

A close-up photograph of a person's upper torso. They are wearing a yellow and black horizontally striped shirt. The background is dark and out of focus.

ORACLE

ORACLE

Base de Datos Convergente: Machine Learning, Spatial and Graph Workshop

Manel Moreno

Andrés Araújo

Daniel Villaverde

Francisco Rivas

Francisco Alvarez

5, 6 y 7 de octubre 2021

Zoom sessions

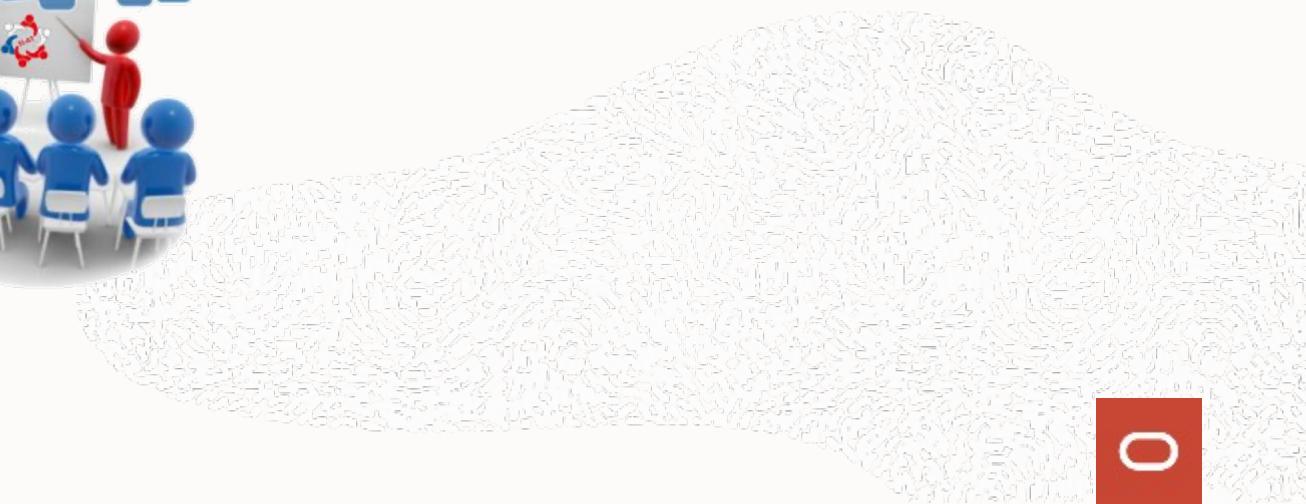
Safe harbor statement

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions.

The development, release, timing, and pricing of any features or functionality described for Oracle's products may change and remains at the sole discretion of Oracle Corporation.



Agenda



Agenda

OVERVIEW

Data Gravity

Moving Data is Slow and Expensive

Leave Data In Place

A diagram featuring a large yellow rock on the left and a red data storage icon on the right. Between them is a red equals sign followed by the text "Data" and a green multiplication sign. To the right of the multiplication sign is a red speedometer-like icon with a needle pointing to zero. Below the equals sign is the equation $e=mc^2$.

Avoid Storing Data in Different
Locations or **Technologies**

It's Easier to **Move Apps to Data**
rather than Moving Data to Apps

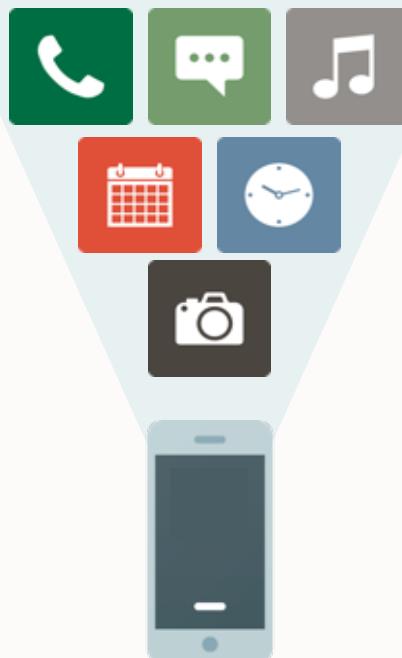
Process Data In Place



Single-Purpose **vs.** Multi-Purpose

Instead of

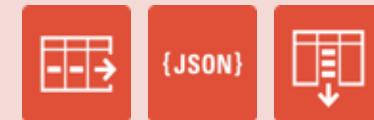
Phones,
Messaging,
Camera, Calendar,
Music, Navigator,
Notepad,
Calculator...



Smart Phone

Instead of

Relational, No-SQL,
JSON, XML,
Transactional,
Analytics, In-Memory,
IoT, ML, Blockchain,
Spatial, Sharding...



Converged Database

Oracle Converged Database

Multi-Model and Multi-Workload

Converged Database

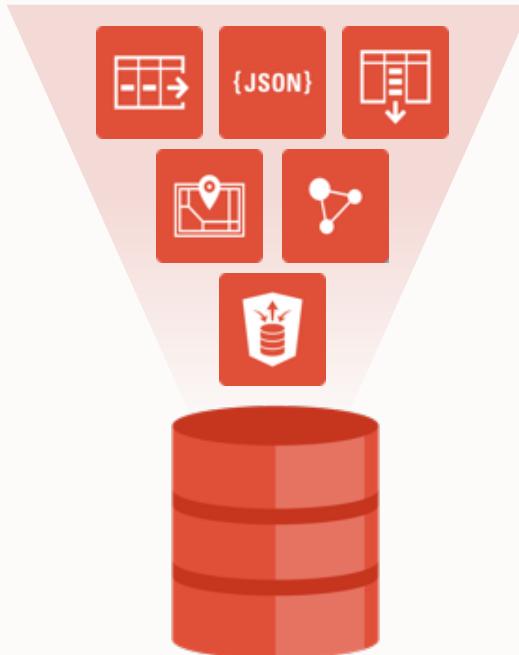
Multi-Model

Multi-Workload

Multiple Data Types **(models and semantics)**

Relational, Document, JSON, XML, OLAP, Spatial, Graph, Object-Oriented, Text, etc.

Multiple Application Types
(workloads and technologies)
Operational, Analytics, **Translytics**, Transactional, IoT, ML, In-Memory, Block-Chain, **HTAP**, etc.

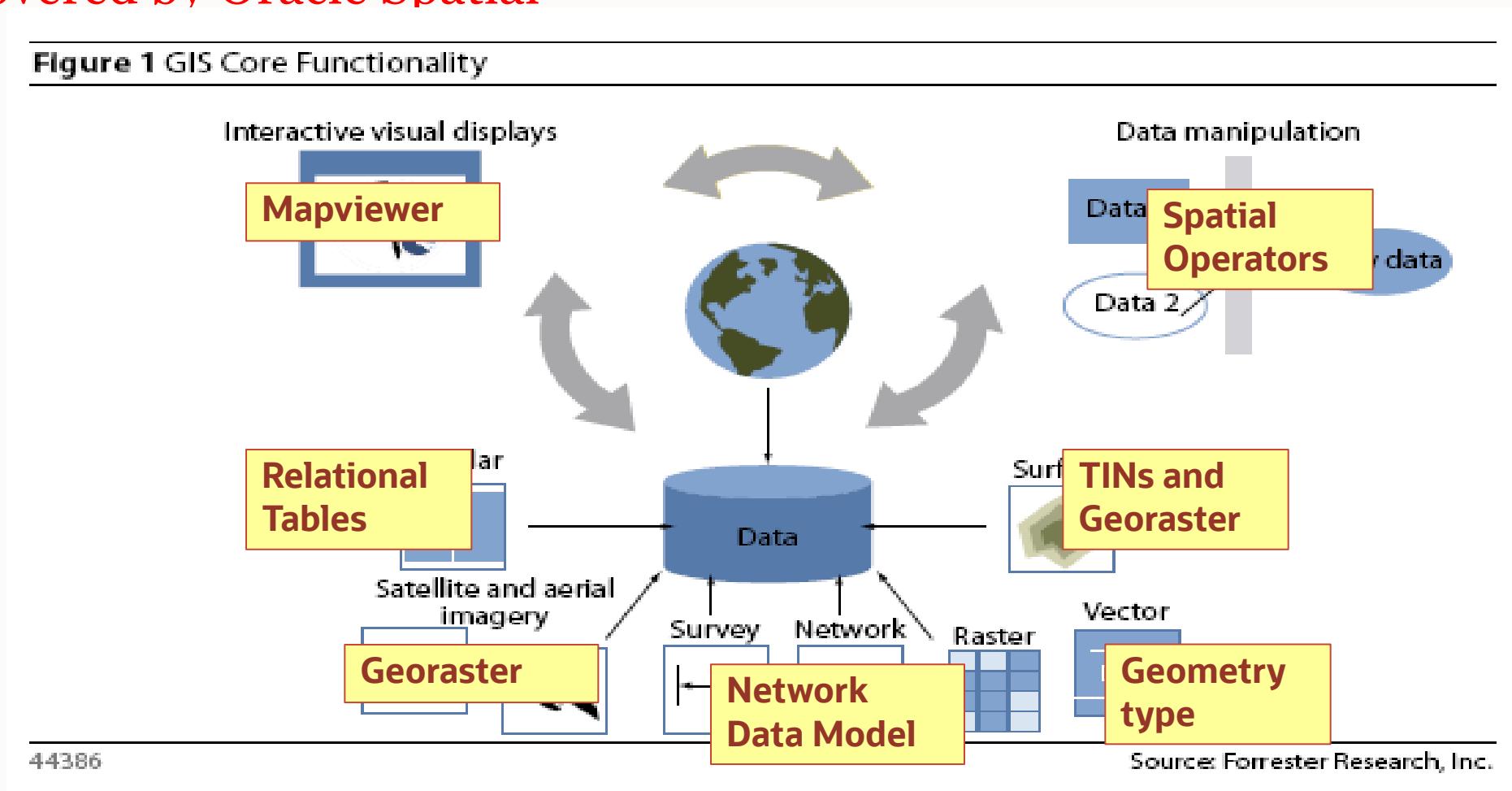


Oracle runs one **Converged Multi-Purpose Database** supporting multiple data types and workloads
Avoid running many **Specialized Single-Purpose Databases** for each data type and workload

