# MATLAB EXPO 2018
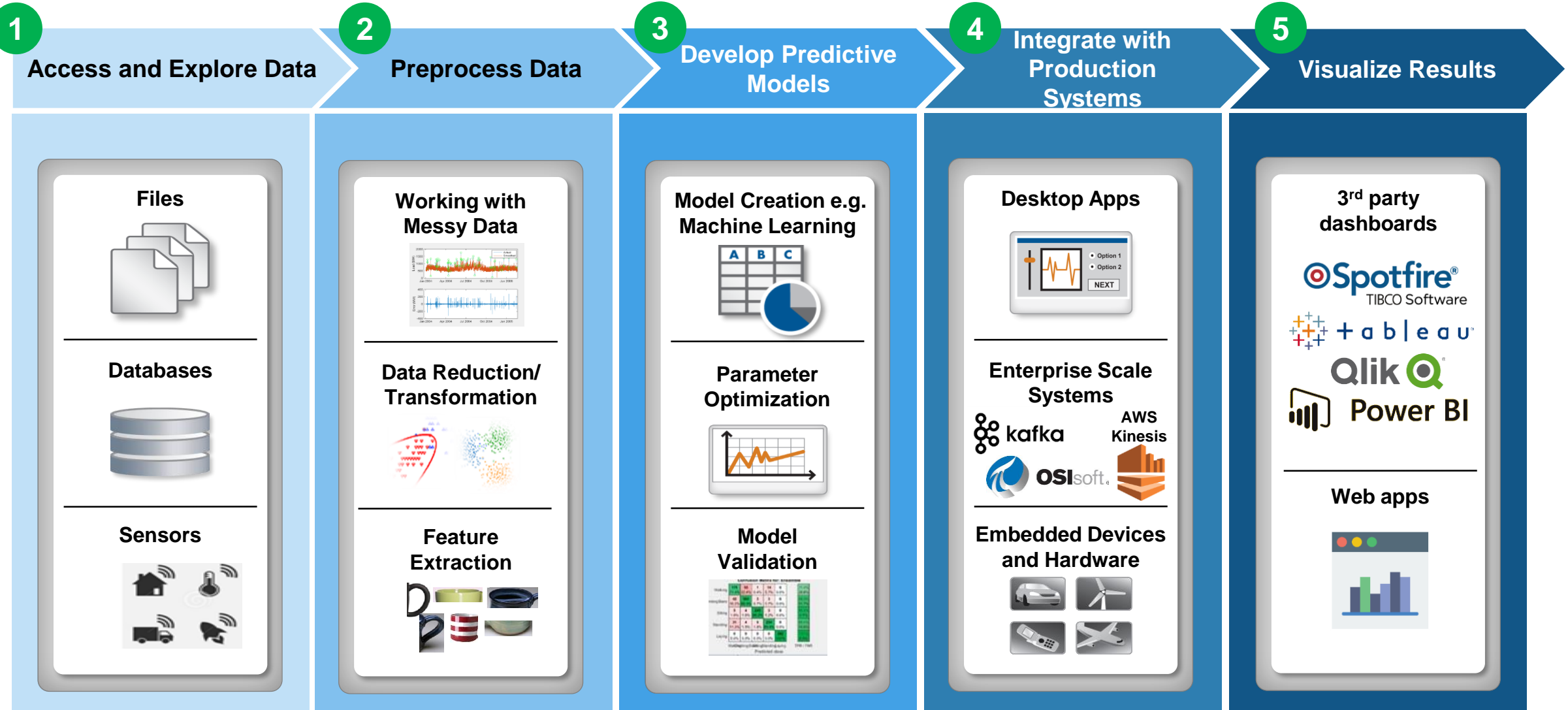
## Scaling up MATLAB Analytics with Kafka and Cloud Services
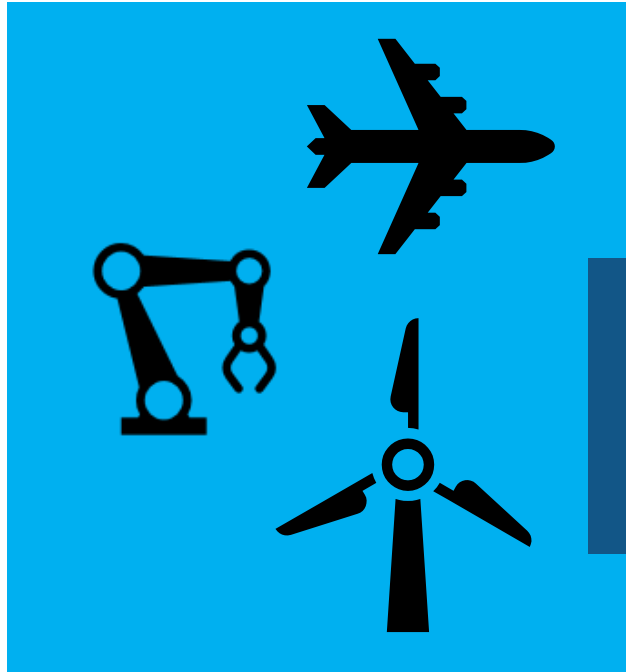
Christoph Stockhammer

