



# Develop AI/ML-Driven Apps with Autonomous Database 23ai



---

**Paul Parkinson**

Architect and Developer Advocate

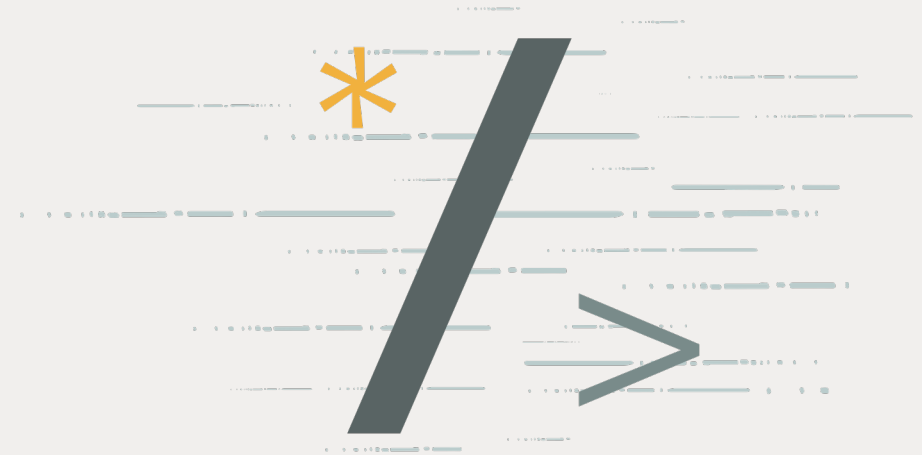
Oracle Database

August 6th 2024

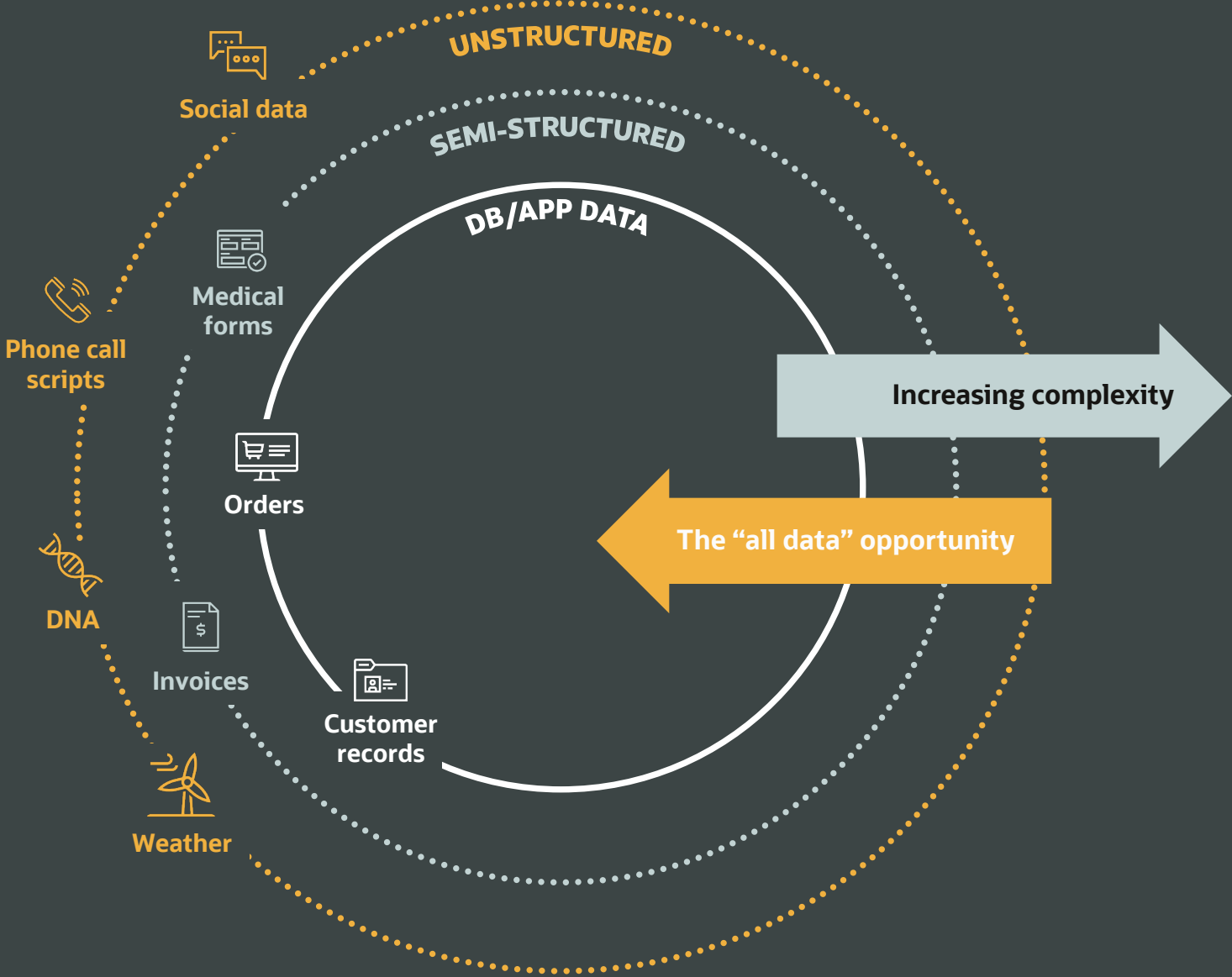
# Safe harbor statement

---

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.



# It all starts with data





# Oracle brings AI to your data

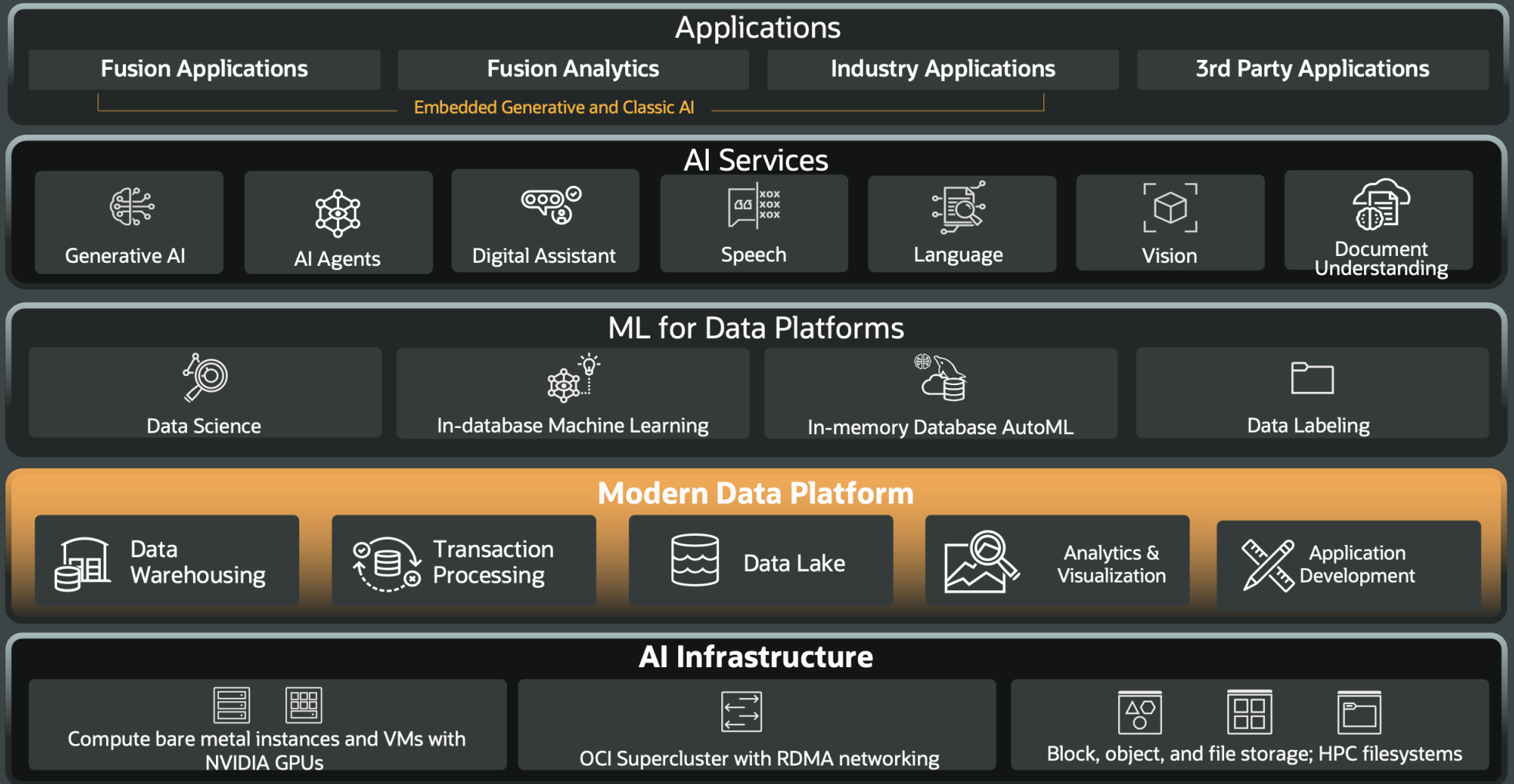


# Agenda

- OCI AI Services Oracle AI Stack
- Oracle Multi-purpose/Converged Database
- Select AI (NL2SQL)
- Vector and RAG
- AI Agents
- AI Playground
- OML4PY
- Spatial AI
- Graph ML and RAG

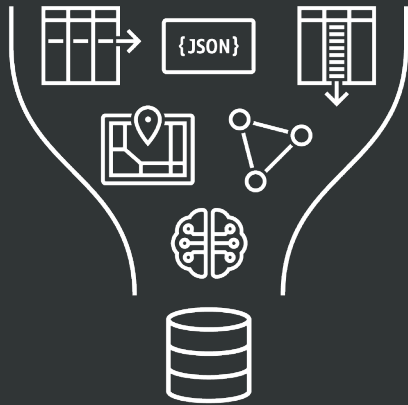


# The Oracle AI stack



# How we deliver the Vision

Complete and Simple Platform for All Data Management Needs



## Converged Database

**Complete:** all modern data types, workloads, and development styles

**Simple:** Add a SQL Statement, not a database to support any need of modern applications



Running on

## Autonomous Database

**Powerful:** All the benefits of converged database running on Exadata

**Simple:** Fully-managed cloud service

# Language and API support

## OCI...

- SDKs for JavaScript, Java, Python, .NET, PL/SQL, as well as REST.

## Oracle Database...

- Support for all languages as well as REST.
- Run Java, JavaScript, and Python within the database itself.
- DBMS\_CLOUD and DBMS\_CLOUD\_AI PL/SQL packages to call other OCI AI, etc. services
- Outbound REST calls to external services including multi-cloud, huggingface, etc.

Most AI apps will involve adding AI to existing apps vs writing from scratch.



What if you could...

Converse in natural language with your database

Combine semantic search on unstructured data with searches on business data

Get better answers from Generative AI by augmenting it with data in your database



# Autonomous Database Select AI

**Simplest** way to get answers about your business



Select AI allows you to use your natural language to query data

- No need to understand where and how your data is stored to gain insights

LLMs infer the intent of our natural language question  
We infer a lot from human language

what are our total streams for each tom hanks movie this month?

total number of  
movie views

breakout views  
by movie

tom hanks is  
an actor

understanding  
of time

# Putting it all together

what are our total streams for each tom hanks movie this month?

what are our total streams for each tom hanks movie this month?