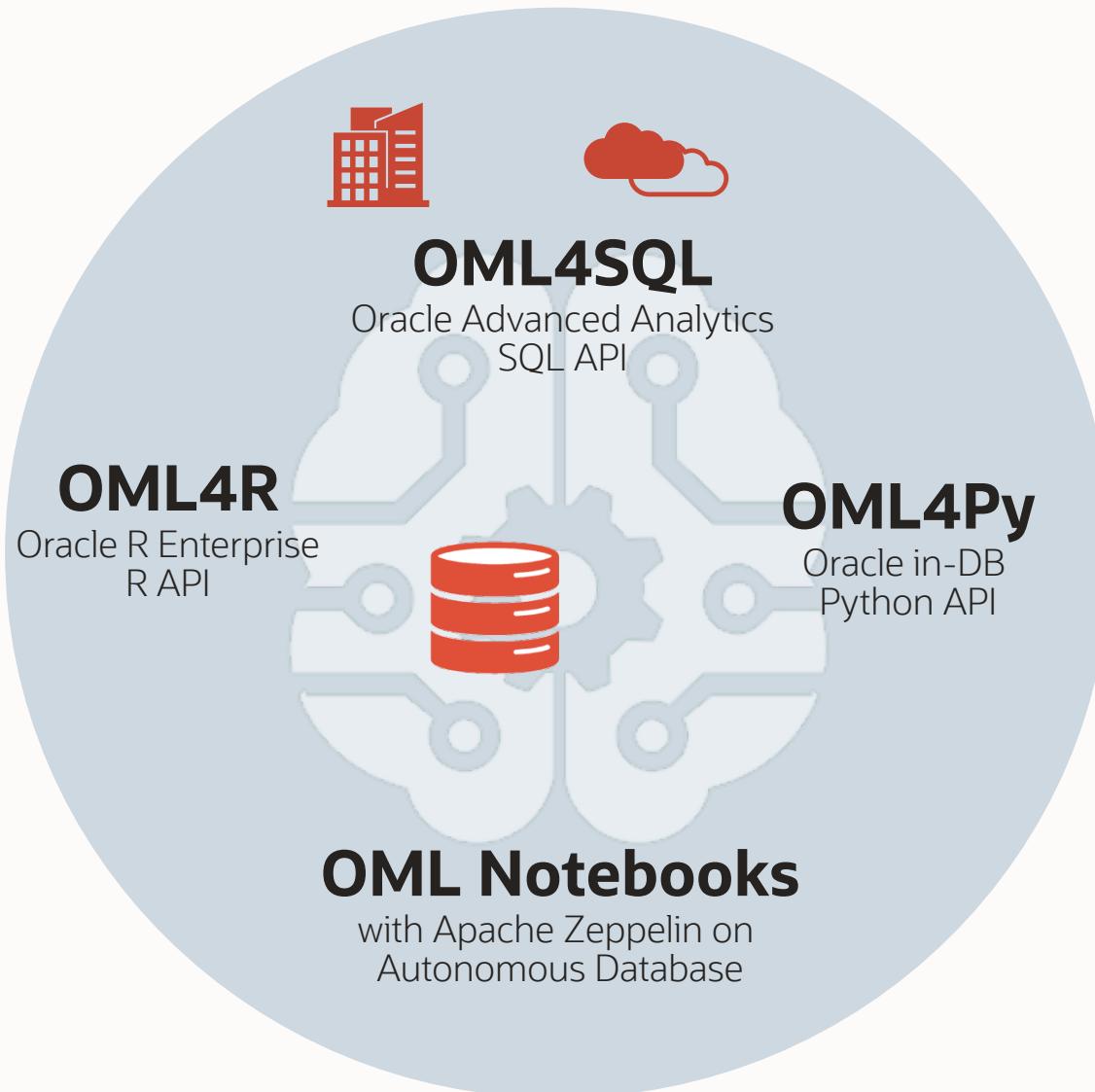
GIS Core Functionality

All Covered by Oracle Spatial

Figure 1 GIS Core Functionality



Oracle Machine Learning

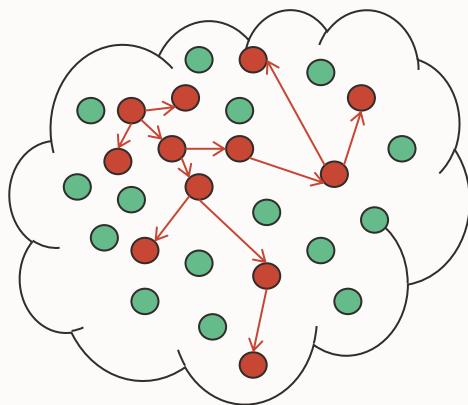


- In-DB Parallel ML Framework
- Python, R or PLSQL
- Cloud Notebook Interface
- Model Lifecycle Management
- Auto-ML and Model Explanation
- Leverage DB Security
- REST and SQL APIs for Scoring



Oracle Machine Learning and Graph Analytics

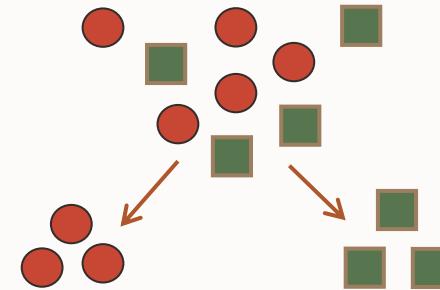
Graph Analytics



Compute graph metric(s)

Explore graph or compute new metrics using ML result

Machine Learning



Build predictive model using graph metric

Use models to score or classify data

Add to structured data

Add to graph

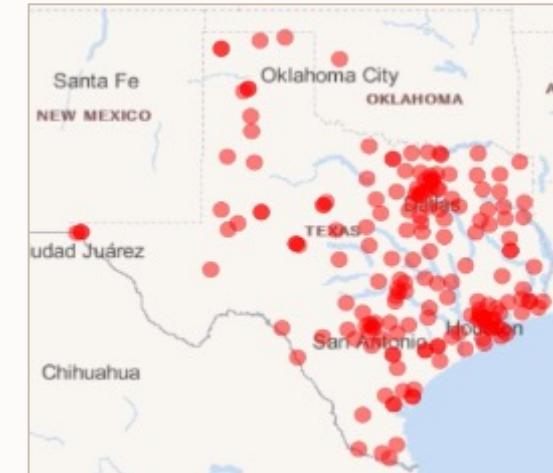


Agenda



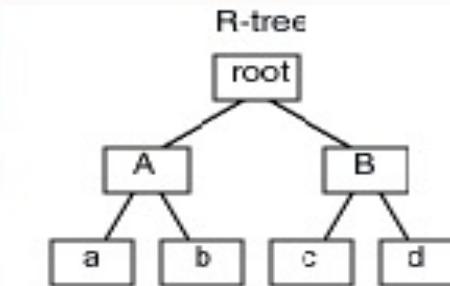
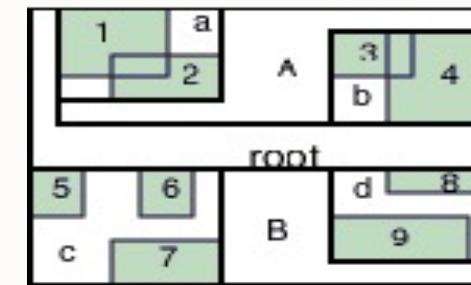
Oracle Spatial– Key Spatial Features

- In-database support for different kinds of geospatial data
- Vector Data (Points, Lines, Linestrings, Areas)
- Geo-referenced Raster Imagery (Orthophotos, Satellite Images, ...)
- 3D Point Cloud Data (Laser scanning, Photogrammetry)
- Network Data (Road Networks, Utility Networks)
- Topology Data (Land management)
- Streaming Point Data (Location tracking)
- Deployable Services
- Map visualization
- Geocoding
- Routing
- Publishing (OGC Web Services)



Database Capabilities for Geospatial Analysis

- Data type to store points, lines, areas, solids, ...
 - In two or three dimensions
 - Taking into account coordinate system
- Topological operators
 - Point-in-polygon, intersecting linestrings, overlapping areas, ...
- Geometric functions
 - Calculating areas, distances, buffer zones, ...
- Spatial indices
 - Fast access to relevant data



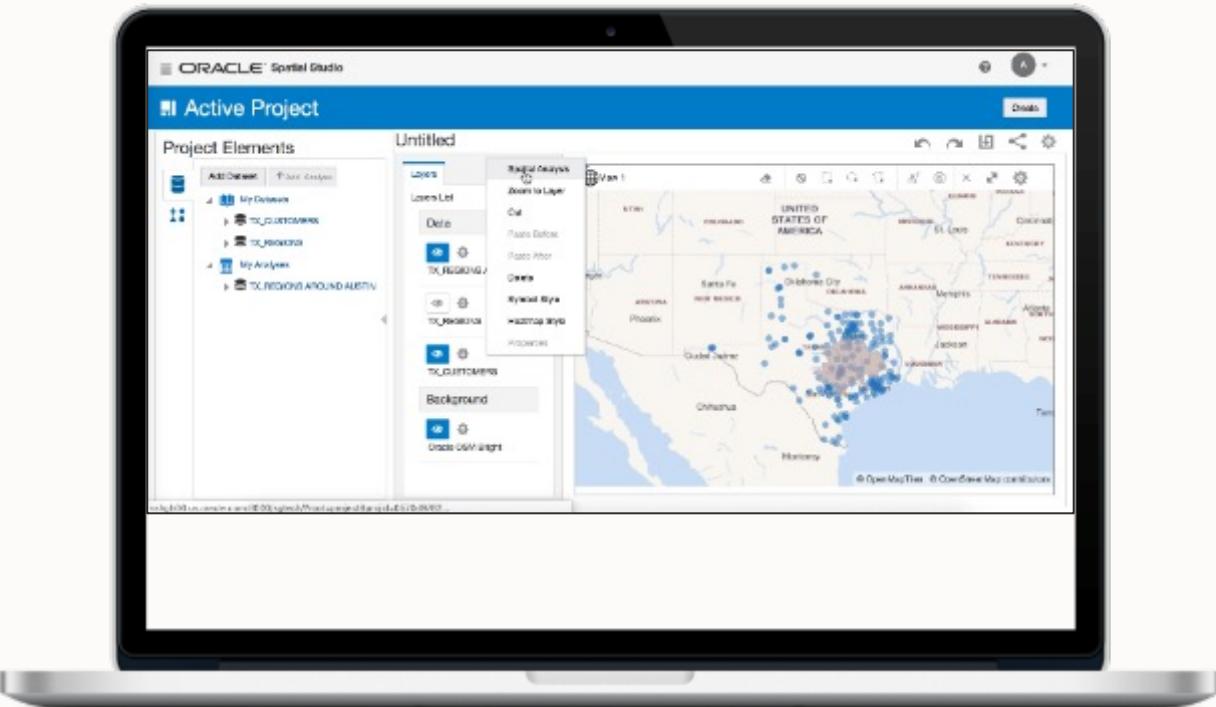
```
SELECT a.owner_name, a.acquisition_status  
FROM properties a, projects b  
WHERE sdo_within_distance (a.property_geom1,  
b.project_geom,  
    'distance = 25 unit = meter') = 'TRUE'  
and b.project_id=189498;
```

Benefits of Managing Spatial Data in Oracle DB

- **Multi-model database, integrating all kinds of data**
 - Relational data, XML or JSON documents, spatial data, images, ...
- **Comprehensive server-side ETL and analytics capabilities**
 - Data integration, geospatial analysis, machine learning, graph analysis, ...
- **Secure datastore**
 - Multi-level access control, encryption, redaction, auditing, ...
- **Highly available, scalable infrastructure**
 - Clustering, parallelization, Maximum Availability Architecture (MAA), ...
- **Core component of data management platform for analytics**
 - Tools integration, standards support, open interfaces, Big Data connectivity, ...