# The Need for Large-Scale Streaming



## Predictive Maintenance

*Increase Operational Efficiency*

*Reduce Unplanned Downtime*

**More applications require near real-time analytics**

Jet engine: ~800TB per day
Turbine:      ~ 2 TB per day

## Medical Devices

*Patient Safety*

*Better Treatment Outcomes*

## Connected Cars

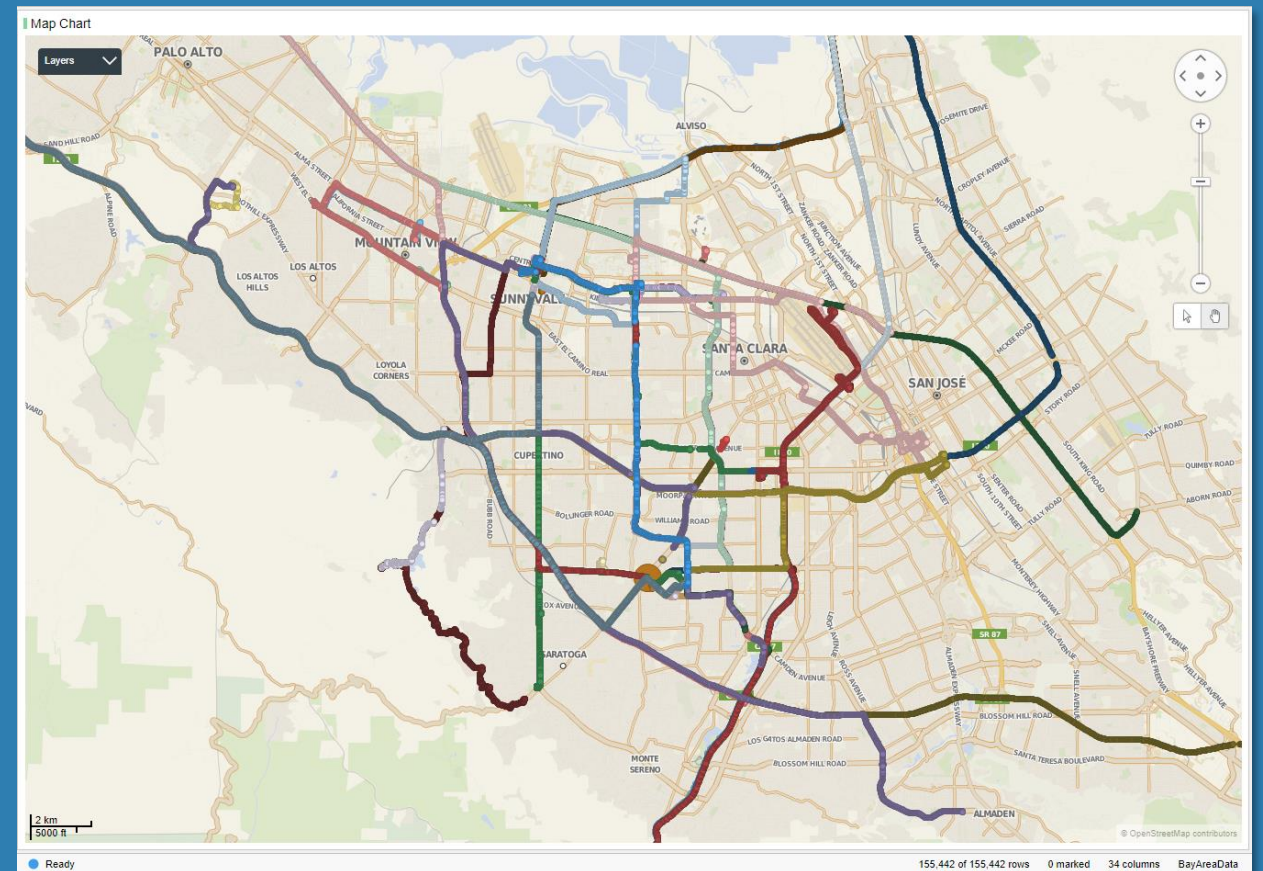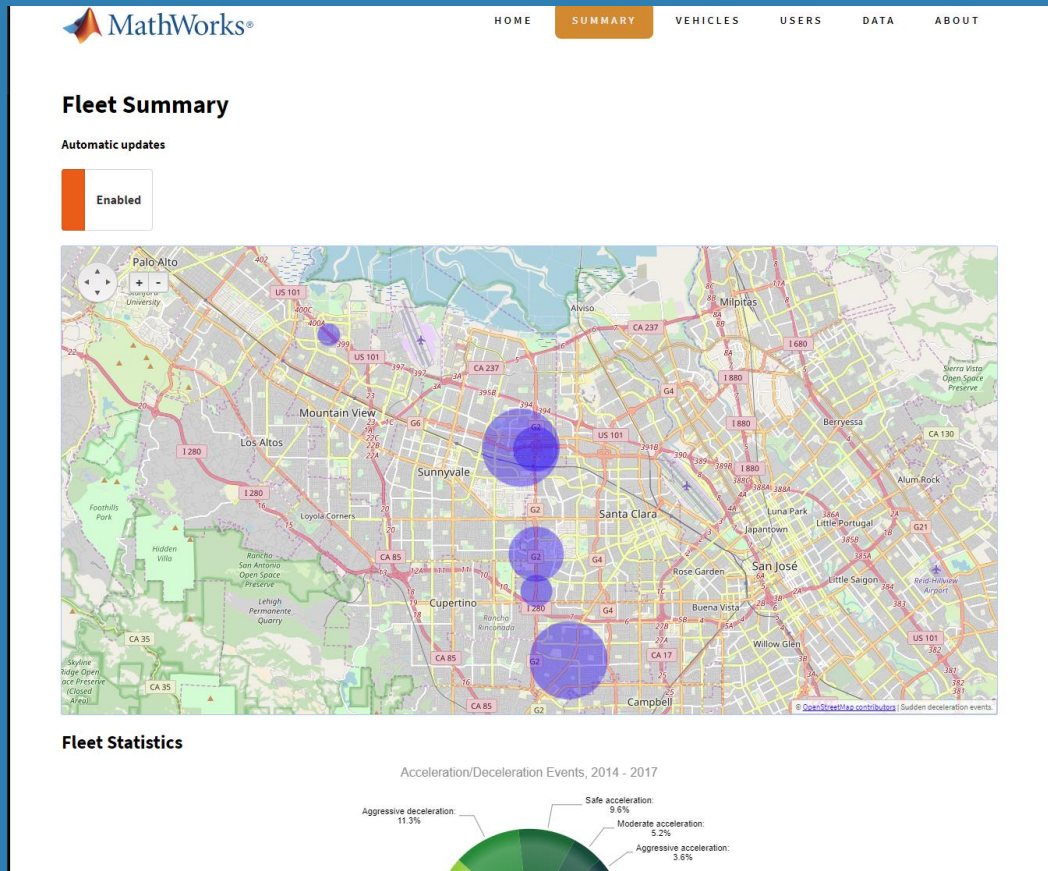*Safety, Maintenance*

