

# Surveillance et débogage d'Apache Spark



Temps estimé nécessaire : 30 minutes

Ce laboratoire vous expliquera comment surveiller et déboguer une application Spark via l'interface utilisateur Web.

## Objectifs

Après avoir terminé ce laboratoire, vous serez capable de :

1. Démarrez un cluster autonome Spark et connectez-vous au shell PySpark.
2. Créez un DataFrame et ouvrez l'interface utilisateur Web de l'application.
3. Débuguez une erreur d'exécution en localisant la tâche ayant échoué dans l'interface utilisateur Web.
4. Exécutez une requête SQL pour surveiller, puis augmentez la capacité en ajoutant un autre worker au cluster.

## Exercice 1 : démarrer un cluster Spark autonome

Dans cet exercice, vous allez initialiser un cluster autonome Spark avec un maître et un travailleur. Ensuite, vous allez démarrer un shell PySpark qui se connecte au cluster et ouvrir l'interface utilisateur Web de l'application Spark pour le surveiller. Nous utiliserons le terminal Theia pour exécuter des commandes et des conteneurs basés sur Docker pour lancer les processus Spark.

### Tâche A : Télécharger des exemples de données pour Spark

1. Ouvrez un terminal Theia en cliquant sur l'élément de menu Terminal -> New Terminal.
2. Utilisez la commande suivante pour télécharger l'ensemble de données que nous utiliserons dans ce laboratoire dans le conteneur exécutant Spark.

```
wget https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-BD0225EN-SkillsNetwork/labs/data/cars.csv
```

### Tâche B : Initialiser le cluster

1. Arrêtez tous les conteneurs précédemment en cours d'exécution avec la commande :

```
for i in `docker ps | awk '{print $1}' | grep -v CONTAINER`; do docker kill $i; done
```

2. Retirez tous les conteneurs précédemment utilisés :

ignorez les erreurs indiquant « Aucun conteneur de ce type »

```
docker rm spark-master spark-worker-1 spark-worker-2
```

3. Démarrez le serveur Spark Master :

```
docker run \
  --name spark-master \
  -h spark-master \
  -e ENABLE_INIT_DAEMON=false \
  -p 4040:4040 \
  -p 8080:8080 \
  -v `pwd`: /home/root \
  -d bde2020/spark-master:3.1.1-hadoop3.2
```

4. Démarrez un Spark Worker qui se connectera au Master :

```
docker run \
  --name spark-worker-1 \
  --link spark-master:spark-master \
  -e ENABLE_INIT_DAEMON=false \
  -p 8081:8081 \
  -v `pwd`: /home/root \
  -d bde2020/spark-worker:3.1.1-hadoop3.2
```

### Tâche C : connecter un shell PySpark au cluster et ouvrir l'interface utilisateur

1. Lancez un shell PySpark dans le conteneur Spark Master en cours d'exécution :

```
docker exec \
  -it `docker ps | grep spark-master | awk '{print $1}'` \
  /spark/bin/pyspark \
  --master spark://spark-master:7077
```