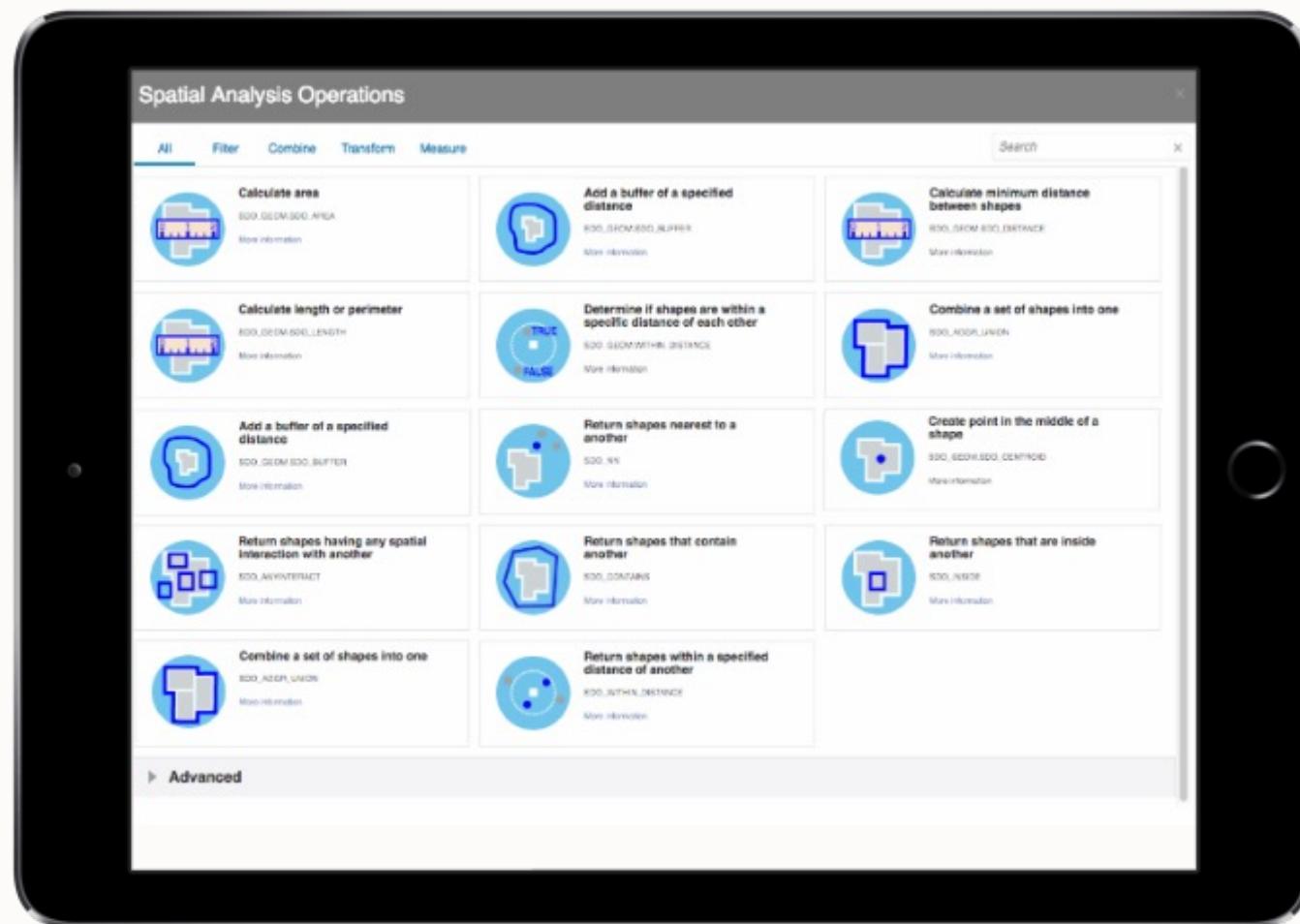
Typical Data Analysis Workflow

- Data ingestion
 - Spatial and non-spatial data
- Data enrichment
 - Address geocoding
 - Converting placenames
- Geospatial processing
 - Creating analytical workflows
- Interactive analysis
 - Map visualization
- Publication of results



Spatial Studio – Self-service spatial analytics

Spatial Studio – Simple Geospatial Analysis



Major New Spatial Features

Ease of Use

- Spatial Studio - Self-service development tool
- Improved JSON and Oracle REST Data Services
- Enhanced Location Tracking Server
- Map Visualization
- Improved web services support (CSW, WFS)
- Georaster enhancements

Performance

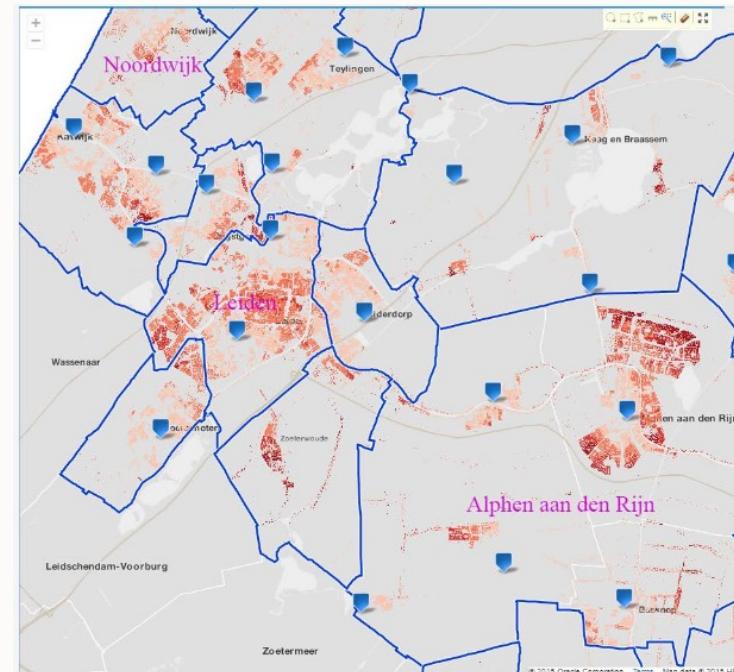
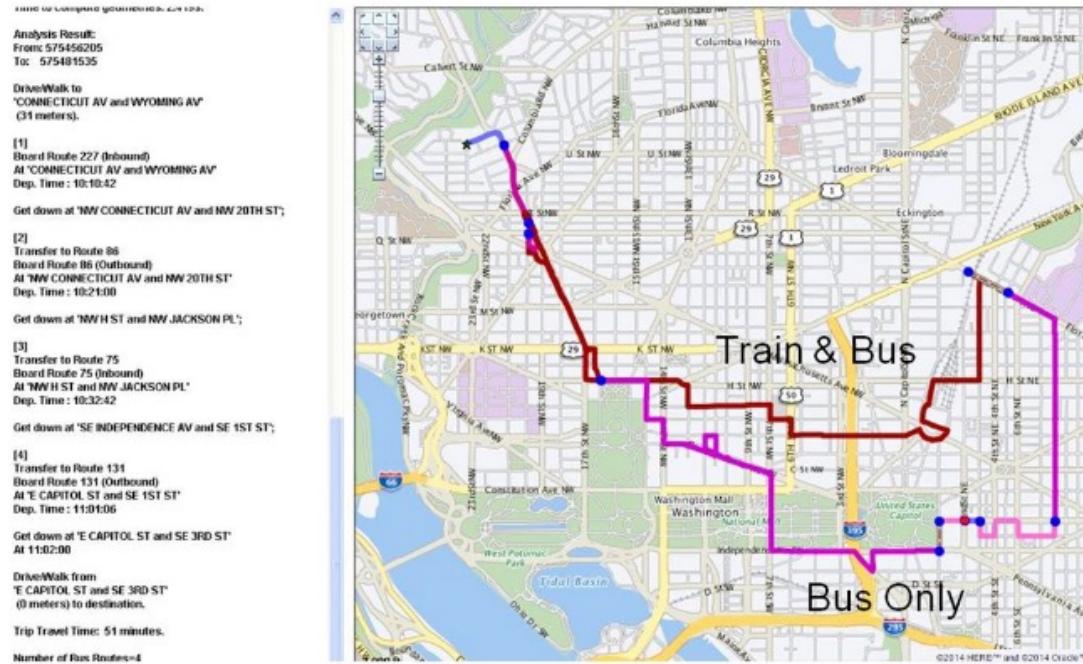
- Spatial index performance improvements
 - 3x faster queries for large point data sets
- Map visualization dynamic tile layer
 - Save storage overhead on large, complex queries

Improved Database Integration

- Spatial support for all partitioning methods
- Spatial support for distributed transactions
- Spatial support for database sharding
- Improved support for queries on external tables