MOVIE_TITLE	TOTAL_STREAMS
Forrest Gump	888.00
Sleepless in Seattle	176.00
Splash	117.00
Angels & Demons	293.00
Saving Private Ryan	494.00
Who Killed the Electric Car?	9.00
Philadelphia	189.00
Cast Away	567.00
Big	189.00
The Great Buck Howard	4.00



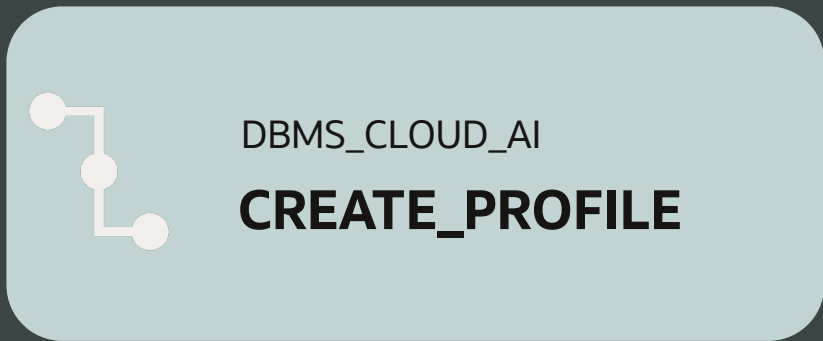
```
SELECT
  m.title AS movie_title,
  COUNT(s.views) AS total_streams
FROM movie m
  JOIN sales_sample s ON m.movie_id = s.movie_id
  JOIN actors a ON m.movie_id = a.movie_id
WHERE a.actor = 'Tom Hanks'
  AND EXTRACT(MONTH FROM s.day_id) =
    EXTRACT(MONTH FROM SYSDATE)
GROUP BY m.title
```

LLMs are remarkable at inferring intent (and getting better)  
But they are not perfect  
It is **very important** to verify the results



# How Select AI Works?

Create a profile describing your LLM & the metadata such as schemas, tables, etc.



```

BEGIN DBMS_CLOUD_AI.CREATE_PROFILE(
  profile_name => 'openai_gpt35',
  attributes => '{
    "provider"      : "openai",
    "credential_name": "OPENAI_CRED",
    "object_list"   : [
      {"owner": "MOVIESTREAM", "name": "movies"},
      {"owner": "MOVIESTREAM", "name": "streams"},
      {"owner": "MOVIESTREAM", "name": "cust_ext"},
      {"owner": "MOVIESTREAM", "name": "pizza_shop"},
      {"owner": "MOVIESTREAM", "name": "actors"},
      {"owner": "MOVIESTREAM", "name": "genre"},
      {"owner": "MOVIESTREAM", "name": "cust_seg"},
      {"owner": "MOVIESTREAM", "name": "cust_contact"}
    ] }' );
END;
/

```



# Configure Select AI

## Perform AI translation with your existing profile



DBMS\_CLOUD\_AI

**Generate**



```
l_response := DBMS_CLOUD_AI.GENERATE(  
  prompt      => '<task + data>',  
  profile_name => 'openai_gpt35',  
  action      => 'chat' / 'showsql' / 'narrate'  
);
```

```
SELECT DBMS_CLOUD_AI.GENERATE(  
  prompt      => 'How many customers do we have',  
  profile_name => 'openai_gpt35',  
  action      => 'showsql'  
) FROM dual;
```

# Adding context to users question

The Prompt Context is compiled and stored into a Page Item - P6\_PROMPT\_CONTEXT

```
:P6_FINAL_PROMPT := apex_string.join_clobs(  
  apex_t_clob(  
    '<CONTEXT>', :P6_PROMPT_CONTEXT||chr(10),  
    '<QUESTION> Based on the context above answer the following question: '||:P6_PROMPT||CHR(10),  
    'If this question cannot be answered based on above context say - "Information not found!"') );
```

--Example query for generating context

```
SELECT 'Overview of the school : '|| OVERVIEW_PARAGRAPH ||chr(10) || chr(13)||  
  'The following Language Courses are taught here : '||LANGUAGE_CLASSES||chr(10) || chr(13)||  
  'The following Advanced Placement Courses are taught : '||ADVANCEDP_COURSES||chr(10) || chr(13)||  
  'The following is the Diversity in Admission Policy for school: '||diadetails||chr(10) || chr(13)||  
  'Extra curricular activities are available: '|| extracurricular_activities|| chr(10) ||chr(13)||  
  'The Public Schools Athletic League (PSAL) sports for boys: '||PSAL_SPTS_BOYS || chr(10) ||chr(13)||  
  'The Public Schools Athletic League (PSAL) sports for girls: '||PSAL_STS_GIRLS || chr(10) ||chr(13)||  
  'Other facilities in this school: '||adttl_info1 || chr(10) ||chr(13)||  
  'The following academic oppurtunities are available : '||academic_opps || chr(10)||chr(13)  
  as prompt_context  
FROM high_schools WHERE id = :P6_ID;
```



## Finding what you're looking for can be challenging...

I have a photo of an artwork, but don't know the artist or work's name

I got this error message ORA-XXXX, but don't know how to resolve it

I heard a catchy tune and can even hum a few bars, but don't know the title or artist

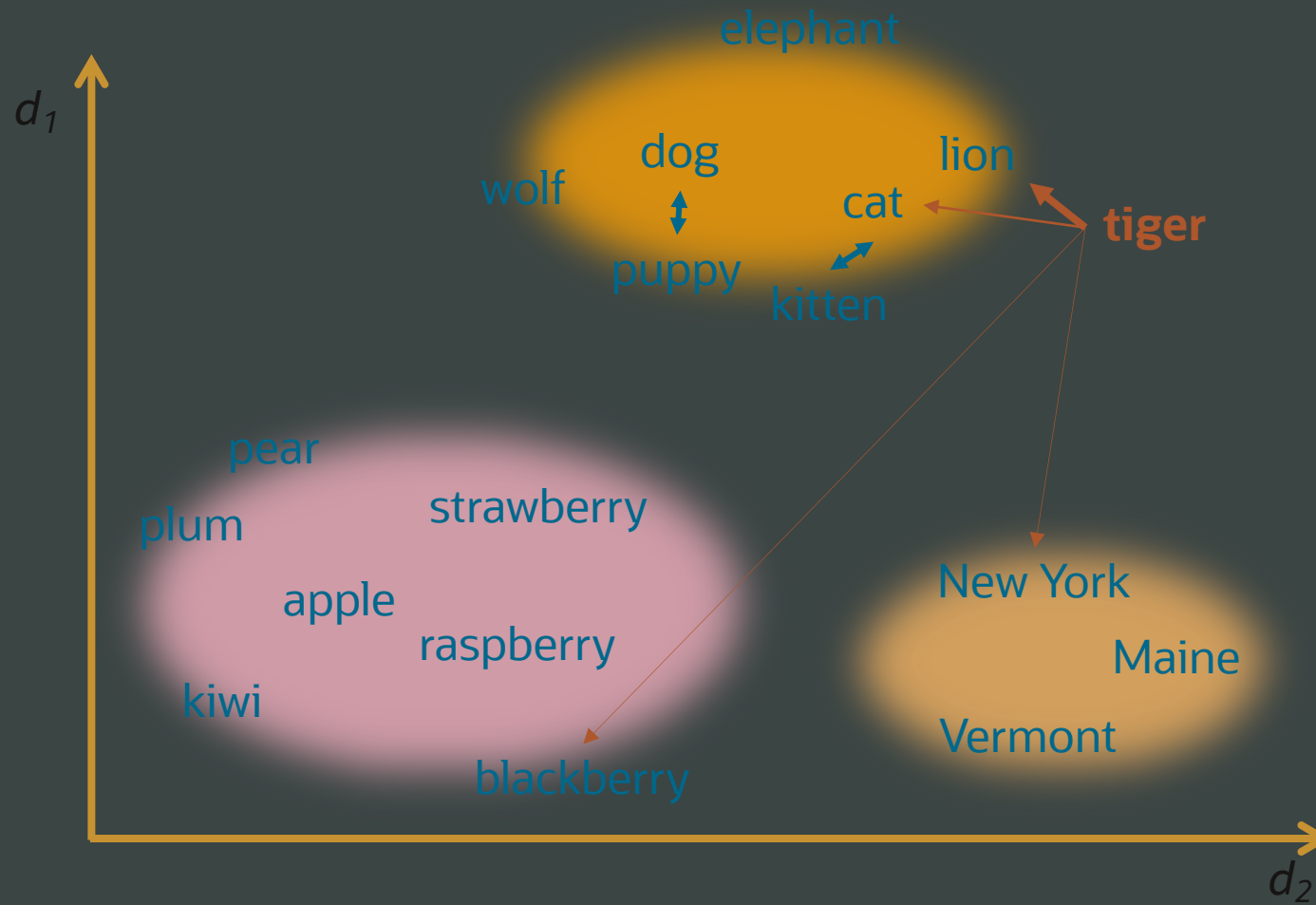


A technology called AI Vector Search enables end-user questions to be mapped to relevant data in your database enabling semantic searches on unstructured data.



# An intuitive example

Word relatedness in two dimensions



# Embeddings for images

Convert an image to one or more vector embeddings and search vector database



0.5, 1.5, 2.6, -1.1, ...

Compare these against other vectors in the database to find similar content





# Introducing Oracle AI Vector Search

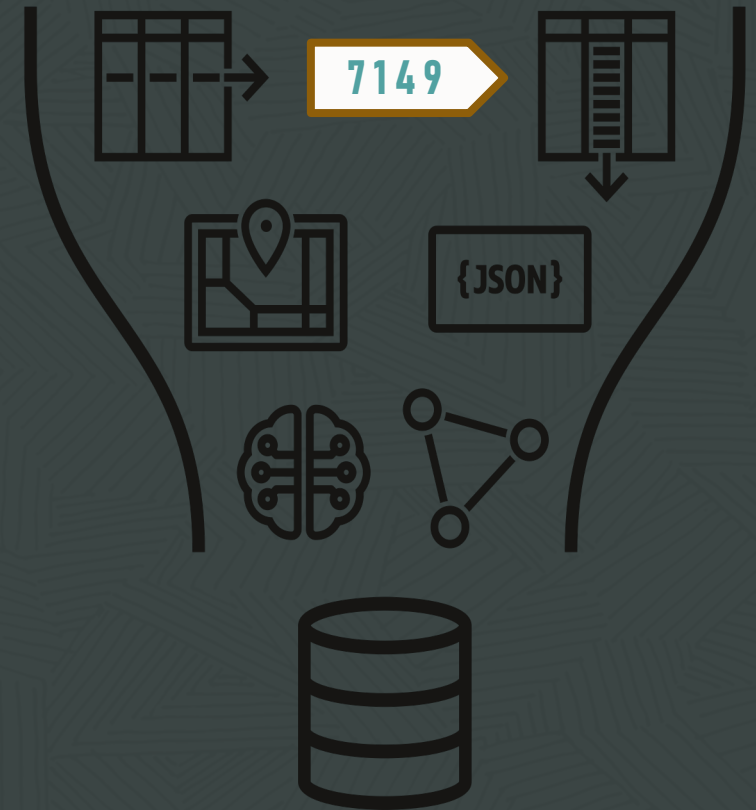
New set of capabilities coming in Oracle Database 23ai

Process vector search and other workloads in same Oracle converged database

Designed to be simple to use and easy to understand

- **New** VECTOR data type for storing vectors
- **New** SQL syntax and functions express similarity search with ease
- **New** Approximate search indexes packaged and tuned for high performance and quality

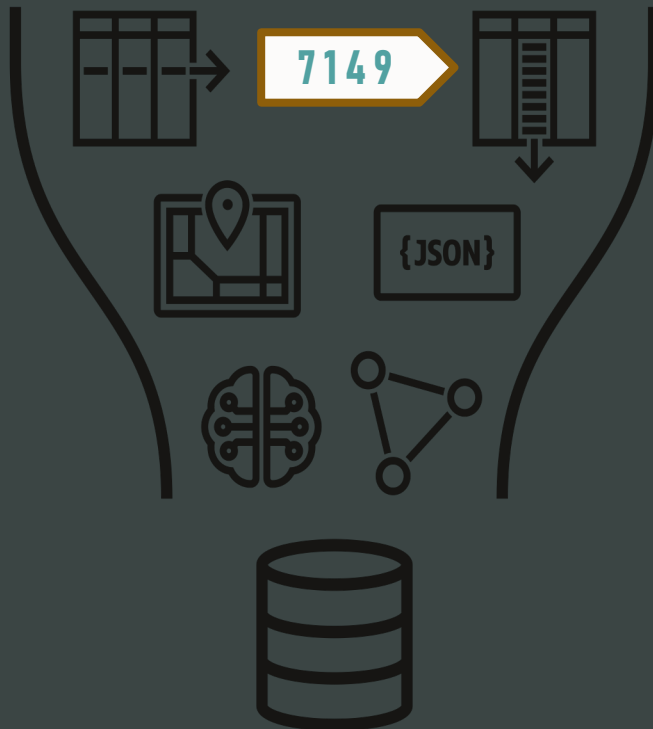
Perform vector search in queries alongside business data about customers and products



Converged Database



# Oracle Database 23ai can store vectors using a new vector data type



```
CREATE TABLE house_for_sale (house_id    number,  
                             price       number,  
                             city        varchar2(400),  
                             house_vector vector  
);
```

# Allows finding data that is semantically similar to an input



Models are pluggable

No ML expertise required

DBAs and Developers can  
learn and use AI Vector  
Search in minutes

Find houses that are similar to this picture



```
SELECT ...  
FROM house_for_sale  
ORDER BY vector_distance(house_vector, :input_vector);
```

# Allows queries that combine AI vector search with business data about customers and products

Combines customer data, product data, and AI search in a few lines of SQL!

