# Enhancing AI Information Retrieval with NVIDIA NIM and DataGemma Models: A Deep Dive into Agentic-RIG

## Introduction

In today's data-driven world, the ability of AI systems to retrieve and generate accurate, context-aware information is more critical than ever. **Agentic-RIG** represents a significant advancement in this domain by enhancing traditional Retrieval-Augmented Generation (RAG) systems with a novel approach called **Retrieval Interleaved Generation (RIG)**. Leveraging the power of **NVIDIA Neural Information Model (NIM)** and **DataGemma models**, and integrated within the **NVIDIA AI Workbench**, Agentic-RIG pushes the boundaries of AI-powered information processing. Through a customizable Gradio Chat application, it offers dynamic, precise, and contextually relevant responses to users.

## The Evolution from RAG to RIG

### Limitations of RAG

Retrieval-Augmented Generation (RAG) combines information retrieval with AI-generated responses. While it marked a significant step forward, RAG has inherent limitations:

- **Single Retrieval Step**: RAG performs a one-time retrieval before the generation process, which may not capture all necessary information, especially for complex or evolving queries.

- **Context Misalignment**: As the response unfolds, the initially retrieved information may become less relevant, leading to gaps or inaccuracies.

- **Difficulty with Complex Queries**: Multi-faceted questions requiring diverse information sources can overwhelm RAG systems, resulting in incomplete or unsatisfactory answers.

### Introducing Retrieval Interleaved Generation (RIG)

**Retrieval Interleaved Generation (RIG)** addresses these challenges by seamlessly integrating retrieval and generation processes.

### Key Features of RIG:

- **Dynamic Retrieval**: Performs multiple retrievals at different stages of response generation, ensuring that new information needs are met as they arise.

- **Contextual Alignment**: Continuously updates the context based on both the user's input and the evolving response, maintaining relevance throughout.

- **Enhanced Accuracy**: By filling knowledge gaps on the fly, RIG reduces omissions and inaccuracies, providing comprehensive and precise responses.

# Model Flexibility: DataGemma and NVIDIA NIM

## DataGemma Models

**DataGemma** is an advanced large language model that embodies the principles of Retrieval Interleaved Generation. It offers:

- **Real-Time Data Integration:** Dynamically accesses and incorporates up-to-date information during response generation.

- **Reduced Hallucinations:** Grounds responses in retrieved data, minimizing the risk of generating incorrect information.

- **Improved Factual Accuracy:** Cross-references and verifies information, leading to highly accurate outputs.

Agentic-RIG can leverage DataGemma models in two ways:

- **Local Deployment with GPUs:** Ideal for high-performance environments with available GPU resources.

- **API Endpoints:** For users without dedicated hardware, DataGemma can be accessed via APIs, such as those provided by Hugging Face.

## Leveraging NVIDIA NIM

The **NVIDIA Neural Information Model (NIM)** offers a robust foundation for implementing RIG due to its advanced capabilities:

- **High Performance**: Optimized for NVIDIA GPUs, NIM provides rapid processing speeds essential for real-time applications.

- **Scalability**: Designed to handle extensive workloads, suitable for both small-scale and enterprise-level deployments.

- **Integration**: Seamlessly fits into NVIDIA AI ecosystems, benefiting from NVIDIA's development tools and frameworks.

## Model Flexibility and Integration with AI Workbench

Agentic-RIG is developed as an **NVIDIA AI Workbench** project, providing flexibility and ease of deployment:

- **Local GPU Deployment**: Users can run Agentic-RIG on their hardware with NVIDIA GPUs, utilizing both DataGemma and NVIDIA NIM models.

- **Cloud-Based APIs**: Integrates with cloud services and APIs, such as Hugging Face and NVIDIA's API Catalog, allowing access to models without dedicated hardware.

# System Overview of Agentic-RIG

### Agentic Workflow and Intelligent Query Routing

Agentic-RIG employs an intelligent routing mechanism where an AI model evaluates incoming queries to determine the most appropriate processing pipeline:

- **RIG Pipeline with DataGemma**: For queries that benefit from dynamic retrieval interleaved with generation, leveraging DataGemma's capabilities.

- **RAG Pipeline with NVIDIA NIM**: For queries suited to traditional retrieval-augmented generation, utilizing NVIDIA NIM for high-performance processing.

- **Web Search Pipeline**: For broader queries requiring up-to-date information from the internet, integrating web search APIs.

