



链滴

数据 ETL 平台 Streamsets 简单应用

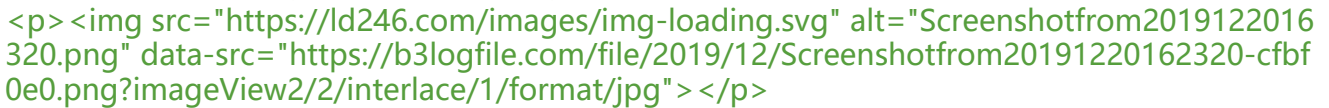
作者: [Gaoshengyue](#)

原文链接: <https://ld246.com/article/1576837366412>

来源网站: [链滴](#)

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

StreamSets 介绍

 Screenshot from 2019122016320.png

StreamSets 是最近兴起的 ETL 平台，它的特点是具有多样性的组件支持，可集成于 CDH 平台最为吸引人的就是可视化的数据流通流程设置，多个 pipelines 的编写，RestApi 形式的自动化支持等等等等，当然选择使用它的最终理由还是因为支持的组件多。这一篇就简单来讲讲 StreamSets 的单使用。

StreamSets 部署

StreamSets 的部署有很多形式，这里列举两个最方便的。

Docker-compose

这里我提供一份写好的 docker-compose 文件，如下

```
version: '3.1'
services:
  streamsets:
    image: streamsets/datacollector
    restart: always
    ports:
      - 18630:18630
    environment:
      HOST_IP: 0.0.0.0
    volumes:
      - streamsets:/opt/steramsets
```

```
-p">:</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">streamsets</span><span class="high
ight-p">:</span><span class="highlight-w">
</span></span></span></code></pre>
<p>这里把端口映射到了 18630 方便等下的登录使用, HOST_IP 改成万能 IP, 不设置访问限制方便
问, 相关数据挂到本地磁盘。 <br>
启动命令</p>
<pre><code class="language-bash highlight-chroma"><span class="highlight-line"><span c
ass="highlight-cl">sudo docker-compose up -d
</span></span></code></pre>
<p>直接复制 yaml 配置创建 docker-compose.yaml 就可以了</p>
<h3 id="K8S">K8S</h3>
<p>因为 StreamSets 内部的数据 collector 重启机制问题, 会导致组件无法成功添加的问题, 所以
里还没有解决, 暂不推荐 K8S, 不过还是把配置文件给出来(组件无法附加版本)</p>
<pre><code class="language-yaml highlight-chroma"><span class="highlight-line"><span c
ass="highlight-cl"><span class="highlight-nt">apiVersion</span><span class="highlight-p
">:</span><span class="highlight-w"> </span><span class="highlight-l">apps/v1</span><span><
pan class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"></span><span class="highlight-nt">kind</span><span class="highlight-p
">:</span><span class="highlight-w"> </span><span class="highlight-l">Deployment</spa
><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"></span><span class="highlight-nt">metadata</span><span class="highli
ht-p">:</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">name</span><span class="highlight
p">:</span><span class="highlight-w"> </span><span class="highlight-l">streamsets-data
collector</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">namespace</span><span class="high
light-p">:</span><span class="highlight-w"> </span><span class="highlight-l">default</s
an><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">labels</span><span class="highlight
p">:</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">app</span><span class="highlight
p">:</span><span class="highlight-w"> </span><span class="highlight-l">streamsets-data
collector</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"></span><span class="highlight-nt">spec</span><span class="highlight-p
">:</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">replicas</span><span class="highlig
t-p">:</span><span class="highlight-w"> </span><span class="highlight-m">1</span><span><s
an class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">selector</span><span class="highlig
t-p">:</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w"> </span><span class="highlight-nt">matchLabels</span><span class="h
ghlight-p">:</span><span class="highlight-w">
```



```
<span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w">  </span><span class="highlight-nt">targetPort</span><span class="high
light-p">:</span><span class="highlight-w"> </span><span class="highlight-m">18630</s
pan><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w">  </span><span class="highlight-nt">nodePort</span><span class="highl
