

CLUDERA

Open Source Data Management for Car-to-Cloud

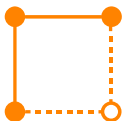
Michael Ger

Managing Director, Manufacturing & Automotive

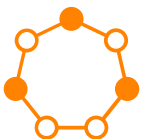
Cludera



CLOUDERA AT A GLANCE



One stop shop for Big Data management for analytics



Unified open source architecture



Hybrid and multi-cloud

The graphic displays the Cloudera SDX logo, which stands for "shared data experience". Above the logo are five icons representing different data management stages: Ingest & Streaming, Operational Database, Data Warehouse, Data Engineering, and Data Science. Below the logo are three logos for major cloud providers: Microsoft Azure, Amazon Web Services, and Google Cloud Platform.

CLOUDERA ENTERPRISE DATA WAREHOUSE

Any Data. Anywhere. From Edge to AI

Traditional Data Warehouse Optimization



- BI Reporting
- Data Warehouse Offload
- ETL Optimization

Operations & Events Data Warehouse



- IOT Data (Sensors, Logs, etc.)
- Connected Vehicles
- Connected Factories (Industry 4.0)
- Etc.

Research & Discovery Data Warehouse



- Unstructured Data (Images, Video, etc.)
- Intelligent Search
- Autonomous Vehicle Research

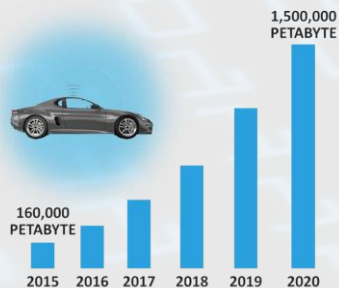
DID YOU KNOW?



10 of the top 10 Automotive OEMs
Trust Cloudera for Big Data-
Enabled Digital Transformations

DATA TRENDS

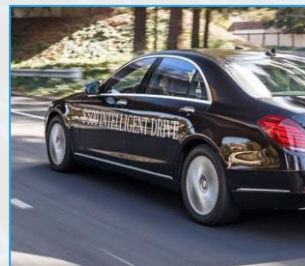
Connected Vehicle Data



10X by 2020*

* Source: Cowen and Company, Gartner

Autonomous Vehicle Data



**TERABYTES
per vehicle,
per day**

Industrial Internet Data



**2X faster
than any
other data
source**

* Source: Wikibon

Real-Time Data



**1.5X faster
than
traditional
data**

* Source: IDC

CAR TO CLOUD – DATA MANAGEMENT ON 5 LEVELS

THE EDGE



MANAGING THE EDGE



ENTERPRISE FLOW/STREAM ANALYTICS



EDGE TO AI



CONNECTED COMMUNITIES



CAR TO CLOUD – DATA MANAGEMENT ON 5 LEVELS

THE EDGE



MANAGING THE EDGE



ENTERPRISE FLOW/STREAM ANALYTICS



EDGE TO AI



CONNECTED COMMUNITIES



EDGE DATA COLLECTION - MINIFI AGENTS

- Extremely small footprint
- Java and C++ agents
 - JRHEL/CentOS, Debian/Ubuntu, Android*
- Key data management functions
 - Filtering, buffering, guaranteed delivery, prioritized queuing
- Secure
 - Encryption, certificate-based authentication
 - Data tagging and provenance
- Execute Machine Learning (ML) models including Tensorflow




GOING SMALLER – FROM NIFI TO MINIFI

Supporting the Need for Smaller Footprints on the Edge



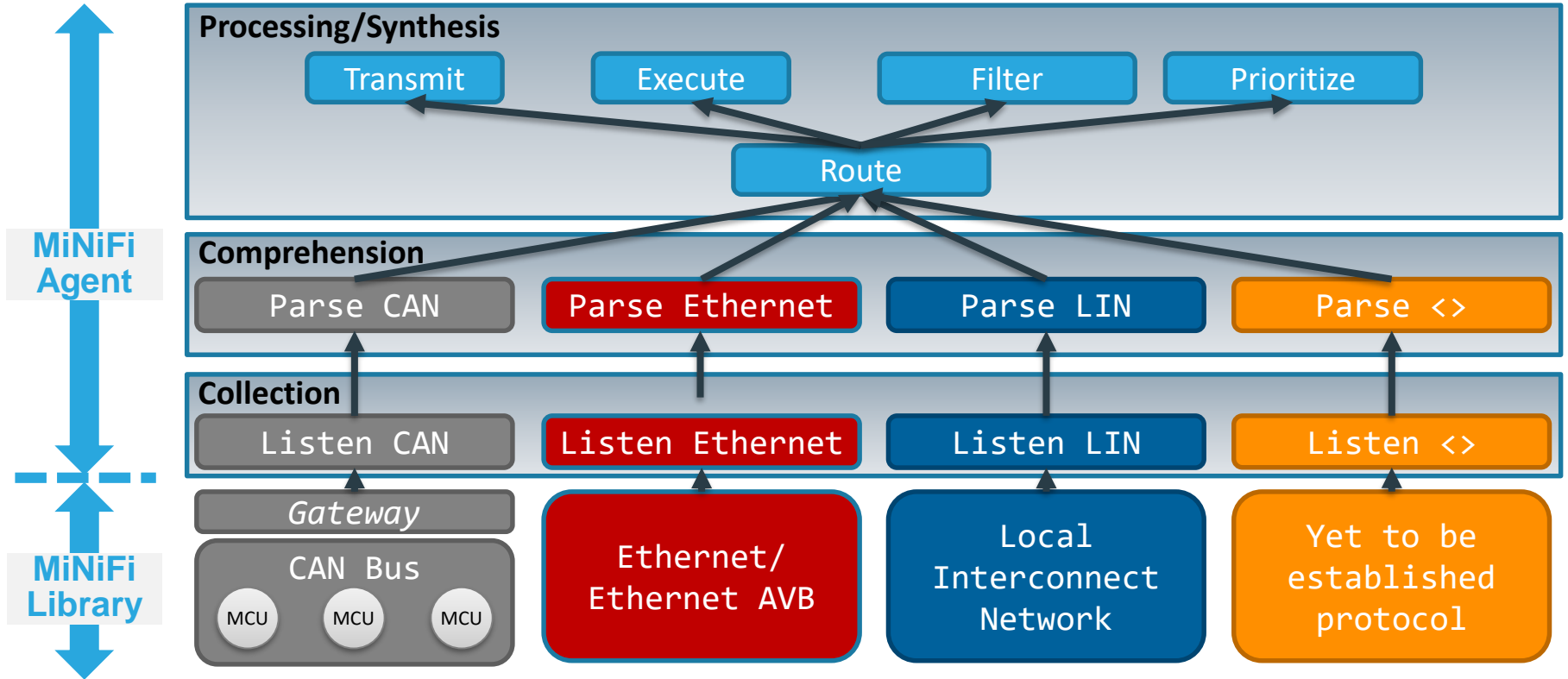
LARGE

A large blue arrow pointing downwards, indicating the progression from larger to smaller footprints.

