

The logo for Spark Summit Europe 2016 features a stylized, circular graphic composed of many thin, white lines radiating from a central point, creating a sunburst or network-like effect. The text "SPARK SUMMIT" is in white, uppercase letters, and "EUROPE 2016" is in black, uppercase letters, both positioned to the right of the graphic.

**SPARK SUMMIT**  
**EUROPE 2016**

# Lambda Architecture in the IoT

Fast Data Analytics with Spark Streaming and MLlib

Bas Geerdink  
ING

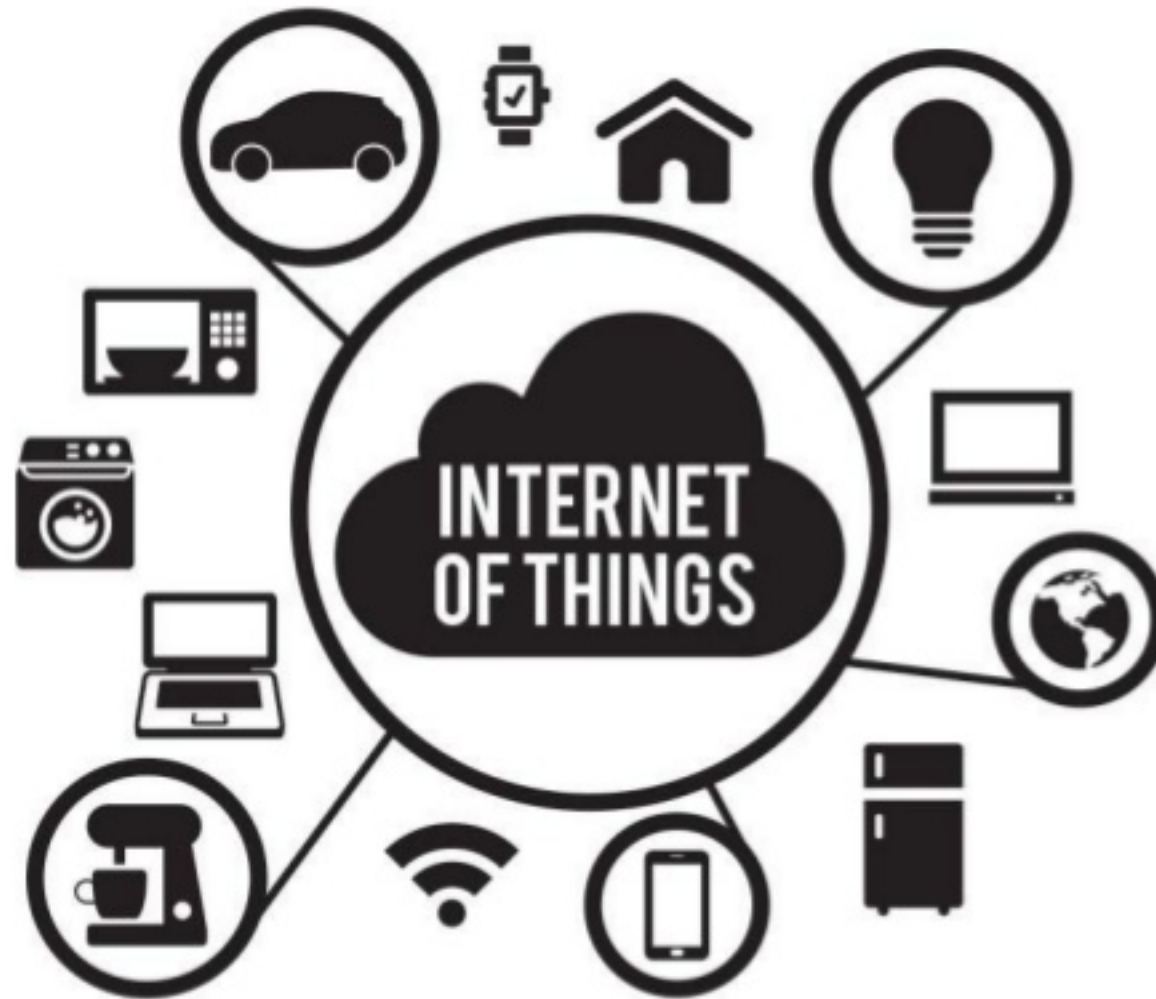
# Who am I?

