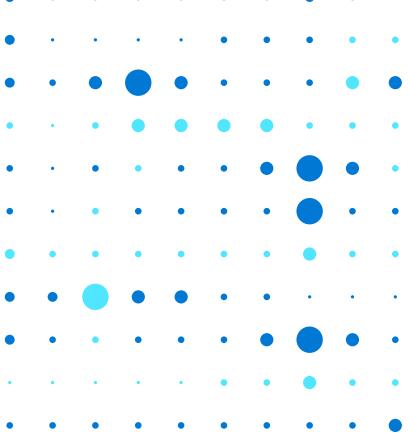




Solving business needs with Delta Lake

Contents

Introduction	3
Scenario 1: Streaming financial stock data analysis	4
Scenario 2: Genomic data analytics	16
Scenario 3: Real-time display advertising attribution	36
Scenario 4: Mobile Gaming Data Event Processing	45
Summary.....	51
About Azure Databricks.....	52
Resources.....	53



Introduction

Data professionals want to help their organizations innovate for competitive advantage by making the most of their data. Good quality, reliable data forms the foundation for success for such analytics and machine learning initiatives.

[Delta Lake](#) is an open source storage layer that brings reliability to data lakes, which are a key component of modern data architectures. It provides:

- ACID transactions
- Scalable metadata handling
- Unified streaming and batch data processing

Delta Lake runs on top of your existing data lake and is fully compatible with Apache Spark™ APIs. Delta Lake on Azure Databricks allows you to configure Delta Lake based on your workload patterns and provides optimized layouts and indexes for fast interactive queries. Seeing how Delta Lake can be used to address various business needs can help you get a good sense of how you can deploy it to help you enable downstream users with reliable data.

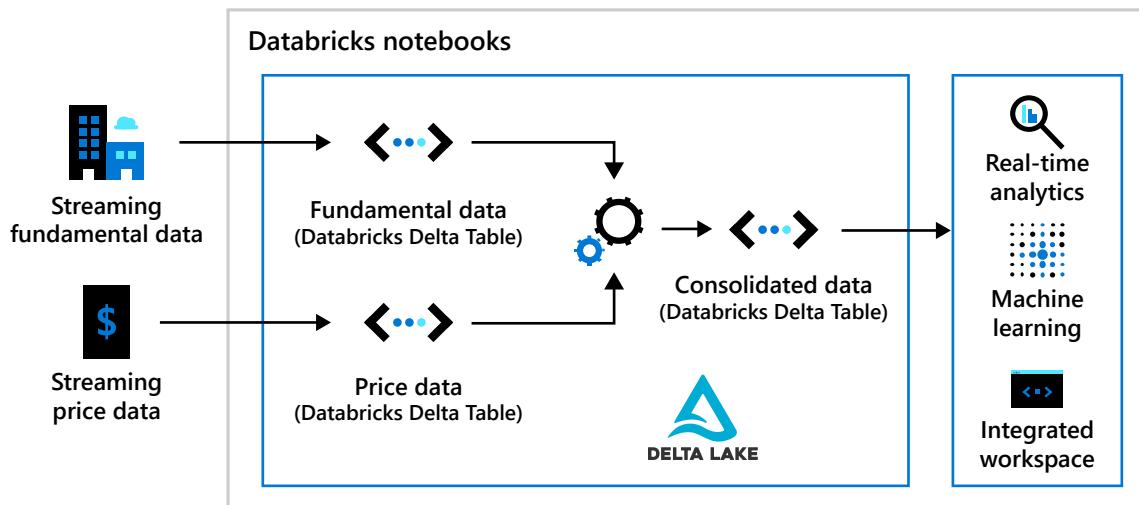
In this guide, we will walk you through four examples of applying Delta Lake to solve business needs. Specifically, they cover building data pipelines for:

- Streaming financial stock data analysis that delivers transactional consistency of legacy and streaming data concurrently
- Genomic data analytics used for analyzing population- scale genomic data
- Real-time display advertising attribution for delivering information on advertising spend effectiveness
- Mobile gaming data event processing to enable fast metric calculations and responsive scaling

Scenario 1: Streaming financial stock data analysis

Real-time analysis of stock data has been a complicated endeavor due to the complexities of maintaining a streaming system and ensuring transactional consistency of legacy and streaming data concurrently. Delta Lake helps solve many of the pain points of building a streaming system to analyze stock data in real-time.

In the following diagram, we provide a high-level architecture to simplify this problem. We start by ingesting two different sets of data into two Delta Lake tables. The two datasets are stock prices and fundamentals. After ingesting the data into their respective tables, we then join the data in an ETL process and write the data out into a third Delta Lake table for downstream analysis.



In this section we will review:

- The typical problems of running real-time stock analysis
- How Delta Lake addresses these problems
- How to implement the system using Azure Databricks (note: the system can be implemented with just open source software)

Scenario 1: Streaming financial stock data analysis

Typical pain points before Delta Lake

The pain points of a traditional streaming and data warehousing solution can be broken into two groups: data lake pains and data warehouse pains.

Data Lake pain points

While data lakes allow you to flexibly store an immense amount of data in a file system, there are many pain points, such as:

- Consolidation of streaming data from many disparate systems is difficult.
- Updating data in a data lake is nearly impossible and much of the streaming data needs to be updated as changes are made. This is especially important in scenarios that involve financial reconciliation and subsequent adjustments.
- Query speeds for a data lake are typically very slow.

Data Warehouse pain points

The power of a data warehouse is that you have a persistent performant store of your data. But the pain points for building modern continuous applications include:

- Constrained to SQL queries; i.e., no machine learning or advanced analytics.
- Accessing streaming data and stored data together is very difficult if not impossible.
- Data warehouses do not scale very well.
- Tying compute and storage together makes using a warehouse very expensive.

Scenario 1: Streaming financial stock data analysis

How Delta Lake Solves These Issues

Because Delta Lake ([Delta Lake Guide](#)) runs on top of your existing data lake and is fully compatible with Apache Spark™ APIs:

- Your streaming/data lake/warehousing solution has ACID compliance.
- Delta Lake, along with Structured Streaming, makes it possible to analyze streaming and historical data together at data warehouse speeds.
- Using Delta Lake tables as sources and destinations of streaming big data makes it easy to consolidate disparate data sources.
- Upserts are supported on Delta Lake tables so changes are simpler to manage.
- You can easily include machine learning scoring and advanced analytics into ETL and queries.
- Compute and storage are decoupled resulting in a completely scalable solution.

Implement Your Streaming Stock Analysis Solution with Delta Lake

Delta Lake and Apache Spark do most of the work for our solution; you can run this Azure Databricks set up notebook so it can automatically download the generated source data and start loading data into a file location and then try out the full Azure Databricks notebook and follow along with the code samples below. There is also an [Azure Databricks setup notebook](#) to help you get started. First run this notebook so it can automatically download the generated source data and start loading data into a file location.

Let's start by enabling Delta Lake.

Scenario 1: Streaming financial stock data analysis

As noted in the preceding diagram, we have two datasets to process – one for fundamentals and one for price data. To create our two Delta Lake tables, we specify the `.format("delta")` against our DBFS locations.

```
# Create Fundamental Data (Databricks Delta table)
dfBaseFund = spark \
    .read \
    .format('delta') \
    .load('/delta/stocksFundamentals')

# Create Price Data (Databricks Delta table)
dfBasePrice = spark \
    .read \
    .format('delta') \
    .load('/delta/stocksDailyPrices')
```

While we're updating the `stockFundamentals` and `stocksDailyPrices`, we will consolidate this data through a series of ETL jobs into a consolidated view (`stocksDailyPricesWFund`). With the following code snippet, we can determine the start and end date of available data and then combine the price and fundamentals data for that date range into DBFS.

```
# Determine start and end date of available data
row = dfBasePrice.agg(
    func.max(dfBasePrice.price_date).alias("maxDate"),
    func.min(dfBasePrice.price_date).alias("minDate")
).collect()[0]
startDate = row["minDate"]
endDate = row["maxDate"]

# Define our date range function
def daterange(start_date, end_date):
    for n in range(
        int((end_date - start_date).days)):
        yield start_date + datetime.timedelta(n)
```

Scenario 1: Streaming financial stock data analysis

```
# Define combinePriceAndFund information by date
def combinePriceAndFund(theDate):
    dfFund = dfBaseFund.where(
        dfBaseFund.price_date == theDate)
    dfPrice = dfBasePrice.where(
        dfBasePrice.price_date == theDate) \
            .drop('price_date')
    # Drop the updated column
    dfPriceWFund = dfPrice.join(dfFund,[“ticker”]) \
        .drop(“updated”)

# Save data to DBFS
dfPriceWFund.write.format(“delta”).mode(“append”) \
    .save(“/delta/stocksDailyPricesWFund”)

# Loop through dates to complete fundamentals
# + price ETL process
for single_date in daterange(startDate,
    (endDate + datetime.timedelta(days=1))):
    start = datetime.datetime.now()
    combinePriceAndFund(single_date)
    end = datetime.datetime.now()
```

The result is a stream of consolidated fundamentals and price data that is being pushed into DBFS in the /delta/ stocksDailyPricesWFund location. We can build a Delta Lake table by specifying .format("delta") against that DBFS location.

```
dfPriceWithFundamentals = spark
.readStream
.format("delta")
.load("/delta/stocksDailyPricesWFund")
# Create temporary view of the data
dfPriceWithFundamentals.
createOrReplaceTempView("priceWithFundamentals")
```

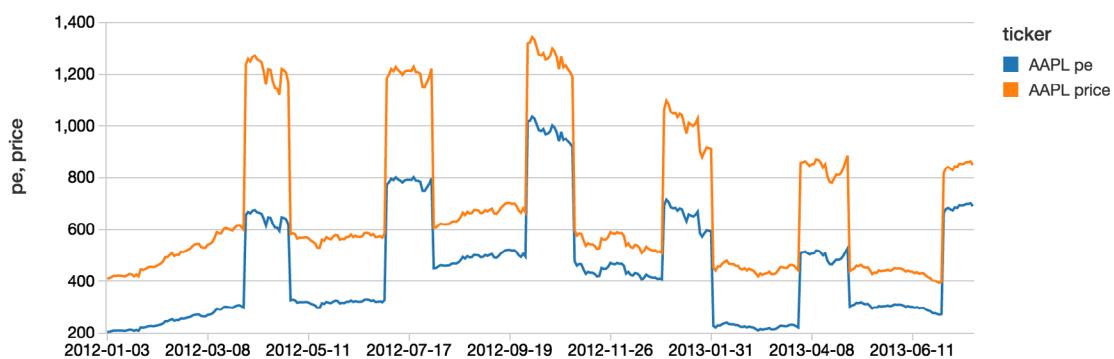
Scenario 1: Streaming financial stock data analysis

Now that we have created our initial Delta Lake table, let's create a view that will allow us to calculate the price/earnings ratio in real time (because of the underlying streaming data updating our Delta Lake table).

```
%sql  
CREATE OR REPLACE TEMPORARY VIEW viewPE AS  
select ticker,  
       price_date,  
       first(close) as price,  
       (close/eps_basic_net) as pe  
from priceWithFundamentals  
where eps_basic_net > 0
```

With the view in place, you can quickly analyze our data [using Spark SQL](#).

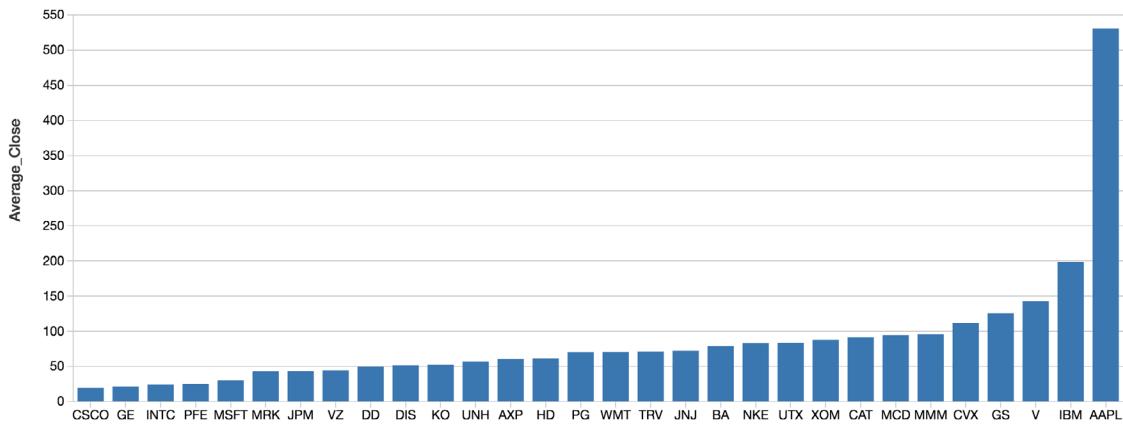
```
%sql  
select *  
from viewPE  
where ticker == "AAPL"  
order by price_date
```



As the underlying source of this consolidated dataset is a Delta Lake table, this view isn't just showing the batch data but also any new streams of data that are coming in as per the following streaming dashboard.