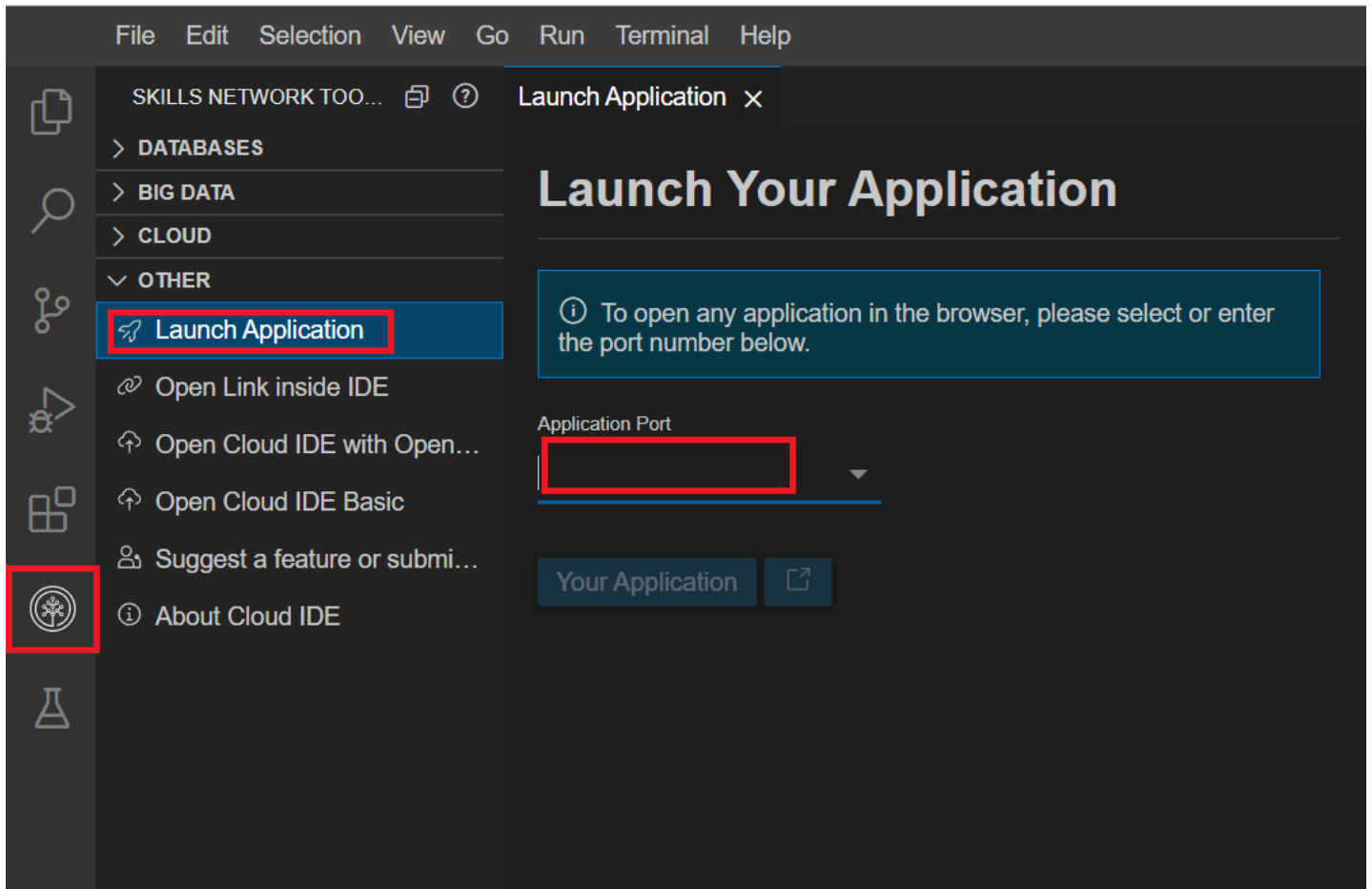
2. Créez un DataFrame dans le shell avec :

