

Oracle Data Lakehouse

Robert Korošec, Oracle

Grega Dvoršak, Qubix

Žiga Vaupot, Qubix

2. junij 2022

Intro to the Data Lakehouse

Typical Approach

Data Warehouse



- Solves the problem of analyzing transactional data
- Focus on curated data with known value that is well understood

Data Lake

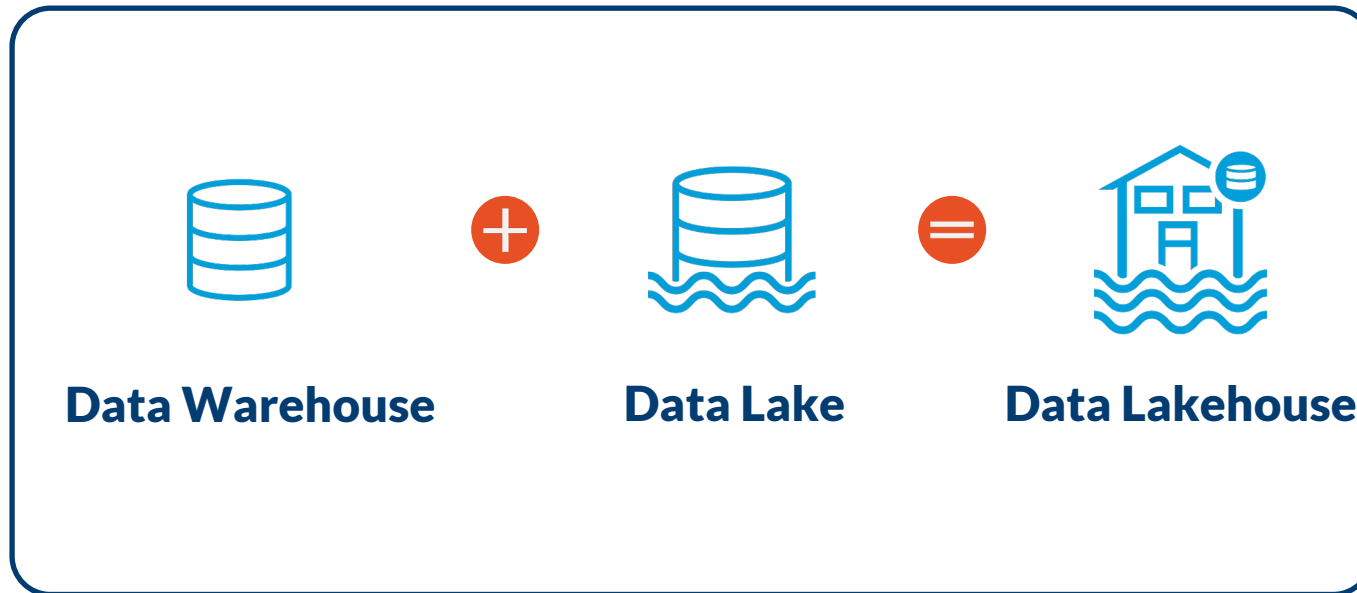


- Handles data that is raw and un-curated. Unknown value or low value
- Open-source tools for processing, analysis and more

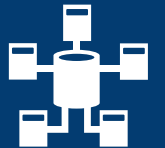


Neither approach alone can answer those questions

Introducing the Data Lakehouse



Integrate the power of a modern **Cloud-based elastic SQL analytic platforms** with modern data design of the **Data Lake** to handle integrated analysis for all data



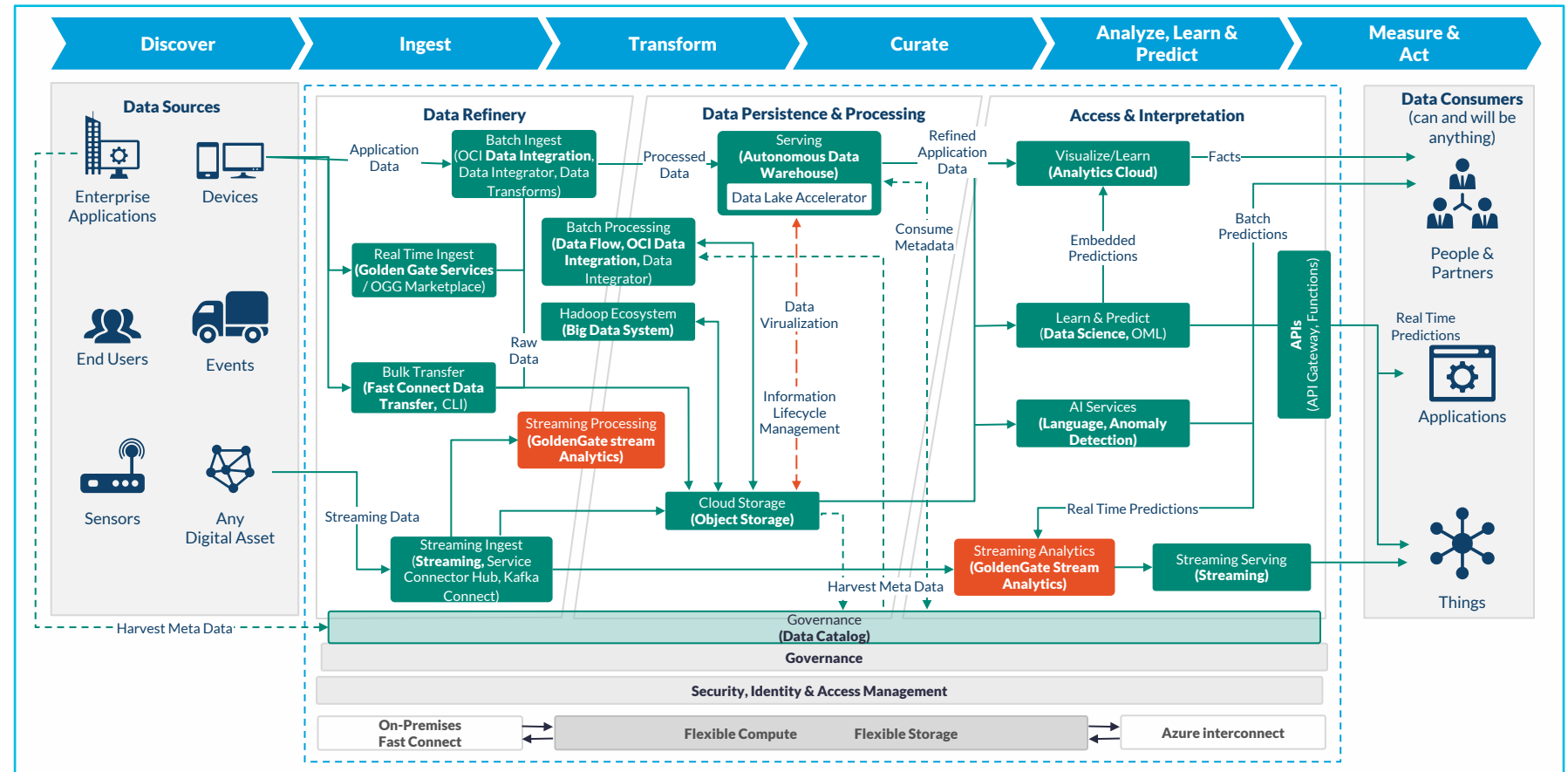
A lakehouse architecture can answer those questions

