



# Beyond SQL

Building hybrid analytics & insights systems  
on top of Apache Spark

---

**Scott Haines**, Principal Software Engineer,  
Voice Insights, **Twilio**



@newfront  
@twilio



@newfrontcreative





# SIGNAL: CUSTOMER AND DEVELOPER CONFERENCE

Explore the intersection of technology, innovation, and communications.



 @newfront @twilio

10% OFF: ENGINEERING10

AUGUST 6 & 7  
MOSCONE WEST, SF CA

2 days of Twilio learning, talks, and networking.



# Laying out the concepts

## SOME BACKGROUND



@newfront @twilio

SQL / NoSQL / NewSQL

Modern Data Architecture

Zero to Hero: SparkSQL

Beyond SQL



# Definitions

---

- **ACID** (Atomicity, Consistency, Isolation, Durability)
- **Atomicity**: Transactions are guaranteed. Eg. All or Nothing
- **Consistency**: Data operations are consistent, rule based fashion
- **Isolation**: Separation of concerns affects visibility of data
- **Durability**: Data state safety. Fault Tolerance



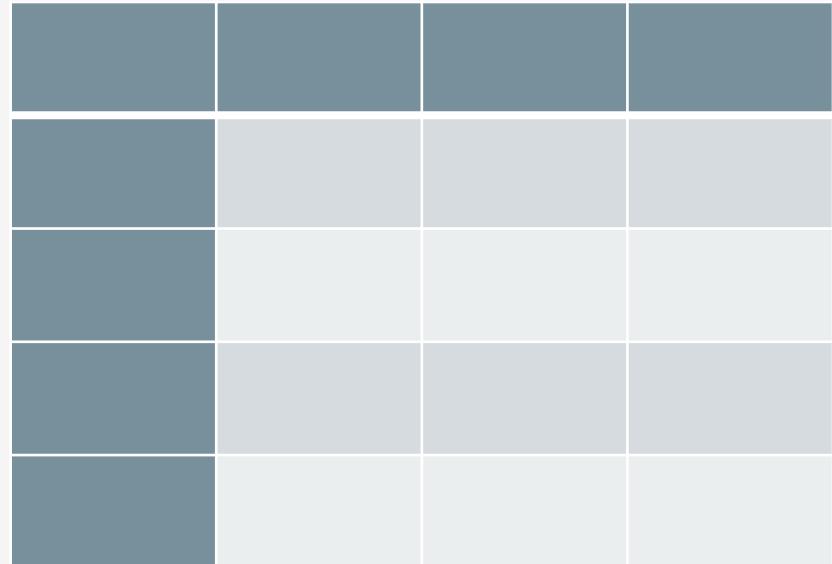
@newfront @twilio



# OG SQL

---

- **Flavors:** MySQL, SQL Server, Oracle, etc
- **Architecture:** Master, Slave
- **Scalability:** Share Everything
- **OLTP:** Online Transaction Processing



@newfront @twilio



# OG SQL

---

- **Structured:** Introduced structured table data via schemas
- **Relational:** Star Pattern: (foreign keys (fk), has and belongs to many, has relationships)
- **Projectable:** `select fieldA, fieldB from table`
- **Filterable:** `select \* from table where conditionA and conditionB and conditionN`




@newfront @twilio



# OLAP SQL

online analytical processing

---

- **Flavors:** Apache Druid, Apache Pinot\*, Apache Kylin, and more
- **Aggregation:** Pre-Aggregation (read optimized), Just-in-time Aggregation (disk optimized, pay in memory)
- **Why:** Allows analytics and insight magic backed by SQL or SQL-like stores



Dashboard magic delivered



@newfront @twilio

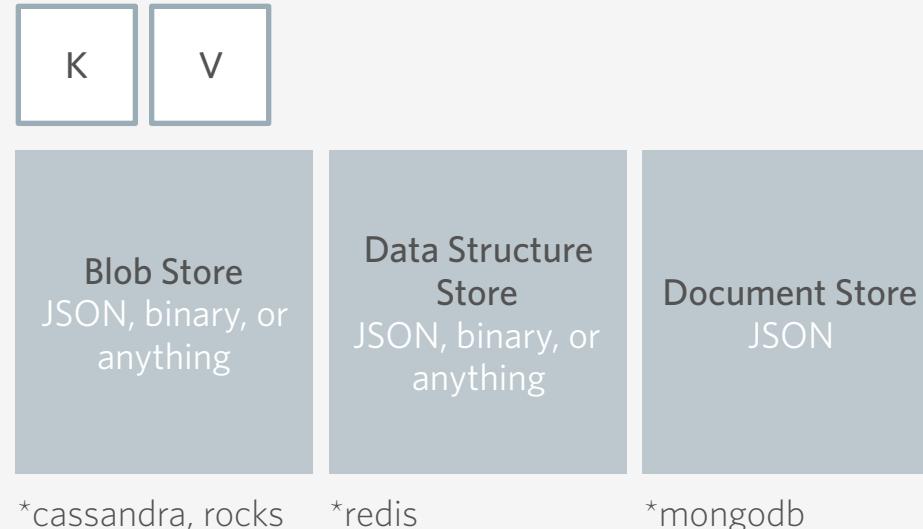


# NoSQL

non-relational

---

- **Flavors:** Redis, Cassandra\*, RocksDB, MongoDB
- **Architectures:** High Availability (HA), Mixed Master (MaMsSS, MaMaSS), Ring/Quorum, Share-Nothing, Eventually Consistent
- **Scalability:** Horizontal via Cluster Sharding