**REMARQUE** : appuyez deux fois sur Entrée pour continuer après avoir exécuté la commande dans le shell

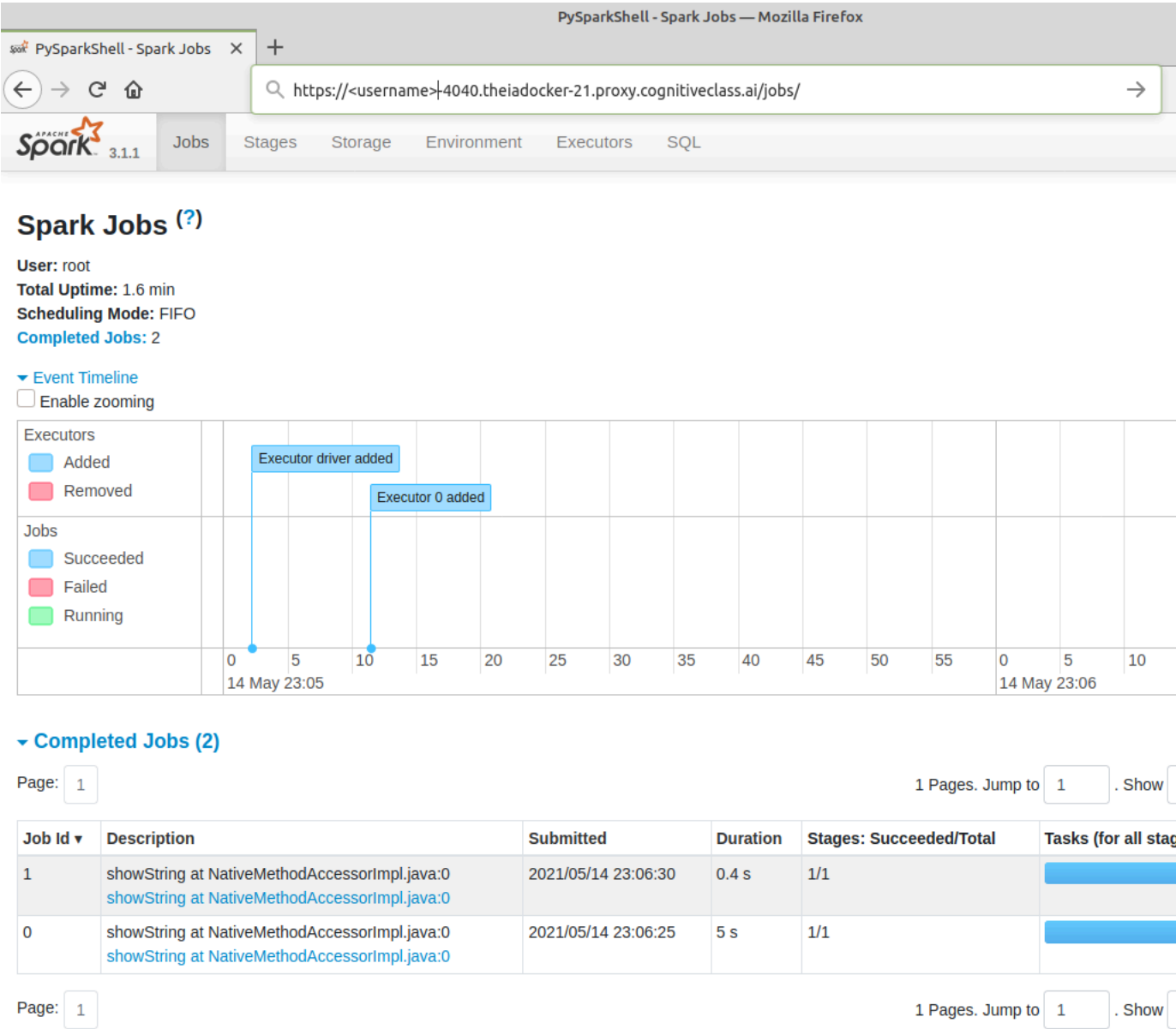
```
df = spark.read.csv("/home/root/cars.csv", header=True, inferSchema=True) \
```

```
.repartition(32) \
.cache()
df.show()
```

3. Cliquez sur le bouton Skills Network à gauche, cela ouvrira la « Skills Network Toolbox ». Cliquez ensuite OTHER sur Launch Application. À partir de là, vous devriez pouvoir saisir le numéro de port 4040 et lancer l'interface utilisateur de l'application Spark dans votre navigateur.



4. Vérifiez que vous pouvez voir la page des offres d'emploi qui devrait ressembler à ce qui suit, bien que pas nécessairement exactement la même :



## Exercice 2 : Exécuter une requête SQL et déboguer dans l'interface utilisateur de l'application

Dans cet exercice, vous allez définir une fonction définie par l'utilisateur (UDF) et exécuter une requête qui génère une erreur. Nous localiserons cette erreur dans l'interface utilisateur de l'application et en trouverons la cause première. Enfin, nous corrigerons l'erreur et réexécuterons la requête.