A single integrated solution, all data is fully consistent

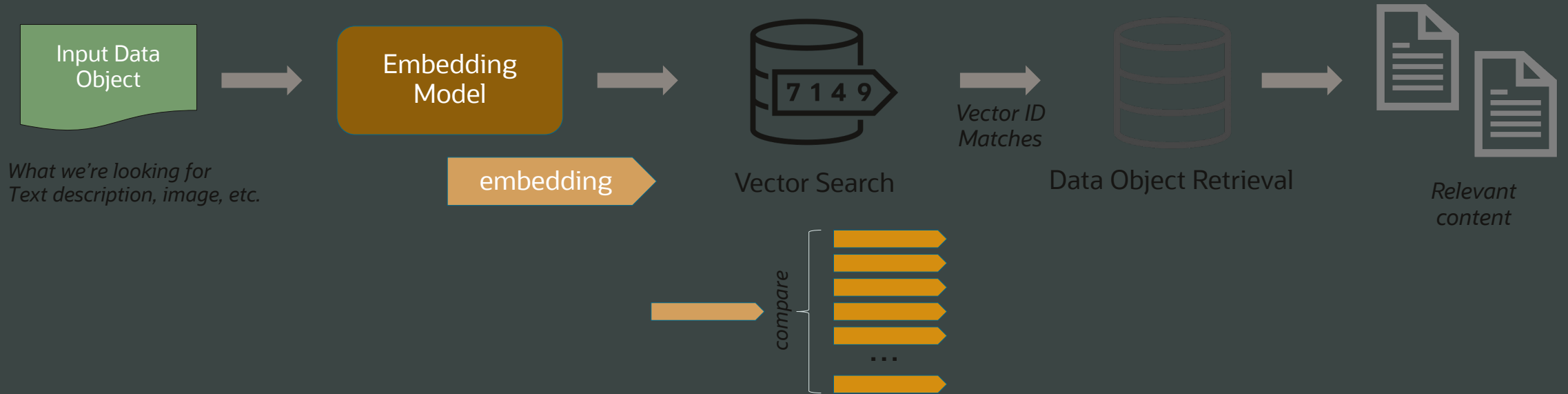
Find houses that are similar to this picture and match the customer's preferred city and budget



```
SELECT ...  
FROM   house_for_sale  
WHERE  price <= (SELECT budget      FROM customer ...)  
AND    city  in (SELECT search_city FROM customer ...)  
ORDER BY vector_distance(house_vector, :input_vector)  
FETCH APPROX FIRST 10 ROWS ONLY  
WITH TARGET ACCURACY 95 PERCENT;
```

# Simple vector search

## Pipeline



## Retrieval Augmented Generation (RAG)

LLMs may have information that is outdated, only public, can not be referenced to a source (provenance), provide hallucinations, ...

AI Vector Search can augment Generative AI by retrieving additional, often private, content needed to answer questions more accurately

RAG addresses this by reusing inference model with additional data but without the need to retrain.

Oracle Database Vector AI allow for the rapid coding of new data, and searches against that data to feed into the LLM.

# Role of vector databases with LLMs

Avoid using sensitive customer data for LLM training and fine-tuning

Cheaper than fine-tuning LLMs, which can be expensive to update

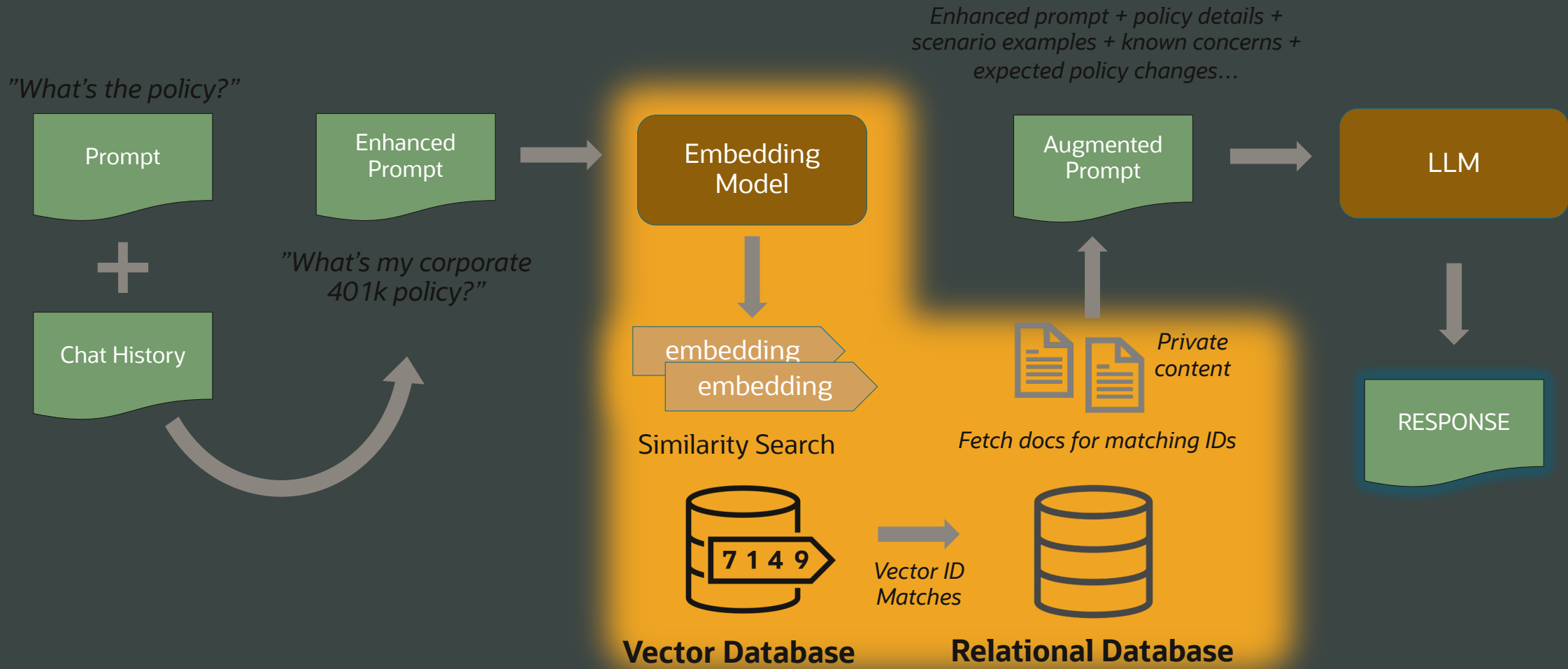
Real-time updated knowledgebase

Cache previous LLM prompts/responses to improve performance and reduce costs



# LLM-based chatbot with “enterprise knowledge”

## Extended pipeline using retrieval architecture

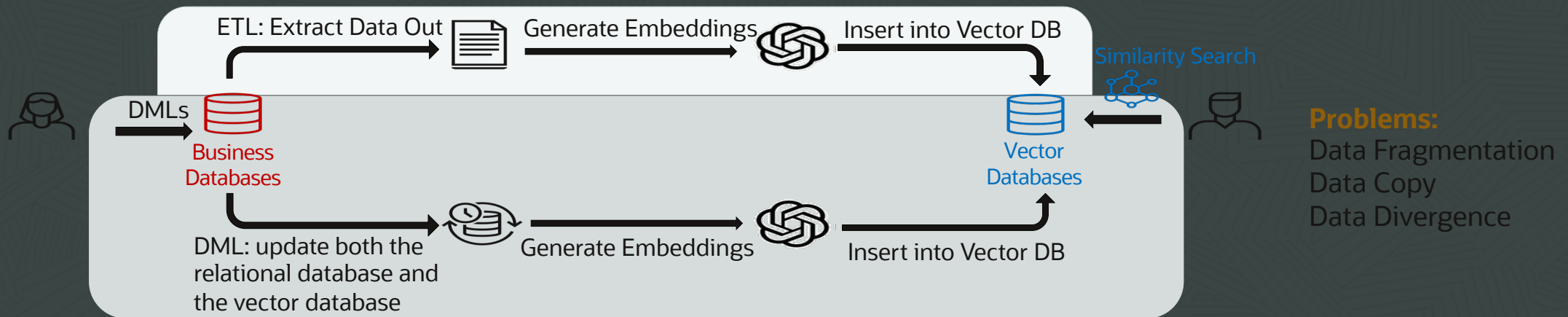




# Oracle's Intrinsic Advantage with Vector Search

Oracle already holds **most of the world's operational enterprise data**

It's an **order of magnitude easier** to **add vector embeddings to enterprise data** than replicating enterprise data to a separate vector database to perform similarity search



It's easy to add vector search capability into a converged, enterprise-grade business database

- *But is it easy to add enterprise-grade relational database functionality into a vector database?*





# AI Vector Search | Complex, Converged SQL

Oracle is a converged database that supports all types of workloads and data models:

- Graph, Text, JSON, Spatial, Relational, etc.

But it also support the full gamut of SQL, including complex operators and functionality:

- Window analytic functions, stored procedures, aggregation

What you get is support for **Complex, Converged SQL**:

*Show me the top 3 photos, grouped by year, over the past 5 years, based on similarity to a provided query image. The photos should have been taken within 20 miles of San Francisco, and have been viewed by at least 100 different people.*

No purpose-built Database (let alone Vector Db) can do this.

Top-3

(top 3 photos per matching group)

Vector Search

(images similar to query image)

Having Clause

Having sum("views" > 100

Group by Sum

(group by year, sum "views")

Spatial

(20 miles from SF)

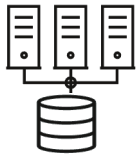
Relational

(last 5 years)



# Oracle AI Vector Search is Fully Integrated

Seamless Integration with core database features for enterprise-grade performance and reliability