@newfront @twilio

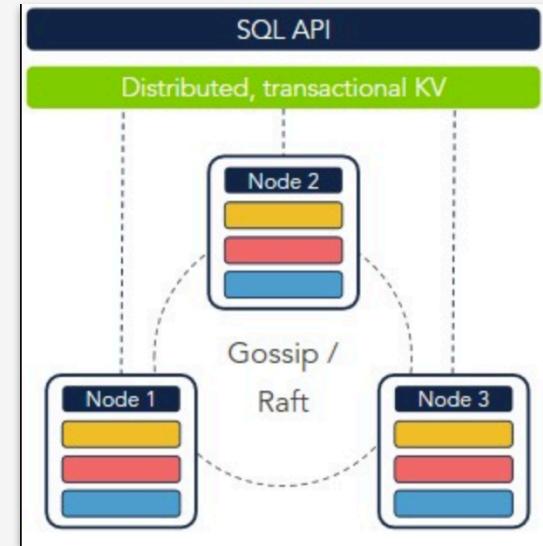


# NewSQL

cloud optimized

---

- **Flavors:** MemSQL, CockroachDB, Google Spanner
- **Architectures:** HA, ACID\*, Distributed, Share-Nothing, Cloud Optimized, All Nodes can serve requests or relay requests.
- **Scalability:** Horizontal, Globally Distributed, Low Latency
- **Supports:** Distributed SQL Queries. Low Latency



Nodes use message passing to resolve distributed queries with low latency



@newfront @twilio

# Musings regarding Data Stores

---

- **Data Warehouses:** OG Shared Data Store
- **Data Lakes:** Bottomless Pits for Good or for Bad



@newfront @twilio



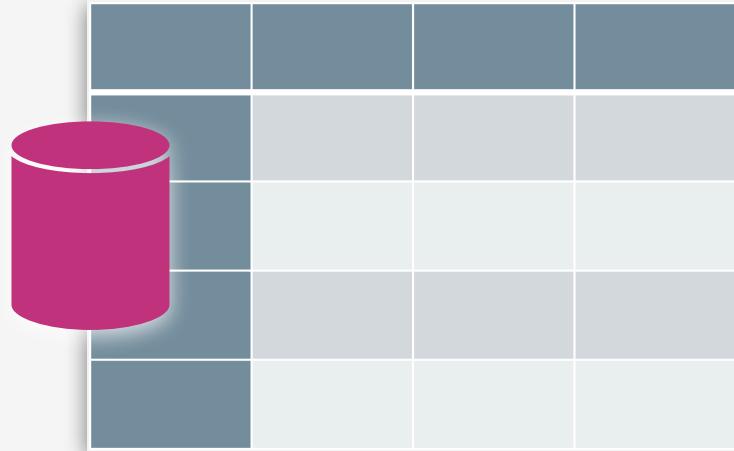
Kimball Group coined 'Warehouse'

# The Data Warehouse

## Pattern One

---

- **How:** Shared Monolithic SQL Database
- **Why:** Organizations need shared access to be able to solve common Business Analytics needs
- **Pros:** Single Shared Database. Uses Grant Permissions to restrict access to Tables
- **Cons:** Single Shared Database. Scalability...



One Silo to Rule Them all.  
One Silo to bring them together...  
Something something Bind them...



@newfront @twilio



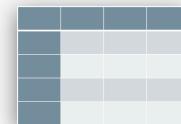
# The Data Lake

## Pattern Two

- **How:** Reliable File Store (s3, hdfs) and connector to hive, looker, etc
- **Why:** Organizations need to be able to store tons of data and analyze in big data fashion.
- **Pros:** Scales out almost indefinitely.
- **Cons:** Schema / San-Schema / Maybe Schema

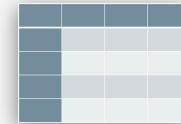


Scales Out. Distributed Storage



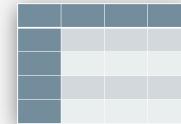
```
{  
  'type': ['json', 'csv']  
}
```

Semi-Structured



```
{  
  'type': ['avro', 'protobuf']  
}
```

Structured  
Binary



```
{  
  'type': ['parquet', 'orc']  
}
```

Columnar



@newfront @twilio

# Musings about Structured Data

---

"Don't drown in your Data Lake"

- JSON <= runtime headache
- **Data Engineering Dos:** Data Contracts and Defined and Versioned Structs.



@newfront @twilio



# Structured Data

---

- JSON has Structure.
- But JSON isn't **strictly** Structured Data.



```
{  
  "type": "employee",  
  "nickname": "scaines",  
  "attributes": [  
    {  
      "height": 6,  
      "unit": "ft"  
    }  
  ]  
}
```



# Structured Data

---

- JSON has **poor runtime guarantees** due to its flexible nature. **Optimize** for compile time guarantees.
- **Debugging** corrupt data in a large distributed system ruins hopes and dreams.

```
{  
  "type": "employee",  
  "nickname": null,  
  "attributes": "bae345"  
}
```

100%

Of people saddened by  
bad data\*



@newfront @twilio

# Structured Data

Versioned & Compiled

---

- **Protocol Buffers** or Avro (arrow, etc) give you types and compile time guarantees
- **Versioning your Data API** now just means sticking to a version of your schema.
- Rely on **release** for versioning.
- Perfect for **Pipeline interoperability**

```
message Employee {  
    uint64 created = 1;  
    string uuid     = 2;  
    string nickname = 3 [default="you"];  
    repeated Attribute attributes = 4;  
}
```

100%

Of people like when  
things work between  
releases!



@newfront @twilio

# Structured Data

columnar

---

- **Parquet is awesome** format for data at rest
- Columnar Data with **Column Stats in partition header\***
- Plays natively with Apache Spark, Presto, Impala and Hive
- \*Schema changes across partitions can be tough

```
spark.read  
.parquet("hdfs://.../year=2019")
```

Spark Encoders Joke ->

100%

Also like not having to define Encoders



@newfront @twilio

# Owning your Data

---

- **What** does your data need to do for you?
- **How** is storing the data going to help your company or organization?
- Does your data meet **GDPR** guidelines?