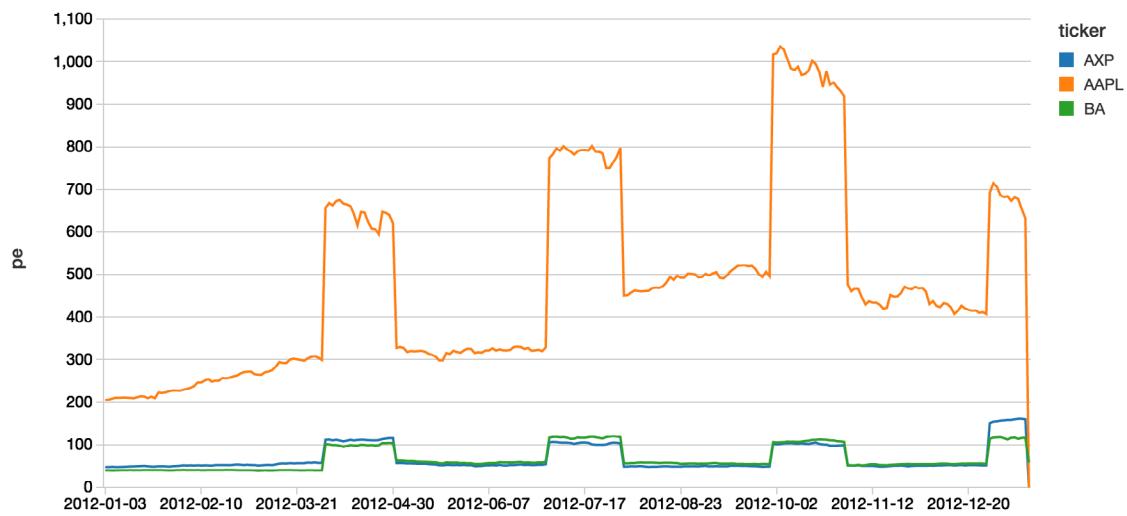
Underneath the covers, Structured Streaming isn't just writing the data to Delta Lake tables but also keeping the state of the distinct number of keys (in this case ticker symbols) that need to be tracked.

Scenario 1: Streaming financial stock data analysis



Because you are using Spark SQL, you can execute aggregate queries at scale and in real-time.

```
%sql
SELECT ticker, AVG(close) as Average_Close
FROM priceWithFundamentals
GROUP BY ticker
ORDER BY Average_Close
```



Scenario 1: Streaming financial stock data analysis

The difference between Delta Lake and Apache Parquet

While the Delta Lake stores data in Apache Parquet format, it includes features that allow data lakes to be reliable at scale. These features include:

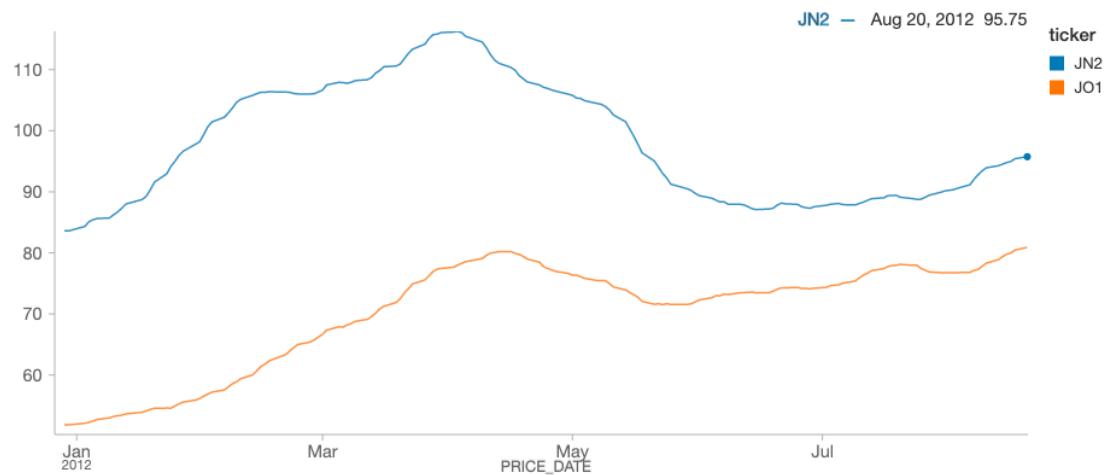
- **ACID Transactions:** Delta Lake ensures data integrity and provides serializability.
- **Scalable Metadata handling:** For big data systems, the metadata itself is often “big” enough to slow down any system that tries to make sense of it, let alone making sense of the actual underlying data. Delta Lake treats metadata like regular data and leverages Apache Spark’s distributed processing power. As a result, Delta Lake can handle petabyte-scale tables with billions of partitions and files at ease.
- **Time travel (data versioning):** Creates snapshots of data, allowing you to access and revert to earlier versions of data for audits, rollbacks or to reproduce experiments.
- **Open format:** All data in Delta Lake is stored in Apache Parquet format, thereby enabling Delta Lake to leverage the efficient compression and encoding schemes that are native to Parquet.
- **Unified batch and streaming source and sink:** A table in Delta Lake is both a batch table, as well as a streaming source and sink. Streaming data ingest, batch historic backfill, and interactive queries all just work out of the box.
- **Schema enforcement:** Delta Lake provides the ability to specify your schema and enforce it. This helps ensure that the data types are correct and required columns are present, preventing bad data from causing data corruption.
- **Schema evolution:** Big data is continuously changing. Delta Lake enables you to make changes to a table schema that can be applied automatically, without the need for cumbersome DDL.
- **100% Compatible with Apache Spark API:** Developers can use Delta Lake with their existing data pipelines with minimal change as it is fully compatible with Spark, the commonly used big data processing engine.

Scenario 1: Streaming financial stock data analysis

How to view the Delta Lake Table for both streaming and batch near the beginning of the notebook

As noted in the [Streaming Stock Analysis with Delta Lake notebook](#), in cell 8 we ran the following batch query:

```
dfPrice = \
spark.read.format("delta").load(deltaPricePath)
display(dfPrice.where(
dfPrice.ticker.isin({'JN2', 'JO1'})))
```



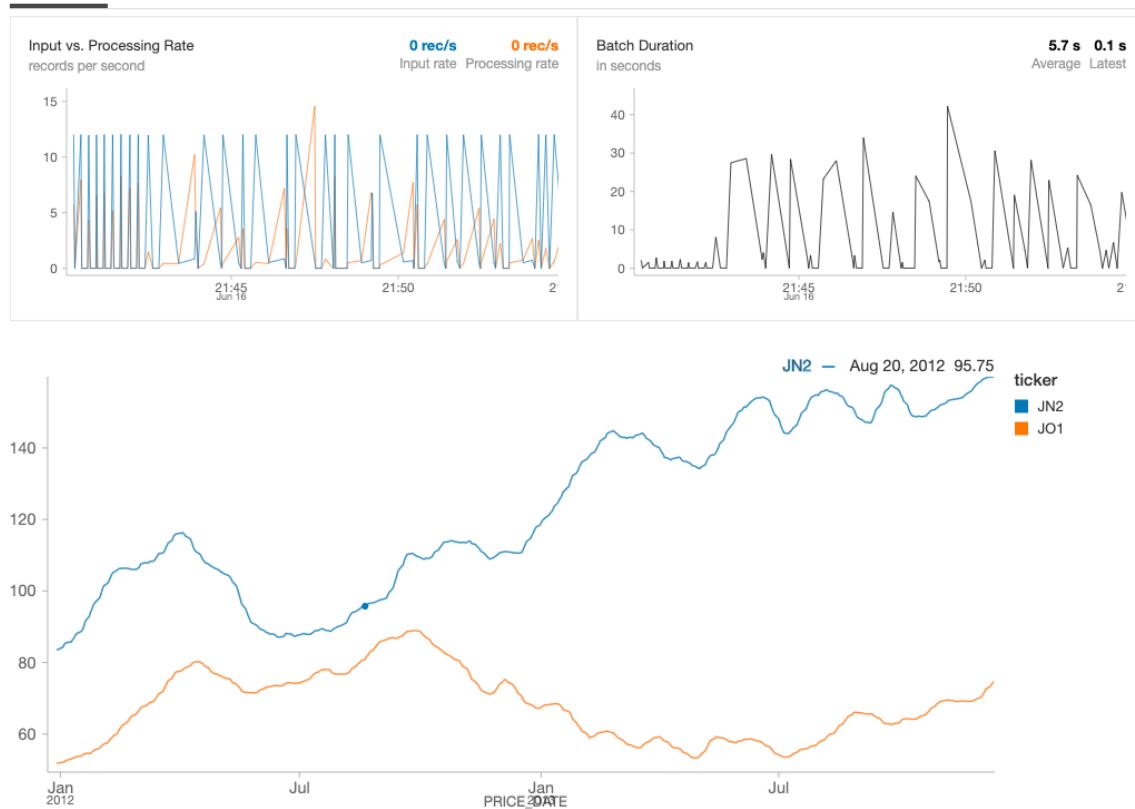
Notice that we ran this query earlier in the cycle with data up until August 20, 2012. Using the same folder path (deltaPricePath), we also created a structured streaming DataFrame via the following code snippet in cell 4:

```
# Create Stream and Temp View for Price
dfPriceStream =
spark.readStream.format("delta").load(deltaPricePath)
dfPriceStream.createOrReplaceTempView("priceStream")
```

Scenario 1: Streaming financial stock data analysis

We can then run the following real-time [Spark SQL](#) query that will continuously refresh.

```
%sql  
SELECT *  
FROM priceStream  
where ticker in ('JN2', 'JO1')
```



Notice that, even though the batch query executed earlier (and ended at August 20, 2012), the structured streaming query continued to process data long past that date (the small blue dot denotes where August 20, 2012 is on the streaming line chart). As you can see from the preceding code snippets, both the batch and structured streaming DataFrames query from the same folder path of `deltaPricePath`.

Scenario 1: Streaming financial stock data analysis

Finding a mistake and correcting it for auditing purposes

Delta Lake has a data versioning feature called Time Travel. It provides snapshots of data, allowing you to access and revert to earlier versions of data for audits and rollbacks or to reproduce experiments. To visualize this, note cells 36 onwards in [the Streaming Stock Analysis with Delta Lake notebook](#). The screenshot (opposite panel) shows three different queries using the VERSION AS OF syntax allowing you to view your data by version (or by timestamp using the TIMESTAMP syntax).

With this capability, you can know what changes to your data were made and when those transactions had occurred.

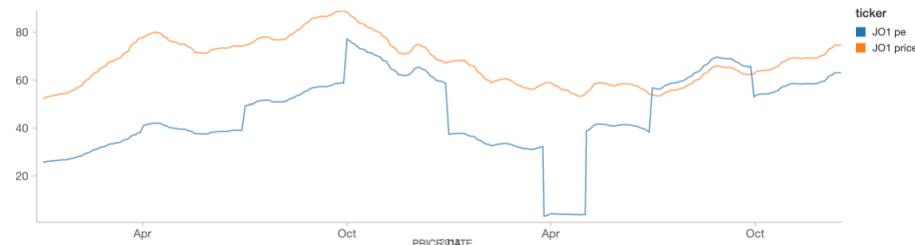
```
%sql  
SELECT ticker, price_date, adj_close as price, (adj_close/eps_basic_net) as pe  
FROM priceWithFundamentalsHistory VERSION AS OF 1  
WHERE ticker == "JO1"
```



```
%sql  
SELECT ticker, price_date, adj_close as price, (adj_close/eps_basic_net) as pe  
FROM priceWithFundamentalsHistory VERSION AS OF 20  
WHERE ticker == "JO1"
```



```
%sql  
SELECT ticker, price_date, adj_close as price, (adj_close/eps_basic_net) as pe  
FROM priceWithFundamentalsHistory  
WHERE ticker == "JO1"
```



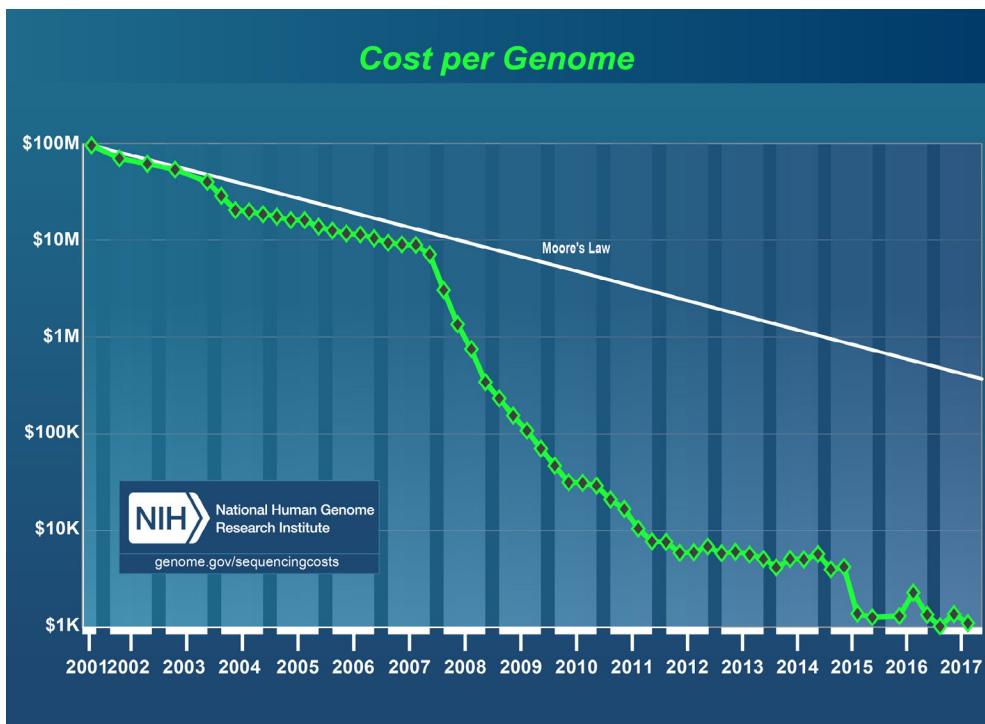
Scenario 1: Streaming financial stock data analysis

Summary

In closing, we demonstrated how to simplify streaming stock data analysis using [Delta Lake](#). By combining Spark Structured Streaming and Delta Lake, we have a scalable solution that has the advantages of both data lakes and data warehouses. This removes the data engineering complexities commonly associated with streaming and transactional consistency enabling data engineering and data science teams to focus on understanding the trends in their stock data.