RAC



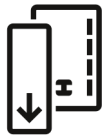
Sharding



Partitioning



Exadata



Transactions



Parallel Execution

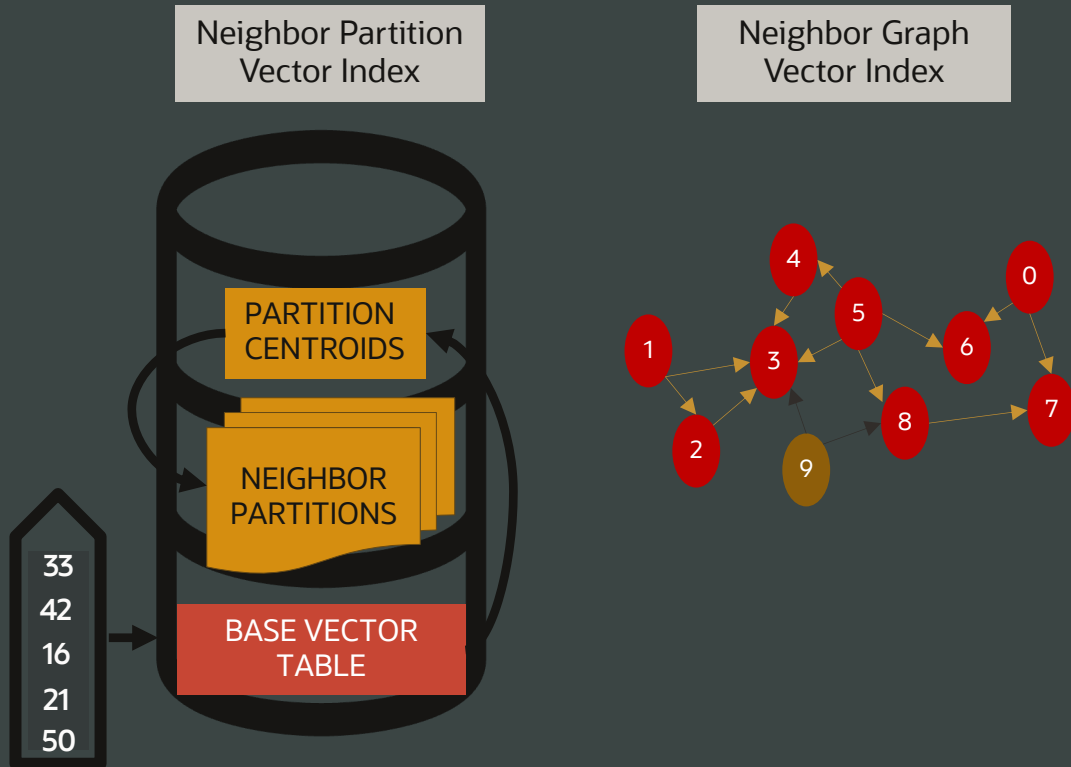


Analytics



Security

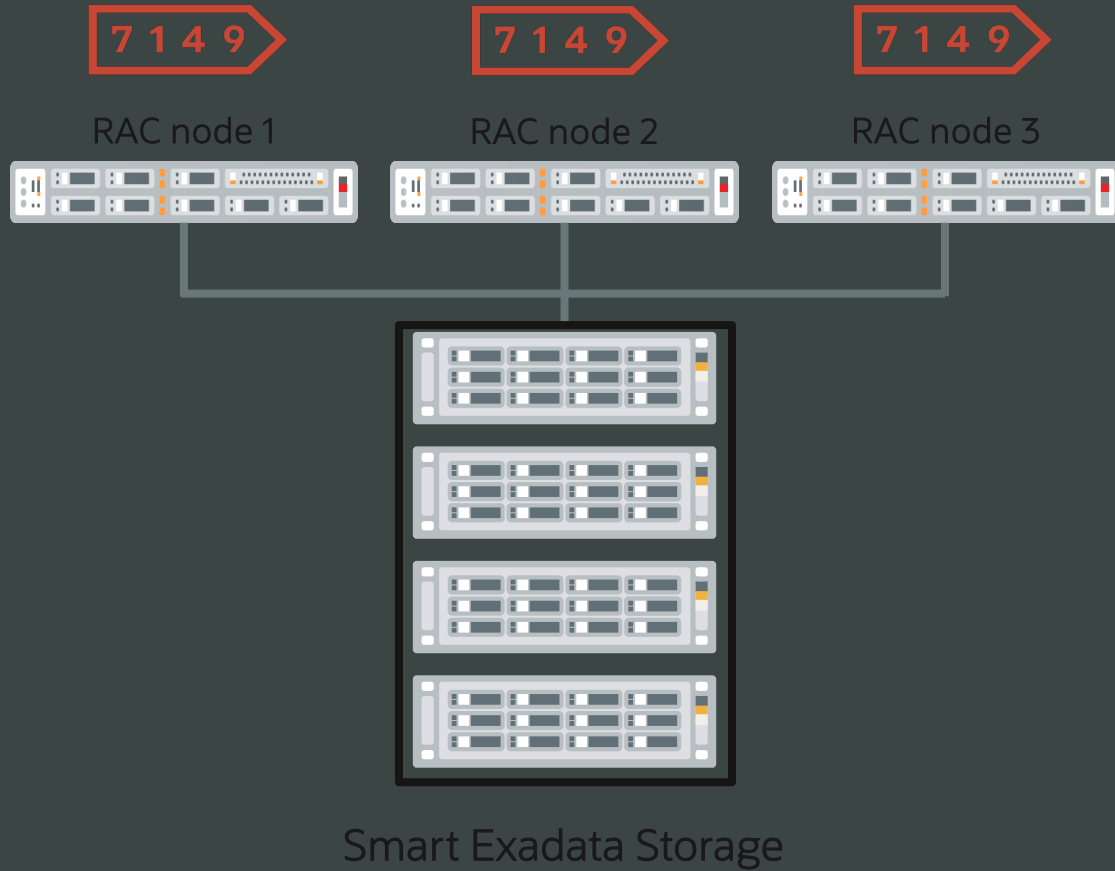
# AI Vector Search | Transactions



Oracle's AI Vector Search Indexes **maintain transactional consistency** with DML activity



# AI Vector Search | Scale-Out with RAC

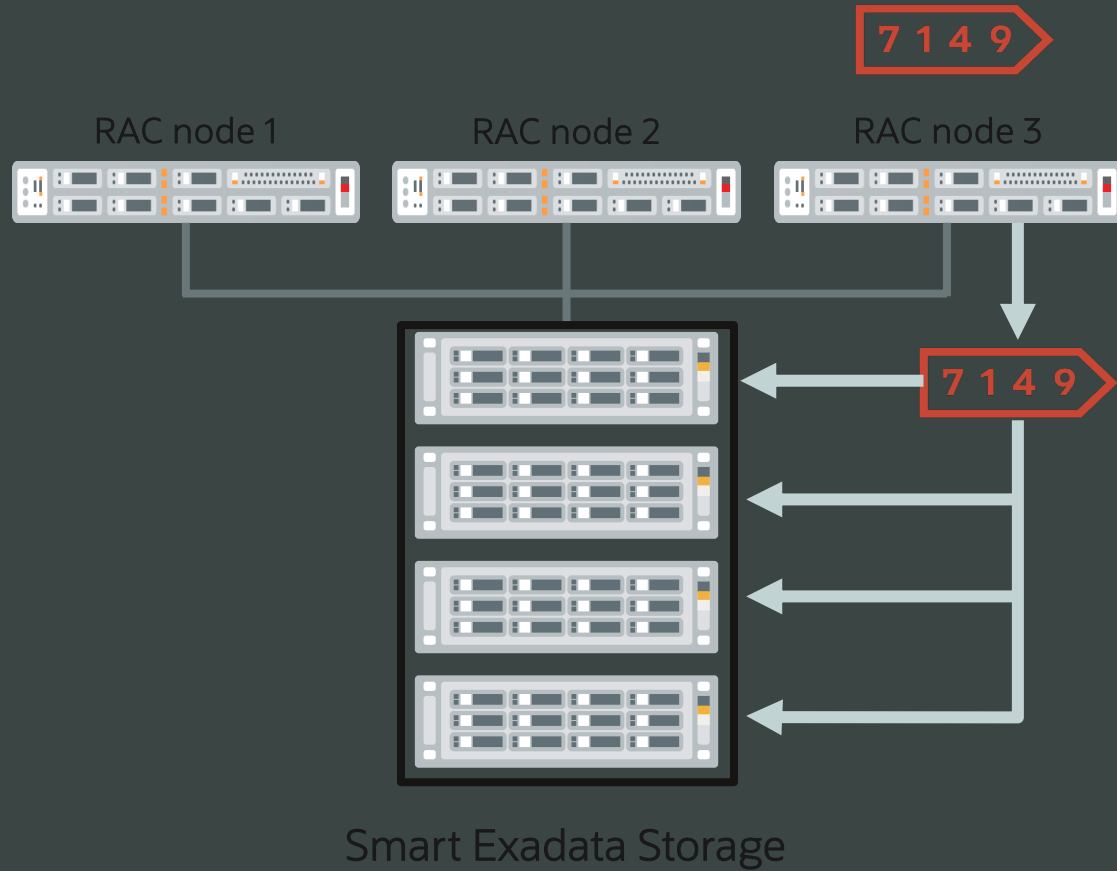


Oracle **transparently scales vector** processing across the computers in a RAC cluster

With full data consistency



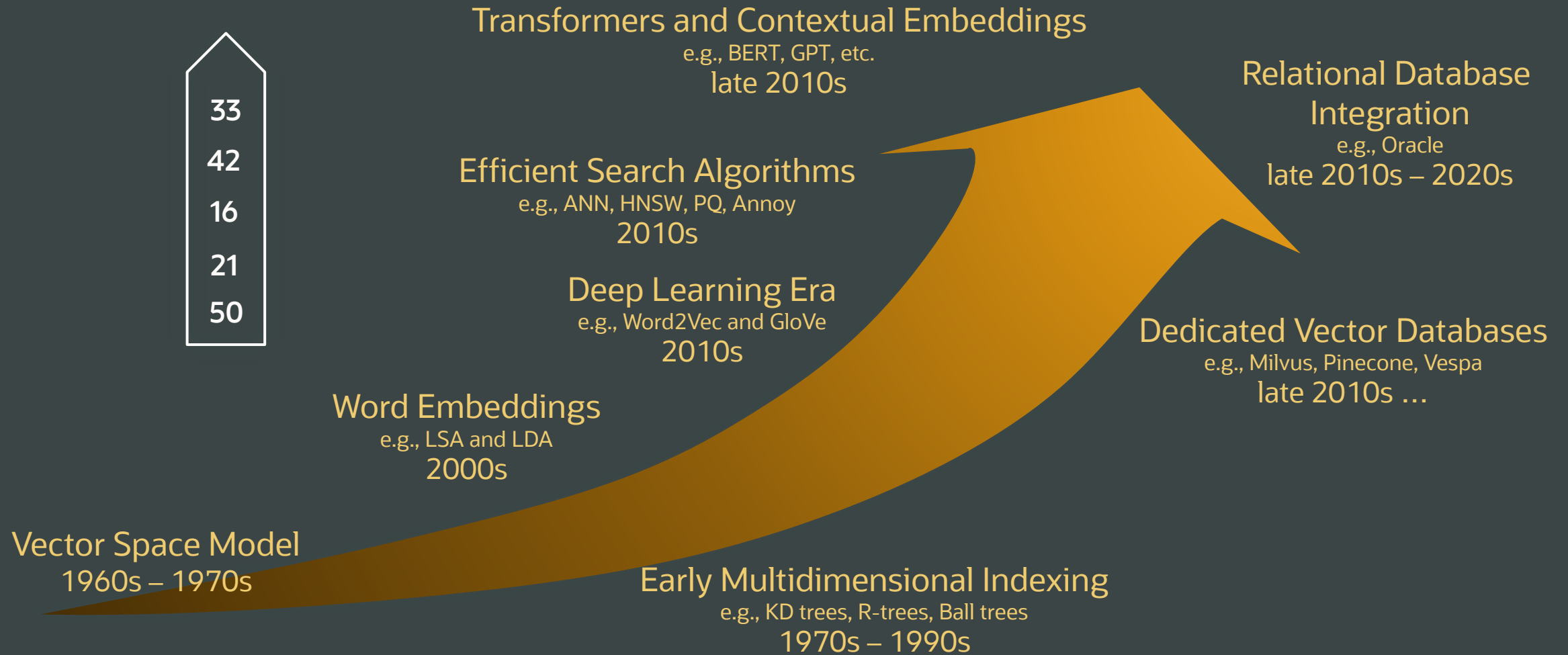
# AI Vector Search | Scale-Out with Exadata



Oracle AI Vector search can be **transparently offloaded** to smart Exadata storage for faster search



# Vector search continues to evolve...



# OCI Generative AI Agents

## Why use OCI Generative AI Agents?

### Converse directly with your enterprise knowledge bases

Leverage a conversational interface that lets anyone query enterprise data stores in natural language.

### Get up-to-date results

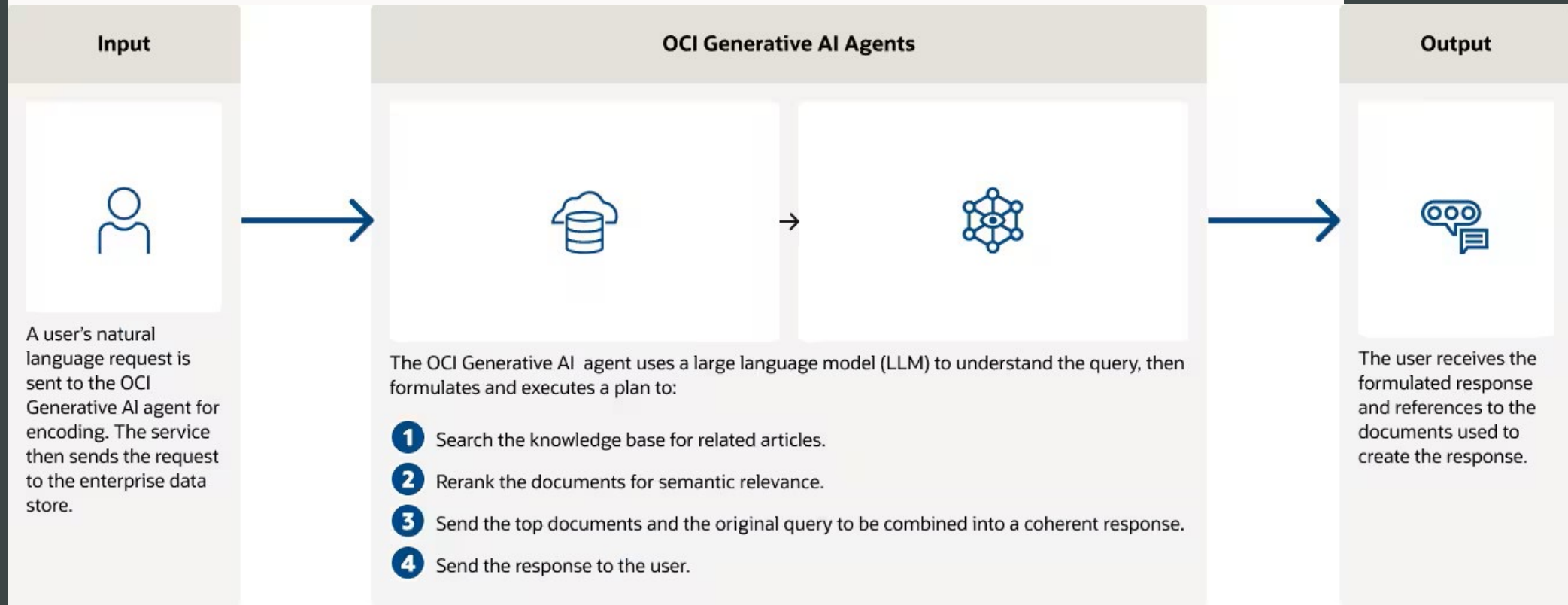
Receive dynamic, real-time responses, even on fast-changing enterprise data stores.

### Embedded for actionable results

Create and embed custom AI agents into enterprise business applications and processes.

### Scale your search

Get faster results through RAG with OpenSearch, with AI Vector Search in Oracle Database 23ai.



# The fallacy of the “quick” chatbot



Many articles present a simple chatbot that can be created very quickly but these are not production-ready – why not?



# It's not as easy as you may think...



**Many AI projects fail miserably.** So today, let's have a look at the recent news on [McDonald's](#) AI fiasco at its drive thru.

An AI-powered voice ordering technology was developed in partnership with [IBM](#) and launched at 100 McDonald's drive-thru locations across the United States.

But a series of viral TikTok videos started appearing recently where the AI failed to comprehend different accents, got confused by background noise and misunderstood and messed up orders.

**Case in point:** it gave one person hundreds of dollars worth of chicken nuggets and gave another customer an ice cream cone topped with bacon. 🤖

Suffice it to say no brand wants this kind of ridicule as publicity.

