AWS re:Invent

ANT332

# METRICS-DRIVEN PERFORMANCE TUNING FOR AWS GLUE ETL JOBS
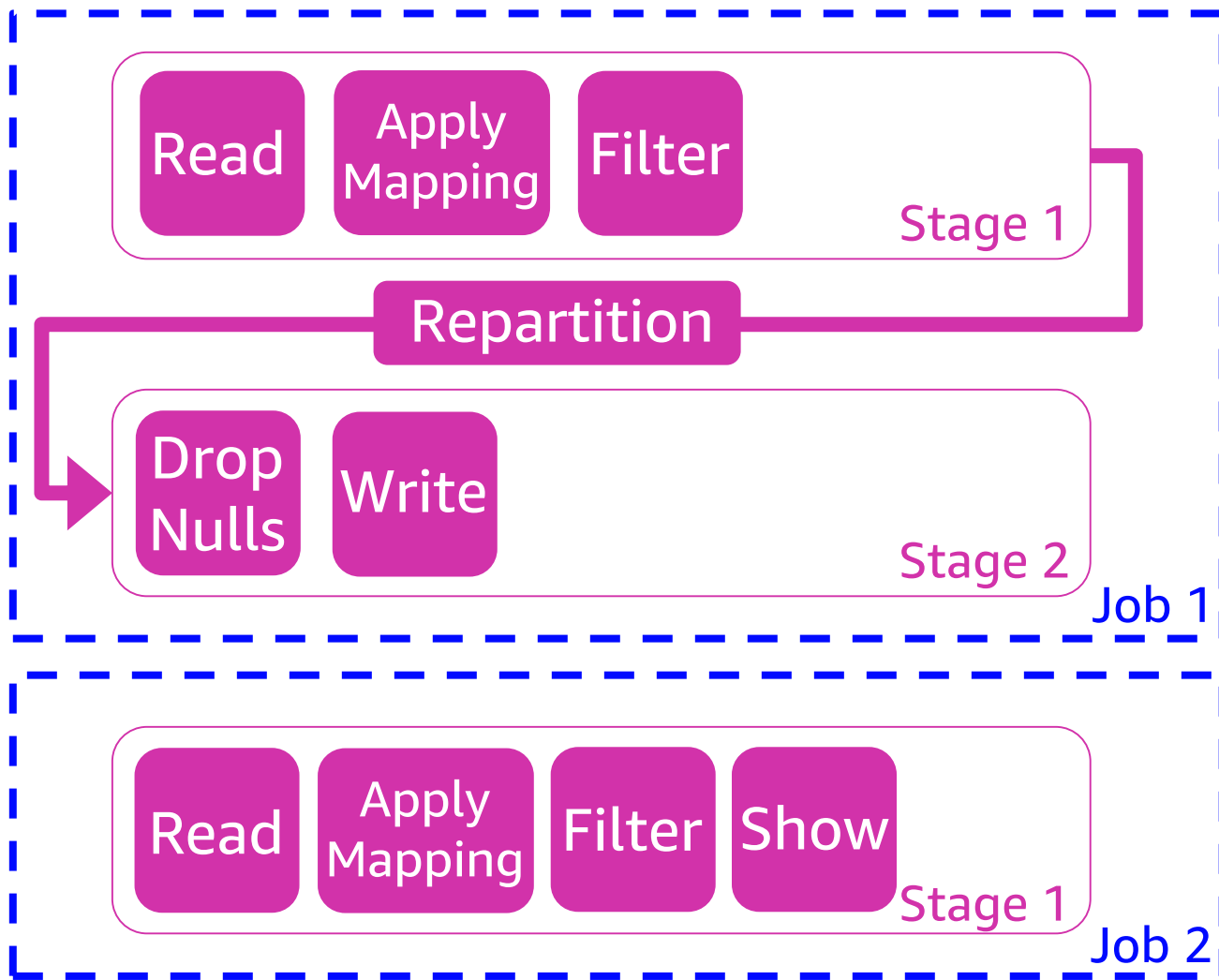
Benjamin Sowell
Principal Engineer
AWS Glue

aws

# Apache Spark and AWS Glue ETL

| | | |
|---|---|---|
| SparkSQL | AWS Glue ETL | Application |
| Spark DataFrames | AWS Glue DynamicFrames | Data Structure |
| Spark Core: RDDs | | Execution |

- Apache Spark is a distributed data processing engine with rich support for complex analytics.
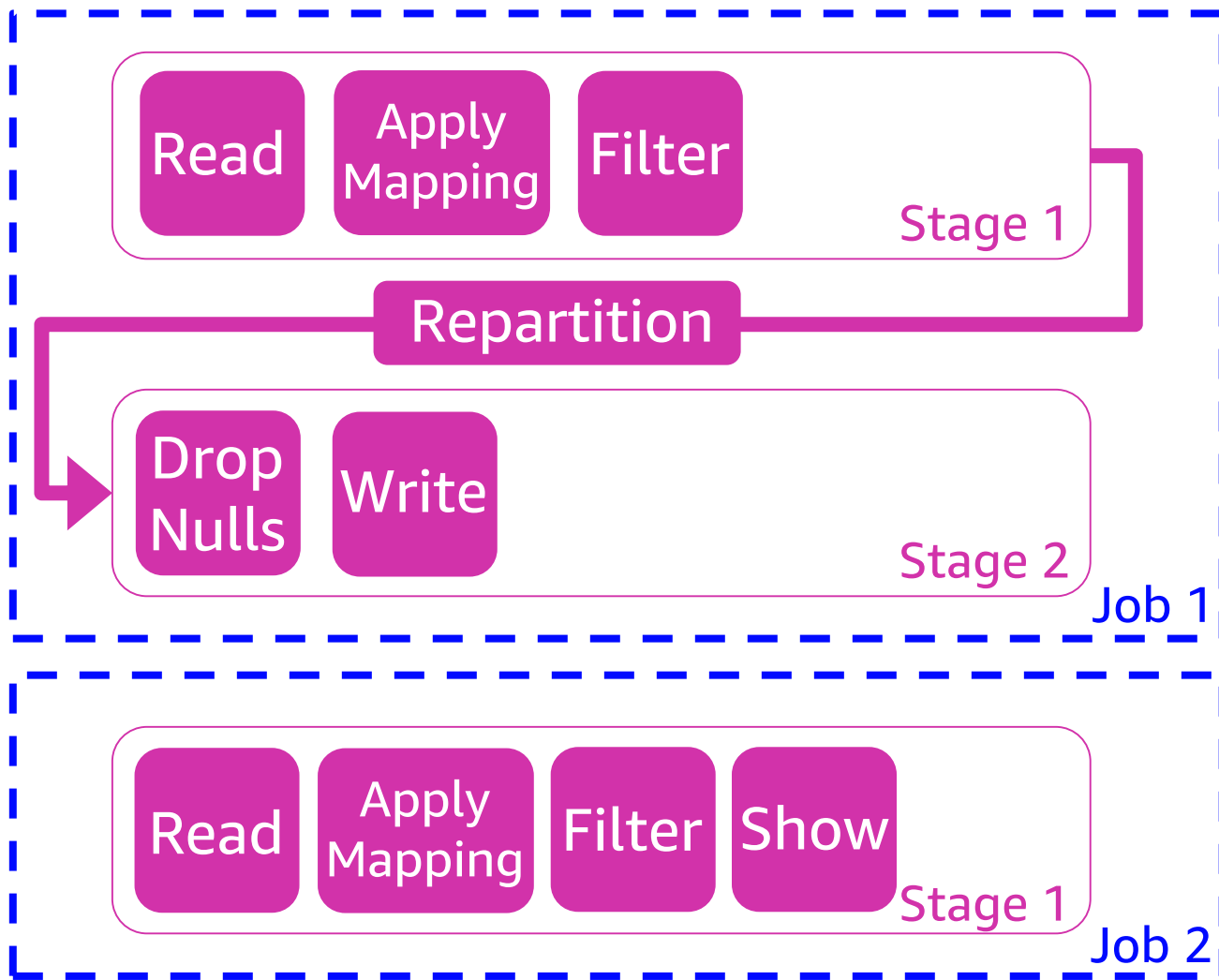- AWS Glue builds on the Apache Spark runtime to offer ETL specific functionality.