Scenario 2: Genomic data analytics

Download the [Simplifying Genomics Pipelines at Scale with Databricks \(Python\)](#) notebook and try it out.



Source: DNA Sequencing Costs: Data

Since the completion of the Human Genome Project in 2003, there has been an explosion in data fueled by a dramatic drop in the cost of DNA sequencing, from \$3B for the first genome to under \$1,000 today.¹

Consequently, the field of genomics has now matured to a stage where companies have started to do DNA sequencing at population-scale. However, sequencing the DNA code

¹ The Human Genome Project was a \$3B project led by the Department of Energy and the National Institutes of Health began in 1990 and completed in 2003.

Scenario 2: Genomic data analysis

is only the first step; the raw data then needs to be transformed into a format suitable for analysis. Typically, this is done by gluing together a series of bioinformatics tools with custom scripts and processing the data on a single node, one sample at a time, until we wind up with a collection of genomic variants. Bioinformatics scientists today spend most of their time building out and maintaining these pipelines. As genomic data sets have expanded into the petabyte scale, it has become challenging to answer even the following simple questions promptly:

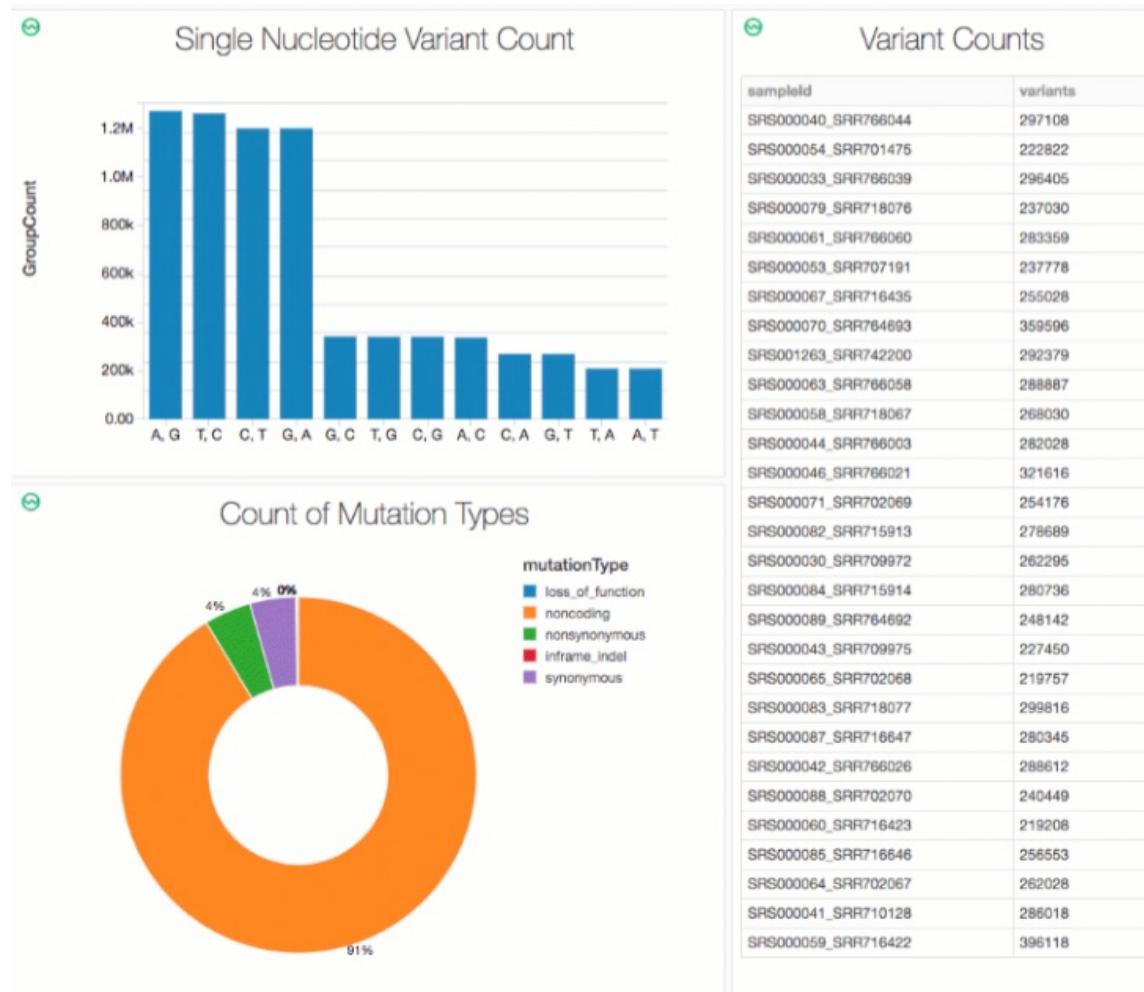
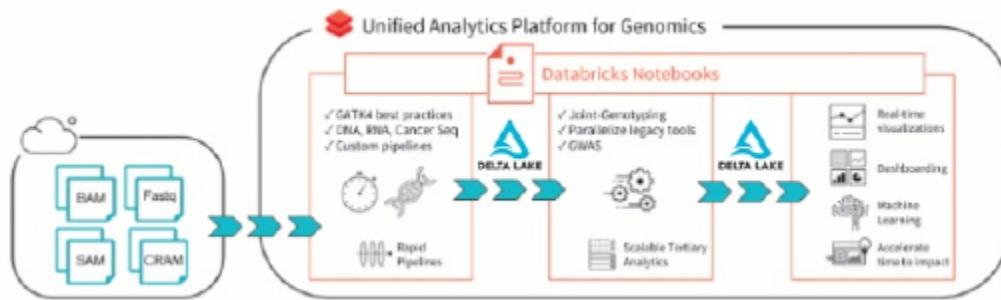
- How many samples have we sequenced this month?
- What is the total number of unique variants detected?
- How many variants did we see across different classes of variation?

Further compounding this problem, data from thousands of individuals cannot be stored, tracked nor versioned while also remaining accessible and queryable. Consequently, researchers often duplicate subsets of their genomic data when performing their analyses, causing the overall storage footprint and costs to escalate. In an attempt to alleviate this problem, today researchers employ a strategy of "data freezes," typically between six months to two years, where they halt work on new data and instead focus on a frozen copy of existing data. There is no solution that incrementally builds up analysis over shorter time frames, causing research progress to slow down.

Scenario 2: Genomic data analysis

Architecture for end-to-end genomics analysis with Azure Databricks

Note that this architecture can also be achieved with open source software.



Scenario 2: Genomic data analysis

There is a compelling need for robust software that can consume genomic data at industrial scale while also retaining the flexibility that scientists need to explore the data, iterate on their analytical pipelines, and derive new insights.

With Delta Lake, you can store all your genomic data in one place and create analyses that update in real-time as new data is ingested. As an example, the following dashboard shows quality control metrics and visualizations that can be calculated and presented in an automated fashion and customized to suit your specific requirements.

In the rest of this section, we will walk through the steps we took to build the quality control dashboard above, which updates in real time as samples finish processing. By using a Delta Lake-based pipeline for processing genomic data, data teams can now operate their pipelines to provide real-time, sample-by-sample visibility. With Azure Databricks notebooks (and integrations such as GitHub and MLflow) they can track and version analysis to ensure their results are reproducible. Their bioinformaticians can devote less time to maintaining pipelines and spend more time making discoveries.

The net effect is to help drive the transformation from ad-hoc analysis to production genomics on an industrial scale, enabling better insights into the link between genetics and disease.

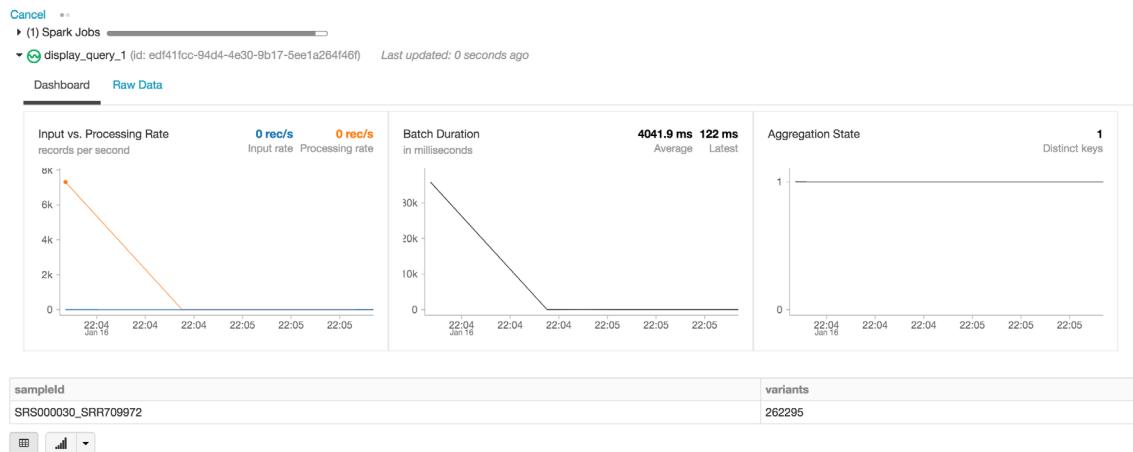
Read Sample Data

Let's start by reading variation data from a small cohort of samples; the following statement reads in data for a specific sampleId and saves it using the Delta Lake format (in the delta_stream_output folder).

```
# Read the initial sample annotation file
spark.read \
    .format("parquet") \
    .load("dbfs:/annotations_etl_parquet/sampleId=" +
          "SRS000030_SRR709972") \
    .write \
    .format("delta") \
    .save(delta_stream_outpath)
```

Note, the annotations_etl_parquet folder contains annotations generated from the 1000 genomes dataset stored in parquet format. The ETL and processing of these annotations were performed using Azure Databricks for Genomics.

Scenario 2: Genomic data analysis



Start Streaming the Delta Lake Table

In the following statement, we are creating the `exomes` Apache Spark DataFrame, which is reading a stream (via `readStream`) of data using the Delta Lake format. This is a continuously running or dynamic DataFrame, i.e., the `exomes` DataFrame will load new data as data is written into the `delta_stream_output` folder.

To view the `exomes` DataFrame, we can run a DataFrame query to find the count of variants grouped by the `sampleId`.

```
# Read the stream of data
exomes = spark.readStream \
    .format("delta").load(delta_stream_outpath)

# Display the data via DataFrame query
display(exomes.groupBy("sampleId").count() \
    .withColumnRenamed("count", "variants"))
```

When executing the `display` statement, the Azure Databricks notebook provides a streaming dashboard to monitor the streaming jobs (as noted in the preceding figure). Immediately below the streaming job are the results of the `display` statement (i.e. the count of variants by `sample_id`).

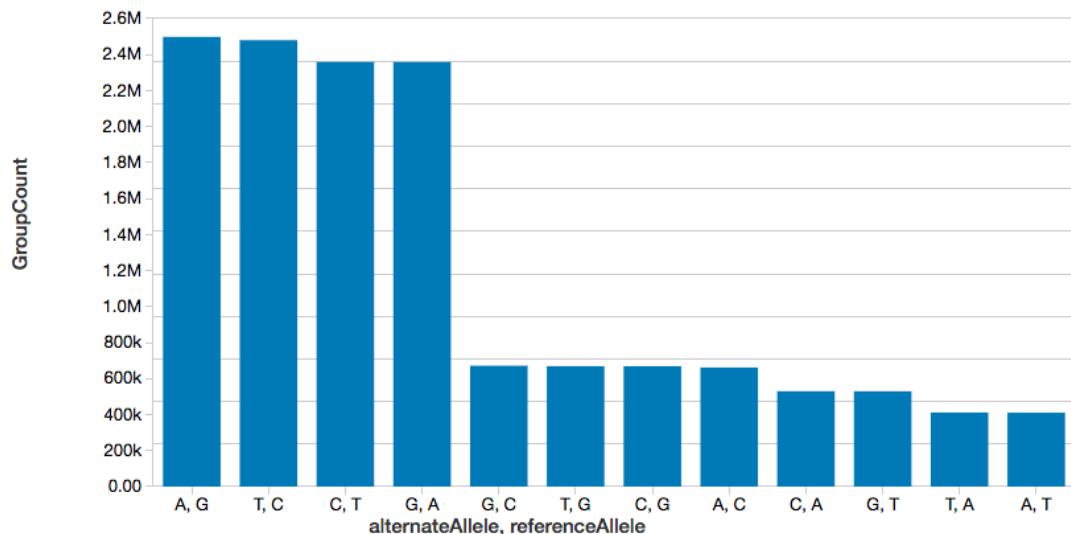
Let's continue answering our initial set of questions by running other DataFrame queries based on our DataFrame.