### Tâche A : Exécuter une requête SQL

1. Définissez un UDF pour afficher le type de moteur. Copiez et collez le code et cliquez sur Enter.

```
from pyspark.sql.functions import udf
import time
@udf("string")
def engine(cylinders):
    time.sleep(0.2) # Intentionally delay task
    eng = {6: "V6", 8: "V8"}
    return eng[cylinders]
```

2. Ajoutez l'UDF en tant que colonne dans le DataFrame

```
df = df.withColumn("engine", engine("cylinders"))
```

3. Regrouper le DataFrame par « cylindres » et agréger les autres colonnes

```
dfg = df.groupby("cylinders")
dfa = dfg.agg({"mpg": "avg", "engine": "first"})
dfa.show()
```

4. La requête a échoué et vous devriez voir de nombreux messages et résultats dans la console.  
La tâche suivante consiste à localiser l'erreur dans l'interface utilisateur de l'application et à déterminer la cause première.

# Tâche B : Débuguer l'erreur dans l'interface utilisateur de l'application

1. Recherchez l'erreur dans l'interface utilisateur de l'application.  
Ouvrez l'interface utilisateur des tâches, regardez la liste des tâches ayant échoué, cliquez sur la première tâche.

APACHE

Spark

3.2.0-SNAPSHOT

Jobs

Stages

Storage

Environment

Executors

SQL

PySparkShell  
application UI

Spark Jobs (?)

User: root

Total Uptime: 1.4 min

Scheduling Mode: FIFO

Completed Jobs: 4

▶ Event Timeline

▼ Completed Jobs (4)

Page: 11 Pages. Jump to 1. Show 100 items in a page. Go

Job Id ▼	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
3	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> <a href="#">showString at NativeMethodAccessorImpl.java:0</a>	2021/07/07 22:08:48	39 ms	1/1 (1 skipped)	1/1 (1 skipped)
2	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> <a href="#">showString at NativeMethodAccessorImpl.java:0</a>	2021/07/07 22:08:47	0.5 s	2/2	2/2
1	<a href="#">csv at NativeMethodAccessorImpl.java:0</a> <a href="#">csv at NativeMethodAccessorImpl.java:0</a>	2021/07/07 22:08:46	1 s	1/1	1/1
0	<a href="#">csv at NativeMethodAccessorImpl.java:0</a> <a href="#">csv at NativeMethodAccessorImpl.java:0</a>	2021/07/07 22:08:44	1 s	1/1	1/1

Page: 11 Pages. Jump to 1. Show 100 items in a page. Go

2. Cela affichera les détails du travail avec une liste des étapes pour ce travail. Dans la liste des étapes ayant échoué, cliquez sur la première étape ayant échoué pour afficher les détails de l'étape avec une liste des tâches pour cette étape.

	<a href="#">+details</a>							
2	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> <a href="#">+details</a>	2021/07/07 22:08:47	0.2 s	1/1	20.6 KiB			23.4 KiB
1	<a href="#">csv at NativeMethodAccessorImpl.java:0</a> <a href="#">+details</a>	2021/07/07 22:08:46	1 s	1/1	20.6 KiB			
0	<a href="#">csv at NativeMethodAccessorImpl.java:0</a> <a href="#">+details</a>	2021/07/07 22:08:44	1 s	1/1	20.6 KiB			

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Skipped Stages (2)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
6	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> <a href="#">+details</a>	Unknown	Unknown	0/1				
4	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> <a href="#">+details</a>	Unknown	Unknown	0/1				

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Failed Stages (1)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write	Failure Reason
7	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> <a href="#">+details</a>	2021/07/07 22:15:19	2 s	0/32 (28 failed) (6	45.5 KiB		8.1 KiB		Job aborted due to stage 7.0 failed 4 tasks. Lost task 3.3 in stage (192.168.1.33 executor org.apache.spark.scheduler.Task) Traceback (most recent

3. Ici, nous voyons de nombreuses tâches échouées. En regardant la première, la colonne la plus à droite montre les détails de l'échec.

