# Vision for Trade Data Lake

Data from new sources such as UPU (postal), ICAO (aviation), vessel tracking, etc, possibly in streams, to be included along with traditional trade statistics plus the underlying administrative data sources (such as customs declarations or shipping manifests) for better in-depth, timely and accurate analysis/forecast of international trade.

The comparative advantage over existing data would be timelier and richer analysis with possibilities of nowcasting and accurate short term forecasting

# Problem statement

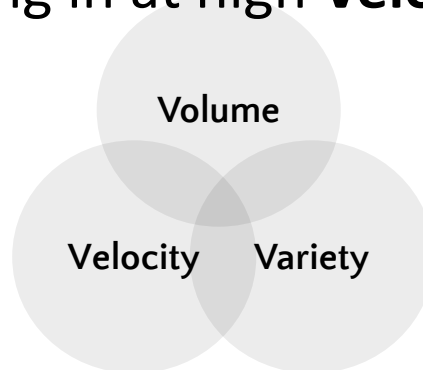The introduction of new heterogeneous streams of data from various new sources demands

- elastic storage capacity that is flexible and more reliable
- real-time analytical capabilities (schema-on-read)
- convenience to do visualizations across these different data sources (structured/unstructured, processed/unprocessed)
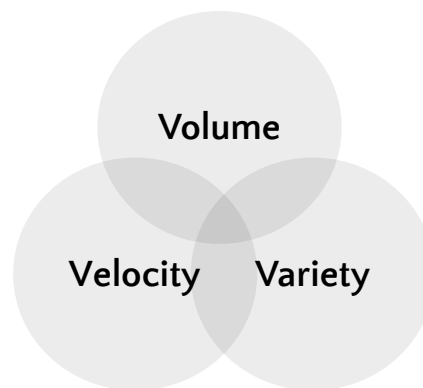- access to data anywhere/anytime through secure channel

# Addressing the challenge

- Reliable elastic storage requirements can be met by moving to cloud-based solutions
- Real-time analytical capability can be realized with interactive query tools (big data analytics)
- Visualization across divergent datasets can possibly achieved through SaaS (Software as a service) tools

# Features of big data analytics

- Possible to store and analyze huge **volumes** of structured/unstructured data that has **Variety** in nature and flowing in at high **Velocity**.

Volume

Velocity    Variety

- When the existing platforms/infrastructures cannot withstand the three Vs, we move to Cloud based platforms such as AWS.
- When the existing analytical solutions cannot withstand the three Vs, we move to Big data analytics such as Hadoop.

# Types of cloud architectures

## IaaS/PaaS

**(Infrastructure as a service/Platform as a service)**

The cloud vendor provides the service up to OS level (IaaS) or application level (PaaS) and we manage the rest.

- The required service is not available in the cloud (such as custom-applications built in house).
- bound to specific technologies/tools

## Serverless (SaaS)

**(Software as a service)**

No servers, no maintenance and no licensing overheads. Desired functionalities are in the form of microservices (or nanoservices). No provisioning of underlying platforms such as OS, application, etc.

- The required service is readily available on the cloud such as data visualizations, api gateway, etc.
- Not bound to any tool/technology.

# Inside the Trade Data Lake

| | |
|---|---|
| CT_Address.CT_Location.code | String |
| CT_Address.CT_Location.name | String |
| CT_m307_supported_messages.T_Item.attributes_version_number | String |
| CT_m307_supported_messages.T_Item.customs_procedure | String |
| CT_m307_supported_messages.T_Item.dangerous_good_hazard_class | String |
| CT_m307_supported_messages.T_Item.dangerous_good_type_e_code | String |
| CT_m307_supported_messages.T_Item.declared_gross_weight | String |
| CT_m307_supported_messages.T_Item.measured_gross_weight | String |
| CT_m307_supported_messages.T_Item.nature_of_transaction_code | String |