@newfront @twilio



## 5 min break & QA then LAYERING CONCEPTS

SQL / NoSQL / NewSQL

Modern Data Architecture

Zero to Hero: SparkSQL

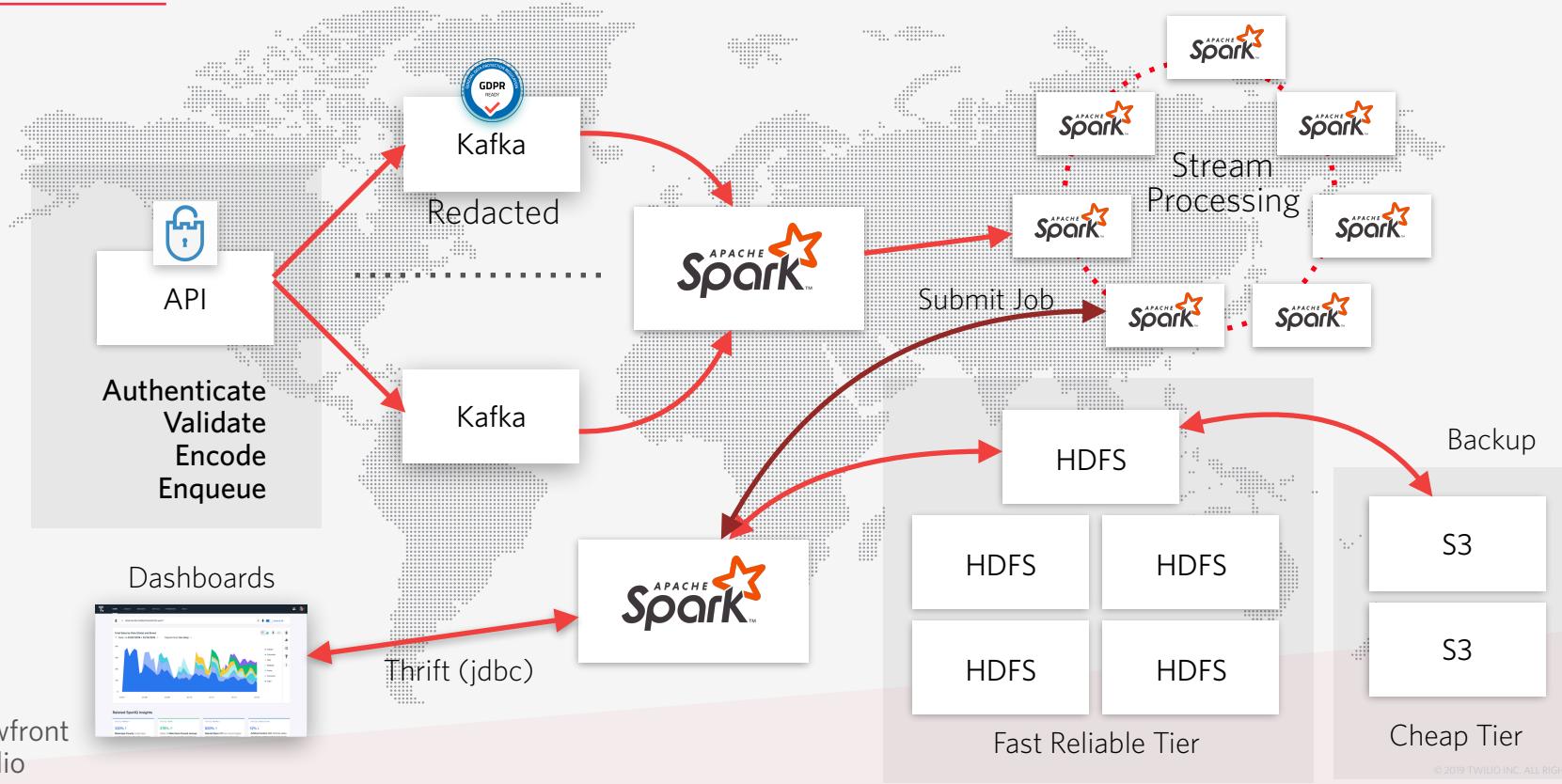
Beyond SQL



@newfront  
@twilio



## MODERN\* DATA ARCHITECTURE STREAMING MICROSERVICE ARCH



@newfront  
@twilio

# Zero to Hero: SparkSQL

---

"Let's look at Supporting OG-style SQL"

- Load all kinds of data
- Transform loaded data
- Project transformed data
- Filtering and more...



@newfront @twilio



# Zero to Hero

## SparkSQL

---

- **The Spark Session:** Literally the GodFather of your application
- Controls and Delegates work
- Jobs, Stages and Tasks report back in and if they don't then we kill em off

```
▶  object ChroniclerApp extends App {  
    ...  
    val log: Logger = LoggerFactory.getLogger(classOf[ChroniclerApp])  
    log.info(s"$args")  
  
    val configPath = if (args.length > 0) args(0) else "conf/application.conf"  
    val config = AppConfig.parse(configPath)  
  
    val sparkConf = new SparkConf()  
        .setAppName(config.sparkAppConfig.appName)  
        .setJars(SparkContext.jarOfClass(classOf[ChroniclerApp]))  
        .setAll(config.sparkAppConfig.core)  
  
    log.info(s"sparkConfig: ${sparkConf.toDebugString}")  
  
    val session = SparkSession.builder  
        .config(sparkConf)  
        .enableHiveSupport()  
        .getOrCreate()  
  
    session.sparkContext.addSparkListener(SparkAppMetricsListener)  
  
    Chronicler(config, session).monitoredRun()  
}
```

# Zero to Hero

## SparkSQL

---

- **Loading** data for batch or streaming can be simple.

```
spark  
  .read  
  .option("inferSchema", "true")  
  .json(path)
```