Application	Memory Footprint	Characteristics
MiNiFi Agent (Java)	10s to 100s of MBs	<ul style="list-style-type: none">• Feature parity and reuse of core NiFi libraries
MiNiFi Agent (C++)	100s kBs to MBs	<ul style="list-style-type: none">• Write once**, run anywhere.• Adaptable to lower level interfaces
MiNiFi Library	10s kBs	<ul style="list-style-type: none">• Write n-many times, embed, run anywhere• Language libraries to support tagging, FlowFile format, Site to Site protocol and provenance generation without a full processing framework• Customizable to environments without disk storage or where threading is prohibitive• Language SDKs, Mobile Platforms

SMALL

MINIFI IN CONNECTED CAR



CAR TO CLOUD – DATA MANAGEMENT ON 5 LEVELS

THE EDGE



MANAGING THE EDGE



ENTERPRISE FLOW/STREAM ANALYTICS



EDGE TO AI

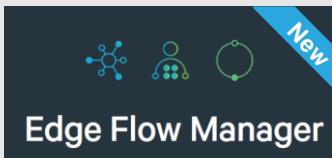


CONNECTED COMMUNITIES



CLOUDERA EDGE MANAGEMENT (CEM)

Edge Management Hub



Edge Agent Management

- Central management of agents
- Collect data from edge device
- Push intelligence to edge

Data

Intelligence

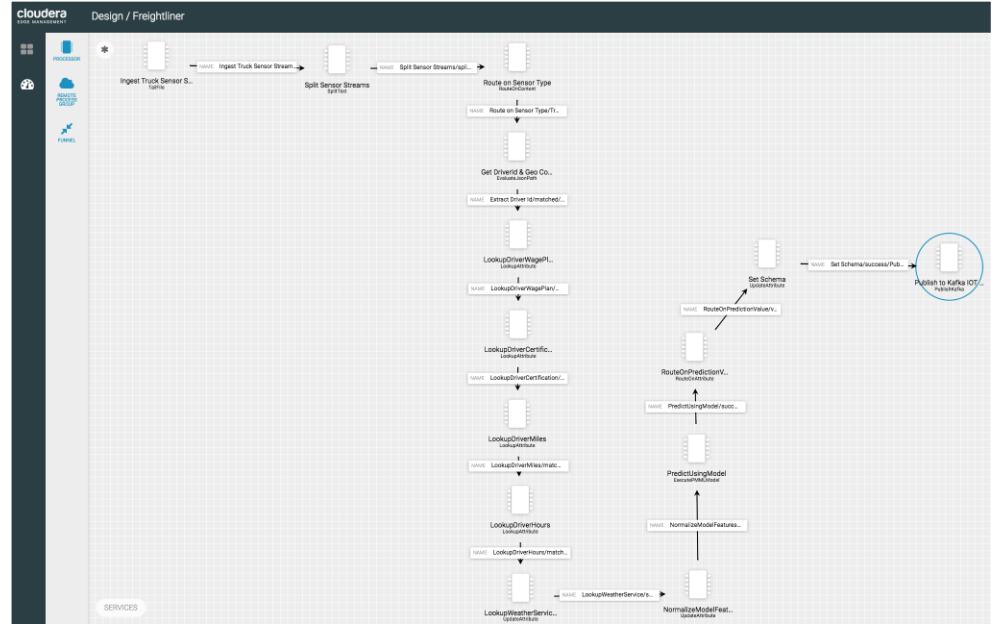
Edge Agents



- Data Collection and Processing at the Edge
- Small, lightweight footprint

EDGE FLOW MANAGER

- **Edge Management Hub**
- Graphical user interface to develop and deploy flows to edge
- Monitor thousands of edge agents
- Deploy updated Machine Learning (ML) models to edge agents
- Integration with NiFi Registry



INTELLIGENCE TO THE EDGE – PUSH ML MODELS TO EDGE

The screenshot displays the Cloudera Edge Management (EFM) Design tool interface for a project named 'Freightliner'. The main workspace shows a data flow diagram with several components: 'Ingest Truck Sensor Stream', 'Split Sensor Streams', 'Route on Sensor Type', 'Truck Speed Event', 'Extract Speed', 'High Speed Filter', 'LookupDriverCertification', 'LookupDriverMiles', 'LookupDriverHours', 'LookupWeatherService', 'NormalizeModelFeatures', 'PredictUsingModel', and 'RouteOnPredictionV...'. A 'Publish Flow' dialog box is open in the foreground, containing the text: 'Publishing this flow will make it available to all agents associated with Freightliner.' Below this, there is a 'CHANGE COMMENTS' section with a text input field containing 'Updated ML Model and capturing speed stream'. The dialog has 'CANCEL' and 'PUBLISH' buttons. A callout box with an orange border and a speech bubble tail points to the 'Publish Flow' dialog, containing the text: 'EFM will version the flow into the NiFi Registry and expose an endpoint for all agents to automatically download and deploy the flow. Powerful Edge Design/Deploy Pattern!'. The Cloudera logo and 'EDGE MANAGEMENT' are visible in the top left corner. The bottom left corner shows 'cloud' and 'SERVICES' buttons. The bottom right corner shows 'served. 14'.

CAR TO CLOUD – DATA MANAGEMENT ON 5 LEVELS

THE EDGE



MANAGING THE EDGE



ENTERPRISE FLOW/STREAM ANALYTICS



EDGE TO AI



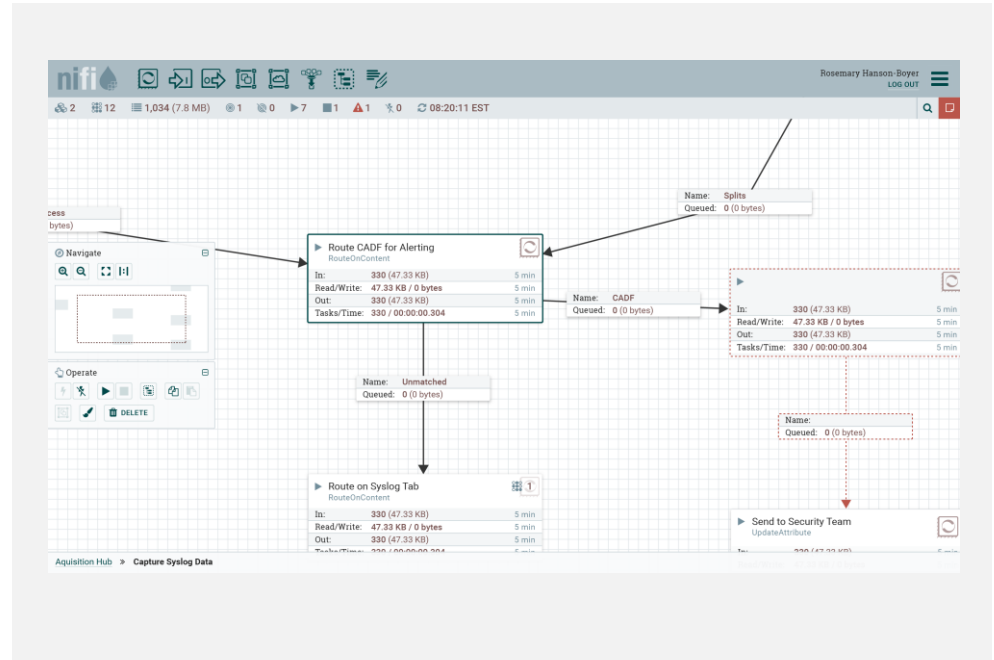
CONNECTED COMMUNITIES



FLOW MANAGEMENT VIA NIFI



- Based on Apache NiFi
- Highly configurable flow creation
- Web-based user interface
- 300+ prebuilt processors
- Secure - fine grained encryption
- Data provenance
- Guaranteed delivery
 - Buffer data during system interruptions
- Designed for extensibility
- NiFi Registry