ght-p">:</span><span class="highlight-w"> </span><span class="highlight-m">31630</s
an><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w">  </span><span class="highlight-nt">selector</span><span class="highlig
ht-p">:</span><span class="highlight-w">
</span></span></span><span class="highlight-line"><span class="highlight-cl"><span cla
s="highlight-w">  </span><span class="highlight-nt">app</span><span class="highlight
p">:</span><span class="highlight-w"> </span><span class="highlight-l">streamsets-data
ollector</span><span class="highlight-w">
</span></span></span></code></pre>
<p>也是把端口暴露在外部 18630<br>
复制下来创建 StreamSets.yaml<br>
启动命令</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight
cl">kubectl apply -f StreamSets.yaml -n data-service
</span></span></code></pre>
<p>这里是使用 data-service 命名空间，K8S 的概念这里不多说了</p>
<h3 id="CDH">CDH</h3>
<p>这个安装方式我就不多说了吧，因为直接在 cloudera-manager 里面添加角色就行了，具体在
个机器上的角色看自己的需求了。</p>
<h2 id="StreamSets应用">StreamSets 应用</h2>
<h3 id="Pipelines创建">Pipelines 创建</h3>
<p>这里我们简单的举一个 HTTP-client 数据获取的例子。<br>
为了方便起见，我们先创建一个 test.txt 文件，内容可以如下：</p>
<pre><code class="language-text highlight-chroma"><span class="highlight-line"><span cl
ss="highlight-cl">name,age
</span></span><span class="highlight-line"><span class="highlight-cl">symoon,18
</span></span><span class="highlight-line"><span class="highlight-cl">willson,50
</span></span></code></pre>
<p>这里用,分割，方便那我们取数据。<br>
当然也可以自己设置测试数据，这里随便写了一个。<br>
接下来我们简单开启一个 python 服务，用来启动一个 api</p>
<pre><code class="language-bash highlight-chroma"><span class="highlight-line"><span c
lass="highlight-cl">python3 -m http.server <span class="highlight-m">8001</span>
</span></span></code></pre>
<p>这个不多说了，大家应该都知道，这条命令在 txt 所在文件夹启动，这时候访问</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight
cl">http://127.0.0.1:8001/test.txt
</span></span></code></pre>
<p>这个 API 就是我们需要的数据。接下来我们到 StreamSets 里面操作。<br>
<br>
账号密码都是默认的 admin<br>
输入正确后，我们进入到主页面<br>
<br>

点击 CREATE NEW PIPELINE 创建一个新的任务，Title 什么的随便填，我直接填了 test<br>

创建成功后我们进入到任务管理界面：<br>

</p>

<p>我们选择拖拽一个 HTTP Client 作为数据源。相关配置如下填写，这里我们用 docker 启动的 StramSets 的化就无法用 127.0.0.1 了，我们使用本机的 IP</p>

<p></p>

<p></p>

<p></p>

<p>这里填写的信息都有说明，相信大家都能看懂就不多说了，接下来测试一下，我们点击一下右上类似眼睛的按钮。<br>

默认信息就好，批量可以从 10 改成 1。结果如图<br>

</p>

<h3 id="Mysql增量采集案例">Mysql 增量采集案例</h3>

<p>这里我们举个简单的增量采集案例，采集 Mysql 指定表中的增量数据到 Kafka 中。<br>

接下来来实战一下，我们先创建一个 pipeline，上面已经说明了，就不多说了。<br>

因为 mysql 的 connector 需要安装 jdbc 组件，所以我们需要进行组件安装。进入 Pipeline 中的界，以如下方式选择。<br>

</p>

<p>点击加号 <strong>Add/Remove Stages</strong>，进入到组件安装界面。<br>

</p>

<p>这里我们找到 JDBC，点击右侧的配置按钮，点击 install。<br>

因为我们需要用到 Kafka 服务，这里我们继续进入 Apache Kafka 中，安装 Kafka-V2.0 版本。安装后会提示重启，根据提示点击 Restart data collector 按钮即可。等待一分钟左右服务重启即可。<p>

<p>重新进入到我们刚才创建的 Pipeline，配置 JDBC 拖动的组件及配置如下：<br>

<br>

这里我们的数据库名字和表名字分别是 test 和 Test，偏移值由 id 字段来计算。初始值设置 0。<br>

因为 jdbc 需要相应的驱动，所以我们需要在 External Libraries 处上传 mysql 的 jdbc 驱动。这里接附上一份。<br>

<a href="https://b3logfile.com/file/2019/12/mysqlconnectorjava8.0.18-c15d56a7.jar">mysqlconnectorjava8.0.18.jar</a><br>

下载好后上传，依旧需要重启服务生效。<br>

后面还需要提供账号密码，这里就不截图多说了。<br>

这时候依旧是点那个眼睛按钮测试一下，首先我们还是需要在数据库里面创建这张表的，具体格式如</p>

```
<pre><code class="highlight-chroma">id|name|age
```



```
1|symoon|18
2|willson|19
3|lake|20
</code></pre>
<p>大概这种表格式，测试后结果如下：

这时候我们已经拿到了增量数据，接下来配置 Kafka 的数据接收端。

</p>
<p>这里我们在上面的接收端选择直接选择 Kafka producer 就可以了，我使用的是本机 docker 启的 Kafka 具体的配置我就不发了，在之前的一键启动 docker 程序里面有，topic 随便输入一个 test 接下来还是点击眼睛测试一下。</p>
<p></p>
<p>kafka 的 input 成功了，接下来我们直接运行起来。

数据流通的性能还有一些相关信息都在这里显示了，这个任务会一直运行着，增量数据都可以流通到 afka，其他的数据库使用其他连接器也可以支持，这里就不多说了，大家可以都——去试试。</p>
```