*Advanced Driving Features*

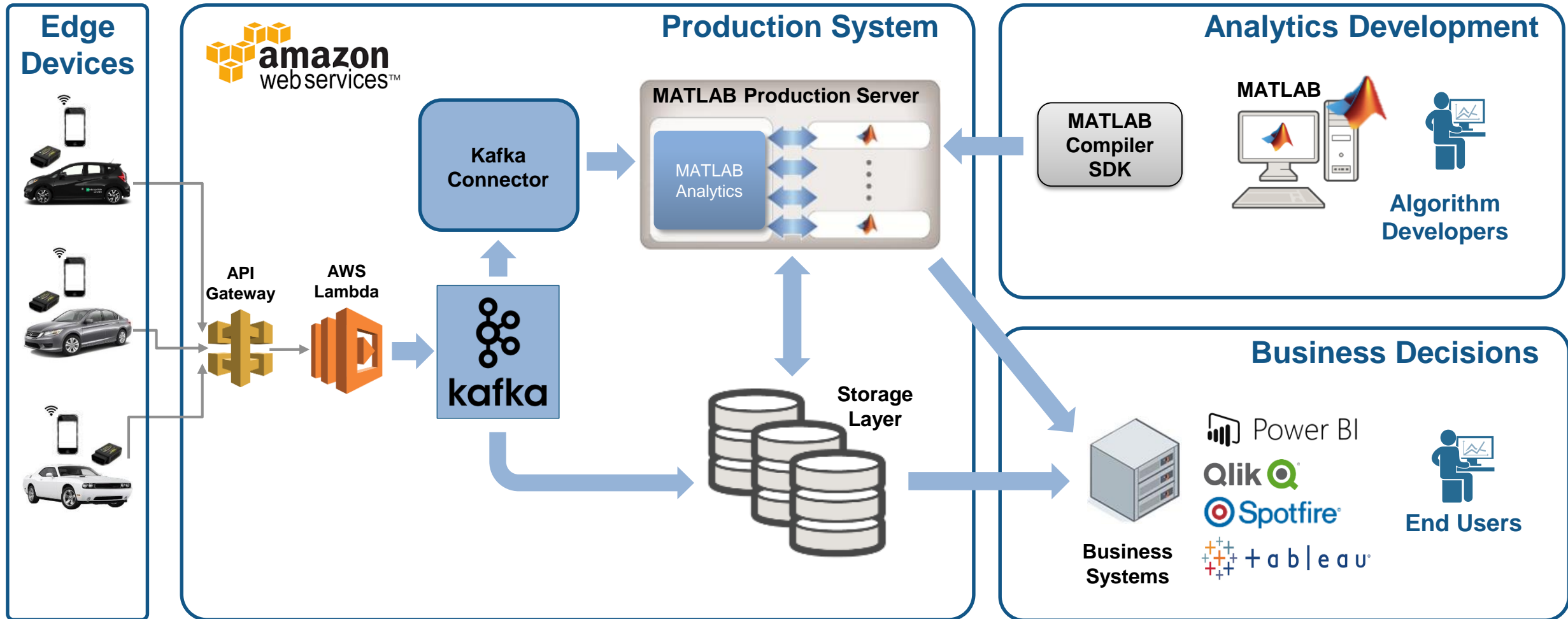Car: ~25 GB per hour

# Example Problem – How's my driving?

- A group of MathWorks employees installed an OBD dongle in their car that monitors the on-board systems

