# Deep Dive: Integrating AI/LLMs/MLOps

This document explores how your enterprise-ready data platform extends its capabilities to support Artificial Intelligence (AI) and Large Language Models (LLMs), all managed under the umbrella of robust Machine Learning Operations (MLOps) principles. While traditional ML focuses on structured data and predictive models, AI/LLMs introduce new dimensions related to unstructured text, natural language understanding, and generative capabilities. MLOps provides the essential operational framework to reliably develop, deploy, and monitor all types of AI and ML solutions.

## 1. Core Concepts & Platform Alignment for AI/LLMs/MLOps

Your existing data platform, designed for enterprise-grade data management, is inherently well-suited to integrate AI/LLM workloads.

- **Data Pipelines for AI/LLMs:**
  - **Ingestion (FastAPI, Kafka, Delta Lake):** Just as with structured data, raw text, documents, or conversation logs are ingested into the platform. This raw data forms the basis for Retrieval Augmented Generation (RAG) systems or fine-tuning datasets for LLMs.
  - **Feature Engineering (Apache Spark):** For traditional ML, Spark processes structured features. For LLMs, Spark can perform text preprocessing, chunking, embedding generation, or preparing structured prompts/responses for fine-tuning.
  - **Data Quality & Governance (OpenMetadata):** Critical for both traditional ML and LLMs. Ensuring the quality of text data used for RAG or fine-tuning, and maintaining lineage for prompt templates, models, and generated outputs.
- **Model Management (MinIO/Delta Lake, conceptual Model Registry):**
  - Trained traditional ML models are versioned in your data lakehouse (MinIO/Delta Lake).
  - LLMs, whether open-source models or proprietary ones accessed via API, also require versioning of their configurations, fine-tuned weights (if applicable), and prompt templates. MinIO/Delta Lake can store these assets.
- **Orchestration (Apache Airflow):**
  - Airflow remains the central orchestrator for complex pipelines:
    - Data preparation for model training/fine-tuning.
    - Automated training/fine-tuning jobs.
    - Batch inference jobs.
    - Model evaluation and validation workflows.
    - Deployment triggers for new model versions or prompt templates.
- **Observability (OpenTelemetry, Grafana Alloy, Grafana):**
  - Monitoring extends beyond system health to model performance. For LLMs, this

includes metrics like:
- **Latency:** Time to generate responses.
- **Token Usage:** Input/output token counts for cost management.
- **Quality Metrics:** Hallucination rates, relevance, coherence (often requires human-in-the-loop or proxy metrics).
- **Abuse/Safety:** Monitoring for inappropriate content generation.
    - Traces can follow LLM calls, and logs can include prompts and responses (with necessary PII masking).

# 2. Interactive How-Tos: Integrating AI/LLMs/MLOps

Let's explore some practical examples of how these integrations can take shape within your platform.

## Basic Use Case: Data Preparation for LLM Fine-tuning/RAG

**Objective:** To demonstrate how Apache Spark can preprocess unstructured text data (e.g., articles, customer support transcripts) from your raw data lake into a structured format suitable for LLM fine-tuning or Retrieval Augmented Generation (RAG) systems.

**Role in Platform:** Transform raw, unstructured data into a consumable format for LLMs, enabling the creation of custom knowledge bases or training datasets.

**Setup/Configuration (Local Environment - Advanced Track):**

1. **Ensure all Advanced Track services are running:** docker compose up --build -d. This includes spark and minio.
2. **Simulate Raw Text Data:** Create a dummy text file raw_docs.json that looks like ingested unstructured data.
    - In your data/minio/raw-data-bucket/ directory, create articles/raw_docs.json.

# data/minio/raw-data-bucket/articles/raw_docs.json
{"id": "doc1", "text": "Apache Spark is a unified analytics engine for large-scale data processing. It provides high-level APIs in Java, Scala, Python and R. Spark also supports a rich set of higher-level tools including Spark SQL for SQL and structured data processing, MLlib for machine learning, GraphX for graph processing, and Structured Streaming for incremental computation and stream processing."}
{"id": "doc2", "text": "Large Language Models (LLMs) are deep learning models trained on vast amounts of text data. They are capable of understanding, generating, and manipulating human language. Popular LLMs include GPT-3, LLaMA, and Gemini. Applications range from chatbots and content generation to code completion and translation."}
{"id": "doc3", "text": "Retrieval Augmented Generation (RAG) is an AI framework for improving the specificity and factual accuracy of generative AI models with information retrieved from external knowledge bases. RAG models combine a retrieval component with a generation component. Instead of generating responses solely based on their training data, RAG models first retrieve relevant documents or passages from a given knowledge base and then use this retrieved information to inform their generated response."}
{"id": "doc4", "text": "MLOps (Machine Learning Operations) is a set of practices that

aims to deploy and maintain machine learning models in production reliably and efficiently. The MLOps lifecycle includes data preparation, model training, deployment, monitoring, and governance. Key tools often include orchestrators like Airflow and observability platforms like Grafana."}

