIVERSITY

```python
def retransform_dict():
    start_time = time.time()

    # 创建临时目录
    os.makedirs(temp_dir, exist_ok=True)

    # 读取映射字典
    print("Reading reference dictionary...")
    referr_dict = pd.read_csv('large_referrer_dict.csv')
    id_to_value = dict(zip(referr_dict["referrer_id"], referr_dict["referrer_value"]))

    # 获取CPU核心数
    num_processes = cpu_count()
    print(f"Using {num_processes} processes for parallel processing")

    # 分块读取数据并并行处理
    print("Processing data chunks...")
    chunk_size = 10000000
    processed_chunks = []

    with Pool(processes=num_processes) as pool:
        results = []
        chunk_num = 0

        # 读取数据并提交处理任务
        for chunk in pd.read_csv(input_file, sep=",", header=None,
                                 names=["web1_id", "web2_id", "similarity"],
                                 chunksize=chunk_size):
            results.append(pool.apply_async(process_chunk,
                                            args=(chunk.copy(), id_to_value, chunk_num)))
            chunk_num += 1

        # 收集处理结果
        for result in results:
            processed_chunks.append(result.get())

    # 合并所有临时文件
    print("Merging temporary files...")
    first_chunk = True
```

```
  ubuntu@VM-0-7-ubuntu: ~                                    —  □  ×

  1  [|||||||||||||||||||||||||100.0%]   Tasks: 97, 671 thr; 11 running
  2  [|                          2.6%]   Load average: 10.13 6.45 3.00
  3  [|||||||||||||||||||||||||100.0%]   Uptime: 01:02:22
  4  [||||                       19.3%]
  5  [|||||||||||||||||||||||||96.7%]
  6  [|||||||                    43.4%]
  7  [|||||||||||||||||||||||||100.0%]
  8  [|||||||                    28.0%]
  9  [|||||||||||||||||||||||||100.0%]
 10  [|                           3.3%]
 11  [|                           5.3%]
 12  [|||||||||||||||||||||||||100.0%]
 13  [|                           2.7%]
 14  [|||||||||||||||||||||||||100.0%]
 15  [|||||||||||||||||         61.8%]
 16  [|||||||||||||||||||||||||100.0%]
 Mem[||||||||||||||||||||7.76G/61.1G]
 Swp[                          0K/0K]

  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
67694 ubuntu     20   0 1167M  374M 23948 R 103.  0.6  0:30.71 python transform.py
67714 ubuntu     20   0 1270M  685M 10368 R 101.  1.1  0:20.17 python transform.py
67715 ubuntu     20   0 1270M  685M 10368 R 100.  1.1  0:18.02 python transform.py
67720 ubuntu     20   0 1269M  685M 10368 R 100.  1.1  0:05.85 python transform.py