- Data is streamed to the cloud where it is aggregated and stored

- We would like to use this data to score the driving habits of participants

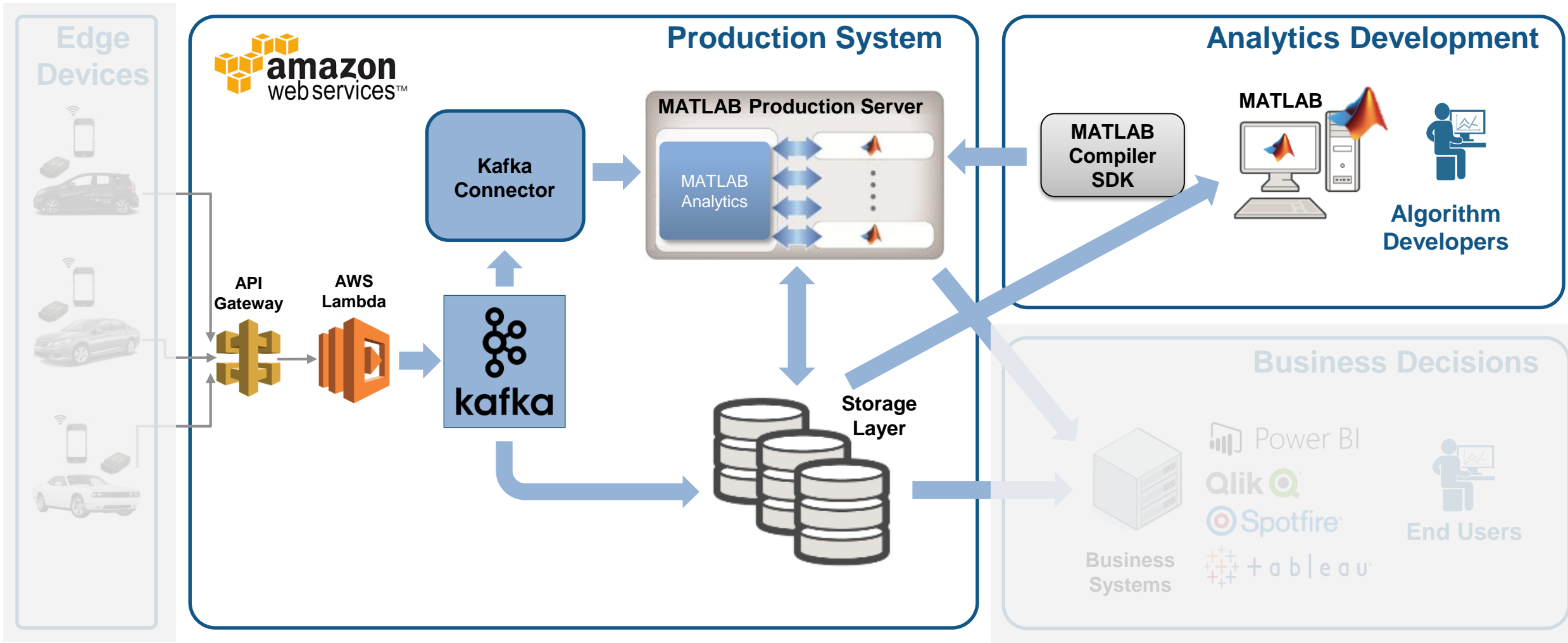# Example: Fleet Analytics with MATLAB

# Fleet Analytics Architecture

MathWorks

**Edge Devices**

**amazon** web services™

**Production System**

Kafka Connector

API Gateway

AWS Lambda

kafka

**MATLAB Production Server**

MATLAB Analytics

Storage Layer

**Analytics Development**

**MATLAB**

MATLAB Compiler SDK

Algorithm Developers

**Business Decisions**

Business Systems

Power BI
Qlik Q
Spotfire
+ableau

End Users

**1** **Access and Explore Data**

# The first step is to clean up the incoming data



**Edge Devices**

**Production System**

**Analytics Development**

amazon web services™

**Kafka Connector**

**MATLAB Production Server**

MATLAB Analytics

**MATLAB Compiler SDK**

**MATLAB**

**Algorithm Developers**

**API Gateway**

**AWS Lambda**

kafka

**Storage Layer**

**Business Decisions**

Power BI

Qlik Q

Spotfire

tableau

**Business Systems**

**End Users**

**Access and Explore Data**

# The Data: Timestamped messages with JSON encoding

```
{
    "vehicles id": {"$oid":"55a3fd0069702d5b41000000"},    Key

    "time" : {"$date":"2015-07-13T18:01:35.000Z"},    Timestamp

    "kc" : 1975.0, "kff1225" : 100.65293, "kff125a" : 110.36619, …    Values
}
```

```
{
    "vehicles_id": {"$oid":"55a3fe3569702d5c5c000020"}
    "time":{"$date":"2015-07-13T18:01:53.000Z"},
    "kc" : 2000.0, "kff1225" : 109.65293, "kff125a" : 115.36619,
      …
}
```

```
{
    "vehicles_id": {"$oid":"55a4193569702d115b000001"}
    "time":{"$date":"2015-07-12T19:04:04.000Z"}
    "kc":2200.0, "kff1225" : 112.65293, "kff125a" : 112.36619,
      …
}
```

