

Generative AI and Retrieval-Augmented Generation (RAG)

Introduction:

This seminar examined two game-changing advancements in artificial intelligence - Generative AI and Retrieval-Augmented Generation (RAG) - that are changing the possibilities and reliability of intelligent systems. The seminar is aimed at final year BCA students and will provide them with conceptual understanding and practical awareness of how modern artificial intelligence systems generate, augment and verify information.

Overview of Generative AI:

Generative AI has emerged as a significant iteration of classic AI. While traditional AI systems classify or predict events, Generative AI models generate new data such as text, images, computer code, and music after learning complex relationships from richer and larger datasets. These models use deep learning techniques and transformer-based architectures and are the fundamental frameworks of applications such as ChatGPT, DALL·E, GitHub Copilot and Gemini. The seminar will unpack the underlying mechanics of token prediction, recalling contexts, and training models, to enable examples of these tools to generate human-like content.

Limitations of Generative AI:

Hallucinations: Producing erroneous or fictitious information

Lack of Real-Time Knowledge: Models rely on previously acquired, pre-trained knowledge

Bias and Inconsistency: Output potentially reflects limitations of the dataset.

Lack of Transparency: Difficulty extricating the use of generated knowledge.

These limitations create a need for a framework that can anchor AI generation in factually anchored, up-to-date information.

Introducing Retrieval-augmented Generation (RAG):

Retrieval-augmented Generation (RAG) is a powerful solution to the limitations of generative AI. It combines the information retrieval capabilities of traditional search systems and generative capabilities of large language models (LLMs). RAG was designed to improve the factuality of AI-generated outputs by retrieving evidence-based information from outside knowledge sources, or documents, or broadly the internet, prior to generating a final response. This helps ensure that generative AI-generated outputs are coherent and grounded in verifiable evidence.

RAG Architecture and Workflow:

The outline of the seminar is the main architecture and workflow of RAG. There are two main components:

- Retriever: finds and retrieves relevant information from external knowledge and vector databases
- Generator: takes the retrieved context and produces accurate, meaningful, and contextually relevant outputs.

Workflow of RAG:

User Query → Retriever (Fetch Data) → Generator (LLM Output) → Factually Grounded Output

This architecture combines creativity and knowledge to produce outputs that leverage AI's strengths in being expressive and trustworthy.

Advantages of using RAG:

- Reduced factual errors and hallucinations
- Provides explainable and traceable outputs
- Real-time updates can be made without retraining the model
- Allows the integration of domain-specific or enterprise databases
- Increases user trust and reliability in AI-generated outputs

Examples of Real-World Applications:

The seminar discusses real-world applications of RAG across multiple domains, such as

- **ChatGPT** (with browsing and custom GPTs) - dynamic retrieval to produce information in real-time
- **Microsoft Copilot** - integrates enterprise data within productivity tools
- **Google Gemini** - has incorporated multimodal inputs while "grounding" its output to those inputs
- **Enterprise AI Chatbots** - can provide accurate responses based on private datasets.

Future Prospects:

The session concludes with an exploration of **emerging trends and innovations**, including:

- **Multimodal RAG** – integrating text, images, and audio retrieval.
- **Knowledge Graph Integration** – enhancing contextual understanding.
- **Privacy-Aware RAG** – ensuring secure access to private data.
- **Autonomous Agentic AI Systems** – enabling reasoning and self-directed information retrieval.

These developments signify a shift toward **intelligent, context-aware, and reliable AI ecosystems**.

Conclusion:

The seminar emphasizes that the future of AI lies in the **fusion of creativity and accuracy**. While Generative AI empowers machines to think creatively, RAG ensures those creations are **factually grounded and contextually relevant**. Together, they mark the next major evolution in artificial intelligence—building systems that are not only intelligent but also **trustworthy, transparent, and aligned with human knowledge**.