67721 ubuntu     20   0 1269M  685M 10368 R 100.  1.1  0:03.11 python transform.py
67717 ubuntu     20   0 1270M  685M 10368 R 99.9  1.1  0:12.75 python transform.py
67719 ubuntu     20   0 1269M  685M 10368 R 99.9  1.1  0:08.63 python transform.py
67718 ubuntu     20   0 1269M  685M 10368 R 99.9  1.1  0:10.69 python transform.py
67716 ubuntu     20   0 1270M  685M 10368 R 99.2  1.1  0:15.08 python transform.py
67722 ubuntu     20   0 1038M  451M  8444 R 34.4  0.7  0:00.52 python transform.py
67728 ubuntu     20   0 1167M  374M 23948 S 19.2  0.6  0:03.37 python transform.py
67713 ubuntu     20   0  650M 67116 10368 S  4.6  0.1  0:21.01 python transform.py
62697 9999       20   0 4721M  661M 29608 S  4.0  1.1  0:22.17 /opt/java/openjdk/b
63014 9999       20   0 4721M  661M 29608 S  3.3  1.1  0:08.00 /opt/java/openjdk/b
 8434 root       20   0  209M  100M 57088 S  3.3  0.2  0:43.90 kube-controller-man
 8503 root       20   0  624M  505M 68784 S  2.6  0.8  1:42.60 kube-apiserver --ad
 8928 root       20   0 3167M 94184 63384 S  2.6  0.1  0:58.34 /usr/bin/kubelet --
 3420 root       20   0  975M 65720 31772 S  1.3  0.1  0:28.51 /usr/local/qcloud/Y
 4536 root       20   0  621M 23772  5824 S  1.3  0.0  0:24.02 barad_agent
 8512 root       20   0 10.1G 81856 26180 S  0.7  0.1  0:43.51 etcd --advertise-cl
68139 ubuntu     20   0 9536   5376  3444 R  0.7  0.0  0:00.14 htop
 3426 root       20   0  975M 65720 31772 S  0.7  0.1  0:08.87 /usr/local/qcloud/Y

F1Help  F2Setup F3Search F4Filter F5Tree  F6SortBy F7Nice -F8Nice +F9Kill  F10Quit
```

# 结果处理——显示最终结果

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ python transform.py
Reading reference dictionary...
Using 16 processes for parallel processing
Processing data chunks...
Merging temporary files...
Success! Processed 383157187 rows in 192.00s
```

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ ls -l
total 35038084
-rw-rw-r-- 1 ubuntu ubuntu     2925518 May 16 21:28 large_referree_dict.csv
-rw-rw-r-- 1 ubuntu ubuntu     1966892 May 16 21:28 large_referrer_dict.csv
-rw-rw-r-- 1 ubuntu ubuntu 23957013430 May 19 16:24 large_result_final.csv
-rw-rw-r-- 1 ubuntu ubuntu 11917070025 May 19 16:12 task
-rw-rw-r-- 1 ubuntu ubuntu        3385 May 19 16:21 transform.py
```

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ wc -l large_result_final.csv
383157188 large_result_final.csv
ubuntu@VM-0-7-ubuntu:~/all_tasks$ head -n 20 large_result_final.csv
