

The Apache Spark User Interface

Objectives

After watching this video, you will be able to:

- Explain how to connect to the Apache Spark user interface web server
- Navigate basic elements of the Spark user interface

Why use the Application UI

The Spark user interface shows information about the running application:
 display app job including job status
 are part of the state of tasks
 tabs

- Shows jobs, stages and tasks
- Storage of persisted RDDs and DataFrames
- Environment configuration and properties
- Executor summary
- SQL information (if SQL queries exist)

Job Details

View the Job Stages timeline and job status

Quickly view stage details based on status

Click the Description link to view completed stage details

Storage Information

details about RDDs that have been cached, or persisted, to memory and/or written to disk

Click the RDD Name link to view task details for that stage.

Executors Information

The Executors tab

Scroll the window to view

- Spark configuration properties
- Resource profiles
- Hadoop properties
- System configuration properties

SQL Query Details - query plan

View the SQL query plan

Access the Application UI

Start your application
 Connect to the application UI using the following URL:

<http://<driver-node>:4040>

The SparkContext starts a web server for the application UI

- The driver program is the host
- The port defaults to 4040
- The UI is available only while the application is active

Jobs Information

The Jobs tab

Jobs summary info

Jobs event timeline

Click the Description link to view completed job details

Stages Information

The Stages tab

Click the Description links to view task details for that stage.

Environment Information

The Environment tab

View:

- Spark configuration properties
- Resource profiles
- Hadoop properties
- System configuration properties

SQL Information

The SQL tab

Click the SQL Description link to view the query details

SQL Query Details - the DAG

View the SQL DAG

Debugging Apache Spark Application Issues

Objectives

After watching this video, you will be able to:

- Identify common Apache Spark application issues
- Debug issues through the application UI
- Describe application log file locations and content

Common Application Issues

Common areas of Spark application issues

- User Code
- Configuration
- Application Dependencies
- Resource Allocation

↓ include

What is User Code

User code is made up of:

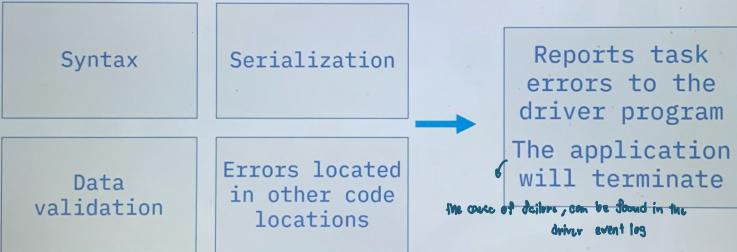
Driver program: code that runs in the driver process and the functions and variables serialized that the executor runs in parallel

Serialized closures contain the code's necessary functions, classes, and variables

The serialized closures are distributed into the cluster to run in parallel by the executors

User Code Issues

- User code can error in the driver



Application dependency issues

Application dependencies include:

- Application files:
 - Source files examples: Python script files, Java JAR files
 - Required data files
- Application libraries

Dependencies must be available for all nodes of the cluster

Application Resource Issues

- CPU and memory resources must be available for tasks to run
- Driver and executor processes require some amount of CPU and memory to run
- Any Worker with free resources can start processes
- Resources are acquired until a process completes
- If resources are not available, Spark retries until a worker is free
- Lack of resources results in task time-outs, also called task starvation!

Examining the Log Files

↳ provide detail that give more insight of failure causes

- Application logs are found in 'work/' directory named as work/<application-id>/<stdout|error>
Spark installation directory
- Spark standalone writes master/worker logs to the 'log/' directory

Understanding Memory Resources

Objectives

After watching this video, you will be able to:

- Describe memory parameters
- Describe Spark memory management
- Describe Spark data persistence
- Explain how to set executor memory on submit
- Explain how to set memory for Spark Standalone memory and cores

Configuring Spark Process Memory

Spark applications allow configuration of driver and executor memory

Upper limit on useable memory enables apps to run without using all available cluster memory

Exceeding memory requirements causes disk spill or out-of-memory errors

Memory Setting Considerations

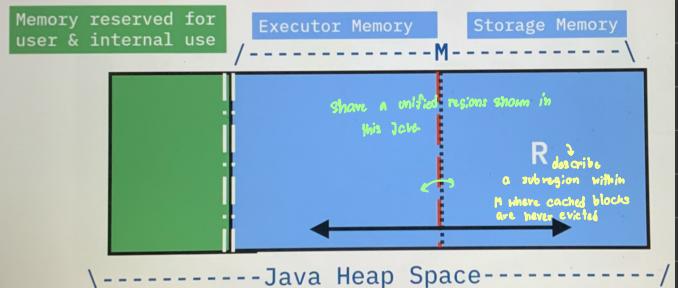
Executor Memory:

- Processing
- Caching
- Excessive caching leads to issues

Driver Memory:

