

# MATLAB EXPO 2018

## 高可扩展的 MATLAB 云端工程 数据分析

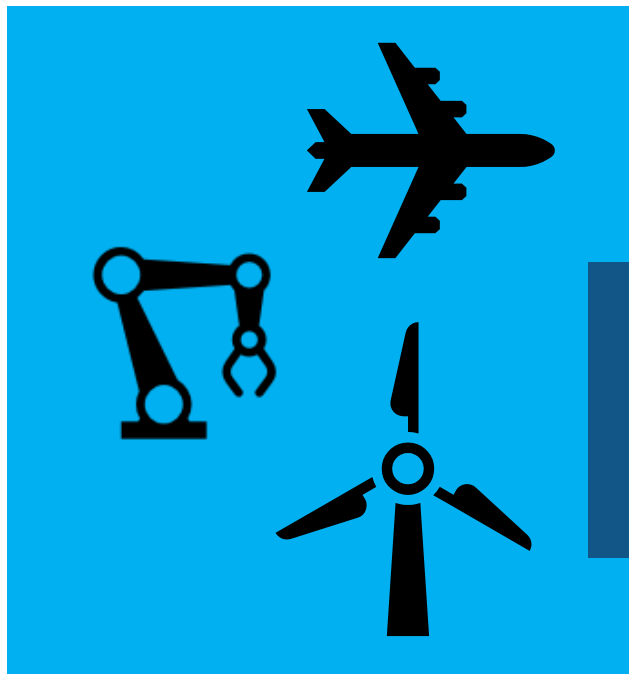
陈建平, MathWorks 中国



# 大纲



# 大规模流处理的需求



## 预测性维护

增加操作效率  
减少计划外的停机

更多的应用需要近乎实时的分析

喷气引擎: ~800TB 每天  
涡轮引擎: ~ 2TB 每天

MATLAB EXPO 2018

## 医疗设备

患者的安全  
更加积极的治疗结果

## 车辆互联

安全  
维护  
先进的驾驶功能



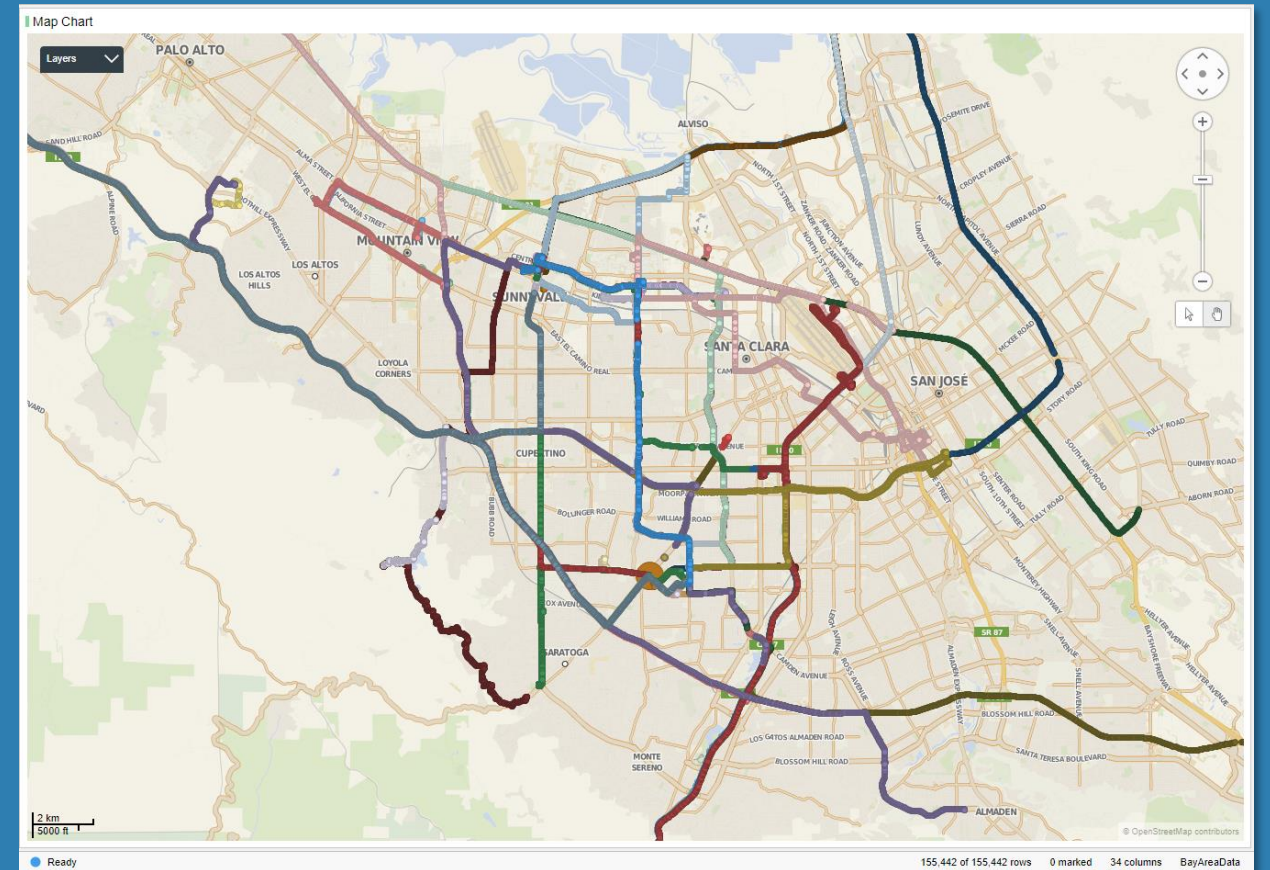
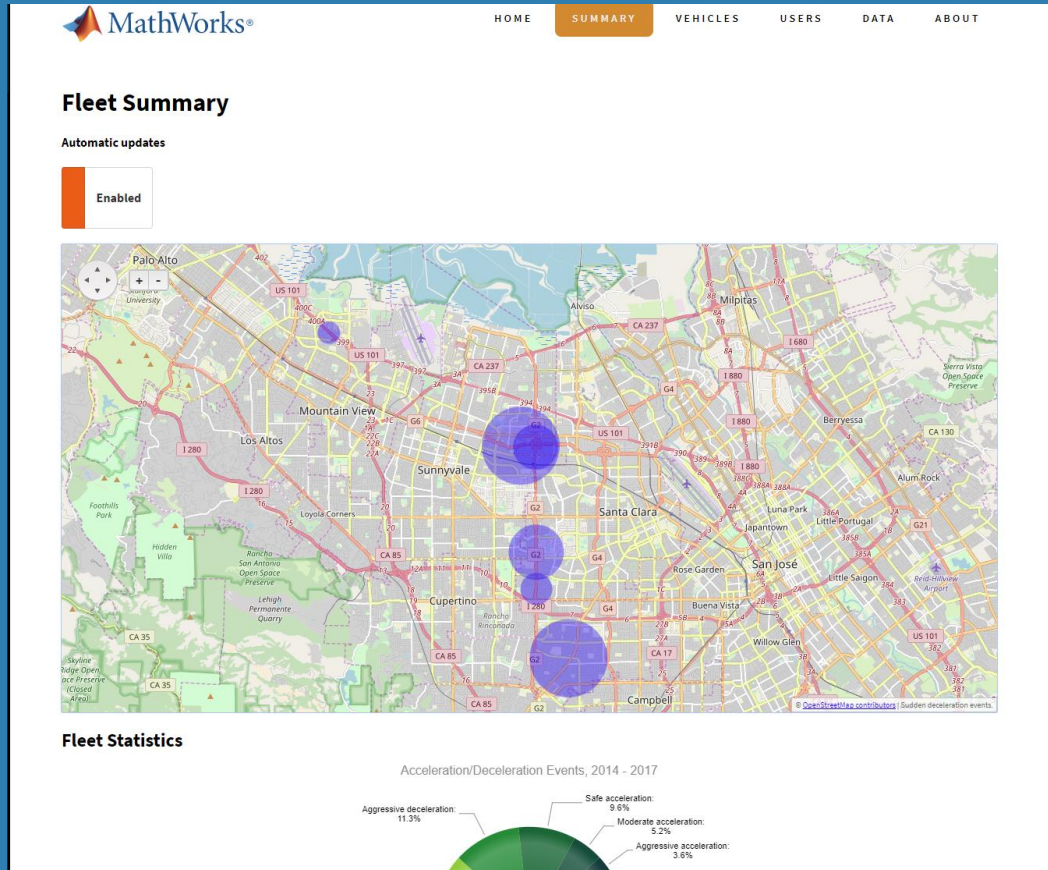
车辆: ~25 GB 每天