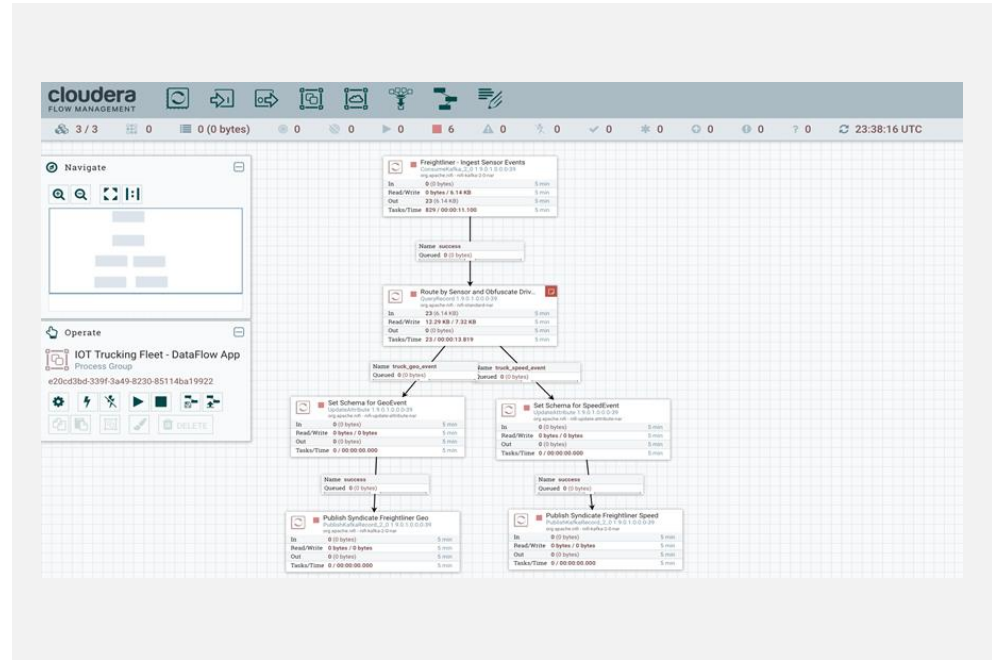


WHAT CAN BE DONE WITH NIFI PROCESSORS?