McDonald's just announced that it will shut off the AI system by the end of July.

# What are some common challenges?

- Choosing the right model
- Context window size
- The **quality** of the corpus – structure, content, formatting, size, presence of distractors
- The preciseness or **ambiguousness** of the question
- The quality of the **prompt**
- Domain-specific specialized language not seen in model pre-training
- Dealing with information after the model pre-training was completed

There are some aspects of working with LLMs that are extremely expensive (in terms of time and/or money)

- Testing/evaluating model performance with **human feedback** or human-created datasets
- **Fine tuning** models
- Creating embeddings for **large corpora**

LLMs tend to “make up” answers if they don’t know the answer – formally called “**hallucination**”



Success factors

Creating a production-ready chatbot

# What are the important choices and success factors?

- Model choice for embedding and inferencing
- Chunking strategy and size
- Vector similarity algorithm
- Preparation of the corpus
- Prompt engineering
- Model parameters
- Memory
- Re-ranking
- Agents
- Evaluation

This is an active area of research in industry and academia and there are many promising emerging techniques including:

- Re-ranking
- Prompt tuning
- Sentence Window Retrieval
- Retrieval Augmented Fine Tuning
- Parameter Efficient Fine Tuning
- Low Rank Adaptation

# Model choice

LLMs are trained for specific tasks and using defined datasets

- Models may be optimized for embeddings, inferencing, summarization, entity extraction, or other tasks
- Some models are **general purpose**, i.e., can perform different tasks

Model Type	Optimized Tasks	Examples
Embedding-focused LLMs	Generating embeddings for downstream tasks	BERT, RoBERTa, ALBERT, DistilBERT
Inference-focused LLMs	Fast and efficient inference	GPT-3, GPT-Neo, T5
General-purpose LLMs	Broad applicability across various tasks	GPT-3, GPT-Neo, T5, BERT, RoBERTa
Multimodal LLMs	Handling text and multimodal inputs	CLIP, ALIGN, LXMERT, UNITER
Multilingual LLMs	Handling multiple languages	mBERT, XLM-R, mT5, MarianMT
Domain-specific LLMs	Tailored for specific domains or tasks	BioBERT (biomedical), SciBERT (scientific text)

# Chunking

To obtain a `ConnectionFactory`, which supports both point-to-point and publish/subscribe operations, use `AQjmsFactory.getConnectionFactory()`. To obtain a `QueueConnectionFactory`, use `AQjmsFactory.getQueueConnectionFactory()`. To obtain a `TopicConnectionFactory`, use `AQjmsFactory.getTopicConnectionFactory()`.

The `ConnectionFactory`, `QueueConnectionFactory`, or `TopicConnectionFactory` can be created using hostname, port number, and SID driver or by using JDBC URL and properties.

### Using JNDI to Look Up ConnectionFactory Objects

A JMS administrator can register `ConnectionFactory` objects in a [Lightweight Directory Access Protocol \(LDAP\)](#) server. The following setup is required to enable [Java Naming and Directory Interface \(JNDI\)](#) lookup in JMS:

1. Register Database  
When the Oracle Database server is installed, the database must be registered with the LDAP server. This can be accomplished using the [Database Configuration](#)

ORACLE 7-2

Chapter 7  
Java Messaging Service Interface for Oracle Transactional Event Queues and Advanced Queuing

Assistant (DBCA). Figure 7-1 shows the structure of Oracle Database Advanced Queuing entries in the LDAP server. `ConnectionFactory` information is stored under `<cn=OracleDBConnections>`, while topics and queues are stored under `<cn=OracleDBQueues>`.

**Figure 7-1 Structure of Oracle Database Advanced Queuing Entries in LDAP Server**

```
graph TD
    A["<cn=acme, cn=com> (administrative context)"] --> B["<cn=OracleContext> (root of oracle RDBMS schema)"]
    B --> C["<cn=db1> (database)"]
    C --> D["<cn=OracleDBConnections> (Connection Factories)"]
    C --> E["<cn=OracleDBQueue> (Queues / Topics)"]
    C --> F["<cn=...> (Other db objects)"]
```

2. Set Parameter `GLOBAL_TOPIC_ENABLED`.

Chunking is more difficult than you'd expect, for example consider this chunk of a PDF

- It spans multiple topics which may not be related
- It includes headings, page headers and footers in the middle of the text
- It includes part of some sentences
- It does not understand the content of images at all

The quality of the chunks has a direct impact on the quality of the context provided to the LLM

# Structured Chunking

For example, using HTML headings to understand the structure of the document

- Captures a single topic per chunk
- Retains the hierarchical structure of the document
- Avoids inclusion of spurious text like page headers, footers, etc.
- May be further split into text chunks if too large for the chosen model

15 Using PL/Scope  
16 Using the PL/SQL Hierarchical Profiler  
17 Using PL/SQL Basic Block Coverage to Maintain Quality  
18 Developing PL/SQL Web Applications  
19 Using Continuous Query Notification (CQN)  
▼ Part IV Advanced Topics for Application Developers  
20 Choosing a Programming Environment  
21 Developing Applications with Multiple Programming Languages  
22 Using Oracle Flashback Technology  
23 Developing Applications with the Publish-Subscribe Model  
24 Using the Oracle Database ODBC Driver  
25 Using the Identity Code Package  
26 Microservices Architecture  
27 Oracle Backend for Spring Boot and Microservices  
28 Developing Applications with Sagas  
29 Using Lock-Free Reservation  
30 Developing Applications with Oracle XA  
31 Understanding Schema Object Dependency  
32 Using Edition-Based Redefinition  
33 Using Transaction Guard  
34 Table DDL Change Notification  
A.1 Appendix: Troubleshooting the Saga Framework  
B.1 Appendix: Troubleshooting UTL\_HTTP  
C.1 Appendix: Recording DML Changes on the Tracked Table

Transaction Failed

An important facet of the Saga pattern is the Saga coordinator. The Saga coordinator tells other participants what to do. The coordinator invokes Saga participants and with every response, the coordinator transitions to the next state. Asynchronous messaging is another important aspect of Sagas, which involves sending interservice messages over a queuing system.

The Saga pattern enables you to build more robust systems because Sagas use a failure management pattern. Every action has a compensating action for rollback, which helps ensure eventual data consistency and correctness across microservices.

### 26.4.2.1 Why Use Sagas?

Use Sagas for the following reasons:

- Perform a group of operations related to different microservices automatically.
- Rely on the Saga pattern to ensure data consistency across microservices and in different databases.
- Reduce the locking period and restrict the data locks to the duration of local transactions for improved concurrency.
- Avoid using Two-phase Commit (2PC) because of its extended locking period and performance constraints.
- Roll back or compensate if one of the transaction operations in the sequence of microservices fails.
- Use microservices that do not support 2PC.

**Why using Sagas is a better option than using 2PC transactions?**

- A 2PC transaction can cause extended locking period and incomplete results for queries until the 2PC transaction is committed or canceled. With Sagas, the data locks are placed only for the duration of the local transaction (that implements a microservice transaction), and not for the entire Saga lifecycle. Reducing the locking period improves the throughput of Saga transactions, resulting in better scalability for applications.
- Sagas are useful for long-lived transactions. 2PC is not ideal for long-lived transactions since the locking of resources for prolonged durations can affect performance and scalability.

When you use Sagas with Oracle Database, lock-free reservation features that are built into the database help improve concurrency and reduce bottlenecks.

### 26.4.2.2 Saga Implementation Approaches

There are two common saga implementation approaches, namely the orchestration and choreography models.

**Orchestration**

All communication between microservices is made through a centralized service called a Saga coordinator. The coordinator service is responsible for receiving the requests and calling the respective services. If any service fails, the coordinator service implements the roll back methods. You can use the orchestration model for complex workflows that need the Saga

# Sentence Window Retrieval

She lay on the operating table unconscious. Jones made a deep cut in her leg with a very sharp knife. He was able to remove the bullet and repair the damage to the blood vessels in the area. Ultimately, this saved her life.

Embedding vectors created from just the sentence to improve search

But the whole context is given to the LLM, not just the one sentence



# Understanding the context window

Each LLM has a specific size context window which determines the largest amount of information that it can work with

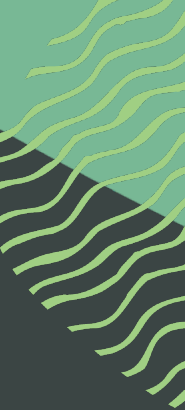
For example, GPT-3's context window is 2000 tokens

Context window
System instructions
Question
Context (documents from the RAG)
Memory (chat history, usually summarized)
Agent work area
Response

Suppose the window was 4k tokens in size, this could be used as follows:

- 512 tokens for chat memory
- 512 tokens for agent work area
- 512 tokens for system instructions, question and response

This leaves enough room for **only three** documents from the RAG if the chunk size is 800 tokens – smaller chunks means more documents but less context



# Prompt engineering

- Models expect specific prompt **formatting** including delimiters and placeholders
- The prompt can give the model **extra information** about the task to be performed
- Prompts have significant impact on the **quality** of responses
- **Significant effort** can go into optimizing the prompt for the intended use case

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>
```

```
You are a friendly, helpful and honest assistant.
```

```
You are given parts of a long document and a question. Use the provided document parts to answer the question.
```

```
If you do not know the answer, just say "I don't know". Do not make up an answer.
```

```
After your answer, tell the human which document parts you used to answer the question.
```