Oracle Cloud Data Lakehouse – Reference Architecture

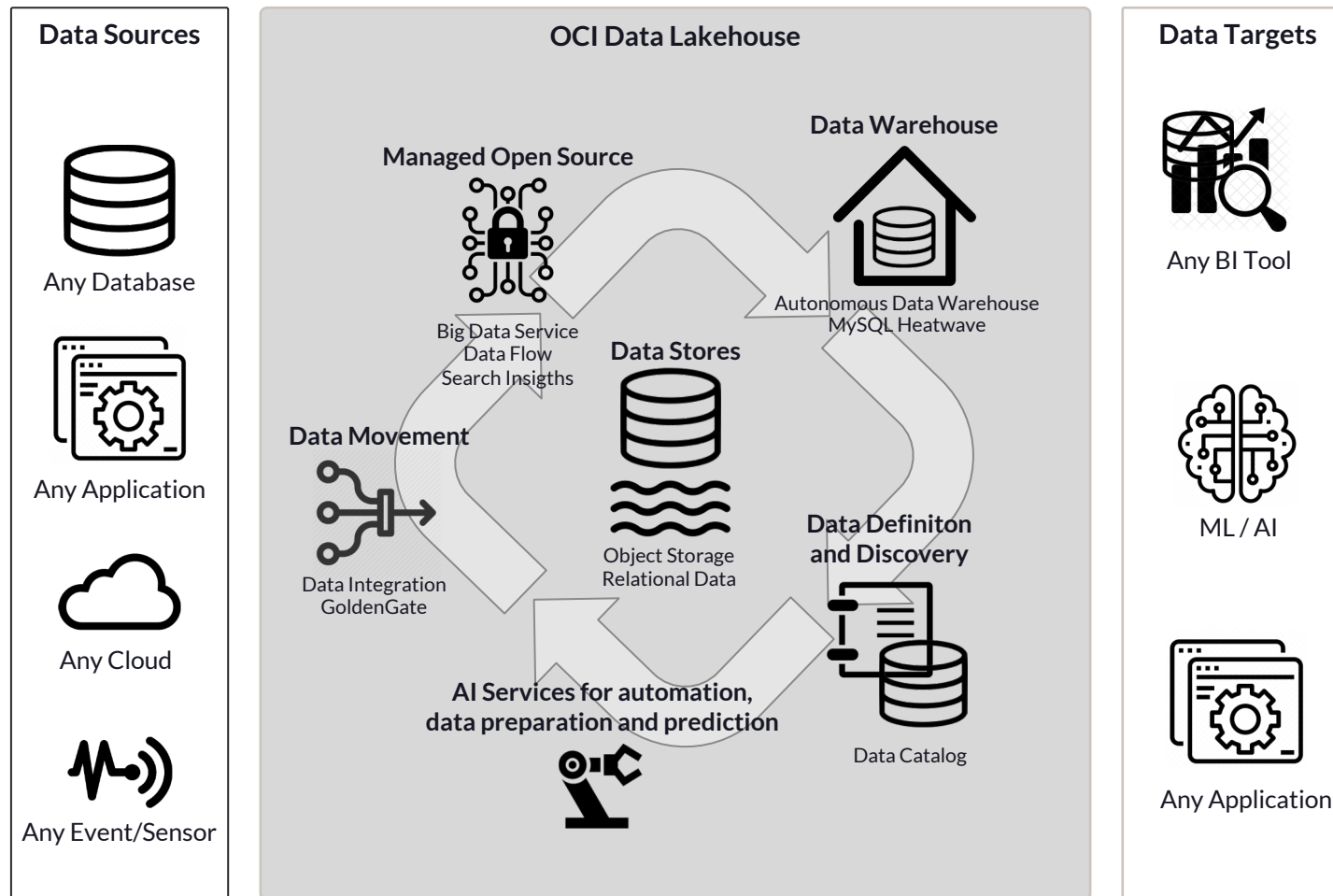
Built on OCI Backbone

Centered around:

- Oracle Autonomous Data Warehouse
- Oracle Machine Learning
- Oracle Cloud Infrastructure Data Integration
- Oracle Cloud Infrastructure Data Flow
- Big Data Service
- Data Catalog
- Oracle Cloud Infrastructure Streaming

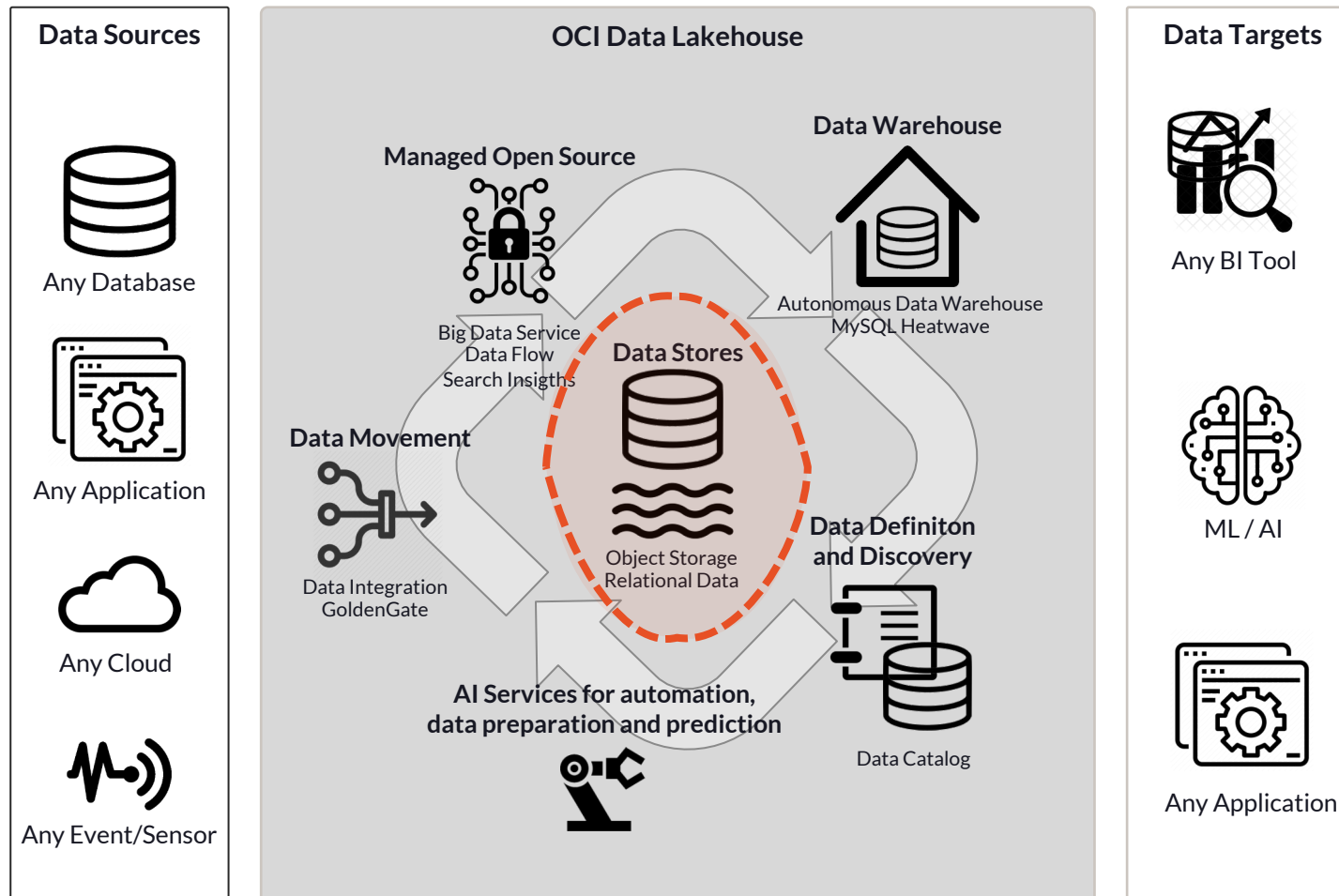


Oracle Data Lakehouse



Data Lakehouse **Object Storage**

Data Stores: Object Storage



- Internet-scale, high-performance storage platform that offers reliable and cost-efficient data durability.
- Stores an unlimited amount of unstructured data of any content type, including analytic data and rich content, like images and videos.

Buckets

Object Storage » Bucket Details



retail_data

Edit Visibility

Move Resource

Bucket Information

Tags

General

Namespace: fryl8pi3k85f

Compartment: [DataLakehouse](#)

Created: Wed, Mar 23, 2022, 12:32:59 UTC

ETag: 3cdb843e-dd88-44f0-b5ff-a23cf893dc87

OCID: ...mqe6xqla [Show](#) [Copy](#)

Usage

Approximate Object Count: 110 objects [i](#)

Approximate Size: 464.9 MIB [i](#)

Uncommitted Multipart Uploads Count: 0 uploads [i](#)

Uncommitted Multipart Uploads Approximate Size: 0 bytes [i](#)

Objects

Location: transactions

Upload

More Actions

Search by prefix

<input type="checkbox"/>	Name	Last Modified	Size	Storage Tier
	↶ ..	-	-	-
<input type="checkbox"/>	transactions_200607.csv	Wed, Mar 23, 2022, 13:25:29 UTC	3.42 MIB	Standard
<input type="checkbox"/>	transactions_200608.csv	Wed, Mar 23, 2022, 13:25:39 UTC	3.54 MIB	Standard
<input type="checkbox"/>	transactions_200609.csv	Wed, Mar 23, 2022, 13:25:49 UTC	3.66 MIB	Standard
<input type="checkbox"/>	transactions_200610.csv	Wed, Mar 23, 2022, 13:25:42 UTC	3.68 MIB	Standard
<input type="checkbox"/>	transactions_200611.csv	Wed, Mar 23, 2022, 13:25:31 UTC	3.54 MIB	Standard
<input type="checkbox"/>	transactions_200612.csv	Wed, Mar 23, 2022, 13:25:55 UTC	3.74 MIB	Standard
<input type="checkbox"/>	transactions_200613.csv	Wed, Mar 23, 2022, 13:25:51 UTC	3.65 MIB	Standard
<input type="checkbox"/>	transactions_200614.csv	Wed, Mar 23, 2022, 13:25:37 UTC	3.7 MIB	Standard

Features

Default Storage Tier: Standard

Visibility: ⚠ Public

Encryption Key: Oracle managed key [Assign](#)

Auto-Tiering: Disabled [Edit](#) [i](#)

Emit Object Events: Disabled [Edit](#) [i](#)

Object Versioning: Disabled [Edit](#) [i](#)