# Access a Sample of Data

**1** Access and Explore Data

## Raw Data

| | timestamp | 1<br>value | 2<br>key |
|---|---|---|---|
| 1 | 15-Jan-2015 22:12:23 | '{ "_id" : { "$oid" : "55a41cb069702d115b059ee0" }, "trip_id" : { "$oid"... | '55a41cb069702d115b059ede' |
| 2 | 15-Jan-2015 22:12:24 | '{ "_id" : { "$oid" : "55a41cb069702d115b059ee1" }, "trip_id" : { "$oid"... | '55a41cb069702d115b059ede' |
| 3 | 15-Jan-2015 22:12:25 | '{ "_id" : { "$oid" : "55a41cb069702d115b059ee2" }, "trip_id" : { "$oid"... | '55a41cb069702d115b059ede' |
| 4 | 15-Jan-2015 22:12:26 | '{ "_id" : { "$oid" : "55a41cb069702d115b059ee3" }, "trip_id" : { "$oid"... | '55a41cb069702d115b059ede' |

✓ **Decode JSON data**
✓ **Create Timetable**

## Timetable

t = 4647×40 timetable

| | trip_id | VIN | kff1001 | kff1005 | kff1006 | kff1220 | kff1221 | kff1222 | kff1223 | kff125a |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 Sun Jul 12 16:18:41 UTC 2015 | 55a3fe356... | 55a3fe356... | 17.1000 | -84.9323 | 45.4704 | NaN | NaN | NaN | NaN | 59.0434 |
| 2 Sun Jul 12 16:18:42 UTC 2015 | 55a3fe356... | 55a3fe356... | 17.1000 | -84.9322 | 45.4704 | NaN | NaN | NaN | NaN | 57.8609 |
| 3 Sun Jul 12 16:18:43 UTC 2015 | 55a3fe356... | 55a3fe356... | 18.9000 | -84.9322 | 45.4705 | NaN | NaN | NaN | NaN | 52.7147 |
| 4 Sun Jul 12 16:18:44 UTC 2015 | 55a3fe356... | 55a3fe356... | 18.9000 | -84.9322 | 45.4705 | NaN | NaN | NaN | NaN | 51.1983 |
| 5 Sun Jul 12 16:18:45 UTC 2015 | 55a3fe356... | 55a3fe356... | 18.0000 | -84.9321 | 45.4706 | NaN | NaN | NaN | NaN | 49.1095 |
| 6 Sun Jul 12 16:19:13 UTC 2015 | 55a3fe356... | 55a3fe356... | 58.5000 | -84.9305 | 45.4686 | NaN | NaN | NaN | NaN | 73.2005 |
| 7 Sun Jul 12 16:19:14 UTC 2015 | 55a3fe356... | 55a3fe356... | 56.7000 | -84.9304 | 45.4685 | NaN | NaN | NaN | NaN | 75.3612 |
| 8 Sun Jul 12 16:19:15 UTC 2015 | 55a3fe356... | 55a3fe356... | 57.6000 | -84.9304 | 45.4683 | NaN | NaN | NaN | NaN | 70.7542 |
| 9 Sun Jul 12 16:19:16 UTC 2015 | 55a3fe356... | 55a3fe356... | 56.7000 | -84.9303 | 45.4682 | NaN | NaN | NaN | NaN | 62.8340 |

**②**
**Preprocess Data**

# Develop a Preprocessing Function

## Timetable

| t = | 4647×40 | *timetable* | | | | | | | | |

| | trip_id | VIN | kff1001 | kff1005 | kff1006 | kff1220 | kff1221 | kff1222 | kff1223 | kff125a |
|---|---------|-----|---------|---------|---------|---------|---------|---------|---------|---------|
| 1 Sun Jul 12 16:18:41 UTC 2015 | 55a3fe356… | 55a3fe356… | 17.1000 | -84.9323 | 45.4704 | NaN | NaN | NaN | NaN | 59.0434 |
| 2 Sun Jul 12 16:18:42 UTC 2015 | 55a3fe356… | 55a3fe356… | 17.1000 | -84.9322 | 45.4704 | NaN | NaN | NaN | NaN | 57.8609 |
| 3 Sun Jul 12 16:18:43 UTC 2015 | 55a3fe356… | 55a3fe356… | 18.9000 | -84.9322 | 45.4705 | NaN | NaN | NaN | NaN | 52.7147 |
| 4 Sun Jul 12 16:18:44 UTC 2015 | 55a3fe356… | 55a3fe356… | 18.9000 | -84.9322 | 45.4705 | NaN | NaN | NaN | NaN | 51.1983 |
| 5 Sun Jul 12 16:18:45 UTC 2015 | 55a3fe356… | 55a3fe356… | 18.0000 | -84.9321 | 45.4706 | NaN | NaN | NaN | NaN | 49.1095 |
| 6 Sun Jul 12 16:19:13 UTC 2015 | 55a3fe356… | 55a3fe356… | 58.5000 | -84.9305 | 45.46?? | NaN | NaN | NaN | NaN | 73.?005 |
| 7 Sun Jul 12 16:19:14 UTC 2015 | 55a3fe356… | 55a3fe356… | 56.7000 | -84.9304 | 45.46? | | | | | |
| 8 Sun Jul 12 16:19:15 UTC 2015 | 55a3fe356… | 55a3fe356… | 57.6000 | -84.9304 | 45.46? | | | | | |
| 9 Sun Jul 12 16:19:16 UTC 2015 | 55a3fe356… | 55a3fe356… | 56.7000 | -84.9303 | 45.46? | | | | | |

