# Gerard Maas

Señor SW Engineer

Lightbend

🐦 @maasg

https://github.com/maasg

https://www.linkedin.com/in/gerardmaas/

https://stackoverflow.com/users/764040/maasg

O'REILLY®

Stream
Processing with
Apache Spark

BEST PRACTICES FOR SCALING AND OPTIMIZING APACHE SPARK

François Garillot & Gerard Maas

O'REILLY®

Designing Fast
Data Application
Architectures

Gerard Maas, Stavros Kontopoulos
& Sean Glover

Lightbend

# Data Pipelines

# Data Pipelines

- Create **Composable** Streaming Applications
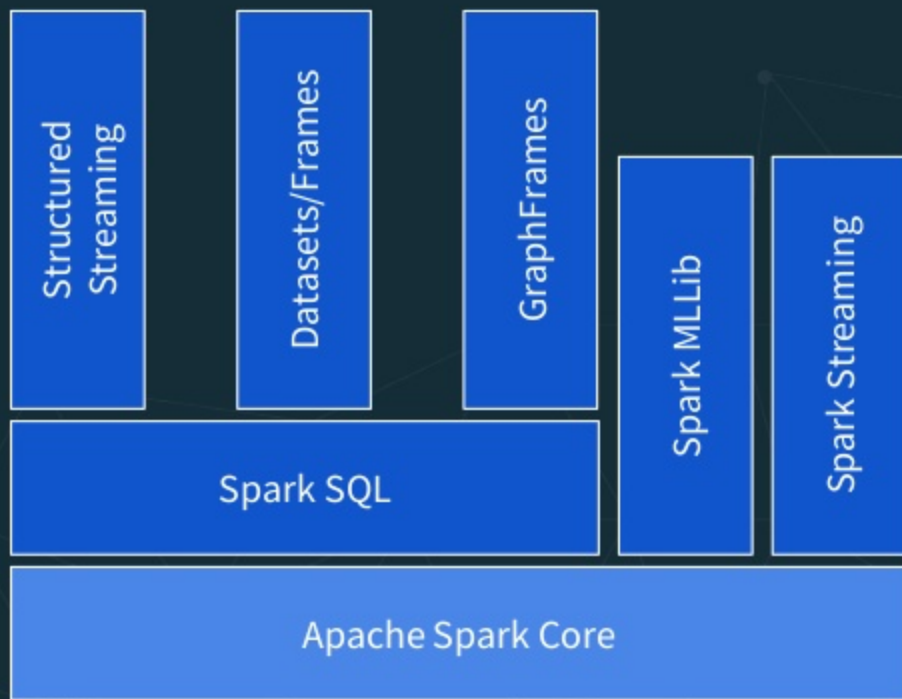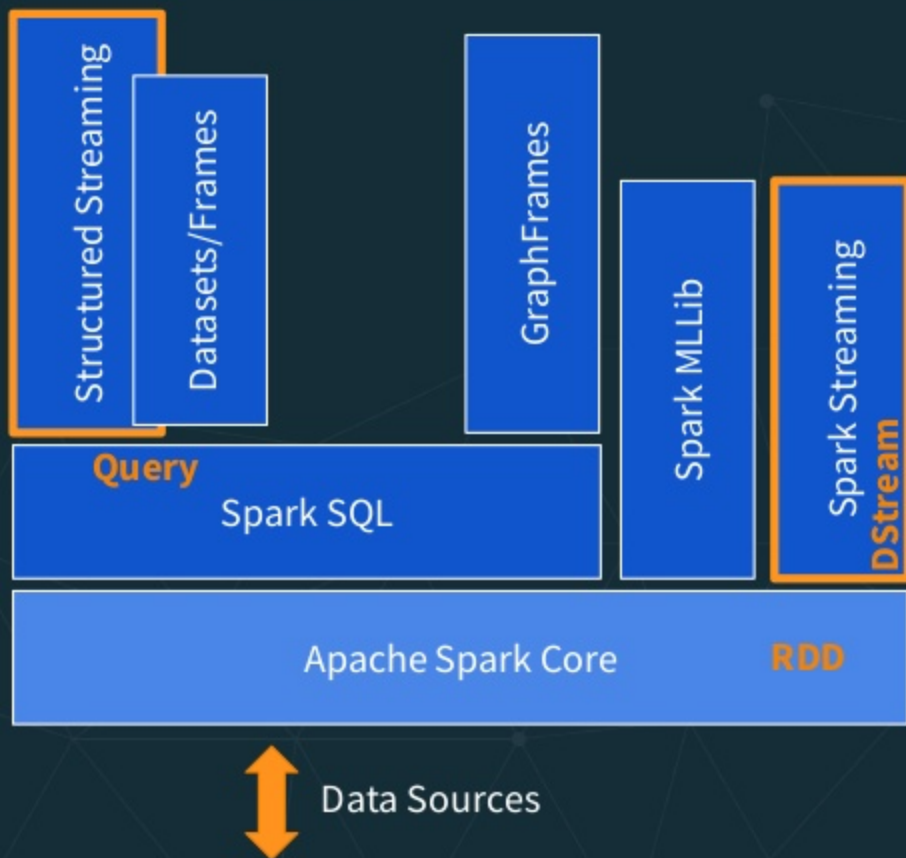- Using the **Best Tool** for the Job
- Generating a **Network Effect**