```
Question: {question}
```

```
Context:
```

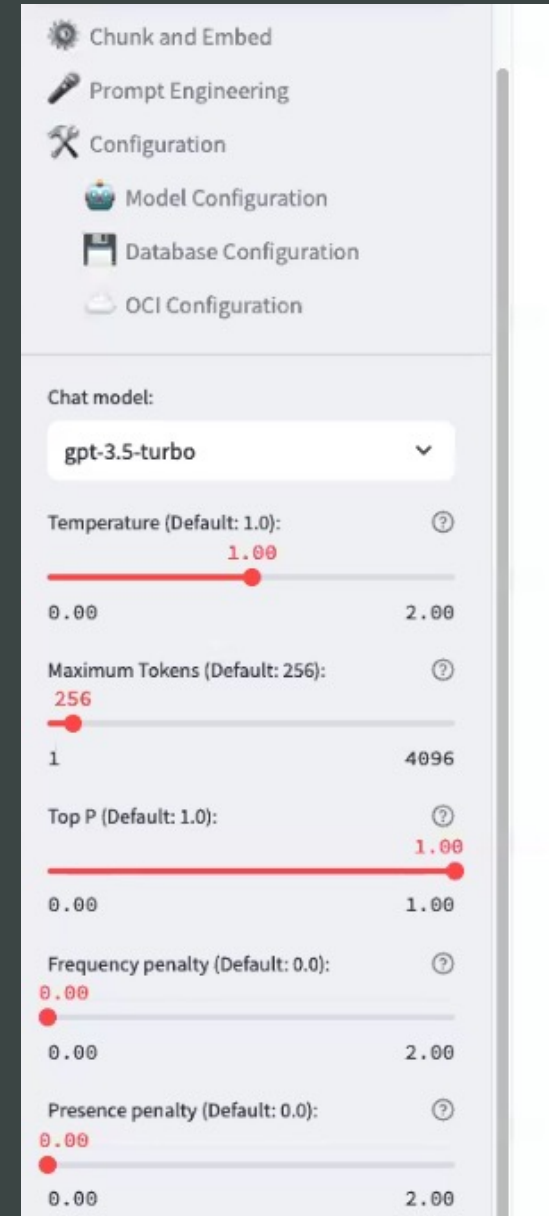
```
{context}<|eot_id|><|start_header_id|>assistant<|end_header_id|>
```

Example of a RAG prompt for Llama 3

# Model parameters

Models have parameters that affect the responses created, including:

- **Temperature** controls how “creative” the response is
- **Top K, Top P, frequency penalty, presence penalty**, and others control how many candidates are considered when choosing the next token
- **Maximum tokens** controls the length of the response





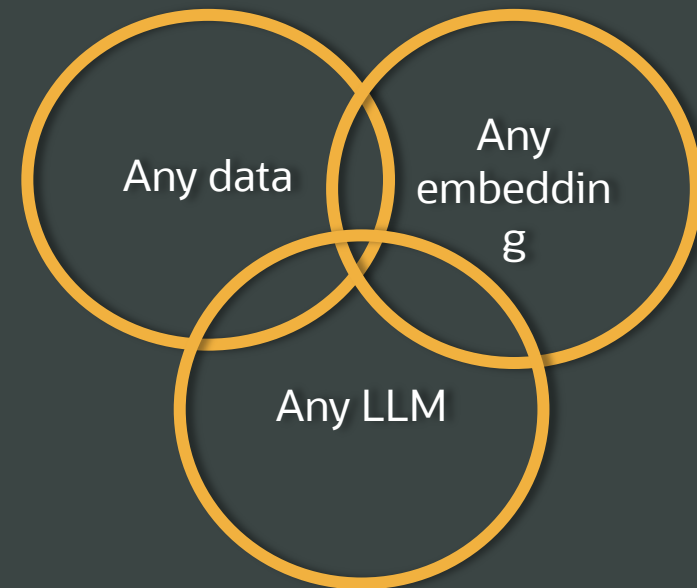
**Is there an approach  
that avoids the downsides?**

The most important success factor in developing a real chatbot is *experimentation* and *iteration*

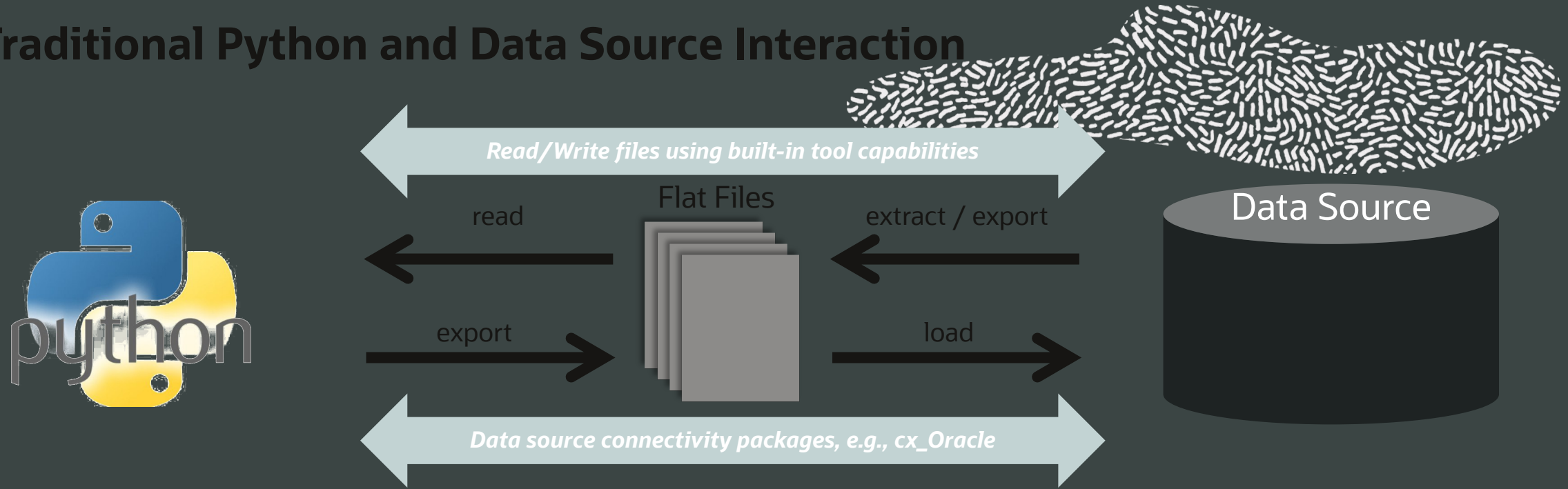
Oracle AI Playground

Developer release coming soon...

# ORACLE AI Playground



# Traditional Python and Data Source Interaction



Access latency

Paradigm shift: Python → *Data Access Language* → Python

Memory limitation – data size, in-memory processing

Single threaded

Issues for backup, recovery, security

Ad hoc production deployment



# Oracle Machine Learning for Python (OML4Py)

Empower data scientists and Python users on Oracle Database and Autonomous Database

Leverage your database as a high-performance compute engine

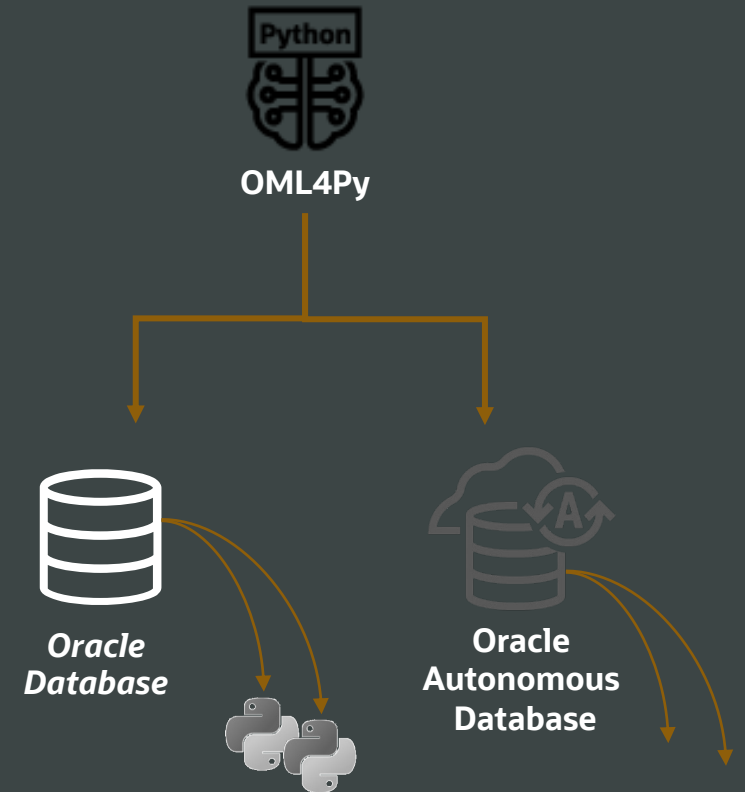
Use in-database optimized ML algorithms

Store and manage Python objects and user-defined functions in your database

Augment solutions with third-party Python packages

Invoke user-defined Python functions from SQL and REST

Use AutoML API with in-database ML algorithms

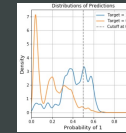


matplotlib

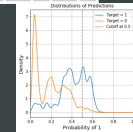


{REST:API}  
SQL

A	B	C	D	E	F



A	B	C	D	E	F



```
def my_function (arg):  
    import sklearn  
    ...  
    return result
```





# Embedded Execution

Example of parallel partitioned data flow using third party package using OML4Py

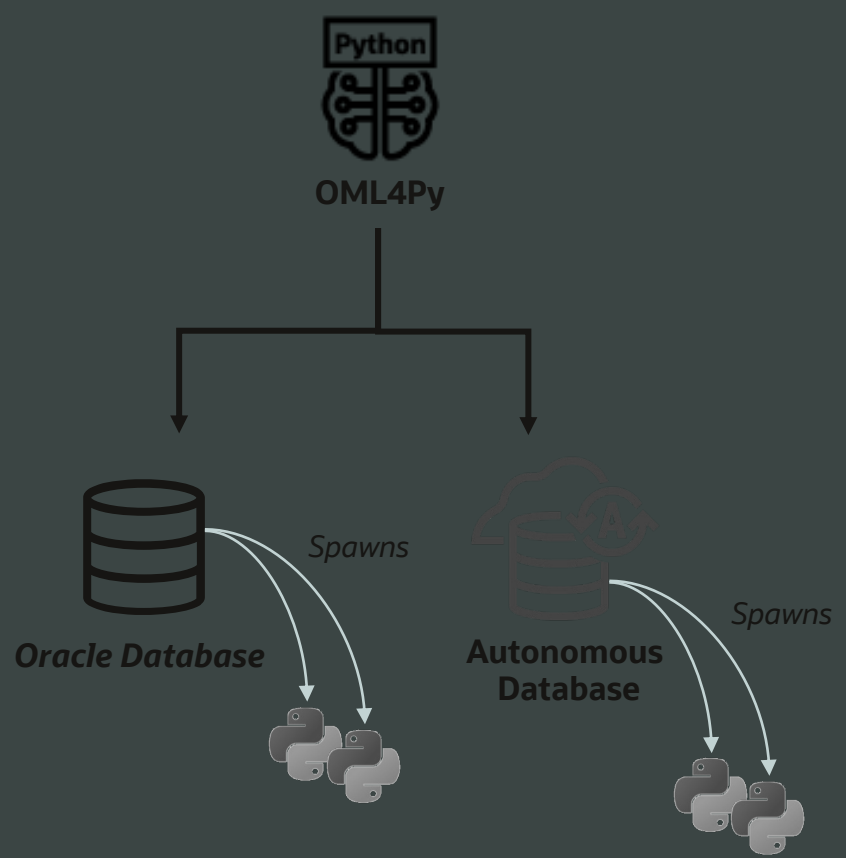


```
# user-defined function using sklearn
def build_lm(dat):
    from sklearn import linear_model
    lm = linear_model.LinearRegression()
    X = dat[['PETAL_WIDTH']]
    y = dat[['PETAL_LENGTH']]
    lm.fit(X, y)
    return lm

# select column(s) for partitioning data
index = oml.DataFrame(IRIS['SPECIES'])

# invoke function in parallel on IRIS table
mods = oml.group_apply(IRIS, index,
                       func=build_lm,
                       parallel=2)

mods.pull().items()
```



# Custom third-party packages on ADB via OML Notebooks

## Expanding ADB as a platform for data science and machine learning

Support 3rd party Python package installation and conda environment creation

- **conda** - open-source package and environment management system
- Admins install third-party packages and manage conda environments
- Users download and activate conda environments from Object Storage
- Environments run in a separate container for security

Use with Python interpreters in OML Notebooks

Use with embedded execution in OML4Py

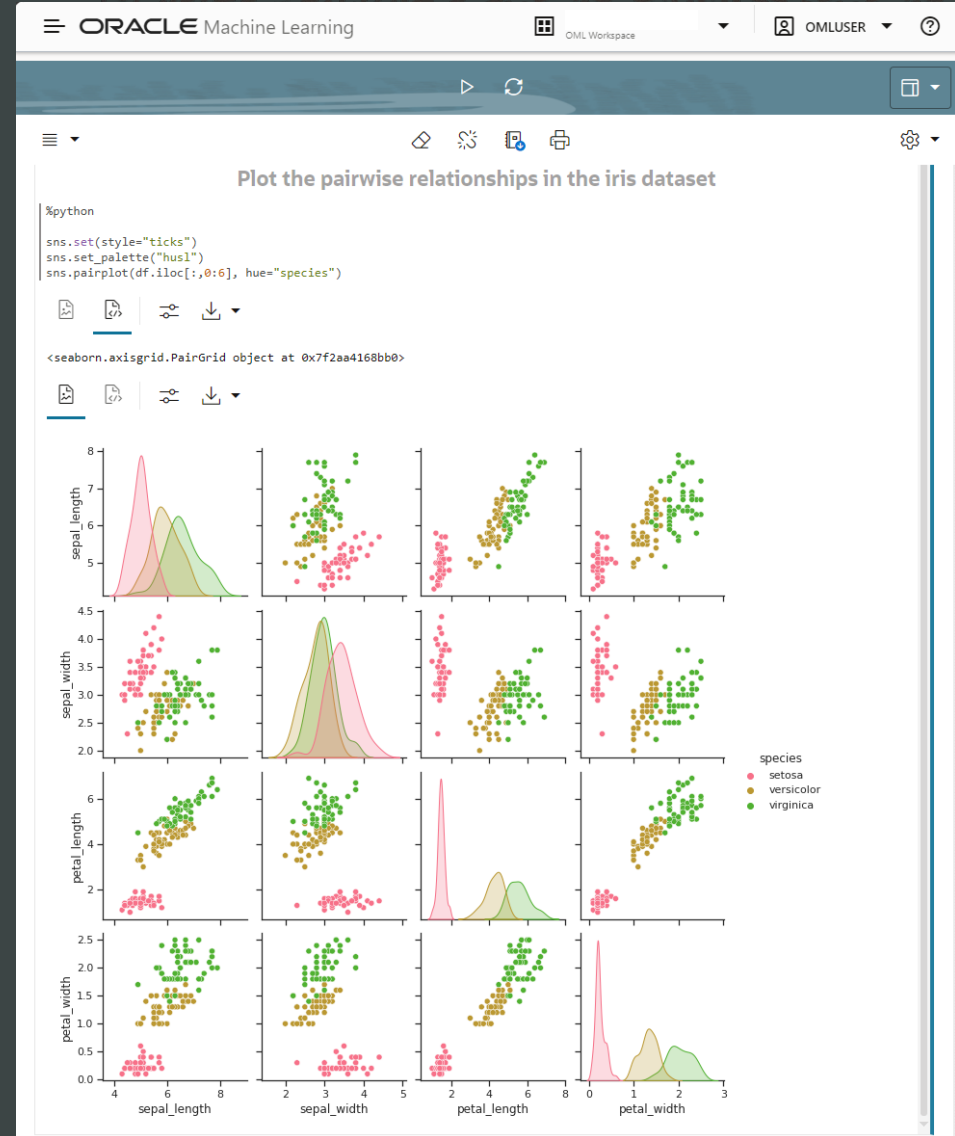


# Conda environment usage via OML Notebooks on ADB

Download and activate the 'mypyenv' environment