Bas Geerdink

- Chapter Lead in Analytics area at ING
- Master degree in Artificial Intelligence and Informatics
- Spark Certified Developer
- @bgeerdink
- <https://www.linkedin.com/in/geerdink>



# The Internet of ...





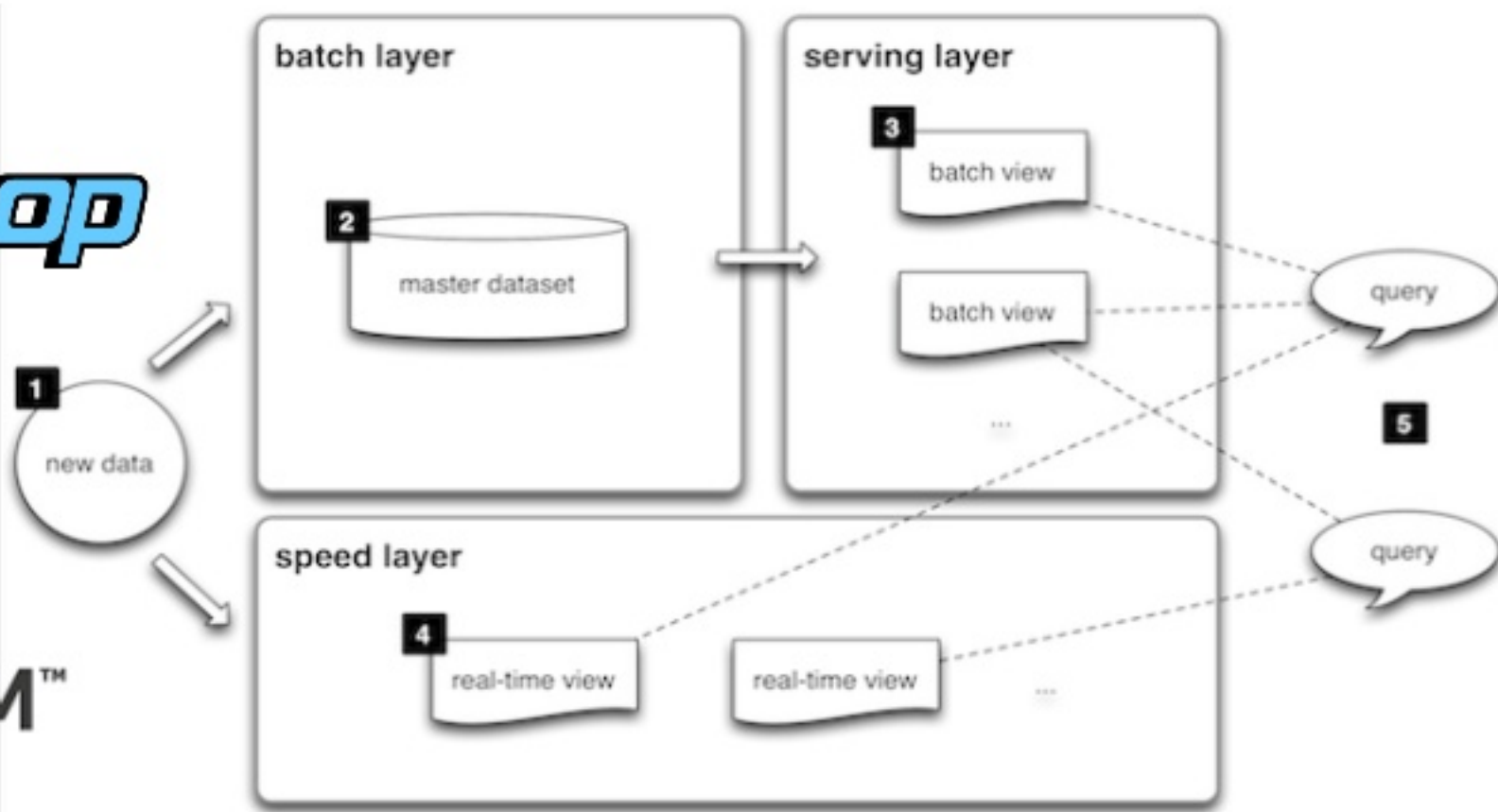
# What's new in the IoT?