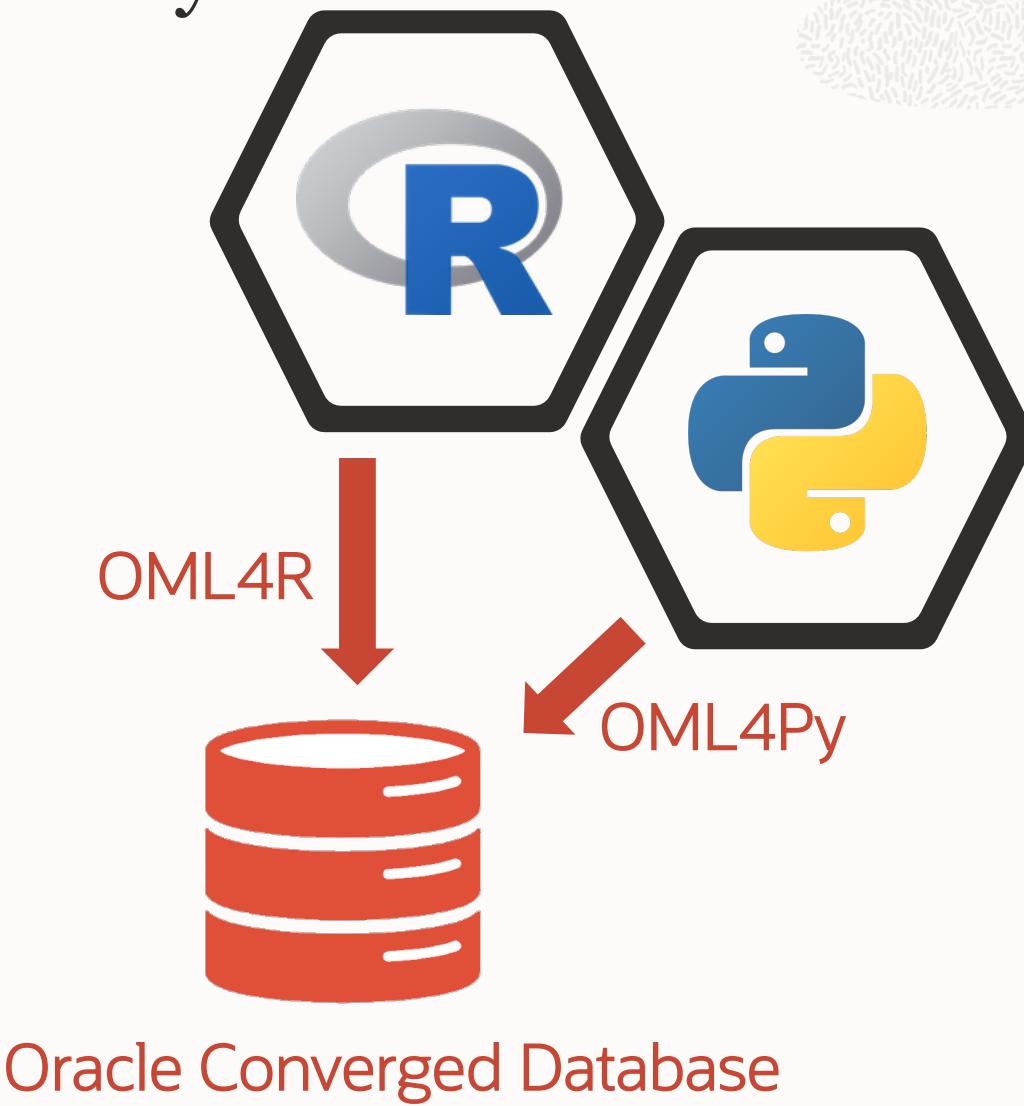
Advanced Spatial Data Models

- **Spatial networks for roads, transport, pipelines, telcos and other geographically connected analysis**
- **Topology for mapping, land management and cadastre applications**

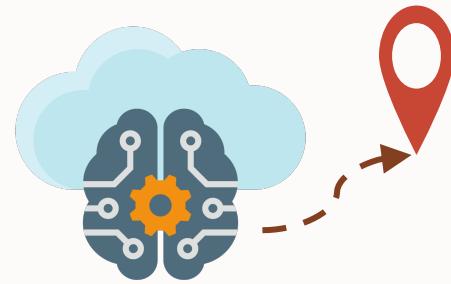


Why data scientists and data analysts use R and Python

- Powerful
- Extensible
- Graphical
- Extensive statistics
- Ease of installation and use
- Rich ecosystem
 - 1000s of open source packages
 - Millions of users worldwide
- Heavily used by data scientists
- Free

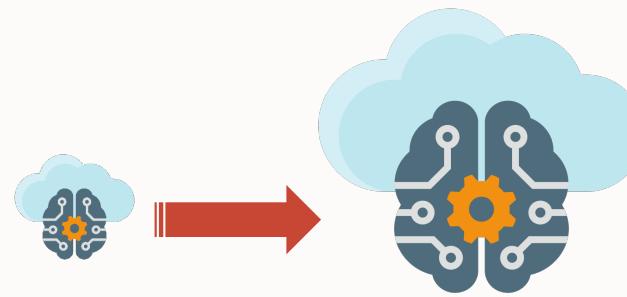


Oracle Machine Learning Key Attributes



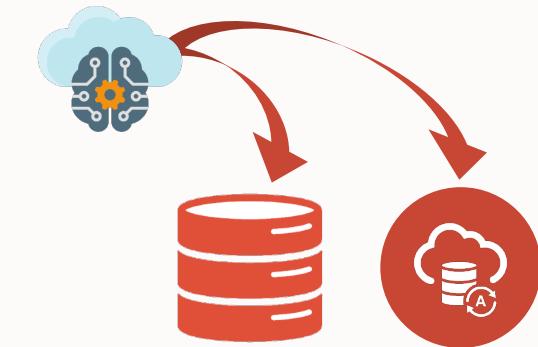
Automated

Get better results faster
with less effort –
even non-expert users



Scalable

Handle big data volumes using
parallel, distributed algorithms –
no data movement

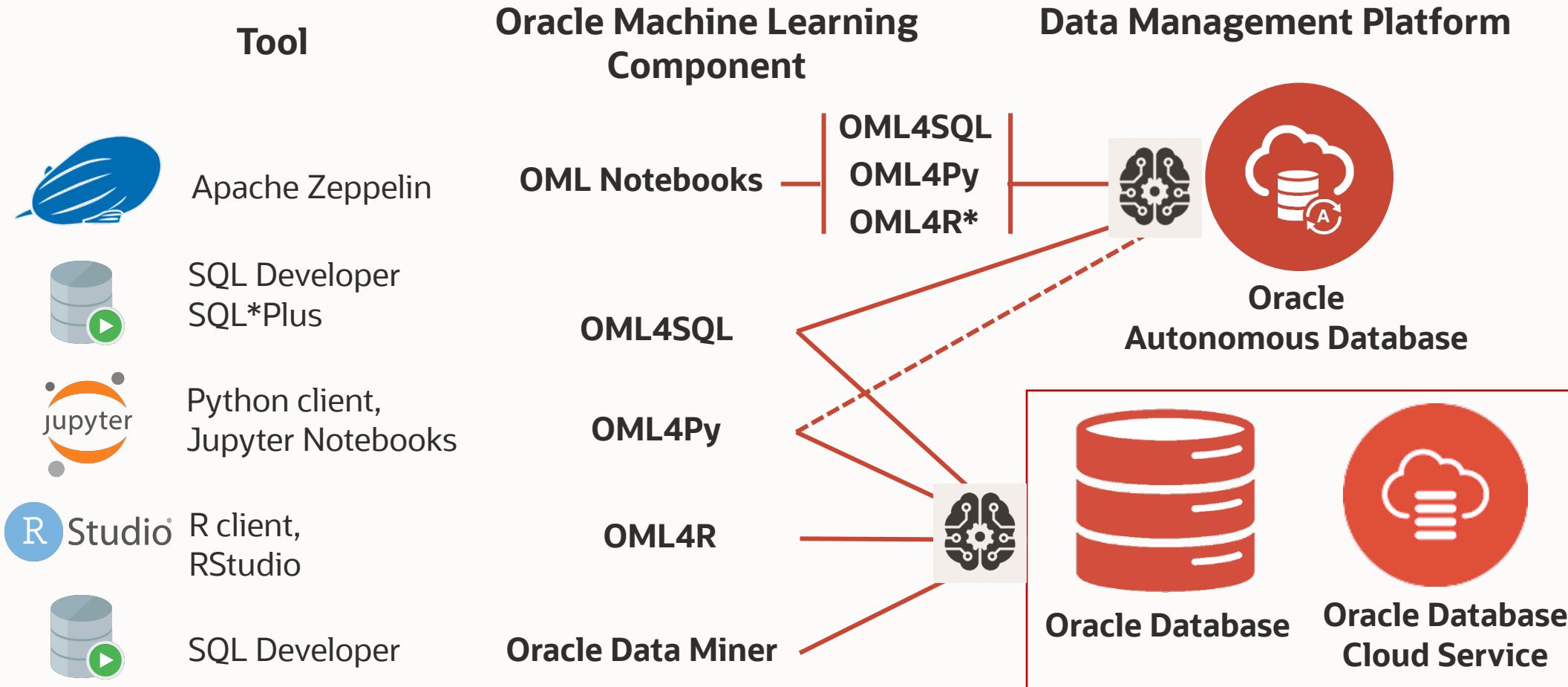


Production-ready

Deploy and update data
science solutions faster with
integrated ML platform

Increase productivity | Achieve enterprise goals | Innovate More

Oracle Machine Learning interfaces to Oracle Database



* coming soon

Oracle Machine Learning Algorithms and Analytics

• CLASSIFICATION

- Naïve Bayes
- Logistic Regression (GLM)
- Decision Tree
- Random Forest
- Neural Network
- Support Vector Machine (SVM)
- Explicit Semantic Analysis

• CLUSTERING

- Hierarchical K-Means
- Hierarchical O-Cluster
- Expectation Maximization (EM)

• ANOMALY DETECTION

- One-Class SVM

• TIME SERIES

- Forecasting - Exponential Smoothing
- Includes popular models
e.g. Holt-Winters with trends,
seasonality, irregularity, missing data

REGRESSION

- Linear Model
- Generalized Linear Model (GLM)
- Support Vector Machine (SVM)
- Stepwise Linear regression
- Neural Network
- LASSO

ATTRIBUTE IMPORTANCE

- Minimum Description Length
- Principal Component Analysis (PCA)
- Unsupervised Pair-wise KL Div
- CUR decomposition for row & AI

ASSOCIATION RULES

- A priori/ market basket

PREDICTIVE QUERIES

- Predict, cluster, detect, features

SQL ANALYTICS

- SQL Windows
- SQL Patterns
- SQL Aggregates

XGBoost
MSET

• FEATURE EXTRACTION

- Principal Comp Analysis (PCA)
- Non-negative Matrix Factorization
- Singular Value Decomposition (SVD)
- Explicit Semantic Analysis (ESA)

• TEXT MINING SUPPORT

- Algorithms support text columns
- Tokenization and theme extraction
- Explicit Semantic Analysis (ESA) for document similarity

• STATISTICAL FUNCTIONS

- Basic statistics: min, max, median, stdev, t-test, F-test, Pearson's, Chi-Sq, ANOVA, etc.