Accessing Object Storage data using SQL

```
1 BEGIN
2 DBMS_CLOUD.CREATE_EXTERNAL_TABLE(
3   table_name => 'RETAIL_TRANSACTIONS',
4   file_uri_list => 'https://objectstorage.eu-frankfurt-1.oraclecloud.com/p/CTuP55zVfoZERhR7Y_fRDh64VXSwJfsHFAGidjRMj19HH0sWG3xyHZ0tB9a@vc_/n/frly8pi3k85f/b/retail_data/o/*.csv',
5   format => json_object('type' value 'csv', 'skipheaders' value '1', 'delimiter' value ','),
6   column_list => ' SHOP_WEEK VARCHAR2 (6) ,
7                   SHOP_DATE VARCHAR2 (8) ,
8                   SHOP_WEEKDAY NUMBER ,
9                   SHOP_HOUR NUMBER ,
10                  QUANTITY NUMBER ,
11                  SPEND NUMBER ,
12                  PROD_CODE VARCHAR2 (10) ,
13                  PROD_CODE_10 VARCHAR2 (7) ,
14                  PROD_CODE_20 VARCHAR2 (8) ,
15                  PROD_CODE_30 VARCHAR2 (6) ,
16                  PROD_CODE_40 VARCHAR2 (6) ,
17                  CUST_CODE VARCHAR2 (14) ,
18                  CUST_PRICE_SENSITIVITY VARCHAR2 (2) ,
19                  CUST_LIFESTAGE VARCHAR2 (2) ,
20                  BASKET_ID NUMBER ,
21                  BASKET_SIZE VARCHAR2 (1) ,
22                  BASKET_PRICE_SENSITIVITY VARCHAR2 (2) ,
23                  BASKET_TYPE VARCHAR2 (20) ,
24                  BASKET_DOMINANT_MISSION VARCHAR2 (20) ,
25                  STORE_CODE VARCHAR2 (10) ,
26                  STORE_FORMAT VARCHAR2 (2) ,
27                  STORE_REGION VARCHAR2 (3) ');
28 END;
```

```
31 SELECT * FROM RETAIL_TRANSACTIONS
32 WHERE SHOP_WEEK = 200607;
```

Query Result Script Output DBMS Output Explain Plan Autotrace SQL History Data Loading

Execution time: 26.881 seconds

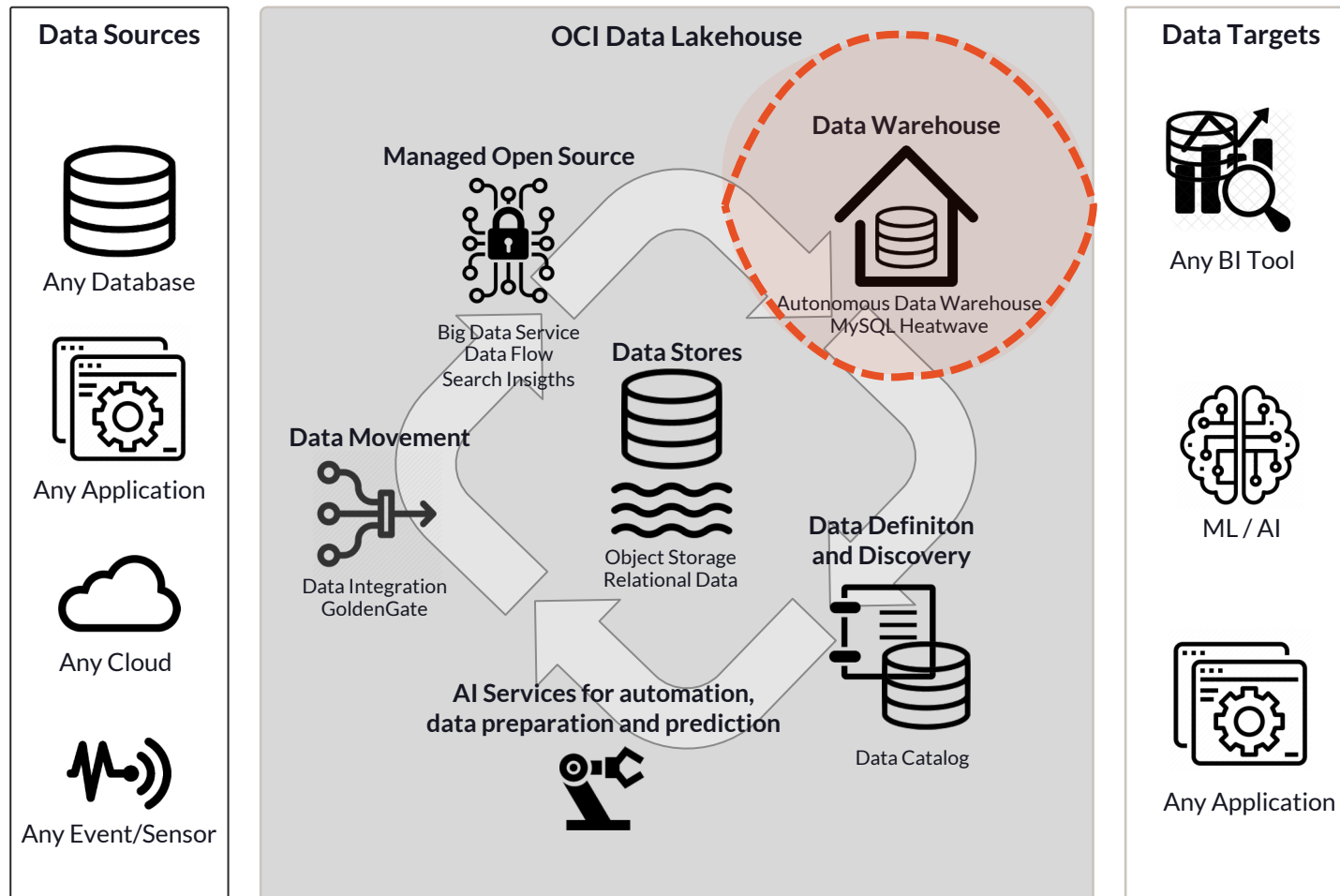
	shop_week	shop_date	shop_weekday	shop_hour	quantity	spend	prod_code	pro
1	200607	20060415	7	19	1	0.93	PRD0900033	CL
2	200607	20060413	5	20	1	1.03	PRD0900097	CL
3	200607	20060416	1	14	1	0.98	PRD0900121	CL
4	200607	20060415	7	19	1	3.07	PRD0900135	CL
5	200607	20060415	7	19	1	4.81	PRD0900220	CL
6	200607	20060412	4	19	1	0.28	PRD0900353	CL
7	200607	20060413	5	18	1	1.56	PRD0900547	CL
8	200607	20060413	5	20	3	0.84	PRD0900550	CL

And finally, bring all into Oracle Analytics



Data Lakehouse
Data Warehouse

Data Warehouse



Autonomous Data Warehouse

- Autonomous Database for Analytics and Data Warehousing provides an easy-to-use, fully autonomous database that scales elastically, delivers fast query performance and requires no database administration.
- Autonomous Database, is deployed on Exadata infrastructure by default

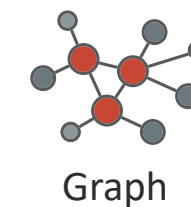
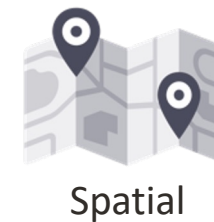
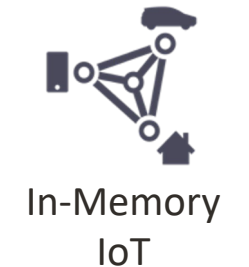
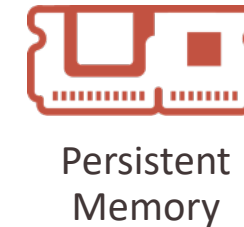
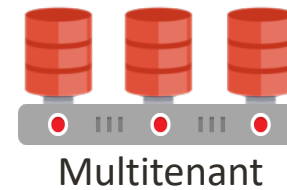
MySQL Heatwave

- HeatWave is a massively parallel, high performance, in-memory query accelerator that accelerates MySQL performance by orders of magnitude for analytics workloads, mixed workloads, and machine learning.

Oracle Autonomous Database

Under the hood – A converged database

- Multitenant for Efficient, Agile Database Clouds
- AutoML for simple integrated Machine Learning
- In-Memory for Database Acceleration
- Native JSON for Document Data
- In-Memory Ingest for Fastest IoT
- Cloud SQL for integrating Object Store Data Lake
- Persistent Memory Store for Lowest Latency
- Spatial and Graph for Mapping and Social Networks