AWS re:Invent

aws

# AWS Glue execution model: jobs and stages



```
df = glueContext.getSource(…)

applyMapping = df.applyMapping(…)

filter = applyMapping.filter(…)

repartition = filter.repartition(10)

dropNulls = repartition.dropNulls()

glueContext.getSink(…)\
    .WriteDynamicFrame(dropNulls)


filter.show()
```
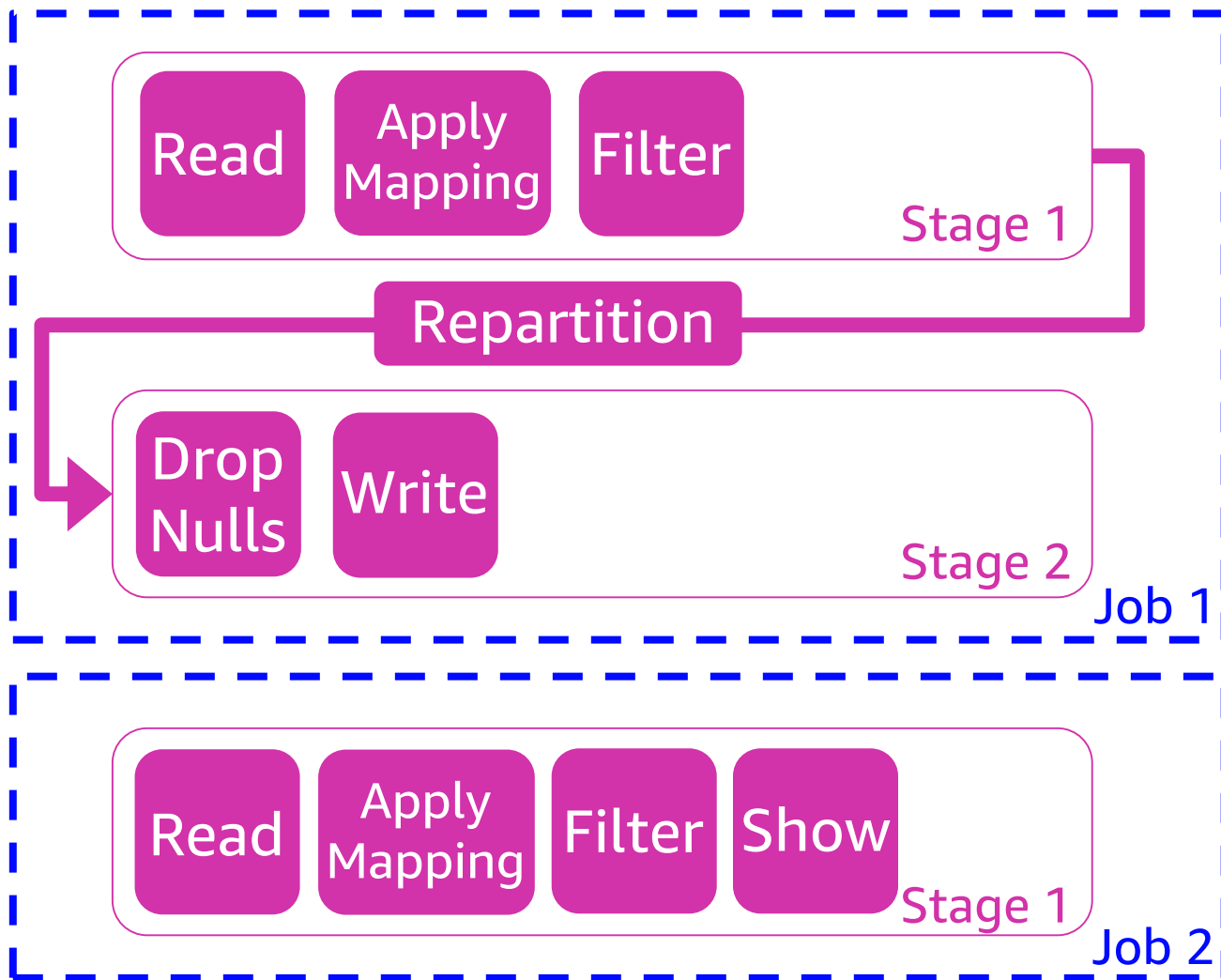
# AWS Glue execution model: jobs and stages



```
df = glueContext.getSource(…)

applyMapping = df.applyMapping(…)

filter = applyMapping.filter(…)

repartition = filter.repartition(10)

dropNulls = repartition.dropNulls()

glueContext.getSink(…)\
    .WriteDynamicFrame(dropNulls)

filter.show()
```
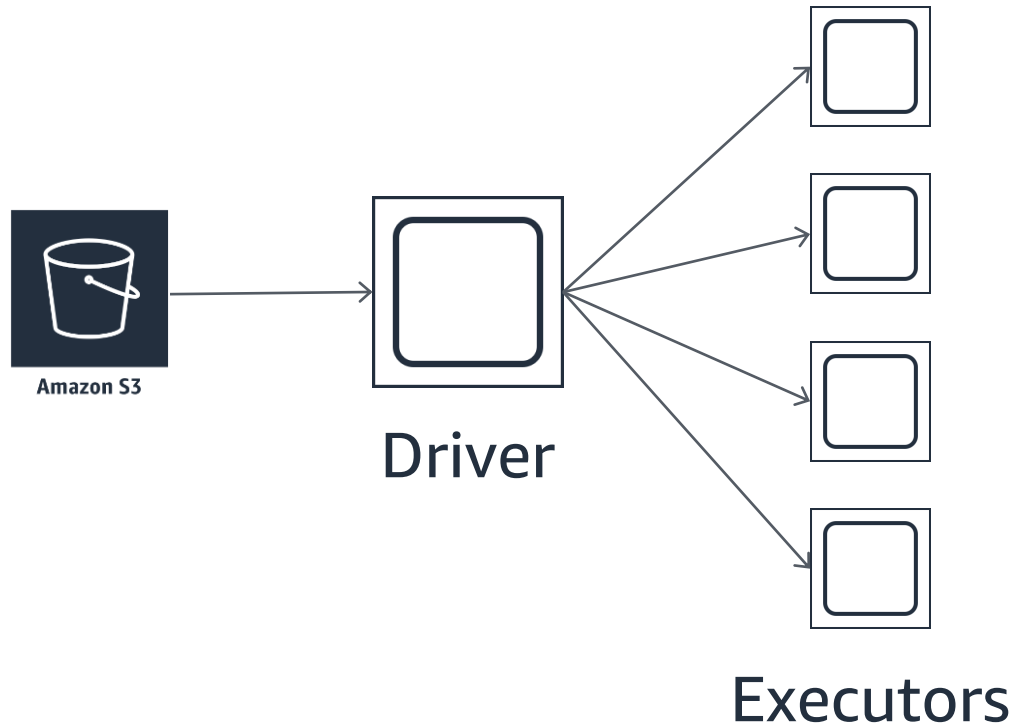
**Actions**

# AWS Glue execution model: jobs and stages



```
df = glueContext.getSource(…)

applyMapping = df.applyMapping(…)

filter = applyMapping.filter(…)

repartition = filter.repartition(10)

dropNulls = repartition.dropNulls()

glueContext.getSink(…)\
    .WriteDynamicFrame(dropNulls)
```

```
filter.show()
```

Jobs

# AWS Glue execution model: data partitions



Driver

Executors

- Apache Spark and AWS Glue are *data parallel*.
- Data is divided into *partitions* that are processed concurrently.
- 1 stage x 1 partition = 1 *task*

Overall throughput is limited by the number of partitions

# AWS Glue performance: key questions

How is your application divided into jobs and stages?