## 实例问题 —— 我的驾驶习惯如何？

- 一组 MathWorks 员工在车内安装了一个 OBD 狗，用于监控车载系统
- 流式数据传输到云端，汇总并存储
- 我们想用这些数据来评估参与者的驾驶习惯

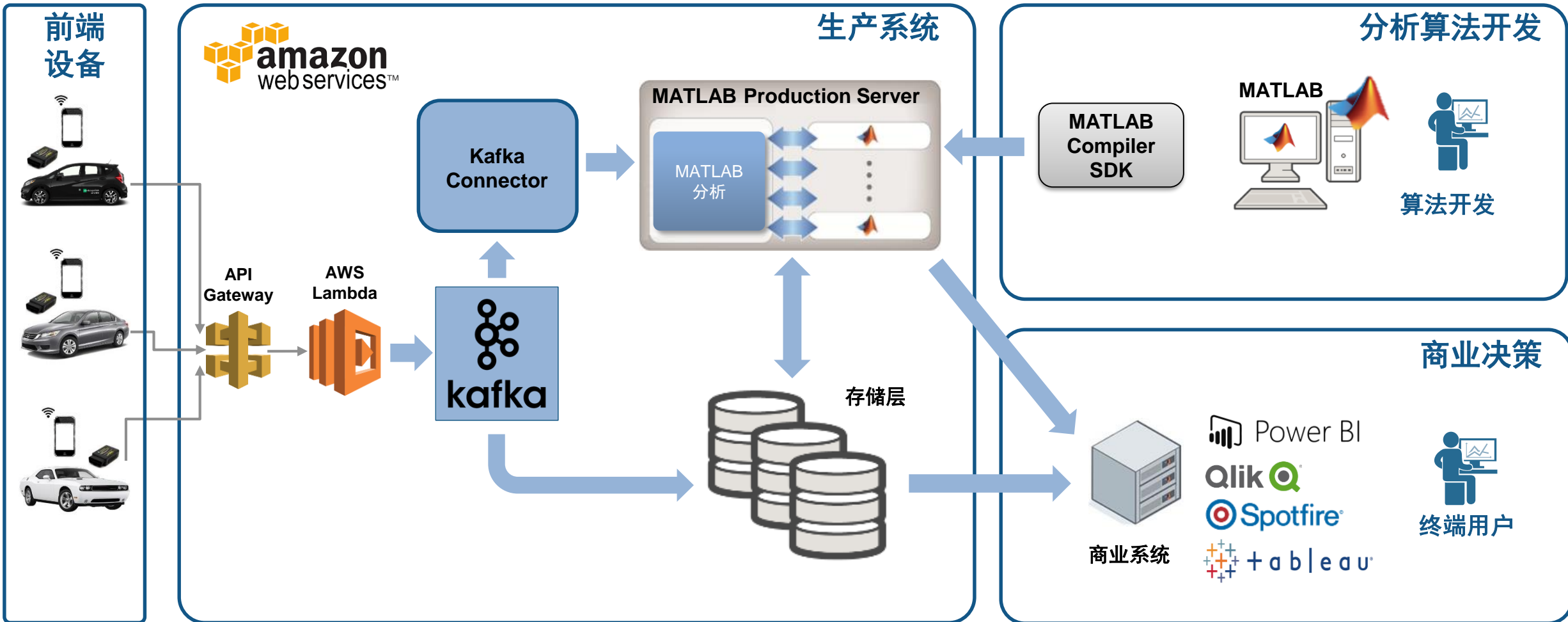


# 示例：MATLAB 车队数据分析





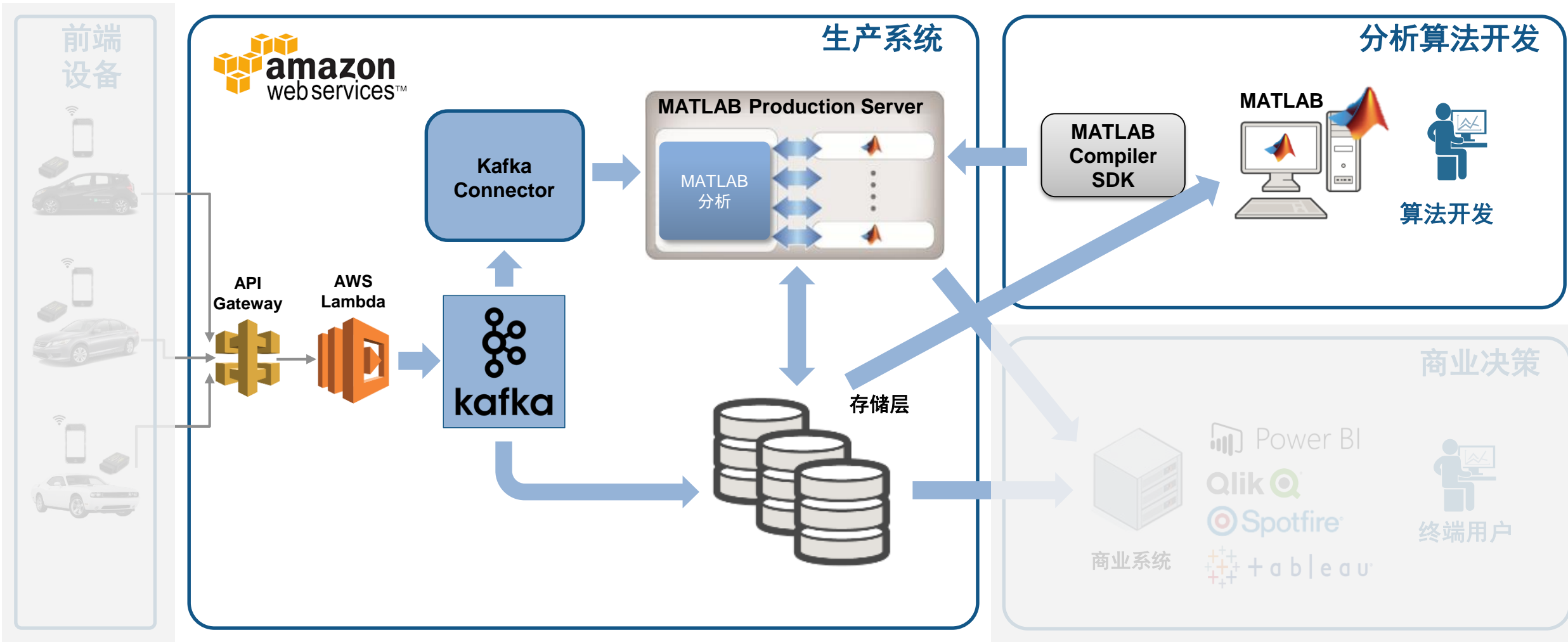
# 车队数据分析架构



1

数据访问和探索

# 输入数据清洗



1

数据访问和探索

## 数据：带时间戳的 JSON 编码消息



```

{
  "vehicles id": {"$oid":"55a3fd0069702d5b4100000"}, Key
  "time" : {"$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,
  ...
}

```



1

数据访问和探索

## 访问数据样点

原始数据

	timestamp	1 value	2 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'

- ✓ JSON 译码
- ✓ 创建 Timetable

Timetable

t = 4647x40 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

数据预处理

## 开发一个预处理函数

## Timetable

t = 4647x40 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	72.2005
7 Sun Jul 12 16:19:14 UTC 2015	55a3fe356...	55a3fe356...	56.7000	-84.9304	45.4686	NaN	NaN	NaN	NaN	72.2005
8 Sun Jul 12 16:19:15 UTC 2015	55a3fe356...	55a3fe356...	57.6000	-84.9304	45.4686	NaN	NaN	NaN	NaN	72.2005
9 Sun Jul 12 16:19:16 UTC 2015	55a3fe356...	55a3fe356...	56.7000	-84.9303	45.4686	NaN	NaN	NaN	NaN	72.2005

## 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);
```

- ✓ 清洗
- ✓ 增强
- ✓ 重建

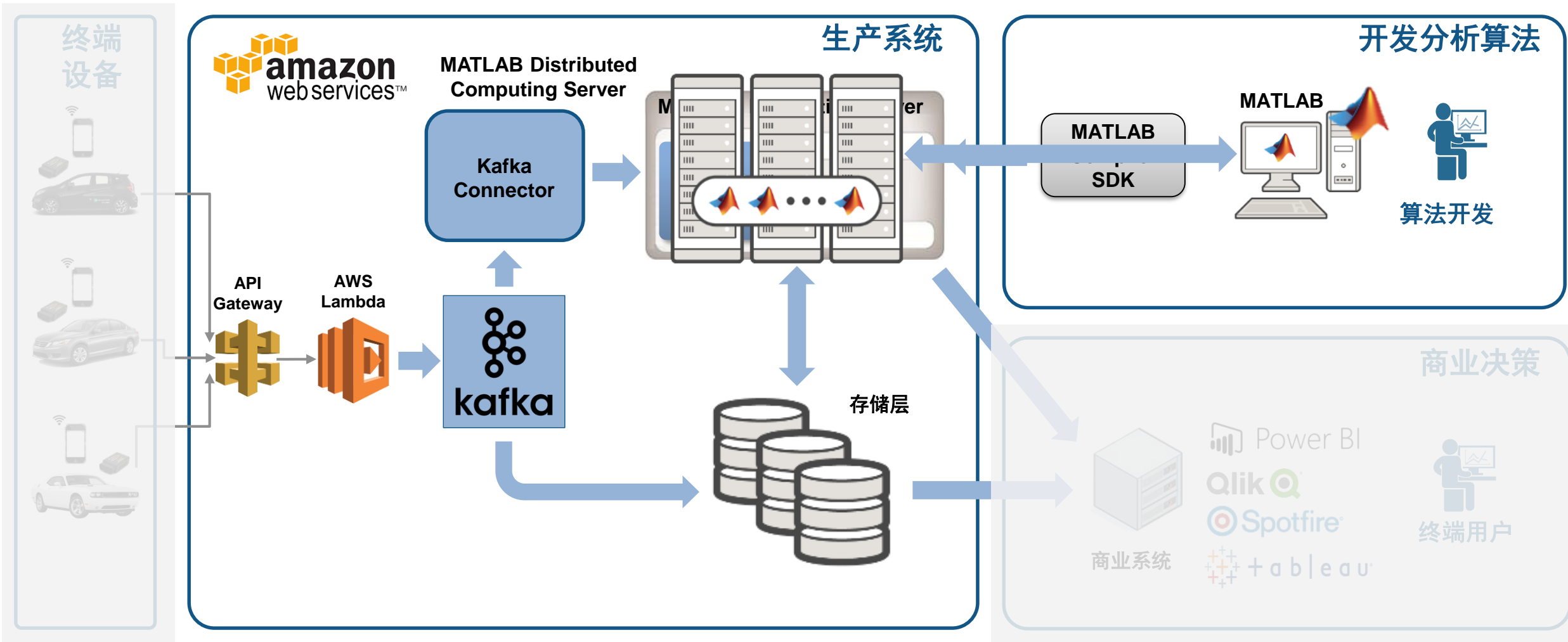
## MATLAB 中随意访问数据