- Data
  - Streaming data from more sources
- Use cases
  - Combining data streams
- Technology
  - Fast processing and scalability
- Challenges
  - Security: encrypt the sensors/network/server

# What do we want in the IoT?

- FAST data: process stream of events (sensory data)
- BIG data: process files/tables in batches (static data)
- ONE analytics engine
- ONE querying/visualization tool
- Scalable environment
- Reactive software

# Lambda Architecture



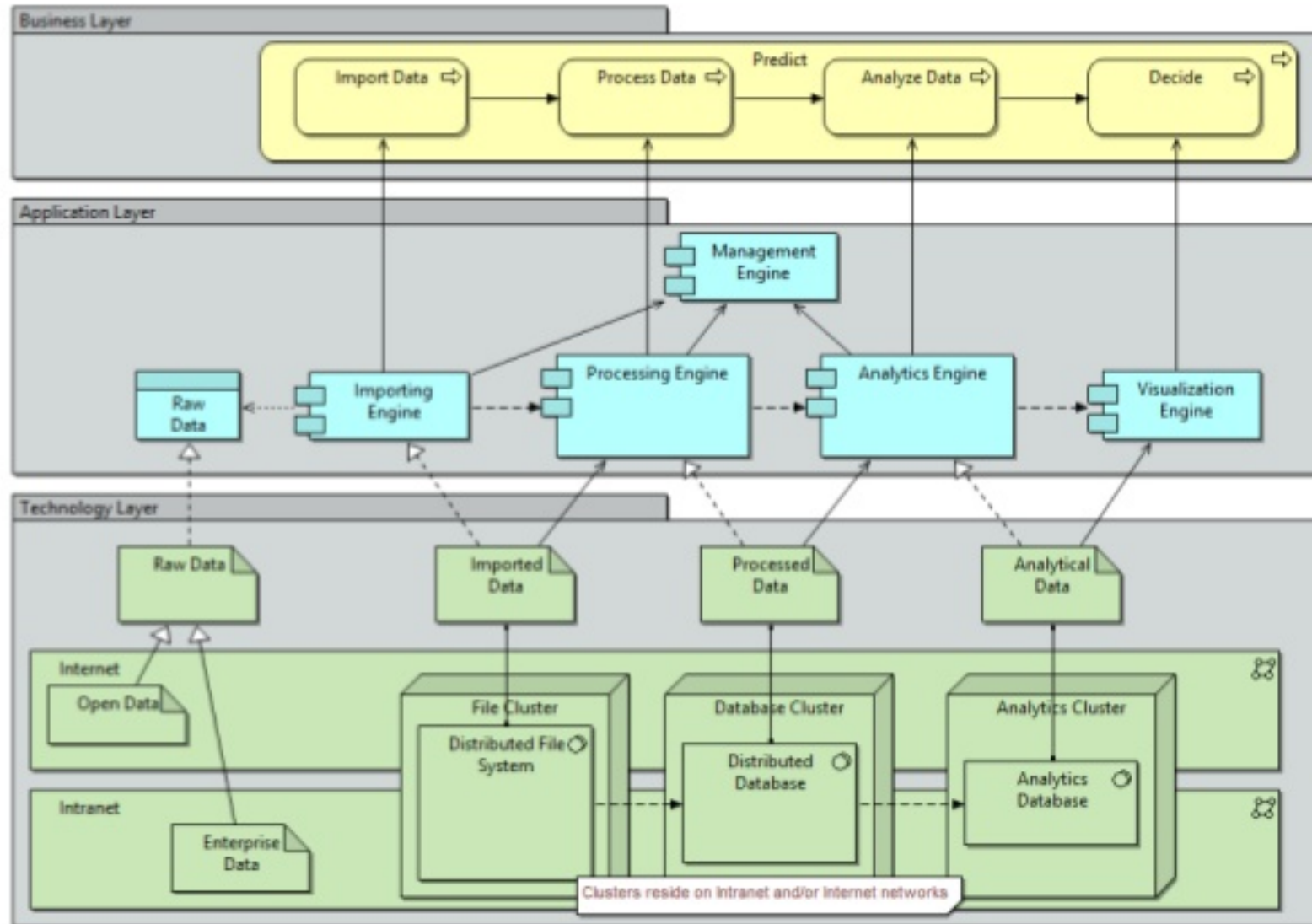
Source: Nathan Marz (2013)

# Gamma/Kappa/Omega/...

## Architecture

- Lambda Criticism:
  - Too complex
  - Two code bases
  - Two data stores
- Alternatives:
  - Treat a stream as a mini-batch
  - Treat a batch as a stream of records
  - Combine batch and stream within one system

# Big Data Reference Architecture



Source: Geerdink (2013)



# Use case: Smart Parking

*Recommend the best car park when driving to Amsterdam*

- Show the car park with the highest score, determined by
- Batch (every x minutes), data from car parks
  - Stream (every x seconds), GPS data of cars



## Parkeren + Reizen (P+R)

Parkeer uw auto voor € 8,00 per 24 uur (na 10.00 uur voor € 1,00 per 24 uur) aan de rand van Amsterdam en reis met openbaar vervoer naar het centrum.