How is your dataset partitioned?
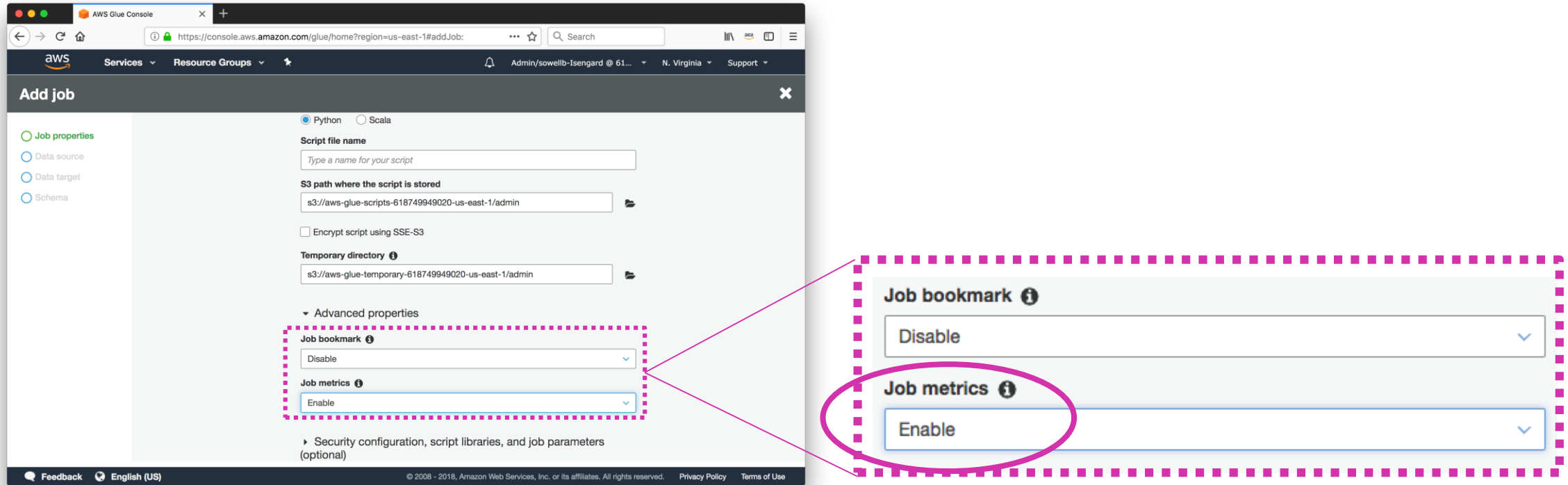
# Example: Processing lots of small files

- One common problem is dealing with large numbers of small files.
  - Can lead to memory pressure and performance overhead.
- Let's look at a straightforward JSON to Parquet conversion job
  - 1.28 million JSON files in 640 partitions:

```
part1/                          part640/
        file1.json                      file1.json
        …                               …
        file2000.json                   file2000.json
```

- We will use AWS Glue *job metrics* to understand the performance.

# Enabling job metrics



- Metrics can be enabled in the CLI/SDK by passing `--enable-metrics` as a job parameter key.

# Example: Processing lots of small files

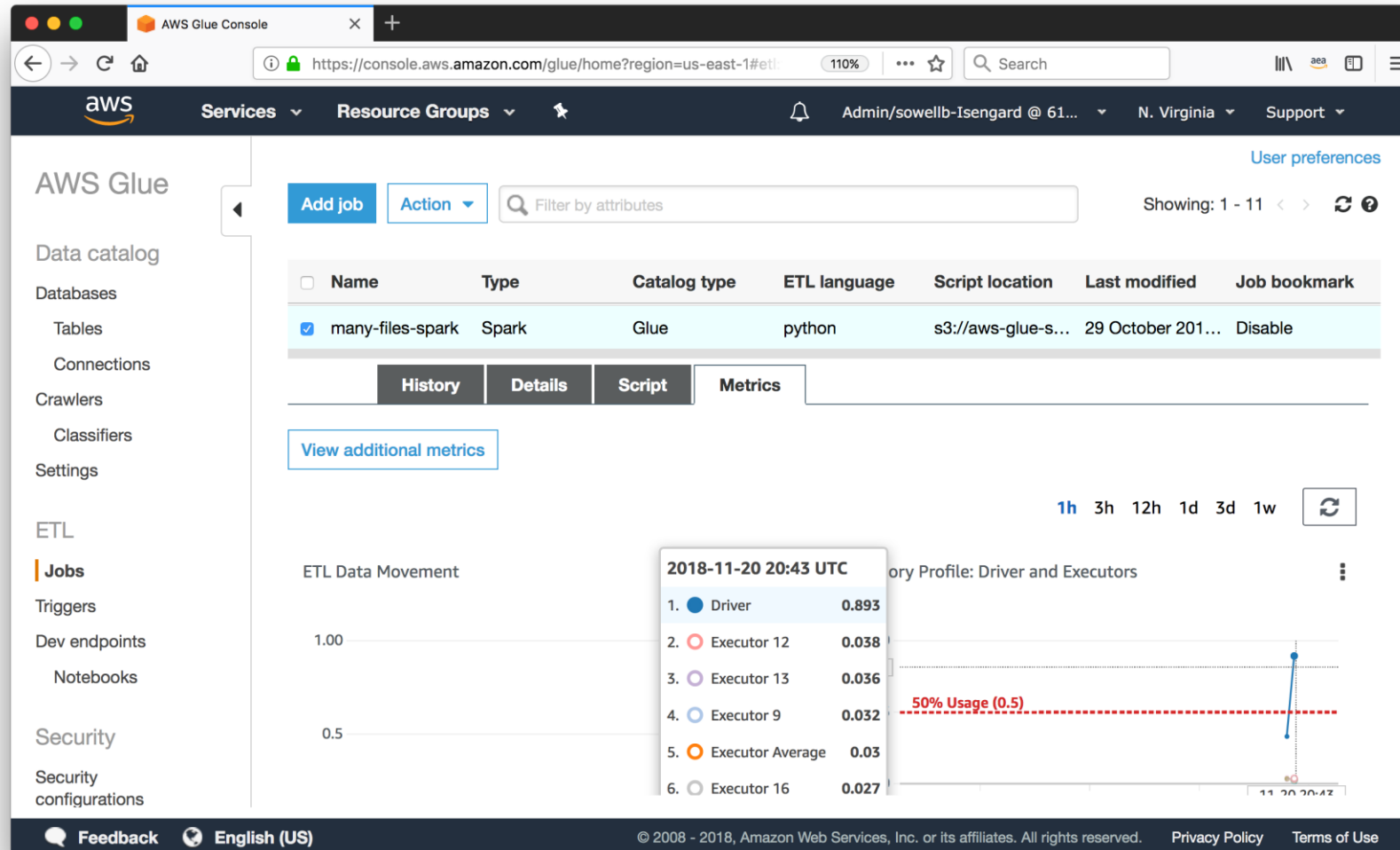- **First try: use a standard SparkSQL job**

```
data = spark.read.format("json").load("<input_location>")
data.write.format("parquet").save("<output_location>")
```

| | History | Details | Script | Metrics | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

View run metrics    ‹ › ⟳      Showing: 1 - 1 ‹ › ⟳

| Run ID | Retry attempt | Run status | Error | Logs | Error logs | Execution time | Timeout | Delay | Triggered by | Start time | End time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ○ jr_baf7b7789ba... | - | Failed | ⚠ Co... | Logs | Error logs | 2 mins | 2880 mins | | | 29 October ... | 29 October ... |

Command failed with exit code 1

aws

# Example: Processing lots of small files

# Example: Processing lots of small files

# Example: Processing lots of small files



- Driver memory use is growing fast and approaching the 5g max.

