

GIS Pipeline acceleration with Apache Sedona



Alihan Zihna Senior Data Scientist, CKDelta (Hutchison Group)



Fernando Ayuso Director of Data Science and Data Engineering, CKDelta (Hutchison Group)

ORGANIZED BY 🗟 databricks

Agenda

What are we sharing today?

- Business context
- Use case
- Technical implementation
- Demo

Business Context







Telco data + GIS

With over **26 billion records** ingested and analysed per day in the platform CKDelta is striving to add **granularity** to telco analysis





+2M Cells in our platform – High Complexity SQL or GIS tools not sufficient

+ granularity & Speed

Combination of python libraries with Spark UDFs on Databricks



Utilities CKDelta a GIS rich industry

CKDelta works with Northumbria Water, UK Power Networks, ISTA or Northern Gas Networks solving problems through data.

Millions of GIS components

Overlaying networks: Water, Gas, Electricity...

Limitation of what GIS can do using only specialised libraries



Use cases: Telco + Utilities

Limits in technology limit the imagination of our solutions



- Water Leakage
- Customer satisfaction
- Electrical Vehicle Charging Station installation and consumption location
- Supply VS Demand

CKDelta





Technical Implementation





Sedona on Databricks

How does Sedona work?

Extending Spark

- Sedona extends the Spark context by new spatial level functions
- Introduces the basic building blocks of GIS to spark such as Polygons, Lines, Spatial Statistics, Spatial Indices



Sedona on Databricks

How does Sedona work?

Extending Spark

- Sedona extends the Spark context by new spatial level functions
- Introduces the basic building blocks of GIS to spark such as Polygons, Lines, Spatial Statistics, Spatial Indices



Sedona on Databricks

How does Sedona work?

Extending Spark

- Sedona extends the Spark context by new spatial level functions
- Introduces the basic building blocks of GIS to spark such as Polygons, Lines, Spatial Statistics, Spatial Indices



Sedona on Databricks

Spatial Functions

Constructors

- ST_GeomFromWKT
- ST_GeomFromGeoJson
- ST_Point

Functions

- ST_Distance
- ST_Length
- ST_Area
- ST_Intersection
- ST_IsValid

Aggregate

- ST_Envelope_Aggr
- ST_Union_Aggr
- ST_Intersection_Aggr

Sedona on Databricks

Spatial Functions

Predicates

- ST_Contains
- ST_Intersects
- ST_Within
- ST_Overlaps

Join Query Optimizers

Other functions can be used to optimize GeoSpatial functions such as joining where shapes:

- intersect
- the distance between is low

Sedona Sample

Read File and Select only valid Shapes

- 1 # Read to RDD and Convert to Spark DataFrame
- 2 rdd = GeoJsonReader.readToGeometryRDD(sc, path)
- 3 Adapter.toDf(rdd, spark).createOrReplaceTempView('temp')

```
4 # Filter only Valid
```

5 spark.sql(

6 """

```
7 SELECT * FROM temp
```

```
8 WHERE ST_IsValid(geometry)
```

9 """

10)

DATA+AI SUMMIT 2022

What does Sedona Introduce

Points, Lines, and Polygons

udt]

2 3 4	data = [
3 4 5	data = [
4	
E	[1, Point(21.0, 52.0)],
Э	[1, Point(23.0, 42.0)],
6	[1, Point(26.0, 32.0)]
7]
8	
9	
10	<pre>gdt = spark.createDataFrame(data</pre>
11	data,
12	[mdex , geom]
14	1
15	gdf.show()
+ ir + 	<pre>(c) gut: Your Your Your Your Your Your Your Your</pre>

1	<pre>from shapely.geometry import LineString</pre>
2	line = $[(40, 40), (30, 30), (40, 20), (30, 10)]$
4	cine = [(40, 40), (30, 30), (40, 20), (30, 20)]
5	data = [
6	<pre>[1, LineString(line)]</pre>
7]
8	
9	gdf = spark.createDataFrame(
10	data,
11	['index', 'geom']
12)
13	
14	gdf.show(1, False)
	(1) Sports Jahr
	(2) Spark JODS
Þ	gdf: pyspark.sql.dataframe.DataFrame = [index: long, geom: ud
+	+
iı	ndex geom
+	+
1	LINESTRING (40 40, 30 30, 40 20, 30 10)
+	++
	\land

1	from shapely.geometry import Polygon
3	polygon = Polygon(
4	[
5	[19,51121, 51,76426],
6	[19.51056, 51.76583],
7	[19.51216, 51.76599],
8	[19.51280, 51.76448],
9	[19.51121, 51.76426]
10]
11)
12	
13	data = [
14	[1, polygon]
15	1
16	
17	gdf = spark.createDataFrame(
18	data,
19	L'index.' . Geom 1
71	,
22	gdf.show(1, False)
	(T) Smark Jahr
1	(e) shair sons
	gdf: pyspark.sql.dataframe.DataFrame = [index: long, geom: udt]
li	ndex geom
+	POLYGON ((19.51121 51.76426, 19.51056 51.76583,
12	

DATA+AI SUMMIT 2022

What does Sedona Introduce

Spatial Indexing and Partitioning

Spatial indexes:

Allow faster access to records

Spatial Partitioning:

 Allow close objects to be in the same partitions







Serialization

What do we need serialization?