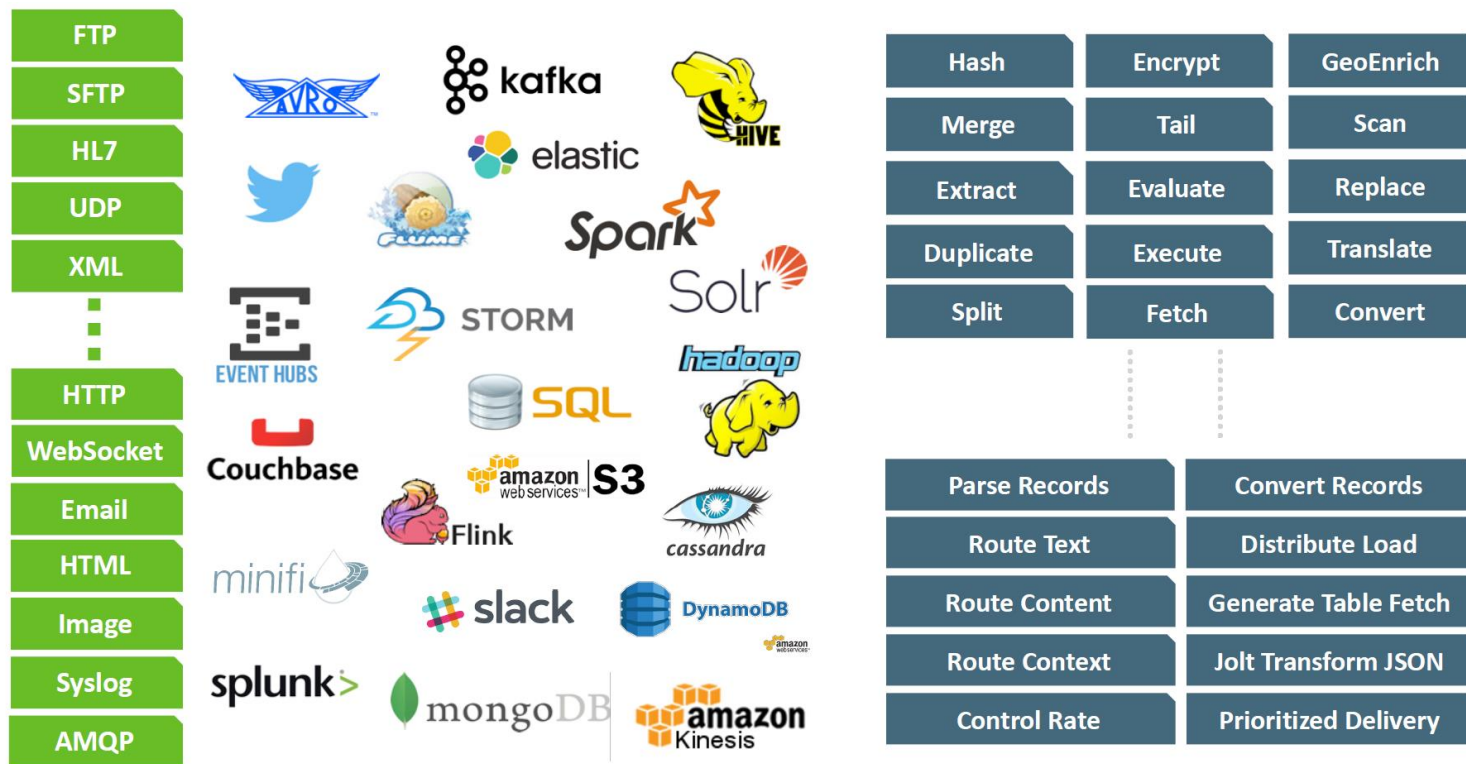
300+ Pre-Built Processors



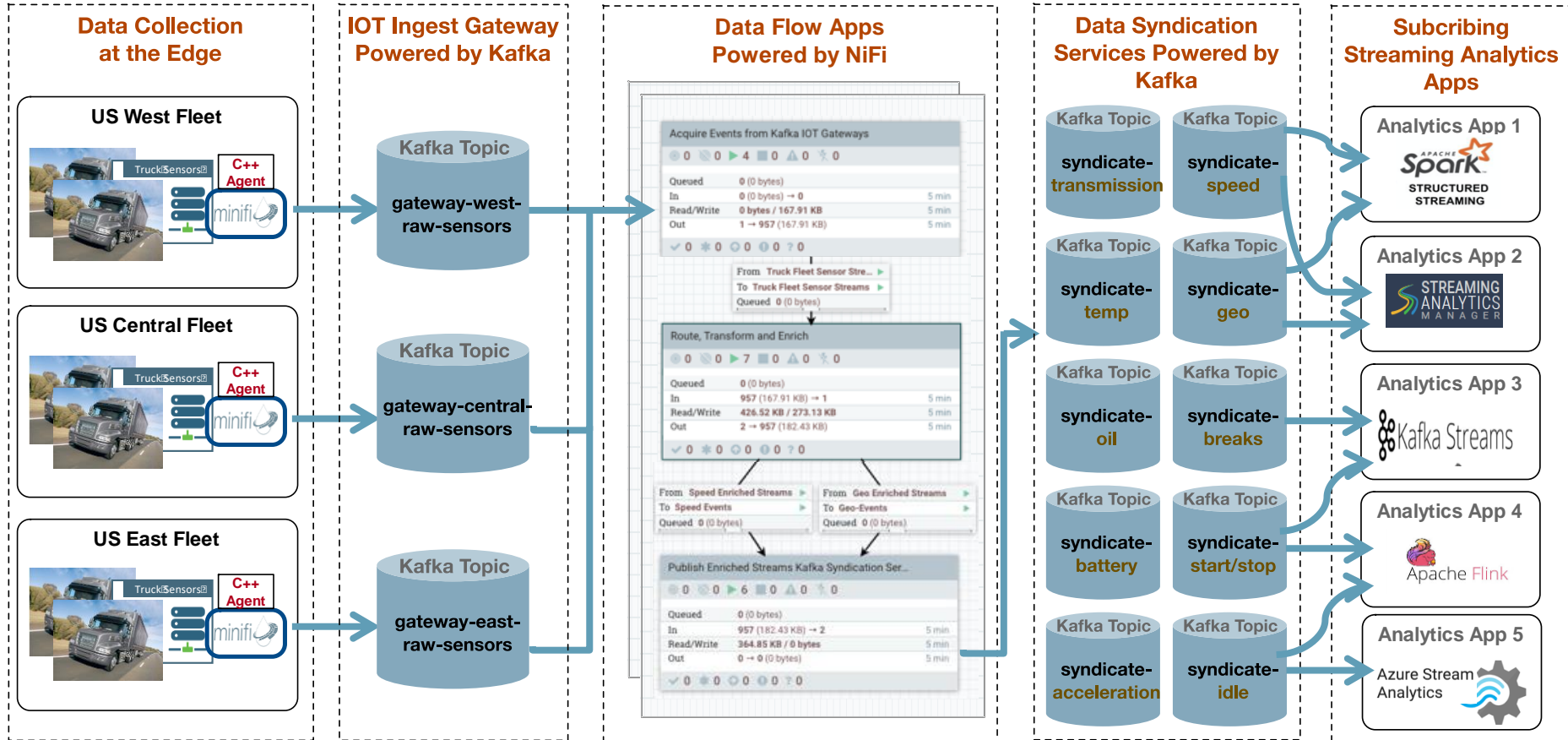
- **Ingestion:** connectors to read/write data from/to several data sources
 - **Protocols:** HTTP (S), AMQP, MQTT, UDP, TCP, CEF, JMS, (S) FTP, AWS IoT, Raw Socket Protocol
 - **Brokers:** Kafka, JMS, AMQP, MQTT etc.
 - **Databases:** JDBC, MongoDB, HBase, Cassandra etc.
- **Extraction** (XML, JSON, Regex, Grok etc.)
- **Transformation**
 - Format conversion (JSON to Avro, CSV to ORC etc.)
 - Compression/decompression, Merge, Split, encryption etc.
- **Data enrichment**
 - Attribute, content, rules etc.
- **Routing**
 - Priority, dynamic/static, based on content or metadata etc.



300+ PROCESSORS FOR DEEPER ECOSYSTEM INTEGRATION



STREAMING ANALYTICS REFERENCE ARCHITECTURE





EXAMPLE - LAS VEGAS CONNECTED VEHICLE PILOT

**City Fleet Vehicle
In Operation**



**Monitor Vehicle
Conditions**



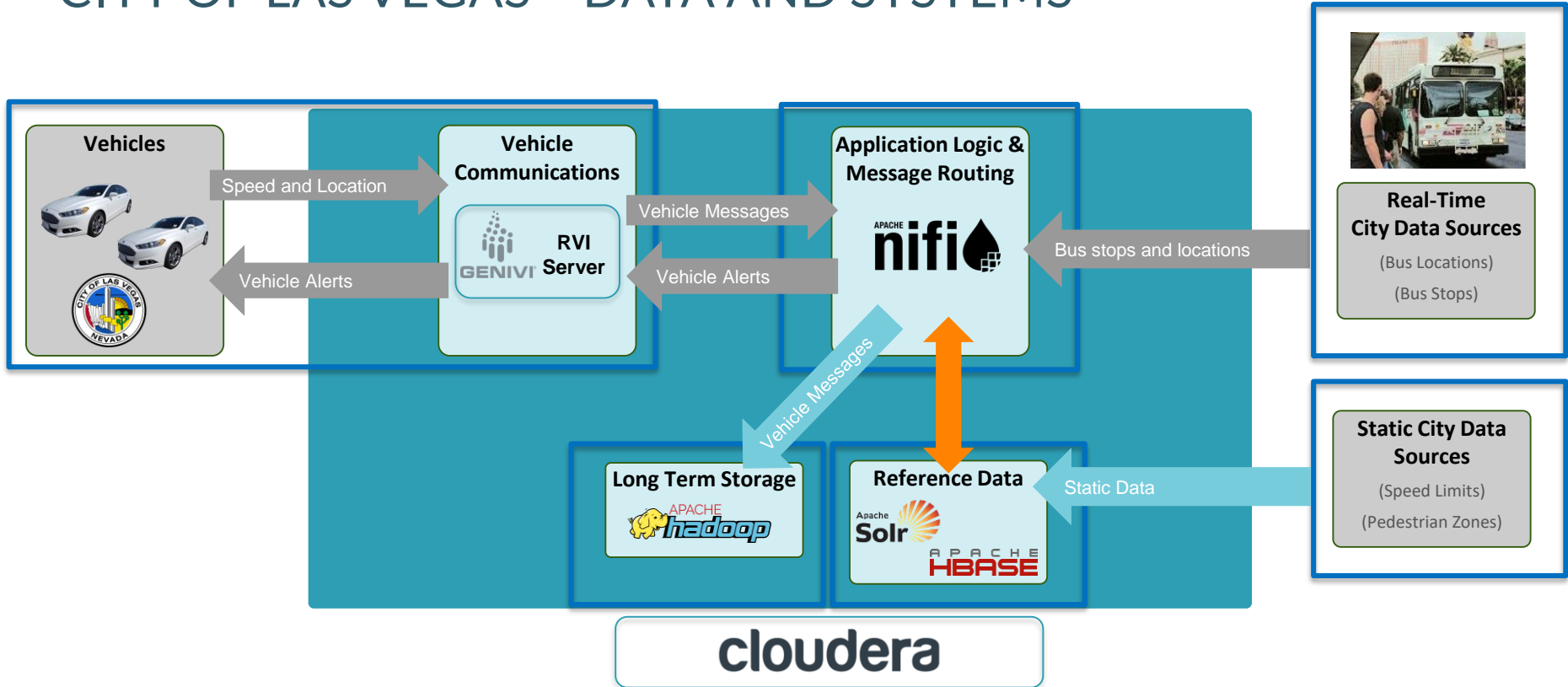
**Issue Alert to
In-Vehicle Display Unit**



**Approaching Active
Bus Stop?**



CITY OF LAS VEGAS – DATA AND SYSTEMS



Navigate

Search, Refresh, Full Screen, Split View

Operate

Receive Vehicle Data
Process Group
de77560c-015b-1000-6f9b-3579589b154d

Settings, Stop, Play, Pause, Refresh, Delete

1. Receive Vehicle Data (ListenHTTP)

Receive RVI Data ListenHTTP	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

2. Fork Data

Route obd or gps Data RouteOnContent	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

MergeContent MergeContent	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

3a. Extract GPS Data

Extract GPS Coordinates EvaluateJsonPath	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

3b. Extract Vehicle Speed Data

Extract OBD Coordinates EvaluateJsonPath	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

PutHDFS PutHDFS	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

6. Store in Hadoop

4. Reformat for each use case

Format Vehicle Data UpdateAttribute	5 min
In 0 (0 bytes)	5 min
Read/Write 0 bytes / 0 bytes	5 min
Out 0 (0 bytes)	5 min
Tasks/Time 0 / 00:00:00.000	5 min