# Example: Processing lots of small files

- Case 2: Run using AWS Glue DynamicFrames.

```
df = glueContext.create_dynamic_frame_from_catalog("<database>","<table>")

glueContext.write_from_options(df,"s3",{"path":"<output_location>"},"parquet")
```

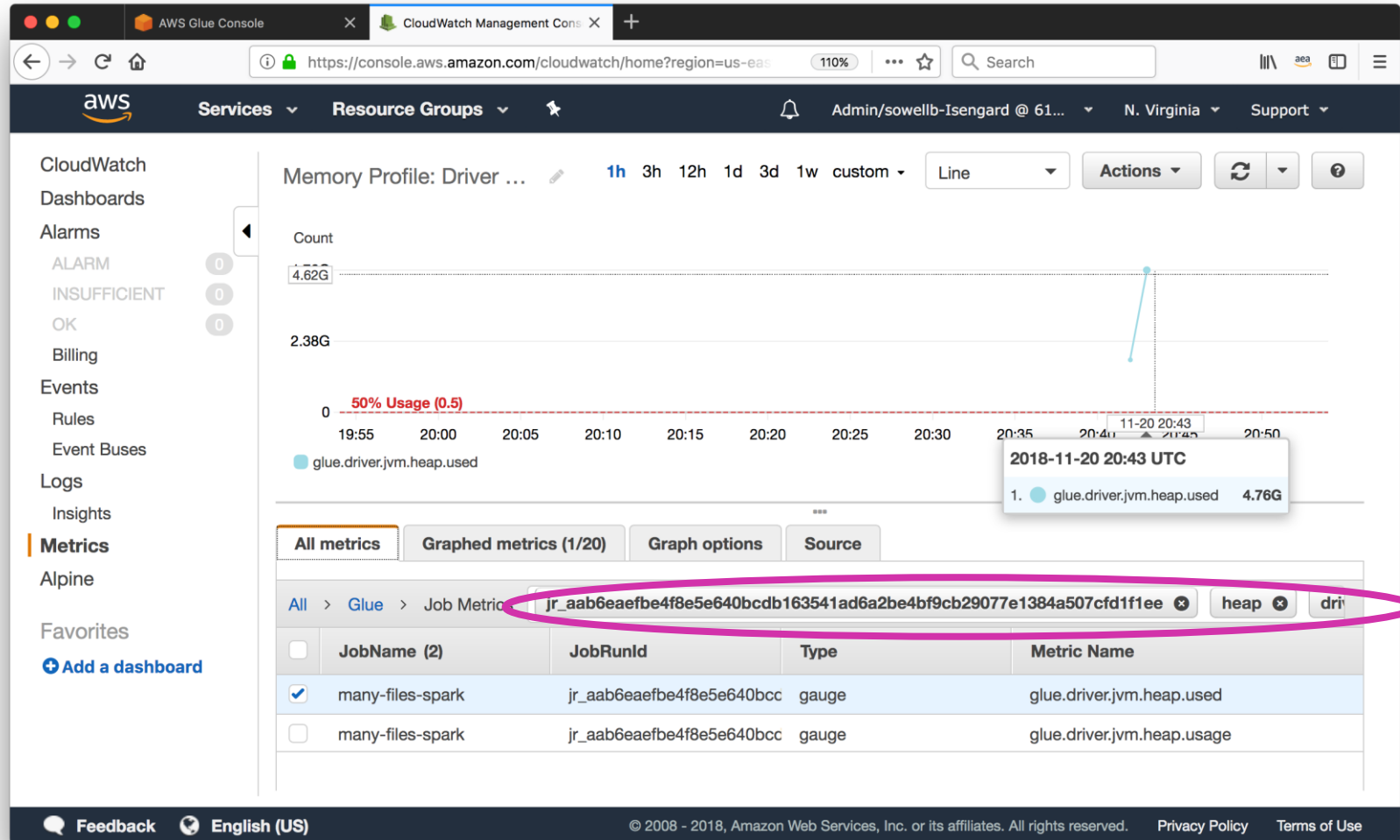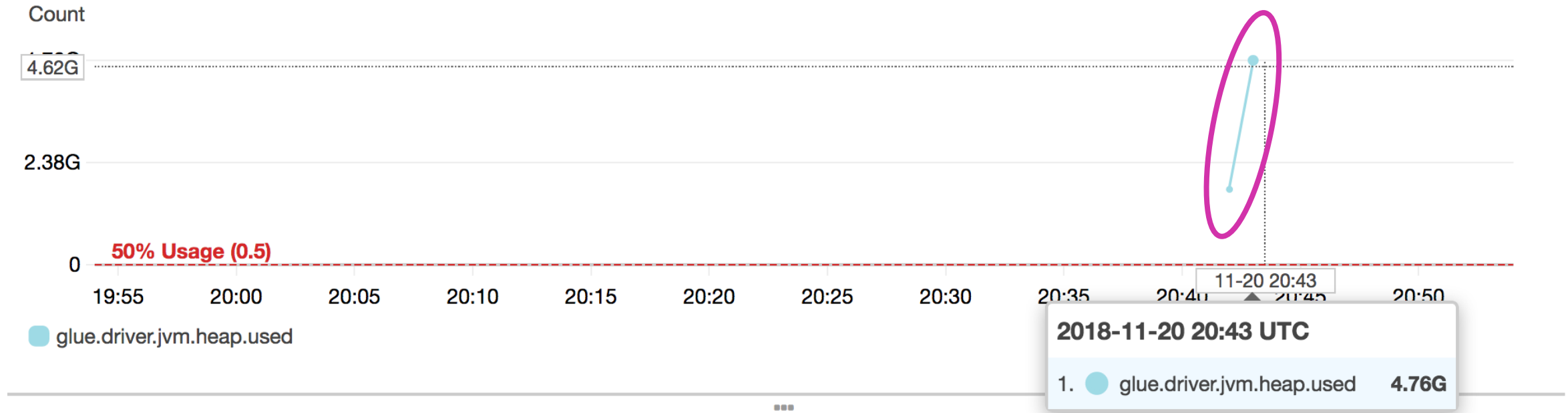| | History | Details | Script | Metrics |
|---|---|---|---|---|

View run metrics    Showing: 1 - 1

| Run ID | Retry attempt | Run status | Error | Logs | Error logs | Execution time | Timeout | Delay | Triggered by | Start time | End time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| jr_671c479fa0... | - | Succeeded | | Logs | | 59 mins | 2880 mins | | | 31 Octobe... | 31 Octobe... |

# Example: Processing lots of small files

# Example: Processing lots of small files



Memory Profile: Driver and Executors

**Driver memory remains below 50% for the entire duration of execution.**

1.00

**50% Usage (0.5)**

0.5

0

19:30   20:00   20:30   21:00   21:30   22:00

# Example: Processing lots of small files

# Example: Processing lots of small files



- **160 max executors: Each partition is assigned 1 task. 4 tasks can be processed on each executor, so there were 640 partitions.**
  - Glue automatically grouped all of the files in each partition.

# Options for grouping files

- groupFiles
  - Set to "inPartition" to enable grouping of files within a partition.
  - Set to "acrossPartition" to enable grouping of files from different partitions. The partition value will not be added to each record!
  - Grouping is automatically enabled when there are more than 50,000 files.

- groupSize
  - Target size of each group in bytes.
  - Default is based on the number of cores in the cluster.

- Let's try increasing the group size.

# Example: Aggressively grouping files

- ## Execution is much slower, but hasn't crashed.

```
glueContext.create_dynamic_frame_from_catalog(
        database = "many_files_dataset",
        table_name="json_2k_million",
        additional_options={"groupSize": 1024 * 1024 * 1024},
                            "groupFiles": "acrossPartition")
```

| History | Details | Script | Metrics |
| --- | --- | --- | --- |

View run metrics    ↻         Showing: 1 - 2 ‹ › ↻

| Run ID | Retry attempt | Run status | Error | Logs | Error logs | Execution time | Timeout | Delay | Triggered by | Start time | End time |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| ○ jr_4b6e297a0f... | - | Running ✕ | | Logs | Error logs | 18 hrs, 53 mins | 2880 mins | | | 31 Octobe... | |
| ○ jr_671c479fa0... | - | Succeeded | | Logs | | 59 mins | 2880 mins | | | 31 Octobe... | 31 Octobe... |

# Example: Aggressively grouping files

Memory Profile: Driver and Executors ✏

1h  **3h**  12h  1d  3d  1w  custom ▾   | Line ▾ |  | Actions ▾ |  ⟳ ▾  ❓

Count
1.00

Executor memory is higher than driver.          Only one executor is active.