Serialization Basics

Required to reduce the footprint when:

- Writing to a file or a database
- Transferring through network (e.g. to other workers)

Spark uses Java Serialization by default.

Kryo vs. Java Serialization

Kryo serializer is more efficient and safer than Java serializer by 20%–50%.

New data objects must be defined individually. If undefined:

- Security implications
- Reduced performance

What does Sedona Introduce

Efficient Serialization

Reducing Memory Impact

Represents the shape objects as bytes using:

- wkb serializer
- shape serializer

Can be configured using spark.configs



Sedona on Databricks

Which version to install?

Depends on Databricks Runtime Version and Spark Version

Cluster Runtime Major	Spark	Sedona version
7 - 9	3.0 - 3.1	1.0.1 - 1.1.0
10 - 11	3.2	1.1.1 - 1.2.0



Sedona on Databricks

How to set-up

Installing	Using UI	
Cuesta Libuaru		

library Source	Library			
Upload	DBFS/ADLS	РуРІ	Maven	CRAN
Library Type	thon Egg F	vthon Wh	d	
		And I red		
Library Name	ę			
Optional				
			Drog	p JAR here
Create	Cancel			

Simple Setup

- Add packages as libraries
- Find the external packages
- Install to the clusters

Drawbacks

- No version control
- Can't benefit from serialization

Sedona on Databricks

How to set-up

Installing Using init scripts

Create init script directory for Sedona mkdir -p /dbfs/FileStore/sedona/

Create init script
cat > /dbfs/FileStore/sedona/sedona-init.sh <<*EOF*
#1/bin/bash</pre>

File: sedona-init.sh
Author: Erni Durdevic
Created: 2021-11-01

On cluster startup, this script will copy the Sedona jars to the cluster's default jar directory. # In order to activate Sedona functions, remember to add to your spark configuration the Sedona extensions: # "spark.sl.v.tortsnions org.spache.sedona.viz.sd, Sedonavižsternstnics.ged, sadenasd[Extensions"

cp /dbfs/FileStore/jars/sedona/1.2.0-incubating/*.jar /databricks/jars

EOF

Spark	Tags	SSH	Logging	Init Scripts	JDBC/ODBC	Permissions	
Init scri	pts O						
	Туре	Fi	le path				
	DBFS	dt	ofs:/FileStore/	sedona/sedona-	initah		

Spark	Tags	SSH	Logging	Init Scripts	JDBC/ODBC	Permissions
Spark c	onfig 🔞					
spark.s org.apa	ql.extensi ache.sedo	ons na.viz.sq	l.SedonaVizE>	xtensions.org.ap	ache.sedona	
.sql.Sec spark.r	donaSqlEx pc.messag	densions ge.maxSi	ze 1024			
.sql.Sec spark.r spark.d spark.k	donaSqlEx pc.messag latabricks ryo.regist	tensions ge.maxSi .delta.pre rator	ze 1024 eview.enabled	l true		

Copy the sample scripts from Sedona setup page

Add the bash file to the cluster's init scripts.

Update the spark configs

DATA+AI SUMMIT 2022

Sedona on Databricks

Configurations

Spark Configs for Spark

.

spark.sql.extensions org.apache.sedona.viz.sql.SedonaVizExtensions,org.apache.sedona.sql.SedonaSqlExtensions

- 2 spark.serializer org.apache.spark.serializer.KryoSerializer
- 3 spark.kryo.registrator org.apache.sedona.core.serde.SedonaKryoRegistrator
- 4 spark.rpc.message.maxSize 1024

Sedona on Databricks

Using in notebook

SedonaRegistrator.registerAll() enables all new functions.

Sedona can now be used with pyspark or SQL functions.

Cmd	1
1	%pip install apache-sedona==1.2.0
Py Co Co Re Ir Su Py	thon interpreter will be restarted. llecting apache-sedona==1.2.0 Downloading apache_sedona=1.2.0-py3-none-any.whl (67 kB) llecting shapely Downloading Shapely-1.8.2-cp38-cp38-manylinux_2_12_x86_64.manylinux2010_x86_64.whl (2.1 MB) quirement already satisfied: attrs in /databricks/python3/lib/python3.8/site-packages (from stalling collected packages: shapely, apache-sedona ccessfully installed apache-sedona-1.2.0 shapely-1.8.2 thon interpreter will be restarted.
Со	nmand took 12.15 seconds by a.zihna@ckhgbdp.onmicrosoft.com at 04/05/2022, 17:11:34 on DS_Sedona_1_2_0
Cmd	2
1 2	<pre>from sedona.register.geo_registrator import SedonaRegistrator SedonaRegistrator.registerAll(spark)</pre>
•	(2) Spark Jobs
0	t[1]: True
Со	nmand took 4.06 seconds by a.zihna@ckhgbdp.onmicrosoft.com at 04/05/2022, 17:11:36 on DS_Sedona_1_2_0