3. **Create a PySpark script for LLM data preparation:** In pyspark_jobs/, create llm_data_prep.py.

```python
# pyspark_jobs/llm_data_prep.py
import sys
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, lit, current_timestamp, split, explode, monotonically_increasing_id, concat_ws
from pyspark.sql.types import StructType, StringType, IntegerType

def create_spark_session(app_name):
    """Helper function to create a SparkSession with Delta Lake packages."""
    return (SparkSession.builder.appName(app_name)
        .config("spark.jars.packages", "io.delta:delta-core_2.12:2.4.0")
        .config("spark.sql.extensions", "io.delta.sql.DeltaSparkSessionExtension")
        .config("spark.sql.catalog.spark_catalog",
"org.apache.spark.sql.delta.catalog.DeltaCatalog")
        .getOrCreate())

# Simple text chunking function (for RAG)
# In a real scenario, this would be more sophisticated (e.g., using NLTK, spaCy, or
Hugging Face tokenizers)
# This UDF is illustrative. For large-scale text, pure Spark functions are preferred where
possible.
# from pyspark.sql.functions import udf
# @udf(returnType=ArrayType(StringType()))
# def chunk_text_udf(text: str, chunk_size: int = 1000):
#    if not text:
#        return []
#    words = text.split(" ")
#    chunks = []
#    current_chunk = []
#    for word in words:
#        current_chunk.append(word)
#        if len(" ".join(current_chunk)) >= chunk_size:
#            chunks.append(" ".join(current_chunk))
#            current_chunk = []
#    if current_chunk:
#        chunks.append(" ".join(current_chunk))
```

```python
    #     return chunks

if __name__ == "__main__":
    if len(sys.argv) != 3:
        print("Usage: llm_data_prep.py <input_raw_text_path> <output_delta_path>")
        sys.exit(-1)

    input_raw_text_path = sys.argv[1]
    output_delta_path = sys.argv[2]

    spark = create_spark_session("LLMDataPreparation")
    spark.sparkContext.setLogLevel("WARN")

    print(f"Reading raw text data from: {input_raw_text_path}")
    # Read raw JSON lines from MinIO/S3
    df_raw = spark.read.json(input_raw_text_path)
    df_raw.printSchema()
    df_raw.show(5, truncate=False)

    print("Performing text preprocessing and chunking (conceptual)...")
    # Example: Simple chunking by paragraph (splitting by double newline for illustration)
    # In a real application, you'd use a more robust text splitter, potentially preserving
context.
    df_chunked = df_raw.withColumn("paragraphs", split(col("text"), "\\n\\n")) \
                .withColumn("chunk", explode(col("paragraphs"))) \
                .filter(col("chunk") != "") # Remove empty chunks

    # Assign a unique ID to each chunk and add metadata
    df_final = df_chunked.withColumn("chunk_id", concat_ws("_", col("id"),
monotonically_increasing_id())) \
                .select(
                    col("chunk_id"),
                    col("id").alias("source_document_id"),
                    col("chunk").alias("text_content"),
                    current_timestamp().alias("processed_at")
                )

    print("Schema of prepared data for LLM:")
    df_final.printSchema()
    df_final.show(5, truncate=False)

    # Write the prepared data to a new Delta Lake table
    print(f"Writing prepared LLM data to: {output_delta_path}")
```

```
df_final.write.format("delta") \
        .mode("overwrite") \
        .option("overwriteSchema", "true") \
        .save(output_delta_path)
print("LLM data preparation job completed.")

spark.stop()
```