DataFrame Loaded

Row1	ColA	ColB	ColC
Row2	ColA	ColB	ColC

# Zero to Hero

## SparkSQL

---

- **Creating Projections** on your data is a breeze (once loaded)
- Creating **temporary views** can **simplify** using SparkSQL

```
spark
.read
.option("inferSchema", "true")
.json(path)
.createOrReplaceTempView("table")

val result = spark.sql(
  "select fieldA, fieldB from table
   where condition>90"
)
```

# Zero to Hero

## SparkSQL

---

- **Loading** data for streaming is just a we bit harder.
- **Must Provide** a Schema (StructType)
- Many **options to throttle** ingest on streaming data

```
def jsonSchema: StructType = {
    val updateToString = Set("start_time", "end_time")

    val baseSchema = summaryEncoder.schema
        .add(StructField("@timestamp", StringType, nullable = true))

    val fields: Array[StructField] = baseSchema.map { field =>
        if (updateToString.contains(field.name)) {
            StructField(field.name, StringType, nullable = true)
        } else field
    }.toArray

    baseSchema.copy(fields)
}

lazy val fileRecoveryStream: DataStreamReader = fileStreamReader()
def fileStreamReader(): DataStreamReader = spark.readStream
    .schema(jsonSchema)
    .option("maxFilesPerTrigger", recoveryMaxFilesPerTrigger)

override def run(): Unit = {
    assert(recoveryInputDir.nonEmpty, "spark.chronicler.recoveryInputDir must be set")
    assert(recoveryOutputDir.nonEmpty, "spark.chronicler.recoveryOutputDir must be set")
    assert(recoveryCheckpointPath.nonEmpty, "spark.chronicler.recoveryCheckpointPath must be set")

    val streamingQuery = RecoveryProcessor(config)
        .process(fileRecoveryStream.json(recoveryInputDir))(spark)
        .writeStream
```

# Zero to Hero

## SparkSQL

---

- **Loading** common data sources comes natively.
- **Data Sources:** Kafka, Flume, FileSystem, TCP
- **Formats:** CSV, JSON, AVRO, PARQUET
- **Bring your Own Format:** Extend ExpressionEncoder

```
val reader = spark.readStream.format("kafka")
```

+

```
override def run(): Unit = {
    // creates a streaming query that will aggregate events from kafka
    // and produce statistical histogram aggregations to be consumed downstream
    val eventProcessor = EventAggregation(config)
    eventProcessor.process(readKafkaStream())(spark)
        .writeStream
        .queryName("prediction.stream")
        .outputMode(eventProcessor.outputMode)
        .trigger(Trigger.ProcessingTime(Duration(config.triggerInterval)))
        .format("kafka")
        .option("topic", "odsc.predictions")
        .options(Map(
            "checkpointLocation" -> config.checkpointPath
        ))
        .start()
}
```

# Zero to Hero

## SparkSQL

---

- **ETL** aka Extract Transform Load
- **Create simplified projections** from complex nested data (parquet supports nesting, so does JSON)
- **Sanitize Data** - aka clean for GDPR & Data Governance

```
val withMetrics = flattenStruct(df, "metrics")
    .drop("metrics").drop("score_data").drop("dimension_hash").drop("iso_time")
    .withColumn("call_sid_enc", base64encode(col("call_sid")))
    .withColumn("account_sid_enc", base64encode(col("account_sid")))
    .drop("call_sid").drop("account_sid")
    .withColumnRenamed("call_sid_enc", "call_sid")
    .withColumnRenamed("account_sid_enc", "account_sid")

val withDimensions = flattenStruct(withMetrics, "dimensions").drop("dimensions")
    .withColumn("dimensions_carrier_sid", base64encode(col("dimensions_carrier_sid")))
    .drop("dimensions_provider_sid").drop("dimensions_account_sid")

expandArray(withDimensions, col("insight_tags")).drop("insight_tags")
    .withColumn("ts", functions.from_utc_timestamp(col("iso_time"), "UTC"))
    .withColumn("year", functions.year(col("ts")))
    .withColumn("month", functions.month(col("ts")))
    .withColumn("day", functions.dayofmonth(col("ts")))
    .withColumn("hour", functions.hour(col("ts")))
    .withColumn("dayofyear", functions.dayofyear(col("ts")))
```

# Zero to Hero

## SparkSQL

---

- Parquet + File System Push Down
- Lazy fetches data to resolve a given query.

```
%spark  
val scoresData = spark.read  
  .option("basePath", "hdfs://localhost:9000/user/ds/odsc/east/realish/callscores")  
  .parquet("hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores")
```

```
scoresData: org.apache.spark.sql.DataFrame = [timestamp: bigint, score: float]
```

Took 7 sec. Last updated by anonymous at June 26 2019, 3:54:04 PM.

```
%spark  
scoresData.inputFiles
```

```
res3: Array[String] = Array(hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores/year=2019/month=1/day=25/part-00000-0dbb1b79-64c4-456e-8650-65752b53844f.c000.snappy.parquet, hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores/year=2019/month=1/day=26/part-00000-10000-30fea43a-140c-4f3c-be88-cd8a386add22.c000.snappy.parquet, hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores/year=2019/month=1/day=27/part-00000-10000-30fea43a-140c-4f3c-be88-cd8a386add22.c000.snappy.parquet, hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores/year=2019/month=1/day=28/part-00000-10000-30fea43a-140c-4f3c-be88-cd8a386add22.c000.snappy.parquet, hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores/year=2019/month=1/day=29/part-00000-10000-30fea43a-140c-4f3c-be88-cd8a386add22.c000.snappy.parquet, hdfs://localhost:9000/user/ds/odsc/east/realish/eda/callscores/year=2019/month=1/day=30/part-00000-10000-30fea43a-140c-4f3c-be88-cd8a386add22.c000.snappy.parquet)
```

Took 0 sec. Last updated by anonymous at June 26 2019, 4:00:28 PM.

```
scoresData.createOrReplaceTempView("scores")
```