Scenario 2: Genomic data analysis

Single Nucleotide Variant Count

To continue the example, we can quickly calculate the number of single nucleotide variants (SNVs), as displayed in the following graph.

```
%sql  
select referenceAllele, alternateAllele, count(1) as GroupCount  
    from snvs  
group by referenceAllele, alternateAllele  
order by GroupCount desc
```



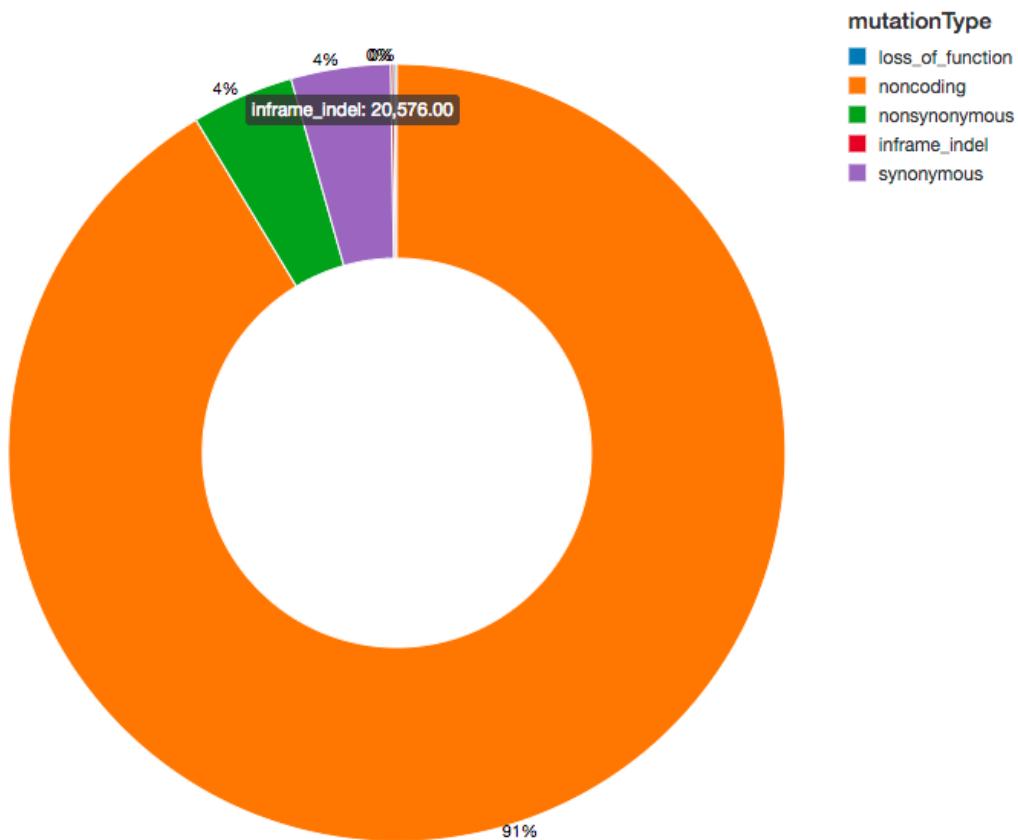
Scenario 2: Genomic data analysis

Variant Count

Since we have annotated our variants with functional effects, we can continue our analysis by looking at the spread of variant effects we see. The majority of the variants detected flank regions that code for proteins, these are known as noncoding variants.

```
display(exomes.groupBy("mutationType").count())
```

Note, the `display` command is part of the Azure Databricks workspace that allows you to view your DataFrame using Azure Databricks visualizations (i.e. no coding required).



Scenario 2: Genomic data analysis

Amino Acid Substitution Heatmap

Continuing with our exomes DataFrame, let's calculate the amino acid substitution counts with the following code snippet. Similar to the previous DataFrames, we will create another dynamic DataFrame (aa_counts) so that as new data is processed by the exomes DataFrame, it will subsequently be reflected in the amino acid substitution counts as well. We are also writing the data into memory (i.e., .format("memory")) and processed batches every 60s (i.e. trigger(processingTime='60 seconds')) so the downstream Pandas heatmap code can process and visualize the heatmap.

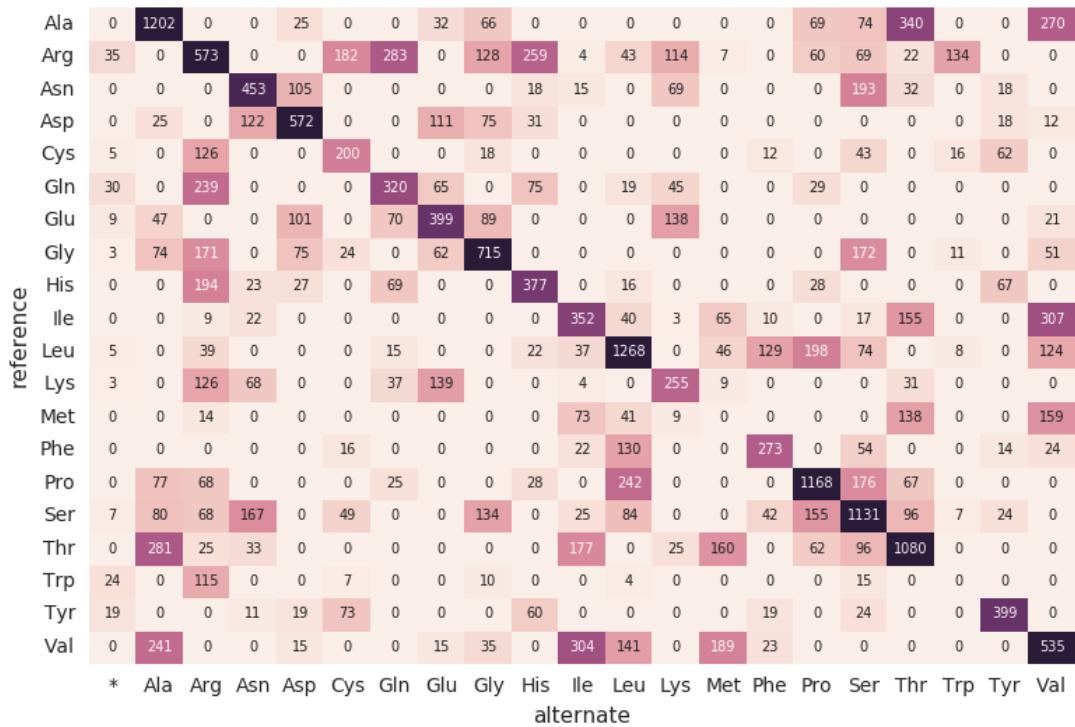
```
# Calculate amino acid substitution counts
coding = get_coding_mutations(exomes)
aa_substitutions =
get_amino_acid_substitutions(coding.select("proteinHgvs"), "proteinHgvs")
aa_counts =
count_amino_acid_substitution_combinations(aa_substitutions)
aa_counts. \
    writeStream. \
    format("memory"). \
    queryName("amino_acid_substitutions"). \
    outputMode("complete"). \
    trigger(processingTime='60 seconds'). \
    start()
```

The following code snippet reads the preceding amino_acid_ substitutions Spark table, determines the max count, creates a new Pandas pivot table from the Spark table, and then plots out the heatmap.

```
# Use pandas and matplotlib to build heatmap
amino_acid_substitutions =
spark.read.table("amino_acid_substitutions")
max_count =
amino_acid_substitutions.agg(fx.max("substitutions")).collect()[0][0]
aa_counts_pd = amino_acid_substitutions.toPandas()
aa_counts_pd = pd.pivot_table(aa_counts_pd,
values='substitutions', index=['reference'],
columns=['alternate'], fill_value=0)
```

Scenario 2: Genomic data analysis

```
fig, ax = plt.subplots()
with sns.axes_style("white"):
    ax = sns.heatmap(aa_counts_pd, vmax=max_count*0.4, cbar=False,
    annot=True, annot_kws={"size": 7}, fmt="d")
plt.tight_layout()
display(fig)
```



Scenario 2: Genomic data analysis

Migrating to a Continuous Pipeline

Up to this point, the preceding code snippets and visualizations represent a single run for a single `sampleId`. But because we're using Structured Streaming and Delta Lake, this code can be used (without any changes) to construct a production data pipeline that computes quality control statistics continuously as samples roll through our pipeline. To demonstrate this, we can run the following code snippet that will load our entire dataset.

```
import time
parquets = "dbfs:/databricks-datasets/genomics/annotations_etl_parquet/"
files = dbutils.fs.ls(parquets)
counter = 0
for sample in files:
    counter += 1
    annotation_path = sample.path
    sampleId = annotation_path.split("/")[4].split("=")[1]
    variants = spark.read.format("parquet"). \
load(str(annotation_path))
    print("running " + sampleId)
    if(sampleId != "SRS000030_SRR709972"):
        variants.write.format("delta"). \
        mode("append"). \
        save(delta_stream_outpath)
time.sleep(10)
```

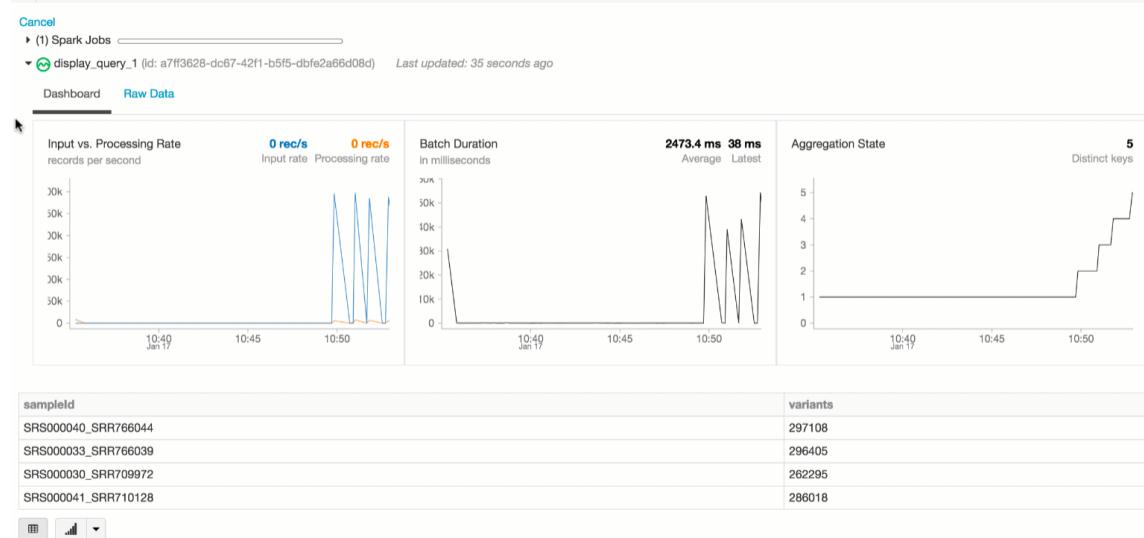
As described in the earlier code snippets, the source of the exomes DataFrame are the files loaded into the `delta_stream_output` folder. Initially, we had loaded a set of files for a single `sampleId` (i.e., `sampleId="SRS000030_SRR709972"`). The preceding code snippet now takes all the generated parquet samples (i.e. parquets) and incrementally loads those files by `sampleId` into the same `delta_stream_output` folder.

Scenario 2: Genomic data analysis

Visualizing Your Genomics Pipeline

When you scroll back to the top of your notebook, you will notice that the `exomes` DataFrame is now automatically loading the new `sampleIds`. Because the structured streaming component of our genomics pipeline runs continuously, it processes data as soon as new files are loaded into the `delta_stream_output` path folder. By using the Delta format, we can ensure the transactional consistency of the data streaming into the `exomes` DataFrame.

```
exomes = spark.readStream.format("delta").load(delta_stream_outpath)
display(exomes.groupBy("sampleId").count().withColumnRenamed("count", "variants"))
```



As opposed to the initial creation of our `exomes` DataFrame, notice how the structured streaming monitoring dashboard is now loading data (i.e., the fluctuating “input vs. processing rate,” fluctuating “batch duration,” and an increase of distinct keys in the “aggregations state”). As the `exomes` DataFrame is processing, notice the new rows of `sampleIds` (and variant counts). This same action can also be seen for the associated group by mutation type query.

Scenario 2: Genomic data analysis

```
import time
parquets = "dbfs:/wbrandler/dnaseq/annotations_etl_parquet/"
files = dbutils.fs.ls(parquets)
counter = 0
for sample in files:
    counter += 1
    annotation_path = sample.path
    sampleId = annotation_path.split("/")[4].split("=")[1]
    variants = spark.read.format("parquet").load(str(annotation_path))
    print("running " + sampleId)
    if(sampleId != "SRS000030_SRR709972"):
        variants.write.format("delta"). \
            mode("append"). \
            save(delta_stream_outpath)
time.sleep(10)
```

Cancel ••• Running command...
▶ (6) Spark Jobs 
▶ 📈 variants: pyspark.sql.dataframe.DataFrame = [contigName: string, start: long ... 12 more fields]
running SRS000030_SRR709972
running SRS000033_SRR766039
running SRS000040_SRR766044

Cmd 26



Scenario 2: Genomic data analysis

With Delta Lake, any new data is transactionally consistent in each step of our genomics pipeline. This is important because it ensures your pipeline is consistent (maintains consistency of your data, i.e., ensures all of the data is “correct”), reliable (either the transaction succeeds or fails completely), and can handle real-time updates (the ability to handle many transactions concurrently and any outage of failure will not impact the data). Thus, even the data in our downstream amino acid substitution map (which had a number of additional ETL steps) is refreshed seamlessly.