```
athenaQuery.mlx x +  
  
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
```

3

开发预测模型

# 开发预测模型



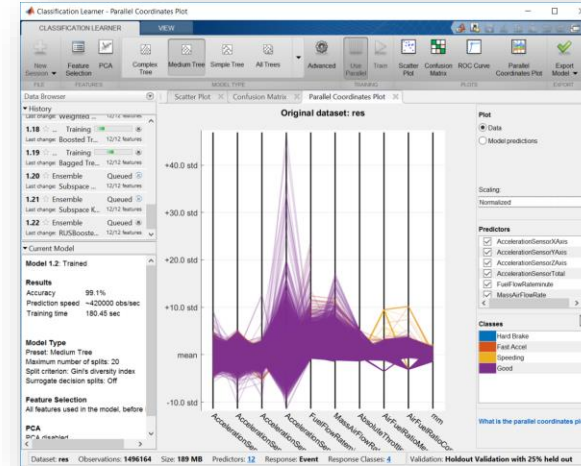
3

开发预测模型

# MATLAB 开发预测模型

time	1 Event	2 SpeedGPS	3 AccelerationSensorXAxis	4 AccelerationSensorYAxis	5 AccelerationSensorZAxis
Mon May 11 04:03:15 UTC 2015	Hard Brake	10.8360	-0.6996	0.6014	0.205
Wed May 06 19:09:48 UTC 2015	Hard Brake	27.8280	0.1419	0.9035	-0.526
Sun May 17 17:09:19 UTC 2015	Hard Brake	6.5520	0.9986	-0.0761	-0.004
Fri Jan 16 20:38:37 UTC 2015	Hard Brake	39.6128	0.0999	0.8000	0.367
Sat May 02 14:00:37 UTC 2015	Hard Brake	61.1280	0.4006	-0.4022	0.663
Mon Apr 27 17:54:27 UTC 2015	Fast Accel	37.7640	0.1527	0.4666	0.857
Sun May 03 21:00:42 UTC 2015	Fast Accel	17.2440	1.0235	0.0815	0.304
Mon May 04 11:30:33 UTC 2015	Fast Accel	19.6560	0.1336	0.8932	-0.578
Wed May 20 10:20:55 UTC 2015	Fast Accel	22.4000	0.2050	0.0054	0.006

标记事件



信号展示

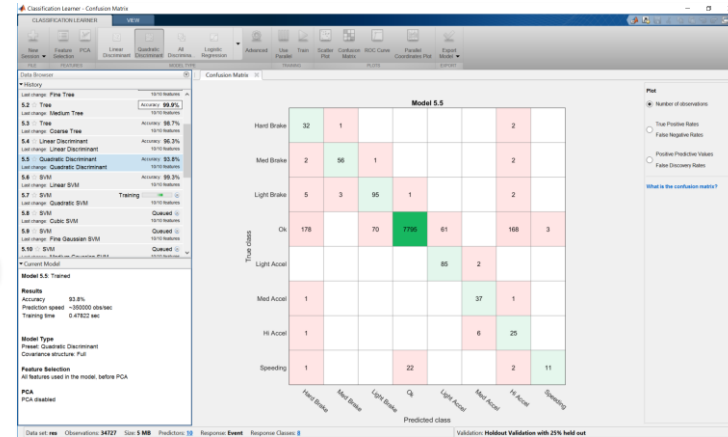
Evaluating tall expression using the Spark Cluster:  
 - Pass 1 of 2: Completed in 11 sec  
 - Pass 2 of 2: Completed in 2.3333 min  
 Evaluation completed in 2.6167 min

```

Scale up
tt = tall(data); % test tall array
model = TreeBagger(50,tt,'Event');

Scale to out of memory data
tt = tall(ds);
tt = preprocessData(tt);
model = TreeBagger(50,tt,'Event');
save machineLearningModel model
    