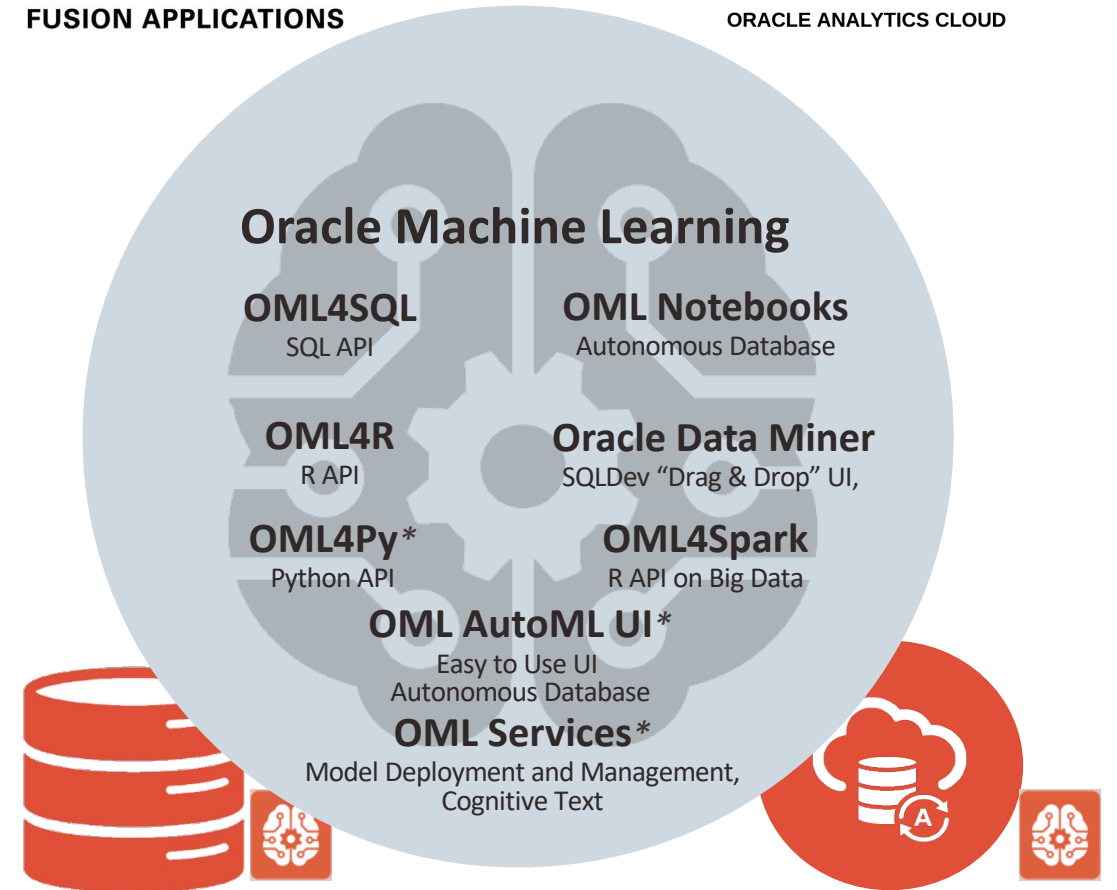


Oracle Machine Learning

- Oracle Machine Learning extends Oracle Database(s) and enables users to build “AI” applications and analytical dashboards
- OML delivers powerful in-database machine learning algorithms, automated ML functionality via SQL APIs and integration with opensource Python* and R.

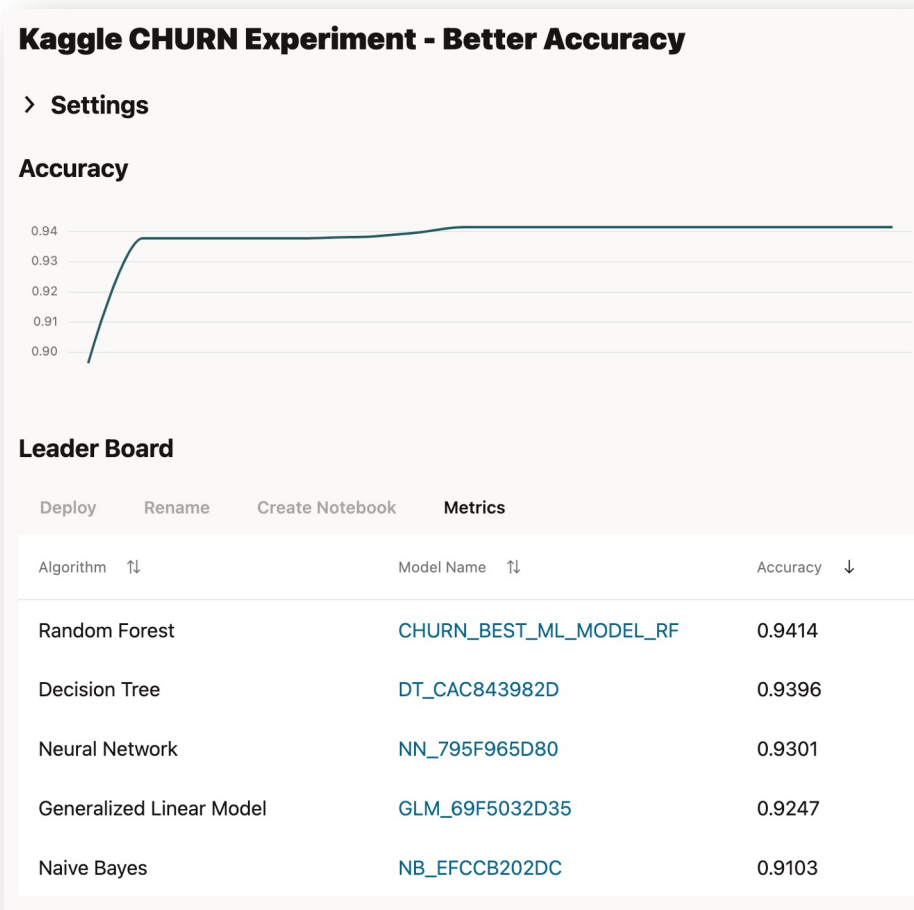
ORACLE®
FUSION APPLICATIONS

ORACLE®
ORACLE ANALYTICS CLOUD



OML AutoML UI

Train the best model using OML Auto UI



Deploy by registering the best model with Oracle Analytics

Select a Model to Register

Search

Type	Name
⌵	AUTOMS_AA8F145E04A4F087
⊕	AUTOMS_E710E48E0EC0326E
⊕	CHURN_BEST_ML_MODEL_RF
⊕	DT_9292407745
⊕	DT_CAC843982D
⊕	DT_CD53399F70
⊕	GLMR_1C906D92AB
⊕	GLMR_E8AF4D7928
⊕	GLM_2B6CCCCF4EA
⊕	GLM_4E3F282E4B

Name: CHURN_BEST_ML_MODEL_RF
Description:

Model Info

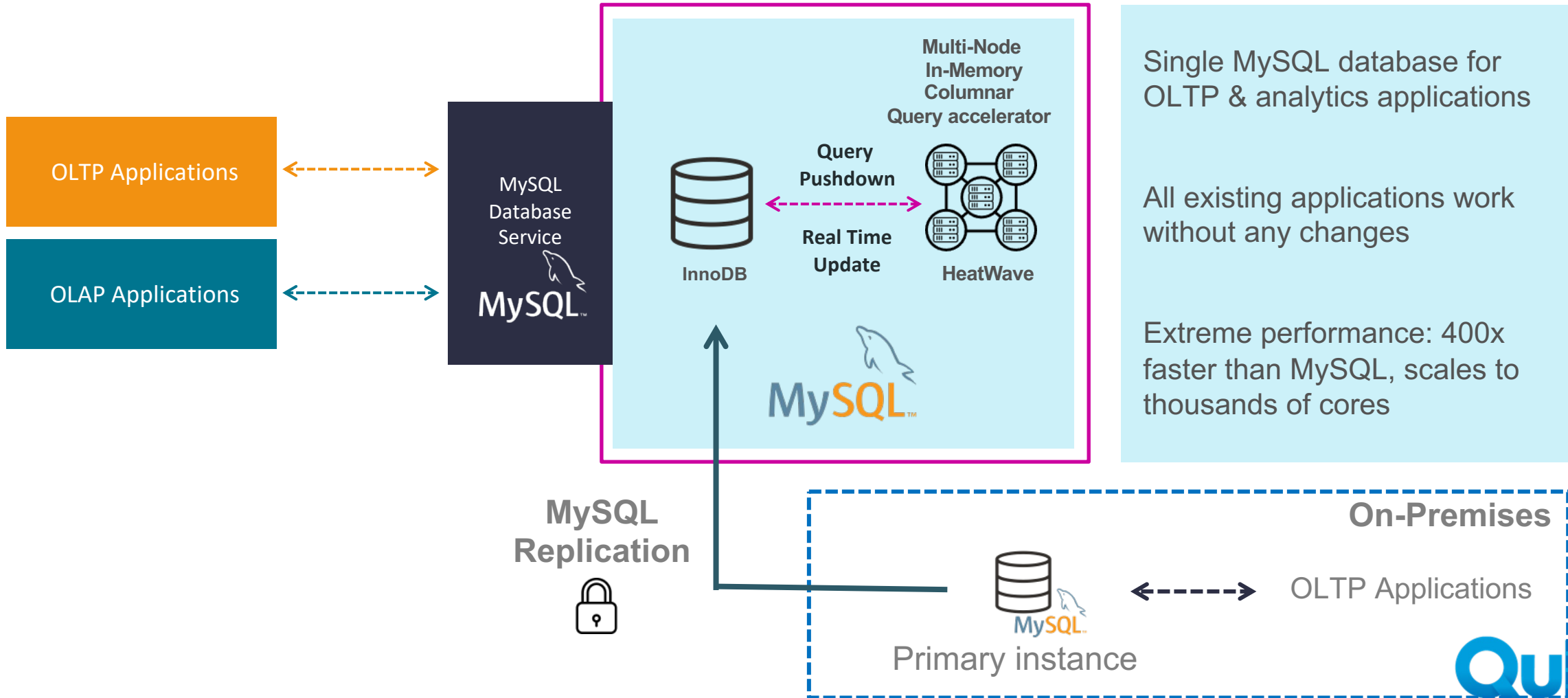
Model Class	CLASSIFICATION
Algorithm	RANDOM_FOREST
DB Model Name	CHURN_BEST_ML_MODEL_RF
DB Model Description	
DB Model Owner	DATALAKEHOUSE
Created On	Yesterday
Target	CHURN