✓ **Clean up**
✓ **Enrich**
✓ **Restructure**

**Preprocess data**

```
t = sortrows(t);
t = rmmissing(t,'MinNumMissing',width(t)-2);
```

**Perform windowed calculations**

```
t.Speed = movmedian(t.SpeedGPS,3);
t.D1 = [0;diff(t.SpeedGPS)];
```

```
[tmin,tmax] = bounds(t.time);
tnew = tmin:seconds(10):tmax;
countsByTime = retime(t(:,'Event'),tnew,@histcounts);
```

# ① Access and Explore Data

# Ad Hoc Access to Data from MATLAB



**athenaQuery.mlx**

## Access the data in S3

### Bring up the AthenaClient

```
athenaClient = aws.athena.Client();
athenaClient.Database = 'trainingdata';
athenaClient.initialize();
```

### Create a query and submit

```
athenaClient.submitQuery('SELECT * FROM "trainingdata"."sampledata" limit 100','s3://fleettrainingdata')
```

### Fetch data as a table for easy analysis

```
ds = datastore('s3://fleettrainingdata/*.csv');
ds.NumHeaderLines = 2;
data = table(ds);
```

### Your usual MATLAB workflow goes here

AWS S3

AWS Athena Service

MATLAB

# Everything you need to develop a predictive model is found in MATLAB

**3** Develop Predictive Models

**Label Events**

**Represent Signals**

**Train Model**

**Validate Model**

**Scale Up**

**Develop Predictive Models**

# Develop a Predictive Model in MATLAB

# A quick Intro to Stream Processing

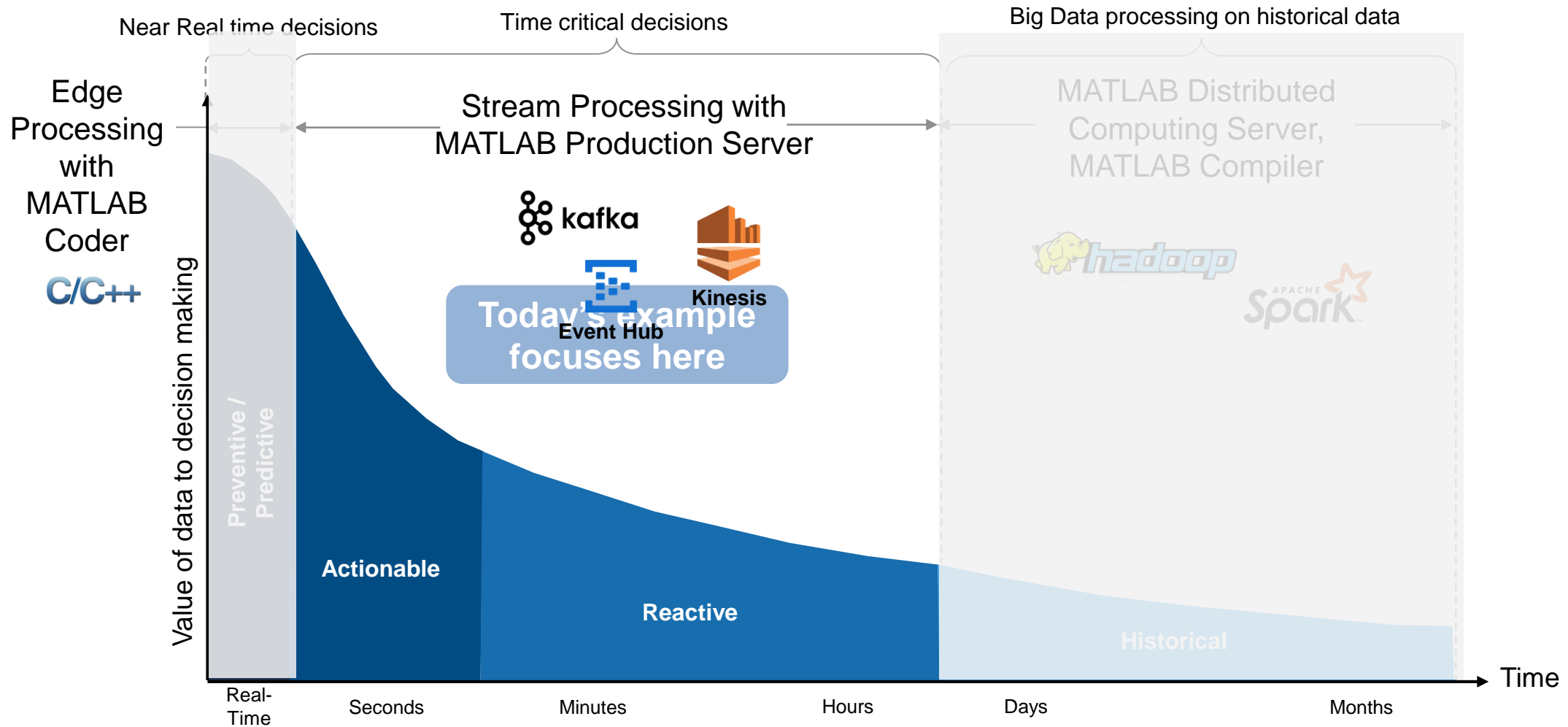- ***Batch Processing*** applies computation to a finite sized historical data set that was acquired in the past

**Historical Data**  **Configure Resources**  **Schedule and Run Job**  **Output Data**

Files

Storage

Files

Storage

- Reporting
- Data Exploration
- Training Models

- ***Stream Processing*** applies computation to an unbounded data set that is produced continuously

