# Apache Spark & Parquet in AWS

## Thaer Khawaja

Hothead Games

✉ tkhawaja@hotheadgames.com

April 25, 2017

# Table of Contents

1 Case Study

2 Tech Overview

3 ETL Integration

4 Documentation

5 Dev Horizon

# Case Study

- First game to encode and log in-mission gameplay event data
    - decodes into processable data inflates data up to $\sim 30x$ the size
    - contains all game events
        - damage taken per character, where and when
        - enemies killed
        - abilities used
        - "ghost" actions (team AI kills, actions, maneuvers, etc.)
        - and more . . .
    - valuable for gameplay parameter tunings
- Analytics data warehouse (Redshift) not a suitable intermediary
    - postprocessing summary required on intermediary decoded data would cause heavy load
    - would impact response time on actively used cluster by analysts and data scientists
    - inflated data size would require us to unnecessarily provision cluster upwards
    - postprocessed summaries still need to be on Redshift (for now)

# Tech Overview

1 Case Study

2 Tech Overview
   - Apache Spark
   - Apache Parquet
   - ETL Pipeline

3 ETL Integration

4 Documentation

5 Dev Horizon

Highlights:

- Supports batch and "real-time" (mini-batch) data processing
  - Guarantees no duplicate processed output
- AWS EMR-friendly (runs on Hadoop YARN)
- Commercially available Databricks solution (runs on Mesos)
- "All the rage" framework with large community support

Lowlights:

- Not a true real-time streaming framework
  - Not an issue for our use case
  - Seasoned alternative: Apache Storm with Trident
  - New alternatives: Apache Flink and Apache Apex
- Databricks DBU units cost \$\$\$
  - One r3.8xlarge instance = 8 DBUs

Highlights:

- Column-oriented format
    - Space-efficient data encodings
    - Queries only the data that is being ask for
    - Supports predicate pushdown
- Really good at dealing with wide and nested JSON data
- Automated schema discovery
    - Not magical and breaks on complex loosely structured schemas
- Integrated support with Apache Spark (and others)
- Benchmarks put Parquet as best performer for our use case

Lowlights:

- ORC seems to outperform on flatter data schemas
    - Possibly because Parquet doesn't have bloom filters
      (PARQUET-41)

# Current ETL Pipeline Tech Stack

**S3** - Data lake

**Redshift** - Data warehouse

**EMR** - ETL workhorse

**Data Pipeline** - Workflow scheduler

**Kinesis** - Data streams

**Python** - Code glue, MapReduce, and business logic

# ETL Integration

Apache Spark &
Parquet in AWS

Thaer Khawaja

Case Study

Tech Overview

ETL Integration
Diagram
Parquet
Conversion
Spark Submit

Documentation

Dev Horizon

1 Case Study

2 Tech Overview

3 ETL Integration
- Diagram
- Parquet Conversion
- Spark Submit

4 Documentation

5 Dev Horizon

# ETL Integration - Diagram

Apache Spark & Parquet in AWS

Thaer Khawaja

Case Study

Tech Overview

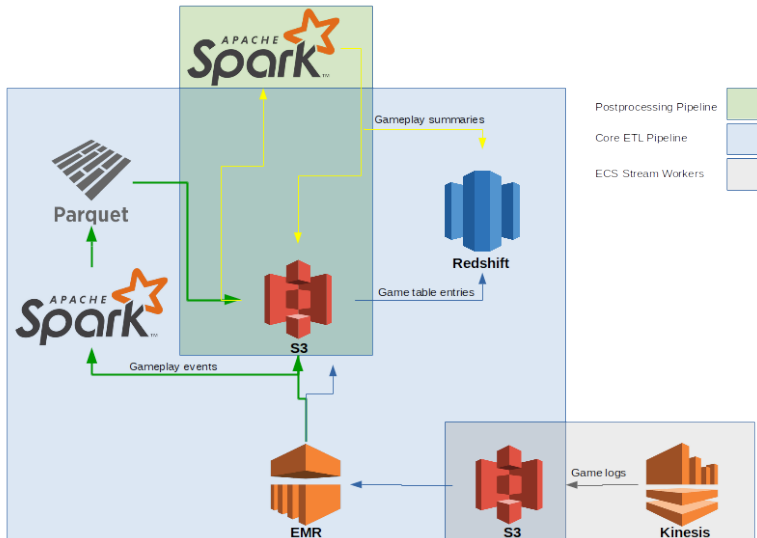ETL Integration
Diagram
Parquet Conversion
Spark Submit