5. Output Data to Use Cases

Check Speeding war...	
--------------------------	--

To UC1

UC2 High-risk pedestrian area...	
-------------------------------------	--

To UC2

UC3 Bus stop warning	
-------------------------	--

To UC3

UC3 – Active Bus Stop Warning

Navigate

Operate

UC3 Bus stop warning
Process Group
e83b2aef-015b-1000-0002-3526ebdc22b1

DELETED

Receive RVI Data

1. Receive UC3 RVI Data (Location)

Queued 0 (0 bytes)

Build Query
UpdateAttribute

In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	0 / 00:00:00.000	5 min

2. Build location based Query

Name success
Queued 0 (0 bytes)

Query Solr for nearby Bus Stops
InvokeHTTP

In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	0 / 00:00:00.000	5 min

3. Query if vehicle approaching bus stop location

Name Response
Queued 0 (0 bytes)

Evaluate Results
Evaluate.JsonPath

In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	0 / 00:00:00.000	5 min

4. Determine if bus stop is ACTIVE

Name matched
Queued 0 (0 bytes)

Post Warning to Driver

5. To Active Bus Stop Warning

Issue Alert to In-Vehicle Display Unit



CAR TO CLOUD – DATA MANAGEMENT ON 5 LEVELS

THE EDGE



MANAGING THE EDGE



ENTERPRISE FLOW/STREAM ANALYTICS



EDGE TO AI

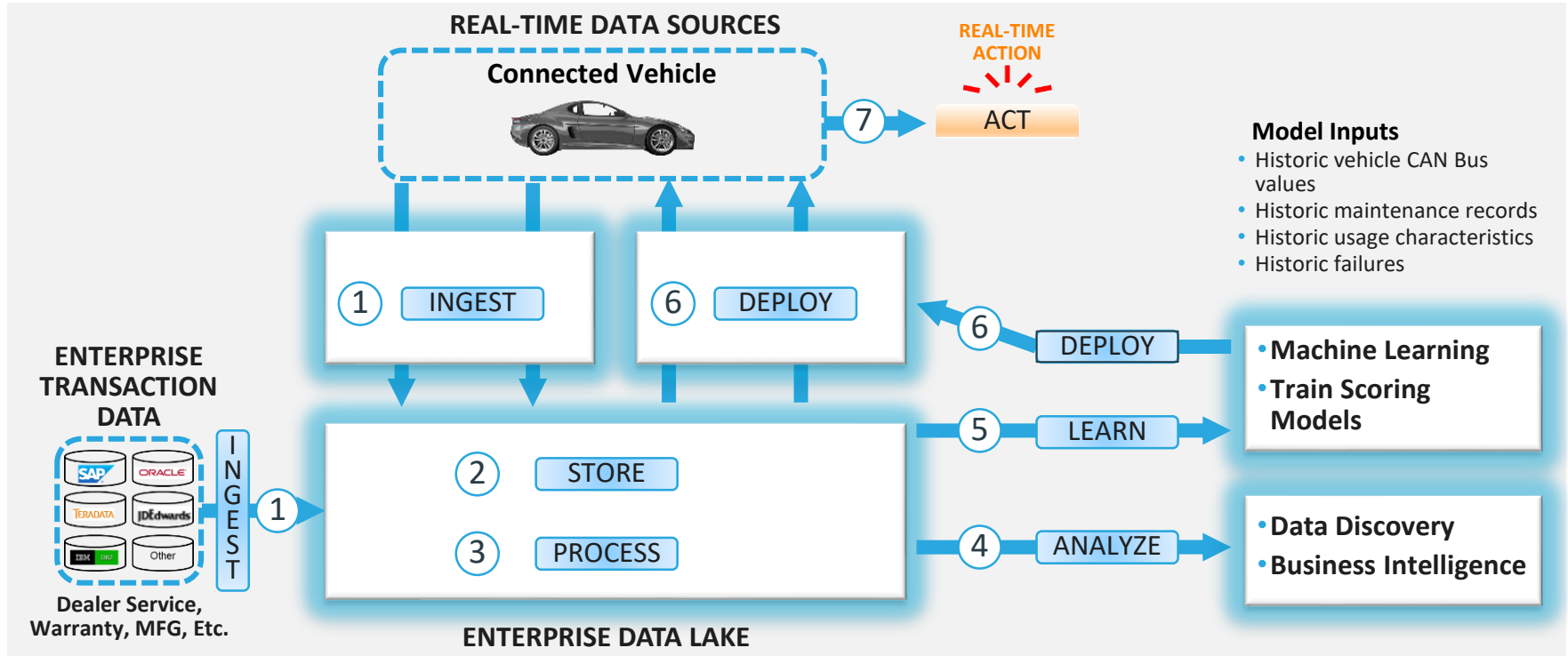


CONNECTED COMMUNITIES

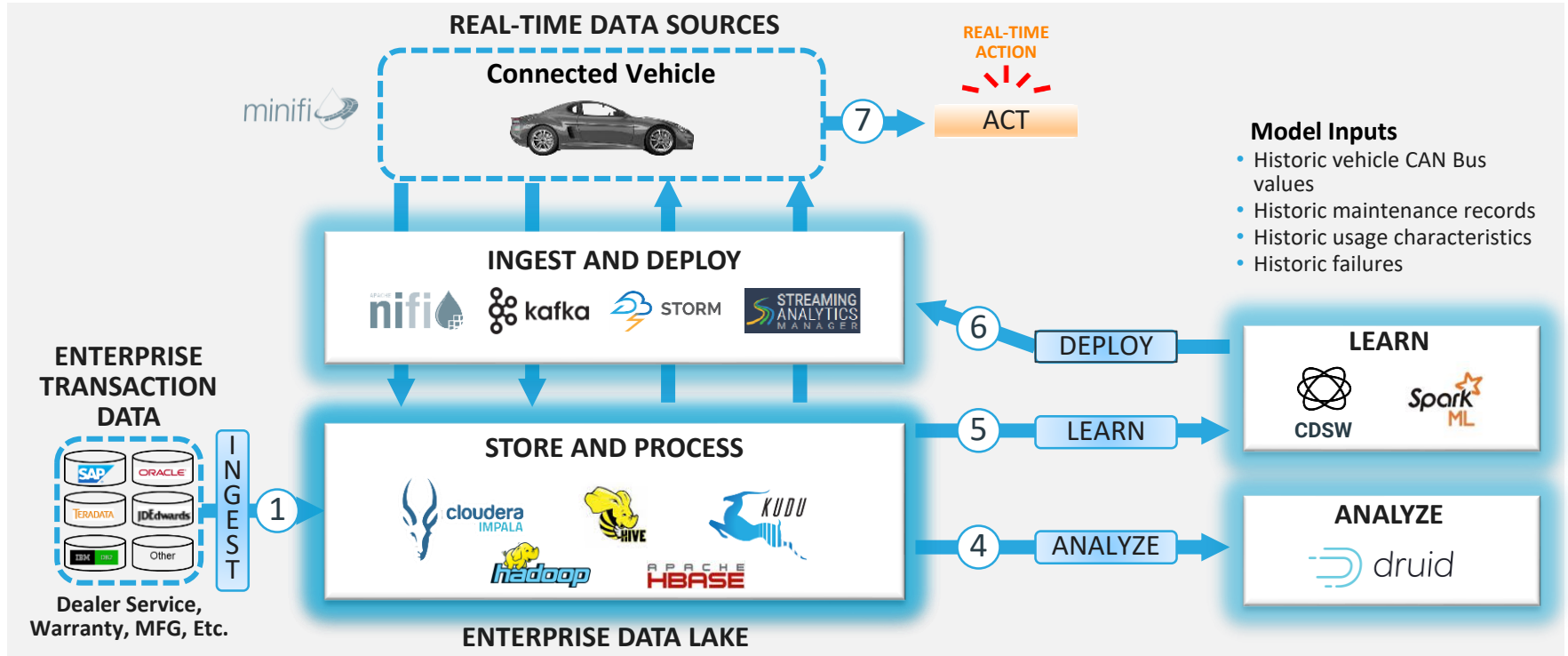


CONNECTED VEHICLE ANALYTICS LIFECYCLE

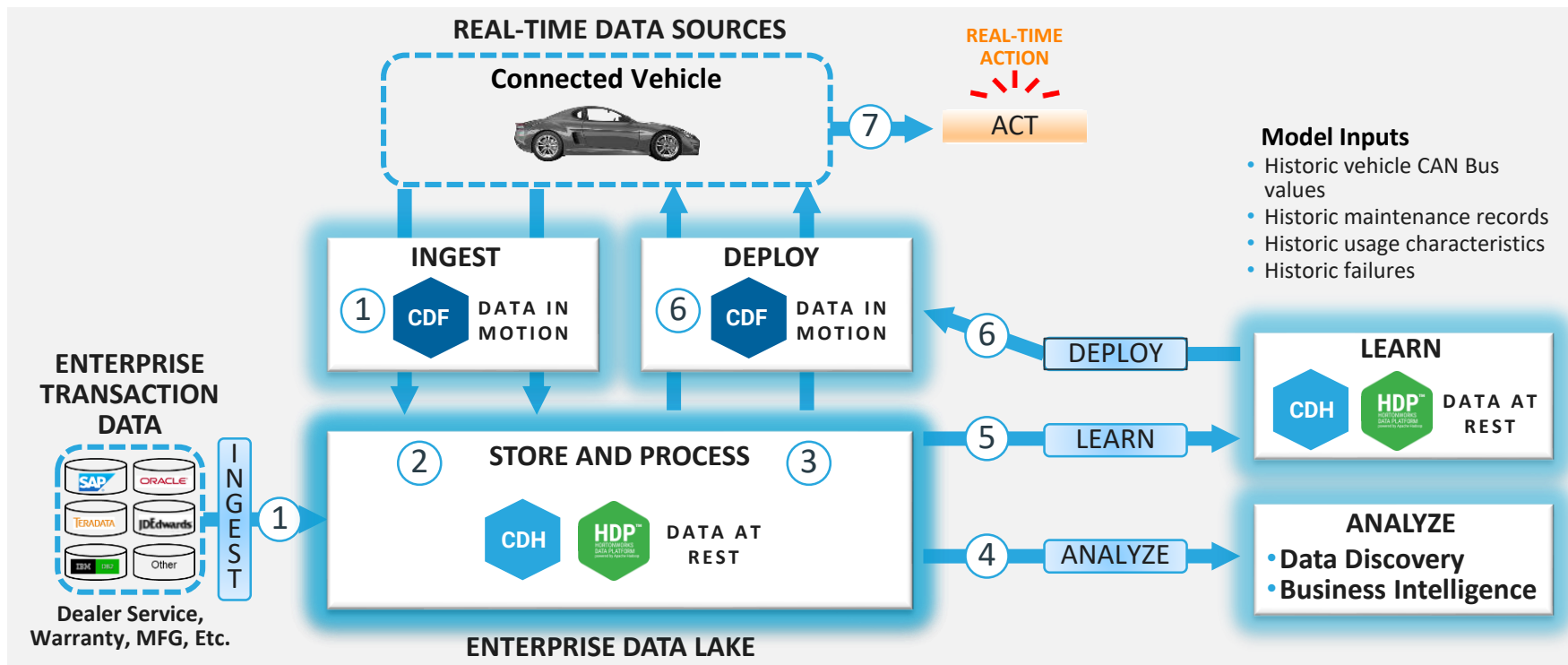
Example: Vehicle Predictive Maintenance



SOLUTION COMPONENT DETAIL



CLOUDERA FOR CONNECTED VEHICLE INNOVATION



CAR TO CLOUD – DATA MANAGEMENT ON 5 LEVELS

THE EDGE



MANAGING THE EDGE



ENTERPRISE FLOW/STREAM ANALYTICS



EDGE TO AI

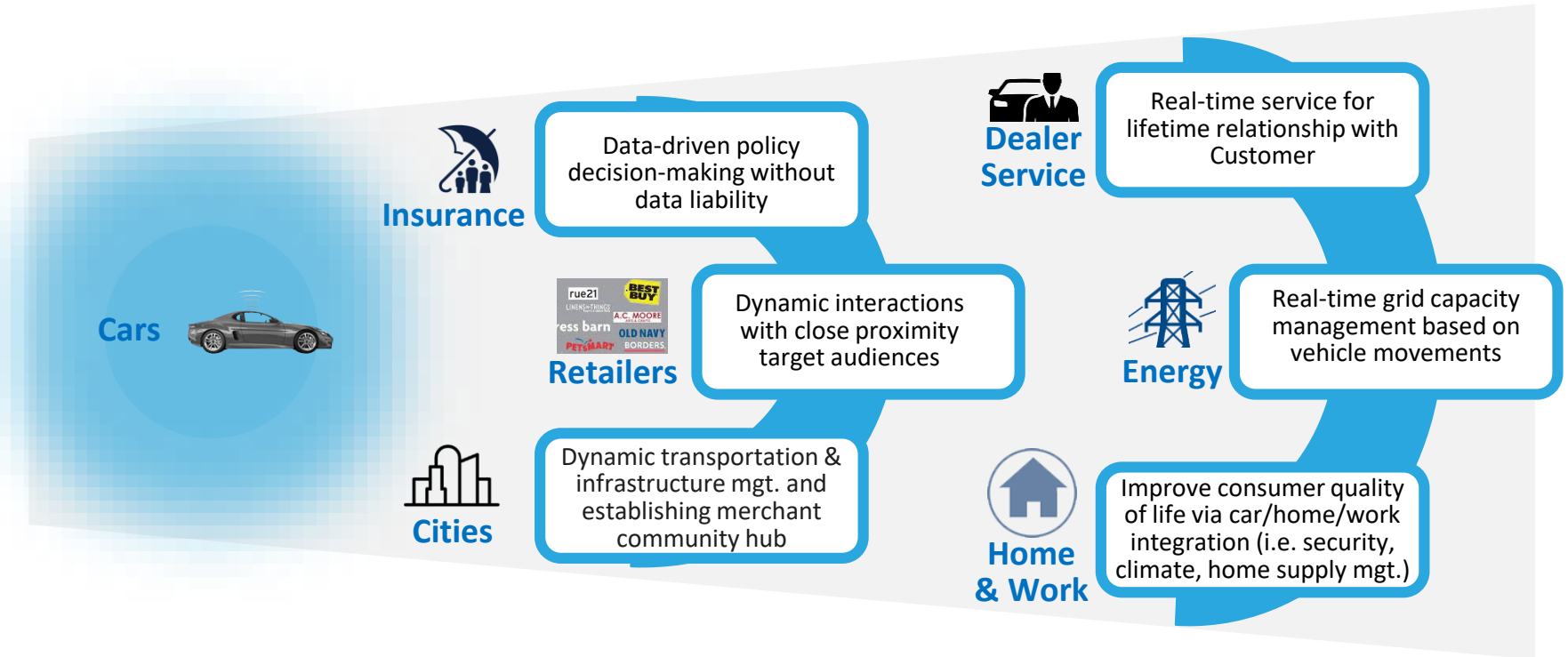


CONNECTED COMMUNITIES



CONNECTED VEHICLE IS LEADING THE TRANSFORMATION

Monetizing the Connected Vehicle



ESTABLISHING TRUST – WHAT’S REQUIRED?