## Postal Data (National)

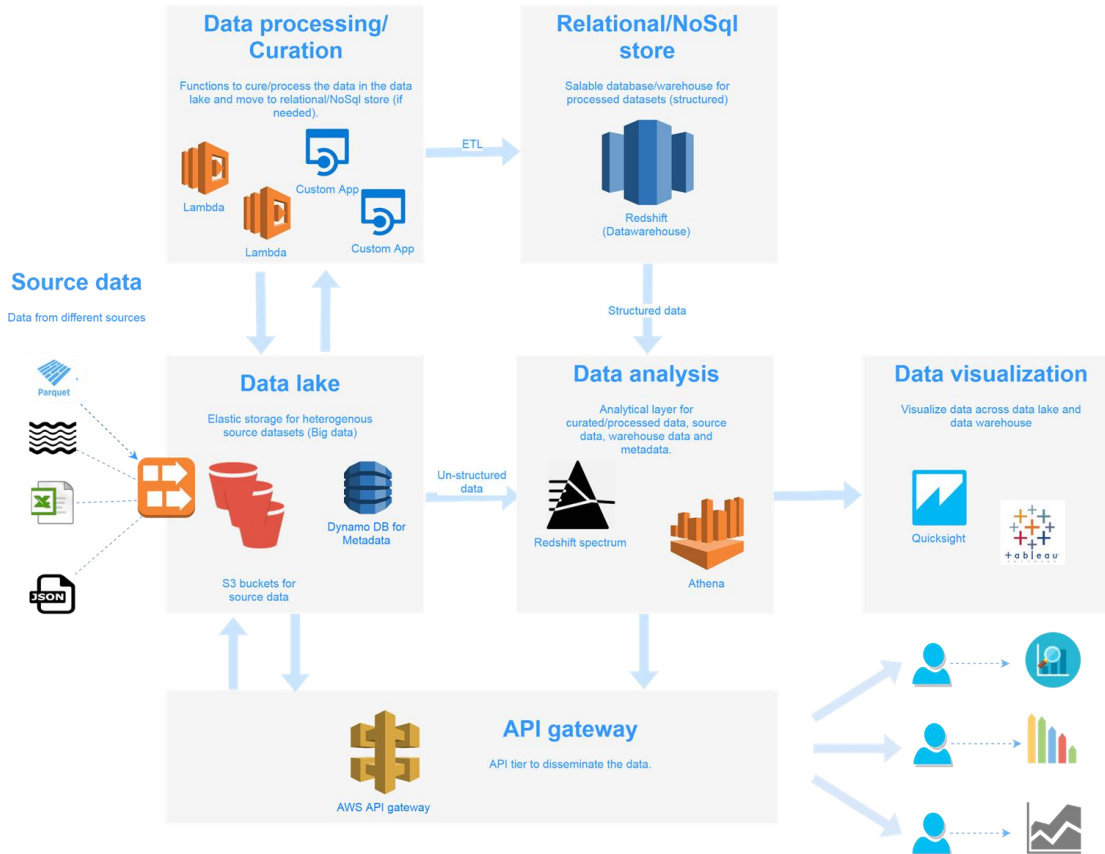| FIELD | DESCRIPTION |
|---|---|
| TIMESTAMP | Time of ship position detection / reception (in UTC) |
| MMSI | Ship's MMSI number sent with the AIS notification |
| Lat | Latitude of the ship position (in decimal degrees) |
| Lon | Longitude of the ship position (in decimal degrees) |
| "Speed over ground" | Speed over ground (in knots) |
| "Course over ground" | Course over ground (in degrees) |
| Heading | True heading (in degrees (0-359)) |
| "IMO number" | Ship's IMO number sent with the AIS notification |
| Shipname | Ship's vessel name sent with the AIS notification |
| Callsign | Ship's Callsign sent with the AIS notification |
| "Type of ship" | Ship type |
| Draught | Maximum Present Static Draught (in meters) |
| Destination | Destination |

BigData
UN Global Working Group

# AIS Vessel Tracking Data

# In addition:

- Freight Aviation Data (ICAO)
- Shipping Manifests
- Advanced Export Declarations
- ASYCUDA records

# Data Architecture based on AWS Cloud solutions

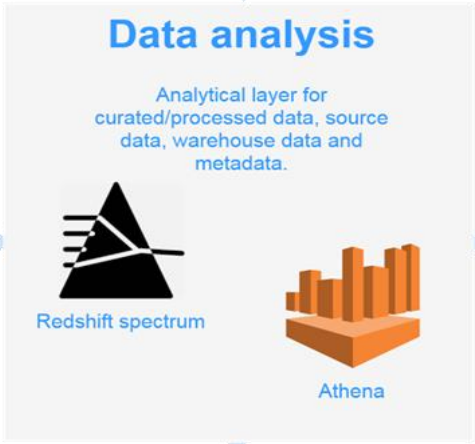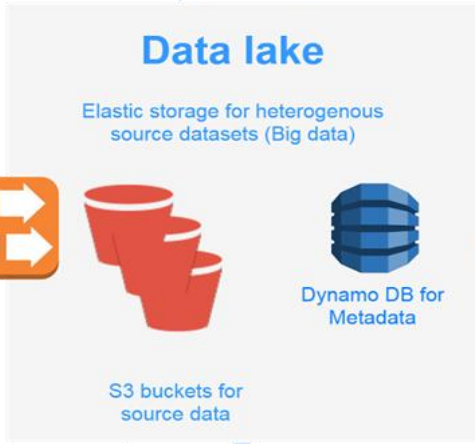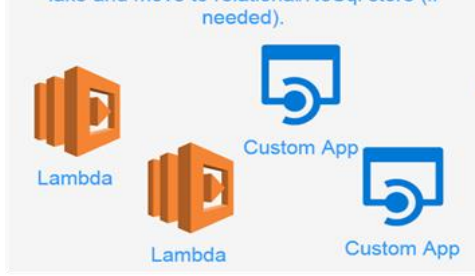Data Lake is AWS's elastic storage layer to dump multi data sources

