# Spark Yoga - saving time & money with lean data pipelines

Vishnu Rao Senior data engineer at Eyeota (ex - Flipkart)

#### Running time is money



"some cloud provider"

# Yoga Asanas





Yep thats me :)

# What will we focus on today ?

4

# What will we focus on today ?

#### 1. Some Asanas for **Spark**



# What will we focus on today ?

- 1. Some Asanas for **Spark**
- 2. And some simple actual Yoga Asanas too at the end :)



• Using Spark in a specific way to save time & money

- Using Spark in a specific way to save time & money
  - Tuning / Using a specific Spark Configuration parameter

- Using Spark in a specific way to save time & money
  - Tuning / Using a specific Spark Configuration parameter
  - Choosing the kind of data lake it interacts with

- Using Spark in a specific way to save time & money
  - Tuning / Using a specific Spark Configuration parameter
  - Choosing the kind of data lake it interacts with
  - Running the spark job in a specific way

- Using Spark in a specific way to save time & money
  - Tuning / Using a specific Spark Configuration parameter
  - Choosing the kind of data lake it interacts with
  - Running the spark job in a specific way
  - Choosing the kind of framework it runs in

- Using Spark in a specific way to save time & money
  - Tuning / Using a specific Spark Configuration parameter
  - Choosing the kind of data lake it interacts with
  - Running the spark job in a specific way
  - Choosing the kind of framework it runs in
  - Balancing Time With Co\$t

Some of my frequently used configs are

1. Spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version = 2 (saves time)

- 1. Spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version = 2 (saves time)
  - a. Or experiment with s3a committers

- 1. Spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version = 2 (saves time)
  - a. Or experiment with s3a committers
- 2. spark.serializer=org.apache.spark.serializer.KryoSerializer (better serialization)

- 1. Spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version = 2 (saves time)
  - a. Or experiment with s3a committers
- 2. spark.serializer=org.apache.spark.serializer.KryoSerializer (better serialization)
- 3. spark.sql.autoBroadcastJoinThreshold=256MB (helps with the join)



- 1. Spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version = 2 (saves time)
  - a. Or experiment with s3a committers
- 2. spark.serializer=org.apache.spark.serializer.KryoSerializer (better serialization)
- 3. spark.sql.autoBroadcastJoinThreshold=256MB (helps with the join)
- 4. spark.hadoop.fs.s3a.connection.ssl.enabled=false (saves time)

- 1. Spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version = 2 (saves time)
  - a. Or experiment with s3a committers
- 2. spark.serializer=org.apache.spark.serializer.KryoSerializer (better serialization)
- 3. spark.sql.autoBroadcastJoinThreshold=256MB (helps with the join)
- 4. spark.hadoop.fs.s3a.connection.ssl.enabled=false (saves time)
- 5. If executor < 32G, Use -XX:+UseCompressedOops (pointer size reduction)

1. S3\* or any other cloud object storage

- 1. S3\* or any other cloud object storage
- 2. HDFS (On Premise / Cloud)

- 1. S3\* or any other cloud object storage
- 2. HDFS (On Premise / Cloud)

The above are the more notable choices. Are they the only choices?

- 1. S3\* or any other cloud object storage
- 2. HDFS (On Premise / Cloud)

Evaluation Criteria:

- 1. S3\* or any other cloud object storage
- 2. HDFS (On Premise / Cloud)

Evaluation Criteria:

- 1. Data access Time
- 2. Maintenance

	S3/Cloud Storage	HDFS (onprem / cloud)
Infra Maintenance	no 🧹	yes 🌓
Access times	Over network - https? 🕜	Local access to node 🕔
Compute	No, Need extra nodes	Yes with hadoop 🗸

	S3/Cloud Storage	HDFS (onprem / cloud)	Local File System + S3
Infra Maintenance	no 🧹	yes 🎁	Yes. local storage might be limited, extend via EBS ?
Access times	Over network - https?	Local access 🕔 to node	Local access of to node
Compute	No, Need extra nodes	Yes with hadoop	Yes, you launch workers on the box

	S3/Cloud Storage	HDFS (onprem / cloud)	Local File System + S3
Infra Maintenance	no 🧹	yes 🎁	Yes. local storage might be limited, extend via EBS ?
Access times	Over network - https?	Local access 🕔 to node	Local access of to node
Compute	No, Need extra nodes	Yes with hadoop	Yes, you launch workers on the box