▶ Input Columns
▶ Output Columns
▶ Parameters

Cancel Register

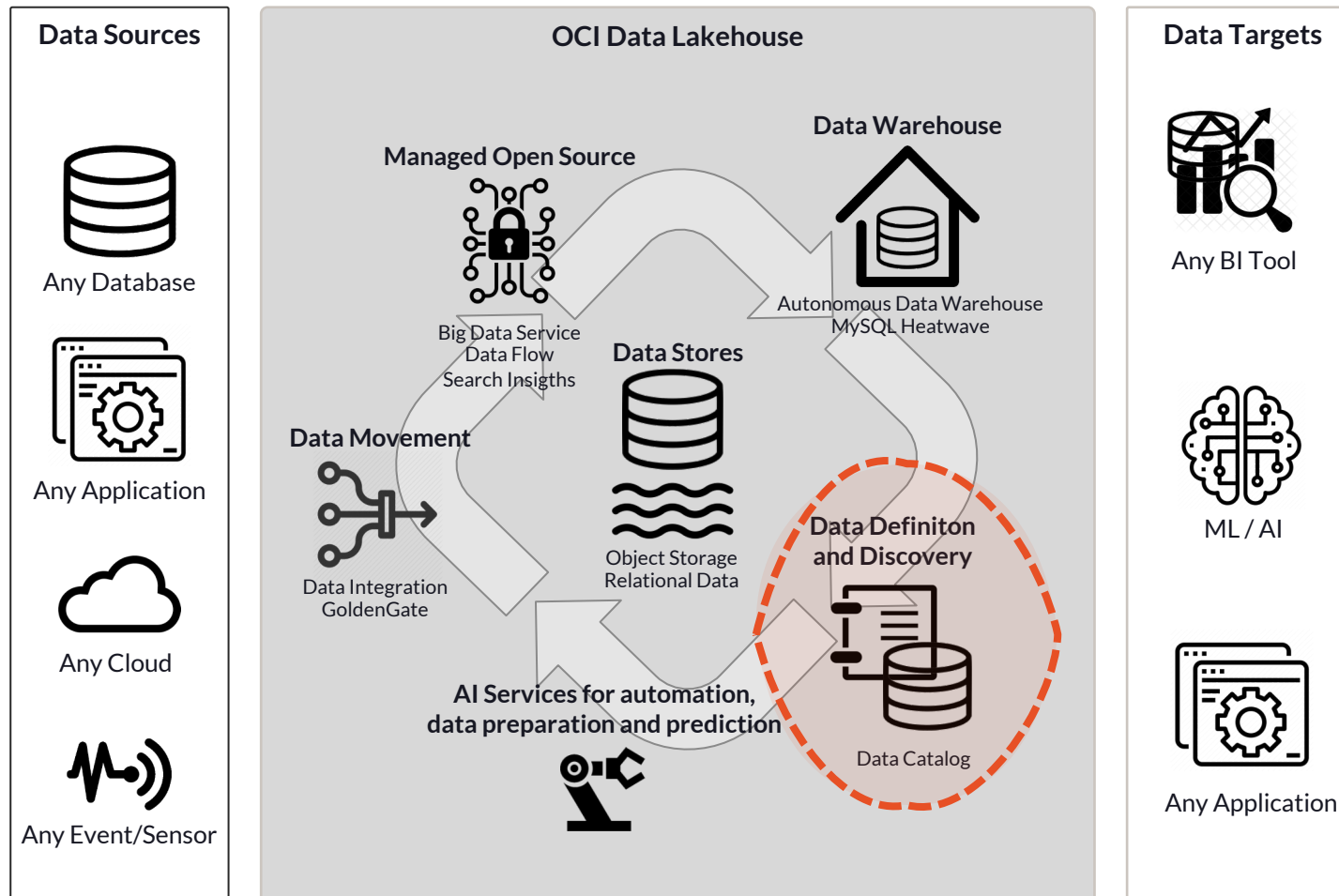
MySQL HeatWave

The only MySQL service with a native massively-scalable query accelerator



Data Lakehouse
OCI Data Catalog

Data Definition and Discovery: Data Catalog



- Metadata management service that helps data professionals discover data and support data governance.
- Designed specifically to work well with the Oracle ecosystem.
- It provides an inventory of assets, a business glossary, and a common metastore for data lakes.

Harvest, Enrich, Search



Harvest metadata from a data source into your catalog. Let Data Catalog discover the available data sources, or you can identify them manually.

[Learn more](#)



Add business context to the technical metadata using business glossary terms and categories, tags, and user defined properties.

[Learn more](#)



Find information about the data available to you by using technical or business names.

[Learn more](#)

Data Catalog

- Discover Data Sources
- Create Data Assets
- Harvesting
- Custom properties
- Glossary and Terms

Search Results for "Churn"

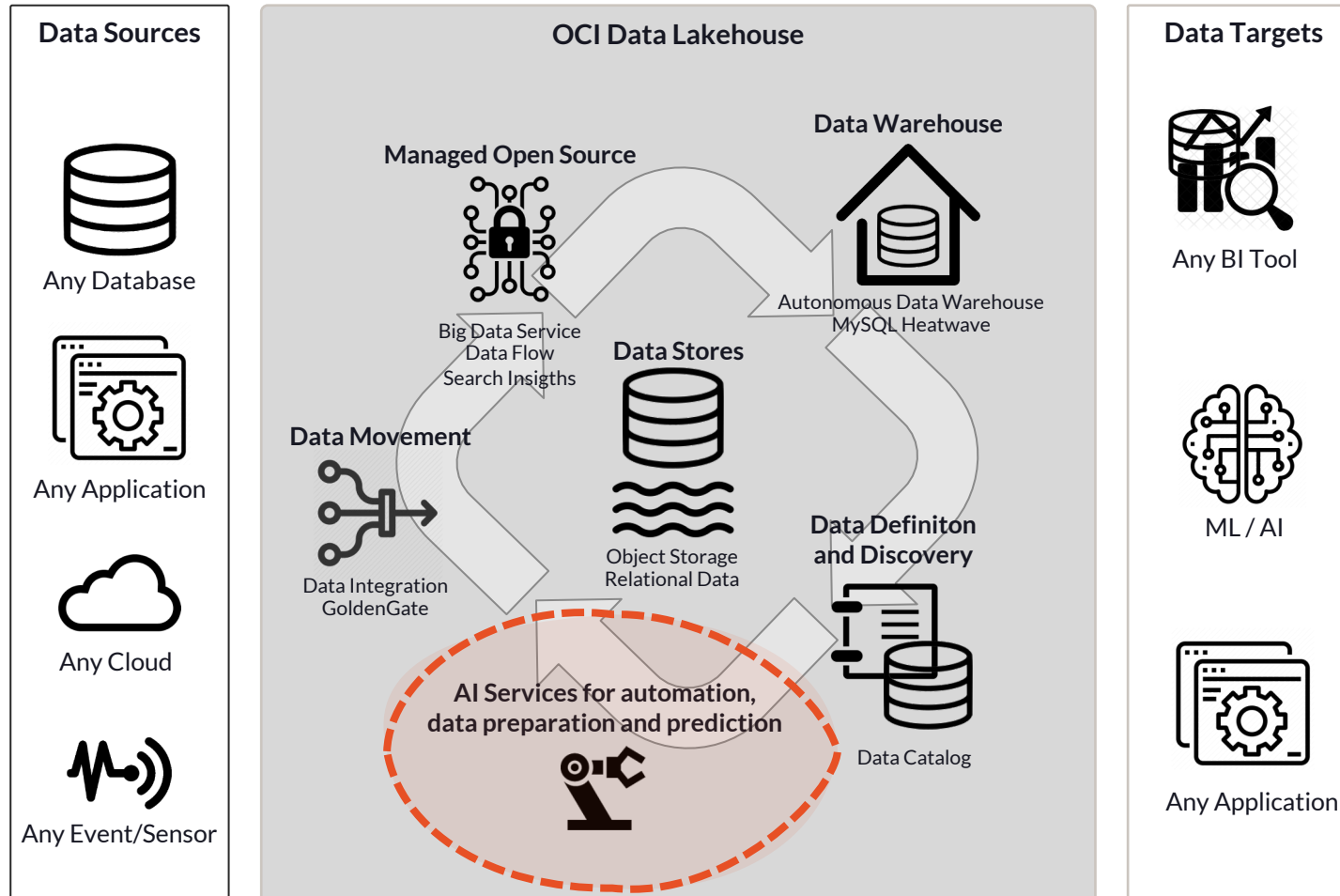
Last updated on Wed, Apr 20, 2022, 12:04:20 PM UTC