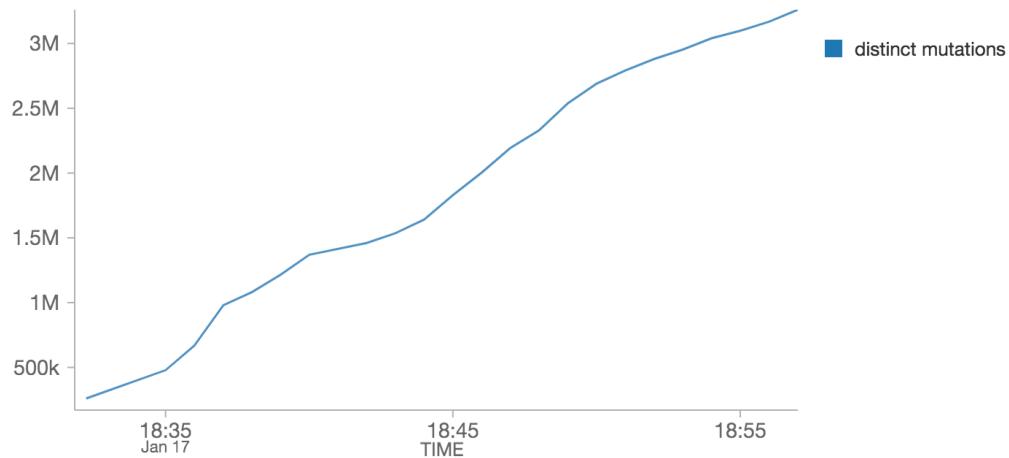
	Ala	0	35306	0	0	912	0	0	1098	1932	0	0	0	0	0	2148	2451	0	10101	0	0	7940			
reference	Arg	977	0	16282	0	0	5088	8095	0	3968	7182	141	1374	3404	279	0	1661	2042	0	726	3662	0	0		
	Asn	0	0	0	12613	2954	0	0	0	468	470	0	1753	0	0	0	5189	0	948	0	462	0	0		
	Asp	0	732	0	3675	16658	0	0	3232	2164	844	0	0	0	0	0	0	0	0	0	0	702	485		
	Cys	164	0	3526	0	0	6226	0	0	585	0	0	0	0	0	429	0	1348	0	0	0	488	1946	0	
	Gln	743	0	7138	0	0	0	9268	2038	0	2390	0	614	1302	0	0	798	0	0	0	0	0	0	0	
	Glu	280	1355	0	0	3146	0	2115	10795	2735	0	0	0	4218	0	0	0	0	0	0	0	0	0	627	
	Gly	120	2128	4860	0	2365	775	0	2204	21045	0	0	0	0	0	0	4776	0	0	0	306	0	1687	0	
	His	0	0	5520	745	665	0	2025	0	0	10857	0	461	0	0	0	776	0	0	0	0	0	1955	0	
	Ile	0	0	184	586	0	0	0	0	0	10653	1190	74	1950	399	0	443	0	4427	0	0	0	8730	0	
	Leu	177	0	1124	0	0	0	510	0	0	629	1032	36255	0	1464	3589	5678	2139	2	0	232	0	0	3495	0
	Lys	64	0	3533	1945	0	0	1178	4138	0	0	139	0	7857	266	0	0	0	995	0	0	0	0	0	0
	Met	0	0	319	0	0	0	0	0	0	2320	1312	270	0	0	0	0	4292	0	0	0	0	0	4841	0
	Phe	0	0	0	0	0	439	0	0	0	596	3859	0	0	8297	0	1694	0	0	0	0	592	530	0	
	Pro	0	2276	1565	0	0	0	695	0	0	945	0	7367	0	0	0	34763	5108	0	1793	0	0	0	0	
	Ser	195	2357	1839	4743	0	1684	0	0	3994	0	719	2438	0	0	0	1463	4398	33024	0	2849	232	685	0	
	Thr	0	8318	775	1153	0	0	0	0	0	4861	0	768	4555	0	1961	2874	0	31761	0	0	0	0	0	
	Trp	781	0	3058	0	0	345	0	0	291	0	0	168	0	0	0	0	262	0	0	0	0	0	1	
	Tyr	445	0	0	312	561	2293	0	0	0	1856	0	0	0	0	688	0	590	1	0	0	0	11149	0	
	Val	0	6930	0	0	412	0	0	559	1209	0	9252	3804	0	5021	611	0	0	0	0	0	0	15683	0	
	*	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Ter	Thr	Trp	Tyr	Val			
																							alternate		

Scenario 2: Genomic data analysis

As the last step of the genomics pipeline, the distinct mutations are being monitored by reviewing the Delta Lake parquet files within DBFS (i.e. increase of distinct mutations over time).

```
display(spark.read.parquet("dbfs:/wbrandler/dnaseq/demo/data/
variant_count_test_parquet"). \
    withColumnRenamed("count", "distinct mutations")
)
```

▶ (3) Spark Jobs



Scenario 2: Genomic data analysis

Part 1 summary

Using Delta Lake, bioinformaticians and researchers can apply distributed analytics with transactional consistency. These abstractions allow data practitioners to simplify genomics pipelines. Here we have created a genomic sample quality control pipeline that continuously processes data as new samples are processed, without manual intervention.

Whether you are performing ETL or performing sophisticated analytics, your data will flow through your genomics pipeline rapidly and without disruption. Try it yourself today by downloading the [Simplifying Genomics Pipelines at Scale with Delta Lake notebook](#).

Part 2: Accurately building genomics cohorts at scale

Azure Databricks can be used to help solve genomics problems so that it is easier for organizations to perform joint-genotyping, the "N + 1" problem, and the challenge of scaling to population-level cohorts. Let's explore how to use Azure Databricks technology to joint genotyping.

Why do people do large scale sequencing? Most people are familiar with the genetic data produced by 23andMe or AncestryDNA. These tests use genotyping arrays, which read a fixed number of variants in the genome, typically ~1,000,000 well-known variants that occur commonly in the normal human population. With sequencing, we get an unbiased picture of all the variants an individual has, whether they are variants we've seen many times before in healthy humans or variants that we've never seen before that contribute to or protect against diseases. Figure 1 demonstrates the difference between variation data produced by genotype arrays, sequencing, and joint genotyping.

Scenario 2: Genomic data analysis

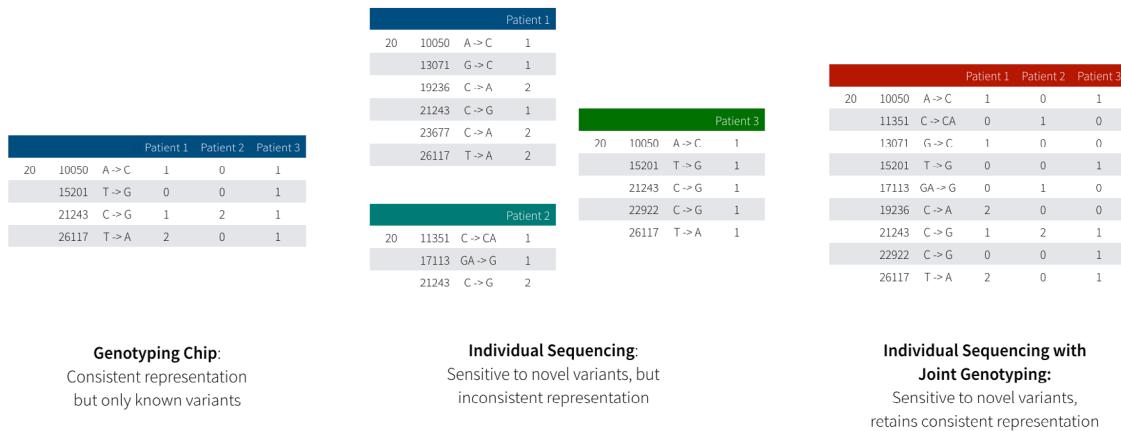


Figure 1: This diagram illustrates the difference between variation data produced by genotype arrays (left) and by sequencing (middle) followed by joint genotyping (right). Genotyping arrays are restricted to "read" a fixed number of known variants, but guarantee a genotype for every sample at every variant. In sequencing, we are able to discover variants that are so rare that they only exist in a single individual, but determining if a novel variant is truly unique in this person or just hard to detect with current technology is a non-trivial problem.

While sequencing provides much higher resolution, it becomes problematic when trying to examine the effect of a genetic variant on many patients. Since a genotyping array measures the same variants across all samples, looking over many individuals is a straightforward proposition: all variants have been measured for all individuals. When working with sequencing data, we have a trickier proposition: if we saw a variant in patient 1, but didn't see that variant in patient 2, what does that tell us? Did patient 2 not have an allele of that variant? Alternatively, when we sequenced patient 2, did an error occur that caused the sequencer to not read the variant we are interested in?

Joint genotyping addresses this problem as follows:

1. Combining evidence from multiple samples enables the rescue of variants that do not meet strict statistical significance needed to be detected accurately in a single sample.
2. As the accuracy of the predictions at each site in the human genome increases, the ability to model sequencing errors and filter spurious variants improves.
3. Joint genotyping provides a common variant representation for all samples that simplifies determining whether a variant in individual X is also present in individual Y.

Scenario 2: Genomic data analysis

Accurately identifying genetic variants at scale with joint genotyping

Joint genotyping works by pooling data from all the individuals in a study when computing the likelihood for each individual's genotype. This provides a uniform representation of how many copies of each variant are present in each individual, a key stepping-stone in looking at the link between a genetic variant and a disease. Computing these new likelihoods also enables the computation of a prior probability distribution for a given variant appearing in a population, which we can use to disambiguate borderline variant calls.

For a more concrete example, table 1 shows the precision and recall statistics for indel (insertions/deletions) and single-nucleotide variants (SNVs) for the sample HG002 called via the GATK variant calling pipeline compared to the [Genome-in-a-Bottle \(GIAB\)](#) high-confidence variant calls in high-confidence regions.

Variant calling accuracy for HG002, processed as a single sample

	Recall	Precision
Indel	96.25%	98.32%
SNV	99.72%	99.40%

By contrast, table 2 shows improvement for indel precision and recall and SNV recall when variants are jointly called across HG002 and two relatives (HG003 and HG004). The halving of this error rate is significant, especially for clinical applications.

Variant calling accuracy for HG002 following joint genotyping with HG003 and HG004

	Recall	Precision
Indel	98.21%	98.98%
SNV	99.78%	99.34%

Scenario 2: Genomic data analysis

Originally, joint genotyping was performed directly from the raw sequencing data for all individuals, but as studies have grown to petabytes in size, this approach has become impractical. Modern approaches start from a genome variant call file (gVCF), a tab-delimited file with all variants seen in a single sample, and information about the quality of the sequencing data at every position where no variant was seen. While the gVCF-based approach touches less data than looking at the sequences, a moderately sized project can still have tens of terabytes of gVCF data. This gVCF-based approach eliminates the need to go back to the raw reads but still requires all $N+1$ gVCFs to be reprocessed when jointly calling variants with a new sample added to the cohort. Our approach uses Delta Lakes to enable incrementally squaring-off the $N+1$ samples in the cohort, while parallelizing regenotyping using Apache Spark™.

The Challenges of Calling Variants at a Population Level

Despite the importance of joint variant calling, bioinformatics teams often defer this step because the existing infrastructure for GATK4 makes these workloads hard to run and even harder to scale. The default implementation of the GATK4's joint genotyping algorithm is single threaded, and scaling this implementation relies on manually parallelizing the joint genotyping kernel using a workflow language and runners like WDL and Cromwell. While GATK4 has support for a Spark-based HaplotypeCaller, it does not support running GenotypeGVCFs parallelized using Spark.

Additionally, for scalability, the GATK4 best practice joint genotyping workflow relies on storing data in GenomicsDB. However, GenomicsDB has limited support for cloud storage systems like Azure Blob Storage, and [studies have demonstrated](#) that the new GenomicsDB workflow is slower than the old CombineGVCFs/GenotypeGVCFs workflow on some large datasets.

Scenario 2: Genomic data analysis

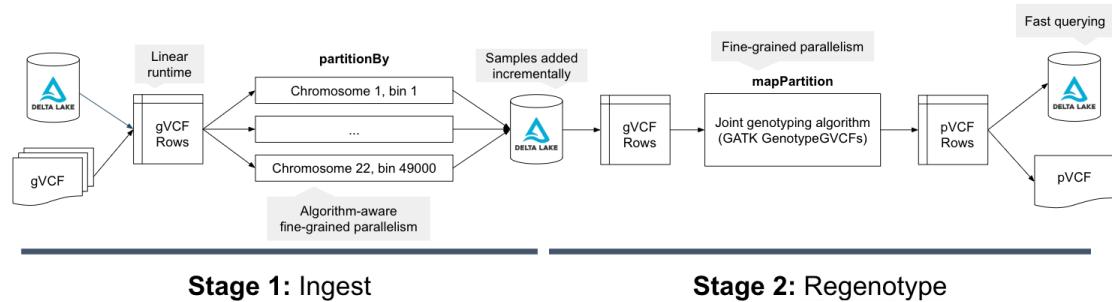


Figure 2: The computational flow of the Azure Databricks joint genotyping pipeline. In stage 1, the gVCF data is ingested into a Delta Lake columnar store in a scheme partitioned by a genomic bin. This Delta Lake table can be incrementally updated as new samples arrive. In stage 2, we then load the variants and reference models from the Delta Lake tables and directly run the core re-genotyping algorithm from the GATK4's GenotypeGVCFs tool. The final squared off genotype matrix is saved to Delta Lake by default, but can also be written out as VCF on some large datasets.

The [Azure Joint Genotyping Pipeline](#) provides a solution for these common needs.

Figure 2 shows the computational architecture of the joint genotyping pipeline. Referred to as an Azure Databricks job, the joint variant calling pipeline is simple to run: the user simply needs to provide their input files and output directory. When the pipeline runs, it starts by appending the input gVCF data to Delta Lake. Delta Lake provides inexpensive incremental updates, which makes it cheap to add an $N+1$ th sample into an existing cohort. When the pipeline runs, it uses Spark SQL to bin the variant calls. The joint variant calling algorithm then runs in parallel over each bin, scaling linearly with the number of variants.