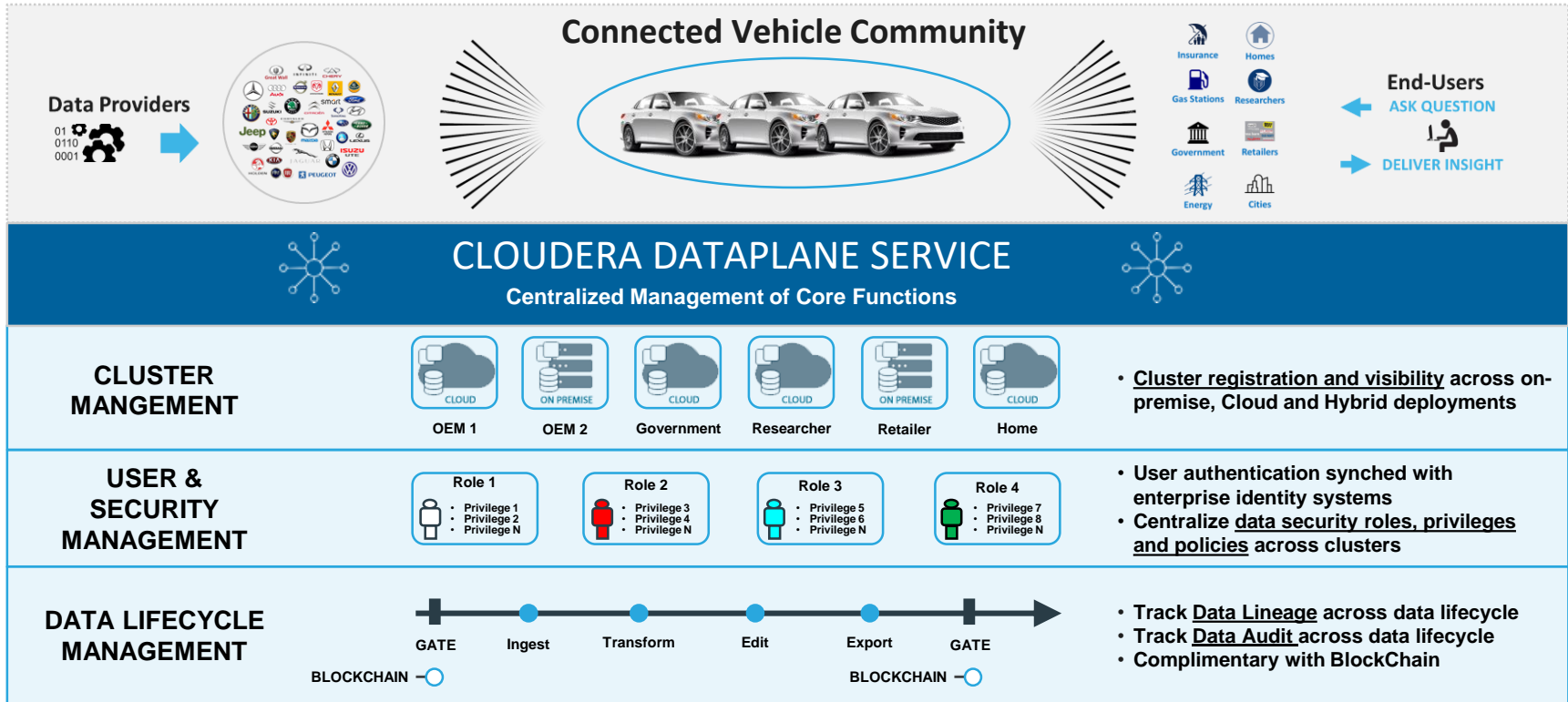
✓ **Data access across systems**

✓ **Controlling who sees what data**

✓ **Data audit across the ecosystem**

- Each data provider shares specific *data assets* and defines specific data sharing rules
- Each data provider precisely defines *which end-users can see their data* throughout the data lifecycle
- Each data provider can track *data usage* across the data lifecycle, including the *data's origin, where it is used, and how it changes over time*

HOW IT WORKS.....



- DATAPLANE ADMIN
- Clusters**
- Users
- Services
- Settings

Admin / Clusters

Single View of All Data Sources



Clusters







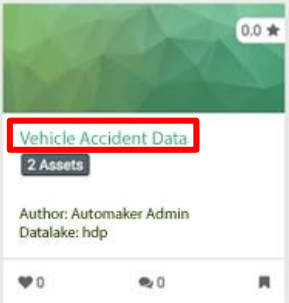
Search ADD CLUSTER

Status	Name	Location	Data Center	Nodes	Uptime	HDFS Used	
	Automaker	Detroit	DC_Detroit	NA	NA	NA/NA	
	Vehicle Owner	Los Angeles	DC_Los Angeles	NA	NA	NA/NA	
	Telematics Provider	Dallas	DC_Dallas	1	a month	76 GB /140 GB	
	Lease Provider	Kansas City	DC_Kansas City	NA	NA	NA/NA	
	City	Sacramento	DC_Sacramento	NA	NA	NA/NA	
	Insurance Provider	Hartford	DC_Hartford	NA	NA	NA/NA	
	Body Shop	Los Angeles	DC_Los Angeles	NA	NA	NA/NA	

Asset Collections

- Asset Collection**
- Bookmarks
- Dashboard
 - hdp, DC_Dallas
 - hdp, DC_Detroit
 - hdp, DC_Hartford
 - hdp, DC_Kansas City
 - hdp, DC_Los Angeles
 - hdp, DC_Los Angeles 2
 - hdp, DC_Sacramento
- Profiler

Type to search Filter by Tags

 <p>Automaker Data 15 Assets Author: Automaker Admin Datalake: hdp</p>	 <p>Vehicle Owner Data 4 Assets Author: Automaker Admin Datalake: hdp</p>	 <p>Telematics Provider Data 3 Assets Author: Telematics Admin Datalake: hdp</p>	 <p>Lease Provider Data 4 Assets Author: Lease Admin Datalake: hdp</p>	 <p>City Data 2 Assets Author: City Admin Datalake: hdp</p>
 <p>Insurance Provider Data 13 Assets Author: Insurance Admin Datalake: hdp</p>	 <p>Vehicle Accident Data 2 Assets Author: Automaker Admin Datalake: hdp</p>			

- DATA STEWARD STUDIO
- Asset Collection**
- Bookmarks
- Dashboard
 - hdp, DC_Dallas
 - hdp, DC_Detroit
 - hdp, DC_Hartford
 - hdp, DC_Kansas City
 - hdp, DC_Los Angeles
 - hdp, DC_Los Angeles 2
 - hdp, DC_Sacramento
- Profiler

Vehicle Accident Data
Data relating to vehicle accidents

Asset Collection Details

Overview Assets

4 Assets
Tables

Search

SOURCE	NAME	DATABASE NAME	OWNER	CREATED TIME
HIVE	Emergency Contacts	OEM_owner_portal	Automaker Admin	-
HIVE	Vehicle Events	OEM_vehicle_events	Automaker Admin	-
HIVE	Emergency Incidents	OEM_emergency_details	Automaker Admin	-
HIVE	Emergency Incident Analysis	OEM_emergency_aggregation	Automaker Admin	-