<input type="button" value="Link Terms and Categories"/>	<input type="button" value="Link Tags"/>	<input type="button" value="Delete"/>	<input type="button" value="Refresh"/>	<input type="text" value="Churn"/>	<input type="text" value="Relevance"/>
<input type="checkbox"/>	CHURN Path: Data Lakehouse ADW_eu-frankfurt-1 / DATALAKEHOUSE / CHURN_TRAIN_CLEAN Object type: Complex Data asset type: ADW Last updated: Wed, Apr 20, 2022, 09:09 AM UTC				Updated by: oracleidentitycloudservice/ziga... ⋮
<input type="checkbox"/>	CHURN_PREDICTED Path: Data Lakehouse ADW_eu-frankfurt-1 / DATALAKEHOUSE Object type: Table Data asset type: ADW Last updated: Wed, Apr 20, 2022, 09:09 AM UTC				Updated by: oracleidentitycloudservice/ziga... ⋮
<input type="checkbox"/>	CHURN_TEST_CLEAN Path: Data Lakehouse ADW_eu-frankfurt-1 / DATALAKEHOUSE Object type: Table Data asset type: ADW Last updated: Wed, Apr 20, 2022, 09:09 AM UTC				Updated by: oracleidentitycloudservice/ziga... ⋮
<input type="checkbox"/>	CHURN_TRAIN_CLEAN Path: Data Lakehouse ADW_eu-frankfurt-1 / DATALAKEHOUSE Object type: Table Data asset type: ADW Last updated: Wed, Apr 20, 2022, 09:21 AM UTC				Updated by: oracleidentitycloudservice/ziga... ⋮
0 selected				Page 1 of 1 (1-4 of 4 items)	<input type="text" value="1"/>

Data Lakehouse

**AI Services for automation,
data preparation and prediction**

Data Science & AI Services

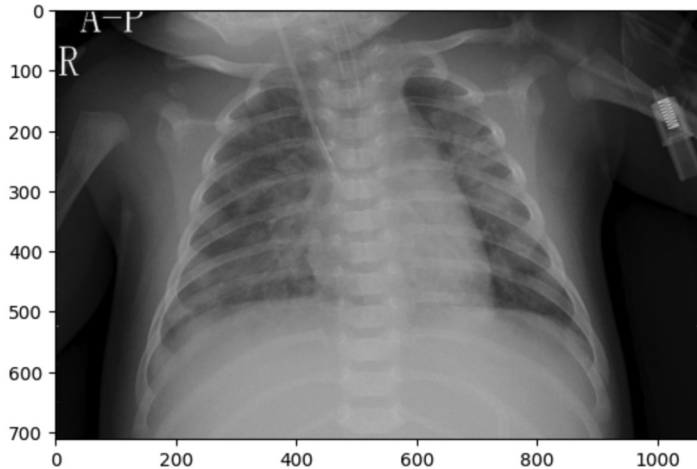


- Data Science Services
 - Data Science
 - Data Labeling
- AI services
 - Vision
 - Speech
 - Language
 - Anomaly Detection
 - Digital Assistant

First, let's start with OCI Data Science

```
[8]: img = cv2.imread("kaggle/input/chest-xray-pneumonia/chest_xray/chest_xray/val/PNEUMONIA/person1946_bacteria_4875.jpeg")
```

```
[9]: plt.imshow(img)  
plt.show()
```



OCI Data Science is a fully managed and serverless platform for data science teams to build, train, and manage machine learning models using Oracle Cloud Infrastructure.

```
[46]: model = Sequential([  
# data_augmentation,  
# layers.Reshape((150,150,1)),  
layers.Conv2D(32, (3,3), strides = 1, input_shape=(150,150,1), padding='same', activation='relu'),  
layers.BatchNormalization(),  
layers.MaxPool2D((2,2), strides = 2, padding = 'same'),  
layers.Conv2D(64, (3,3), strides = 1, padding = 'same', activation = 'relu'),  
layers.Dropout(0.1),  
layers.BatchNormalization(),  
layers.MaxPool2D((2,2), strides = 2, padding = 'same'),  
layers.Conv2D(64, (3,3), strides = 1, padding = 'same', activation = 'relu'),  
layers.BatchNormalization(),  
layers.MaxPool2D(pool_size = (2,2), strides = 2, padding = 'same'),  
layers.Conv2D(128, kernel_size = (3,3), strides = 1, padding = 'same', activation = 'relu'),  
layers.Dropout(0.2),  
layers.BatchNormalization(),  
layers.MaxPool2D(pool_size = (2,2), strides = 2, padding = 'same'),  
layers.Conv2D(256, kernel_size = (3,3), strides = 1, padding = 'same', activation = 'relu'),  
layers.Dropout(0.2),  
layers.BatchNormalization(),  
layers.MaxPool2D(pool_size = (2,2), strides = 2, padding = 'same'),  
layers.Flatten(),  
layers.Dense(units = 128, activation='relu'),  
layers.Dropout(0.2),  
layers.Dense(units = 1, activation = 'sigmoid')  
)  
  
METRICS = [  
keras.metrics.BinaryAccuracy(name='accuracy'),  
keras.metrics.Precision(name='precision'),  
keras.metrics.Recall(name='recall'),  
keras.metrics.AUC(name='auc'),  
)  
)  
  
model.compile(optimizer='rmsprop',  
loss='binary_crossentropy',  
metrics=METRICS)
```

```
[73]: print("Testing set los: \t\t" + str(res_loss))  
print("Testing set accuracy: \t" + str(res_acc*100) + "%")  
print("Testing set precision: \t" + str(res_prec*100) + "%")  
print("Testing set recall: \t" + str(res_rec*100) + "%")  
print("Testing set AUC: \t\t" + str(res_auc*100) + "%")
```

```
Testing set los: 0.24827460944652557  
Testing set accuracy: 92.62820482254028%  
Testing set precision: 91.54589176177979%  
Testing set recall: 97.1794843673706%  
Testing set AUC: 97.02881574630737%
```

Vision AI Service

Object detection

Identifies objects and their location within an image along with a confidence score

Image source

Local file Object storage

Upload image

Drop a file or [select one...](#)

Upload image



Results ⓘ

Labels Raw text

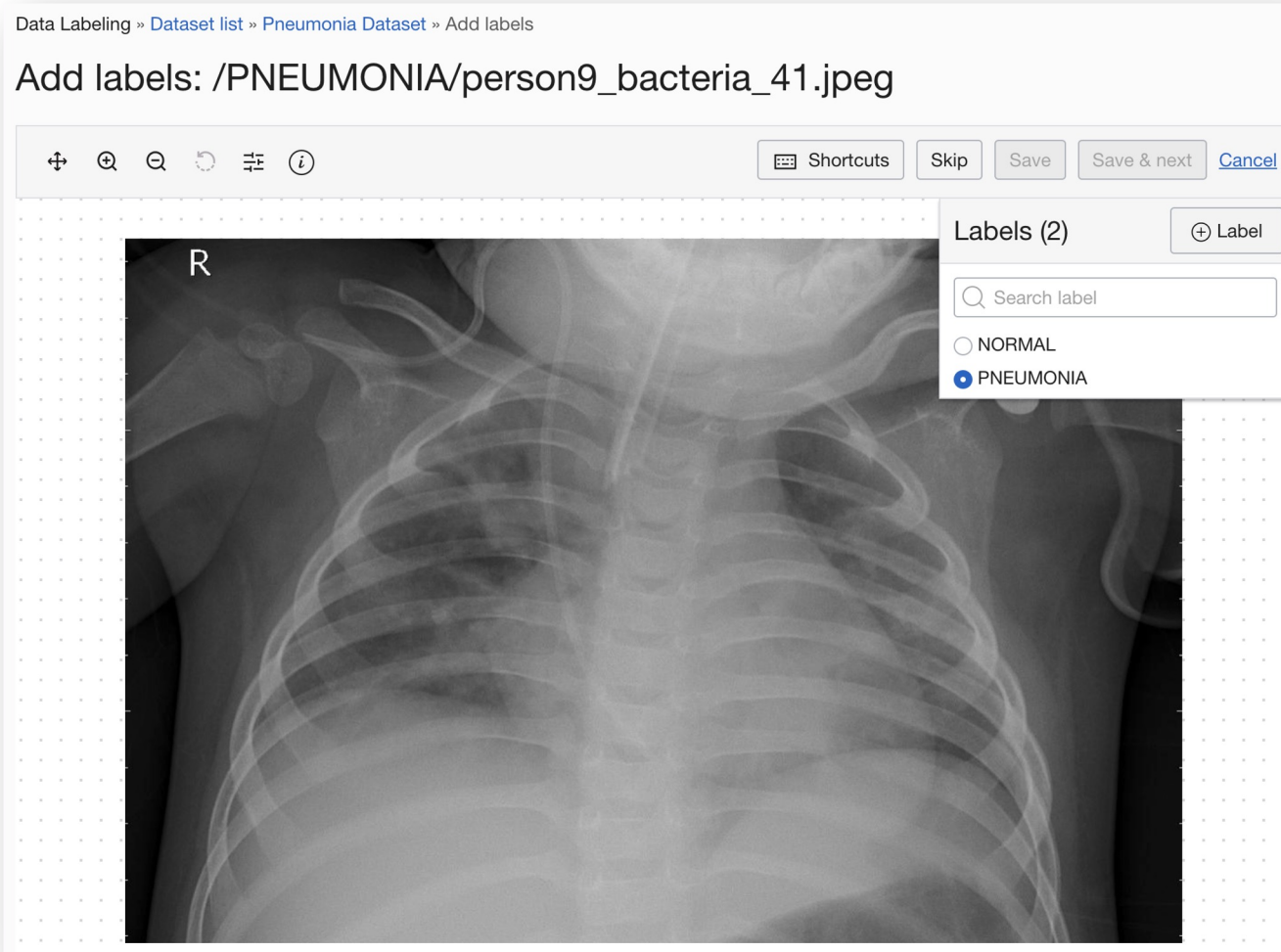
Label	Confidence
Motorcycle	98.90%
Person	98.67%
Helmet	98.13%
Car	97.51%
Glove	96.94%
Boot	96.05%

