# Retrieval Augmented Generation (RAG)

RAG is a technique that enhances the capabilities of Large Language Models (LLMs) by grounding their responses in external, contextually relevant data. In its basic form, RAG involves these core steps:

1. **Retrieval:**
   - A user's query is used to retrieve relevant documents or chunks of text from an external knowledge base (e.g., a vector database).
   - This retrieval is typically based on semantic similarity, using embeddings to find contextually related information.
2. **Augmentation:**
   - The retrieved context is combined with the user's query.
   - This augmented prompt is then fed into an LLM.
3. **Generation:**
   - The LLM generates a response based on both the user's query and the retrieved context, providing more accurate and grounded answers.

## Key Advantages:

- Reduced Hallucinations**:** RAG minimizes LLM hallucinations by providing factual, external data.
- Increased Accuracy: Answers are more likely to be accurate and relevant to the provided context.
- Knowledge Updates: The external knowledge base can be updated independently of the LLM, allowing for dynamic information integration.

## Advanced RAG:

Advanced RAG builds upon the foundational RAG framework to address its limitations and improve performance. It incorporates more sophisticated techniques for retrieval and context processing

1. **Query Transformation:**
   - Techniques like query rewriting, sub-question generation, and query expansion are used to improve retrieval accuracy.
2. **Chunking and Indexing Strategies:**
   - Advanced chunking methods (e.g., semantic chunking, sliding windows) and indexing structures are employed to optimize context retrieval.
3. **Context Re-ranking and Filtering:**
   - Retrieved documents are re-ranked based on relevance, and irrelevant context is filtered out to reduce noise.
4. **Metadata and Knowledge Graph Integration:**
   - Metadata and knowledge graphs are incorporated to provide structured information and enhance context understanding.

5. **Hybrid Retrieval:**
   - Combines vector search with keyword search or other methods.
6. **Context Compression:**
   - Techniques to reduce the size of the context provided to the LLM, to improve performance and reduce cost.

## Key Improvements:

- Enhanced Relevance: More precise retrieval of relevant context.
- Improved Context Understanding: Better processing and utilization of retrieved information.
- Reduced Noise: Filtering and re-ranking minimize the impact of irrelevant context.
- Increased Efficiency: Context compression and optimized retrieval reduce computational overhead.

# Cache RAG

Cache RAG focuses on optimizing the performance and efficiency of RAG systems by introducing caching mechanisms. It addresses the issue of redundant retrieval and LLM processing for repeated

1. **Caching Retrieval Results:**
   - The retrieved context for a given query is stored in a cache.
   - If the same query is repeated, the cached context is used, avoiding redundant retrieval.
2. **Caching LLM Responses:**
   - The LLM's response for a given query and context is also cached.
   - Repeated queries with the same context can be served from the cache, bypassing LLM processing.

## Key Benefits:

- Reduced Latency: Caching significantly reduces response times for repeated queries.
- Lower Costs: Caching minimizes the number of retrieval and LLM calls, reducing computational costs.
- Increased Throughput: Caching allows RAG systems to handle higher query loads. **In essence:**