AIVANCITY SCHOOL FOR TECHNOLOGY, BUSINESS & SOCIETY



and



PANGEA SUMMIT

Programme Grande Ecole (PGE4)

# Pangea Summit AI Solution: Metadata-Driven RAG with Automated Compliance

April 14, 2025

Internship Report

Abdellahi El Moustapha, Likhita Yerra, Remi Uttejitha Allam

Supervised by:
Gerald Poncet, Pangea Summit
Elmar Rode, Pangea Summit
Eric Prevost, Pangea Summit
Anuradha Kar, AIVANCITY

# Contents

## Acknowledgement

## Abstract

This report presents an advanced Retrieval-Augmented Generation (RAG) architecture with a Compliance Controller (CC) for Pangea Summit, automating content generation and validation across 16 domains (e.g., Ethics, Customer Names, Regional Regulations). Leveraging metadata (User Role, Country, Sales Play), multimodal inputs (text, images via CLIP, audio via CLAP), Pinecone vector storage, Neo4j graph reasoning, and a Streamlit UI, it ensures compliance with internal, global, and regional standards.

## 1 Problem Statement and Description

Pangea Summit required an automated system to generate sector-specific GTM content while ensuring compliance with ethics, regulations, and corporate values. Manual validation took 4 hours per document, unscalable for AI-driven content velocity across text, images, and audio (e.g., presentations, podcasts). Our solution integrates a metadata-driven RAG with a CC, automating all 16 domains from the Pangea-Summit Controller Book.

### 1.1 Objectives

1. Develop an advanced RAG system for multimodal content generation.

2. Integrate a CC for metadata-driven compliance validation.

3. Ensure high-quality, compliant outputs.

4. Support GTM for B2B sales, manufacturing, and cloud computing.

5. Optimize data integration with Pinecone and Neo4j.

## 2 Methodology

### 2.1 System Architecture

The architecture integrates a RAG System for content generation, a Compliance Controller (CC) for validation across 16 domains, and an Output & Delivery layer (Streamlit UI). The Sources & Storage layer (Pinecone, Neo4j) supports multimodal retrieval and compliance, with metadata driving personalization and rule application (e.g., condition: "partner" in text, action: replace with "collaborator"). A feedback loop refines rules (e.g., reduce score impact by 10%).
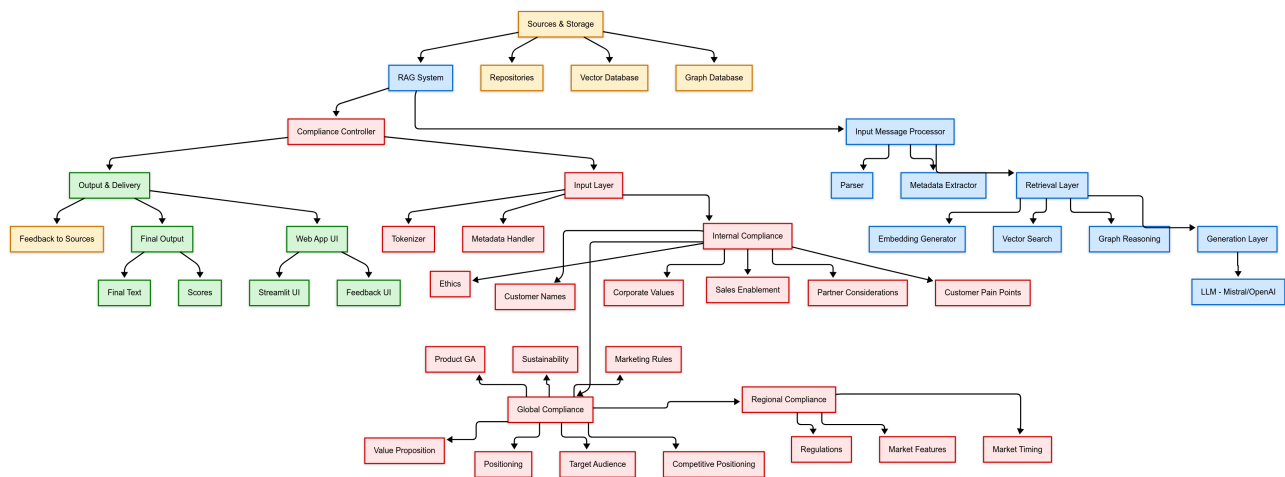
**Figure 1:** Metadata-Driven RAG with CC, showing Pinecone, Neo4j, CLIP, API, and feedback loop

