# The Oracle AI Microservices Sandbox for RAG Rapid Prototyping

Corrado De Bari, Mark Nelson, John Lathouwers

# Future & Past TechCasts:

**Dec 12th**
Unlocking Insights: Mastering Data Storytelling with Oracle Analytics
Presented by **Philip Godfrey**

**Jan 9th**
The Oracle AI Microservices Sandbox for RAG Rapid Prototyping
Presented by **Corrado De Bari, Mark Nelson, & John Lathouwers**

**Jan 23rd**
Leveraging Vector Search for RAG in Generative AI
Presented by **Kai Yu**

## TechCast Archive

| 2024 | 2023 | 2022 | 2021 | 2020 | 2019 |
|------|------|------|------|------|------|

| Date | Title | Presenter(s) | Replay | Download(s) |
|------|-------|--------------|--------|-------------|
| Nov 21 | Backup, Cloning and DR for Oracle Analytics Cloud | Jason Lester | Video | Slides |
| Nov 7 | Gimme a Vector, Victor: Leveraging Vector Datatypes for Practical Generative AI Applications | Jim Czuprynski | Video | Slides |
| Oct 17 | Our Favorite New Features in OAC | Dan Vlamis, Wayne Van Sluys, Cathye Pendley, Tim Vlamis, Mystery Guest: Gautam Pisharam | Video | Slides |
| Oct 3 | Harnessing Oracle GoldenGate 23ai with AI Automation: A New Era of Intelligent Monitoring, Building, and Diagnostics | Bobby Curtis | Video | Slides |
| Sept 19 | Discovering Oracle Fusion Data Intelligence | Peter Koutroubis & Jai Gangwani | Video | Slides |
| Sept 5 | Oracle CloudWorld 2024 Session Analysis and Expert Agenda Recommendations | Roger Cressey, Dan Vlamis, Jim Czuprynski, Tim Vlamis, Cathye Pendley | Video | Slides |

Submit a topic to share at **https://andouc.org/techcasts/**

Analytics and Data
ORACLE USER COMMUNITY

**www.andouc.org**

Oracle Spatial & Graph SIG

# Let's Connect

**Website**
http://andouc.org/

**Chat with the Experts**
https://bit.ly/Join-ANDOUC-Slack

**Watch Previous TechCasts**
https://bit.ly/3qmGgHN

**@AnalyticAndData**

https://www.facebook.com/AnDOracleUserCommunity

https://www.linkedin.com/company/analytics-and-data-oracle-user-community

**Spatial + Graph SIG**
bit.ly/Spatial-Graph-LinkedIn

**The Oracle Analytics**

# Data Visualization Challenge 2025

Get recognized for your Oracle Analytics skills!

Running February 3rd – 28th 2025

**Registration open NOW**

*Save the Date*

Analytics and Data
Summit 2025

April 8-10, 2025
Oracle Conference Center
Redwood Shores, California

https://andouc.org/andsummit2025/

Registration is
now **open**!

ORACLE

# The Oracle AI Microservices Sandbox for RAG Rapid Prototyping

**Mark Nelson, John Lathouwers, Corrado De Bari**
**Developer Evangelists**
**Oracle**

# 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.

Everybody is building chatbots
But what defines enterprise quality?

# 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?

# 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"

**ORACLE**

# AI Microservices Sandbox

*A Zero-Code platform to build, test & deploy a RAG Chatbot with:*

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

- Any data
- Any embedding
- Any chat LLM

# Oracle AI Microservices Sandbox



**Open Source Project in**
*Developer Preview*

https://github.com/oracle-samples/oaim-sandbox

# Building a RAG+ Chatbot

- **Playground env for RAG agent configuration: 0-Code approach**

- **Any LLM model (embeddings/chat):**
  - **PRIVATE**: On-premises with OLLAMA, HF TEI, HF TGI
  - **PUBLIC**: service providers like Cohere, OpenAI, Perplexity

- **Vector Store ingestion for DB23ai**:
  - Split & Model configuration
  - Documents: HTML, Txt, Pdf docs
  - From: local/OCI Object Storage/Web URL

- **Prompt configuration**

- **Test framework**:
  - Auto Generation for massive tests
  - Self Evalution

- **REST API Server** (OpenAI API compatible)