0.741

50% Usage (0.5)
0.5

0

19:00    19:15    19:30    19:45    20:00    20:15    20:30    11-01 20:45    21:00    21:15    21:30    21:45

■ Driver  ■ Executor Average  ■ Executor 1  ■ Executor 2  ■ Executor 3  ■ Executor 4  ■ Executor 5  ■ Executor 6    ■ Executor 9  ■ Executor 10  ■ Executor 11
■ Executor 12  ■ Executor 13  ■ Executor 14  ■ Executor 15  ■ Executor 16  ■ Executor 17

**2018-11-01 20:46 UTC**

| | | |
|---|---|---|
| 1. ⭕ Executor 1 | 0.43 |
| 2. ⭕ Executor Average | 0.43 |
| 3. ⭕ Driver | 0.354 |
| 4. ⭕ Executor 2 | - |
| 5. ⭕ Executor 3 | - |
| 6. ⭕ Executor 4 | - |
| 7. ⭕ Executor 5 | - |
| 8. ⭕ Executor 6 | - |

| All metrics | **Graphed metrics (19)** | Graph options | Source |

➕ Add a math expression ❓                                      age ▾   **Period: 1 Minute ▾**   **Remove all**

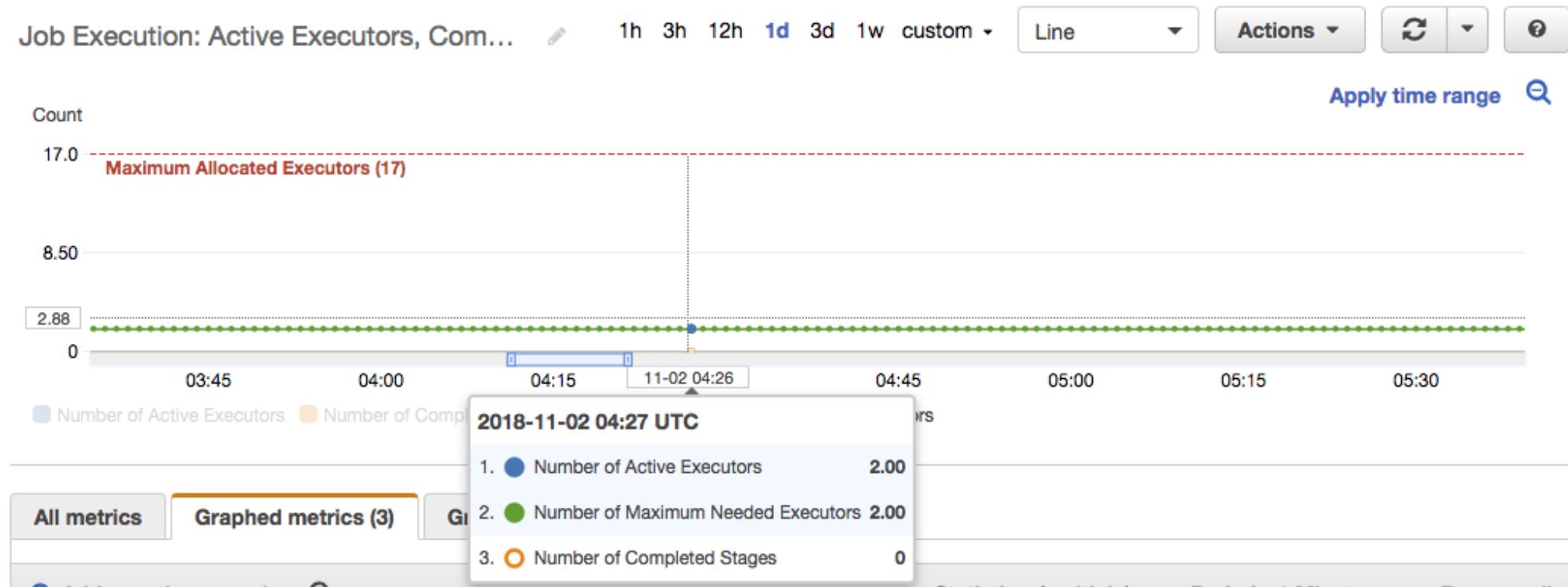| ✔ | | Label | Details | | Period | Y Axis | Actions |
|---|---|---|---|---|---|---|---|
| ✔ | ■ | Driver | Glue • glue.driver.jvm.heap.usage • Type: gauge • Jc | | 1 Minute | ‹ › | 🔔 ⎘ ✖ |
| ✔ | ■ | Executor Average | Glue • glue.ALL.jvm.heap.usage • Type: gauge • Jot | | 1 Minute | ‹ › | 🔔 ⎘ ✖ |

# Example: Processing a few large files

- Processing large text files can also cause problems if they cannot be split.
- How the data is compressed makes a big difference in performance.
  - Files that are *uncompressed* or *bzip2* compressed can be split and processed in parallel.
  - Files that are *gzip* compressed cannot be split.
- Let's see how this looks on a sample dataset of 5 large csv files.
- Each file is
  - 12.5 GB uncompressed
  - 1.6 GB gzip
  - 1.3 GB bzip2
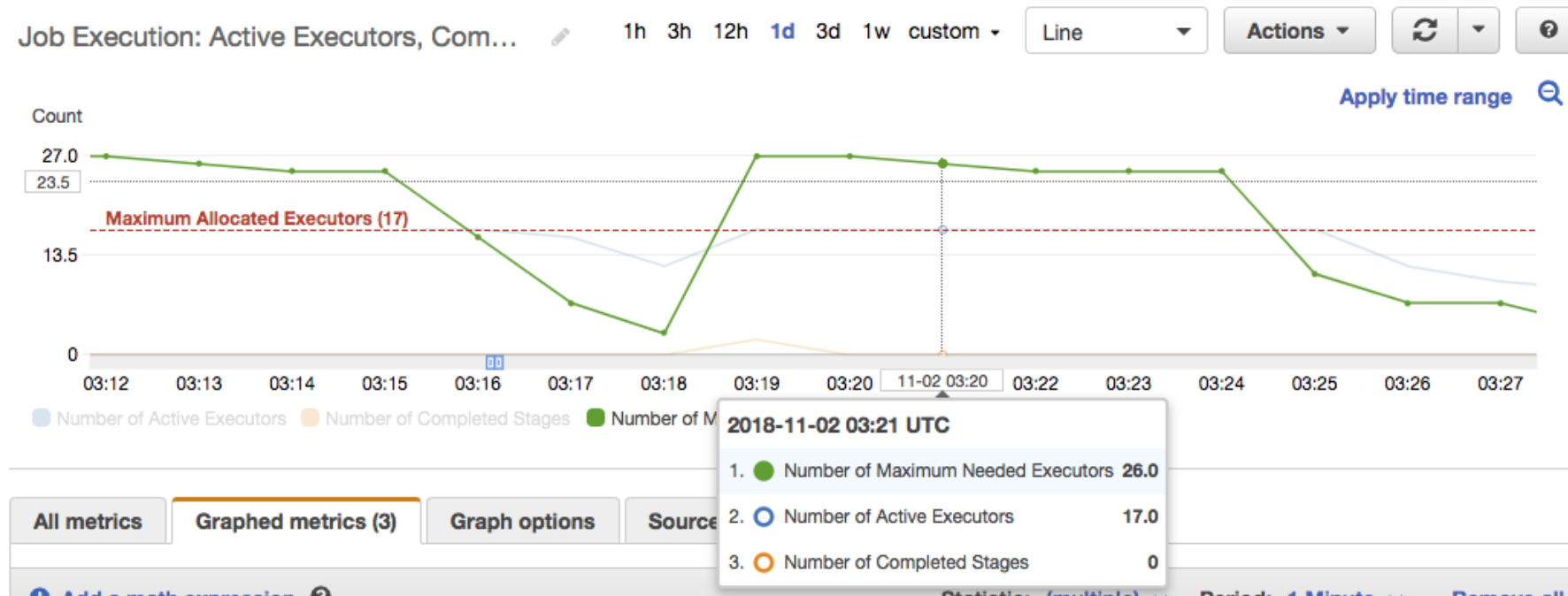- Script converts data to Parquet.