**Continuous Data**  **Messaging Service**  **Stream Analytics**

Connected Devices

$f(x)$

Dashboards

Alerts

Storage

- Reporting
- Real Time Decision Support

MATLAB EXPO 2018

# Streaming data is treated as an unbounded Timetable

## Input Table

| Event Time | Vehicle | RPM | Torque | Fuel Flow |
|---|---|---|---|---|
| 18:01:10 | 55a3fd | 1975 | 100 | 110 |
| 18:10:30 | 55a3fe | 2000 | 109 | 115 |
| 18:05:20 | 55a3fd | 1980 | 105 | 105 |
| 18:10:45 | 55a3fd | 2100 | 110 | 100 |
| 18:30:10 | 55a419 | 2000 | 100 | 110 |
| 18:35:20 | 55a419 | 1960 | 103 | 105 |
| 18:20:40 | 55a3fe | 1970 | 112 | 104 |
| 18:39:30 | 55a419 | 2100 | 105 | 110 |
| 18:30:00 | 55a3fe | 1980 | 110 | 113 |
| 18:30:50 | 55a3fe | 2000 | 100 | 110 |
| … | … | … | … | … |

State

MATLAB Function

State

MATLAB Function

State

MATLAB Function

State

## Output Table

| Time window | | Vehicle | Score |
|---|---|---|---|
| … | | … | … |
| 18:00:00 | 18:10:00 | 55a3fd | 5 |
| | | 55a3fe | … |
| | | 55a419 | … |
| 18:10:00 | 18:20:00 | 55a3fd | 7 |
| | | 55a3fe | 3 |
| | | 55a419 | … |
| 18:20:00 | 18:30:00 | 55a3fd | … |
| | | 55a3fe | 4 |
| | | 55a419 | … |
| 18:30:00 | 18:40:00 | 55a3fd | … |
| | | 55a3fe | 5 |
| | | 55a419 | 8 |

**4** Integrate with Production Systems

# Introducing MATLAB Production Server

## Data

### Databases

DynamoDB

Microsoft SQL Server

mongoDB

Cassandra

Cosmos DB

### Cloud Storage

Azure Blob

S3

### Streaming

AWS Kinesis

OSIsoft PI System

kafka

MQTT ORG

Azure IoT Hub

## Analytics

**MATLAB Production Server**

Request Broker

## Business System

### Dashboards

Qlik Q

tableau

Microsoft Power BI

Spotfire

### Web

Microsoft IIS

Apache Tomcat

WebSphere

### Custom Apps

## Platform

Google Cloud Platform

Azure

amazon web services

rackspace

openstack

vmware

# Connecting MATLAB Production Server to Kafka

- Kafka client for MATLAB Production Server feeds topics to functions deployed on the server

- Configurable batch of messages passed as a MATLAB Timetable

- Each consumer process feeds one topic to a specified function

- Drive everything from a simple config file
  - No programming outside of MATLAB!