0.0 ★

♥ 1 💬 0 📄

Created By
admin

Datalake
hdp

Tags
customer prospect lead

System Tags
CURATED

Created On
Jun 19, 2016, 12:08 AM

Last Modified
Jun 19, 2016, 12:08 AM

Data Steward / Asset Collections / Details / Asset Details (OEM_owner_portal) Customer Emergency Event Data

Tagging of Sensitive Information

Vehicle Identification Number (VIN)	string	
Vehicle License Plate Number	string	
Primary Driver Name	string	name
Primary Driver DOB	date	
Primary Driver Address	string	
Primary Driver Phone Number	string	telephone
Primary Driver Driver License Number	string	
Primary Driver SSN	string	ssn
Primary Driver Blood Type	string	
Primary Driver EMERGENCY MEDICAL INSTRUCTIONS	string	
Primary Driver Emergency Contact Name	string	name
Primary Driver Emergency Contact Phone Number	string	telephone
Emergency Contact Relationship to Primary Driver	string	
2nd Driver Name	string	name
2nd Driver DOB	date	
2nd Driver Address	string	
2nd Driver Phone Number	string	telephone
2nd Driver Driver License Number	string	
2nd Driver SSN	string	ssn
2nd Driver Blood Type	string	
2nd Driver EMERGENCY MEDICAL INSTRUCTIONS	string	
2nd Driver Emergency Contact Name	string	name
2nd Driver Emergency Contact Phone Number	string	telephone
Emergency Contact Relationship to 2nd Driver	string	

* Approximate values as being computed using HLL algorithm.

- Asset Collection
- Bookmarks
- Dashboard
 - hdp, DC_Dallas
 - hdp, DC_Detroit
 - hdp, DC_Hartford
 - hdp, DC_Kansas City
 - hdp, DC_Los Angeles
 - hdp, DC_Los Angeles 2
 - hdp, DC_Sacramento
- Profiler

OEM_emergency_aggregation

HIVE

Adding Security Policies

OVERVIEW SCHEMA **POLICY** AUDIT

Tag Based Policies

POLICIES

USERS

Policy ID	Policy Name	Status	Audit Logging	Group	Users
26	access: CITY_emergency_analysis_table	ENABLED	ENABLED	public	-
28	mask: name	ENABLED	ENABLED	govt_analyst, city_analyst, insurance_analyst	-
34	mask: ssn show first 4	ENABLED	ENABLED	govt_analyst, city_analyst, insurance_analyst	-
35	mask: telephone	ENABLED	ENABLED	govt_analyst, city_analyst, insurance_analyst	-

Resource Based Policies *Table and Field Level Controls*

This asset has 0 policies associated with it.

- Asset Collection
- Bookmarks
- Dashboard
 - hdp, DC_Dallas
 - hdp, DC_Detroit
 - hdp, DC_Hartford
 - hdp, DC_Kansas City
 - hdp, DC_Los Angeles
 - hdp, DC_Los Angeles 2
 - hdp, DC_Sacramento
- Profiler

OEM_emergency_details

Data Lineage (For Table)

OVERVIEW SCHEMA POLICY AUDIT

50001
Number of Rows

39
Number of Columns

5
Sensitive Columns



50001
Number of Rows

39
Number of Columns

Lineage

Customer Emergency Contact Information

/hive_data/OEM_owner_portal

/hive_data/OEM_vehicle_events

Vehicle Event Data

create external table

OEM_emergency_details

OEM_emergency_aggregation

Accident Detail
(For Use by City
Emergency Dispatch)

Aggregated Detail
(For Use by Various
Research Entities)

Lineage Impact

Top 10 Users

Table Properties

Owner: hive

Qualified Name: hortoniabank.us_customers@hdp

Table Type: EXTERNAL_TABLE

Database: hortoniabank

Created On: 13 Jun 2018

Last Modified: 12 Jul 2018

System Tags

email name ssn telephone

Profilers

Ranger Audit Profiler	Last Run
Active	-
Hive Column Profiler	Last Run
Active	a day ago
Sensitive Profiler	Last Run
Active	24 minutes ago
Hive Metastore Profiler	Last Run
Active	-

OEM_emergency_aggregation

HIVE

Audit Trail

OVERVIEW SCHEMA POLICY **AUDIT**

Event Time

Event Type

Access Type: ALL SELECT UPDATE CREATE DROP ALTER INDEX READ WRITE Result: ALL

Policy ID	Event Time	User	Resource Type	Access Type	Result	Acc	Client IP
21	07/17/2018 16:11:49 GMT	kate_hr	@column	SELECT	ALLOWED	ranger	127.0.0.1
28	07/17/2018 16:11:49 GMT	kate_hr	@column	SELECT	ALLOWED	ranger	127.0.0.1
21	07/17/2018 16:08:37 GMT	kate_hr	@column	SELECT	ALLOWED	ranger-acl	127.0.0.1
28	07/17/2018 16:08:37 GMT	kate_hr	@column	SELECT	ALLOWED	ranger-acl	127.0.0.1
28	07/17/2018 16:06:44 GMT	mark_bizdev	@column	SELECT	DENIED	ranger-acl	127.0.0.1
28	07/17/2018 16:05:26 GMT	diane_csr	@column	SELECT	DENIED	ranger-acl	127.0.0.1
28	07/17/2018 16:02:46 GMT	eti_user	@column	SELECT	DENIED	ranger-acl	127.0.0.1
21	07/17/2018 16:02:05 GMT	kate_hr	@column	SELECT	ALLOWED	ranger-acl	127.0.0.1
28	07/17/2018 16:02:05 GMT	kate_hr	@column	SELECT	ALLOWED	ranger-acl	127.0.0.1
28	07/12/2018 14:14:36 GMT	jermy_contractor	@column	SELECT	DENIED	ranger-acl	127.0.0.1
28	07/12/2018 14:13:52 GMT	diane_csr	@column	SELECT	DENIED	ranger-acl	127.0.0.1
28	07/12/2018 14:10:36 GMT	jermy_contractor	@column	SELECT	DENIED	ranger-acl	127.0.0.1
28	07/12/2018 14:07:29 GMT	ivanna_eu_hr	@column	SELECT	DENIED	ranger-acl	127.0.0.1
28	07/12/2018 14:05:54 GMT	joe_analyst	@column	SELECT	ALLOWED	ranger-acl	127.0.0.1
23	07/12/2018 14:05:54 GMT	joe_analyst	@column	MASK_NULL	ALLOWED	ranger-acl	127.0.0.1
34	07/12/2018 14:05:54 GMT	joe_analyst	@column	MASK_SHOW_LAST_4	ALLOWED	ranger-acl	127.0.0.1
39	07/12/2018 14:05:54 GMT	joe_analyst	@column	CUSTOM	ALLOWED	ranger-acl	127.0.0.1
38	07/12/2018 14:05:54 GMT	joe_analyst	@column	CUSTOM	ALLOWED	ranger-acl	127.0.0.1
35	07/12/2018 14:05:54 GMT	joe_analyst	@column	MASK_SHOW_FIRST_4	ALLOWED	ranger-acl	127.0.0.1
23	07/12/2018 14:05:54 GMT	joe_analyst	@column	MASK_HASH	ALLOWED	ranger-acl	127.0.0.1

CONCLUSION

- Connected Vehicle Data is Key to Industry Transformation
- Evolving Intelligence at Edge and Connected Community Use Cases
- Big Data Management Challenges on Multiple Levels
- Open Source Data Management Innovation Can Help!

THANK YOU

CLOUDERA