Documentation

Dev Horizon

Gameplay summaries

Game table entries

Gameplay events

Game logs

Parquet

Redshift

S3

EMR

Kinesis

Postprocessing Pipeline

Core ETL Pipeline

ECS Stream Workers

# Parquet Conversion - TSV and JSON

Apache Spark &
Parquet in AWS

Thaer Khawaja

Case Study

Tech Overview

ETL Integration
  Diagram
  **Parquet
  Conversion**
  Spark Submit

Documentation

Dev Horizon

```python
54 # DataFrames need explicit typing since we're reading split lines from a file
55 def tsv_tables_to_parquet(sc, sqc, table_imports, s3_import_path, parquet_path):
56     game_tables = [
57         (csv_import.redshift_table.table, csv_import.import_path)
58         for csv_import in table_imports
59     ]
60
61     for table, name in game_tables:
62         import_path = '{}{}/'.format(s3_import_path, name)
63         if not path_exists(sc, import_path):
64             continue
65         table_lines = sc.textFile(import_path)
66
67         # Spark freaks out if mission_table is referenced directly (namedlist attribute not pickleable)
68         validator = table.validator
69         table_rows = table_lines.map(lambda l: validator(l.split('\t')))
70         struct_type = to_struct_type(table)
71         df = sqc.createDataFrame(table_rows, struct_type)
72         create_or_append_parquet(df, parquet_path, name)
73
74
75 # Depends on Parquet's schema discovery
76 def transforms_to_parquet(sqc, transform_imports, s3_import_path, parquet_path):
77     for transform_import in transform_imports:
78         import_path = '{}{}/'.format(s3_import_path, transform_import)
79         # TODO-TK: parallelize s3 load
80         df = sqc.read.json(import_path)
81         create_or_append_parquet(df, parquet_path, transform_import)
```

```python
28  struct_types = {
29      BigInteger: LongType,
30      Boolean: BooleanType,
31      Float: FloatType,
32      Integer: IntegerType,
33      String: StringType,
34      StringDateTime: StringType # TODO-TK: deal with funkiness in TimestampType
35  }
36
37
38  def to_struct_type(table):
39      return StructType(
40          [
41              StructField(column.name, struct_types[column.data_type.__class__]())
42              for column in table.columns
43          ]
44      )
45
46
47  def create_or_append_parquet(df, parquet_path, table_name):
48      df.write.parquet(
49          '{}{}/'.format(parquet_path, table_name),
50          mode='append' if is_incremental(table_name) else 'overwrite'
51      )
```

```python
 8  def pyspark_step(script, custom_args, py_files, spark_args):
 9      return {
10          'Jar': 'command-runner.jar',
11          'Args': [
12              'spark-submit',
13              '--master',
14              'yarn',
15              '--deploy-mode',
16              'cluster',
17              '--conf',
18              'spark.executorEnv.PYSPARK_PYTHON=python2.7',
19              '--conf',
20              'spark.yarn.appMasterEnv.PYSPARK_PYTHON=python2.7',
21          ] + spark_args + ['--py-files', py_files, script] + custom_args
22      }
23
```

```python
47  def job_flow_step(script, custom_args, py_files=None, main_class=None, spark_args=[]):
48      if re.search('.jar$', script):  # h4x
49          return jar_step(script, custom_args, main_class)
50      else:
51          return pyspark_step(script, custom_args, py_files, spark_args)
52
53
54  def submit_spark_job(cluster_id, script, name, custom_args, py_files=None, main_class=None, spark_args=[]):
55      conn = boto3.Session().client('emr')
56      resp = conn.add_job_flow_steps(
57          JobFlowId=cluster_id,
58          Steps=[
59              {
60                  'Name': name,
61                  'ActionOnFailure': 'CONTINUE',
62                  'HadoopJarStep': job_flow_step(script, custom_args, py_files, main_class, spark_args)
63              }
64          ]
65      )
66      track_progress(conn, resp['StepIds'][0], cluster_id)
```

# Documentation

Apache Spark &
Parquet in AWS

Thaer Khawaja

Case Study
Tech Overview
ETL Integration
Documentation
Dev Horizon

PySpark docs: https://spark.apache.org/docs/2.1.0/api/python/pyspark.html
Boto3 docs: https://boto3.readthedocs.io/en/latest/

# Dev Horizon

1 Case Study

2 Tech Overview

3 ETL Integration

4 Documentation

5 Dev Horizon

Apache Spark &
Parquet in AWS

Thaer Khawaja

Case Study
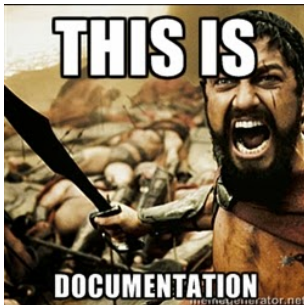
Tech Overview

ETL Integration

Documentation

Dev Horizon

- Add a Parquet copy of our data lake
    - ad-hoc queries for a given user or set of users
    - Partitioning by day/week/month
    - integration tests with randomly sampled data
- Leverage Spark Streaming to improve ETL data freshness
- Migrate post-processing queries from Redshift to Spark
- Alternatives away from Redshift as a data warehouse (???)
    - AWS Athena supports ANSI SQL and has JDBC driver
    - Spark SQL is mature and more testing-friendly

Apache Spark &
Parquet in AWS

Thaer Khawaja

Case Study
Tech Overview
ETL Integration
Documentation
Dev Horizon

Questions?