```

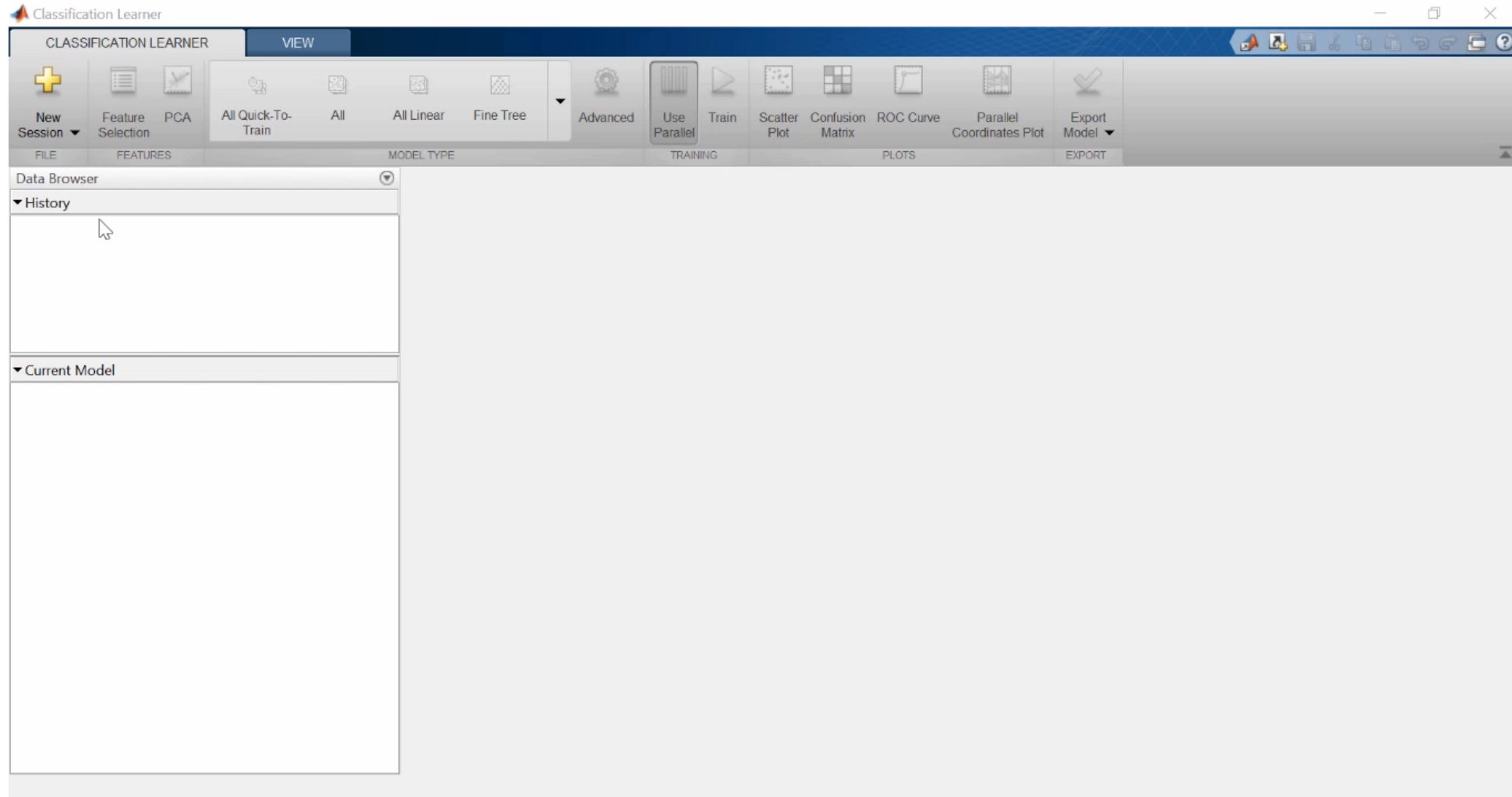
扩展



训练模型

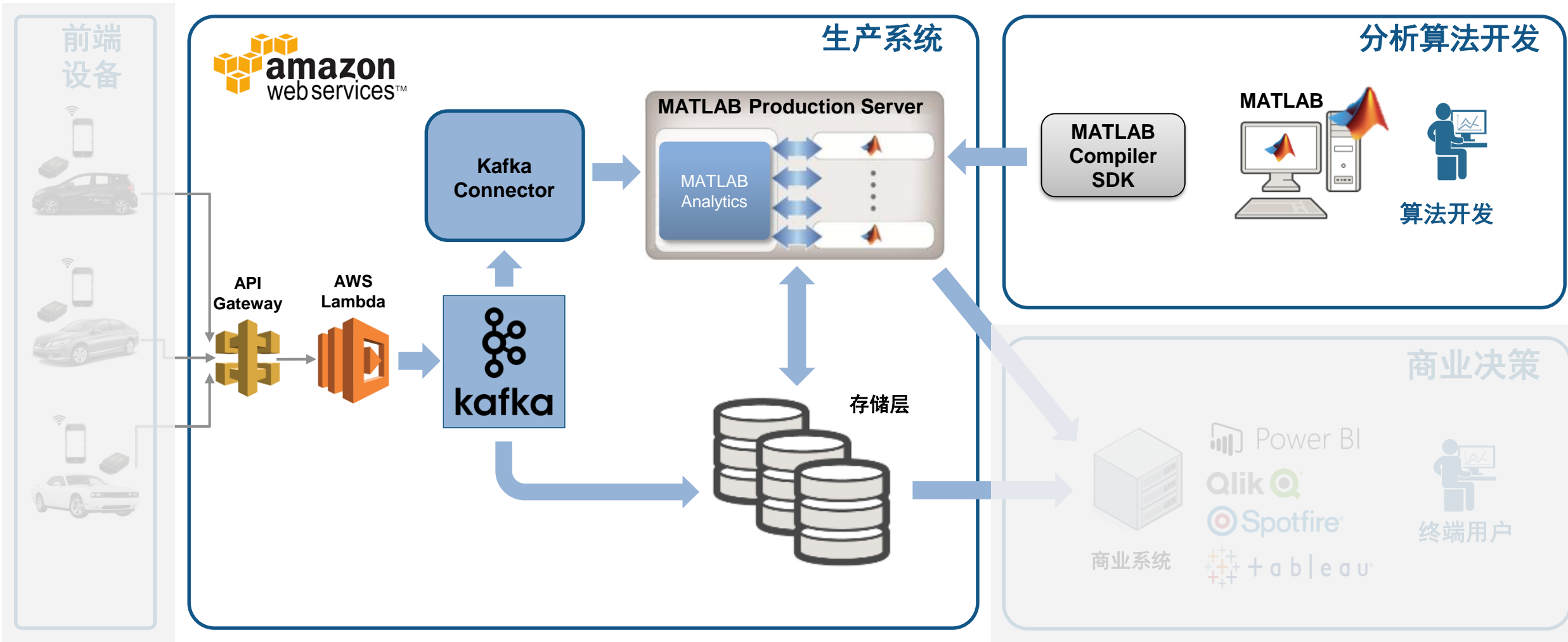
验证模型

## 在 MATLAB 中开发预测模型





# 在生产系统中集成分析算法



## 流处理快速入门

- **批处理:** 应用计算于*过去*采集的*有限大小*的历史数据集

历史数据



Files



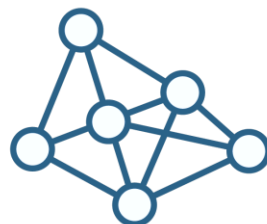
Storage



资源设置



调度和任务计算



输出数据



Files



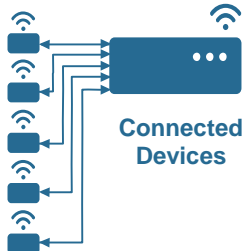
Storage



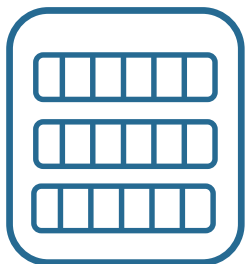
- 报告
- 数据探索
- 训练模型

- **流处理:** 应用计算于*连续*产生的*无限*数据集

连续数据

Connected  
Devices

消息服务



流式分析

 $f(x)$ 

Dashboards



Alerts



Storage

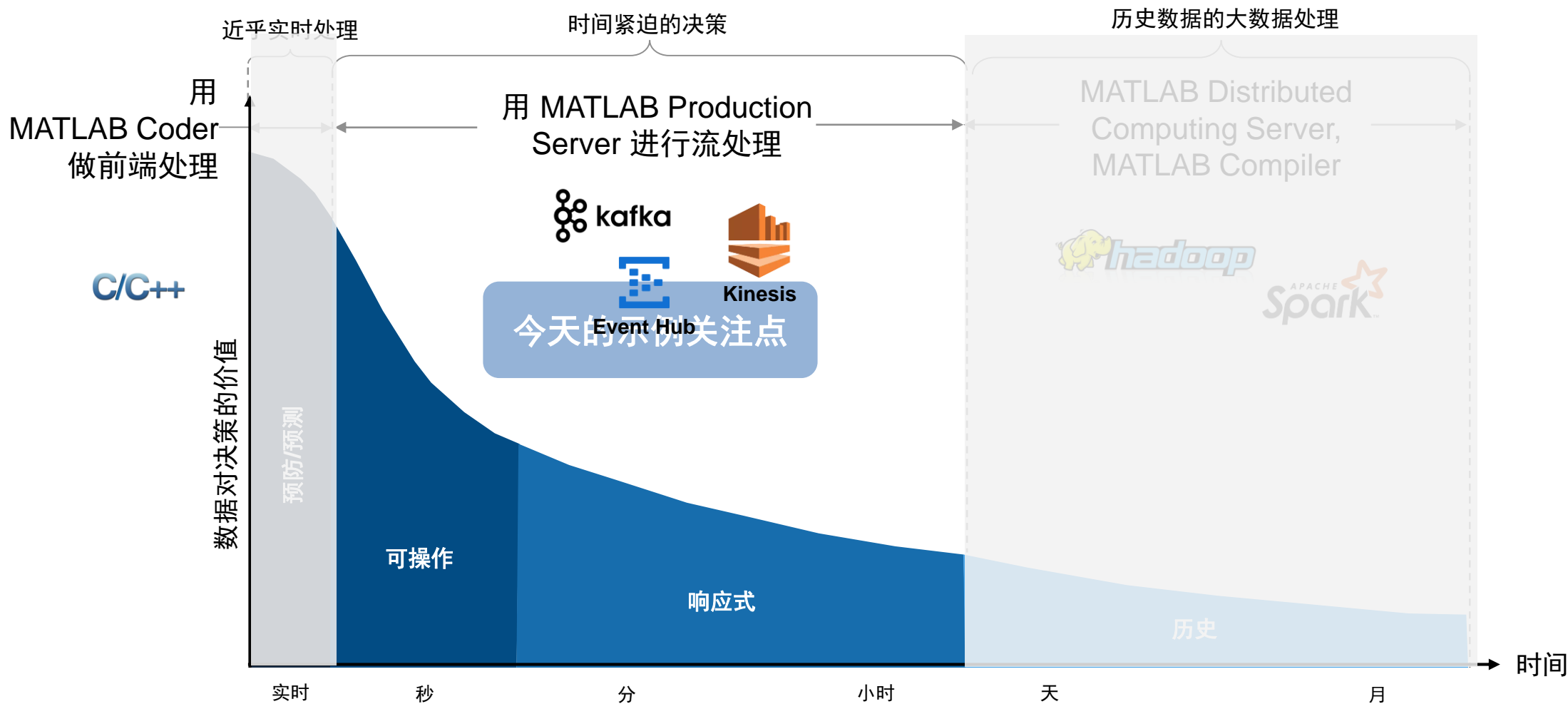


- 报告
- 实时
- 决策支持

4

生产系统集成

# 流处理



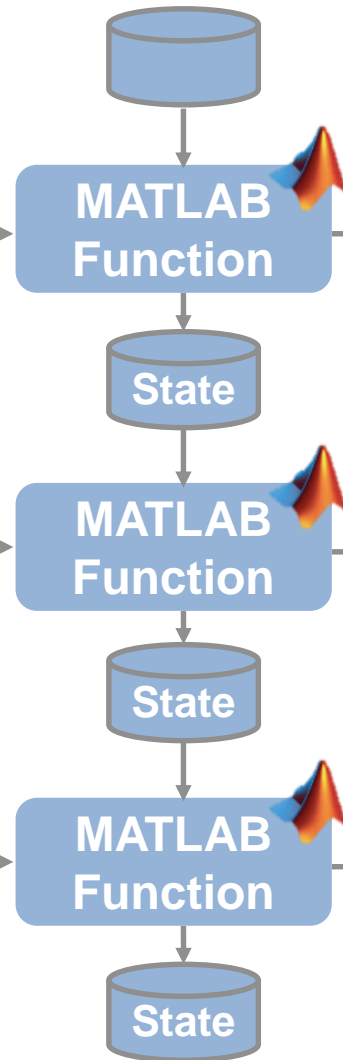
# 流数据被视为无限的时间表

输入表

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