Pro tip: 'aws cp' to local file system is faster than spark interacting with s3

### Example of Orchestration

	Airflow	DAGs	Data Profiling 🗸	Browse 🗸	Admin 🗸	Docs 🗸	About 🗸			
DAG: eyeota_dag_seg_counting_on_part0000										
	# Graph View	🕈 Tre	e View <b>"II</b> Task	Duration 🌓 T	Task Tries	★ Landing T	īmes 📃 Gantt	i Details	4 Code	• Trigger DAG
	success Base da	te: 2021	-10-10 00:00:07	Number of runs	s: 25 ~	Run: man	ual_2021-10-10T00	:00:06+00:00	<ul> <li>✓ Layout:</li> </ul>	Left->Right v
BashOperator										
	print_date_p	uses and (aws	ownload_input_s3_t Ansible to do cli) in paralle	o_local → run wnload I to all hos	<u>_spark_job_c</u>	ounts_genera	ation upload_ou	itput_local_to_s	3 <b>→</b> Clean	Lub

Note: while using local file system(file://), you need to download the entire dataset to each host

	S3/Cloud Storage	HDFS	Local File System + S3
Infra Maintenance	no 🧹	yes 🎁	Yes. local storage might be <b>limited</b> , extend via EBS ?
Access times	Over network -	Local access to node	Local access of to node
Compute	No, Need extra nodes	Yes with hadoop	Yes, you launch vorkers on the box

Sometimes usage of Local File System along with S3,

o- Can lead to less running times as its a local data access pattern

o- Yes, we launch & maintain our own cluster, which might not be a burden as long as cost is low.

We always tend to go the spark-client mode / run in a cluster .

We always tend to go the spark-client mode / run in a cluster.

How about the Spark - Local mode ??

Spark.master = local[\*]

We always tend to go the spark-client mode / run in a cluster.

How about the Spark - Local mode ??

Spark.master = local[\*]

• Valid way of running jobs. One might think it's for only running tests.

We always tend to go the spark-client mode / run in a cluster.

How about the Spark - Local mode ??

Spark.master = local[\*]

- Valid way of running jobs. One might think it's for only running tests.
- I say why not ? You get benefit of spark DSL while solving your problem.
We always tend to go the spark-client mode / run in a cluster.

How about the Spark - Local mode ??

Spark.master = local[\*]

- Valid way of running jobs. One might think it's for only running tests.
- I say why not ? You get benefit of spark DSL while solving your problem.
- You can always scale horizontally by going to a cluster
  - Cluster is not always needed !

We always tend to go the spark-client mode / run in a cluster.

How about the Spark - Local mode ??

Spark.master = local[\*]

- Valid way of running jobs. One might think it's for only running tests.
- I say why not ? You get benefit of spark DSL while solving your problem.
- You can always scale horizontally by going to a cluster
  - Cluster is not always needed !
- Also you have the option of running in a **kubernetes** pod !

Data preparation

Data preparation

Repartition the data for better parallelization of work.

Ex: we wish to increase our write throughput to the db



Data preparation - Sharding Data



Data preparation - Sharding Data



Data preparation - Sharding Data

• Divide & Conquer principle



Data preparation - Sharding Data

- Divide & Conquer principle
  - Say 100 shards/buckets/folders



#### Data preparation - Sharding Data

- Divide & Conquer principle
  - Say 100 shards/buckets/folders
  - 100 spark jobs on 100 smaller datasets



#### Data preparation - Sharding Data

- Divide & Conquer principle
  - Say 100 shards/buckets/folders
  - 100 spark jobs on 100 smaller datasets
- Resource allocation now based on shard size , not large dataset => smaller boxes/cluster



#### Data preparation - Sharding Data

- Divide & Conquer principle
  - Say 100 shards/buckets/folders
  - 100 spark jobs on 100 smaller datasets
- Resource allocation now based on shard size , not large dataset => smaller boxes/cluster
- Isolated failures & ReRun rerun jobs only for failed buckets



#### Data preparation - Sharding Data

- Divide & Conquer principle
  - Say 100 shards/buckets/folders
  - 100 spark jobs on 100 smaller datasets
- Resource allocation now based on shard size , not large dataset => smaller boxes/cluster
- Isolated failures & ReRun rerun jobs only for failed buckets
- Debugging failures / Improving performance focussed wrt small bucket



Data preparation - Sharding data

Does it help always?

- No, Say you shard using user-id/order-id & a spark job needs to cut across all the shards (group by country rather than group by user id)
- We might have to aggregate/merge the results obtained in each job to get the final result. This might be complex.