- **Export as SpringAI Microservice**

# The Oracle DB 23ai as knowledge base:

NEW vector indexes provide 100x faster and highly accurate similarity searches

private models

Entry Point

Oracle database accelerates AI vector search using sophisticated vector indexes

Query Vector

**Graph Vector Index** (e.g., **HNSW** Index)

**Partition Vector Index** (e.g., **IVF_FLAT** index)

# Experiment a Chatbot Application (in Sandbox)

# An Example of improving a Chatbot application (Before)

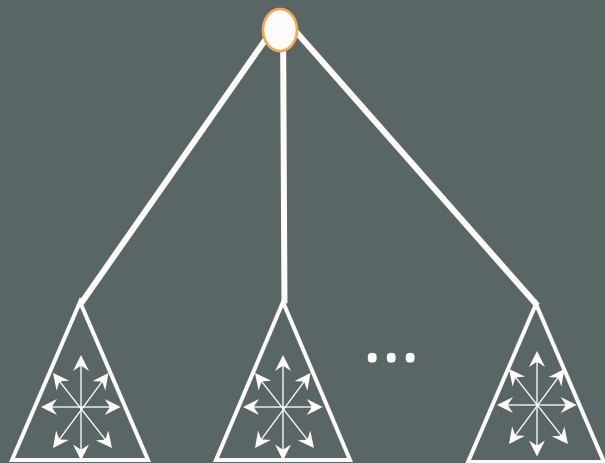**Terminal**

You: How do I determine the accuracy of my vector indexes?

Bot: To determine the accuracy of your vector indexes, you can calculate various evaluation metrics like precision, recall, and F1-score…

You: How do I do it in Oracle Database 23ai?

Bot: I'm sorry, but it seems you mentioned Oracle Database 23ai. As of my last update, the latest version of Oracle Database is Oracle Database 19c…

*Knowledge Cutoff*

You: In Oracle Database 23ai, how do I determine the accuracy of my vector indexes?

Bot: **To determine the accuracy of your vector indexes** in Oracle Database 23ai, you can use the `DBMS_DATA_MINING.DETECT_INDEX` procedure…

*Hallucinations*

You: Please give me more details.

Bot: Of course! More details about what specifically? Feel free to provide more information so I can better assist you.

*No Memory*

gpt-3.5-turbo

# Experiment a Chatbot Application (in Sandbox)

**Terminal**

You: How do I determine the accuracy of my vector indexes?

Bot: To determine the accuracy of your vector indexes **you can utilize the index accuracy reporting feature** provided by the DBMS_VECTOR package.

◁ *Live Knowledge*

You: How do I do it in Oracle Database 23ai?

Bot: In Oracle Database 23ai, **you can determine the accuracy of your vector indexes**...

◁ *Memory/Context*

You: In Oracle Database 23ai, how do I determine the accuracy of my vector indexes?

Bot: In Oracle Database 23ai, you can determine the accuracy of your vector indexes **using the DBMS_VECTOR.INDEX_ACCURACY_QUERY** procedure...

◁ *Grounded Response*

You: Please give me more details.

Bot: ...This information **is based on the documentation in the oracle-ai-vector-search-users-guide.pdf** regarding index accuracy reporting features in Oracle Database 23ai...

◁ *RAG Attribution*

# RAG+ Testset Generation & Evaluation

Generation



Reference
Document

Contexts Collection

independent
LLM for
Q&A generation

**Question:**
**Reference Answer**:
Context:
Metadata:
**Topic**:
Ref:

Synthetic
Q&A collection

List of topics

# RAG+ Testset Generation & Evaluation

Evaluation

**Overall Correctness Score**

**Question:**
Reference Answer:
Context:
Metadata:
        Topic:
        Ref:

*Q&A fix*

*Synthetic*
*Q&A collection*

*AI Microservices*
*Sandbox Agent*

*Update:*
*Knowledge base/*
*Agent Conf*

*Vector Store*
*+*
*RAG conf*

Question:
Reference Answer:
**Actual Answer:**
Context:
Metadata:
        Topic:
        Ref:

*Q & Agent's Answer*
*collection*

*Independent*
*LLM as a judge*

Question:
Reference Answer:
Actual Answer:
**Correctness**:[T/F]
**Correctness Reason**:
Context:
Metadata:
        Topic:
        Ref:

*Score for topics*

# Oracle Backend for Microservices and AI v1.3
Working with AI in Microservices – Get Starters, Build, Optimize, & Deploy

**1** — Get Spring AI and LangChain starters,
Get LLM API Keys
Export from OIAM sandbox

**2** — Build AI microservices with REST
endpoints in IDE/Python.
**SpringAI**

**3** — Optimize AI microservices
in AI Sandbox

**4** — Deploy AI microservices
In App Containers

Oracle Backend for Microservices
and AI
Platform Services Containers

**LLM services**
OpenLLM
OLLAMA
TEI - HF

App Container   App Container   App Container

**Kubernetes**

**Oracle Database**

**AI Developers**

**DevOps**

Available **free: on-premise, OCI and Azure.**
*support through an Oracle DBMS contract*

# Summary

- Rag platform for GenAI Chatbot building
- Zero-code
- Base on Oracle DB 23ai high scalable Vectorstore
- Any LLM for Embeddings & Chat
- Well suitable for Private and/or On-premises deployement
- Free and Open Source

It's now time for Q&A

Got any questions?

Thank you!

# Helpful Links –

**ORACLE ANALYTICS VIDEOS:**
https://www.youtube.com/@OracleAnalytics/videos

**OAC SEPTEMBER NEW FEATURES VIDEOS BY ORACLE:**  https://bit.ly/OACSept24Features

**OAC NEW FEATURES  DOCUMENTATION BY ORACLE:**
https://docs.oracle.com/en/cloud/paas/analytics-cloud/acswn/index.html#GUID-CFF90F44-BCEB-49EE-B40B-8D040F02D476

**ORACLE ANALYTICS COMMUNITY:**
https://community.oracle.com/products/oracleanalytics

**ORACLE ANALYTICS LIBRARY/EXAMPLES:**
https://www.oracle.com/business-analytics/data-visualization/examples/

**ORACLE ANALYTICS LIVE DEMOS:**
https://www.oracle.com/business-analytics/data-visualization/demos/

# Future & Past TechCasts:

**Dec 12th**
**Unlocking Insights: Mastering Data Storytelling with Oracle Analytics**
Presented by **Philip Godfrey**

**Jan 9th**
**The Oracle AI Microservices Sandbox for RAG Rapid Prototyping**
Presented by **Corrado De Bari, Mark Nelson, & John Lathouwers**

**Jan 23rd**
**Leveraging Vector Search for RAG in Generative AI**
Presented by **Kai Yu**

## TechCast Archive

| 2024 | 2023 | 2022 | 2021 | 2020 | 2019 |
|------|------|------|------|------|------|

| Date | Title | Presenter(s) | Replay | Download(s) |
|------|-------|--------------|--------|-------------|
| Nov 21 | Backup, Cloning and DR for Oracle Analytics Cloud | Jason Lester | Video | Slides |
| Nov 7 | Gimme a Vector, Victor: Leveraging Vector Datatypes for Practical Generative AI Applications | Jim Czuprynski | Video | Slides |
| Oct 17 | Our Favorite New Features in OAC | Dan Vlamis, Wayne Van Sluys, Cathye Pendley, Tim Vlamis, Mystery Guest: Gautam Pisharam | Video | Slides |
| Oct 3 | Harnessing Oracle GoldenGate 23ai with AI Automation: A New Era of Intelligent Monitoring, Building, and Diagnostics | Bobby Curtis | Video | Slides |
| Sept 19 | Discovering Oracle Fusion Data Intelligence | Peter Koutroubis & Jai Gangwani | Video | Slides |
| Sept 5 | Oracle CloudWorld 2024 Session Analysis and Expert Agenda Recommendations | Roger Cressey, Dan Vlamis, Jim Czuprynski, Tim Vlamis, Cathye Pendley | Video | Slides |

Submit a topic to share at **https://andouc.org/techcasts/**

**Analytics and Data**
ORACLE USER COMMUNITY

www.andouc.org

Oracle Spatial & Graph SIG

# *Save the Date*

## Analytics and Data
## Summit 2025

April 8-10, 2025
Oracle Conference Center
Redwood Shores, California

https://andouc.org/andsummit2025/

*Registration is now **open**!*