These data can be curated and stored at AWS data warehouse service (Redshift)

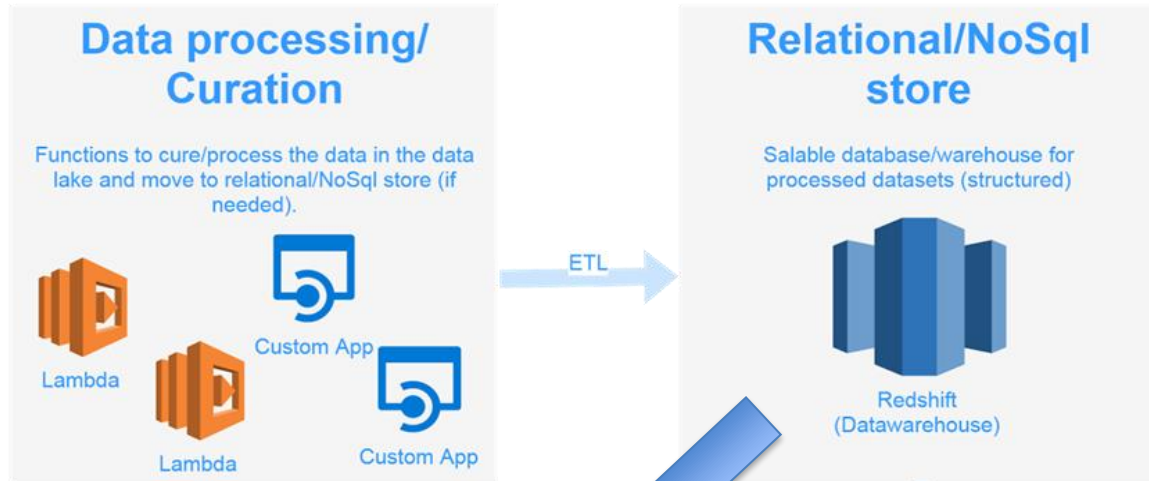Or can be directly analyzed in analytical layer

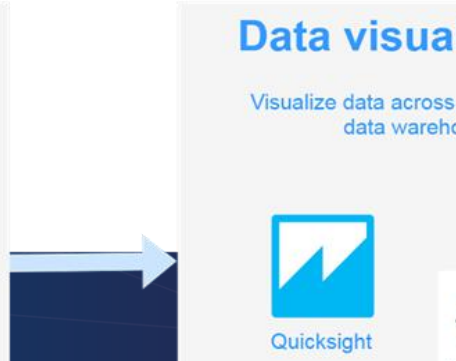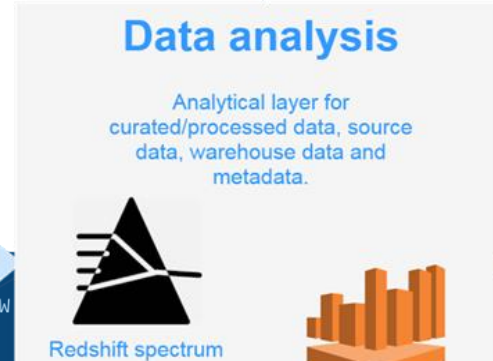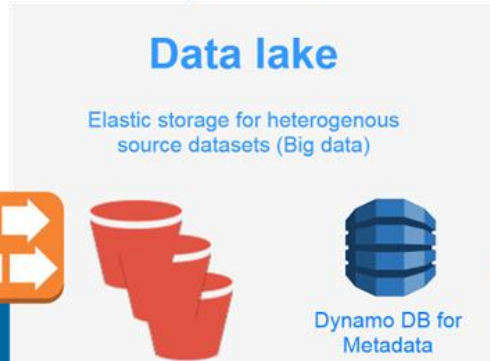Raw data at Data lake and Refined data from analytical layer can be distributed through API or be visualized

**Source data**

Data from different sources

Parquet

JSON

lake and move to relational nosql store (if needed).

Lambda

Lambda

Custom App

Custom App

ETL

processed datasets (structured)

Redshift (Datawarehouse)

Structured data

## Data lake

Elastic storage for heterogenous source datasets (Big data)

Dynamo DB for Metadata

S3 buckets for source data

Un-structured data

## Data analysis

Analytical layer for curated/processed data, source data, warehouse data and metadata.

Redshift spectrum

Athena

## Data visuali

Visualize data across da data warehous

Quicksight

BigData
UN Global Working Group

**API gateway**

Elastic storage for heterogenous
source datasets (Big data)

Un-structured
data

Analytical layer for
curated/processed data, source
data, warehouse data and
metadata.

Visualize data across data lake and
data warehouse

Working Group

Dynamo DB for
Metadata

Redshift spectrum

Athena

Quicksight

+++
++
+ + +
+ tableau

S3 buckets for
source data

# API gateway

API tier to disseminate the data.

AWS API gateway

# Next steps…

**Establish continuous data feed to Data Lake**

**Prep Data Lake for data analysis**

**Make it accessible publicly**