...	...	...	...	...

输出表

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

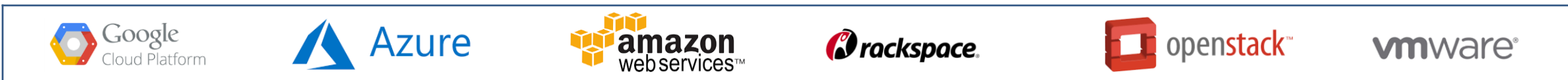
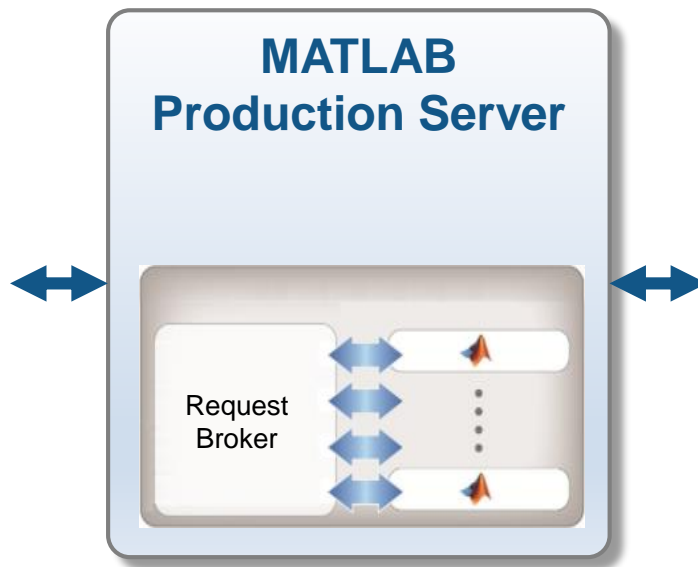
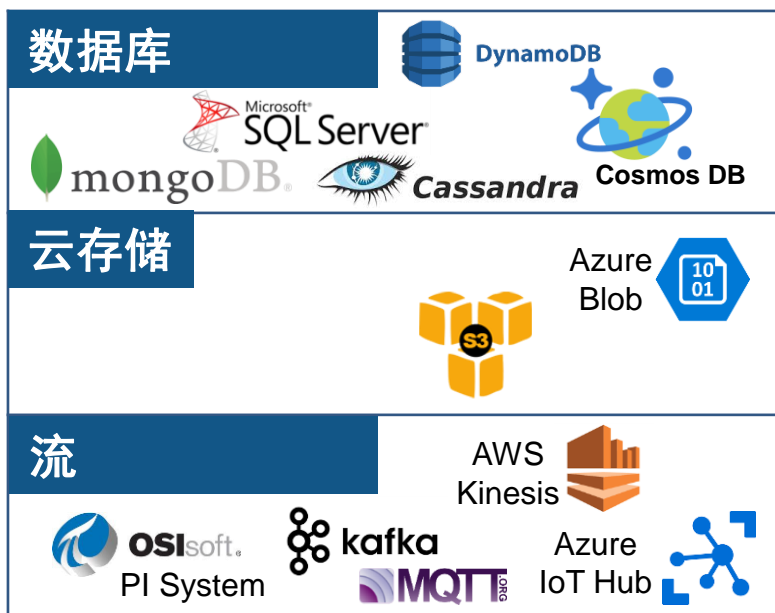


4 生产系统集成 **MATLAB Production Server 介绍**

数据

分析

商业决策

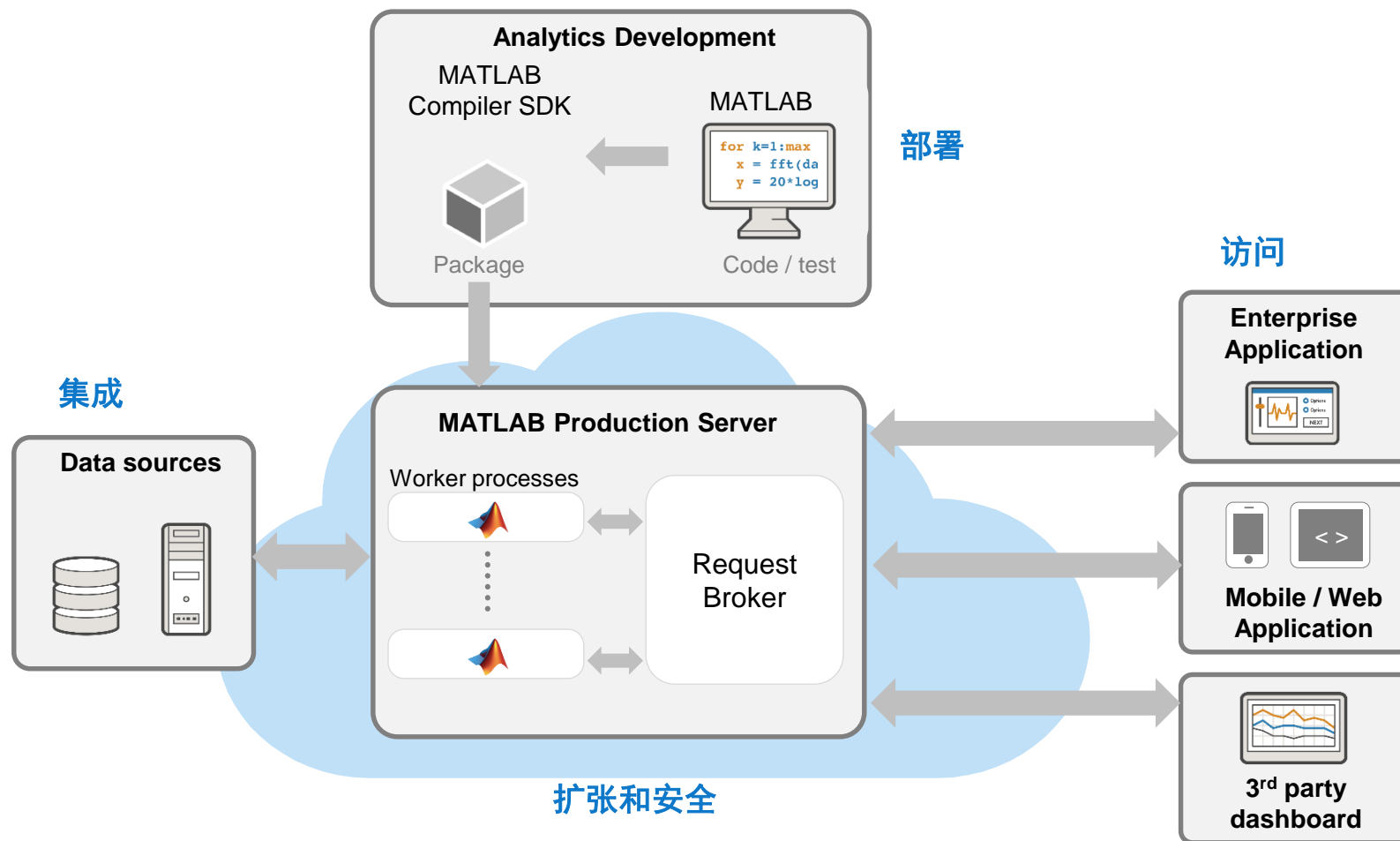


平台

4

生产系统集成

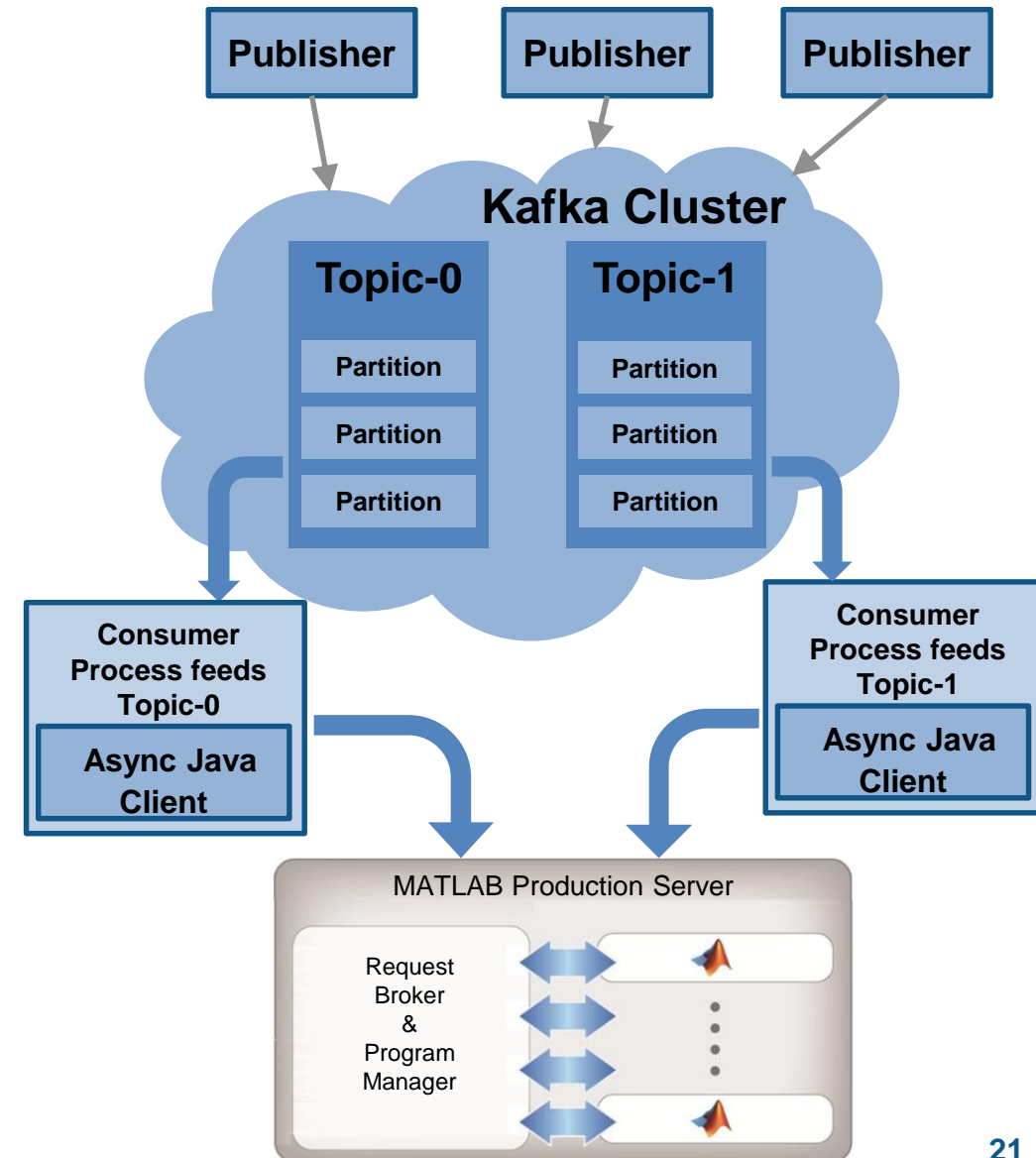
# MATLAB Production Server: MATLAB 代码发布成 API 的应用服务器



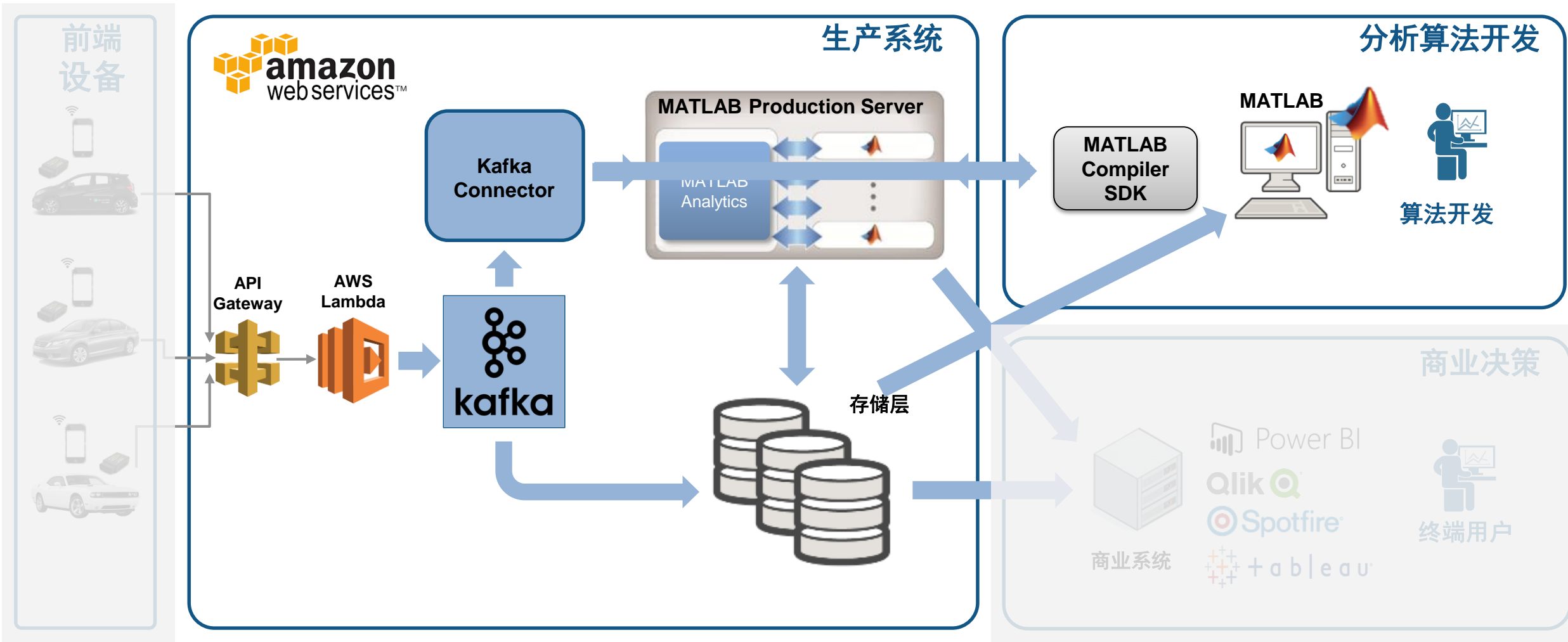


## 连接 MATLAB Production Server 到 Kafka

- Kafka 客户端把数据发送给部署在服务器中的函数
- 可以把一批消息配置成 MATLAB 时间表传递
- 每个消费者进程把一个主题送给指定的 MATLAB 函数
- 从一个简单的配置程序驱动
  - 除了 MATLAB 无需其他编程



# 开发和部署流处理函数



## 在 MATLAB 中开发流处理函数

数据到达时处理每一个数据窗