aws

# Example: Processing a few large gzip files

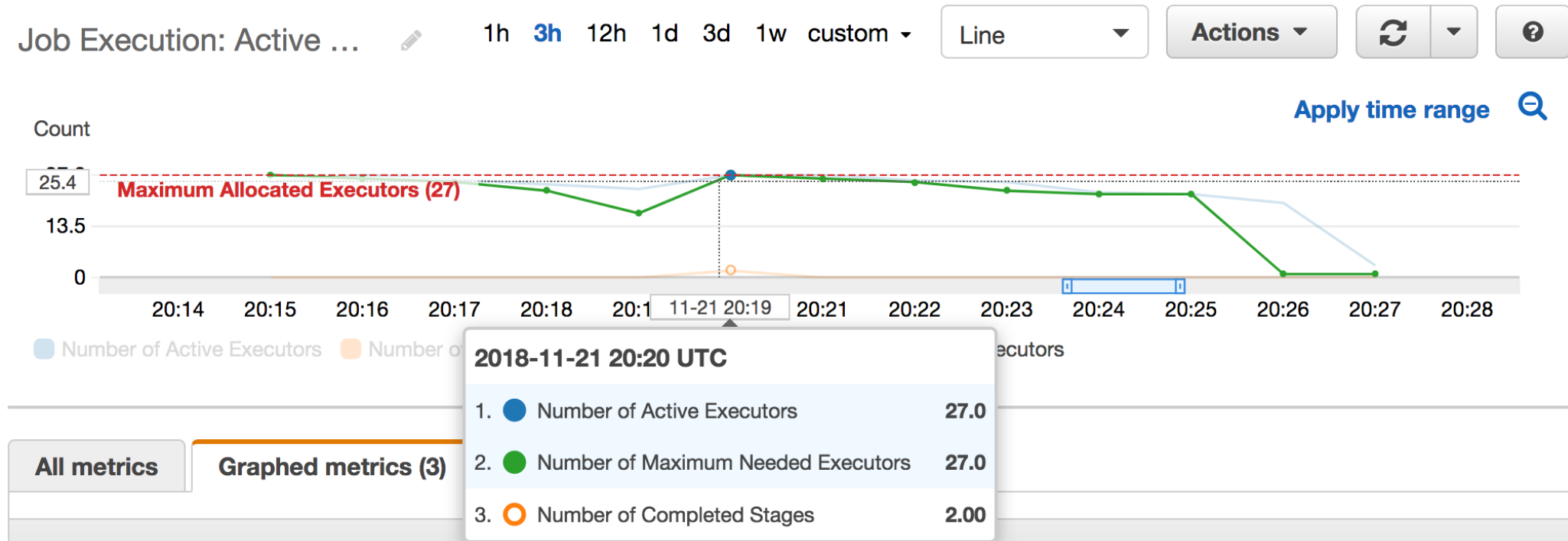- We only have 5 partitions – one for each file.
- Job fails after 2 hours.

# Example: Processing a few large bzip2 files

- Bzip2 files can be split into blocks, so we see up to 104 tasks.
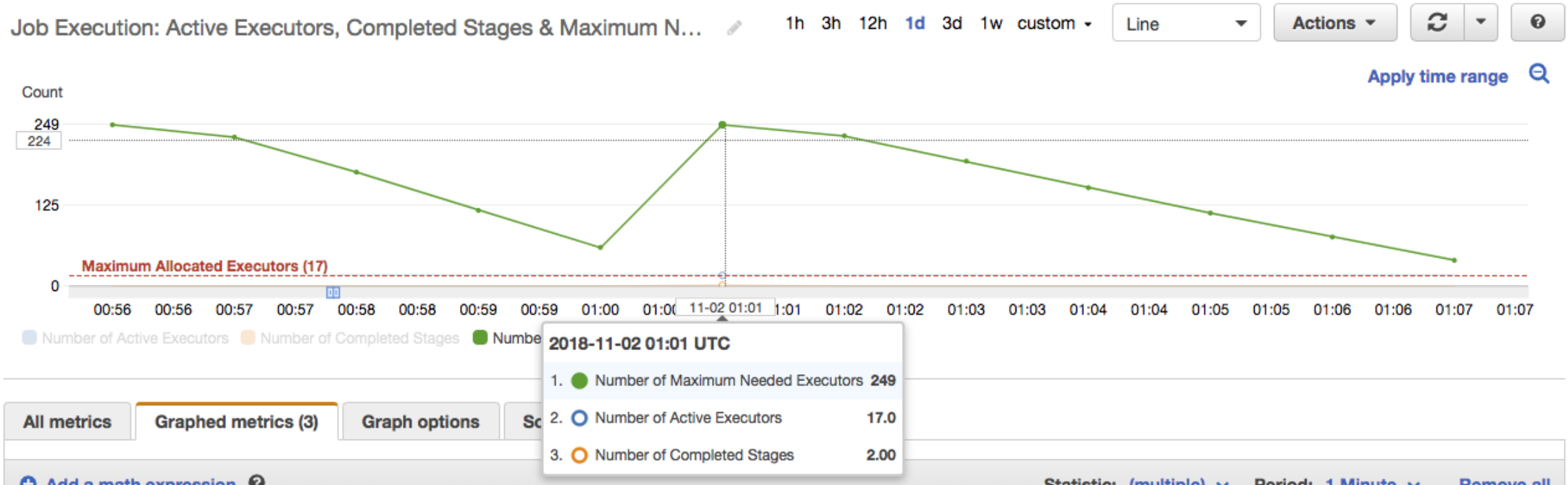- Job completes in 18 minutes.

# Example: Processing a few large bzip2 files

- With 15 DPU, the number of active executors closely tracks the maximum needed number of executors.

# Example: Processing a few large uncompressed files

- Uncompressed files can be split into lines, so we construct 64MB partitions.
- Job completes in 12 minutes.

aws

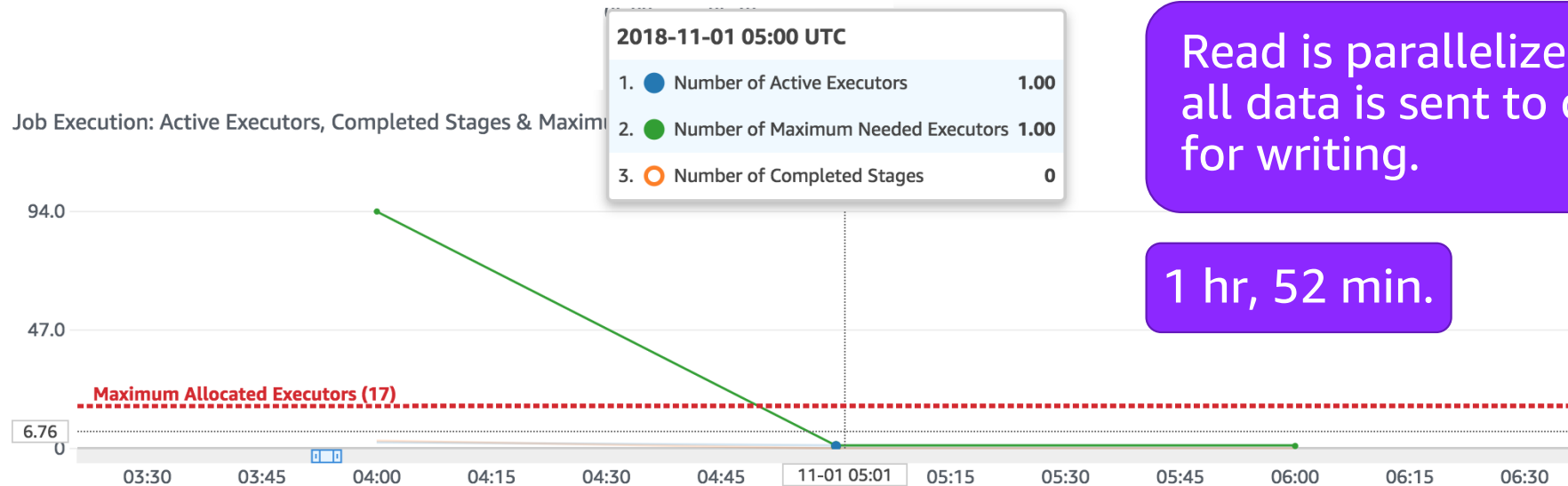# Example: Processing a few large files

- If you have a choice of compression type, prefer bzip2.
- If you are using gzip, make sure you have enough files to fully utilize your resources.
- Bandwidth is rarely the bottleneck for AWS Glue jobs, so consider leaving files uncompressed.

# Example: Coalescing a dataset

- Common use case: reduce the number of files (32 GB compressed)

```
df = spark.read.format("json").load(...)
df2 = df.coalesce(1)
df2.write.format("json").save(...)
```

Job Execution: Active Executors, Completed Stages & Maxim...

**2018-11-01 05:00 UTC**

| | | |
|---|---|---|
| 1. ● Number of Active Executors | **1.00** |
| 2. ● Number of Maximum Needed Executors | **1.00** |
| 3. ○ Number of Completed Stages | **0** |

Read is parallelized and then all data is sent to one executor for writing.