The parallel joint variant calling step is implemented through Spark SQL. Specifically, you bin all the input genotypes/ reference models from the gVCF files into contiguous regions of the reference genome. Within each bin, data is sorted by reference position and sample ID. You then directly invoke the joint genotyping algorithm from the GATK4's GenotypeGVCFs tool over the sorted iterator for the genomic bin. You then save this data out to a Delta table, and optionally as a VCF file.

Scenario 2: Genomic data analysis

Benchmarking

The benchmark for the approach was to use low-coverage WGS data from the [1000 Genomes](#) project for scale testing, and data from the [Genome-in-a-Bottle consortium](#) for accuracy benchmarking. For generating input gVCF files, the DNASEq pipeline was used to aligned and call variants. Figure 3 demonstrates that the approach is efficiently scalable with both dataset and cluster size. With this architecture, it is possible to jointly call variants across the 2,504 sample of whole genome sequencing data from the 1000 Genomes Project in 79 hours on 13 c5.9xlarge machines. Since then, customers have scaled this pipeline across projects with more than 3,000 whole genome samples.

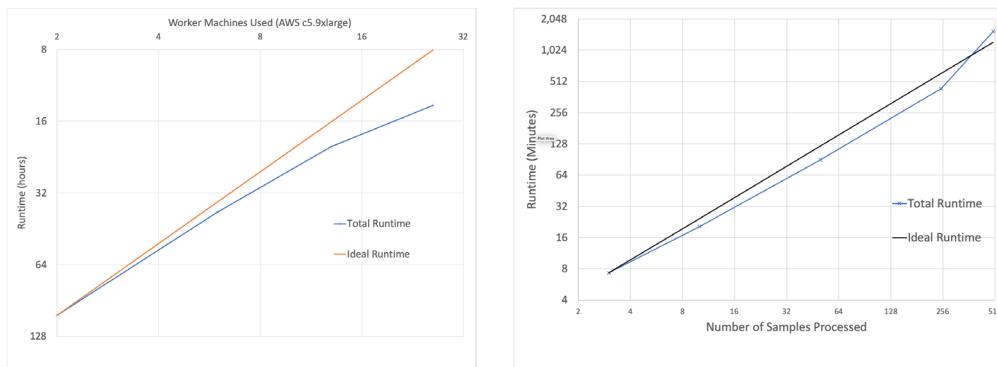


Figure 3: Strong scaling (left) was evaluated by holding the input data constant at 10 samples and increasing the number of executors. Weak scaling (right) was evaluated by holding the cluster size fixed at 13 i3.8xlarge workers and increasing the number of gVCFs processed.

Scenario 2: Genomic data analysis

Running with default settings, the pipeline is highly concordant with the GATK4 joint variant calling pipeline on the HG002, HG003 and HG004 trio. Table 3 describes concordance at the variant and genotype level when comparing the pipeline against the “ground truth” of the GATK4 WDL workflow for joint genotyping. Variant concordance implies that the same variant was called across both tools; a variant called by our joint genotyper only is a false positive, while a variant called only by the GATK GenotypeGVCFs workflow is a false negative. Genotype concordance is computed for all variants called by both tools. A genotype call is treated as a false positive relative to the GATK if the count of called alternate alleles increased in the pipeline, and as a false negative if the count of called alternate alleles decreased.

Concordance statistics comparing open-source GATK4 vs. Databricks for joint genotyping workflows

	Recall	Precision
Variant	99.9985%	99.9982%
Genotype	99.9988%	99.9992%

The discordant calls generally occur at locations in the genome that have a high number of observed alleles in the input gVCF files. At these sites, the GATK discards some alternate alleles to reduce the cost of re-genotyping. However, the alleles eliminated in this process depends on the sequence that variants are read. An effective approach is to list the alternate alleles in the lexicographical order of the sample names prior to pruning. This approach ensures that the output variant calls are consistent given the same sample names and variant sites, regardless of how the computation is parallelized.

Of additional note, the Azure Databricks joint genotyping implementation exposes a configuration option that “squares off” and re-genotypes the input data without adjusting the genotypes by a prior probability which is computed from the data. This feature addresses concerns about the deflation of rare variants caused by the prior probability model used in the GATK4.

Part 2 summary

Using Delta Lake and Spark SQL, it is possible to develop an easy-to-use, fast, and scalable joint variant calling framework on Azure Databricks.

Scenario 3: Real-time display advertising attribution

In digital advertising, one of the most important things to be able to deliver to clients is information about how their advertising spend drove results. The more quickly they receive this information, the better.

To tie conversions or engagements to the impressions served in an advertising campaign, companies must perform attribution. Attribution can be a fairly expensive process, and running attribution against constantly updating datasets is challenging without the right technology. Traditionally, this has not been an easy problem to solve as there because of questions like:

- How do we make sure the data can be written at low latency to a read location without corrupting records?
- How can we continuously append to a large, query-able dataset without exploding costs or loss of performance?
- And where and how should I introduce a join for attribution?



Fortunately, this is easy with Structured Streaming and Delta Lake. In this section,, you can see how to use the DataFrame API to build Structured Streaming applications on top of [Azure EventHubs](#), [Apache Kafka on HDInsight](#), or [Azure Cosmos DB integration](#), and use Delta Lake to query the streams in near-real-time. You will also see how to use the BI tool of your choice to review your attribution data in real- time.

Scenario 3: Real-time display advertising attribution

Define streams

The first thing that needs to be done is establishing the impression and conversion data streams. The impression data stream provides a real- time view of the attributes associated with those customers who were served the digital ad (impression) while the conversion stream denotes customers who have performed an action (e.g. click the ad, purchased an item, etc.) based on that ad.

You can quickly plug into the stream as Databricks supports direct connectivity to Apache Kafka on Azure HDInsight (this is for impressions, repeat this step for conversions) and Azure Event Hubs in the following code snippet.

```
// Read impressions stream
val kinesis = spark.readStream
  .format("azure-event-hubs")
  .option("streamName", kinesisStreamName)
  .option("region", kinesisRegion)
  .option("initialPosition", "latest")
  .option("awsAccessKey", $awsAccessKeyId$)
  .option("awsSecretKey", $awsSecretKey$)
  .load()
```

Next, it's time to create data streams schema as noted in the following snippet.

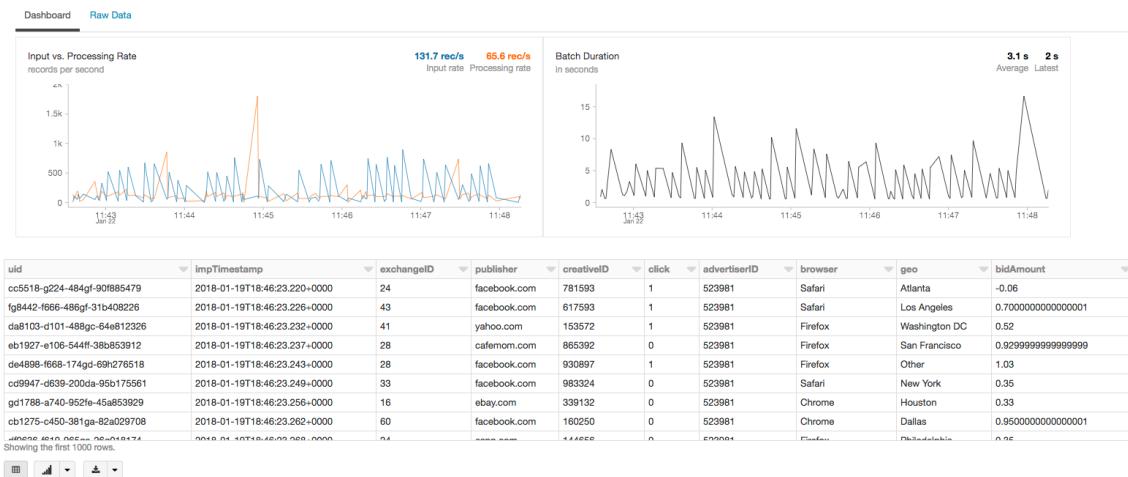
```
// Define impressions stream schema
val schema = StructType(Seq(
    StructField("uid", StringType, true),
    StructField("impTimestamp", TimestampType, true),
    StructField("exchangeID", IntegerType, true),
    StructField("publisher", StringType, true),
    StructField("creativeID", IntegerType, true),
    StructField("click", StringType, true),
    StructField("advertiserID", IntegerType, true),
    StructField("browser", StringType, true),
    StructField("geo", StringType, true),
    StructField("bidAmount", DoubleType, true)
))
```

Scenario 3: Real-time display advertising attribution

The final step is to create the streaming impressions DataFrame. With the Databricks display command, both the data and the input/processing rate can be seen in real-time alongside the data.

```
// Define streaming impressions DataFrame
val imp =
  azure-event-hubs.select(from_json('data.cast("string"), schema)
as "fields").select($"fields.*")

// View impressions real-time data
display(imp)
```



Scenario 3: Real-time display advertising attribution

Sync streams to Delta Lake

The impression (imp) and conversion (conv) streams can be synced directly to Delta Lake, allowing a greater degree of flexibility and scalability for this real-time attribution use-case. You can quickly write these real-time data streams into Parquet format on Blob Storage while enabling users to read from the same directory simultaneously without the overhead of managing consistency, transactionality and performance yourself. In the following code snippet, the raw records are being captured from a single source and written into its own Delta Lake table.

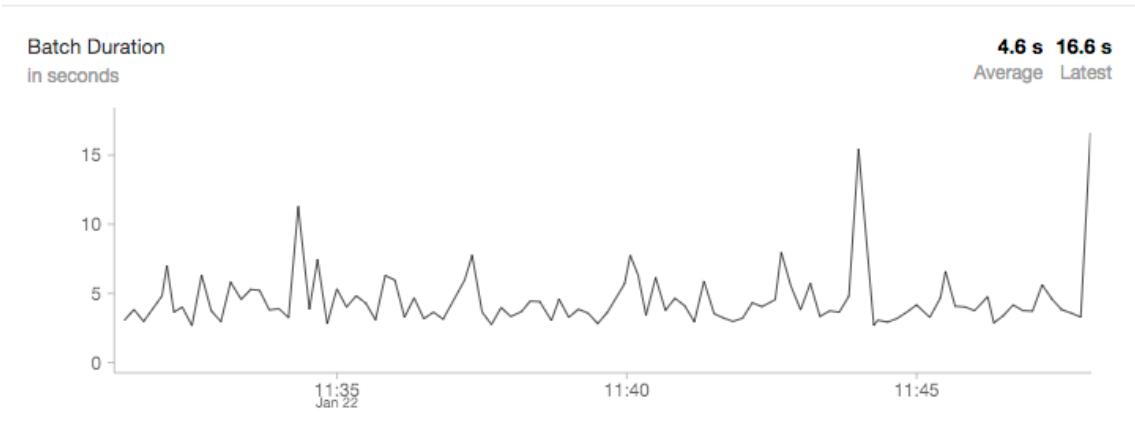
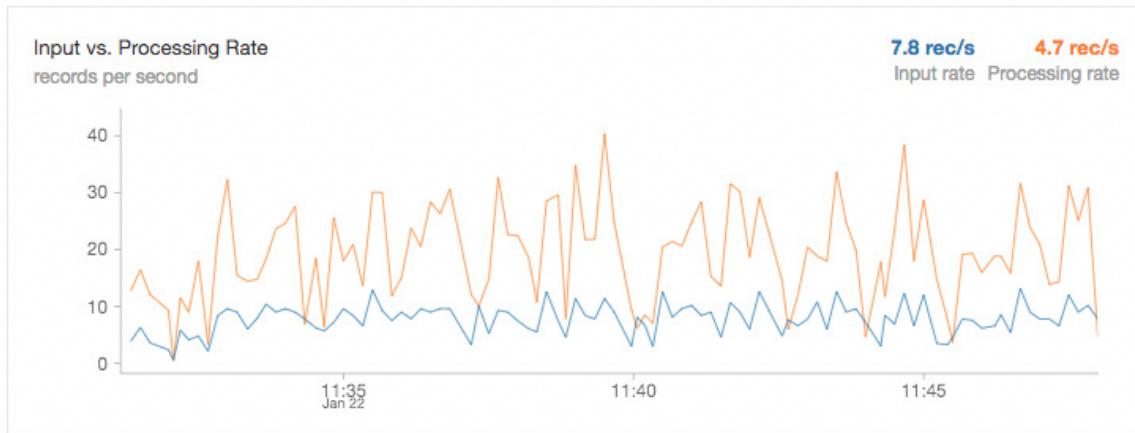
```
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.expressions.Window

// Persist Impression `imp` data to Databricks Delta
imp.withWatermark("impTimestamp", "1 minute")
    .repartition(1)
    .writeStream
    .format("delta")
    .option("path", "/tmp/adtech/impressions")
    .option("checkpointLocation", "/tmp/adtech/impressions-checkpoints")
    .trigger(org.apache.spark.sql.streaming.Trigger.ProcessingTime("10
seconds")).start()
```

Scenario 3: Real-time display advertising attribution

It is important to note that with Delta Lake, you can also:

- Apply additional ETL, analytics, and enrichment steps at this point.
- Write data from different streams or batch process and different sources into the same table.



Scenario 3: Real-time display advertising attribution

Delta Lake views for ad hoc reporting

Now that the impression and conversion Delta Lake tables are done, the next step is to create named views for easily executing the joins in Spark SQL, as well as make this data available to query from your favorite BI tool. First comes creating the Delta Lake Views.

```
%sql  
use adtech;  
create or replace view impressionsDelta  
as select *  
from delta.`/tmp/adtech/impressions`;  
create or replace view conversionsDelta  
as select *  
from delta.`/tmp/adtech/conversions`;
```

Now it's possible to calculate the last touch attribution on the view and then calculate weighted attribution on the view.

