

Base-line Apache Spark performance with Apache Spark history server in Azure Synapse Analytics

Note:

You are not required to complete the processes, tasks, activities, or steps presented in this example. The various samples provided are for illustrative purposes only and it's likely that if you try this out you will encounter issues in your system.

The Apache Spark history server can be used to debug and diagnose completed and running Apache Spark applications. You can use the Apache Spark history server web UI from the Azure Synapse Studio environment.

Once you launch it, there are several tabs that you can use in order to monitor the Apache Spark application:

- Jobs
- Stages
- Storage
- Environment
- Executors
- SQL

The Apache Spark history server is the web user interface known as the *Spark UI*, and is used to view completed and running Apache Spark applications.

If you want to navigate to the Apache Spark History server, you can navigate to the Azure Synapse Analytics Studio environment and go to the **Monitor** tab. In the Monitor tab, you can select **Apache Spark Applications**.

Integration

- Home
- Data
- Develop
- Integrate
- Monitor
- Manage

Activities

- Pipeline runs
- Trigger runs
- Integration runtimes
- Apache Spark applications
- SQL requests
- Data flow debug

Apache Spark applications > Livy ID 0

Synapse_NBS_sparkpoolmod_1605597928
Completed tasks 16 of 16 Status Cancelled Total duration 33m 31s

Cancel Refresh Spark history server

Attempts 1 of 1
All job IDs View Progress Playback 0ms / 7s 571ms

Stage 0
Job ID: 0
Task 1
Duration: 34.464ms
Run ID: 0
Data read: 0Byte
Data written: 0Byte
[View details](#)

Stage 1
Job ID: 1
Task 4
Duration: 39.722ms
Run ID: 0
Data read: 0Byte
Data written: 0Byte
[View details](#)

Stage 2
Job ID: 2
Task 11
Duration: 133ms
Run ID: 0
Data read: 0Byte
Data written: 0Byte
[View details](#)

Diagnostics Logs

- Failed jobs
All 3 jobs were completed successfully.
- Data skew
No data skew detected.
- Time skew
No time skew detected.
- Executor utilization

Details

Summary

Application

Application ID
application_1605598015663_0001

Queued duration
1s

Running duration
33m 30s

Livy ID
0

Submitter
kawa@je@microsoft.com

Submit time
11/11/2020 7:25:28 AM

Executors
2

Spark pool

Name
sparkpoolmod

Auto-pause
On

Number of minutes idle
15

Node size
Medium (8 vCores / 64 GB)

Monitor Apache Spark Application Overview

If you are familiar with Apache Spark, you can find the standard Apache Spark history server UI by selecting **Open Spark UI**.

Another way to open the Apache Spark History server is to navigate to the **Data** tab. If you create a notebook and read a DataFrame then you can go to the bottom of the page and find the Spark History Server known as the Spark UI.

1. From your Azure Synapse Studio notebook, select **Open Spark UI** from the job execution output cell or from the status panel at the bottom of the notebook document.

We use optional cookies to provide a better experience. [Learn more](#)

Accept Reject More options

Data Workspace Linked

Filter resources by name

Databases 1

Notebook 1

+ Cell Run all Publish Attach to: sparkpoolmod Language: PySpark (Python) Preview Features

Ready

```
1 new_rows = [{"CA", 22, 45000}, {"WA", 35, 65000}, {"WA", 50, 85000}]
2 demo_df = spark.createDataFrame(new_rows, ['state', 'age', 'salary'])
3 demo_df.show()
```

Command executed in 10s 187ms by kawa@je on 11-20-2020 12:27:40:494 +00:00

Job execution Succeeded Spark 2 executors 16 cores


[View in monitoring](#) [Open Spark UI](#)

ID	Description	Status	Stages	Tasks	Submission Time	Duration
Job 0	showString at NativeMethodAccessorImpl.java:0	Succeeded	1/1	1/1 succeeded	12:27:32 PM, 11/20/20	3 sec
Job 1	showString at NativeMethodAccessorImpl.java:0	Succeeded	1/1	4/4 succeeded	12:27:36 PM, 11/20/20	3 sec
Job 2	showString at NativeMethodAccessorImpl.java:0	Succeeded	1/1	11/11 succeeded	12:27:39 PM, 11/20/20	0 sec

```
state|age|salary|
-----+-----+-----
CA| 22| 45000|
WA| 35| 65000|
WA| 50| 85000|
```

Open Spark UI through Data Tab

2. Select **Spark UI** from the slide out panel to be redirected to the Spark monitoring tab where you will land in the **Jobs** tab.