```
calculateScores.mlx x +
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
```

处理结果

前一个状态

当前待处理数据

## 在 MATLAB 中开发流处理函数

```
calculateScores.mlx x +  
  
Develop a Streaming Function  
function new_state = calculateScores(car_id, current_data)  
  
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
```

```
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)];
```

应用预处理算法

## 在 MATLAB 中开发流处理函数

使用分类学习 App 习得的模型

```
calculateScores.mlx x +
```

**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
```

```
function current_data = predictEvents(current_data)
% Predict events for current data based on machine learning model
predictorNames = {'kff1005', 'kff1006', 'kff125a', 'k10', 'kff1249', 'Speed', 'D1', 'D2', ...
                  'kff1001', 'kff1220', 'kff1221', 'kff1222', 'kff1223', ...
                  'k47', 'kff124d'};
predictors = current_data(:, predictorNames);
mdl = load('machineLearningModel.mat');
current_data.Event = predict(mdl.model, predictors);
end
```

## 在 MATLAB 中开发流处理函数

```
calculateScores.mlx x +
```

### 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

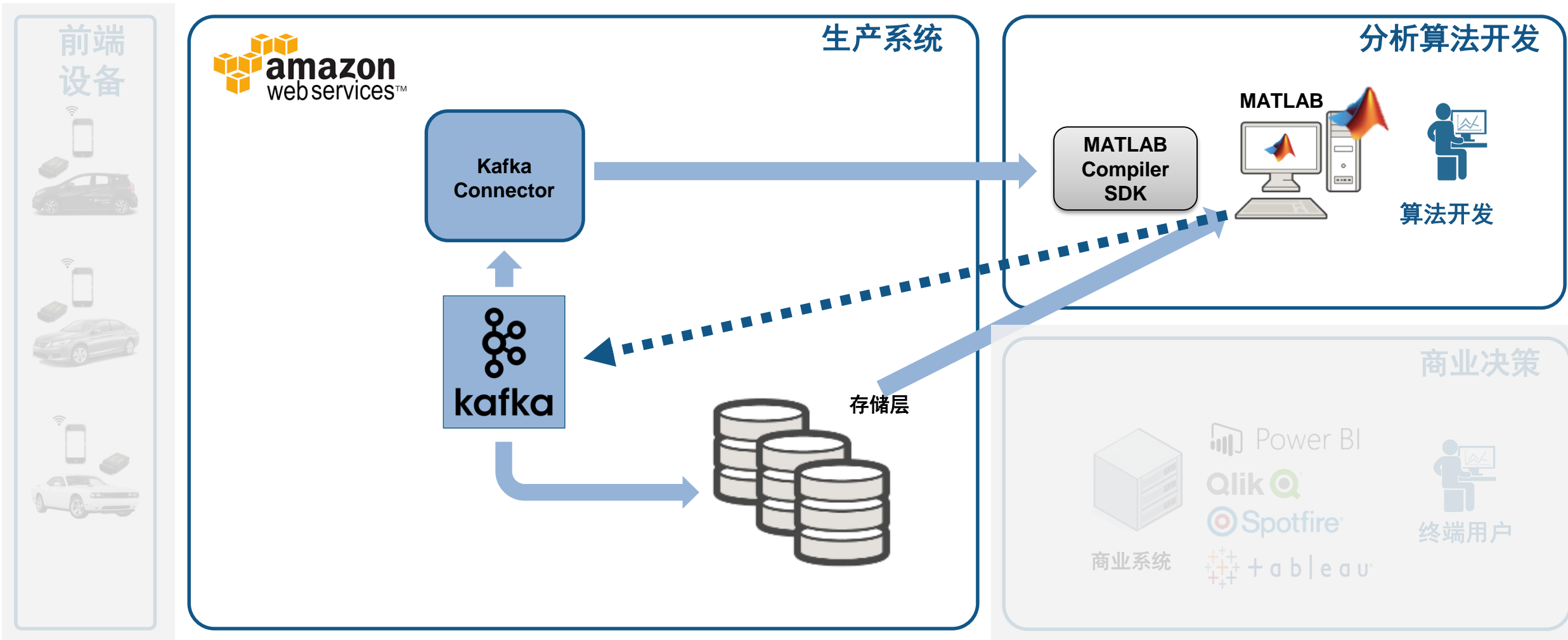
```
new_state = updateState(countsByTime, old_state);  
end
```