5 min break & QA then  
OLAP ANALYTICS

SQL / NoSQL / NewSQL

Modern Data Architecture

Zero to Hero: SparkSQL

Beyond SQL



@newfront  
@twilio

# Zero to Hero: Analytics

---

"Supporting OLAP with SparkSQL"

- Aggregation Functions
- Windowing and Watermarking
- Extending Aggregations (MutableAggregationBuffers)

```
verageSeizureRatio extends UserDefinedAggregateFunction {  
  
    val mainKey: String = "completed"  
  
    The actual fieldName doesn't matter here, it is used as an example.  
    def inputSchema: StructType = {  
        type(StructField("row", StringType) :: Nil)  
  
        ...  
    }  
  
    // aggregation buffer schema  
    def bufferSchema: StructType = {  
        ...  
        field("answered", LongType) :: StructField("total", LongType) :: Nil  
    }  
  
    // returned value  
    def dataType: DataType = DoubleType  
  
    // deterministic or not  
    def deterministic: Boolean = true  
  
    // initialize the aggregation buffer  
    def initialize(buffer: MutableAggregationBuffer): Unit = {  
  
        ...  
        buffer.setLong(0, 0L)  
    }  
  
    // update the aggregation buffer with a new row  
    def update(buffer: MutableAggregationBuffer, input: Row): Unit = {  
        ...  
        buffer.setLong(0, buffer.getLong(0) + 1L)  
    }  
  
    // merge two aggregation buffers  
    def merge(buffer1: MutableAggregationBuffer, buffer2: Row): Unit = {  
        ...  
        buffer1.setLong(0, buffer1.getLong(0) + buffer2.getLong(0))  
        buffer1.setLong(1, buffer1.getLong(1) + buffer2.getLong(1))  
    }  
  
    // evaluate the aggregation function  
    def evaluate(buffer: Row): Double = {  
        ...  
        buffer.getDouble(0).toDouble  
    }  
}
```

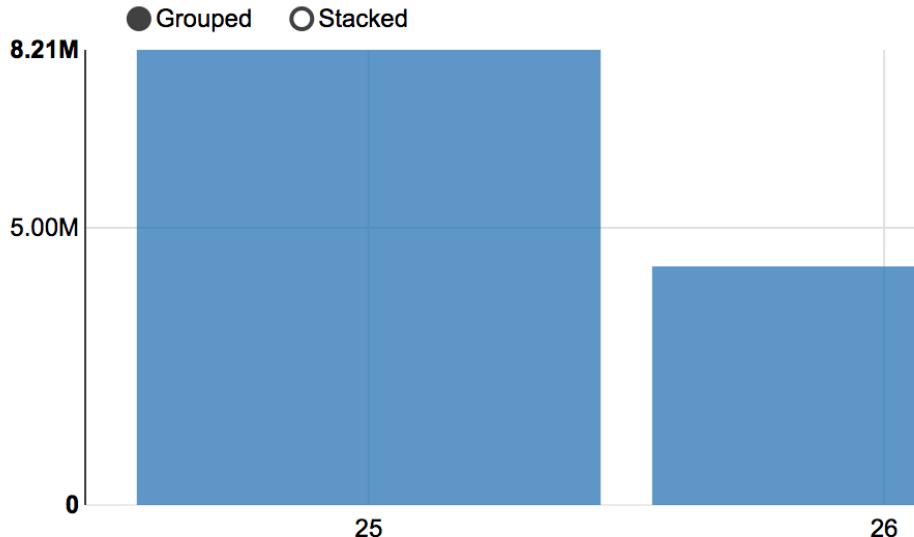


# Zero to Hero

## SparkSQL: Analytics

- Batch or Streaming queries are treated more or less the same\*
- Use the DSL or query interpreter

```
%sql  
select count(*) as total, day from scores  
where year=2019 and month=1 and day in (25, 26, 27)  
group by day
```



Took 11 sec. Last updated by anonymous at June 26 2019, 3:55:26 PM. (outdated)

# Zero to Hero

## SparkSQL: Analytics

---

- Hate SQL. Write SQL as a DAG!

```
%spark
val aggs = scoresData
    .select("account_sid", "dimensions_country_code", "score")
    .filter(not(col("dimensions_country_code") === ""))
    .groupBy("account_sid", "dimensions_country_code")
    .agg(
        count("*") as "total",
        min("score") as "min_score",
        avg("score") as "avg_score",
        max("score") as "max_score"
    )
    .sort(desc("total"))

aggs.createOrReplaceTempView("score_aggs")

aggs: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row]
... 4 more fields]
```

Took 2 sec. Last updated by anonymous at June 26 2019, 4:08:13 PM.