web1,web2,similarity
1068958236763207711807,11329054788938503295 8,0.0252707581227436
10081612426229044507 0,1095052271976405891 04,0.0121155638397017
1163379757684562028 28,1173447869847882054 59,0.0611620795107033
1134538746896065945 46,1163758478902925878 30,0.0298755186721991
1013202356080312379 06,1163854894664044089 87,0.00836820083682
1013202356080312379 06,1099749504873420424 28,0.0166666666666666
1013202356080312379 06,1154653805834757179 8,0.0101010101010101
1155972963007639620 23,1182789728162018655 35,0.0303030303030303
1005786456846821967 43,1015706468090059715 41,0.0348623853211009
1013202356080312379 06,1177130692871572690 96,0.0108695652173913
1013202356080312379 06,1145535282887443904 01,0.0714285714285714
1013202356080312379 06,1102079535203995578 67,0.0454545454545454
1013202356080312379 06,1154644773765397822 6,0.0018315018315018
1011363134520223728 34,1012148419672818090 81,0.0445168295331161
1005578190524698266 67,1080877767762015713 4,0.1247947454844006
1042430537321227193 47,1079742983398679180 05,0.0580645161290322
1103991783546465073 11,1182128753257668286 0,0.0353356890459364
1068958236763207711807,1071694751350453611 47,0.0151057401812688
1068958236763207711807,1158046716735687562 7,0.0056818181818181
```

```
ubuntu@VM-0-7-ubuntu:~/all_tasks$ tail -n 20 large_result_final.csv
10035138078242150883 2,1124438724363424694 63,0.0133333333333333
10035138078242150883 2,1074862434672825012 80,0.0010649627263045
10702856501799133051 7,1120827485651010458 22,0.009009009009009
11136675618757671207 8,1145166492444456189 45,0.0011286681715575
11288907087442037174 5,1138322825346741623 27,0.0045662100456621
11140941745229992076 9,1168775821618765718 00,0.0031645569620253
10350680097300259097 6,1054111820227399041 20,0.001081081081081
10842008215203564988 8,1158750566302369475 29,0.0158730158730158
10009346888040944943 6,1020502730561355204 52,0.0022123893805309
10187583940076015871 2,1057316998471742303 47,0.0032051282051282
10187583940076015871 2,1139302163202530874 85,0.0294117647058823
10187583940076015871 2,1036277219813782043 76,0.0011792452830188
10187583940076015871 2,1112062260759498773 58,0.0034013605442176
10187583940076015871 2,1043596262679067916 01,0.0163934426229508
10187583940076015871 2,1099264737832086350 57,0.03125
10187583940076015871 2,1027182747918896106 73,0.0294117647058823
10187583940076015871 2,1109828844915502894 66,0.0188679245283018
10187583940076015871 2,1123121477697203491 03,0.0083333333333333