**Steps to Exercise:**
1.  **Place raw_docs.json:** Ensure the raw_docs.json file is in data/minio/raw-data-bucket/articles/.
2.  **Submit the Spark LLM Data Prep Job:**
    ```
    docker exec -it spark spark-submit \
        --packages io.delta:delta-core_2.12:2.4.0 \
        --conf spark.sql.extensions=io.delta.sql.DeltaSparkSessionExtension \
        --conf
    spark.sql.catalog.spark_catalog=org.apache.spark.sql.delta.catalog.DeltaCatalog \
        --conf spark.hadoop.fs.s3a.endpoint=http://minio:9000 \
        --conf spark.hadoop.fs.s3a.access.key=minioadmin \
        --conf spark.hadoop.fs.s3a.secret.key=minioadmin \
        --conf spark.hadoop.fs.s3a.path.style.access=true \
        /opt/bitnami/spark/jobs/llm_data_prep.py \
        s3a://raw-data-bucket/articles/raw_docs.json \
        s3a://curated-data-bucket/llm_prepared_docs
    ```

3.  **Monitor Spark Job:** Observe the console output for Spark logs, confirming data reading, processing, and writing.
4.  **Verify Prepared Data in MinIO:** Access the MinIO Console (http://localhost:9001). Navigate to curated-data-bucket/llm_prepared_docs/. You should see new Delta Lake files.
5.  **Query Prepared Data via Spark SQL:**
    ```
    docker exec -it spark spark-sql \
        --packages io.delta:delta-core_2.12:2.4.0 \
        --conf spark.sql.extensions=io.delta.sql.DeltaSparkSessionExtension \
        --conf
    spark.sql.catalog.spark_catalog=org.apache.spark.sql.delta.catalog.DeltaCatalog \
        --conf spark.hadoop.fs.s3a.endpoint=http://minio:9000 \
        --conf spark.hadoop.fs.s3a.access.key=minioadmin \
        --conf spark.hadoop.fs.s3a.secret.key=minioadmin \
        --conf spark.hadoop.fs.s3a.path.style.access=true \
        -e "SELECT * FROM delta.\`s3a://curated-data-bucket/llm_prepared_docs\` LIMIT 10;"
    ```

**Verification:**
● **Spark Job Completion:** The job runs successfully, indicated by console output.

- **MinIO Contents:** The llm_prepared_docs Delta Lake table is created with new .parquet files and _delta_log.
- **Spark SQL Query:** The query results show the text_content column, demonstrating that raw text has been transformed into a structured format with unique chunk_ids, ready for use in RAG or fine-tuning datasets.

# Advanced Use Case 1: Integrating LLM Inference into an Application (Conceptual)

**Objective:** To conceptually demonstrate how a FastAPI service could integrate with an external LLM API (e.g., Gemini API) to enrich incoming data or generate responses. This shows how your platform's API layer can become a gateway for AI capabilities.

**Role in Platform:** Enable real-time AI-powered features, such as intelligent routing of support tickets, sentiment analysis on customer feedback, or dynamic content generation.