### P+R locaties

Meer informatie per P+R locatie en uitleg hoe het werkt:

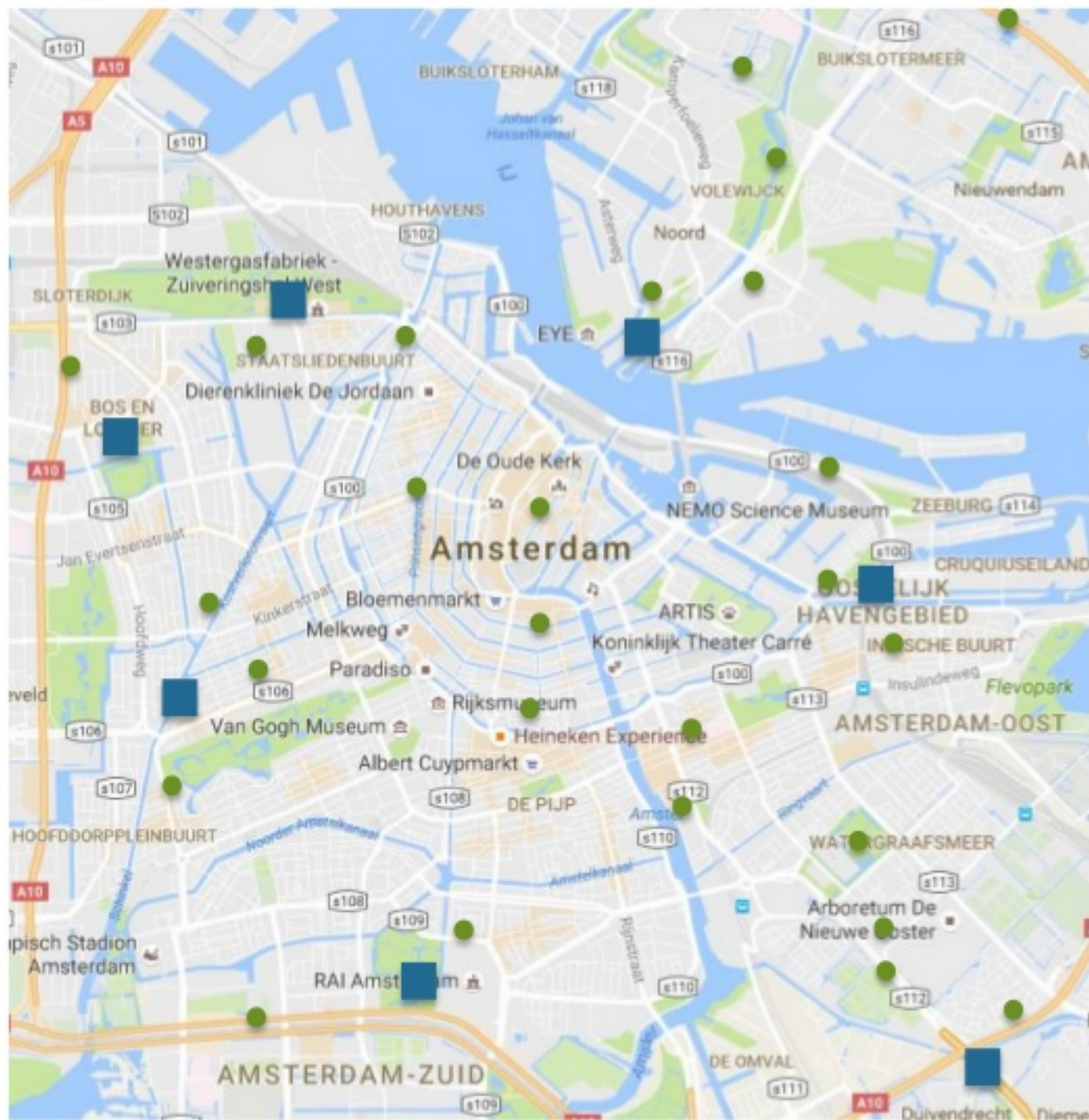
- P+R ArenA
- P+R Zeeburg I en II
- Weekend P+R VUmc
- P+R Bos en Lommer
- P+R RAI
- P+R Sloterdijk
- P+R Olympisch Stadion
- Welke P+R kies ik?