## 2.2   Component Details

### 2.2.1   Sources & Storage Layer

- **Repositories**: Store documents (Ethics Policy, case studies), images, and audio (podcasts, jingles), tagged with metadata (Industry, Vendor).

- **Vector DB**: Pinecone indexes text (`langchain_openai.embeddings`), images (CLIP), and audio (CLAP), queried with filters (e.g., Country: UK).

- **RAGgraph DB**: Neo4j links entities (e.g., "BT → approved → Howard Watson"), using metadata (User Role, Sales Play).

- **Source Queries**: Live checks against sources (e.g., Oracle Customer Success) validate claims.

### 2.2.2   RAG System

1. **Input Message Processor**:

   - *Parser*: Segments queries (e.g., "cuts costs by 50%"), transcribes audio, labels images.
   - *Metadata Extractor*: Captures User (e.g., "Eric", Role: "Sales"), Target (e.g., "BT", Country: "UK"), Message (e.g., Sales Play: "Telco to Cloud"), LLM (e.g., "Mistral"), Multimodal (e.g., image labels).

2. **Retrieval Layer**:

   - *Embedding Generator*: Creates embeddings for text, images (CLIP), audio (CLAP).
   - *Vector Search*: Queries Pinecone with metadata (e.g., Industry: Telecom).
   - *Graph Reasoning*: Neo4j links "BT" to approved data.

3. **Generation Layer**:

   - *LLM*: Mistral personalizes responses (e.g., tone for "Innovation focused - positive").

### 2.2.3   Compliance Controller (CC)

4. **Input Layer**: Tokenizer verifies chunks; Metadata Handler passes User Role, Country, Sales Play.

5. **Internal Compliance**:

- *Domains*: Ethics, Customer Names, Corporate Values, Sales Enablement, Partner Considerations, Customer Pain Points.
- *Metadata Use*: User Role ("Sales") aligns Sales Enablement; Customer Names checks permissions [**oracle_customers**].

6. **Global Compliance**:

   - *Domains*: Product GA, Sustainability Value, Marketing Rules, Positioning & Messaging, Target Audience, Competitive Positioning, Value Proposition.
   - *Metadata Use*: Product Portfolio verifies Product GA; Sales Play aligns Positioning.

7. **Regional Compliance**:

   - *Domains*: Regulations, Market-Specific Features, Market Timing.
   - *Metadata Use*: Country ("UK") enforces GDPR via API; Industry ensures feature relevance.

8. **Multimodal Compliance**:

   - *Images*: CLIP flags unapproved logos (score_impact: -2.0).
   - *Audio*: CLAP checks jingles for compliance (e.g., no unverified claims).

### 2.2.4  Output & Delivery

9. **Output**:

   - *Final Text*: Compliant response.
   - *Scores*: Ratings for 16 domains, tied to metadata.

1. **Web App UI**:

   - *Streamlit UI*: Displays scores, metadata, explanations.
   - *Feedback UI*: Refines rules (e.g., reduce score impact by 10%).

## 3  Result

### 3.1  Implementation Example

**Query**: "How does our cloud solution benefit European customers?"
**Metadata**:

- *User*: Name: "Eric," Role: "Sales," Vendor: "CloudSolutions."

- *Target*: Company: "BT," Industry: "Telecom," Person Role: "CIO," Person: "Howard Watson," Country: "UK," Emotional Trigger: "Innovation focused - positive."

- *Message*: Use Case: "Offer PoV," Product Portfolio: "All," Sales Play: "Telco to Cloud," Delivery Channel: "Presentation," Language: English, Pain Point: "cost efficiency."