1 hr, 52 min.

94.0

47.0

Maximum Allocated Executors (17)

6.76

0

03:30  03:45  04:00  04:15  04:30  04:45  11-01 05:01  05:15  05:30  05:45  06:00  06:15  06:30

aws re:Invent

aws

# Example: Coalescing with grouping

- You can control this further with grouping
- Group size 500 GB.

Job Execution: Active Executors, Completed Stages & Maximum Needed Executo

**2018-11-06 23:45 UTC**

| 1. | ⭘ Number of Active Executors | **1.00** |
| 2. | ⭘ Number of Maximum Needed Executors | **1.00** |
| 3. | ⭘ Number of Completed Stages | **0** |

17.0 ····································································································
**Maximum Allocated Executors (17)**

**Only a single executor is used for the entire job**

8.50

1.00
0
22:50  22:55  23:00  23:05  23:10  23:15  23:20  23:25  23:30  23:35  2 11-06 23:42  23:50  23:55  00:00  00:05  00:10  00:15

**1 hr, 27 min.**

AWS re:Invent

aws

# Example: Coalescing with grouping

- You can control this further with grouping
- Group size 2 GB.

2018-11-07 00:40 UTC

| | | |
|---|---|---|
| 1. 🔵 Number of Active Executors | 4.00 |
| 2. 🟢 Number of Maximum Needed Executors | 4.00 |
| 3. ⭕ Number of Completed Stages | 0 |

Job Execution: Active Executors, Completed Stages & Maximum Nee

17.0 -------------------------------------------------
**Maximum Allocated Executors (17)**

8.50

4.56

0

00:20    00:25    00:30    00:35    11-07 00:39    00:45    00:50    00:55

Four executors were used to process 16 partitions of ~2 GB each.

26 min.

aws

# Example: Baseline performance without coalescing



**2018-11-06 22:35 UTC**

| | | |
|---|---|---|
| 1. ⬤ Number of Maximum Needed Executors | 105 | |
| 2. ⬤ Number of Active Executors | 17.0 | |
| 3. ⬤ Number of Completed Stages | 0 | |

Job E ... & Maximum Needed Executors

105
79.6
52.5

Maximum Allocated Executors (17)

0

22:34   22:35   22:35   22:36   22:   11-06 22:36   37   22:37   22:38   22:38   22:39   22:39   22:40   22:40   22:41

- Job is fully parallelized and completes in 6 minutes.
- Output data files are around 500 MB.

# DynamicFrames and schema computation

- DynamicFrames are Spark RDDs of *self-describing* records.

- An overall schema is not required for basic ETL operations.
  - I can drop the field "age" without looking at other records.

- Some operations do require a complete schema.
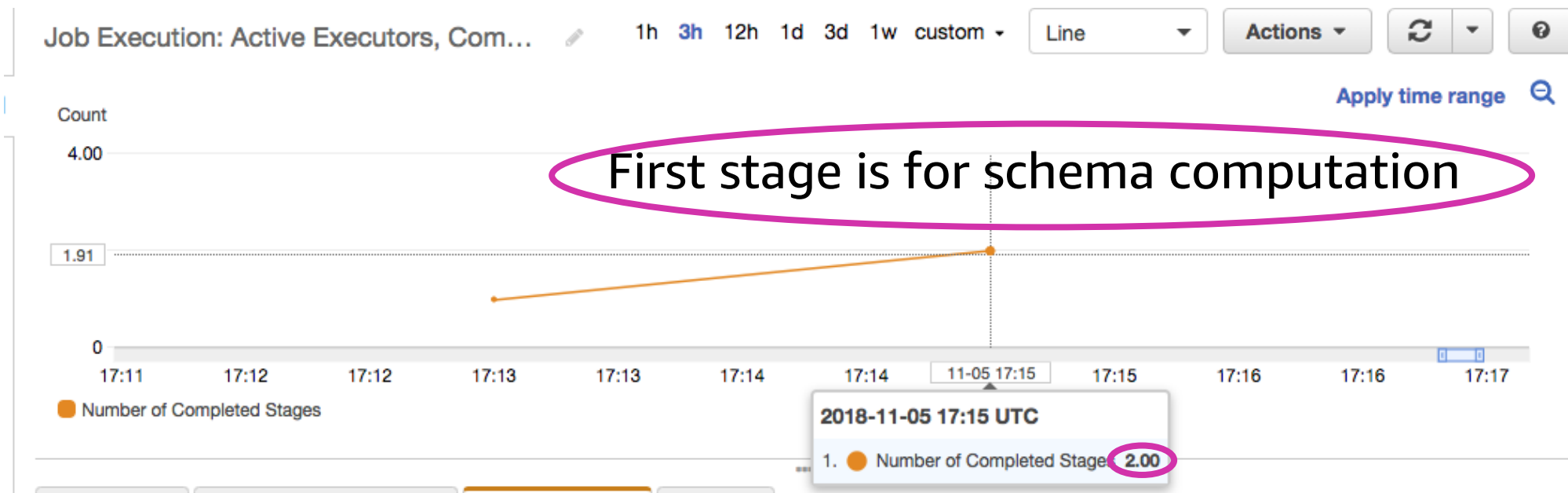  - This can force an extra job in your application.



Transforms that require a schema:
- DropNullFields
- Relationalize
- ResolveChoice without specifying columns
...

# DynamicFrame Schema Example: DropNullFields

```
df = glueContext.create_dynamic_frame_from_catalog(…)
df2 = DropNullFields.apply(df)
glueContext.write_dynamic_frame_from_options(df2, …)
```
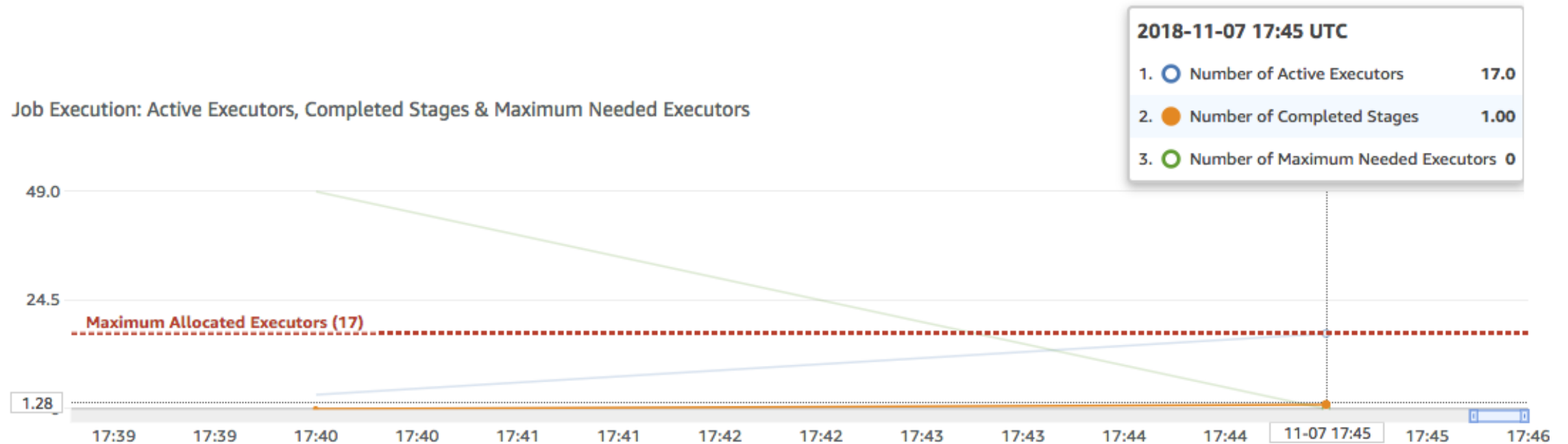


First stage is for schema computation

# DynamicFrame Schema Example: DropNullFields

```
df = glueContext.create_dynamic_frame_from_catalog(…)
df2 = df.applyMapping([
    ('col1', 'col1', 'string'),
    ('col2', 'col2', 'string'),
    ('col3', 'col3', 'string'),
    ('col4', 'col4', 'string'),
    ('col5', 'col5', 'boolean')
])
df3 = DropNullFields.apply(df2)
glueContext.write_dynamic_frame_from_options(df3, …)
```

ApplyMapping sets the schema without an additional pass.

