

Enhanced RAG Techniques

Advanced Chunking And Preprocessing

pdf, txt, Word, json



Document Splitter

Document

Recursive Character

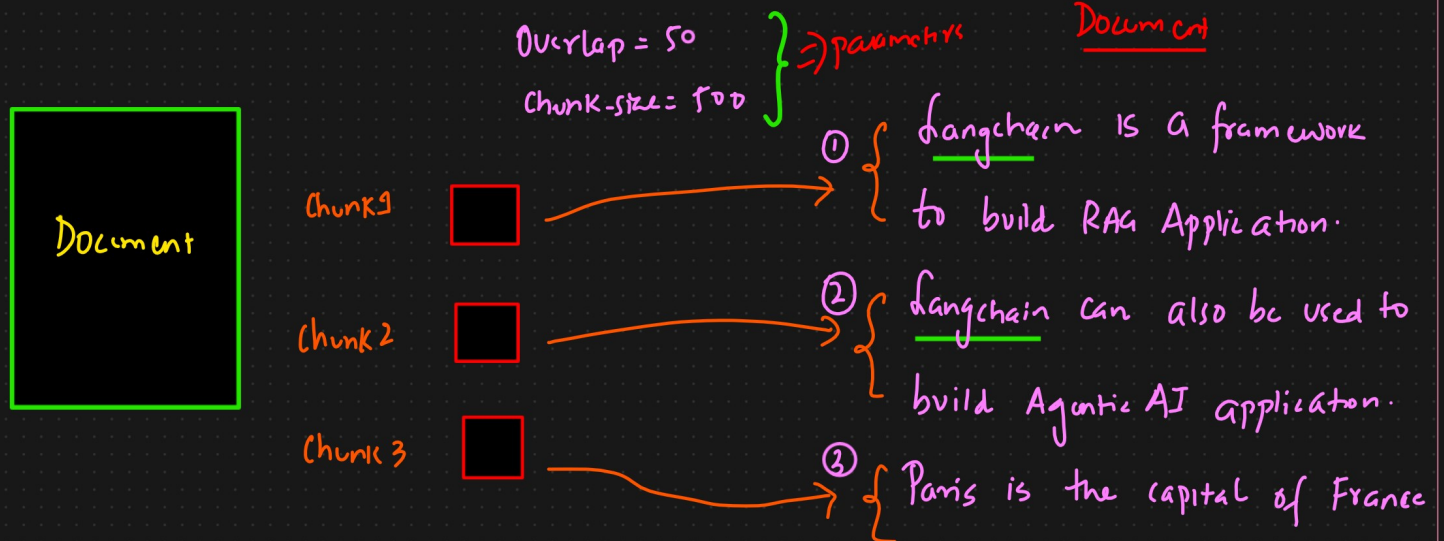
1) Semantic Chunking

Semantic chunking is the process of splitting a document into meaningful units (chunks) based on semantic similarity — not just by number of tokens or lines.

This is important in RAG systems because:

Better chunks → better retrieval → better grounding → better answers.

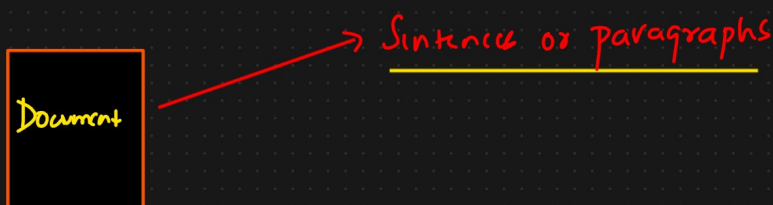
Chunks should be self-contained, contextually rich, and logically separated.



Contextual rich logically separated chunks.

How does It Work?

1) Document Segmentation



2) Sentence Embedding : Each sentence is converted into a Vector representation

OpenAI, Hf \Rightarrow Sentence Transformers

3) Semantic Similarity Check : Cosine Similarity between adjacent embedding

$S_1 \leftrightarrow S_2 \Rightarrow \text{Similar} \Rightarrow 0.95$

$S_2 \leftrightarrow S_3$

$S_3 \leftrightarrow S_4$

$\left\{ \begin{array}{l} \underline{\text{Threshold}} \Rightarrow 0.8 \\ \Rightarrow 0.75 \\ \Rightarrow 0.95 \end{array} \right\}$

4) Merging of Sentences : Merge adjacent sentences if they are semantically similar

5) Form chunks : Merge S_1 and $S_2 \Rightarrow [S_1, S_2] \Rightarrow \text{chunk 1}$
 $\Rightarrow [S_3] \Rightarrow \text{chunk 2}$

Chunks

1. LangChain is a framework for building LLM-powered apps.
2. It integrates with tools like OpenAI and Pinecone.
3. The Eiffel Tower is located in Paris.
4. France is a popular tourist destination.



Semantic chunking

o/p

["LangChain is a framework...", "It integrates with tools..."]
["The Eiffel Tower is located in Paris."]
["France is a popular tourist destination."]