# Zero to Hero

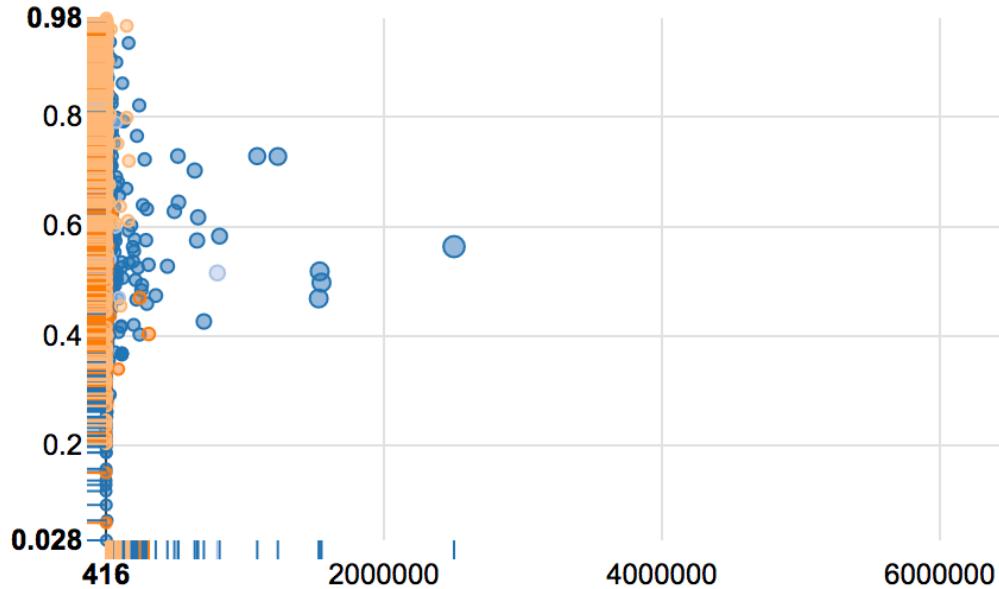
## SparkSQL: Analytics

- Not Strictly SQL.
- Not Strictly MapReduce.
- Use what works.

```
%sql  
select  
account_sid, dimensions_country_code, total, min_score, avg_score,  
from score_aggs  
where dimensions_country_code in ("US", "MX", "GB", "CA")
```



settings ▾



# Zero to Hero

## SparkSQL: Analytics

- Aggregate and groupBy with ease.

```
%spark
val carrierEdgebyCountry = scoresData
    .filter($"dimensions_edge_type" === "carrier_edge")
    .filter(not($"dimensions_country_code" === ""))
    .groupBy("dimensions_country_code", "dayofyear")
    .agg(
        count(*) as "total",
        avg($"score") as "mean_score",
        sum("tag_high_packet_loss") as "pktloss_sum",
        sum("tag_silent") as "silent_sum",
        sum("tag_high_pdd") as "high_pdd_sum",
        sum("tag_high_jitter") as "high_jitter_sum",
        sum("tag_high_latency") as "high_latency_sum",
        avg($"metrics_connect_duration") as "avg_connect_duration"
    )
carrierEdgebyCountry.createOrReplaceTempView("carrier_by_country")

carrierEdgebyCountry: org.apache.spark.sql.DataFrame = [dimensions_country_code: string, dayofyear: int ... 8 more fields]
```

# Zero to Hero

## SparkSQL: Analytics

- Extend SparkSQL Aggregations with the **MutableAggregationBuffer**
- Chain Complex Queries into a Simple Flow (Spark will do the heavy lifting -> LogicalPlan)

```
val report = reporter.generate()
    .withColumn("is200",
        when(col("last_sip_response_number").equalTo(200)).otherwise(0))
    .groupBy("parent_account_sid",
        .agg(
            count("call_sid") as "total",
            sum("is200") as "normal",
            sum("connect_duration") as "total_connect",
            bround(asr.col("call_state")))
        .bround(aloc.col("connect_duration"))
        .insight_tags(col("insight_tags"))
    )
    .withColumn("connected_minutes",
        col("total_connected_seconds")
        , 2)
    )
    .withColumn("normal_sessions", bround(lit(100.0) * (col("normal") / col("total"))))
    .drop("normal")
    .na.fill("", Seq("parent_account_sid"))
    .withColumn("reportId", lit(reportId))
```

/\*  
Answer-seizure ratio  
ASR = 100.\* ( answered / total ) - eg.. completion rate  
\*/  
**def asr:** AverageSeizureRatio.type = {  
 AverageSeizureRatio  
}  
  
/\*  
Average Length of Call  
sum(connect\_duration) / N  
\*/  
**def aloc:** AverageLengthOfCall.type = {  
 AverageLengthOfCall  
}

# Zero to Hero

## SparkSQL: Analytics

---

- Looking at creating an  
Distributed Aggregation -  
aka the Monoid

```
object AverageSeizureRatio extends UserDefinedAggregateFunction {  
  
    private val mainKey: String = "completed"  
  
    // note: The actual fieldName doesn't matter here it is the intent.  
    // - when you pass in (asr(col("call_state")) the schema is struct of string type  
    override def inputSchema: StructType = {  
        StructType(StructField("row", StringType) :: Nil)  
    }  
  
    // DataTypes of the values in the aggregation buffer  
    override def bufferSchema: StructType = {  
        StructType(  
            StructField("answered", LongType) :: StructField("total", LongType) :: Nil  
        )  
    }  
  
    // DataType of returned value  
    override def dataType: DataType = DoubleType  
  
    override def deterministic: Boolean = true  
  
    override def initialize(buffer: MutableAggregationBuffer): Unit = {  
        buffer(0) = 0L  
        buffer(1) = 0L  
    }  
  
    override def update(buffer: MutableAggregationBuffer, input: Row): Unit = {  
        if (!input.isNullAt(0)) {  
            if (input.getString(0) == mainKey) {  
                buffer(0) = buffer.getLong(0) + 1L  
            }  
            buffer(1) = buffer.getLong(1) + 1L  
        }  
    }  
  
    override def merge(buffer1: MutableAggregationBuffer, buffer2: Row): Unit = {  
        buffer1(0) = buffer1.getLong(0) + buffer2.getLong(0)  
        buffer1(1) = buffer1.getLong(1) + buffer2.getLong(1)  
    }  
}
```

# Zero to Hero

## SparkSQL: Analytics

---

- Looking at creating an  
Distributed Aggregation -  
aka the Monoid

```
object AverageSeizureRatio extends UserDefinedAggregateFunction {  
  
    private val mainKey: String = "completed"  
  
    // note: The actual fieldName doesn't matter here it is the intent.  
    // - when you pass in (asr(col("call_state")) the schema is struct of string type  
    override def inputSchema: StructType = {  
        StructType(StructField("row", StringType) :: Nil)  
    }  
  
    // DataTypes of the values in the aggregation buffer  
    override def bufferSchema: StructType = {  
        StructType(  
            StructField("answered", LongType) :: StructField("total", LongType) :: Nil  
        )  
    }  
  
    // DataType of returned value  
    override def dataType: DataType = DoubleType  
  
    override def deterministic: Boolean = true  
  
    override def initialize(buffer: MutableAggregationBuffer): Unit = {  
        buffer(0) = 0L  
        buffer(1) = 0L  
    }  
  
    override def update(buffer: MutableAggregationBuffer, input: Row): Unit = {  
        if (!input.isNullAt(0)) {  
            if (input.getString(0) == mainKey) {  
                buffer(0) = buffer.getLong(0) + 1L  
            }  
            buffer(1) = buffer.getLong(1) + 1L  
        }  
    }  
  
    override def merge(buffer1: MutableAggregationBuffer, buffer2: Row): Unit = {  
        buffer1(0) = buffer1.getLong(0) + buffer2.getLong(0)  
        buffer1(1) = buffer1.getLong(1) + buffer2.getLong(1)  
    }  
}
```