This routing ensures that each query is handled by the most suitable process, enhancing response quality and relevance.

### Interactive Interface with Gradio Chat App

Agentic-RIG features a user-friendly interface built with **Gradio**, enabling users to interact with the AI system seamlessly:

- **Real-Time Interaction**: Users can input queries and receive responses in an intuitive chat format.

- **Visual Workflow Indicator**: A diagram displays the agentic workflow, providing visual feedback on the processing stages of each query.

### User Configurable Settings and Flexibility

#### Model Settings

Users can customize the AI models used for various components of the pipeline:

- **Component Selection**: Configure models for the router, generator, and graders, choosing between DataGemma, NVIDIA NIM, or other models.

- **Custom Prompts**: Adjust the prompts for each component to tailor the AI's behavior and focus.

#### Document Settings

Agentic-RIG allows users to upload their own documents, enhancing personalization:

- **Document Upload**: Embed webpages or PDFs into a locally running **Chroma vector database**.

- **Contextual Relevance**: The system retrieves relevant information from these documents during response generation.

**Monitoring Tools**

Built-in monitoring features provide transparency into the AI's decision-making process:

- **Action Console**: Logs the agent's actions when processing queries.

- **Detailed Trace**: Offers in-depth insights into the retrieval and generation steps for each response.

# Integration with Third-Party Services

Agentic-RIG enhances its capabilities by integrating with various third-party services:

- **Tavily Search API**: Provides web search functionalities for queries requiring external information.

- **Hugging Face**: Accesses pre-trained models like DataGemma and inference endpoints.

- **OpenAI and LangSmith**: Utilizes additional language processing and analytics tools.

These integrations expand the system's functionality and adaptability, allowing users to leverage a wide range of AI resources.

# Real-World Application: Medical Research Assistant

### Scenario

A team of medical researchers is investigating potential treatments for a rare disease. They require up-to-date information from medical journals, clinical trials, and their own research documents.

### Agentic-RIG in Action

1 **User Interaction**

- Researchers upload their documents into the system and configure the model settings to focus on medical topics.

2 **Dynamic Retrieval and Generation**

- **RIG Pipeline with DataGemma**: For complex queries requiring dynamic retrieval, DataGemma models provide detailed, context-aware responses.

- **Web Search Pipeline**: For the latest information, the system performs web searches to gather recent studies and trial results.

3 **Interleaved Generation**

- The system interleaves retrieval and generation, incorporating new data as needed to provide comprehensive insights.

4 **Monitoring and Feedback**

- Researchers utilize monitoring tools to understand how the AI processed their queries and refine the system's settings accordingly.

**Impact**

- **Accelerated Research**: Provides timely and relevant information, enabling researchers to make informed decisions faster.

- **Enhanced Accuracy**: Reduces the risk of overlooking critical data by dynamically retrieving and integrating information.

# Challenges and Solutions

### Model Integration Flexibility

**Challenge**: Ensuring seamless operation across different models like DataGemma and NVIDIA NIM.

**Solution**:

- **Modular Architecture**: Developed an architecture that allows easy switching between models based on availability and suitability.

- **Dynamic Adjustment**: The system detects available resources (e.g., GPUs) and selects the optimal model accordingly.

### Harmonizing Retrieval and Generation

**Challenge**: Ensuring that retrieval and generation processes work seamlessly without introducing latency.

**Solution**:

- **Efficient Orchestration**: Implemented coordination mechanisms that optimally schedule tasks.

- **Asynchronous Processing**: Allows retrieval and generation to proceed without unnecessary delays.

## Conclusion

**Agentic-RIG**, powered by **DataGemma models** and **NVIDIA NIM**, and integrated within the **NVIDIA AI Workbench**, represents a transformative step in AI information retrieval and generation. By interleaving retrieval with generation and providing a customizable, user-friendly interface, it overcomes the limitations of traditional RAG systems, offering dynamic, accurate, and context-aware responses.

This project demonstrates how advanced AI techniques can be made accessible and practical, allowing users to tailor the system to their specific needs and enhance their productivity.

## Call to Action

We invite researchers, developers, and organizations to explore **Agentic-RIG** and collaborate with us to unlock new possibilities in AI-powered information systems. Together, we can drive innovation and shape the future of intelligent information retrieval.

*Note: This blog post is intended for the HACKAI Challenge and showcases our project utilizing NVIDIA technologies and DataGemma models.*