更新 Mongo 数据库

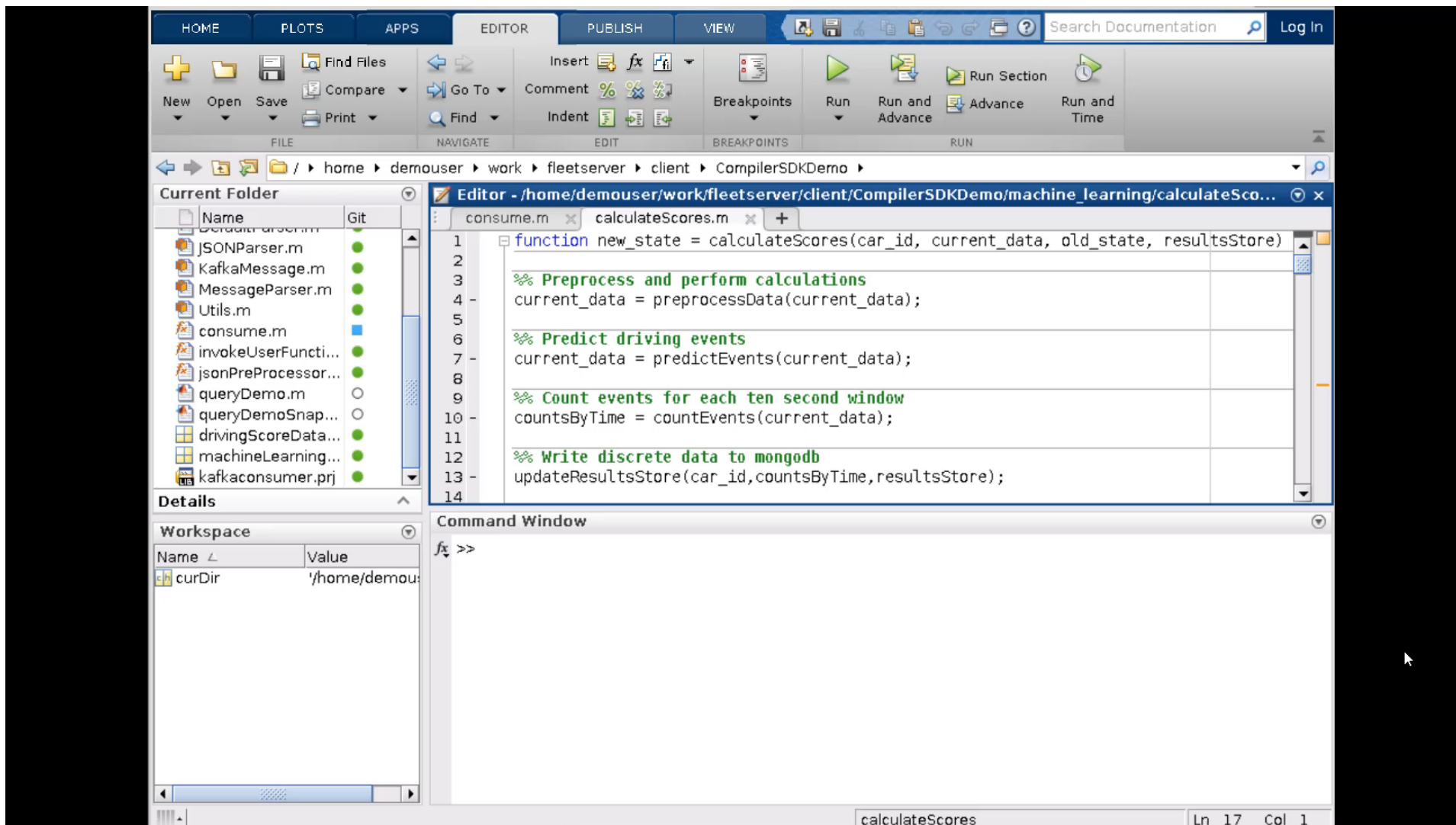
- 针对类型和位置统计事件
- 驾驶习惯估计



# 在 MATLAB 中调试流处理函数



## 在 MATLAB 中调试流处理函数



The screenshot displays the MATLAB IDE interface. The Editor window shows the following code:

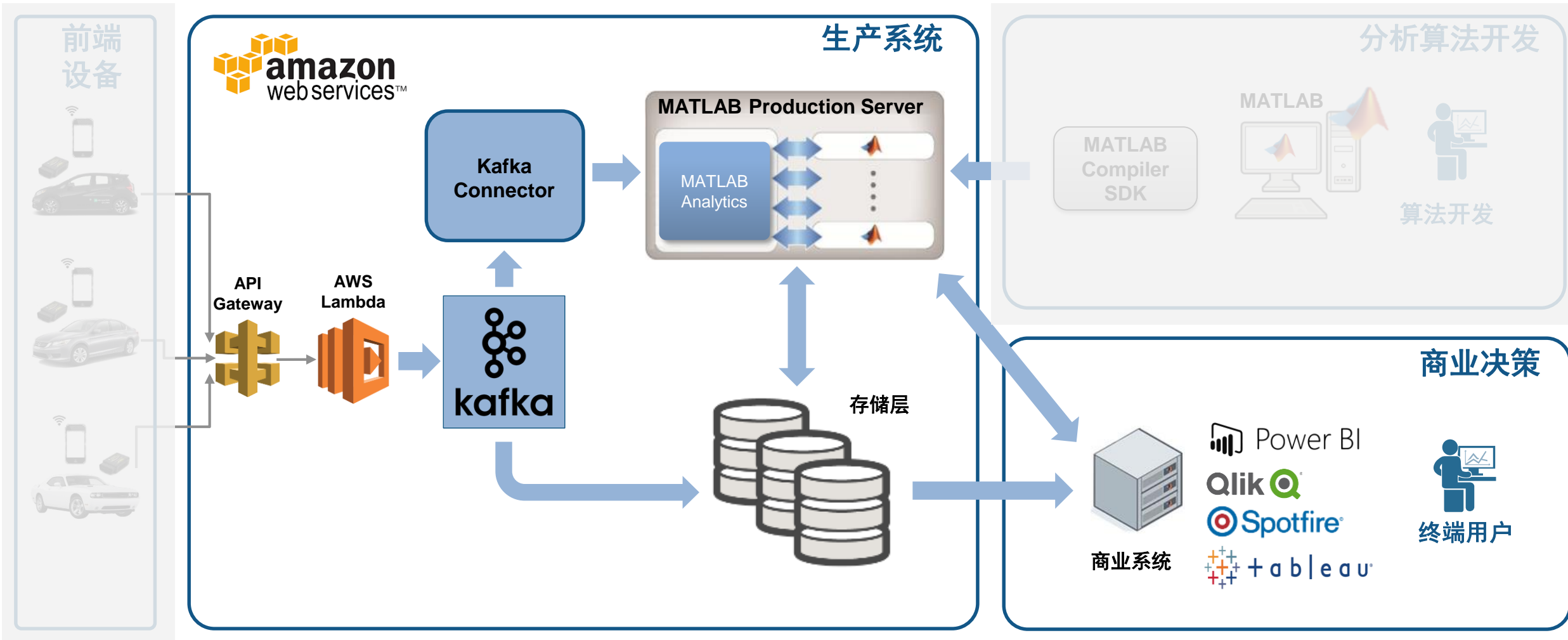
```
1 function new_state = calculateScores(car_id, current_data, old_state, resultsStore)
2
3 %% Preprocess and perform calculations
4 current_data = preprocessData(current_data);
5
6 %% Predict driving events
7 current_data = predictEvents(current_data);
8
9 %% Count events for each ten second window
10 countsByTime = countEvents(current_data);
11
12 %% Write discrete data to mongodb
13 updateResultsStore(car_id, countsByTime, resultsStore);
14
```

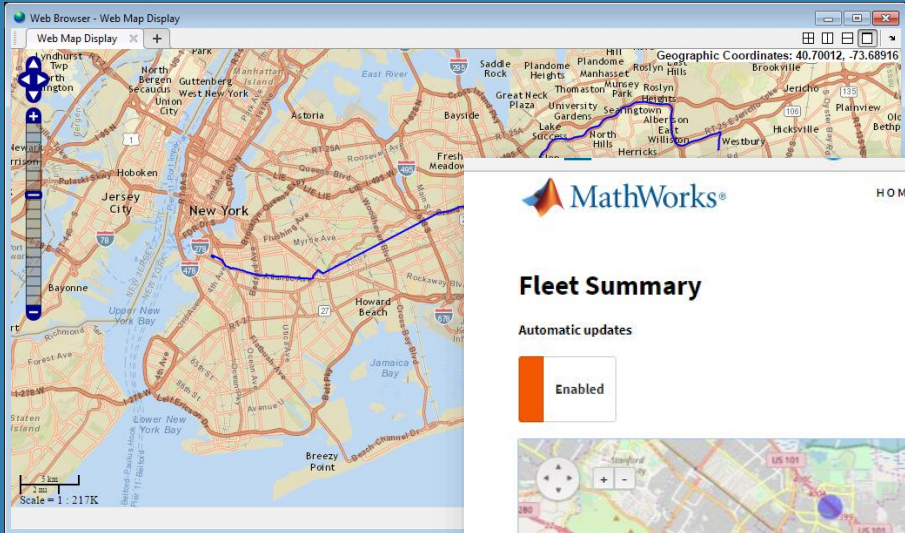
The Command Window shows the prompt `fx >>`. The Workspace shows the current directory `curDir` as `'/home/demou...`.

4

生产系统集成

# 绑定您的仪表盘应用

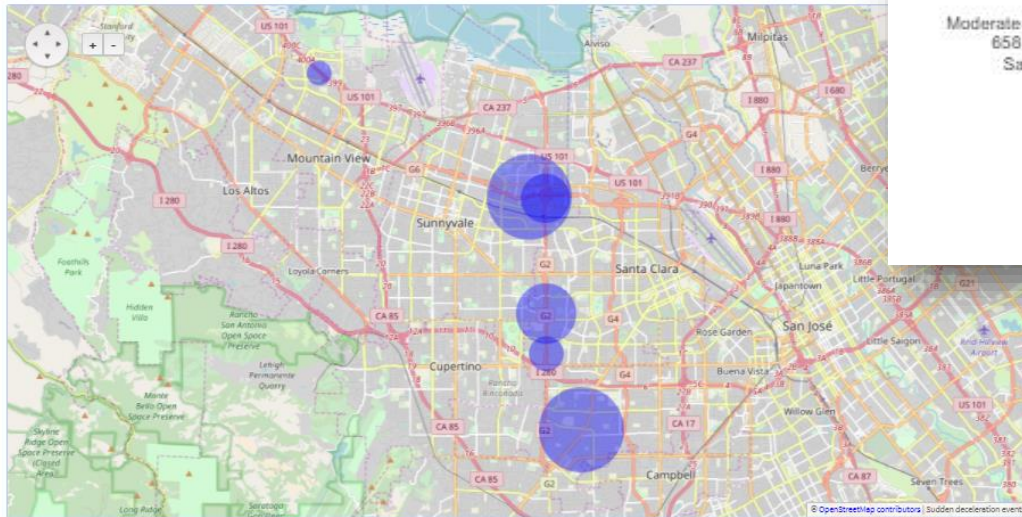




HOME SUMMARY VEHICLES USERS TRIPS REPORT

### Fleet Summary

Automatic updates

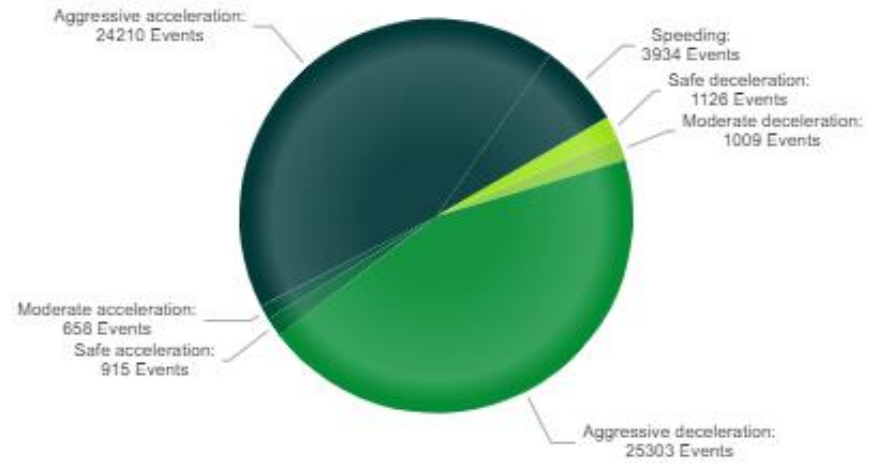


### Fleet Statistics

Total Events:

193351

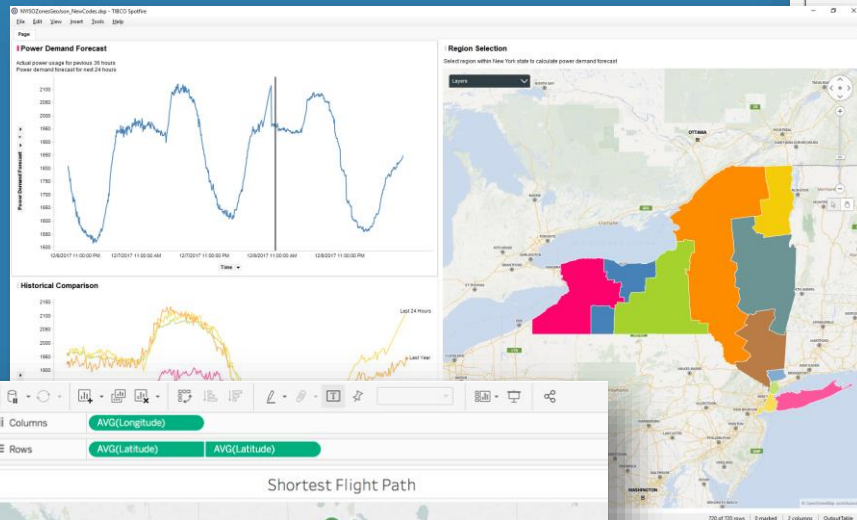
Acceleration/Deceleration Events, 2014 - 2017



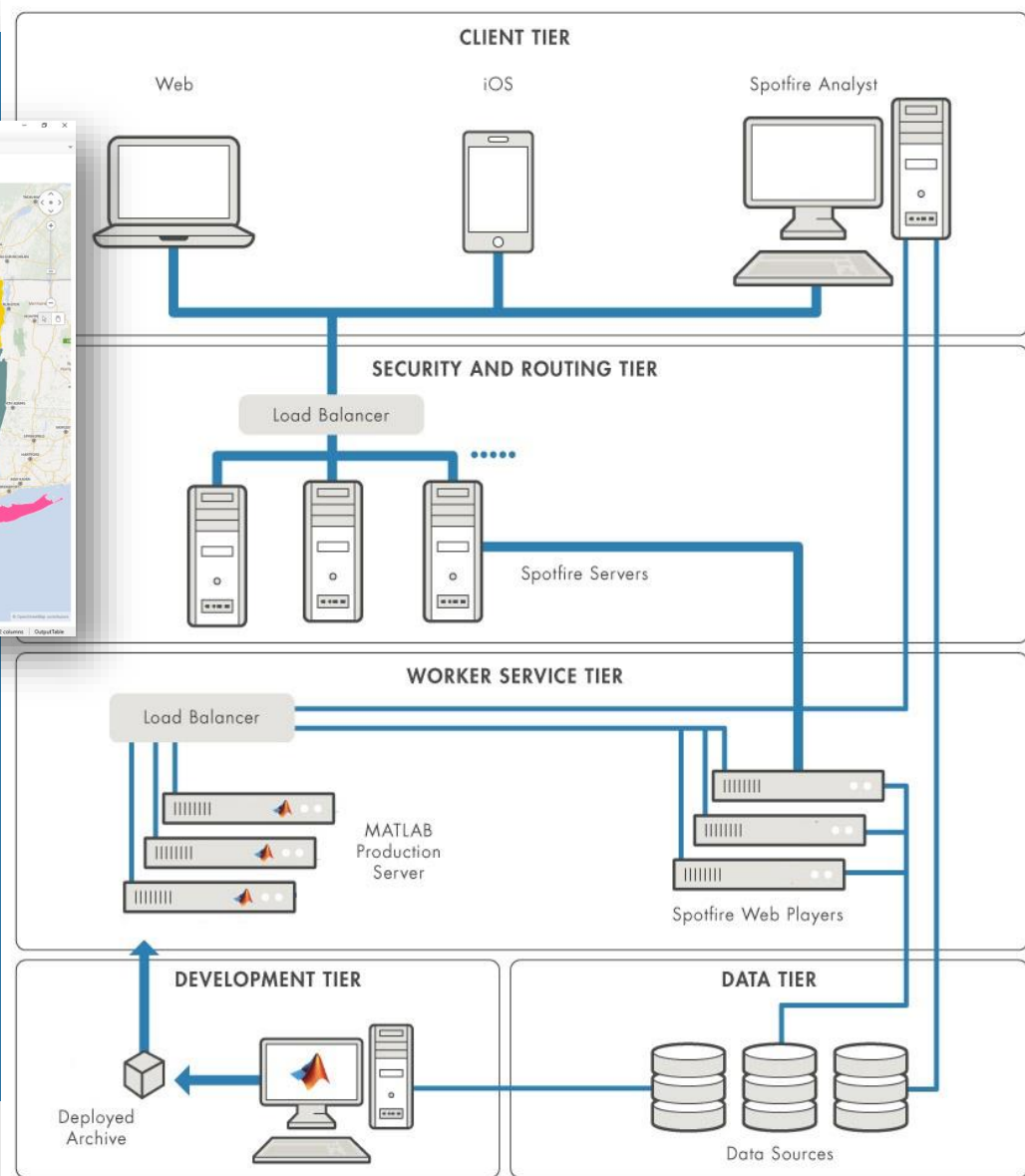
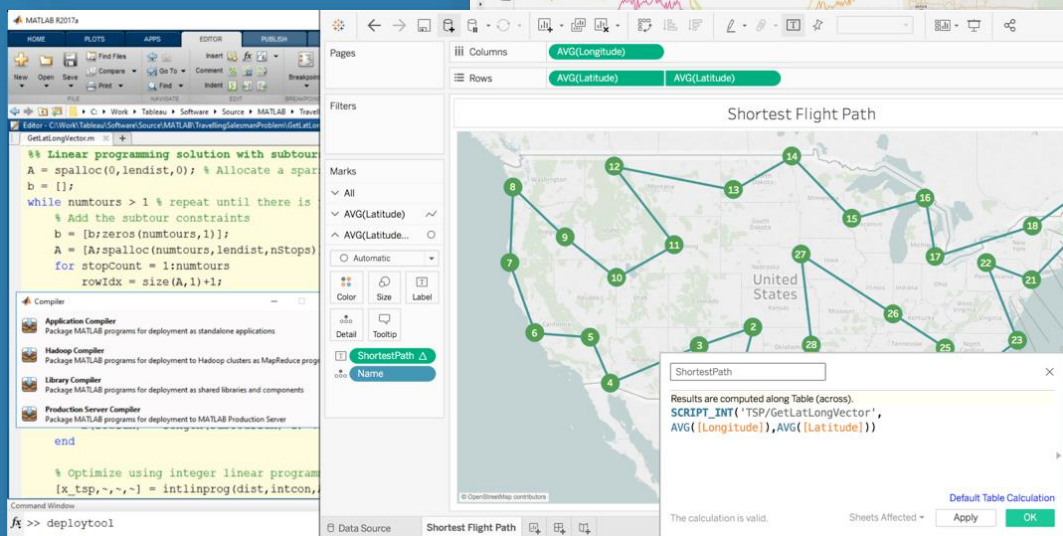


# 协同企业 BI 工具的扩展分析

## TIBCO Spotfire



## Tableau



## 要点回顾

- MATLAB 直接连接到您的数据之中，方便快速设计和验证算法
- MATLAB 语言和 apps 支持快速设计迭代
- MATLAB 产品服务器可以轻松将算法与企业生产系统集成
- 您的时间将聚焦于数据理解和算法设计

# 学习资源和快速起步

- [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](#)

The screenshot shows a MathWorks webpage titled "Reference Architecture" for "Scalable Analytics with TIBCO Spotfire and MATLAB Production Server". The page includes a navigation menu with "Products", "Solutions", "Academia", "Support", "Community", and "Events". A search bar is present with the text "Search MathWorks.com".

The main content area features the title "Scalable Analytics with TIBCO Spotfire and MATLAB Production Server" and a subtitle "Resources and Spotfire extension to scale MATLAB analytics for use with Spotfire applications".

Below the text, there are two call-to-action buttons: "Request the Extension & Getting Started Guide" (with a "Submit request" button) and "Download the Free Technical Brief" (with a "Download now" button).

On the right side, there is a detailed architecture diagram. It shows a "TIBCO SPOTFIRE SERVER" at the bottom, which connects to a "TIBCO SPOTFIRE WEB PLAYER" (containing "Mathworks MPSExtension"). This web player is accessed by "MOBILE", "WEB", and "DESKTOP" clients (the desktop client also includes "Mathworks MPSExtension"). A "Load Balancer" sits between the clients and the web player. The web player connects to a "MATLAB PRODUCTION SERVER" (containing "MATLAB Analytics") via another "Load Balancer".