**Publisher**   **Publisher**   **Publisher**

**Kafka Cluster**

**Topic-0**   **Topic-1**

Partition   Partition

Partition   Partition

Partition   Partition

**Consumer Process feeds Topic-0**

**Async Java Client**

**Consumer Process feeds Topic-1**

**Async Java Client**

MATLAB Production Server

Request Broker & Program Manager

## 4 Integrate with Production Systems

# Develop and Deploy a Stream Processing Function



**Edge Devices**

**Production System**

amazon web services™

**API Gateway**

**AWS Lambda**

kafka

**Kafka Connector**

**MATLAB Production Server**

MATLAB Analytics

**Storage Layer**

**Analytics Development**

**MATLAB Compiler SDK**

**MATLAB**

**Algorithm Developers**

**Business Decisions**

Power BI

Qlik Q

Spotfire

+ableau

**Business Systems**

**End Users**

# Develop a Stream Processing Function in MATLAB



**Process each window of data as it arrives**

**Current score**

**Previous state**

**Current window of data to be processed**

# Develop a Stream Processing Function in MATLAB

calculateScores.mlx ✕ +

## Develop a Streaming Function

```matlab
function new_state = calculateScores(car_id, current_data
```

**Preprocess and perform calculations**

```matlab
current_data = preprocessData(current_data);
```

**Predict driving events**

```matlab
current_data = predictEvents(current_data);
```

**Count events for each ten second window**

```matlab
countsByTime = countEvents(current_data);
```

**Write discrete data to mongodb**

```matlab
updateResultsStore(car_id,countsByTime,resultsStore);
```

**Update new state**

```matlab
new_state = updateState(countsByTime,old_state);
end
```
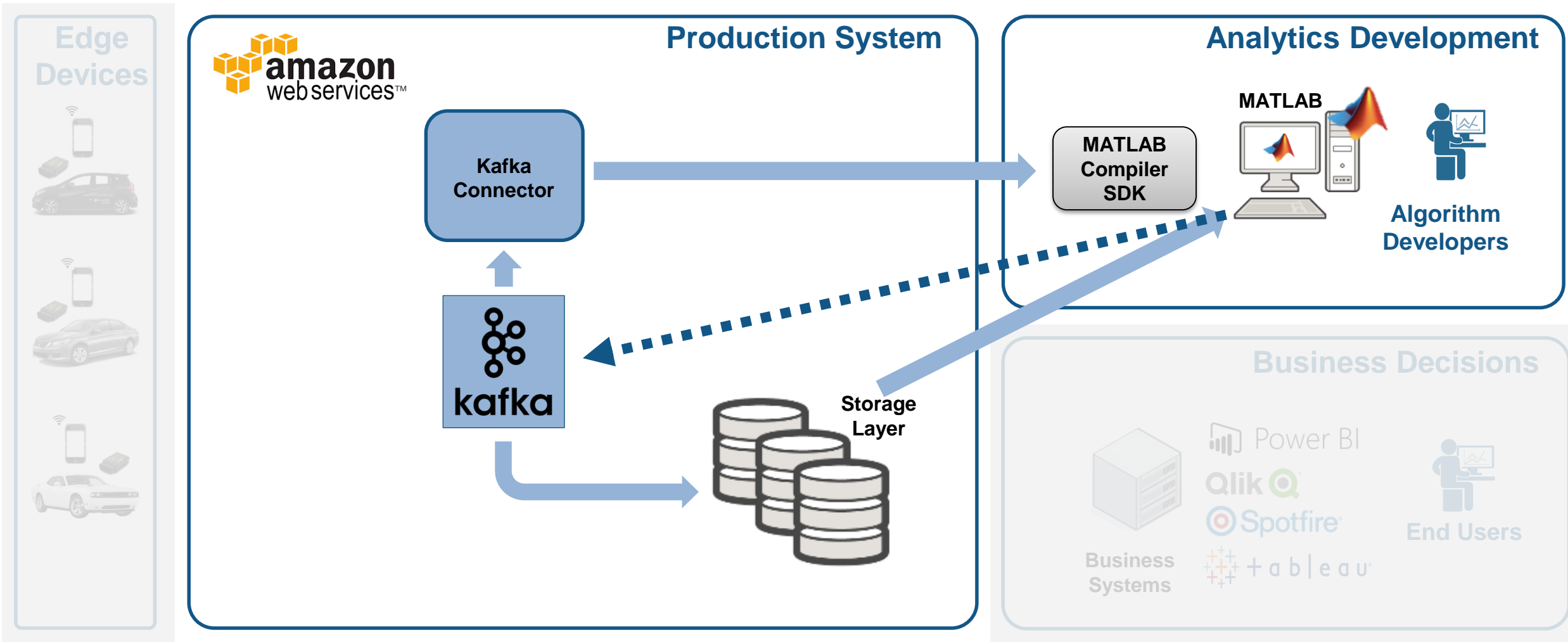
```matlab
function current_data = preprocessData(current_data)
% Preprocess and perform calculations

% Remove records with all missing data
current_data = rmmissing(current_data,'MinNumMissing',width(current_data)-1);

% Smooth and calculate approximate gradients
current_data.Speed = movmedian(current_data.kff1001,5);
current_data.D1 = [0;diff(current_data.kff1001)];
current_data.D2 = [0;0;diff(current_data.kff1001,2)];
```

**Apply your pre-processing algorithm**

# Develop a Stream Processing Function in MATLAB



**Use the model you created with Classification Learner App**

# Develop a Stream Processing Function in MATLAB

**calculateScores.mlx** ✕ +

## Develop a Streaming Function

```
function new_state = calculateScores(car_id, current_data, old_state, resultsStore)
```

**Preprocess and perform calculations**

```
current_data = preprocessData(current_data);
```

**Predict driving events**

```
current_data = predictEvents(current_data);
```

**Count events for each ten second window**

```
countsByTime = countEvents(current_data);
```

**Write discrete data to mongodb**

```
updateResultsStore(car_id,countsByTime,resultsStore);
```

**Update new state**

```
new_state = updateState(countsByTime,old_state);
end
```

**Update Mongo database**
- **Count of events by type and location**
- **Results of driver scoring**

# Debug a Stream Processing Function in MATLAB

MathWorks

**Edge Devices**

**Production System**

amazon web services™

**Kafka Connector**

kafka

**Storage Layer**

**Analytics Development**

**MATLAB Compiler SDK**

**MATLAB**

**Algorithm Developers**

**Business Decisions**

Power BI

Qlik Q

Spotfire

+ableau

**Business Systems**

**End Users**

# Debug a Stream Processing Function in MATLAB

**4** Integrate with Production Systems

# Tie in your Dashboard Application

Production System

Analytics Development

Edge Devices

amazon web services™

**MATLAB Production Server**

MATLAB Analytics

MATLAB Compiler SDK

MATLAB

Algorithm Developers

**Kafka Connector**

API Gateway

AWS Lambda

kafka

Storage Layer

**Business Decisions**

Power BI

Qlik Q

Spotfire

+ableau

Business Systems

End Users

**Visualize Results**

# Complete Your Application

# Key Takeaways

➢ MATLAB connects directly to your data so you can quickly design and validate algorithms

➢ The MATLAB language and apps enable fast design iterations

➢ MATLAB Production Server enables easy integration of your MATLAB algorithms with enterprise production systems

➢ You to spend your time understanding the data and designing algorithms

# Resources to learn and get started

- [Data Analytics with MATLAB](#)

- [MATLAB Production Server](#)
- [MATLAB Compiler SDK](#)
- [Statistics and Machine Learning Toolbox](#)
- [Database Toolbox](#)
- [Mapping Toolbox](#)

- [MATLAB with TIBCO Spotfire](#)
- [MATLAB with Tableau](#)
- [MATLAB with MongoDB](#)