# Agenda

Creating a Fast Data Pipeline with
**Structured Streaming** and
**Spark Streaming**

# Spark Structured Streaming

```scala
val lines = spark.readStream
  .format("socket")
  .option("host", "localhost")
  .option("port", 9999)
  .load()
val words= lines.as[String].flatMap(_.split(" "))
val wordCounts = words.groupBy("value").count()

val query = wordCounts.writeStream
  .outputMode("complete")
  .format("console")
  .start()
```

# Spark Streaming

```scala
val ctx= new StreamingContext(conf,Seconds(1))

val lines = ssc.socketTextStream("localhost", 9999)

val words = lines.flatMap(_.split(" "))

val pairs = words.map(word => (word, 1))
val wordCounts = pairs.reduceByKey(_ + _)

wordCounts.print()

ctx.start()
```

Lightbend

@maasg

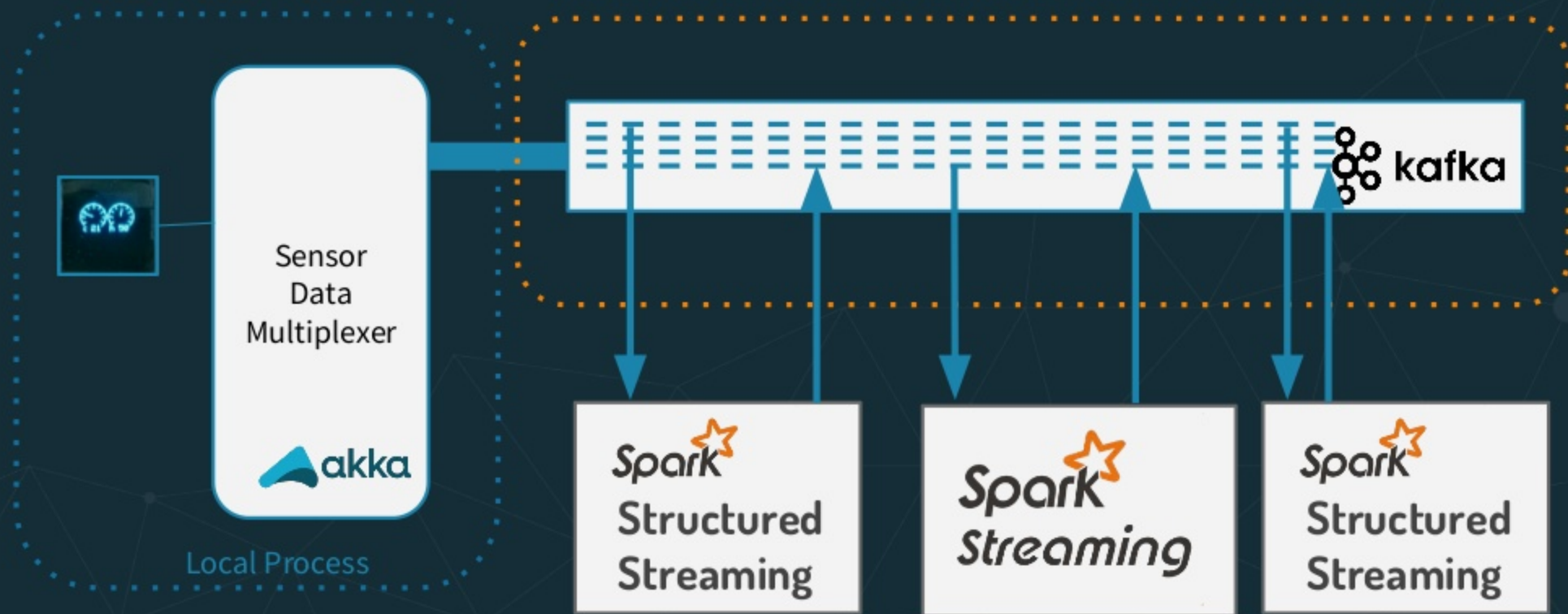| | Structured Streaming | Spark Streaming |
|---|---|---|
| Time | Abstract (Processing Time, Event Time) | Fixed to microbatch Streaming Interval |
| Execution | Fixed Micro batch, Best Effort MB, Continuous (NRT) | Fixed Micro batch |
| Abstraction | DataFrames/Dataset | DStream, RDD |