**Setup/Configuration:**

1. **Ensure FastAPI is running:** Your fastapi_ingestor service should be active.
2. **Conceptual LLM API Key:** Assume an API key for the LLM service exists (e.g., in environment variables).
3. **Modify fastapi_app/app/main.py:** Add a new endpoint that calls a placeholder LLM API.

```
# fastapi_app/app/main.py (conceptual additions for LLM integration)
# ... existing imports
import httpx # For making async HTTP requests
# ... existing OpenTelemetry setup, Kafka producer, Pydantic models, etc.

# --- LLM Integration Configuration ---
LLM_API_URL = os.getenv("LLM_API_URL",
"https://generativelanguage.googleapis.com/v1beta/models/gemini-2.0-flash:generateC
ontent")
LLM_API_KEY = os.getenv("LLM_API_KEY", "") # IMPORTANT: In production, use AWS
Secrets Manager or similar!

# --- New Pydantic Models for LLM Interaction ---
class LLMQueryRequest(BaseModel):
    text_input: str = Field(..., example="Explain the concept of data lineage in simple
terms.")
    context: Optional[str] = Field(None, example="Consider a financial data pipeline.")

class LLMResponse(BaseModel):
    original_input: str
    generated_text: str
    tokens_used: Optional[int] = None
    processing_time_ms: Optional[float] = None
```

```python
        error: Optional[str] = None


# --- New FastAPI Endpoint for LLM Interaction ---
@app.post("/llm-query/", response_model=LLMResponse, tags=["AI/LLM"])
async def llm_query(request: LLMQueryRequest):
    start_time = time.time()
    prompt_text = request.text_input
    if request.context:
        prompt_text = f"Context: {request.context}\n\nQuery: {request.text_input}"

    try:
        # Prepare chat history for the LLM API call
        chatHistory = []
        chatHistory.push({ "role": "user", "parts": [{ "text": prompt_text }] })

        # Prepare the payload for the LLM API
        payload = { "contents": chatHistory }
        # If you want structured response, add generationConfig
        # payload["generationConfig"] = {
        #     "responseMimeType": "application/json",
        #     "responseSchema": {
        #         "type": "OBJECT",
        #         "properties": {
        #             "summary": { "type": "STRING" },
        #             "keywords": { "type": "ARRAY", "items": { "type": "STRING" } }
        #         }
        #     }
        # }

        # Make the async HTTP request to the LLM API
        async with httpx.AsyncClient() as client:
            response = await client.post(
                f"{LLM_API_URL}?key={LLM_API_KEY}",
                json=payload,
                timeout=60 # Adjust timeout as needed for LLM response
            )
            response.raise_for_status() # Raise an exception for bad status codes

            llm_result = response.json()
            generated_text = "No response from LLM."
            tokens_used = None # LLM API might provide token usage
```

```python
        if llm_result.get("candidates") and len(llm_result["candidates"]) > 0:
            first_candidate = llm_result["candidates"][0]
            if first_candidate.get("content") and first_candidate["content"].get("parts"):
                # Extract the text from the LLM's response
                generated_text = first_candidate["content"]["parts"][0]["text"]

                # Conceptual token usage (replace with actual API response parsing)
                # if llm_result.get("usage_metadata"):
                #     tokens_used = llm_result["usage_metadata"].get("total_token_count")

        processing_time_ms = (time.time() - start_time) * 1000

        # Log the LLM interaction (with PII masking if necessary)
        print(f"LLM Query: '{request.text_input[:50]}...', Response:
'{generated_text[:50]}...', Time: {processing_time_ms:.2f}ms")

        return LLMResponse(
            original_input=request.text_input,
            generated_text=generated_text,
            tokens_used=tokens_used,
            processing_time_ms=processing_time_ms
        )

    except httpx.HTTPStatusError as e:
        error_msg = f"LLM API HTTP Error: {e.response.status_code} - {e.response.text}"
        print(error_msg)
        raise HTTPException(status_code=e.response.status_code, detail=error_msg)
    except httpx.RequestError as e:
        error_msg = f"LLM API Request Error: {e}"
        print(error_msg)
        raise HTTPException(status_code=500, detail=error_msg)
    except json.JSONDecodeError:
        error_msg = "Invalid JSON response from LLM API."
        print(error_msg)
        raise HTTPException(status_code=500, detail=error_msg)
    except Exception as e:
        error_msg = f"An unexpected error during LLM query: {e}"
        print(error_msg)
        raise HTTPException(status_code=500, detail=error_msg)
```

*Note:* You'll need to add httpx to fastapi_app/requirements.txt and rebuild the FastAPI image.

**Steps to Exercise:**

1. **Add httpx to fastapi_app/requirements.txt:**
   # fastapi_app/requirements.txt
   ...
   httpx

2. **Rebuild FastAPI image and restart container:**
   docker compose build fastapi_ingestor
   docker compose restart fastapi_ingestor

3. **Send a query to the new LLM endpoint:**
   curl -X POST -H "Content-Type: application/json" \
       -d '{
           "text_input": "What is the capital of France?",
           "context": "Answer briefly."
         }' \
       http://localhost:8000/llm-query/

   If you have a real Gemini API key, ensure LLM_API_KEY is set as an environment variable for the fastapi_ingestor service in docker-compose.yml.
   If not, the request will likely fail, but the internal logic of parsing and attempting the call will be demonstrated.

**Verification:**
- **HTTP Response:** You will receive a JSON response from the FastAPI endpoint. If the LLM API call was successful (e.g., if you provided a valid API key), generated_text will contain the LLM's response. If not, it will show an error message.
- **FastAPI Logs:** The fastapi_ingestor logs will show the "LLM Query..." message, confirming the endpoint was hit and the conceptual interaction occurred.

# Advanced Use Case 2: MLOps for LLMs - Monitoring & Evaluation (Conceptual)

**Objective:** To discuss and conceptually outline how LLM-specific metrics and evaluation results can be collected using OpenTelemetry and visualized in Grafana, crucial for understanding LLM performance and cost in production.