- Loads data, broadcasts variables
- Handles results, such as collections

Spark Unified Memory

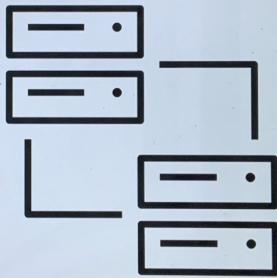


- This design ensures several preferable properties
 - 1. app. that do not use caching can use the entire space for executor memory
 - 2. app. that do not use caching can reserve a minimum storage space

Spark Data Persistence

Spark data persistence:

- Store intermediate calculations for reuse
- Persist to memory/disk the same data
- Less computation
- Faster iterations over data when training a model



Data Persistence Example

Caching a DataFrame with a features column of random values

```
# Define DataFrame with feature column
df = spark.range(100).withColumn("features", rand()).cache()

# Features are generated then cached first time action is applied
print(df.filter(col("features") > 0.5).count())

# Cached DataFrame with features is reused in subsequent calls
print(df.filter(col("features") < 0.5).count())
```

Setting Memory on Submit

Set the executor memory for a Spark standalone cluster:

```
$ ./bin/spark-submit \
--class org.apache.spark.examples.SparkPi \
--master spark://<spark-master-URL>:7077 \
--executor-memory 10G \
/path/to/examples.jar \
1000
```

Set each executor to use up to 10GB of memory

Setting Worker Node Resources

Set the Spark Standalone worker memory and cores:

```
# Start standalone worker with MAX 10Gb memory, 8 cores
$ ./sbin/start-worker.sh \
spark://<spark-master-URL> \
--memory 10G --cores 8
```

Default

- All available memory minus 1GB
- All available cores

Understanding Processor Resources

Objectives

After watching this video, you will be able to:

- Explain how Apache Spark uses processor cores
- Describe how to use Spark to configure cores for an application

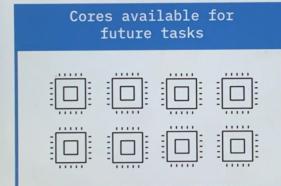
CPU Cores as a Spark Resource

- Spark assigns CPU cores to driver and executor processes



CPU Cores as a Spark Resource

- After processing, CPU cores become available for future tasks



CPU Cores as a Spark Resource

- Workers in the cluster contain a limited number of cores
- If no cores are available to an application, the application must wait for currently running tasks to finish

Setting Executor Cores on Submit

- Specifies the number of executor cores for a Spark standalone cluster per executor process

```
$ ./bin/spark-submit \
--class org.apache.spark.examples.SparkPi \
--master spark://<spark-master-URL>:7077 \
--executor-cores 8 \
/path/to/examples.jar \
1000
```

Default CPU Core Usage

- Spark queues tasks and waits for available executors and cores for maximized parallel processing
- Parallel processing tasks mainly depend on the number of data partitions and operations
- Application settings will override default behavior

Setting Executor Cores on Submit

- Specifies the executor cores for a Spark standalone cluster for the application

```
$ ./bin/spark-submit \
--class org.apache.spark.examples.SparkPi \
--master spark://<spark-master-URL>:7077 \
--total-executor-cores 50 \
/path/to/examples.jar \
1000
```

Setting Worker Node Resources

```
# Start standalone worker with MAX 10Gb memory, 8 cores
$ ./sbin/start-worker.sh \
spark://<spark-master-URL> \
--memory 10G --cores 8
```

- By default, Spark uses all available memory minus 1GB and all available cores

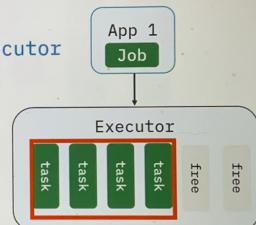
Core Utilization Example

A Spark standalone cluster with one worker node and six cores

- Submit App 1
- The application requests 4 cores per executor

```
$ ./bin/spark-submit \
--master spark://<spark-master-URL>:7077 \
--executor-cores 4 \
examples/src/main/python/pi.py \
1000
```

- App 1 occupies four cores

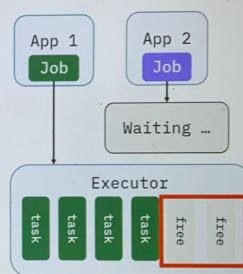


Core Utilization Example

- Submit App 2 using the same cluster
- The application requests 4 cores per executor

```
$ ./bin/spark-submit \
--master spark://<spark-master-URL>:7077 \
--executor-cores 4 \
examples/src/main/python/pi.py \
1000
```

- Two cores are available
- App 2 must wait until two additional cores are available



Core Utilization Example

- Start App 2 and request two cores

```
$ ./bin/spark-submit \
--master spark://<spark-master-URL>:7077 \
--executor-cores 2 \
examples/src/main/python/pi.py \
1000
```

- App 1 and App 2 now run simultaneously