### Actuele beschikbaarheid P + R parkeerplekken

Laatste update: 25-okt-2016 00:21 (elke minuut)

P+R locatie	Beschikbaarheid	Parkeerplekken
P+R ArenA	Vrij	1.440
P+R Zeeburg 2	Vrij	322
P+R Olympisch Stadion	Vrij	161
P+R Zeeburg 1	Vrij	60
P+R Sloterdijk	Geen informatie	
P+R Bos en Lommer	Vrij	21
Weekend P+R VUmc	Geplaatst	







# Stream Process

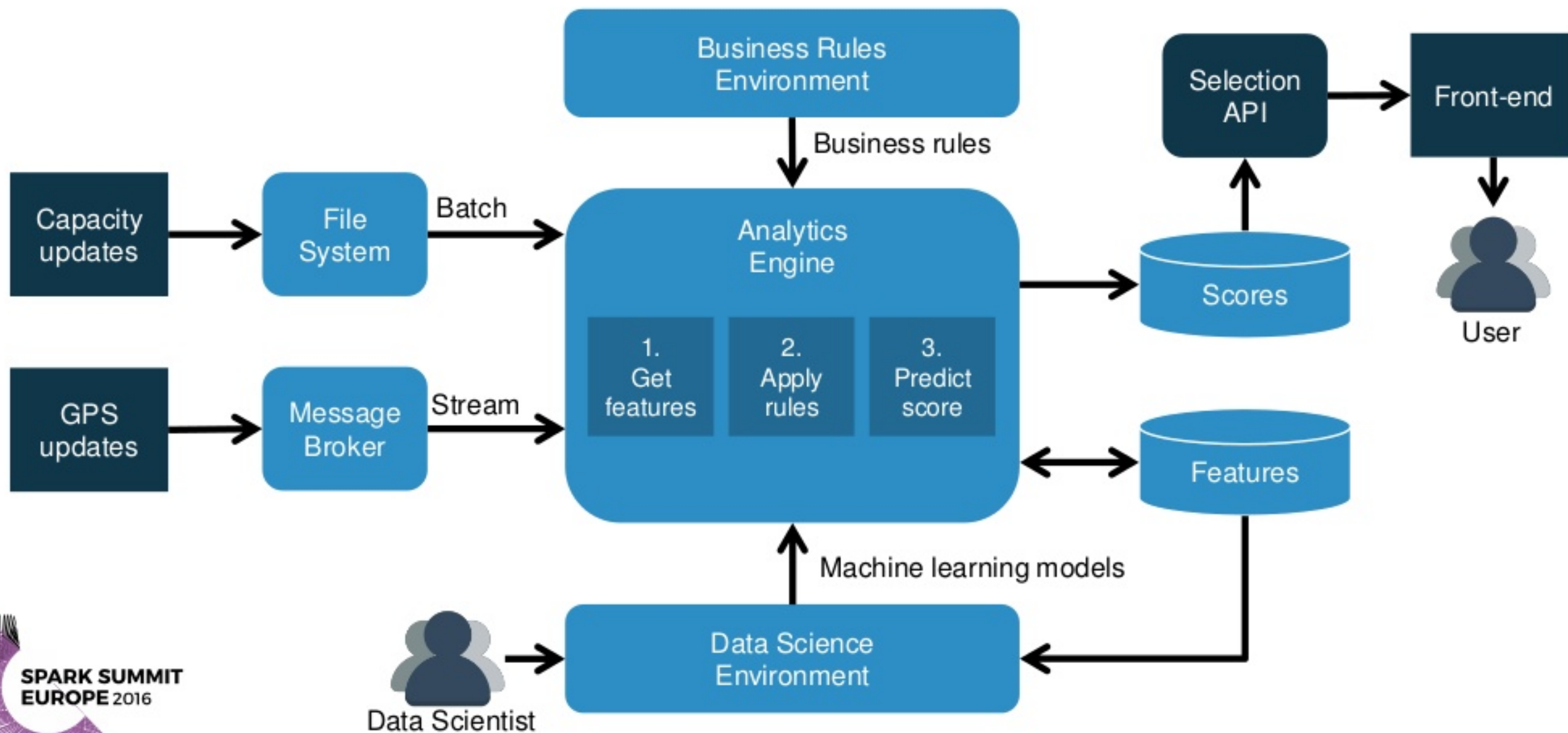
- Get car events (GPS data)
- Filter events (business rules)
- Store events
- For each car park in the neighborhood:
  - Get feature set (location, capacity, usage, ...)
  - Combine event with previous events (running sum)
  - Update feature set
  - Predict score (machine learning) and update database