**Role in Platform:** Extend observability to cover LLM-specific KPIs, enabling proactive monitoring for issues like high latency, excessive token usage, or potential "hallucinations."

**Setup/Configuration:**
1. **Ensure OpenTelemetry is configured for FastAPI:** Your fastapi_app/app/main.py should have the OpenTelemetry setup (as detailed in "Highlighting OpenTelemetry" document), with OTLPMetricExporter sending to grafana-alloy.
2. **Ensure Grafana Alloy and Grafana are running.**

**Steps to Exercise (Conceptual Discussion):**

1. **Define Custom LLM Metrics (within main.py where LLM call is made):**
   - **Latency:** llm.request.duration_ms (Histogram with attributes like model_name, endpoint, success).
   - **Token Usage:** llm.input.tokens_total, llm.output.tokens_total (Counters with attributes like model_name).
   - **Cost:** llm.estimated.cost_usd (Counter, if you can estimate cost per token).
   - **Error Rate:** llm.api.errors_total (Counter with error_type, model_name).
   - **Conceptual Quality Metrics (Proxy):** llm.response.length_chars, llm.response.keywords_count (Gauges, or derived metrics from text post-processing).

```python
# fastapi_app/app/main.py (conceptual additions for LLM metrics)
# ... existing OpenTelemetry setup ...
# From Advanced Use Case 1 in OpenTelemetry:
# meter = metrics.get_meter("fastapi.ingestion.app")

# Create LLM-specific metrics
llm_request_duration = meter.create_histogram(
    "llm.request.duration",
    description="Duration of LLM API requests",
    unit="ms",
    boundaries=[50, 100, 200, 500, 1000, 2000, 5000, 10000] # Adjust for expected
LLM latencies
)
llm_input_tokens = meter.create_counter(
    "llm.input.tokens_total",
    description="Total input tokens sent to LLM",
    unit="1"
)
llm_output_tokens = meter.create_counter(
    "llm.output.tokens_total",
    description="Total output tokens received from LLM",
    unit="1"
)
llm_error_count = meter.create_counter(
    "llm.api.errors_total",
    description="Total errors from LLM API calls",
    unit="1"
)

# --- Inside @app.post("/llm-query/") endpoint ---
# ...
    # After successful LLM call:
    llm_request_duration.record(processing_time_ms, {"model_name":
"gemini-2.0-flash", "success": True})
    # if actual_input_tokens:
    #    llm_input_tokens.add(actual_input_tokens, {"model_name":
```

```
"gemini-2.0-flash"})
      # if actual_output_tokens:
      #     llm_output_tokens.add(actual_output_tokens, {"model_name":
"gemini-2.0-flash"})
# ...
# In error handling block (e.g., in except):
#     llm_error_count.add(1, {"error_type": "http_error" if isinstance(e,
httpx.HTTPStatusError) else "other", "model_name": "gemini-2.0-flash"})
#     llm_request_duration.record(processing_time_ms, {"model_name":
"gemini-2.0-flash", "success": False, "error_type": "timeout"}) # Example
```

2. **Visualize in Grafana:**
   ○ After adding metrics and generating traffic, access http://localhost:3000.
   ○ Create a new dashboard or panel.
   ○ **PromQL Queries:**
     ■ LLM Request Rate:
       rate(llm_request_duration_count{model_name="gemini-2.0-flash"}[1m])
     ■ LLM P99 Latency: histogram_quantile(0.99, sum by(le, model_name)
       (rate(llm_request_duration_bucket[1m])))
     ■ Total Input/Output Tokens:
       llm_input_tokens_total{model_name="gemini-2.0-flash"} (and for output)
     ■ Error Rate: rate(llm_api_errors_total{model_name="gemini-2.0-flash"}[5m])
3. **Discuss LLM Evaluation & Monitoring for Drift:**
   ○ **Data Drift:** Use OpenMetadata's data profiling on the llm_prepared_docs Delta
     table. Monitor for changes in text length distribution, vocabulary, or topic
     distribution of incoming text data.
   ○ **Model Drift (LLM Output):** This is harder.
     ■ **Proxy Metrics:** Monitor metrics derived from LLM output (e.g.,
       llm.response.length_chars). Sudden changes could indicate drift.
     ■ **Human Evaluation:** For critical LLM applications, a sample of responses
       might require periodic human review.
     ■ **Embeddings:** Generate embeddings for LLM inputs/outputs and monitor
       the distribution of these embeddings over time (e.g., using dimensionality
       reduction and clustering in Grafana for anomaly detection).
   ○ **Logging Prompts/Responses:** Crucial for post-hoc analysis and debugging.
     Ensure sensitive data is masked. OpenTelemetry LoggingInstrumentor can help
     link these logs to traces.