▶ Request

▶ Response

- OCI Vision is a serverless, multi-tenant service, accessible using the Console, or over REST APIs.
- Upload images to detect and classify objects in them.
- Process images in batch using asynchronous API endpoints

But before anything else, data can be labeled using Data Labeling service



- Label images one by one manually or
- Run bulk labeling utility to label images all in one go

```
import dls_list_records
import dls_create_annotation
import sys
from config import *
import datetime
import labeling_schemes.first_letter as first_letter
import labeling_schemes.first_match as first_match

def main():
    num_records = list_records_limit
    count_batches=1
    count_records_total=0
    while num_records == list_records_limit: # if num_records < list_records_limit,
that would indicate the last loop i.e. batch
        names, ids, num_records = dls_list_records.main()
        count_records_in_batch=0
        for n in names:
            if labeling_algorithm == "first_match":
                first_match.main(name=n,
record_id=ids[count_records_in_batch])
            elif labeling_algorithm == "first_letter":
                first_letter.main(name=n,
record_id=ids[count_records_in_batch])
        count_records_in_batch+=1
        count_records_total+=1
        print("current time: " + str(datetime.datetime.now()))
        print("current batch #: " + str(count_batches))
        print("# records labeled in current batch: " +
str(count_records_in_batch))
        count_batches+=1
        count_batches-=1
        print("----")
        print("TOTALS:")
        print("current time: " + str(datetime.datetime.now()))
        print("# batches: " + str(count_batches))
        print("# records labeled: " + str(count_records_total))
        print()
main()
```

... and then simply train and use the model

Pneumonia Image Classification Model

Edit Move Resource Delete

Model information Tags

General information

Description:

Model type: IMAGE_CLASSIFICATION

OCID: ...uocra [Show](#) [Copy](#)

Training data asset: ...lp3jxabotz34igu3yq [Show](#) [Copy](#)

Created by: oracleidentitycloudservice/ziga.vaupot@qubix.com

Created: Fri, May 20, 2022, 14:38:40 UTC

Max Training duration: Recommended mode (24 hours)

Training metrics

Precision: 0.9737

Recall: 0.9652

F1 score: 0.9694

Total images: 5060

Test images: 506

Trained duration in hours: 5.5525

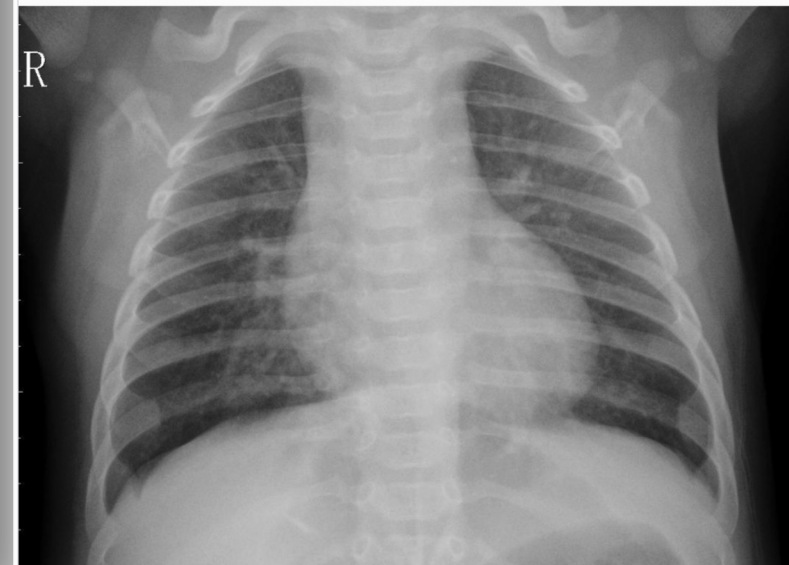
Analyze

Analyze an image to test the newly trained model.

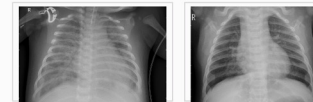
Image source

Local file Object storage

Enter Image URL



Images



Selected: person1000_virus_1681.jpeg

Labels

Label	Confidence
PNEUMONIA	98.20%
NORMAL	1.40%

Request

JSON [Copy](#)

```
{
  "compartmentId": "ocid1.compartment.oc1.",
  "image": {
    "source": "OBJECT_STORAGE",
    "namespaceName": "frly8pi3k85f",
    "bucketName": "Pneumonia",
    "objectName": "train/PNEUMONIA/person1000_virus_1681.jpeg"
  },
  "features": [
    {
      "modelId": "ocid1.aivisionmodel.oc1.",
      "featureType": "IMAGE_CLASSIFICATION",
      "maxResults": 5
    }
  ]
}
```

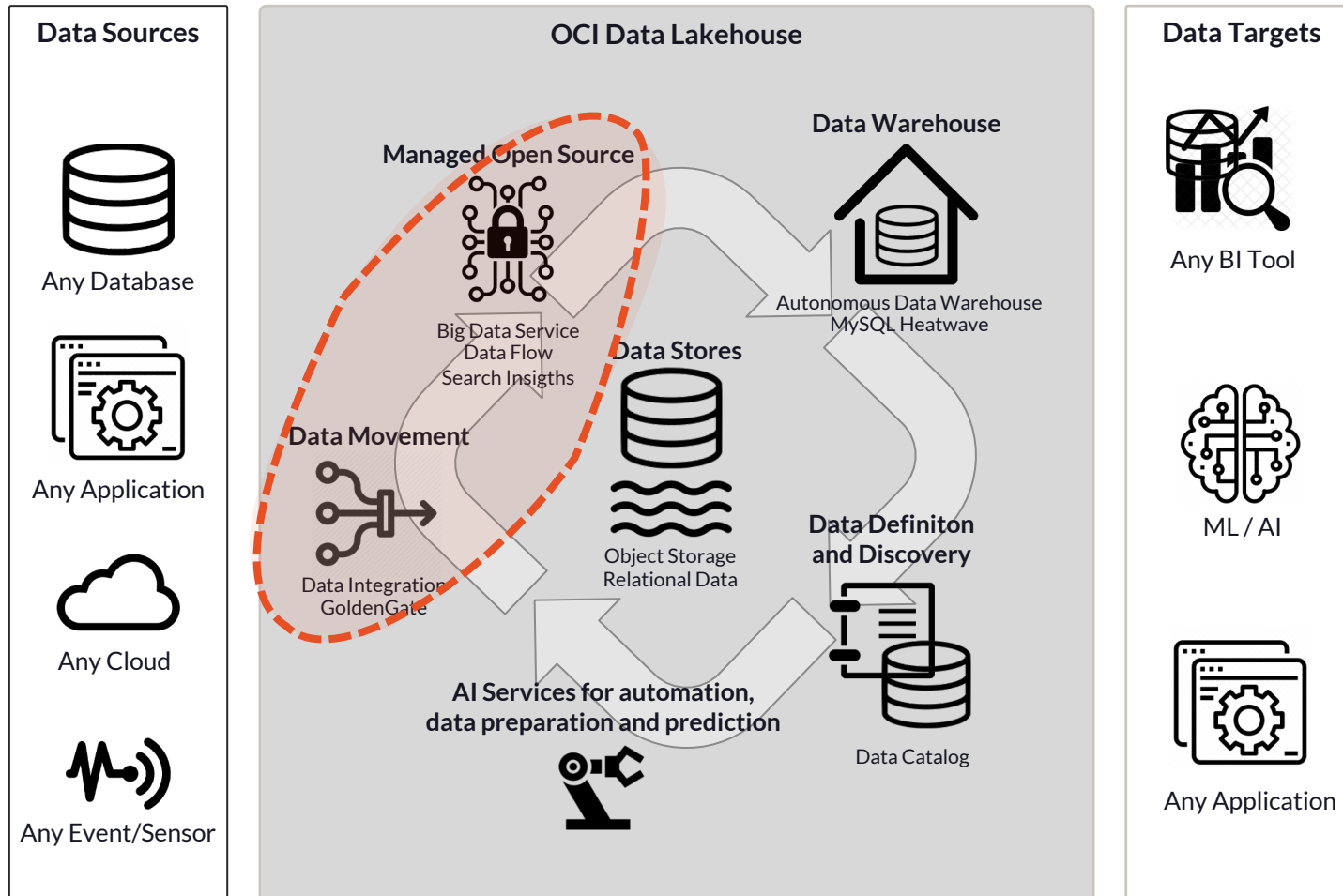
Response

JSON [Copy](#)

```
{
  "imageObjects": null,
  "labels": [
    {
      "name": "PNEUMONIA",
      "confidence": 0.9820319
    }
  ]
}
```