⭐ Access to the scheduler

Lightbend

@maasg

# Agenda

## Hands On with Spark:
Creating a Fast Data Pipeline with Structured Streaming and Spark Streaming

Lightbend

@maasg

# Sensor Anomaly Detection Pipeline

**Data Exploration**

[Structured Streaming]

Sensor
Data
Multiplexer

akka

**Data Preparation**

[Structured Streaming]

**Online Model
Creation +
Training**

[Spark Streaming]

**Anomaly Detection**

[Structured Streaming]

Lightbend

@maasg

# Sensor Anomaly Detection

# Sensor Anomaly Detection



Senso
Data
Multipl

Local Proce

kafka

Spark
Structured
Streaming

Lightbend

@maasg

🔴 **Live**

# Sensor Anomaly Detection Pipeline

**Data Exploration**

[Structured Streaming]

**Kafka Source
Memory Sink
SQL Operations**

Sensor
Data
Multiplexer

**akka**

**Data Preparation**

[Structured Streaming]

**Online Model
Creation +
Training**

[Spark Streaming]

**Anomaly Detection**

[Structured Streaming]

Lightbend

@maasg

# Sensor Anomaly Detection Pipeline

**Data Exploration**

[Structured Streaming]

Kafka Source
Memory Sink
SQL Operations

Sensor
Data
Multiplexer

akka

**Data Preparation**

[Structured Streaming]

**Online Model
Creation +
Training**

[Spark Streaming]

Event Time
Windows
Watermark
Kafka Sink

**Anomaly Detection**

[Structured Streaming]

Lightbend

@maasg

# Sensor Anomaly Detection Pipeline

Sensor Data Multiplexer
akka

Data Exploration

[Structured Streaming]

Kafka Source
Memory Sink
SQL Operations

Data Preparation

[Structured Streaming]

Online Model Creation + Training

[Spark Streaming]

RDD Programming
Local vs Distributed
Use Spark SQL
Kafka Source + Sink

Event Time
Windows
Watermark
Kafka Sink

Anomaly Detection

[Structured Streaming]

Lightbend

@maasg

# Sensor Anomaly Detection Pipeline

|  | **Structured Streaming** | **Spark Streaming** |
|---|---|---|
| Time | Abstract (Processing Time, Event Time) | Fixed to microbatch Streaming Interval |
| Execution | Fixed Micro batch, Best Effort MB, Continuous (NRT) | Fixed Micro batch |
| Abstraction | DataFrames/Dataset | DStream, RDD |

⭐ Access to the scheduler

Lightbend

@maasg

# Resources

**Notebooks used today:**

https://github.com/maasg/spark-notebooks/tree/master/streaming-anomaly-detection

**Pipelines:**

https://www.reactivesummit.org/2018/schedule/taking-the-pain-out-of-deploying-streaming-applications
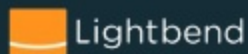
**Structured Streaming + Spark Streaming:**

https://www.reactivesummit.org/2018/schedule/processing-fast-data-with-apache-spark-the-tale-of-two-streaming-apis

**Fast Data:**

https://www.lightbend.com/products/fast-data-platform

Lightbend

# Gerard Maas

Señor SW Engineer

 Lightbend

 @maasg

 https://github.com/maasg

 https://www.linkedin.com/in/gerardmaas/

 https://stackoverflow.com/users/764040/maasg



Lightbend