**Verification (Conceptual):**
● **Grafana Dashboards:** Successfully visualize LLM-specific operational metrics like
  request rate, latency, and token usage.
● **Architectural Understanding:** Demonstrate a clear understanding of how to set up
  monitoring for LLM applications, extending beyond traditional infrastructure metrics to
  capture critical AI-specific KPIs.

# Advanced Use Case 3: Versioning and Deploying Models/Prompts

# (Conceptual MLOps Workflow)

**Objective:** To illustrate how MLOps principles (versioning, automated pipelines, controlled deployments) apply to both traditional ML models and LLM prompts/configurations within your platform.

**Role in Platform:** Ensure reproducibility, auditability, and reliable deployment of all AI/ML artifacts, from data to models to prompts.

**Setup/Configuration (Conceptual Discussion):**

1. **Existing CI/CD Pipeline:** Refer to .github/workflows/release.yml for the overall CI/CD structure.
2. **MinIO/Delta Lake for Artifact Storage:** MinIO acts as the artifact repository.
3. **Airflow for Orchestration:** Refer to airflow_dags/ for pipeline orchestration.

**Steps to Exercise (Conceptual/Discussion):**

1. **Model Versioning in Data Lakehouse:**
   - **ML Models:** After training (model_training_job.py), the model is saved to a specific versioned path in MinIO/Delta Lake (e.g., s3a://models-bucket/financial_fraud_model/v1.0.0/).
   - **LLM Prompts/Configs:** Store prompt templates, RAG knowledge base configurations, or LLM fine-tuning datasets in versioned paths within MinIO/Delta Lake as well (e.g., s3a://llm-artifacts-bucket/prompts/sentiment_analysis_v1.0/prompt_template.txt). This ensures traceability of the exact prompt used with a given model.
   - **OpenMetadata Role:** OpenMetadata can catalog these model and prompt artifacts, linking them to their source data and lineage.

2. **Automated Training/Fine-tuning Pipelines (Airflow):**
   - **Airflow DAGs:** Create Airflow DAGs that:
     - Trigger feature_engineering_job.py.
     - Trigger model_training_job.py (or a fine-tuning job for LLMs).
     - Run model evaluation (e.g., on a test dataset) and store metrics/artifacts.
     - Conditionally push the new model/prompt version to the artifact store.
     - *Conceptually:* Trigger a deployment pipeline if metrics are satisfactory.
   - **Example Trigger:** An Airflow DAG could be scheduled daily, or triggered by a data quality anomaly on features.

3. **Deployment of Inference Services:**
   - **CI/CD:** The existing CI/CD (GitHub Actions) can be extended.
     - **Traditional ML:** The FastAPI inference service (or a dedicated Spark Streaming inference job) would be deployed/updated to point to the *new model version* in MinIO. This could involve updating an environment variable for the FastAPI service or a Spark job parameter.
     - **LLM Integration:** If the FastAPI service directly integrates with an LLM API, prompt changes might be deployed by updating the FastAPI application code itself (e.g., new prompt templates in a config file). Or, if prompts are externalized to a data store, the deployment ensures the service picks up

the latest version.
- ○ **Blue/Green or Canary Deployments:** In a production cloud environment, you'd apply blue/green or canary deployments for inference services to minimize downtime and risk during model/prompt rollouts.
4. **Rollback Strategy:**
   - ○ Because models/prompts are versioned in MinIO/Delta Lake, a rollback is as simple as re-deploying the inference service to point to a *previous, known-good version* of the model or prompt. This makes ML deployments robust.

**Verification (Conceptual):**
- ● **MLOps Lifecycle Understanding:** Demonstrates a comprehensive understanding of how data platform components support a robust MLOps lifecycle, including versioning of models and prompts, automated pipelines, and controlled deployments, which are essential for reliable AI/ML initiatives in an enterprise.

This concludes the deep dive into Integrating AI/LLMs/MLOps within your data platform.