Details for Stage 7 (Attempt 0)

Resource Profile Id: 0  
Total Time Across All Tasks: 17 s  
Locality Level Summary: Node local: 28; Process local: 6  
Input Size / Records: 45.5 KiB / 22  
Shuffle Read Size / Records: 8.1 KiB / 136  
Associated Job Ids: 4

- DAG Visualization
- Show Additional Metrics
- Event Timeline

Summary Metrics for 0 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
--------	-----	-----------------	--------	-----------------	-----

Aggregated Metrics by Executor

Tasks (34)

Show20entries

Search:

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Input Size / Records
0	5	0	FAILED	PROCESS_LOCAL	0	192.168.1.33	<a href="#">stdout</a> <a href="#">stderr</a>	2021-07-07 22:15:19	1 s		2.1 KiB / 1
0	20	1	FAILED	PROCESS_LOCAL	0	192.168.1.33	<a href="#">stdout</a> <a href="#">stderr</a>	2021-07-07 22:15:20	0.4 s		2.1 KiB / 1
0	31	2	KILLED	PROCESS_LOCAL	0	192.168.1.33	<a href="#">stdout</a> <a href="#">stderr</a>	2021-07-07 22:15:20	0.4 s		2.1 KiB / 1
1	6	0	FAILED	PROCESS_LOCAL	0	192.168.1.33	<a href="#">stdout</a> <a href="#">stderr</a>	2021-07-07 22:15:19	1 s	19.0 ms	2 KiB / 1
1	23	1	FAILED	PROCESS_LOCAL	0	192.168.1.33	<a href="#">stdout</a> <a href="#">stderr</a>	2021-07-07 22:15:20	0.6 s		2 KiB / 1
1	37	2	KILLED	PROCESS_LOCAL	0	192.168.1.33	<a href="#">stdout</a> <a href="#">stderr</a>	2021-07-07 22:15:20	0.6 s		2 KiB / 1

Cliquez pour développer les détails.

APACHE

spark

3.2.0-SNAPSHOT

stderr log page for app-20210707220734-0000/0

[Back to Master](#)

Showing 102400 Bytes: 39045 - 141445 of 141445

```
return f(*args, **kwargs)
File "<ipython-input-2-12484af4df99>", line 8, in engine
KeyError: 4

at
org.apache.spark.api.python.BasePythonRunner$ReaderIterator.handlePythonException(PythonRunner.scala:517)
at org.apache.spark.sql.execution.python.PythonUDFRunner$$anon$2.read(PythonUDFRunner.scala:84)
at org.apache.spark.sql.execution.python.PythonUDFRunner$$anon$2.read(PythonUDFRunner.scala:67)
at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.hasNext(PythonRunner.scala:470)
at org.apache.spark.InterruptibleIterator.hasNext(InterruptibleIterator.scala:37)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:489)
at scala.collection.Iterator$$anon$10.hasNext(Iterator.scala:458)
at scala.collection.Iterator$$anon$10.hasNext(Iterator.scala:458)
at
org.apache.spark.sql.catalyst.expressions.GeneratedClass$GeneratedIteratorForCodegenStage1.sort_addToSorter_0$(Unknown Source)
at
org.apache.spark.sql.catalyst.expressions.GeneratedClass$GeneratedIteratorForCodegenStage1.processNext(Unknown Source)
```

Faites défiler un peu vers le bas jusqu'à ce que vous puissiez voir la dernière partie de l'erreur Python qui montre la cause. Vous devriez pouvoir voir que cela a été causé par une KeyError dans notre UDF engine().

Vous pouvez également visualiser ces erreurs en regardant la colonne qui contient des liens vers les journaux et en cliquant sur « std err » pour afficher le journal des erreurs standard.

Fermez l'onglet du navigateur PySpark.

4. Dans le terminal Theia, corrigez l'UDF en ajoutant une entrée au dictionnaire des types de moteurs et fournissez une valeur par défaut pour tous les autres types. Copiez et collez ce code et cliquez sur Enter.

```
@udf("string")
def engine(cylinders):
    time.sleep(0.2) # Intentionally delay task
    eng = {4: "inline-four", 6: "V6", 8: "V8"}
    return eng.get(cylinders, "other")
```

5. Relancez la requête. Vous devrez ajouter à nouveau la colonne « moteur » et saisir la requête puisque nous avons modifié l'UDF.

```
df = df.withColumn("engine", engine("cylinders"))

dfg = df.groupby("cylinders")

dfa = dfg.agg({"mpg": "avg", "engine": "first"})

dfa.show()
```