```
%conda  
download mypyenv  
  
activate mypyenv
```

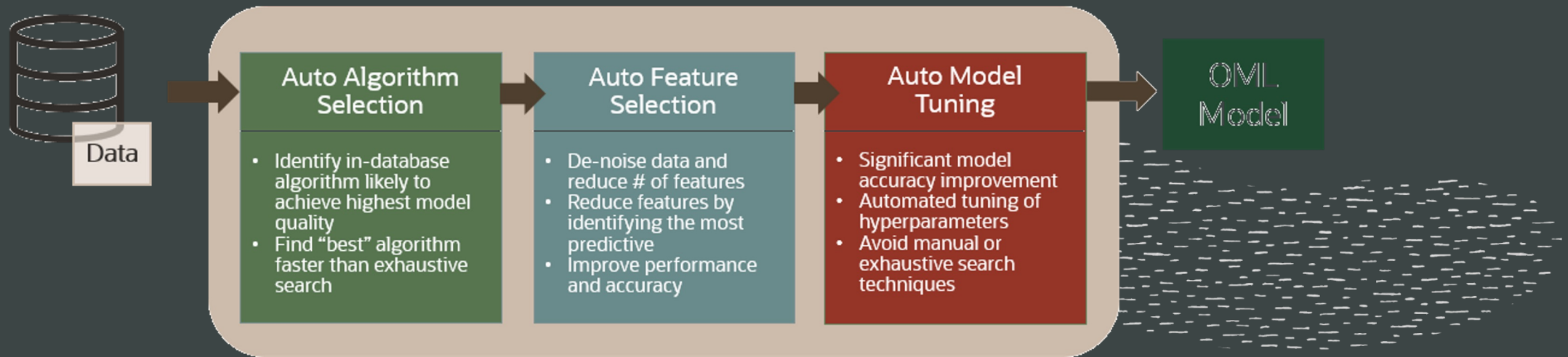
Downloading conda environment mypyenv  
Download successful for conda environment mypyenv



# OML4Py AutoML

Enhance data scientist productivity and enable Python users with ML

- Eliminate repetitive tasks of model building/evaluation to increase user productivity
- Enable non-expert users to leverage machine learning
- Apply ML to the ML process to reduce algorithm and hyperparameters search space and reduce compute time and cost



# OML4Py AutoML API



## Algorithm Selection

```
%python
as_wine_cl = automl.AlgorithmSelection(mining_function='classification',
                                       score_metric='accuracy', parallel=2)

wine_alg_ranking_cl = as_wine_cl.select(WINE_X_cl, WINE_y_cl, k=4)

print("Ranked algorithms:\n", wine_alg_ranking_cl)

selected_wine_alg_cl = next(iter(wine_alg_ranking_cl))
print("Best algorithm: ", selected_wine_alg_cl)

Ranked algorithms:
[('svm_gaussian', 0.98255159474672114714554), ('rf', 0.949547038327114714554)]
Best algorithm: svm_gaussian
```

## Feature Selection

```
%python
fs_wine_cl = automl.FeatureSelection(mining_function = 'classification',
                                     score_metric = 'accuracy', parallel=2)

selected_wine_features_cl = fs_wine_cl.reduce(selected_wine_alg_cl,
                                              WINE_X_cl, WINE_y_cl)

WINE_X_reduced_cl = WINE_X_cl[:,selected_wine_features_cl]

print("Selected columns:", WINE_X_reduced_cl)
print("Number of columns:")
print("{} reduced to {}".format(WINE_X_cl.shape[1], WINE_X_reduced_cl.shape[1]))

Selected columns: ['alcohol', 'ls', 'color_intensity', 'hue']
Number of columns:
'13 reduced to 9'
```

## Model Tuning

```
%python
mt_wine_cl = automl.ModelTuning(mining_function = 'classification', parallel=2)

results_cl = mt_wine_cl.tune(selected_wine_alg_cl, WINE_X_reduced_cl, WINE_y_cl)
tuned_model_cl = results_cl['best_model']
tuned_model_cl

Algorithm Name: Support Vector Machine

Mining Function: CLASSIFICATION
```



# Generate notebooks with Python code using AutoML UI

Enhance data scientist productivity and enable non-expert data professionals

Accelerate ML projects with no-code AutoML

Automate repetitive and time-consuming tasks

Compare multiple models using multiple metrics

Generate editable notebooks for desired models based on OML4Py

Deploy ML models immediately using SQL or to OML Services as REST endpoints

The screenshot displays the Oracle Machine Learning AutoML UI interface. The main window shows the 'Wine Cultivators' experiment, which is completed. A modal window indicates the completion of various steps: Algorithm Selection, Adaptive Sampling, Feature Selection, Model Tuning, Neural Network, and Support Vector Machine (Linear). Below this, a 'Leader Board' table compares three models: Neural Network, Support Vector Machine (Linear), and Random Forest. The Neural Network model shows the highest performance with a Balanced Accuracy of 0.9956 and an F1 Micro score of 0.9946. The interface also includes a sidebar with experiment details and a bottom section with Python code for data preparation and model building.

Algorithm	Model Name	Balanced Accuracy	Accuracy	F1 Micro
Neural Network	NN_82861092F4	0.9956	0.9946	0.9946
Support Vector Machine (Linear)	SVML_6CC1183B39	0.9944	0.9941	0.9941
Random Forest	RF_A5C5B15A1D	0.9744	0.9715	0.9715

```
%python
import warnings
import oml
warnings.simplefilter(action='ignore', category=FutureWarning)

columns = ["total_phenols", "acidity_of_ash", "alcohol"]
schema = "OMLUSER02"
table = "WINE"

column = ','.join(columns)
query = 'SELECT ' + column + ' FROM ' + schema + '.' + table + ';'

build_data = oml.sync(query=query)
z.show(build_data)
```

```
%python
X_train = build_data.drop('target')
y_train = build_data[:, 'target']
```

```
%python
svm_settings = {
    'SVMS_COMPLEXITY_FACTOR': '10', 'SVMS_KERNEL_FUNCTION': 'SVMS_GAUSSIAN',
    'SVMS_STD_DEV': '2.121320343596424', 'CLAS_WEIGHTS_BALANCED': 'OFF',
    'SVMS_NUM_PIVOTS': '200', 'ODMS_DETAILS': 'ODMS_ENABLE',
    'ODMS_TEXT_POLICY_NAME': 'AML_svm_19c8cc7b864947a591d6'
}