AWS re:Invent
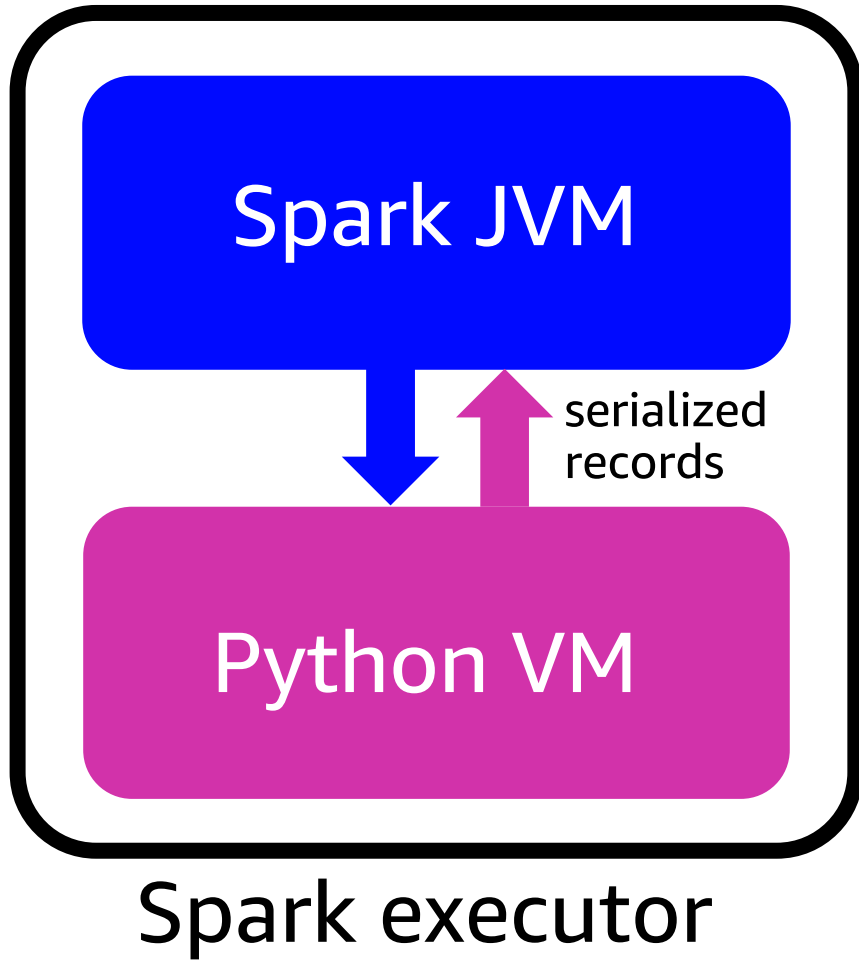
aws

# DynamicFrame Schema Example: DropNullFields

- There is only one stage in the application.

# Optimizing DynamicFrames

- Use ApplyMapping where appropriate to set a schema.

- Be thoughtful about where to add transformations like DropNullFields.
  - The default generated scripts are designed to be safe, but you may be able to optimize if you know your data.

- Try your old scripts!
  - We've increased data conversion speed by 4x in the past year.

# Python performance in AWS Glue



**Spark executor**

- Using map and filter in Python is expensive for large data sets.
  - All data is serialized and sent between the JVM and Python.
- Alternatives
  - Use AWS Glue Scala SDK.
  - Convert to DataFrame and use Spark SQL expressions.
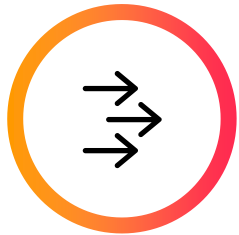
# AWS Lake Formation

Build a secure data lake in days

Register existing data or load new data using blueprints. Data stored in Amazon S3.

Secure data access across multiple services using single set of permissions.

No additional charge. Only pay for the underlying services used.

**Quickly build data lakes**

Move, store, catalog, and clean your data faster. Use ML transforms to de-duplicate data and find matching records.
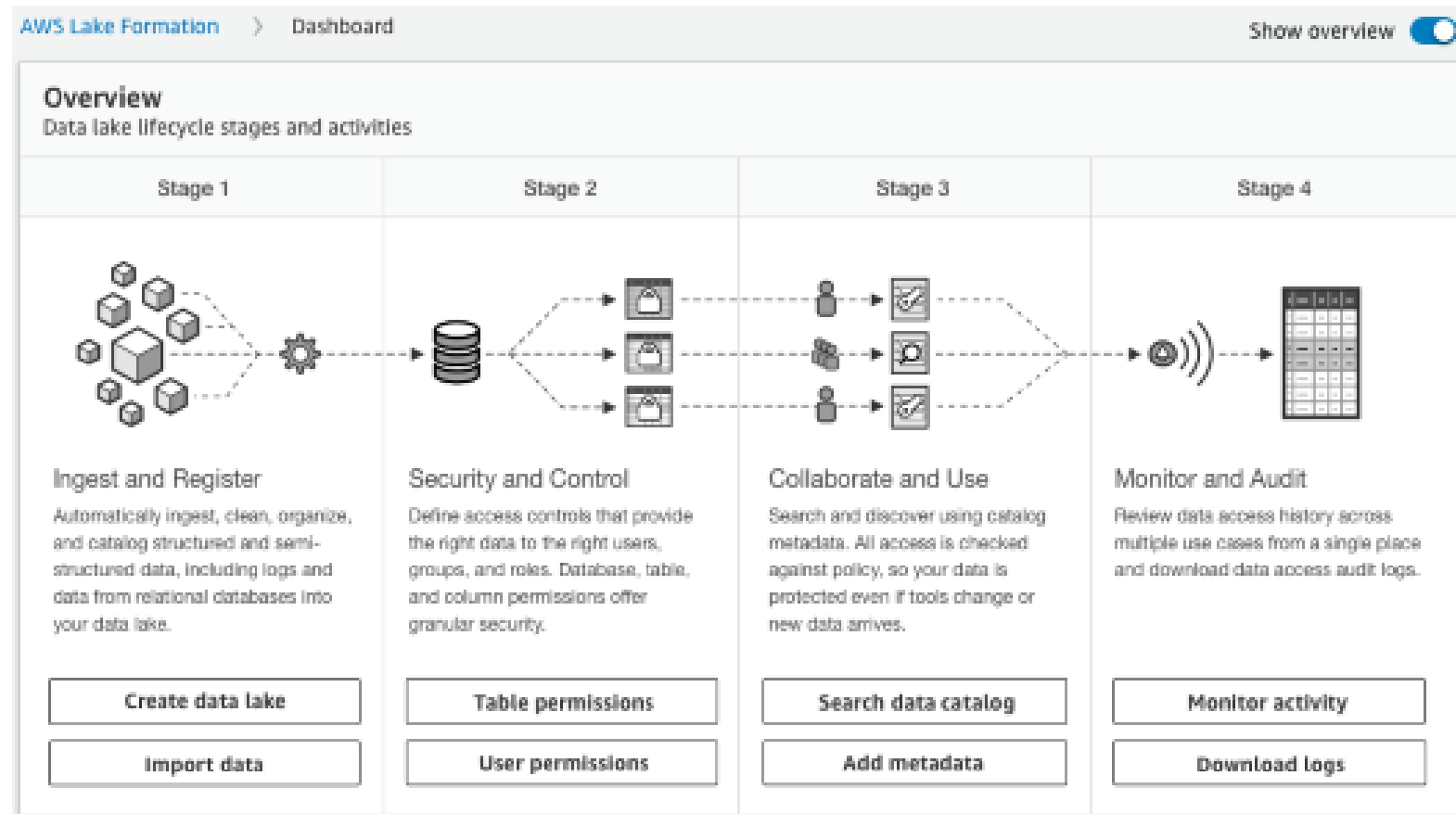
**Easily secure access**

Centrally define table and column-level data access and enforce it across Amazon EMR, Amazon Athena, Amazon Redshift Spectrum, Amazon SageMaker, and Amazon QuickSight

**Share and collaborate**

Use data catalog in Lake Formation to search and find relevant data sets and share them across multiple users and accounts

AWS re:Invent

# How it works

# Thank you!

Benjamin Sowell

aws

Please complete the session survey in the mobile app.