5 min break & QA then  
INSIGHTS



@newfront  
@twilio

SQL / NoSQL / NewSQL

Modern Data Architecture

Zero to Hero: SparkSQL

**Beyond SQL**

# Zero to Hero: Insights

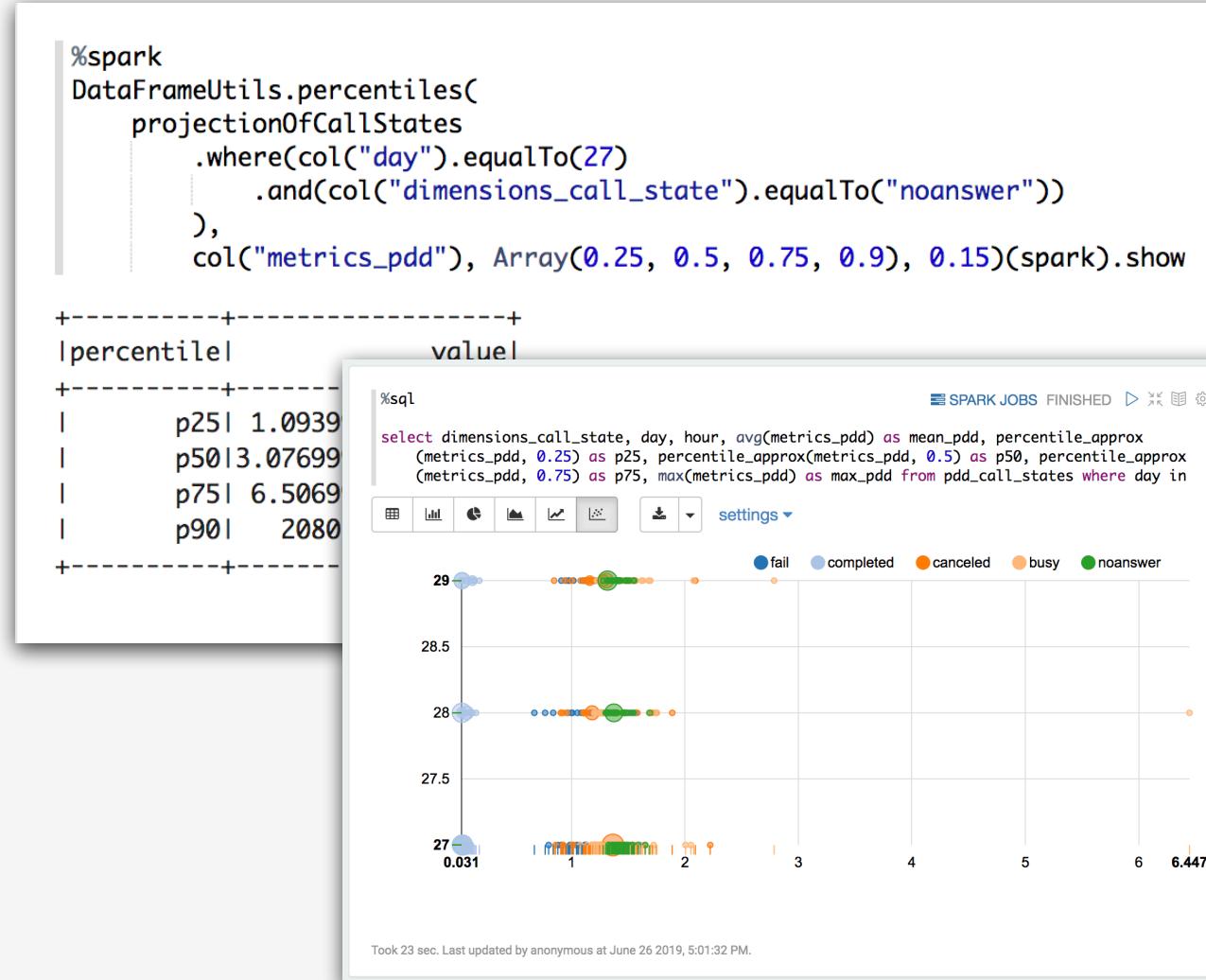
"Insights are the Why behind your Analytics"

- Find patterns in your data.
  - Deliver value to customers

```
val report = reporter.generate()  
  .withColumn("is200", when(col("is200").isTrue(), 1).otherwise(0))  
  .groupBy("parent_account_sid")  
  .agg(  
    count("call_sid") as "total",  
    sum("is200") as "normal",  
    sum("connect_duration") as "duration",  
    bround(asr(col("call_state")) * 100) as "asr",  
    bround(aloc(col("connect_duration")) / 1000) as "avg_connect_duration",  
    insight_tags(col("insight_tags")) as "tags"  
  )  
  
  .withColumn("connected_minute", col("total").divide(lit(60)))  
  .withColumn("total_connected_second", col("duration").divide(lit(1000)))  
  .withColumn("normal_sessions", col("normal").divide(lit(2)))  
  
  .fill("", Seq("parent_account_sid", "reportId", "tags"))  
  .column("reportId", lit("report"))
```