svm_mod = oml.svm(**svm_settings)
svm_mod = svm_mod.fit(X_train, y_train)
```



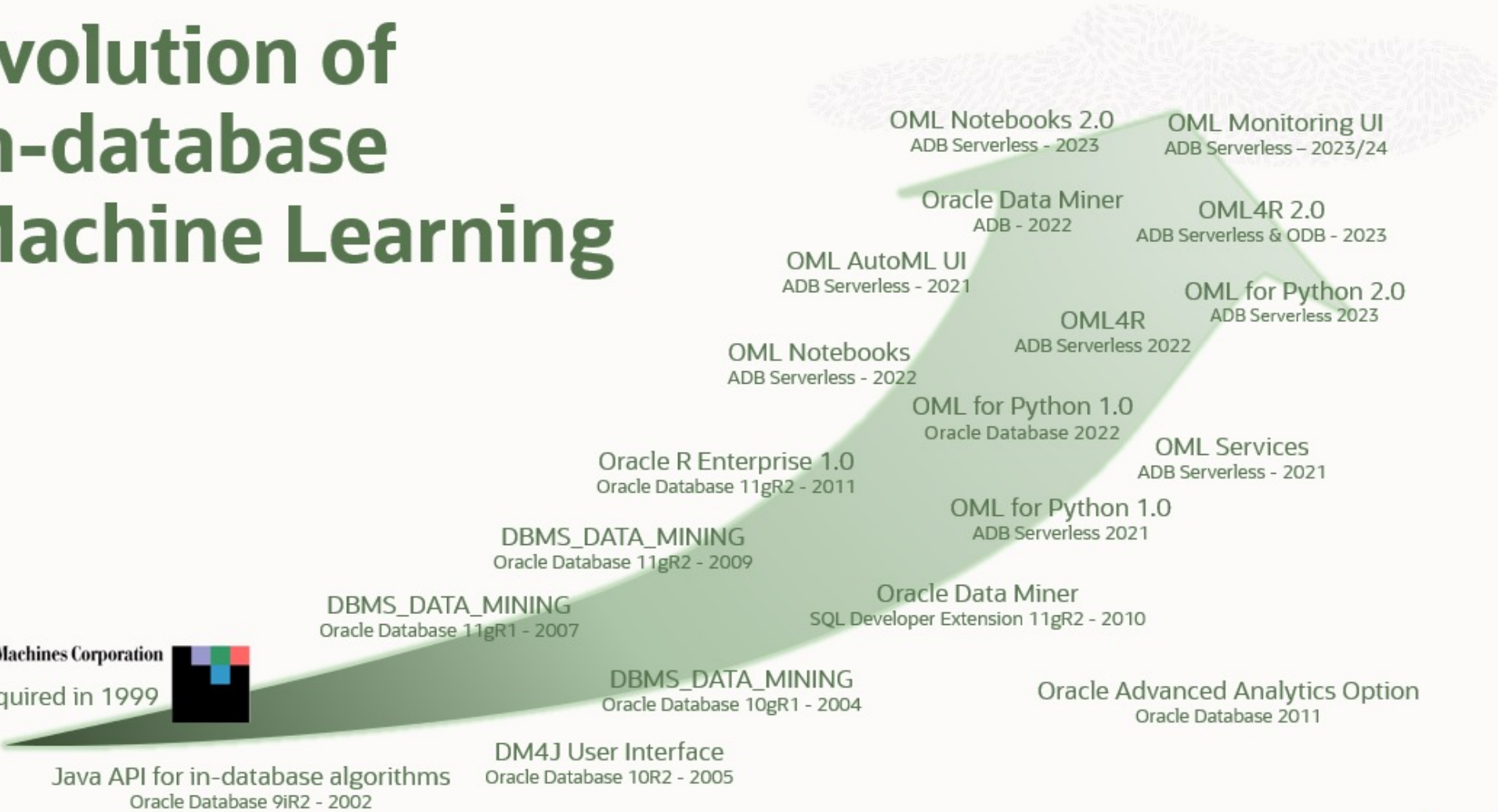


# Evolution of in-database Machine Learning

Thinking Machines Corporation



Acquired in 1999



# Oracle Spatial

## Native spatial type and spatial analysis in Autonomous Database

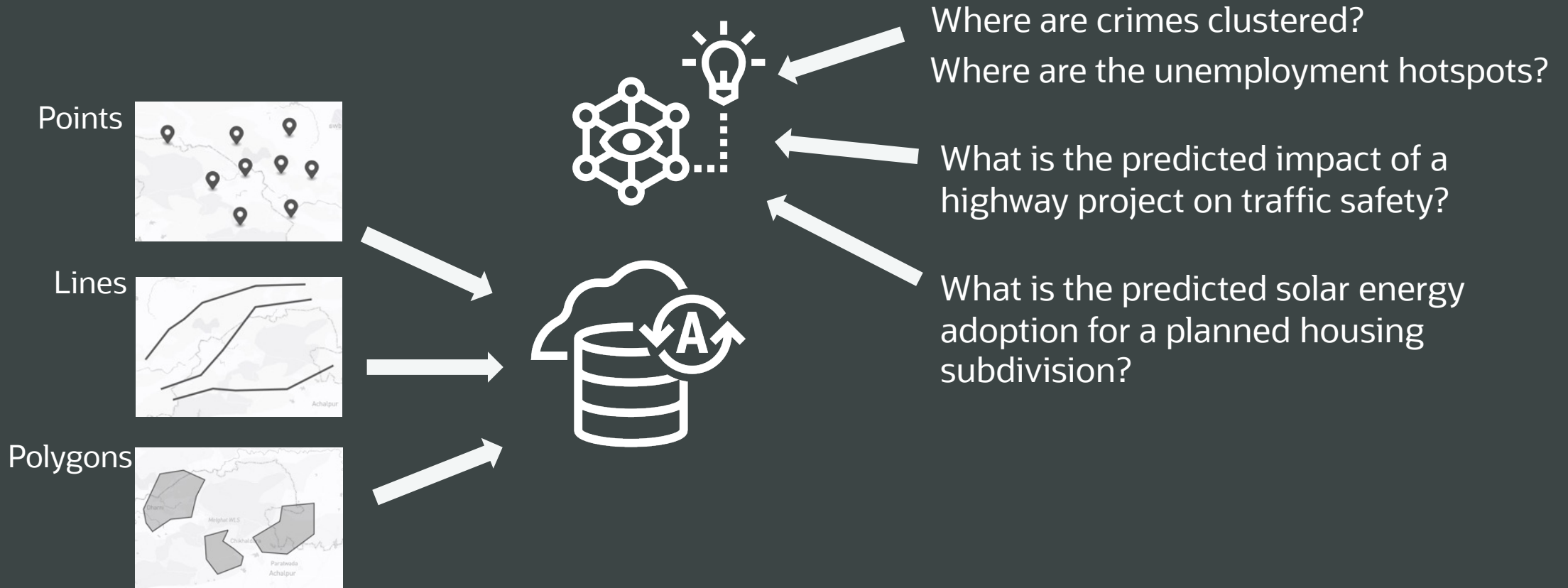


Also 3D, point clouds, raster imagery, spatial networks and more



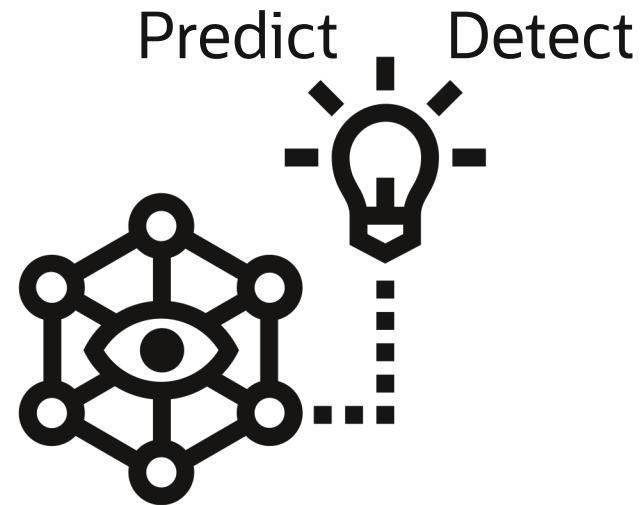
# Oracle Spatial AI

## Spatial algorithms in Oracle Machine Learning for Python (OML4Py)



# Spatial AI

“Derive insights from spatial data”



# Common scenarios addressed by spatial algorithms

Scenario	Jargon	Example
Predictive influence of nearby values	“Spatial dependence”	Property values influenced by nearby values
Regional variation of predictive factors	“Spatial heterogeneity”	Predictive factors for vehicle purchase vary regionally
Geographic concentration of observations	“Spatial clustering”	Crime clusters, unemployment hotspots
Geographic alignment of observations	“Spatial colocation”	Successful sandwich shops tend to be in close proximity to desert shops

# Oracle Spatial AI

Oracle Machine Learning for Python

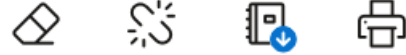
Spatial AI

Autonomous Database Serverless

- Python API for common Spatial SQL operators and functions
- Spatial pre/post processing
  - Feature engineering
  - Gap filling
  - Spatial weights/spatial lag
- Spatial algorithms
  - Regression
  - Classification
  - Clustering
  - Anomaly detection
  - Colocation
- Spatial pipeline
- Python, REST, SQL APIs to invoke models

## Detect density clusters in data from Oracle Spatial

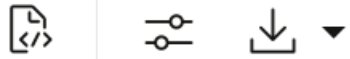
Versioning



low Zeppelin

%python

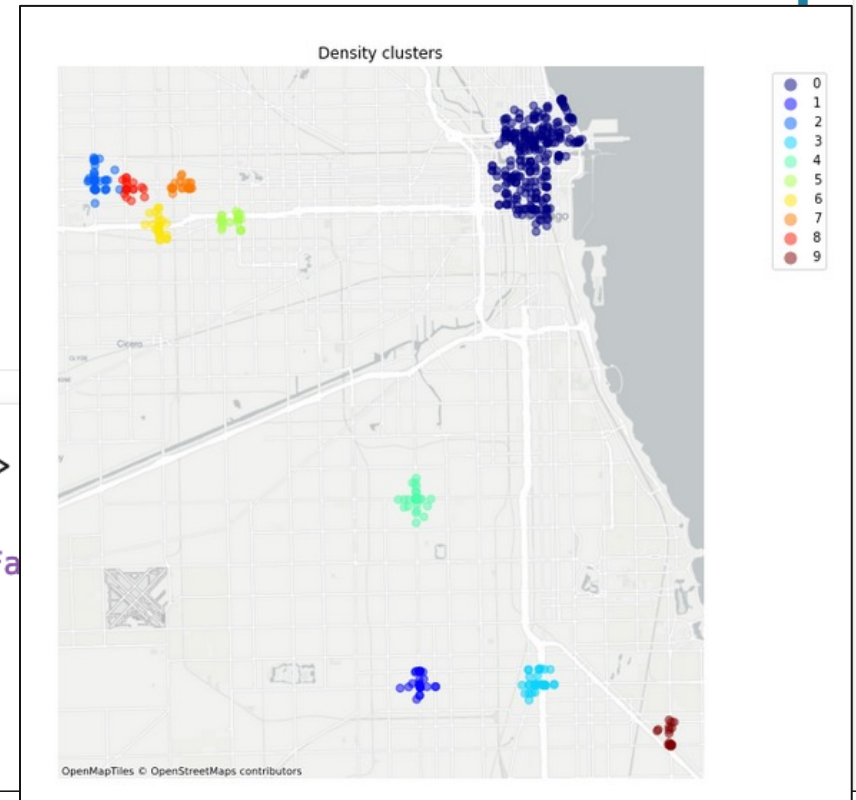
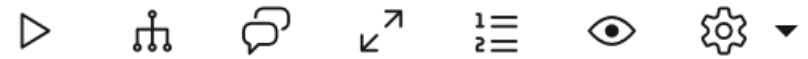
```
X = accidents_injury_sdf[['geometry']].to_crs("EPSG:3857")
labels = dbscan.fit_predict(SCoordTransformer().transform(X))
np.unique(labels)
```



array([-1, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9])

%python

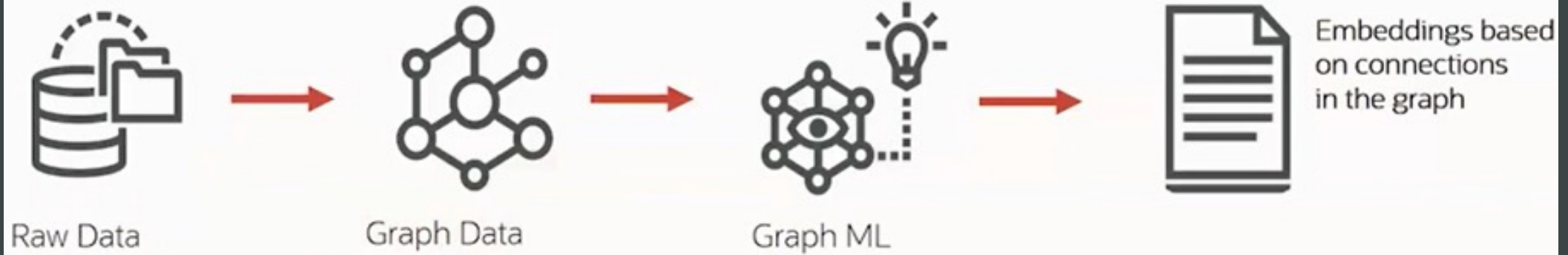
```
_, axs = plt.subplots(figsize=(10, 10))
plot_clusters(X, labels, title='Density clusters', with_noise=False, with_bounds=False)
```



# Spatial AI

- Library added to OML4Py on ADB-S
- Python API for common Oracle Spatial operations
- Spatial pre/post-processing
- Spatial ML algorithms and APIs

# Graphs and Machine Learning

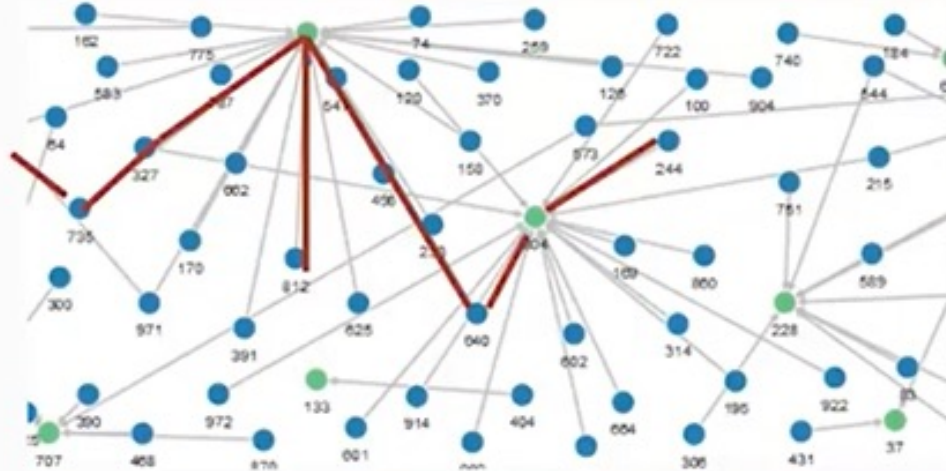


## Embeddings capture

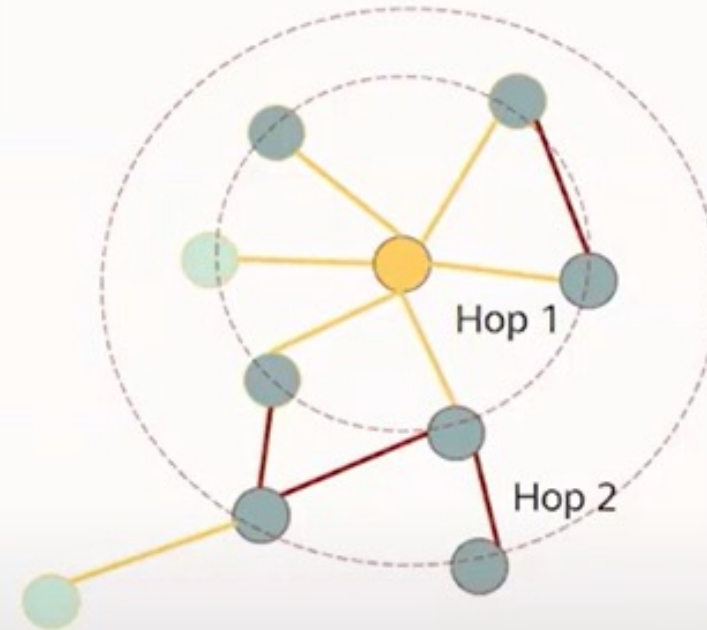
- Graph topology
- Node and edge features
- How information propagates in a graph: How node and edge features are influenced by neighbors

# Graphs and Machine Learning

DeepWalk: Follow walks from each node



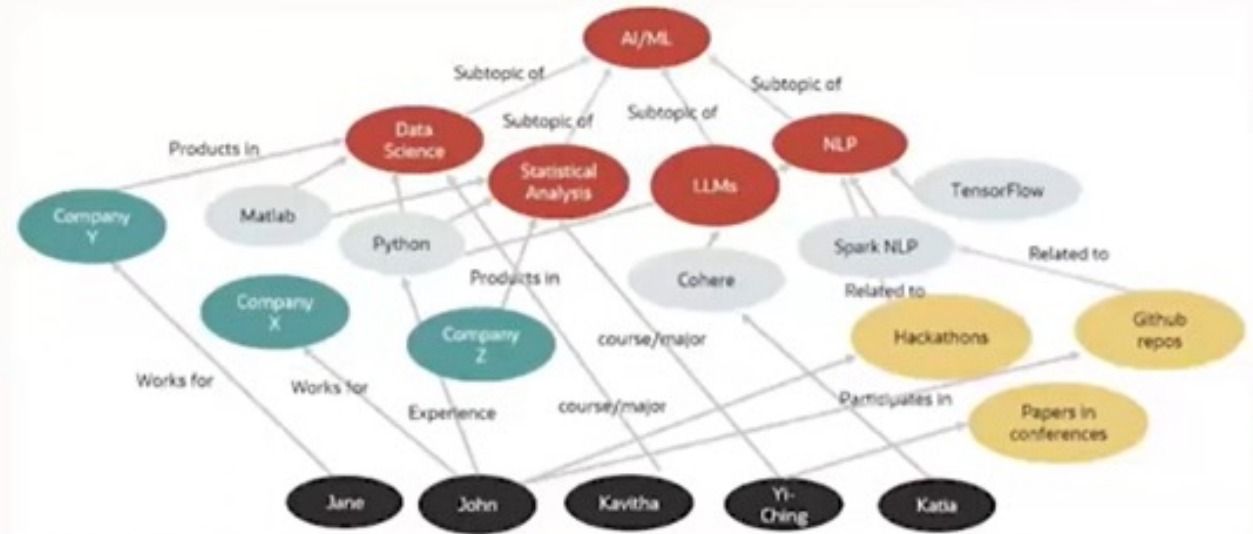
GraphWise: Sample neighbors recursively to learn a function





# Knowledge Graphs and Retrieval Augmented Generation (RAG)

- Knowledge graphs capture expert information about a domain
- Generate embeddings from the knowledge graph using Graph ML techniques
- Use vector search for fast retrieval of top results
- Combine with user query in a prompt to LLMs



# Recap, what we covered...

- OCI AI Services Oracle AI Stack
- Oracle Multi-purpose/Converged database
- Select AI (NL2SQL)
- Vector and RAG
- AI Agents
- AI Playground
- OML4PY
- Spatial AI
- Graph ML and RAG



# Thanks for Attending!

Oracle Database 23ai



AI Solutions Hub



LiveLabs workshops



Paul Parkinson  
Architect and Developer Advocate  
Oracle Database



## Develop AI/ML-Driven Apps with Autonomous Database 23ai