Une fois la requête terminée sans erreur, vous devriez voir un résultat similaire à celui-ci.

```
+-----+-----+-----+
|cylinders|      avg(mpg)|first(engine)|
+-----+-----+-----+
|        6|19.985714285714288|          V6|
|        3|          20.55|          other|
|        5|27.366666666666664|          other|
|        4|29.286764705882348|    inline-four|
|        8|14.963106796116506|          V8|
+-----+-----+-----+
```

## Exercice 3 : Surveiller les performances des applications avec l'interface utilisateur

Maintenant que nous avons exécuté notre requête avec succès, nous allons faire évoluer notre application en ajoutant un worker au cluster. Cela permettra au cluster d'exécuter davantage de tâches en parallèle et d'améliorer les performances globales.

### Tâche A : Ajouter un travailleur au cluster

1. Affichez l'onglet Étapes, puis cliquez sur l'étape comportant 32 tâches. À ce stade, notre UDF est appliqué à chaque partition du DataFrame.

Stages for All Jobs

Completed Stages: 8  
Skipped Stages: 7  
Failed Stages: 1

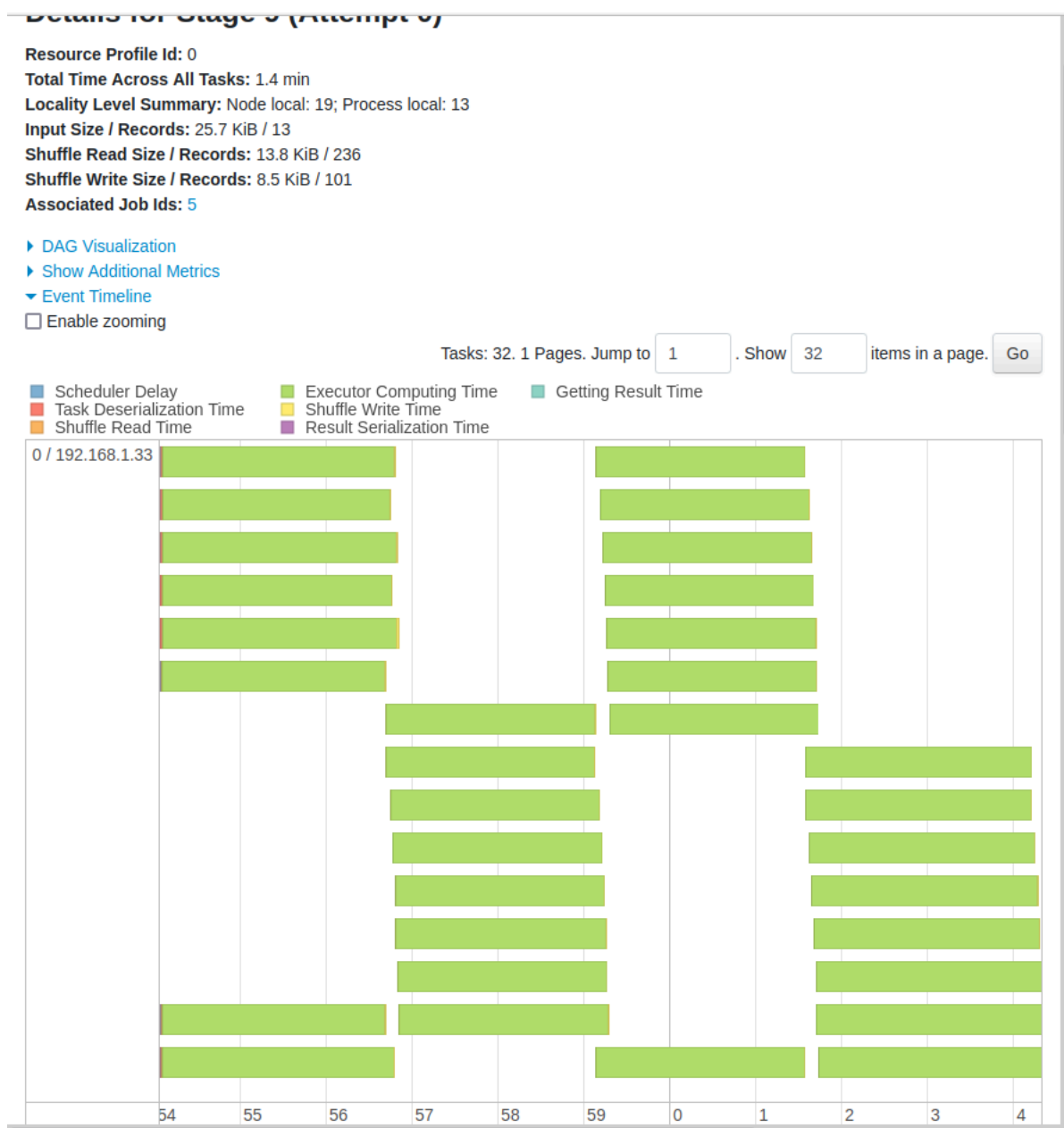
Completed Stages (8)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Stage Id ▾	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
15	showString at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:39:04	47 ms	3/3			6.0 KiB	
12	showString at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:39:04	0.2 s	1/1			2.5 KiB	
9	showString at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:38:54	10 s	32/32	25.7 KiB		13.8 KiB	8.5 KiB
5	showString at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:08:48	33 ms	1/1			778.0 B	
3	showString at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:08:48	0.3 s	1/1			829.0 B	
2	showString at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:08:47	0.2 s	1/1	20.6 KiB			23.4 KiB
1	csv at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:08:46	1 s	1/1	20.6 KiB			
0	csv at NativeMethodAccessorImpl.java:0 <a href="#">+details</a>	2021/07/07 22:08:44	1 s	1/1	20.6 KiB			

