

Custom MCP Server For Oracle 26ai

1. Executive Summary

This project establishes a secure, high-performance bridge between Large Language Models (LLMs) and the **Oracle 26ai** database. By architecting a custom **Model Context Protocol (MCP)** server in Python, we eliminate the high costs and vendor lock-in associated with cloud-managed middleware. This approach grants us full control over the AI's interaction with our data, allowing for specialized **AI Vector Search** tools that are perfectly aligned with our unique enterprise data schema and security requirements.

2. System Architecture

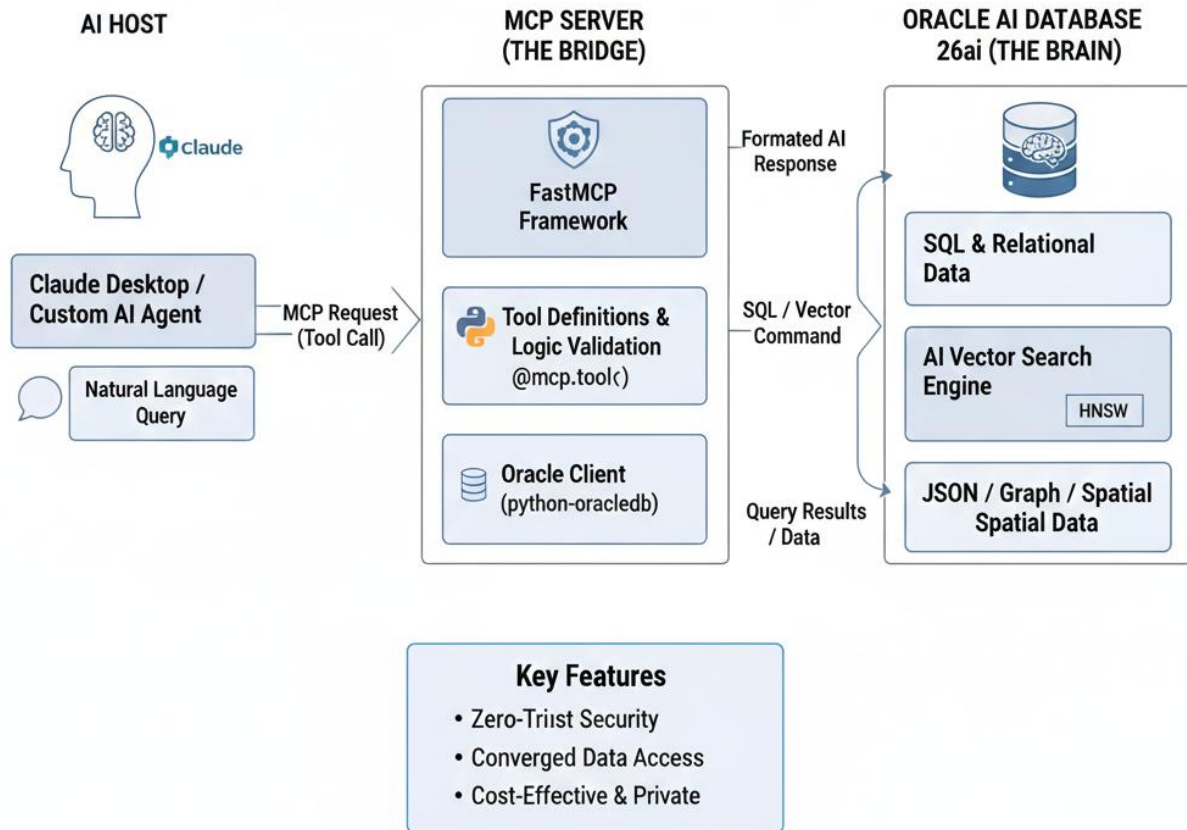
The system is built on a **three-tier architecture** designed for maximum isolation and performance. At the center sits the **MCP Server**, which functions as an intelligent "Interpreter." This layer is critical because it ensures the AI Host never has direct access to the database credentials or the raw storage engine. Instead, the AI only interacts with a set of pre-defined, secure "Tools."

The **AI Host** (such as Claude Desktop or a custom agent) initiates requests, the **MCP Server** (built with Python and FastMCP) processes the logic and security checks, and the **Oracle 26ai Database** (The Brain) performs the heavy lifting of relational queries and high-speed vector similarity searches.

1. **Host (The AI):** Claude Desktop or a custom Python agent.
2. **Server (MCP Bridge):** A Python application using the **FastMCP** framework.
3. **Database (The Brain):** Oracle 26ai (Local or Cloud) utilizing **AI Vector Search**.