R AND PYTHON PACKAGES

- Third-party R and Python Packages through Embedded Execution
- Spark MLlib algorithm integration



Oracle Machine Learning Notebooks



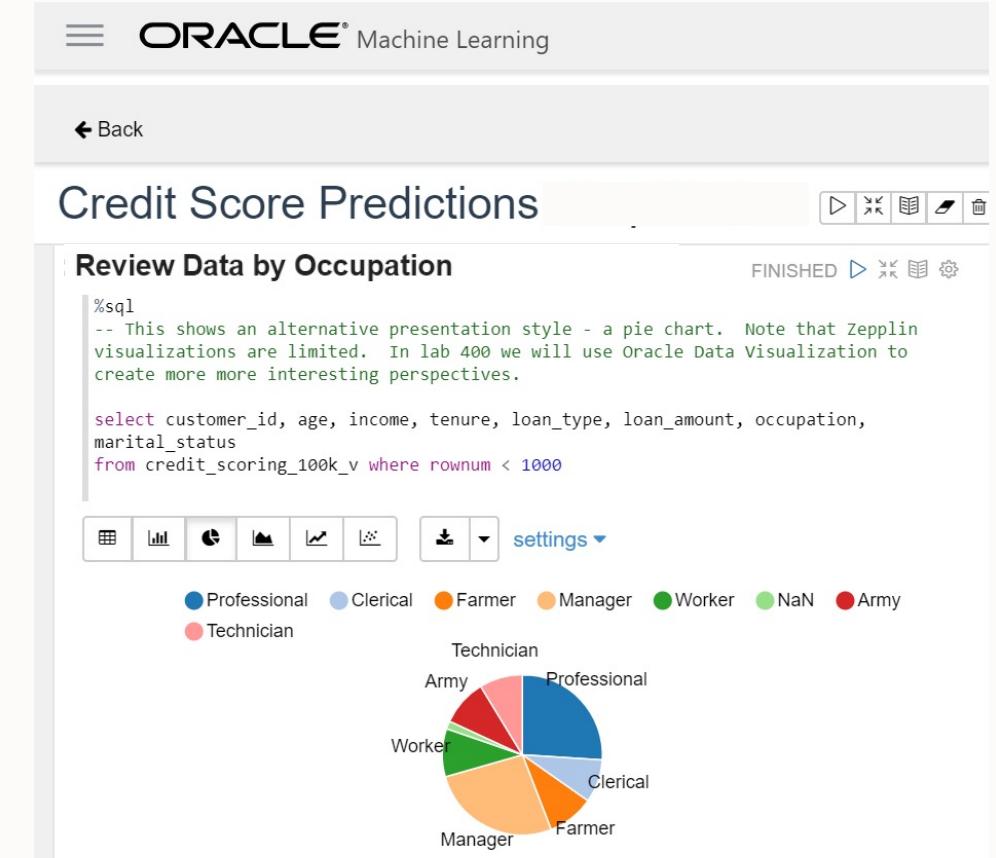
Autonomous Database as a Data Science Platform

- **Collaborative UI**

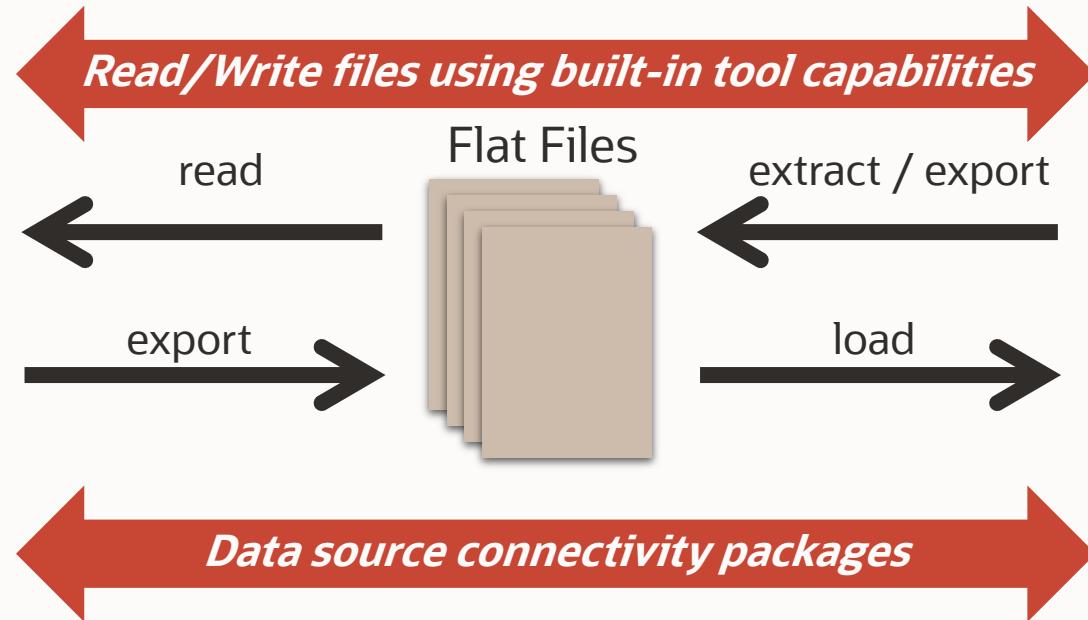
- Based on Apache Zeppelin
- Supports data scientists, data analysts, application developers, DBAs
- Easy sharing of notebooks and templates
- Edits made in one notebook immediately appear in other open shared notebooks
- Permissions, versioning, and execution scheduling

- **Included with Autonomous Database**

- Automatically provisioned, managed, backed up
- In-database SQL algorithms and analytics functions
- Explore and prepare, build and evaluate models, score data, deploy solutions
- Python available!



Traditional Analytics and Data Source Interaction



- **Access latency**
- **Paradigm shift: R/Python → *Data Access Language* → R/Python**
- **Memory limitation – data size, in-memory processing**
- **Single threaded**
- **Issues for backup, recovery, security**
- **Ad hoc production deployment**

Oracle Machine Learning for R and Python

- **Transparency layer**

- Leverage proxy objects so data remain in database
- Overload native functions translating functionality to SQL
- Use familiar R / Python syntax on database data

- **Parallel, distributed algorithms**

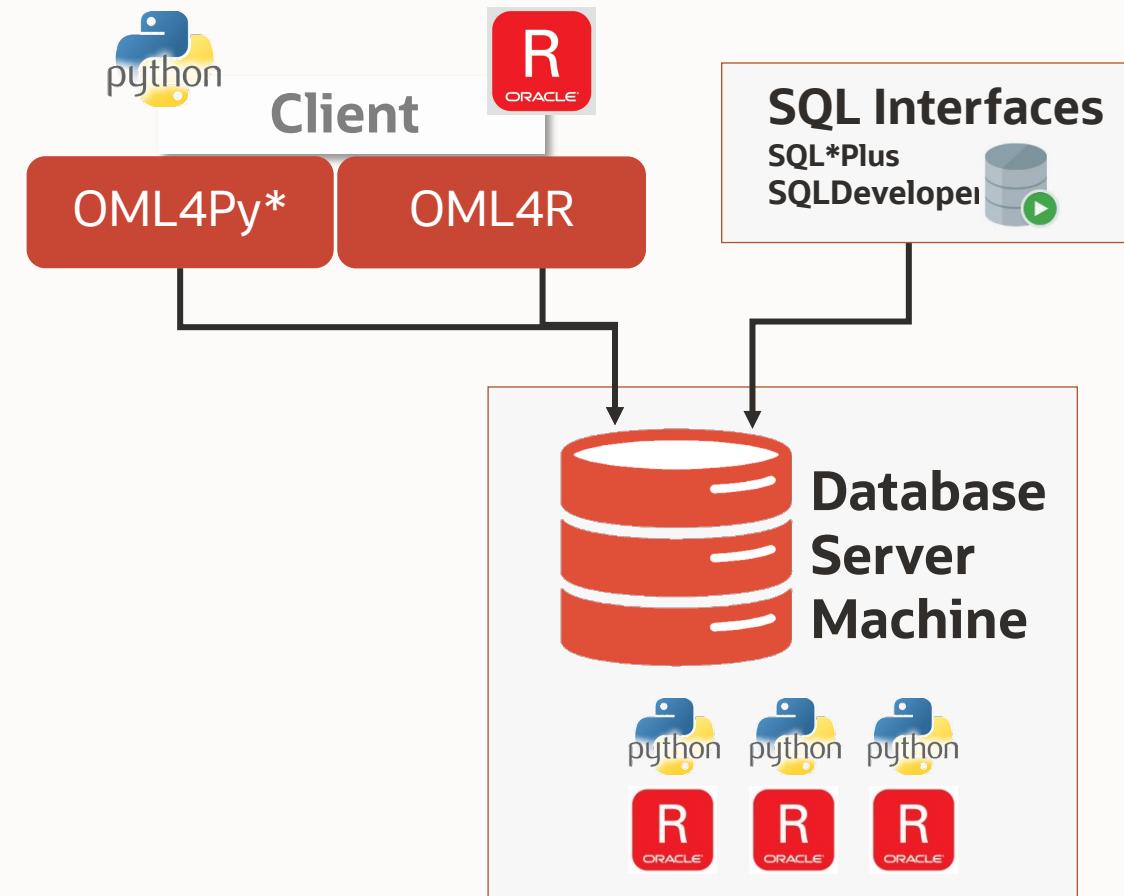
- Scalability and performance
- Exposes in-database algorithms available from OML4SQL

- **Embedded execution**

- Manage and invoke R or Python scripts in Oracle Database
- Data-parallel, task-parallel, and non-parallel execution
- Use open source packages to augment functionality