# Beyond SQL Insights

- **Percentiles** help provide a sense of what has been\*
- **Explain** why things changed by triggering additional jobs
  - ML / Reports / etc
- **Triangulation:** Find a quorum in your data to answer a question



# Beyond SQL

## BYO-Logic

---

- **Window:** [start-end time]  
for a grouped operation
- **Watermark:** [throw away  
late arriving data]
- Use Spark to coordinate,  
you manage data in  
StateStore

```
// 1. Convert from the kafka data frame to an instance of our Metric
val groupedMetrics = df.mapPartitions(
    _.flatMap { kd =>
        bytesToMetricRow(kd.getAs[Array[Byte]]("value"), windowInterval)
    }
)(metricEncoder)
// 1b. apply watermark to ditch updates to said record state and allow for e
// 2. groupBy binary key
// 3. aggregate the value of the metric
// 4. output the final aggregations
groupedMetrics
    .withWatermark("timestamp", s"${config.watermarkIntervalMinutes} minutes")
    .groupByKey(_.metricAggregation)
    .flatMapGroupsWithState[Array[Byte], MetricAggregation](
        outputMode,
        GroupStateTimeout.EventTimeTimeout()
    )(metricAggregator)
```

# FROM ZERO TO DATA HERO: HOW TO EFFECTIVELY USE VOICE INSIGHTS DATA TO SAVE THE DAY

How's it going? Learn from the experts on the Voice Insights team just how to answer common Business Intelligence questions over the course of their workshop. We will detail how organizations can start using the Voice Insights APIs in order to start answering their core Business Intelligence questions (in a streaming fashion) while leveraging expert insights from the team behind the data.

**Tuesday, Aug 6**  
4:55 PM to 5:25 PM

Intermediate

Talk

Best Practices

Product Innovation

Voice

Developer / Engineer

Product Management

Solutions or Systems Architect

Systems Administrator



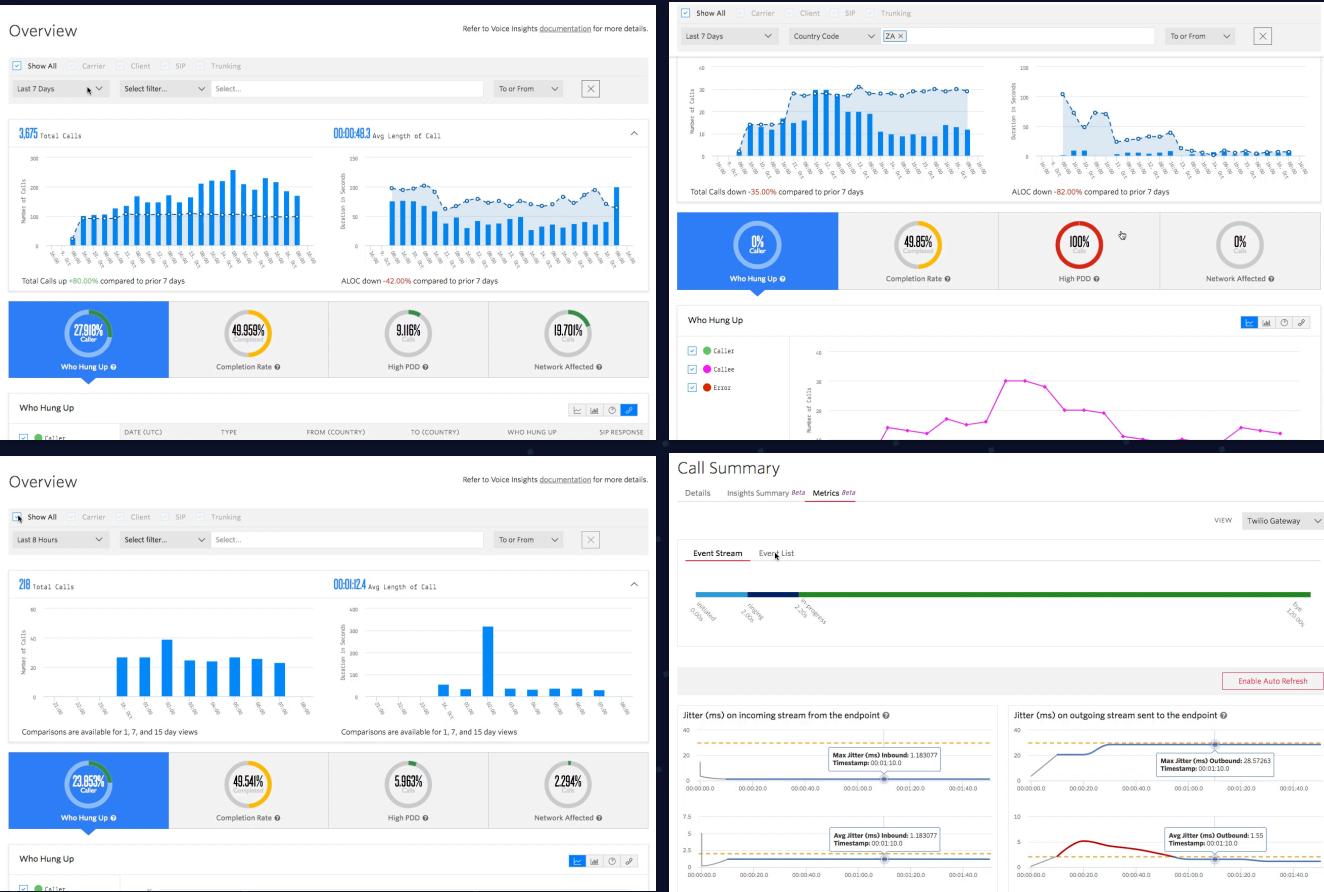
# What I Do Voice Insights

## Voice Analytics Platform

Millions of complex time series aggregations daily across billions of raw events

Supports internal and external customers.

\*Near Real-Time



twitter: [@newfront](#)

blog: [blog.twilio.com](http://blog.twilio.com)

# THANK YOU