En regardant la chronologie, vous pouvez voir qu'il existe un seul travailleur avec un identifiant 0 / <ip-address>qui peut exécuter jusqu'à un certain nombre de tâches en parallèle à la fois. L'ajout d'un autre travailleur permettra d'exécuter des tâches supplémentaires en parallèle.





2. Ouvrez un nouveau terminal Theia en cliquant sur l'élément de menu Terminal -> New Terminal.

3. Ajoutez un deuxième travailleur au cluster avec la commande dans le nouveau terminal :

```
docker run \
  --name spark-worker-2 \
  --link spark-master:spark-master \
  -e ENABLE_INIT_DAEMON=false \
  -p 8082:8082 \
  -d bde2020/spark-worker:3.1.1-hadoop3.2
```

4. Si la commande réussit, il y aura une sortie unique indiquant l'ID du conteneur :

```
theia@theiadocker-user:/home/project$ docker run \
> --name spark-worker-2 \
> --link spark-master:spark-master \
> -e ENABLE_INIT_DAEMON=false \
> -p 8082:8082 \
> -d bde2020/spark-worker:3.1.1-hadoop3.2
1935a71827668ae3476e6a16f0bebcd4c2a342a21271dc22be487aa1b1731708
theia@theiadocker-user:/home/project$
```

5. Cliquez sur le premier terminal sur lequel le shell PySpark est ouvert pour continuer.

## Tâche B : Réexécuter la requête et vérifier les performances

1. Réexécutez la requête, cette fois nous pouvons simplement appeler `show()` à nouveau :

```
dfa.show()
```

2. Lancez l'application sur le numéro de port 4040en suivant le même processus que ci-dessus, pour ouvrir le navigateur PySpark. Accédez à l'onglet **Étapes** et consultez l'ID de l'étape la plus récente.