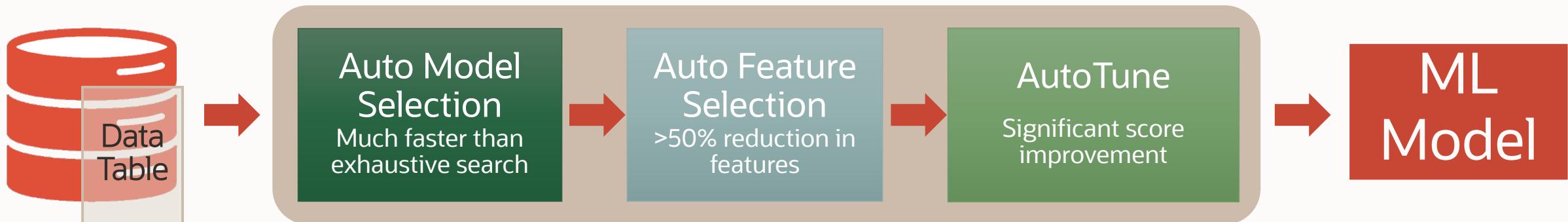
- **OML4Py AutoML**

- Model selection, feature selection, hyper-parameter tuning
- Supports Classification and Regression



AutoML – *new* with OML4Py

Increase data scientist productivity – reduce overall compute time



Auto Model Selection

- Identify in-database algorithm that achieves highest model quality
- Find best model faster than with exhaustive search

• Auto Feature Selection

- Reduce # of features by identifying most predictive
- Improve performance and accuracy

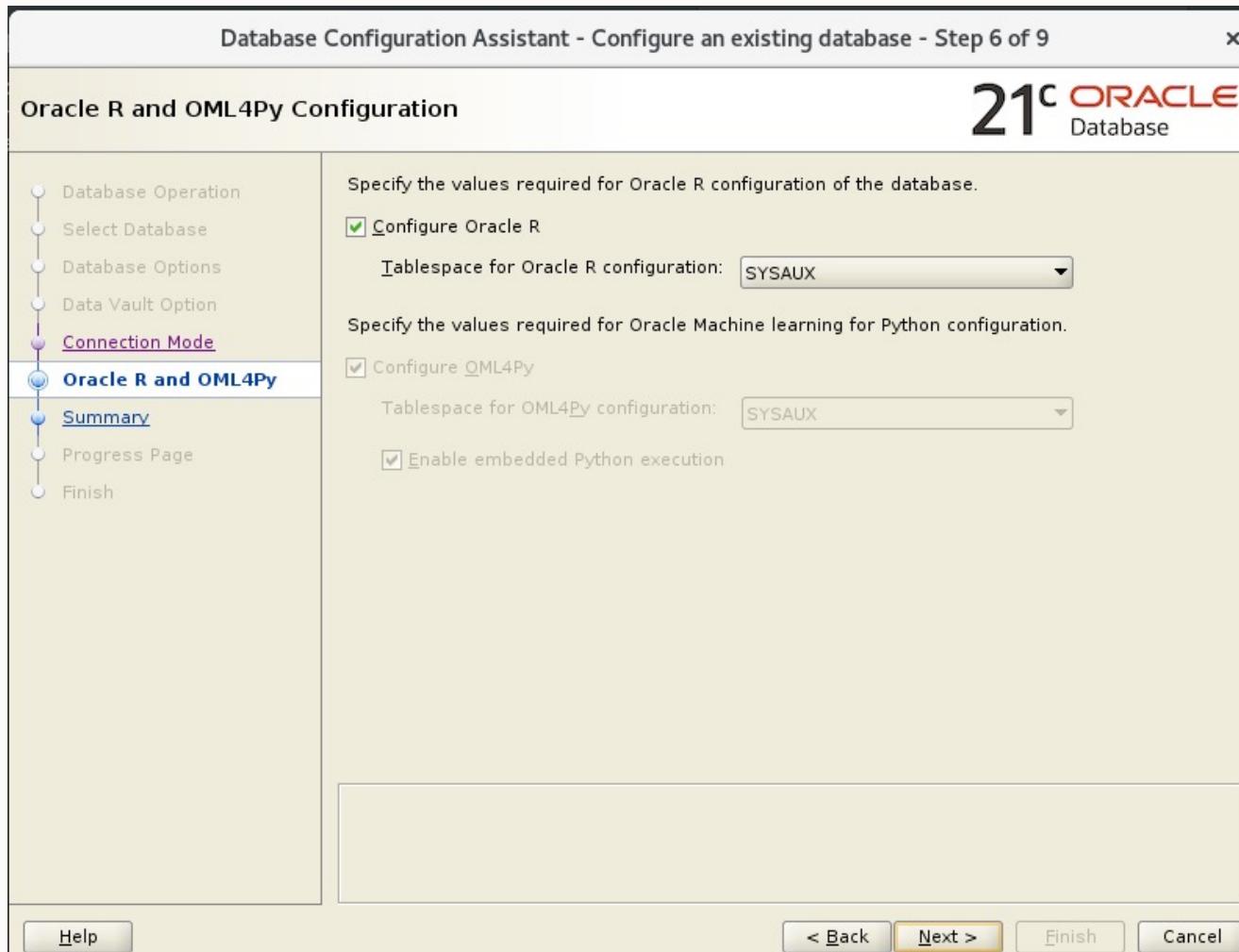
Auto Tune Hyperparameters

- Significantly improve model accuracy
- Avoid manual or exhaustive search techniques

Enables non-expert users to leverage Machine Learning

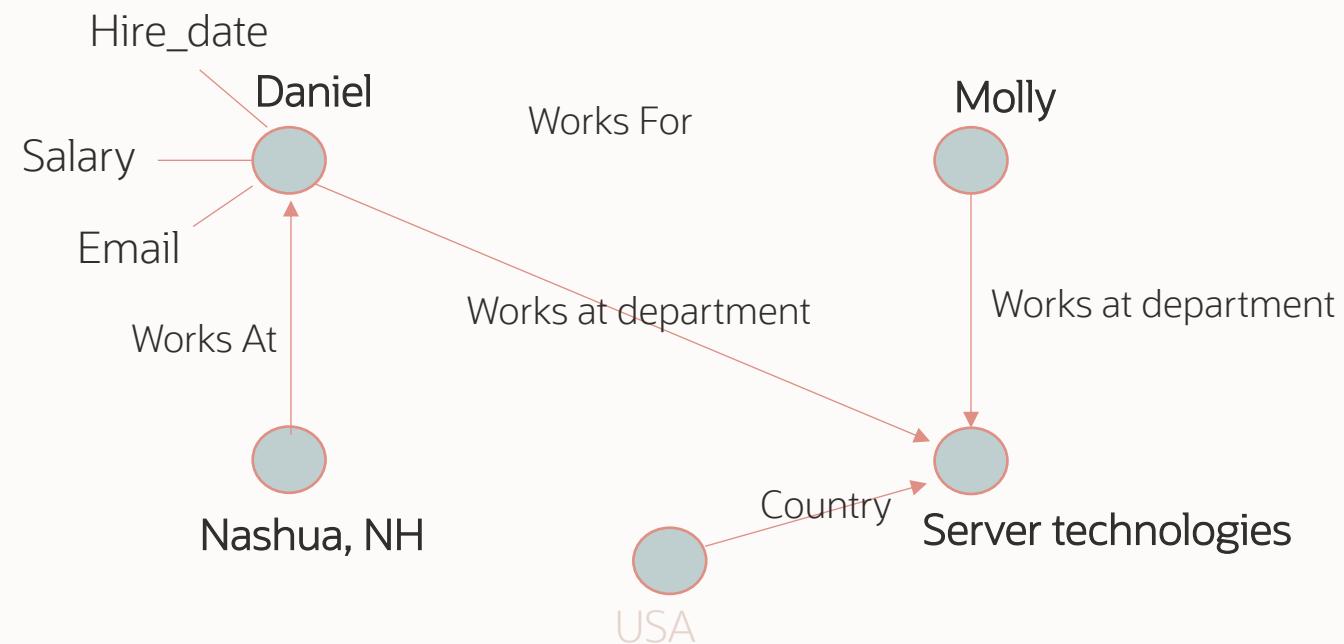
Easy to set up Now in 21c!

Configure R and Python from Database Configuration Assistant



Graph Analytics

- Analytics based on **connections** and **relationships** between data entities



What is Graph Analytics?

A **labeled-property** graph model is represented by a set of nodes, edges, properties, and labels.

What is a graph?

Data model representing entities as vertices and relationships as edges

Optionally including attributes

What are typical graphs?

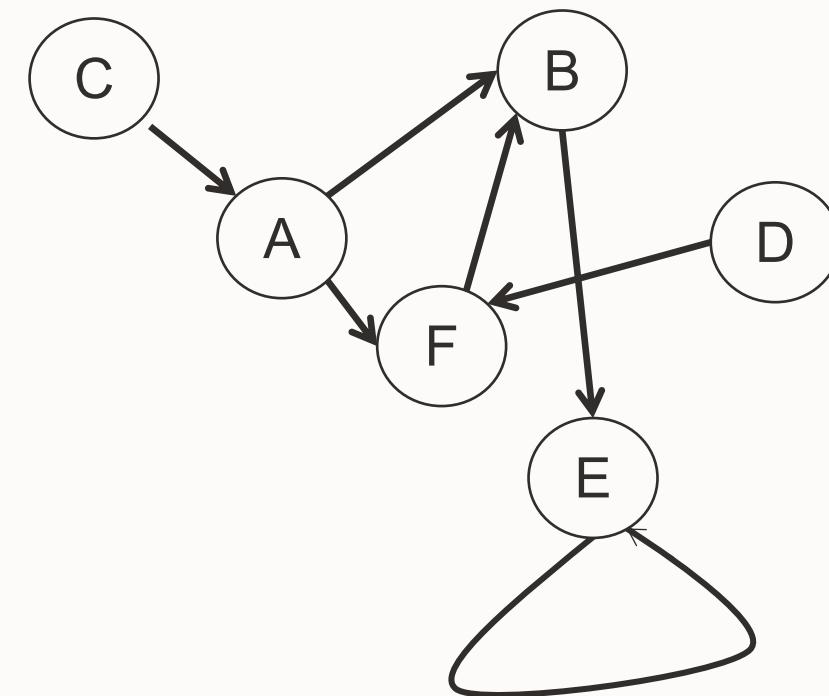
Social Networks

LinkedIn, Facebook, Google+, Twitter, ...