10322166658431155693 1,1090328688149257042 18,0.00093896713615023
11412256727096216971 9,1150116464129566791 15,0.0033003300330033
```

# 3.问题处理

# 问题处理

跑flink作业时发现空闲节点
TaskManager节点只有一个

原因：
k8s对于JobManager和TaskManager
随机分配
该节点被JobManager调度占用，不进
行计算

Master节点被自动打上污点，TM不会
分配到该节点；或master节点资源不
足，TM的pod一直处于pending状态

# 解决办法

```
     taskManager:
35     resource:
36       memory: "51200m"
37       cpu: 12 # 每个taskmanager分配12个CPU
38     replicas: 3 # 启动2个taskmanager pod, 调度到2个节点
39     podTemplate:
40       spec:
41         nodeSelector:
42           flink-role: taskmanager
43         tolerations:
44           - key: "node-role.kubernetes.io/master"
45             operator: "Exists"
46             effect: "NoSchedule"
47         containers:
48           - name: flink-main-container
49             volumeMounts:
50               - name: my-jar
51                 mountPath: /opt/flink/largejob
52         volumes:
53           - name: my-jar
54             hostPath:
55               path: /home/ubuntu/tableapp/upload # 本地JAR所在目录
56               type: Directory
57     job:
58       jarURI: local:///opt/flink/largejob/inputfromfiledemo-1.0.jar
59       entryClass: cn.edu.shu.large_result # 填写作业的主类
       parallelism: 36
```

移除master节点污点限制
kubectl taint node vm-0-9-ubuntu node-role.kubernetes.io/master:NoSchedule-

若想强制选择pod运行的节点
可以在yaml文件中添加节点选择标签，
只有当节点标签为taskmanager时才分配为TM

为节点打上标签
kubectl label nodes vm-0-14-ubuntu flink-role=taskmanager
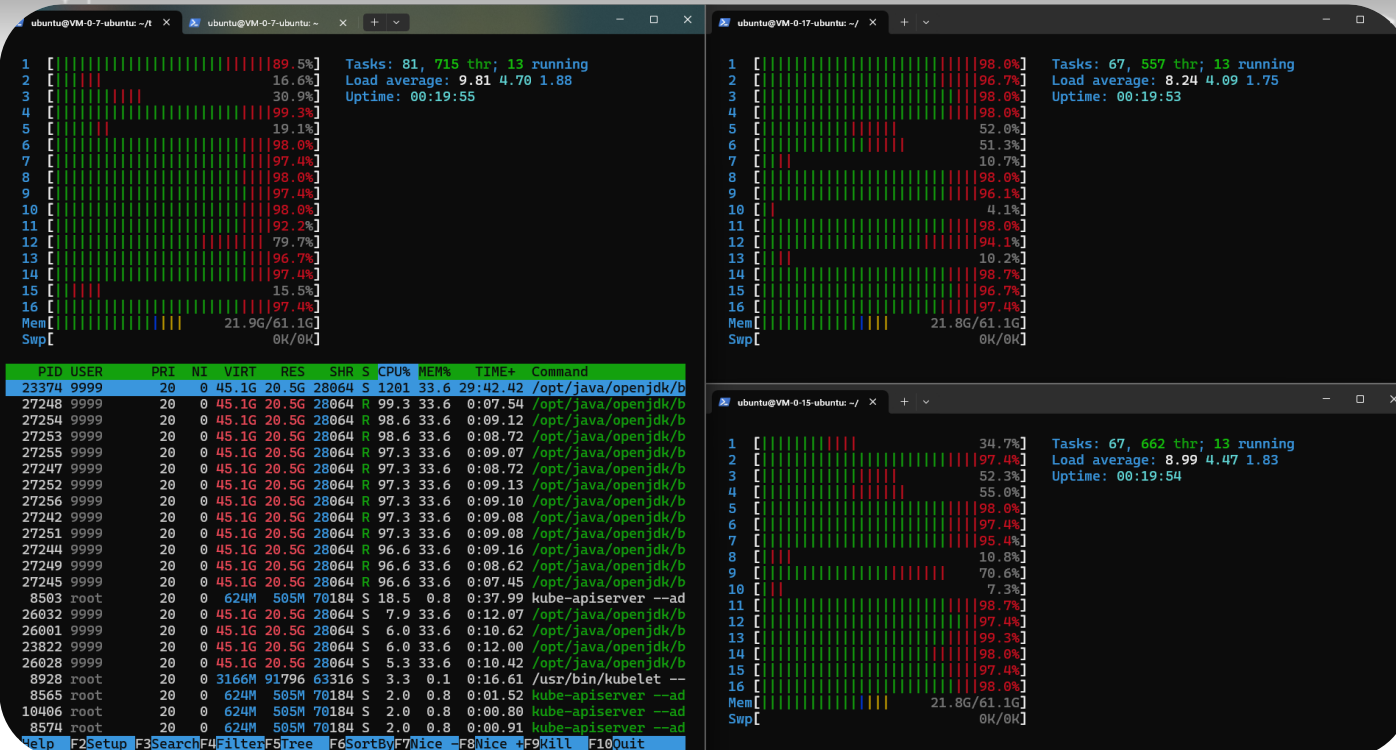显示节点标签
kubectl get nodes --show-labels

```
ubuntu@VM-0-7-ubuntu:~/tableapp$ kubectl get nodes --show-labels
NAME            STATUS   ROLES    AGE   VERSION   LABELS
vm-0-15-ubuntu  Ready    <none>   15m   v1.18.0   beta.kubernetes.io/arch=amd64,beta.kubernetes.io/os=linux,flink-role=taskmanager,kubernetes.io/arch=amd64,kubernetes.io/
hostname=vm-0-15-ubuntu,kubernetes.io/os=linux
vm-0-17-ubuntu  Ready    <none>   15m   v1.18.0   beta.kubernetes.io/arch=amd64,beta.kubernetes.io/os=linux,flink-role=taskmanager,kubernetes.io/arch=amd64,kubernetes.io/
hostname=vm-0-17-ubuntu,kubernetes.io/os=linux
vm-0-7-ubuntu   Ready    master   17m   v1.18.0   beta.kubernetes.io/arch=amd64,beta.kubernetes.io/os=linux,flink-role=taskmanager,kubernetes.io/arch=amd64,kubernetes.io/
hostname=vm-0-7-ubuntu,kubernetes.io/os=linux,node-role.kubernetes.io/master=
ubuntu@VM-0-7-ubuntu:~/tableapp$
```

# 让TaskManager跑在3个节点上

# THANK YOU!

汇报人： 22121630 汪江豪
22120721 阮金桐