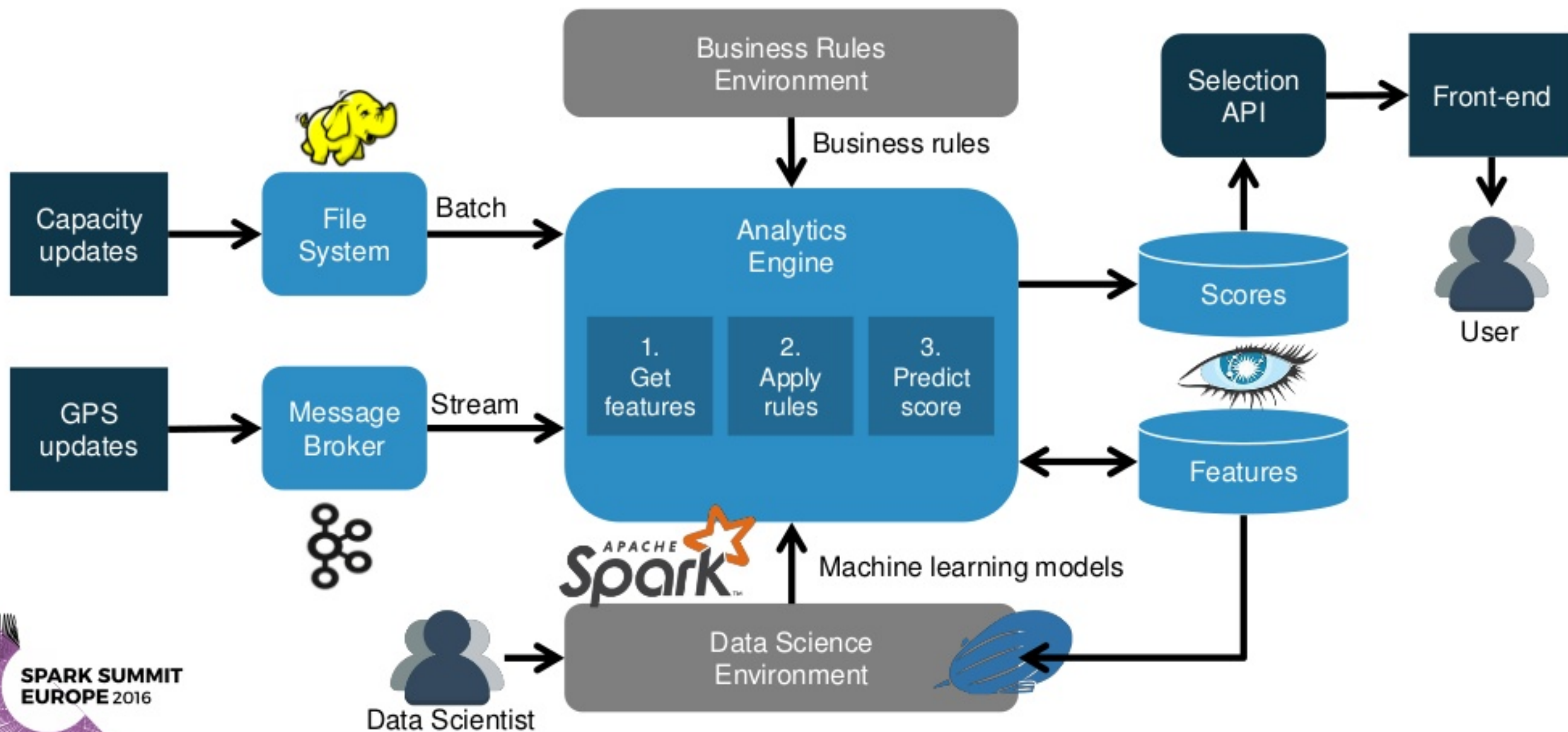
# Batch Process

- Get car park data
- Clean up (remove old car data)
- For each car park in the data set:
  - Get recent set of car data (running sum)
  - Update feature set
  - Predict score (machine learning) and update database