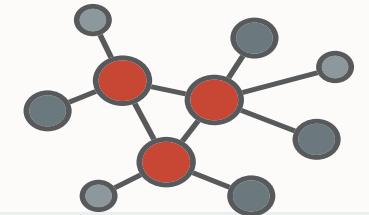
Physical networks, Supplier networks,...

Knowledge Graphs

Apple SIRI, Google Knowledge Graph, ...



From tables to Property Graphs



PRODUCT_ID	BOUGHT_WITH
0	1
0	2
0	4
1	0
1	12
1	23
...	...

PGQL DDL SYNTAX:

```
CREATE PROPERTY GRAPH products
```

VERTEX TABLES (

```
    PRODUCTS KEY(PRODUCT_ID) PROPERTIES (PRODUCT_ID)  
    )
```

EDGE TABLES(

```
    SOURCE KEY(PRODUCT_ID) REFERENCES PRODUCTS  
    DESTINATION KEY(BOUGHT_WITH) REFERENCES PRODUCTS
```

```
)
```

- Every product id is a vertex
- Two vertices in one row are connected by an edge
- (“bought_with” relationship)

Property Graph Product Overview

Store, manage, query and analyze graphs

- **Enterprise capabilities:** Built on Oracle Infrastructure
- Manageability, fine-grained security, high availability, integration and more

High scalable

- In-memory query and analytics and in-database query
- 10s of billions of edges and vertices

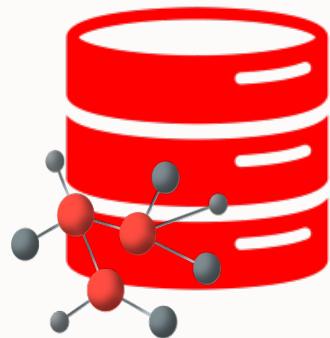
PGQL: Powerful SQL-like graph query language

Analytics Java API: 50+ pre-built graph analytics algorithms

Visualization:

- Light-weight web application: UI accessible from a browser

Oracle Database as a Graph Store

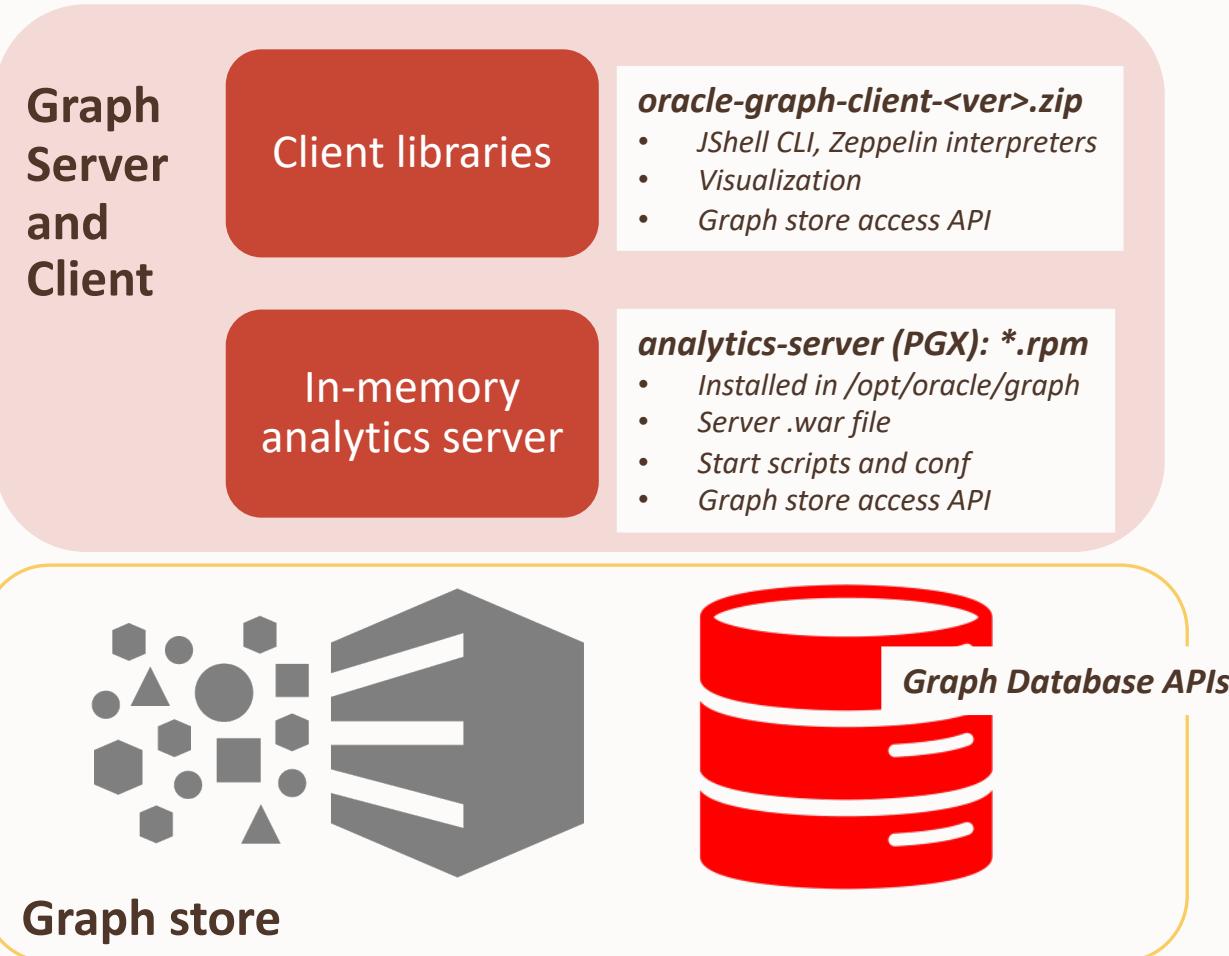


Database stores and manages Graph Nodes, Edges and Properties

Database provides graph traversal and query language and API's

- Java API to develop applications
- Command-line submission of graph queries
- Graph visualization tool
- APIs to update graph store
- PGQL language for Property Graph
- SPARQL language for RDF Triple Store

Now: Graph Server, Client and Storage



- **Graph Server and Client kit**
 - Separate download from e-delivery and oracle.com
(not shipped with \$ORACLE_HOME)
 - 20.1 (first kit) released **Jan 2020**
 - **21.1 latest release**
- Graph Server and Client works with **both Database and Big Data**

PGQL Graph Query Language

Graph pattern matching

(person)-[:works_for]->(person)

Basic patterns and reachability patterns

Can we reach from A to B with an arbitrary number of hops?

Familiarity of SQL users

- Similar language construct and syntax

SELECT ... WHERE ...

GROUP BY ... ORDER BY ...

- 'Result set' (table) as output

PGQL Graph Query

```
1 SELECT n, n0, n1, e0, e1, e2, n.pageRank, n0.pageRank, n1.pageRank  
2 MATCH (n)-[e0]-(n0)-[e1]-(n1), (n)-[e2]-(n1)  
3 WHERE ID(n0) = 'IRON MAN/TONY STARK'  
4 ORDER BY n.pageRank DESC, n0.pageRank DESC, n1.pageRank DESC LIMIT 30
```

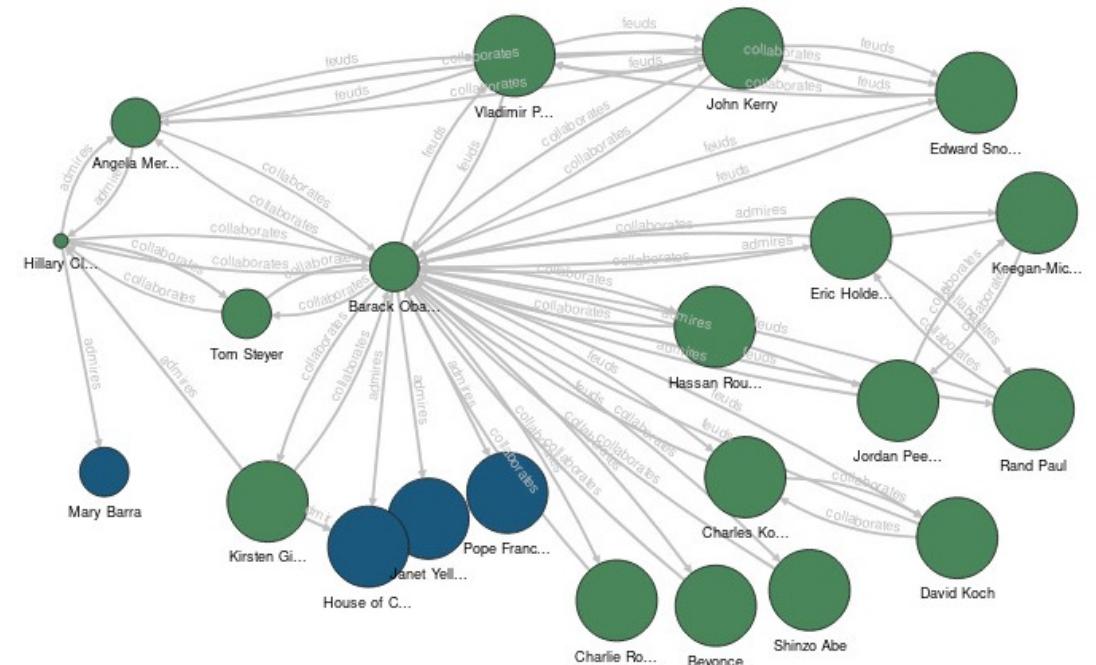
Graph

GraphViz Tool

PGQL Graph Query

```
1 SELECT m,n,e  
2 FROM MATCH (m)-[e]->(n)  
3 where n.distance<=2 and m.distance<=2
```

Graph Parallelism ?
CONNEC... ▾ 0 ▾ ▲ ⏪

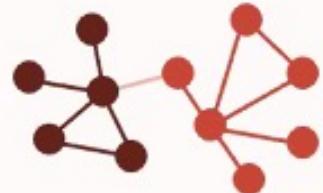


Visualize PGQL query results

- Pre-loaded and Published graphs
- Themes, styles, layouts
- Interactive Graph manipulation