Demo



Conclusion

- Apache Sedona extends Spark
- Minimal changes required for migration
- Easy setup on Databricks from Sedona's docs
- Significant performance upgrade

DATA+AI SUMMIT 2022

Thank you



Alihan Zihna Senior Data Scientist

Fernando Ayuso Director of Data Science and Data Engineering

Appendix





We are CKHutchison's Group data innovation company



Performance Comparison

Cmd 1	
Data • Two files with 3G and 4G cell shapes • 75,000 cells in both files	Markdown 🖬 🗸 — 🗙
Cluster Standard_DS3_v2 with 2 workers	
Task	
Find the overlapping cells	
Cmd 2	
Apache Sedona	
1 dbutils.notebook.run('sedona_trial', 3600)	
Notebook job #705177502235475	
Command took 30.65 seconds by a.zihna@ckhgbdp.onmicrosoft.com at 30/06/2022, 15:50:27 on DS_Sedona_1_2_0	
Cmd 3	
GeoPandas	
1 dbutils.notebook.run('geopandas_trial', 3600)	
Notebook job #364503638060240	

Command took 1.18 minutes -- by a.zihna@ckhgbdp.onmicrosoft.com at 30/06/2022, 15:46:15 on DS_Sedona_1_2_0

DAIA+AI SUMMIT 2022

GeoPandas Implementation

Cmd 2	
Import Python Modules	Python 🕨 🗕 🗙
1 import geopandas as gpd	
/local_disk0/.ephemeral_nfs/envs/pythonEnv-19d88e75-85df-4291-b273-22c9d2ecb8a5/lib/python3.8/site-packages/geopandas/_compat.py:111: U version (3.10.2-CAPI-1.16.0) is incompatible with the GEOS version PyGEOS was compiled with (3.10.1-CAPI-1.16.0). Conversions between b warnings.warn(serWarning: The Shapely GEOS oth will be slow.
Command took 0.44 seconds by a.zihna@ckhgbdp.onmicrosoft.com at 13/06/2022, 16:22:26 on DE_General	
Cmd 3	
Read Files	
<pre>1 df_3g = gpd.read_file('/dbfs/dais/vor_dais_3g.geojson') 2 df_4g = gpd.read_file('/dbfs/dais/vor_dais_4g.geojson')</pre>	
Command took 1.56 minutes by a.zihna@ckhgbdp.onmicrosoft.com at 13/06/2022, 16:22:26 on DE_General	
Cmd 4	
Get the intersections of polygons	
1 gpd.overlay(df_3g, df_4g, keep_geom_type=False).head()	

Sedona Implementation - 1

Cmd 2

Import Sedona Modules and Register

- 1 from sedona.core.formatMapper import GeoJsonReader
- 2 **from** sedona.register.geo_registrator **import** SedonaRegistrator
- 3 **from** sedona.utils.adapter **import** Adapter
- 4
- 5 SedonaRegistrator.registerAll(spark)
- 6 sc = spark.sparkContext
- 7 sc.setSystemProperty("sedona.global.charset", "utf8")

Cmd 3

Read in the data as Spatial RDDs

- 1 rdd_3g = GeoJsonReader.readToGeometryRDD(sc, 'dbfs:/dais/vor_dais_3g.geojson')
- 2 df_3g = Adapter.toDf(rdd_3g, spark).createOrReplaceTempView('df3g')
- 4 rdd_4g = GeoJsonReader.readToGeometryRDD(sc, 'dbfs:/dais/vor_dais_4g.geojson')
- 5 df_4g = Adapter.toDf(rdd_4g, spark).createOrReplaceTempView('df4g')



Sedona Implementation – 2



Result Comparison

GeoPandas

	ci3g	ci4g
0	10-89-13810-6247	10-89-0-313681-2
1	10-89-13810-6247	10-89-0-313681-2
2	10-89-13810-6247	10-89-0-313681-2
3	10-89-13810-3093	10-89-0-313681-2
4	10-89-13810-3093	10-89-0-313681-2

Apache Sedona

	ci3g 🔺	ci4g 🔺
1	510-89-10942-3032	10-89-10699-110096
2	510-89-10942-3032	10-89-10699-110096
3	510-89-10942-3702	10-89-10699-110096
4	510-89-10942-3702	10-89-10699-110096
5	510-89-10942-3032	10-89-10699-110096
6	510-89-10942-3702	10-89-10699-110096
7	510-89-10942-3218	10-89-10699-110096