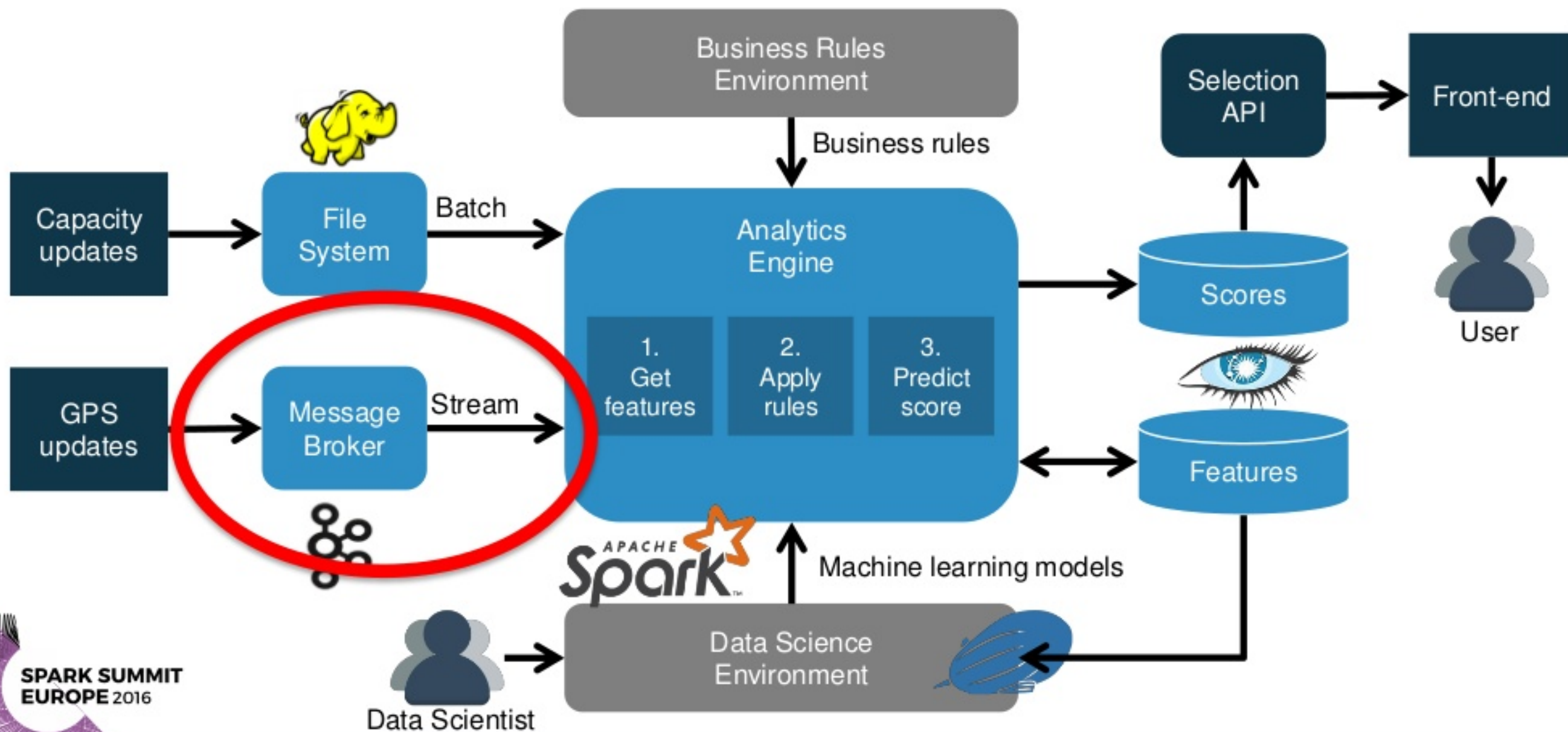
# Lambda Architecture



# Lambda Architecture



# Lambda Architecture





# Event Processing with Kafka

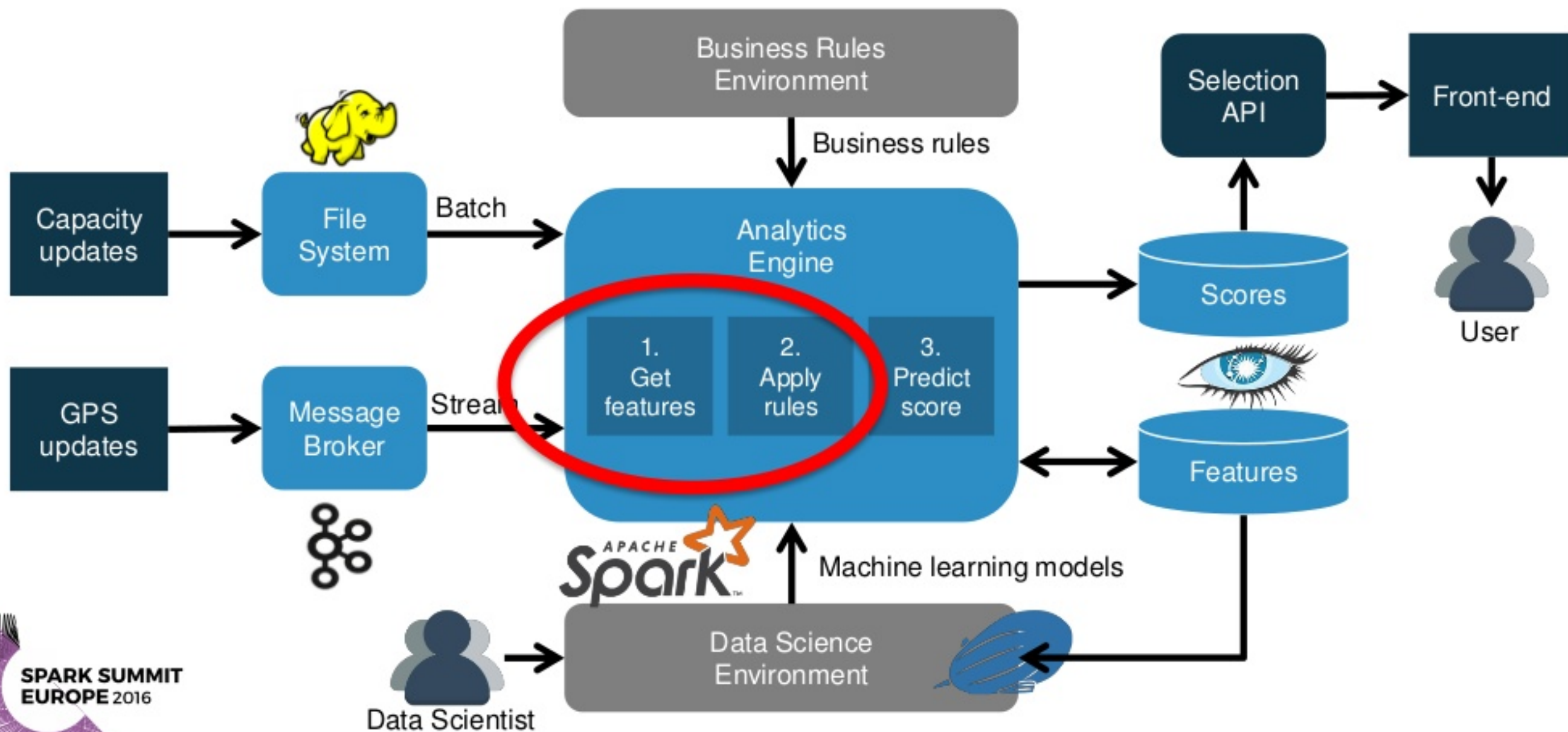
```
import org.apache.spark.streaming.{Seconds, StreamingContext}
import org.apache.spark.streaming.kafka010._
import org.apache.spark.streaming.kafka010.LocationStrategies.PreferConsistent
import org.apache.spark.streaming.kafka010.ConsumerStrategies.Subscribe

// initialize Spark Streaming
val conf = new SparkConf().setAppName("fast-data").setMaster("local[*]")
val ssc = new StreamingContext(conf, Seconds(1)) // batch interval = 1 sec

// set parameters for Kafka connection
val topics = Array("cars")
val kafkaParams = Map[String, Object]("bootstrap.servers" -> "localhost:9092")

// subscribe to stream -> create Spark DStream
val stream = KafkaUtils.createDirectStream[String, String](
  ssc,
  PreferConsistent,
  Subscribe[String, String](topics, kafkaParams))
```

# Lambda Architecture



# Stream: Data Preparation

```
// change raw data to business events  
val events = stream  
| .map(event => CarLocationHelper.createCarLocation(event.value))  
  
// apply business rules  
val filtered = events  
| .filter(Location.filterLocalArea)    // only select cars in local area  
  
// store car locations (update or create)  
filtered  
| .foreachRDD(rdd => rdd.foreach(cle => { CassandraHelper.insertCarLocation(cle) })))
```



# Lambda Architecture