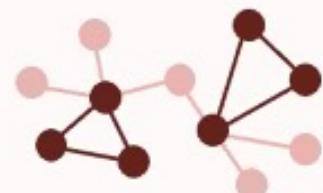
Graph Analytics: 50+ Pre-built Algorithms

Detecting Components and Communities



Strongly Connected Components,
Weakly Connected Components,
Label Propagation,
Conductance Minimization,

Evaluating Structures

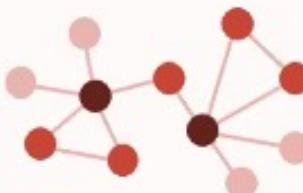


Adamic-Adar Index, Conductance,
Cycle Detection, Degree Distribution,
Eccentricity, K-Core, LCC, Modularity,
Reachability Topological Ordering,
Triangle Counting

Link Prediction

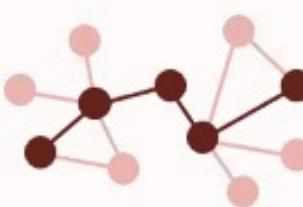
WTF (Who to follow)

Ranking and Walking



PageRank, Personalized PageRank,
Degree Centrality, Closeness Centrality,
Vertex Betweenness Centrality,
Eigenvector Centrality, HITS, SALSA,
Random Walk with Restart

Path-Finding



Shortest Path (Bellman-Ford, Dijkstra,
Bidirectional Dijkstra), Fattest Path,
Compute Distance Index,
Enumerate Simple Paths,
Fast Path Finding, Hop Distance

Others

Minimum Spanning-Tree,
Matrix Factorization

Interaction with the Property Graph

- Access through APIs
 - Implementation of Apache Tinkerpop Blueprints APIs
 - Based on Java, REST plus SolR Cloud/Lucene support for text search
- Scripting
 - Groovy, Python, JavaScript, ...
 - Apache Zeppelin integration, JavaScript (node.js) language binding
- Graphical UIs
 - Cytoscape, plug-in available
 - Commercial tools such as TomSawyer Perspectives
 - Vis.js and D3 among others



Agenda

Converged Database Workshop Series



Converged Database Workshop Series

Marketing Campaign on going

1. Oracle Converged Database: Multitenant, Multimodel, In-Memory

- For DBAs, Solutions Architects and Developers, including CTOs
- One day in Torre Picasso or
- Remote Zoom sessions (3 consecutive days, 2 hours each)



Converged Database Workshop Series

Marketing Campaign on going

1. Oracle Converged Database: Multitenant, Multimodel, In-Memory

- For DBAs, Solutions Architects and Developers, including CTOs
- One day in Torre Picasso or
- Remote Zoom sessions (3 consecutive days, 2 hours each)



2. Oracle Converged Database: Multicloud ECX with Autonomous DB

- For Data Engineers and Cloud Solutions Architects
- One morning in Torre Picasso or
- Remote Zoom sessions (3 consecutive days, 3 hours each)



Converged Database Workshop Series

Marketing Campaign on going

1. Oracle Converged Database: Multitenant, Multimodel, In-Memory

- For DBAs, Solutions Architects and Developers, including CTOs
- One day in Torre Picasso or
- Remote Zoom sessions (3 consecutive days, 2 hours each)



2. Oracle Converged Database: Multicloud ECX with Autonomous DB

- For Data Engineers and Cloud Solutions Architects
- One morning in Torre Picasso or
- Remote Zoom sessions (3 consecutive days, 3 hours each)



3. Oracle Converged Database: Spatial, Graph & ML with Python and R

- For Data Engineers, Data Scientists, Business Analysts and Solutions Architects
- One day in Torre Picasso or
- Remote Zoom sessions (3 consecutive days, 2 hours each)



Agenda

Machine Learning, Spatial and Graph



Oracle Converged Database: Machine Learning, Spatial and Graph Workshop

Hands On Labs

HOL0 – To Be Uploaded (Spatial Web Services) – Not included

HOL1 – Spatial and Spatial Studio

HOL2 – Machine Learning with Python and R

HOL3 - Graph



Materials and HOL manuals

Search or jump to... Pull requests Issues Marketplace Explore

OracleDataManagementSpain / ConvergedDatabase Public Watch 3 Star 4 Fork 4

Code Issues Pull requests 1 Actions Projects 1 Wiki Security Insights

master ConvergedDatabase / MLSpatialGraph / Go to file Add file ...

fralra Update Readme.md 666548c 13 hours ago History

..

Readme.md Update Readme.md 13 hours ago

WORKSHOP_ML_Spatial_Graph_HOL1.pdf Add files via upload 14 hours ago

omlwls.zip Add files via upload 13 hours ago

Readme.md

Machine Learning, Spatial and Graph Workshop

Information

Manuals and all information related to this workshop, that covers:

- Geospatial analysis using *Spatial Studio*
- Machine Learning with R and Python executed in-database
- Property Graph analytics and visualization using PGX

Hands-on Labs Breakout Rooms



Andrés

Room 1



Daniel

Room 2



Manel

Room 3



Francisco

Room 4



Predictive equipment maintenance

Many companies have lots of data about the condition and operation of their industrial equipment.

Analytics can help by providing insight so companies can predict the remaining optimal life of their systems and components, ensuring that their assets operate at optimum production efficiency.

Potential issues can be discovered by analyzing both structured data (equipment year, make, and model) and multi-structured data (log entries, sensor data, error messages, engine temperature, and other factors). With this data, manufacturers can maximize parts and equipment uptime and deploy maintenance more cost effective

Challenges

Machine, log, and sensor data from different types of equipment comes in varying formats.

Integrating all of this data can be difficult. Moreover, the data needs to be analyzed quickly and put into operation to effectively prevent downtime.

Fraud and compliance

When it comes to security, it's not just a few rogue hackers. Every industry is up against entire expert teams. While security landscapes and compliance requirements are constantly evolving. Using data, companies can identify patterns that indicate fraud and aggregate large volumes of information to streamline regulatory reporting.

Challenges

This data requires the integration of different transaction datasets with additional information, such as interaction events and customer behavior. To identify potential fraud patterns, companies will need to sift through a large volume of data.

Production optimization

Optimizing production lines can decrease costs and increase revenue. Data can help manufacturers understand the flow of items through their production lines and see which areas can benefit. Data analysis will reveal which steps lead to increased production time and which areas are causing delays.

Challenges

Optimizing production requires manufacturers to analyze their production equipment data, material use, and other factors. Combining the different kinds of data can pose a challenge.

Product development

Data can help you anticipate customer demand. By classifying key attributes of past and current products and then modeling the relationship between those attributes and the commercial success of the offerings, you can build predictive models for new products and services. Dig deeper by using the data and analytics from focus groups, social media, test markets, and early store rollouts to plan, produce, and launch new products.

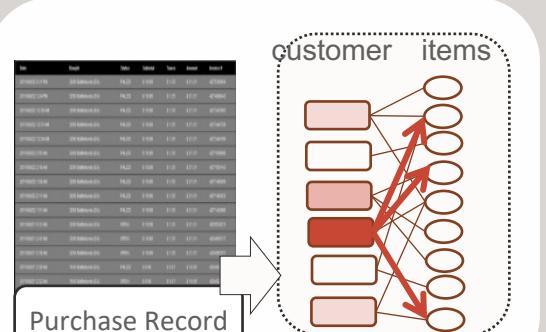
Challenges

Companies will have to analyze what can be a high volume of data coming in varying formats, and then create segments according to customer behavior. They will also have to identify sophisticated use patterns and behavior and map them to potential new offerings.

Common Graph Analysis Use Cases

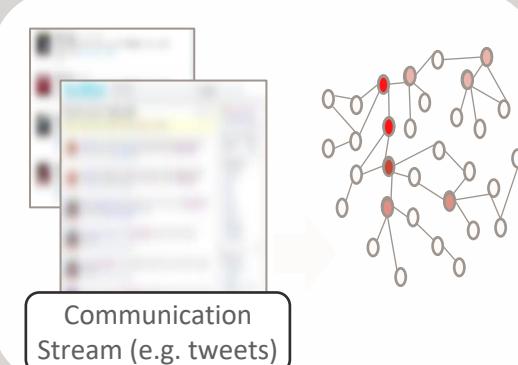
Recommend the most *similar* item purchased by *similar* people

Product Recommendation



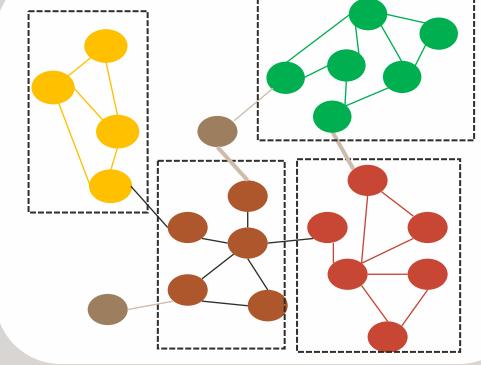
Find out people that are *central* in the given network – e.g. influencer marketing

Influencer Identification



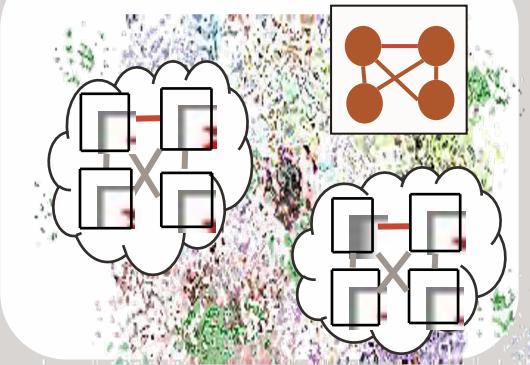
Identify group of people that are close to each other – e.g. target group marketing

Community Detection



Find out all the sets of entities that match to the given pattern – e.g. fraud detection

Graph Pattern Matching



Oracle Cloud Free Tier

Build, test and deploy apps
on Oracle Cloud - for free!

Start Now

Always Free

Services you can use for unlimited time

+

30 Day Free Trial

Up to 400€

Follow this link →

<https://bit.ly/2wG4gPK>

O

ORACLE

ENCUESTA

Workshop Virtual BD Convergente: ML, Spatial y Graph



¡Tu opinión es muy importante!

Escanea el QR para responder o entra aquí:

<https://bit.ly/3Fjo5dS>

Inspiration & Innovation



Our mission is to help people
see data in new ways, discover
insights, unlock endless possibilities.



Thank you!

Oracle Spain