Why This Architecture is Cost-Free

- Oracle 26ai Free: No license fees (limited to 2 CPU, 2GB RAM, 12GB data)
- MCP: Open-source protocol by Anthropic
- Python/Node.js SDKs: Free and open-source
- All Components: Can run on a single machine



3. Implementation Roadmap

Phase 1: Environment Preparation

The first step involves configuring a lightweight Python environment. We utilize the **python-oracledb** driver, which is the gold standard for high-performance Oracle connectivity. Crucially, we use it in "Thin Mode," which removes the need for bulky Oracle Instant Client installations. Alongside the driver, we install the **MCP SDK** to handle the standardized communication protocol that allows the AI to discover and use our database tools.

Phase 2: Database Layer (Oracle 26ai)

In this stage, we prepare the database to handle semantic meaning. We create a specialized knowledge base table that utilizes the native VECTOR data type introduced in **Oracle 26ai**. This allows the database to store mathematical representations of text or images.

```
CREATE TABLE ai_knowledge_base (
  id NUMBER GENERATED BY DEFAULT AS IDENTITY PRIMARY KEY,
  content_text CLOB,
  metadata JSON,
  embedding VECTOR(1024, FLOAT32) -- Oracle 26ai Native Vector Type);
```

Phase 3: The Custom MCP Server (Python)

The heart of the system is the `oracle_mcp.py` script. This script defines the "Tools" that the AI is allowed to use. For example, a `semantic_search` tool allows the AI to ask the database for "concepts" rather than just exact word matches. The server connects to Oracle using secure environment variables, executes the specialized `VECTOR_DISTANCE` SQL syntax, and returns only the most relevant results back to the AI.

```
from mcp.server.fastmcp import FastMCP
import oracledb
import os

# Initialize FastMCP
mcp = FastMCP("Oracle_Enterprise_Bridge")

# Database Connection Logic
def get_db():
    return oracledb.connect(
        user="SYSTEM",
        password="your_password",
        dsn="localhost:1521/FREEPDB1")

@mcp.tool()
def semantic_search(user_query: str):
    """Searches Oracle 26ai for concepts related to the user query."""
    conn = get_db()
    cursor = conn.cursor()

    # Note: In production, use an embedding model to convert query_text to a vector
    # Here we show the Oracle 26ai VECTOR_DISTANCE syntax
    sql = """
    SELECT content_text
    FROM ai_knowledge_base
    ORDER BY VECTOR_DISTANCE(embedding, :v, COSINE)
    FETCH FIRST 3 ROWS ONLY

    """

    # result = cursor.execute(sql, [generated_vector])
    return "Results from Oracle 26ai matching your query..."

if __name__ == "__main__":
    mcp.run()
```

4. Technical Features & Capabilities

The most significant feature of this system is **Zero-Trust Security**. Because the connection strings and sensitive logic are hidden inside the Python server, the LLM only sees the names of the functions it can call. Furthermore, this system enables **Converged Data Access**. Unlike standard RAG systems that require moving data to separate vector stores, our system can join live sales data with AI-driven sentiment analysis in a single, atomic operation. This results in extremely **low latency** and high data integrity.

5. Deployment Guide

Integrating the system is straightforward. We configure the AI client by pointing it to our Python script within a local configuration file. This setup allows the AI to "boot up" our server as a background process whenever it needs to interact with Oracle. By defining the database DSN and credentials within the environment variables of this configuration, we maintain a clean separation between the code and our sensitive access keys.

```
{
  "mcpServers": {
    "oracle_custom": {
      "command": "python",
      "args": ["C:/projects/oracle_mcp.py"],
      "env": {
        "ORACLE_HOME": "/opt/oracle",
        "DB_DSN": "localhost:1521/FREEPDB1"
      }
    }
  }
}
```

6. Why "Build" instead of "Buy"?

Choosing to build a custom server is a strategic decision driven by **cost and privacy**. Most cloud-based AI gateways charge per-user or per-query fees that can scale rapidly. By self-hosting this lightweight Python bridge, those costs are eliminated. More importantly, **Data Privacy** is guaranteed; since the MCP server lives on our local network or private cloud, our sensitive enterprise data never travels to a third-party middleman for "processing" before reaching the database. This provides the ultimate level of flexibility to sanitize or audit every single request the AI makes.

7. Conclusion and Final Remarks

The architecture established in this framework provides a high-performance, secure, and cost-effective solution for bridging enterprise data with Large Language Models. By utilizing the Model Context Protocol (MCP) as an intelligent interpreter, the system maintains a robust layer of isolation between the AI host and the raw database engine.

The shift from external middleware to a custom, self-hosted Python bridge offers several critical advantages:

- **Security and Privacy:** The Zero-Trust model ensures that sensitive connection strings and logic are encapsulated within the server, preventing the LLM from accessing raw storage or credentials.
- **Data Sovereignty:** Enterprise data remains within the private cloud or local network, eliminating the risks associated with third-party processing.
- **Architectural Efficiency:** The use of the python-oracledb driver in "Thin Mode" streamlines deployment by removing the requirement for bulky Instant Client installations.
- **Converged Access:** The system enables the seamless integration of live relational data with AI Vector Search, allowing for low-latency, atomic operations that maintain high data integrity.