Calculate last touch attribution on view

Calculating the real-time window attribution, as noted in the preceding sections, requires joining two different Delta Lake streams of data: impressions and conversions. The code snippet shows the definition of the first Delta Lake-based impressions and conversions and the window specification that will be used by the following `dense_rank()` statement. The window and rank define the attribution logic.

```
# Define needed Impression data  
val imps = spark.sql("select uid as impUid, advertiserID as impAdv,  
* fromsparksummit.imps")  
.drop("advertiserID")  
  
// Define needed Conversions Data  
val convs = spark.sql("select * from sparksummit.convs")  
  
// Define Spark SQL window ordered by Impression  
// Timestamp partitioned by Impression Uid and Impression Advertisor  
val windowSpec = Window.partitionBy("impUid","impAdv")  
.orderBy(desc("impTimestamp"))
```

Scenario 3: Real-time display advertising attribution

```
// Calculate real-time attribution by joining
// impression.impUid == conversion.uid to ensure
// impression time happened before conversion time
// filtering via dense_rank
val windowedAttribution = convs.join(imps,
  imps.col("impUid") === convs.col("uid") &&
  imps.col("impTimestamp") < convs.col("convTimestamp")
  && imps.col("impAdv") === convs.col("advertiserID"))
    .withColumn("dense_rank", dense_rank()
      .over(windowSpec))
    .filter($"dense_rank" === 1)

// Create global temporary view
windowedAttribution.createGlobalTempView("realTimeAttribution")
```

To calculate the real-time window attribution, you join the impression and conversion data but filter for the most recent impression user id (as defined by `impUid`). A global temporary view (`.createGlobalTempView`) enables this real-time view to be accessible by downstream systems. For example, you can view your real-time data using Spark SQL in the following code snippet.

```
%sql
select * from global_temp.realTimeAttribution
```

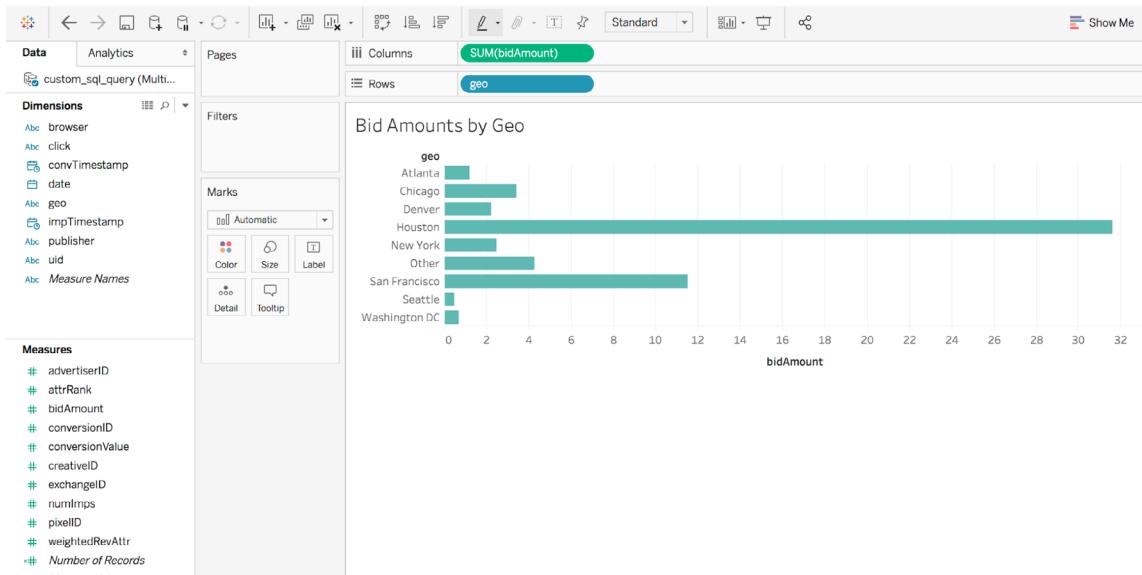
(7) Spark Jobs																	
uid	convTimestamp	conversionID	advertiserID	pixelID	conversionValue	impUid	impAdv	uid	impTimestamp	exchangeID	publisher	creativeID	click	advertiserID	browser		
gc1932-e367-422af-78a153015	2018-01-19T18:55:13.645+0000	89169	523981	103952	57.54	gc1932-e367-422af-78a153015	523981	gc1932-e367-422af-78a153015	2018-01-19T18:56:59.297+0000	24	facebook.com	442957	1	523981	Internet Explorer		
fe3576-e384-785ec-92g730761	2018-01-19T18:55:08.565+0000	50320	523981	103952	22.080000000000002	fe3576-e384-785ec-92g750761	523981	fe3576-e384-785ec-92g750761	2018-01-19T18:55:14.349+0000	35	ebay.com	481457	0	523981	Chrome		
fe3576-e384-785ec-	2018-01-19T18:55:03.565+0000	93331	523981	103952	37.269999999999996	fe3576-e384-785ec-	523981	fe3576-e384-785ec-	2018-01-19T18:55:14.349+0000	35	ebay.com	481457	0	523981	Chrome		

Scenario 3: Real-time display advertising attribution

You can also plug in your favorite BI tool such as Power BI to perform ad-hoc analysis of your data.

The screenshot shows a BI tool interface with a sidebar on the left containing 'Connections' (demo.cloud.databricks.com, Spark SQL) and 'Schema' (default). The main area displays a table titled 'custom_sql_query (Multiple Connections)' with the following columns:

#	custom_sql_query.exchangeID	Abc custom_sql_query.publisher	#	custom_sql_query.creativeID	Abc custom_sql_qu... click	Abc custom_sql_query.uid	Abc custom_sql_query.browser	Abc custom_sql_query.geo	#	custom_sql_query.bidAmount	custom_sql_query.impTimestamp
19	espn.com	953,766	0	fg4597-g002-054ge...	Safari	Seattle	0.33000	6/7/2018 11:09			
19	espn.com	953,766	0	fg4597-g002-054ge...	Safari	Seattle	0.33000	6/7/2018 9:35*			
40	espn.com	335,653	0	fg4597-g002-054ge...	Safari	Seattle	0.75000	6/7/2018 9:22*			
42	vice.com	804,280	0	fg4597-g002-054ge...	Safari	Seattle	0.86000	6/7/2018 9:03!			
60	yahoo.com	251,928	0	fg4597-g002-054ge...	Safari	Seattle	0.49000	5/19/2018 4:04			
39	espn.com	543,581	0	fg4597-g002-054ge...	Safari	Seattle	0.46000	5/19/2018 2:26			
21	ebay.com	587,442	1	fg4597-g002-054ge...	Safari	Seattle	0.32000	5/8/2018 3:46:			
26	facebook.com	291,589	1	fd0018-f693-135ha...	Firefox	New York	0.35000	5/18/2018 2:31			
39	facebook.com	756,743	1	fd0018-f693-135ha...	Firefox	New York	0.43000	5/16/2018 12:4			



Scenario 3: Real-time display advertising attribution

Calculate Weighted Attribution on View

The preceding case demonstrated a naive model — isolating all of a user's impressions prior to conversion, selecting the most recent, and attributing only the most recent impression prior to conversion. A more sophisticated model might apply attribution windows or weight the impressions by time such as the code snippet below.

```
// Define attribution window for impressions
val attrWindow = Window \
    .partitionBy("uid")
    .orderBy($"impTimestamp".desc)
val attrRank = dense_rank().over(attrWindow)

// Define ranked window by taking attribution window
// and partition by conversionID
val rankedWindow = Window.partitionBy("conversionID")
val numAttrImps = \
    max(col("attrRank")).over(rankedWindow)

// Reference impression table
val imps = spark \
    .sql("select * from sparksummit.imps") \
    .withColumn("date", $"impTimestamp" \
        .cast(DateType)) \
    .drop("advertiserID")

// Reference conversion table
val convs = spark \
    .sql("select * from sparksummit.convs") \
    .withColumnRenamed("uid", "cuid")

// Join impression and conversions
val joined = imps.join(convs,
    imps.col("uid") === convs.col("cuid") &&
    imps.col("impTimestamp") <
    convs.col("convTimestamp")).drop("cuid")

// Calculate weighted attribution
val attributed = joined \
    .withColumn("attrRank",attrRank) \
    .withColumn("numImps",numAttrImps) \
    .withColumn("weightedRevAttr", $"conversionValue"/$"numImps")

// Create attributed temporary view
attributed.createOrReplaceTempView("attributed")
```

Scenario 3: Real-time display advertising attribution

```
// Display data  
display(spark.sql("select sum(weightedRevAttr),  
advertiserID from attributed group by advertiserID"))
```

sum(weightedRevAttr)	advertiserID
7846749.570946295	702394
4431214.068630882	489251
2200343.4445629856	295381
10097961.92226806	523981

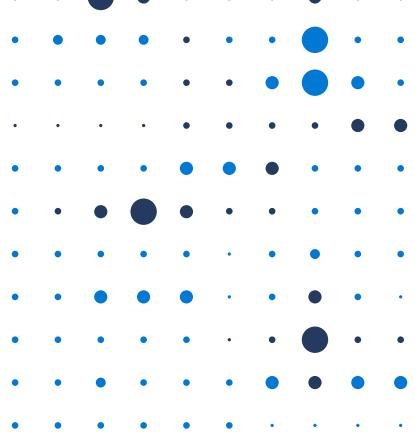


Summary

This section reviewed how Delta Lake provides a simplified solution for a real-time attribution pipeline. The advantages of using Delta Lake to sync and save your data streams include the ability to:

- Save and persist your real-time streaming data like a data warehouse because Delta Lake maintains a transaction log that efficiently tracks changes to the table.
- Run your queries and perform your calculations in real-time and completing in seconds.
- Multiple writers can simultaneously modify a dataset and still see consistent views.
- Writers can modify a dataset without interfering with jobs reading the dataset.

Together with your streaming framework and Azure Databricks, you can quickly build and use your real-time attribution pipeline with Delta Lake to solve your complex display advertising problems in real-time.



Scenario 4: Mobile Gaming Data Event Processing

The world of mobile gaming is fast paced and requires the ability to scale quickly. With millions of gamers around the world generating millions of events per second, you will need to calculate key metrics (score adjustments, in-game purchases, in-game actions, etc.) in real-time. Just as important, a popular game launch or feature will increase event traffic by orders of magnitude and you will need infrastructure to handle this rapid scale.

With complexities of low-latency insights and rapidly scalable infrastructure, building data pipelines for high volume streaming use cases like mobile game analytics can be complex and confusing. Developers who are tasked with this endeavor will encounter a number of architectural questions:

- What set of technology should they consider that will reduce their learning curve and that integrate well?
- How scalable will the architecture be when built?
- How will different personas in an organization collaborate?

Ultimately, they will need to build an end-to-end data pipeline comprised of these three functional components:

- Data ingestion/streaming
- Data transformation (ETL)
- Data analytics and visualization

In this section, we will explore how to:

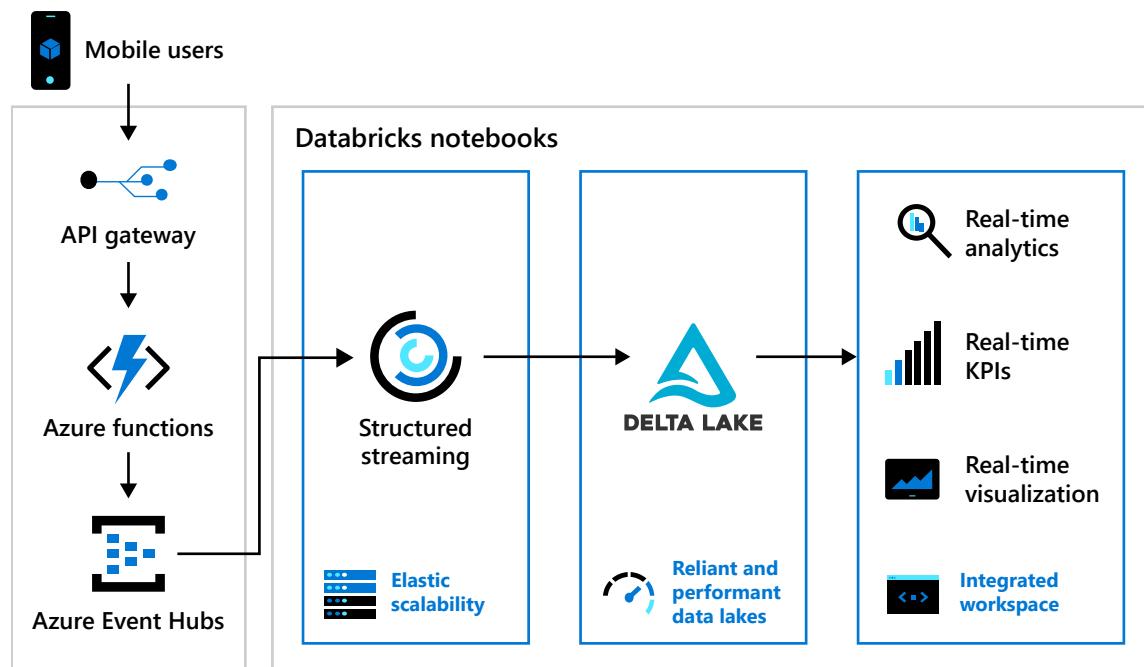
- Build a mobile gaming data pipeline using Azure Event Hubs _____
- Build a stream ingestion service using Spark Structured Streaming
- Use Delta Lake as a sink for our streaming operations
- Explore how analytics can be performed directly on this table, minimizing data latency
- Illustrate how Delta Lake solves traditional issues with streaming data

Scenario 4: Mobile Gaming Data Event Processing

High Level Infrastructure Components

Building mobile gaming data pipelines is complicated by the fact that you need rapidly scalable infrastructure to handle millions of events by millions of users and gain actionable insights in real-time. Thankfully, Azure Event Hubs throughput units can be dynamically re-provisioned to handle increased loads, and Azure Databricks automatically scales out your cluster to handle the increase in data.