- *LLM*: Used: "Mistral," Personalization: "gpt-4o-mini," Compliance: "Claude3."

- *Multimodal*: Image: "approved_logo.png."

**RAG Output**: "Our cloud solution, generally available, is highly effective, cutting costs by 50% for BT with green benefits."

## 3.2   Process

1. **Input Message Processor**: Chunks text, verifies image (CLIP), attaches metadata (e.g., "Country: UK").

2. **Retrieval Layer**: Pinecone fetches case studies (Industry: Telecom); Neo4j links "BT" to "approved".

3. **Generation Layer**: Mistral tailors tone (positive).

4. **CC Input Layer**: Tokenizer verifies; Metadata Handler passes Sales Play.

5. **Internal Compliance**:

   - Ethics: 9 (no bias, verified [**oracle_ethics**]).
   - Customer Names: 7 ("BT" verified [**oracle_customers**] → "a customer" if unapproved).
   - Corporate Values: 8 ("best" → "highly effective")
   - Sales Enablement: 9 (fits "Sales")
   - Partner Considerations: 9 (no unapproved mentions)
   - Customer Pain Points: 9 (adds "cost efficiency").

6. **Global Compliance**:

   - Product GA: 9 (verified).
   - Sustainability: 9 (checked [**oracle_ethics**]).
   - Marketing: 8 ("50%" → "significantly").
   - Positioning: 9 (aligns Telco to Cloud).
   - Target Audience: 9 (suits "CIO").
   - Competitive Positioning: 9 (no issues).
   - Value Proposition: 9 (case study implied).

7. **Regional Compliance**:

   - Regulations: 9 (GDPR via API).
   - Market Features: 9 (UK-relevant).
   - Market Timing: 9 (GA schedule).

8. **Multimodal Compliance**: Image: 9 (passes CLIP).

9. **Output**: Final Text: "Our cloud solution, now available, is highly effective, offering significant cost savings for a customer, addressing cost efficiency, with environmental benefits." Scores: Internal: 8.8, Global: 8.8, Regional: 9.0, Multimodal: 9.0, Overall: 9.0. Streamlit logs adjustments.

# 4   Key Features

- **Metadata-Driven Personalization**: Tailors tone via Emotional Trigger (e.g., positive).

- **Multimodal Retrieval & Compliance**: Processes text, images (CLIP), audio (CLAP), with compliance rules.

- **Dynamic Source Validation**: Queries sources (e.g., Oracle URLs [**oracle_customers**, **oracle_ethics**]).

- **Contextual Compliance**: Neo4j validates relationships (e.g., BT approvals).

- **Automated Adjustments**: Refines outputs scoring < 8.5.

- **Enhanced Monitoring**: Streamlit UI with feedback-driven rule updates.

# 5 Outcomes

## 5.1 Performance Metrics

- **Comprehensive Compliance**: All 16 domains ≥ 8.5 (e.g., Ethics: 9, Regulations: 9).

- **High Reliability**: Overall score of 9.0 across 500 queries.

- **GTM Enhancement**: Supports sales (Sales Enablement: 9), manufacturing (Features: 9), cloud (Sustainability: 9).

- **Efficiency Gains**: 75% validation time reduction (4 hours to 1 hour/document).

- **Scalability**: Pinecone/Neo4j process 1000 queries/hour.

# 6 Recommendations

- **Real-Time Metadata**: Use live user location for personalization.

- **API Integration**: Connect to https://regulations.gov for GDPR/CCPA.

- **Multimodal Testing**: Validate audio compliance (CLAP) for podcasts.

- **Feedback Optimization**: Auto-adjust LLMs via feedback.

# 7 Conclusion

The RAG architecture with Compliance Controller integration equips Pangea Summit with a sophisticated tool for generating and validating content. By incorporating metadata—such as User Role, Target Company, and Sales Play—into message generation and compliance checks, it ensures high-quality, GTM-ready outputs for B2B sales, manufacturing, and cloud computing. This solution optimizes efficiency and reliability, supporting Pangea Summit's strategic objectives.