3.1.1

Jobs

Stages

Storage

Environment

Exe

## Stages for All Jobs

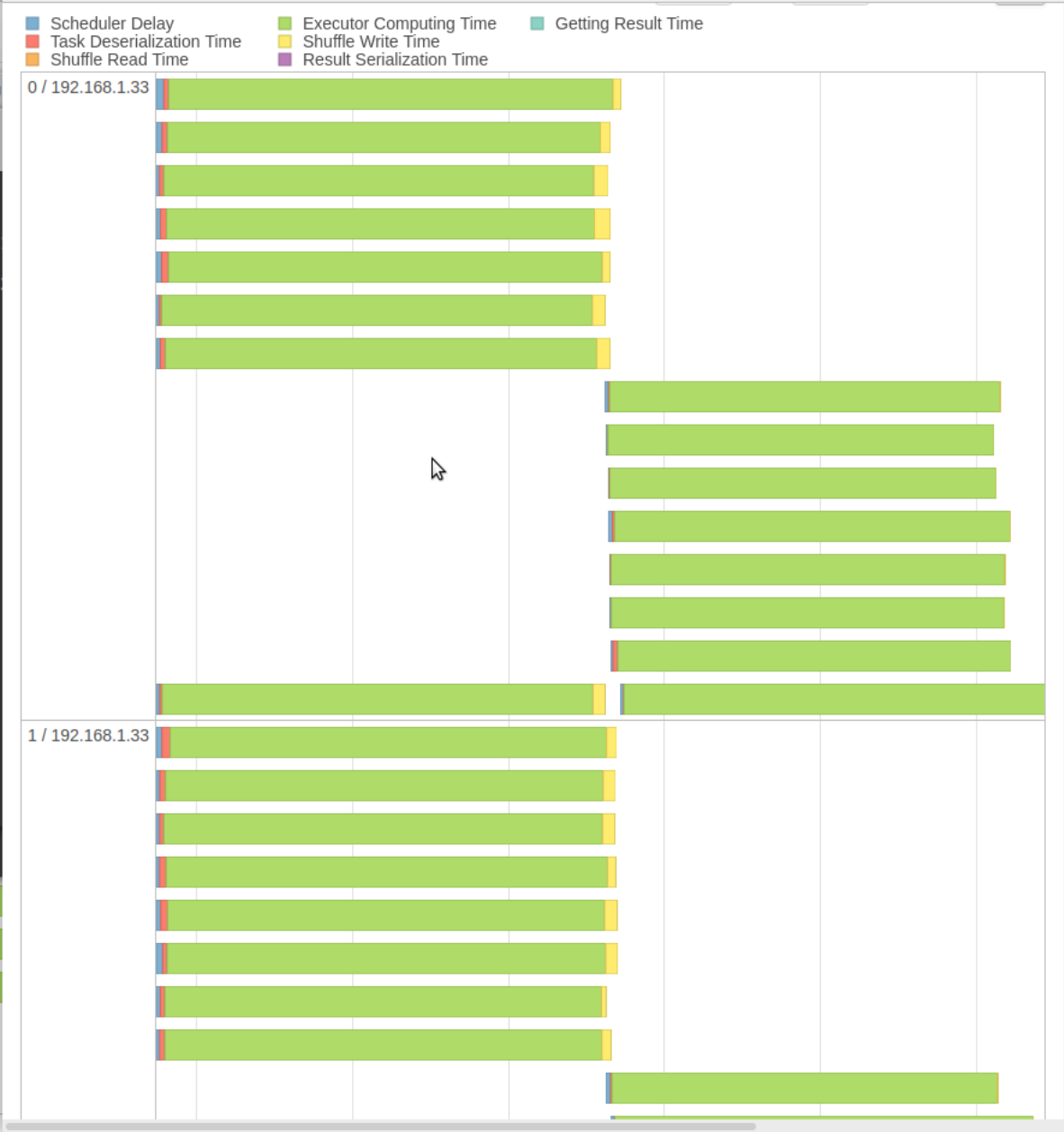
**Completed Stages:** 23**Skipped Stages:** 30**Failed Stages:** 1

### ▼ Completed Stages (23)

Page: 1

Stage Id ▼	Description	Submitted
53	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> +details	2022/01/03 06:27:21
50	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> +details	2022/01/03 06:27:19
47	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> +details	2022/01/03 06:27:18
44	<a href="#">showString at NativeMethodAccessorImpl.java:0</a> +details	2022/01/03 06:27:17

3. Vous verrez que le travailleur supplémentaire avec l'ID 1 / <ip-address> est répertorié et permet désormais d'exécuter davantage de tâches en parallèle. La chronologie des tâches doit ressembler à ce qui suit.



**Auteur(s)**

Aije

**Autre(s) contributeur(s)**

Lavande

© IBM Corporation 2021. Tous droits réservés.