Event Hubs throughput units are provisioned by throughput, so you can provision as many units as necessary to handle your expected data throughput. For more information about Event Hubs throughput units, check out [this documentation](#). Random PartitionKeys are important for even distribution if you have more than one shard.



Note: The above architecture can also be achieved with open source software.

Scenario 4: Mobile Gaming Data Event Processing

Ingesting from Event Hubs using structured streaming

Ingesting data from an Event Hub stream is straightforward. In a production environment, you will want to setup the appropriate IAM role policies to make sure your cluster has access to your Event Hub stream. The minimum permissions look like this:

```
{  
  "Version": "2012-10-17",  
  "Statement": [  
    {  
      "Effect": "Allow",  
      "Action": [  
        "azure-event-hubs:DescribeStream",  
        "azure-event-hubs:GetRecords",  
        "azure-event-hubs:GetShardIterator"  
      ],  
      "Resource": "ARN_FOR_YOUR_STREAM"  
    }  
  ] }
```

Real-time data pipelines using Delta Lake

Now that you have our streaming DataFrame defined, you can do some simple transformations. Event data is usually time-series based, so it's best to partition on something like an event date. This particular incoming stream does not have an event date parameter, however, so you'll make your own by transforming the `eventTime` column. And, there will also be a check to make sure the `eventTime` is not null:

```
base_path = '/path/to/mobile_events_stream/'  
eventsStream = gamingEventDF \  
    .filter(gamingEventDF.eventTime.isNotNull()) \  
    .withColumn("eventDate", \  
               to_date(gamingEventDF.eventTime))  
    .writeStream \  
    .partitionBy('eventDate') \  
    .format('delta') \  
    .option('checkpointLocation', base_path + \  
           '/_checkpoint') \  
    .start(base_path)
```

Scenario 4: Mobile Gaming Data Event Processing

You can also take this opportunity to define your table location.

```
CREATE TABLE
IF NOT EXISTS mobile_events_delta_raw
USING DELTA
location '/path/to/mobile_events_stream/';
```

NOTE: DDL is part of the roadmap and will be available in Delta Lake in the near future.

Real-Time Analytics, KPIs, and Visualization

Now that data streaming is live in your Delta Lake table, you can review some KPIs. Traditionally, companies would only look at these once a day, but with Structured Streaming and Delta Lake, you have the capability to visualize these in real time all within your Azure Databricks notebooks.

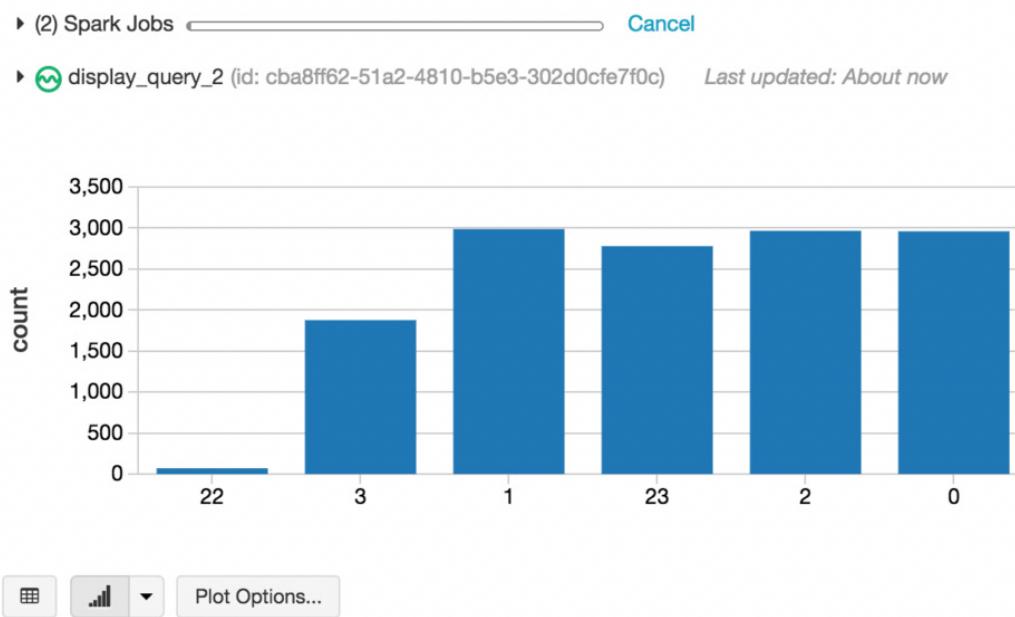
You can start with a simple one. How many events have I seen in the last hour?

```
countsDF = gamingEventDF \
    .withWatermark("eventTime", "180 minutes") \
    .groupBy(window("eventTime", "60 minute")) \
    .count()
countsQuery = countsDF \
    .writeStream \
    .format('memory') \
    .queryName('incoming_events_counts') \
    .start()
```

Scenario 4: Mobile Gaming Data Event Processing

You can then visualize this in your notebook as say, a bar graph:

```
display(countsDF.withColumn('hour', hour(col('window.start'))))
```



Maybe you can make things a little more interesting. How much money have I made in the last hour? Let's examine bookings.

Understanding bookings per hour is an important metric because it can be indicative of how the application/production systems are doing. If there were a sudden drop in bookings right after a new game patch was deployed, for example, you immediately know something is wrong.

Scenario 4: Mobile Gaming Data Event Processing

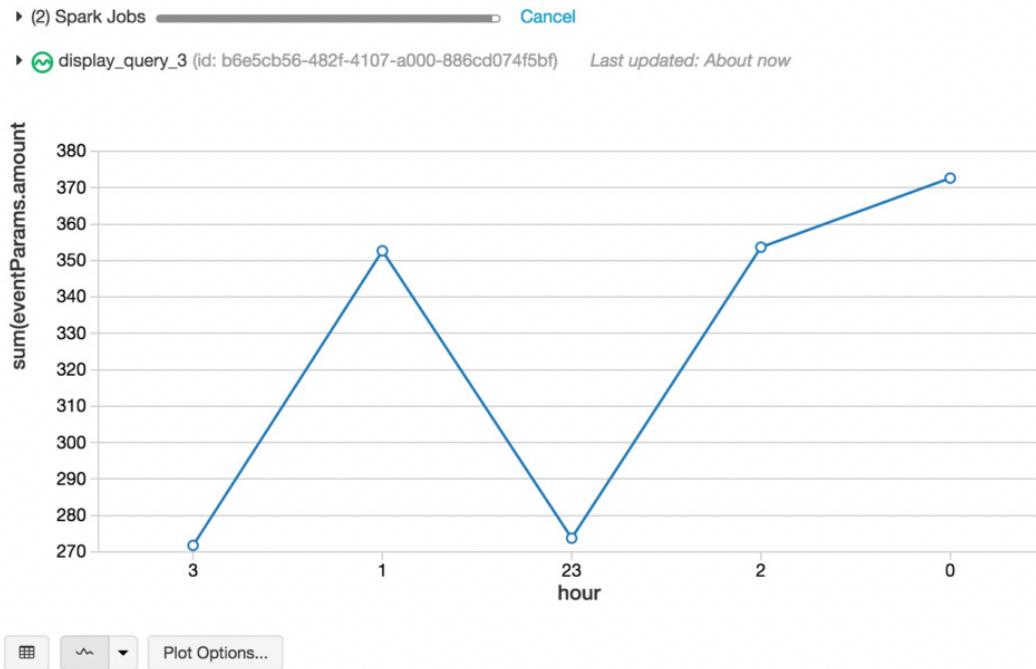
You can take the same DataFrame, but filter on all purchaseEvents, grouping by a window of 60 minutes.

```
bookingsDF = gamingEventDF \
    .withWatermark("eventTime", "180 minutes") \
    .filter(gamingEventDF.eventName == 'purchaseEvent') \
    .groupBy(window("eventTime", "60 minute")) \
    .sum("eventParams.amount")

bookingsQuery = bookingsDF \
    .writeStream \
    .format('memory') \
    .queryName('incoming_events_bookings') \
    .start()
```

A line graph is used to visualize this one:

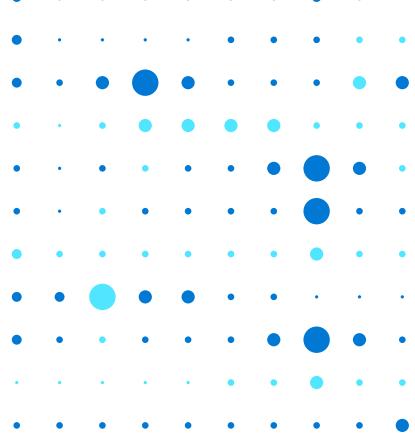
```
display(bookingsDF.withColumn('hour', hour(col('window.start'))))
```



Scenario 4: Mobile Gaming Data Event Processing

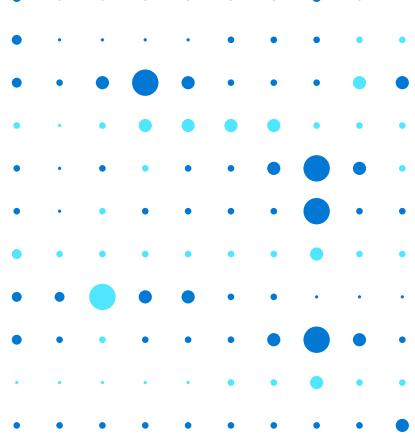
For the SQL enthusiasts, you can query the Delta Lake table directly. Here is a simple query to show the current daily active users (DAU). You know you're actually looking at device id because your sample set doesn't contain a user id, so for the sake of example, you can assume that there is a 1-1 mapping between users and devices (although, in the real world, this is not always the case).

```
select count (distinct eventParams.device_id)
as DAU
from mobile_events_delta_raw
where to_date(eventTime) = current_date;
```



Summary

This eBook has demonstrated how you can build a data pipeline's three functional components: data ingestion/streaming, data transformation, and data analytics and visualization. We've illustrated different ways that you can extrapolate key performance metrics from this real-time streaming data, as well as solve issues that are traditionally associated with streaming. The combination of Spark Structured Streaming and Delta Lake reduces the overall end-to-end latency and availability of data, enabling data engineering, data analytics, and data science teams to respond quickly to events like a sudden drop in bookings, or an increased error-message events, that have direct impact on revenue. Additionally, by removing the data engineering complexities commonly associated with such pipelines, this enables data engineering teams to focus on higher-value projects.

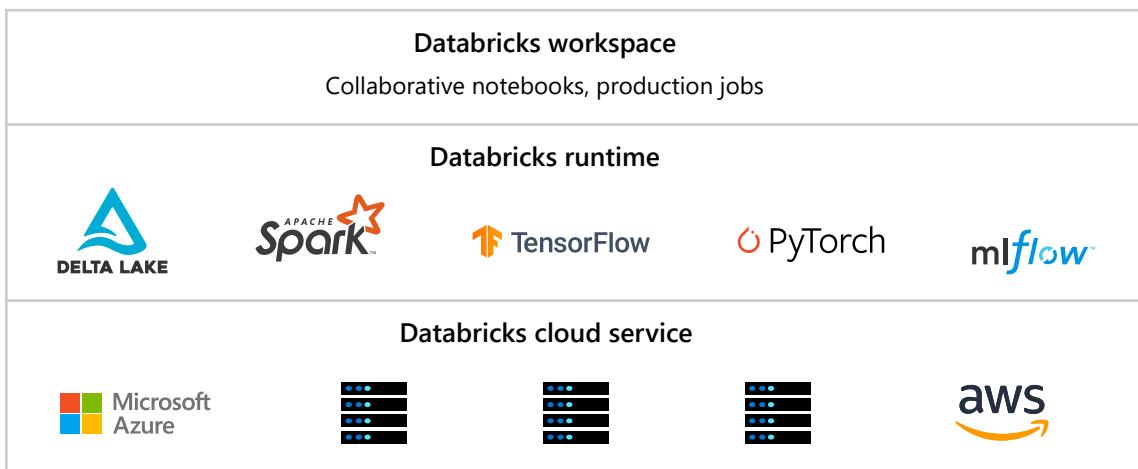


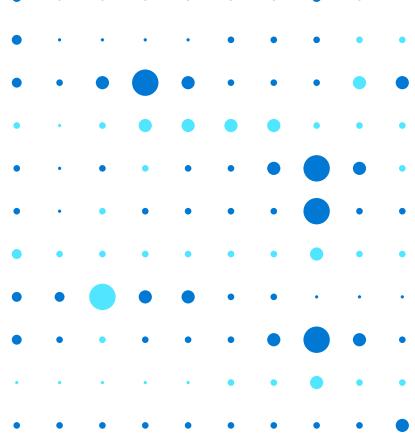
About Azure Databricks

Azure Databricks enables data professionals to manage the entire data analytics lifecycle from data ingestion and preparation to model development, serving and analytics.

Innovations such as Delta Lake, for data lake reliability, and MLflow, for managing the complete machine learning lifecycle, are key enablers of the platform.

Databricks Unified Analytics Platform





Resources

Invent with purpose



Microsoft Azure