<div>  2.4.4.2.6.99.201-25973884 </div> <div> Jobs Stages Storage Environment Executors SQL </div> <div> Synapse_sparkpoolmod_1606115314 application UI </div>					
<h3>Spark Jobs (?)</h3> <p> User: trusted-service-user Total Uptime: 4.7 min Scheduling Mode: FIFO Completed Jobs: 3 Event Timeline Completed Jobs (3) </p>					
Job Id (Job Group) *	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
2 (3)	Job group for statement 3 showString at NativeMethodAccessorImpl.java:0	2020/11/23 07:11:20	0.1 s	1/1	<div>11/11</div>
1 (3)	Job group for statement 3 showString at NativeMethodAccessorImpl.java:0	2020/11/23 07:11:17	3 s	1/1	<div>4/4</div>
0 (3)	Job group for statement 3 showString at NativeMethodAccessorImpl.java:0	2020/11/23 07:11:13	4 s	1/1	<div>1/1</div>

Open Spark UI

3. Within the Jobs tab of the Spark History server, you can view:

- Job ID
- Job description
- the time when the job was submitted
- the duration
- the stages
- and the task (for all stages).

If you select a job ID, go to the **description** and select the **url** to be redirected to the following screen:

Details for Job 2

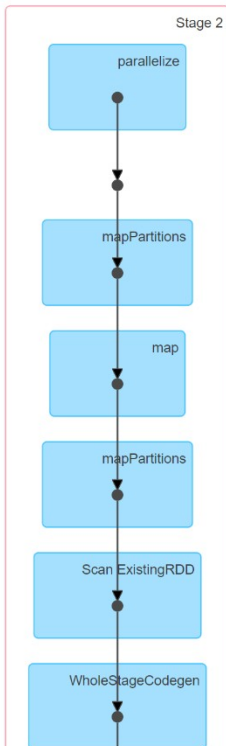
Status: SUCCEEDED

Job Group: 3

Completed Stages: 1

▶ Event Timeline

▼ DAG Visualization



DAG visualization

It will give you a DAG Visualization of the stage.

4. If you select the **Stages** tab, you'll find all the completed stages. If you want to dig deeper into the stages, you can select the **description url** as shown below:

Disk Size: 0.0 B

Block Name	Storage Level	Size in Memory	Size on Disk	Executors
rdd_22_0	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f00337162244:36557
rdd_22_1	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f004f3689053:41025
rdd_22_10	Memory Deserialized 1x Replicated	720.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f00337162244:36557
rdd_22_11	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f004f3689053:41025
rdd_22_12	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f00337162244:36557
rdd_22_13	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f004f3689053:41025
rdd_22_14	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f00337162244:36557
rdd_22_15	Memory Deserialized 1x Replicated	720.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f004f3689053:41025
rdd_22_2	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f00337162244:36557
rdd_22_3	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f004f3689053:41025
rdd_22_4	Memory Deserialized 1x Replicated	24.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f00337162244:36557
rdd_22_5	Memory Deserialized 1x Replicated	720.0 B	0.0 B	bf5e6d37de2247e7847bed207a7b5e8f004f3689053:41025

10. Finally, there is the **SQL** tab. Here you can find the completed queries and details within the IDs for the queries.

SQL

Completed Queries: 8

Completed Queries (8)

ID	Description		Submitted	Duration	Job IDs
-4	PeregrinePlanLogEvent Spark SQL Plans with annotations.	+details	2020/11/23 07:31:45	21 ms	
3	count at NativeMethodAccessorImpl.java:0 org.apache.spark.sql.Dataset.count(Dataset.scala:2843) sun.reflect.NativeMethodAccessorImpl.invoke(Native Method) sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62) sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43) java.lang.reflect.Method.invoke(Method.java:498) py4j.reflection.MethodInvoker.invoke(MethodInvoker.java:244) py4j.reflection.ReflectionEngine.invoke(ReflectionEngine.java:357) py4j.Gateway.invoke(Gateway.java:282) py4j.commands.AbstractCommand.invokeMethod(AbstractCommand.java:132) py4j.commands.CallCommand.execute(CallCommand.java:79) py4j.GatewayConnection.run(GatewayConnection.java:238) java.lang.Thread.run(Thread.java:748)	+details	2020/11/23 07:31:44	0.6 s	[5]
-3	PeregrinePlanLogEvent	+details	2020/11/23 07:30:58	12 ms	
2	count at NativeMethodAccessorImpl.java:0	+details	2020/11/23 07:30:58	0.3 s	[4]
-2	PeregrinePlanLogEvent	+details	2020/11/23 07:28:43	26 ms	
1	count at NativeMethodAccessorImpl.java:0	+details	2020/11/23 07:28:42	0.6 s	[3]
-1	PeregrinePlanLogEvent	+details	2020/11/23 07:11:20	0.1 s	
0	showString at NativeMethodAccessorImpl.java:0	+details	2020/11/23 07:11:12	7 s	[0][1][2]

SQL Tab Spark UI

Now we have been through all the different tabs using the Spark UI, enabling you to optimize and baseline through the Apache Spark performance and settings.