Data Lakehouse **Streaming**

Streaming



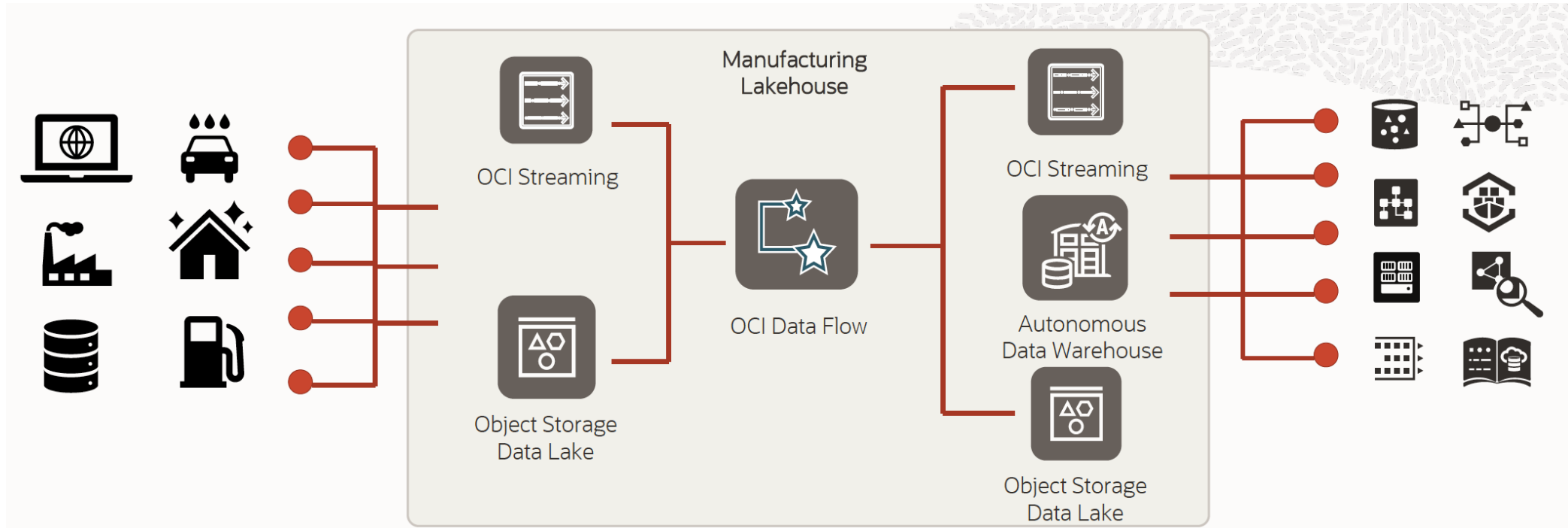
- **OCI Data Flow**

- Fully managed serverless Apache Spark service that supports Spark Streaming.
- Customers can now use OCI Data Flow to do cloud scale ETL on their continuously produced stream data.

- **GoldenGate Stream Analytics**

- Creation of custom operational dashboards that provide real-time monitoring and analyses of event streams in an Apache Spark-based system.
- Customers can identify events of interest in their Apache Spark-based system, execute queries against those event streams in real time and drive operational dashboards or raise alerts based on that analysis.

Oracle Streaming & Data Flow



Demo

Machine Learning on Streaming Data with OCI Data Flow

About the demo

- Use case: predicting remaining useful life of equipment
 - achieved with Data Flow (fully managed serverless Spark as a service)
 - streaming (continuous data streams)
- Other similar use cases:
 - real time fraud detection (financial services)
 - churn prediction (telecommunications)

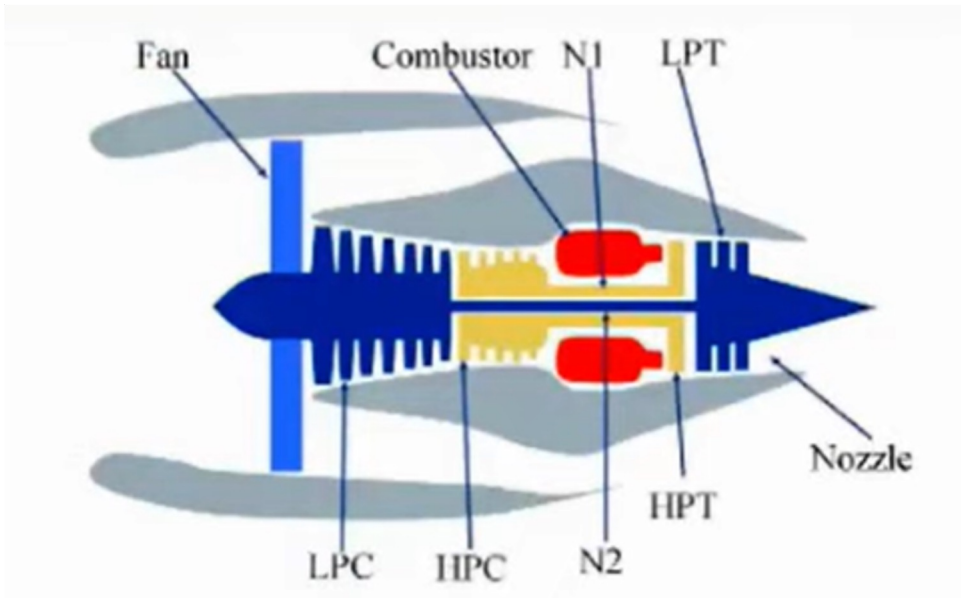
Application details

- Three applications used:
 - model trainer (constructs and trains a machine learning model)
 - simulator (generates sensor data, sends via stream)
 - predictor (uses simulated data from stream and trained model to predict remaining useful life)
- Predictions can be used in further streams or recorded into a database

Color	RUL	Status, Action
Red	0 - 30 days	Degraded, perform maintenance
Yellow	30 - 60 days	Degrading, watch
Green	> 60 days	Healthy, no action needed

Machine learning model and dataset

- ML model: Spark ML Survival Regression Model - Accelerated Failure Model (AFT)
 - available in the Spark ML library
- Dataset: Turbofan Engine Simulation Degradation Dataset
 - from NASA, contains data from a simulation of engine degradation of rocket parts



spark 3.2.1 Overview Programming Guides API Docs Deploying More Search the docs

MLlib: Main Guide

- Basic statistics
- Data sources
- Pipelines
- Extracting, transforming and selecting features
- Classification and Regression
- Clustering
- Collaborative filtering
- Frequent Pattern Mining
- Model selection and tuning
- Advanced topics

MLlib: RDD-based API Guide

- Data types
- Basic statistics
- Classification and regression
- Collaborative filtering
- Clustering
- Dimensionality reduction
- Feature extraction and transformation

Survival regression

In `spark.ml`, we implement the *Accelerated failure time (AFT)* model which is a parametric survival regression model for censored data. It describes a model for the log of survival time, so it's often called a log-linear model for survival analysis. Different from a *Proportional hazards* model designed for the same purpose, the AFT model is easier to parallelize because each instance contributes to the objective function independently.

Given the values of the covariates \mathbf{x}' , for random lifetime t_i of subjects $i = 1, \dots, n$, with possible right-censoring, the likelihood function under the AFT model is given as:

$$L(\beta, \sigma) = \prod_{i=1}^n \left[\frac{1}{\sigma} f_0 \left(\frac{\log t_i - \mathbf{x}'\beta}{\sigma} \right) \right]^{\delta_i} S_0 \left(\frac{\log t_i - \mathbf{x}'\beta}{\sigma} \right)^{1-\delta_i}$$

Where δ_i is the indicator of the event has occurred i.e. uncensored or not. Using $\epsilon_i = \frac{\log t_i - \mathbf{x}'\beta}{\sigma}$, the log-likelihood function assumes the form:

$$\ell(\beta, \sigma) = \sum_{i=1}^n [-\delta_i \log \sigma + \delta_i \log f_0(\epsilon_i) + (1 - \delta_i) \log S_0(\epsilon_i)]$$

Where $S_0(\epsilon_i)$ is the baseline survivor function, and $f_0(\epsilon_i)$ is the corresponding density function.

The most commonly used AFT model is based on the Weibull distribution of the survival time. The Weibull distribution for lifetime corresponds to the extreme value distribution for the log of the lifetime, and the $S_0(\epsilon)$ function is:

$$S_0(\epsilon_i) = \exp(-e^{\epsilon_i})$$

the $f_0(\epsilon_i)$ function is:

$$f_0(\epsilon_i) = e^{\epsilon_i} \exp(-e^{\epsilon_i})$$

The log-likelihood function for AFT model with a Weibull distribution of lifetime is:

$$\ell(\beta, \sigma) = - \sum_{i=1}^n [\delta_i \log \sigma - \delta_i \epsilon_i + e^{\epsilon_i}]$$

Qubix

Enabling the Predictive Enterprise