## The Setup

S3 Bucket for Raw data with 100 parts/shards/folders



### The Setup

- Dynamically create 'n' Airflow dags.
- Run 'n' Airflow Dags for 'n' shards/folders

#### Sample source code taken from

https://www.astronomer.io/quides/dynamically-generating-dags

```
Airflow
                 DAGs
                         Security-
                                              Admin-
                                                                Astronomer-
                                    Browse-
                                                        Docs-
DAGs
           Active 1 Paused 3
  All (4)
                                                                          Fi
 0
    DAG
                                                                   Owner
 loop_hello_world_1
                                                                   airflow
 loop_hello_world_2
                                                                   airflow
 loop_hello_world_3
                                                                   airflow
```

from airflow import DAG				
<pre>from airflow.operators.python_operator import PythonOperator</pre>				
from datetime import datetime				
def exerts des (des id				
der create_dag(dag_10,				
dag number				
default area):				
uerau((_arys)).				
def hello world pv(*aros):				
<pre>print('Hello World')</pre>				
<pre>print('This is DAG: {}'.format(str(dag number)))</pre>				
<pre>dag = DAG(dag_id,</pre>				
<pre>schedule_interval=schedule,</pre>				
<pre>default_args=default_args)</pre>				
with dag:				
t1 = PythonOperator(				
task_id='hello_world',				
python_callable=hello_world_py)				
return dag				
# huild a dag for each number in range(10)				
# build a dag for each number in range(10)				
dog id = lloop hallo world {}! format(str(n))				
default args = {'owner': 'airflow'.				
'start date': datetime(2021, 1, 1)				
}				
schedule = '@daily'				
dag_number = n				
<pre>globals()[dag_id] = create_dag(dag_id,</pre>				
schedule,				
dag_number,				
default_args)				

Reusable datasets

• Are there multiple actions on same dataset in the lineage?



Reusable datasets

• Are there multiple actions on same dataset in the lineage?

Dataset1 = spark.load(...)

Dataset2 = complex\_transform\_(Dataset1); // takes lot of time.

Dataset2.count() //action 1

Dataset2.save() // action 2

Reusable datasets

• Are there multiple actions on same dataset in the lineage?

Dataset1 = spark.load(...)

Dataset2 = complex\_transform\_(Dataset1); // takes lot of time.

Dataset2.persist(DISK\_ONLY); // try with memory if you have it

Dataset2.count() //action 1

Dataset2.save() // action 2

• EMR/Yarn - choice for most people

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?
  - Yes you maintain the cluster
  - But sometimes you can be smart with it

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?
  - Yes you maintain the cluster
  - But sometimes you can be **smart** with it
    - I/O intensive jobs ex: writing to cassandra

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?
  - Yes you maintain the cluster
  - But sometimes you can be **smart** with it
    - I/O intensive jobs ex: writing to cassandra
      - => Less CPU used
      - => Not all cores used to full potential

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?
  - Yes you maintain the cluster
  - But sometimes you can be smart with it
    - I/O intensive jobs ex: writing to cassandra
      - => Less CPU used
      - => Not all cores used to full potential
      - Can we increase Threads -> **YES** !

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?
  - Yes you maintain the cluster
  - But sometimes you can be smart with it
    - I/O intensive jobs ex: writing to cassandra
      - => Less CPU used
      - => Not all cores used to full potential
      - Can we increase Threads -> YES !
      - Box has N cores -> you configure SPARK\_WORKER\_CORES to M cores (M > N)

- EMR/Yarn choice for most people.
- How about Standalone Cluster ?
  - Yes you maintain the cluster
  - But sometimes you can be smart with it
    - I/O intensive jobs ex: writing to cassandra
      - => Less CPU used
      - => Not all cores used to full potential
      - Can we increase Threads -> YES !
      - Box has N cores -> you configure SPARK\_WORKER\_CORES to M cores (M > N)
      - This is called Core-Oversubscribing
        - EMR might not give this option\*. Only increase number of boxes => more\$

### Oversubscribing



Spark Master at spark:// 7077 URL: spark://1-------7077 Alive Workers: 2 Cores in use: 384 Total, 0 Used Memory in use: 48.0 GiB Total, 0.0 B Used **Resources in use:** Applications: 0 Running, 193 Completed Two C5.4X large boxes Drivers: 0 Running, 0 Completed (16core, 32G) Status: ALIVE - Workers (2)

Worker Id	Address	State	Cores
worker-20210506014111-1000000000000000000000000000000	69015	ALIVE	192 (0 Used)
worker-20210506014121-1-46247	1.46247	ALIVE	192 (0 Used)

Spark on Kubernetes

Spark on Kubernetes

+

Using Spot instances

Spark on Kubernetes

+

Using Spot instances

=

Cheaper way to run workloads

• At the end of day \$ does matter with an SLA (time) in mind.

• At the end of day \$ does matter with an **acceptable** SLA (time) in mind.

- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost

- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost
  - Compute cluster cost i.e. how long boxes are used ?

- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost
  - Compute cluster cost i.e. how long boxes are used ?
  - Storage Cost during computation
- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost
  - Compute cluster cost i.e. how long boxes are used ?
  - Storage Cost during computation
- Sometimes we can decrease \$ with acceptable increase in SLA (time)
  - $\circ$  Ex: Say running time

- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost
  - Compute cluster cost i.e. how long boxes are used ?
  - Storage Cost during computation
- Sometimes we can decrease \$ with acceptable increase in SLA (time)
  - $\circ$  Ex: Say running time
    - on Two c5.4x.large is 8 hours => X\$ for 8 hour run time

- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost
  - Compute cluster cost i.e. how long boxes are used ?
  - Storage Cost during computation
- Sometimes we can decrease \$ with acceptable increase in SLA (time)
  - Ex: Say running time
    - on Two c5.4x.large is 8 hours => X\$ for 8 hour run time
    - Theoretically on Two t3.2x.large is ~16 hours => Y\$ for 16 hour run time (Y<X)</p>

- At the end of day \$ does matter but with acceptable SLA (time) in mind.
- Cost
  - Compute cluster cost i.e. how long boxes are used ?
  - Storage Cost during computation
- Sometimes we can decrease \$ with acceptable increase in SLA (time)
  - Ex: Say running time
    - on Two c5.4x.large is 8 hours => X\$ for 8 hour run time
    - Theoretically on Two t3.2x.large is ~16 hours => Y\$ for 16 hour run time (Y<X)</li>
    - If you are ok with the SLA of 16 hours, why not go for it?



Dataset records = spark.session.orc.load("s3a://...").select(few columns); // 15gb in s3 folder

Dataset lowercase = records.map(.....transform to lowercase....);

lowercase.count(); //some action

```
Dataset Uppercase = records.map(.....transform to uppercase...);
```

Uppercase.count(); //some action

Running Time: 12.5 min

(spark.hadoop.fs.s3a.connection.ssl.enabled=false)

Dataset records = spark.session.orc.load("s3a://...").select(few columns); // 15gb in s3 folder

Dataset lowercase = records.map(.....transform to lowercase....);

lowercase.count(); //some action

Dataset Uppercase = records.map(.....transform to uppercase...);

Uppercase.count(); //some action

Running Time: 11.4min 🕔

(spark.hadoop.fs.s3a.connection.ssl.enabled=false)

Dataset records = spark.session.orc.load("s3a://...").select(few columns); // 15gb in s3 folder

records.persist(DISK\_ONLY); // if you have lot of free ram, try with memory

Dataset lowercase = records.map(.....transform to lowercase....);

lowercase.count(); //some action

```
Dataset Uppercase = records.map(.....transform to uppercase...);
```

Uppercase.count(); //some action



(aws s3 cp data to local disk & read from there- took 2min)

Dataset records = spark.session.orc.load("file://...").select(few columns); // 15gb in s3 folder

Dataset lowercase = records.map(.....transform to lowercase....);

*lowercase.count(); //some action* 

Dataset Uppercase = records.map(.....transform to uppercase...);

Uppercase.count(); //some action

Running Time: **5.9min** (*including 2min for downloading data*)

#### (aws s3 cp data to local disk & read from there)

Dataset records = spark.session.orc.load("file://...").select(few columns); // 15gb in s3 folder

#### records.persist(DISK\_ONLY); // if you have lot of free ram, try with memory

```
Dataset lowercase = records.map(.....transform to lowercase....);
```

```
lowercase.count(); //some action
```

```
Dataset Uppercase = records.map(.....transform to uppercase...);
```

Uppercase.count(); //some action

Total Running Time: **6.8min** (includes downloading data to disk. The best was 5.9min)

#### Exercise (Results)

Scenario	Minutes
Program as it is (from s3)	12.5
Program as it is (from s3) + ssl disable	11.4
Persisting disk (from s3)	8.6
Persist disk (from s3) + ssl disable	8.5
Aws cp to local	<mark>5.9</mark>
Aws cp to local + Persist (disk)	6.8
Aws cp to local + Persist (mem_ser)	
xms=6g	6.9
Aws cp to local + Persist (mem)	
xms=6g	6.7



\*Aws s3 cp took 2 minutes

### Note: Please do run your own tests !

Explore your own code.

#### Finally, It's about the Balance



# So try out these Asanas for a cost-effective Spark Life



## Thank you

linkedin.com/in/213vishnu

twitter.com/sweetweet213

https://mash213.wordpress.com/