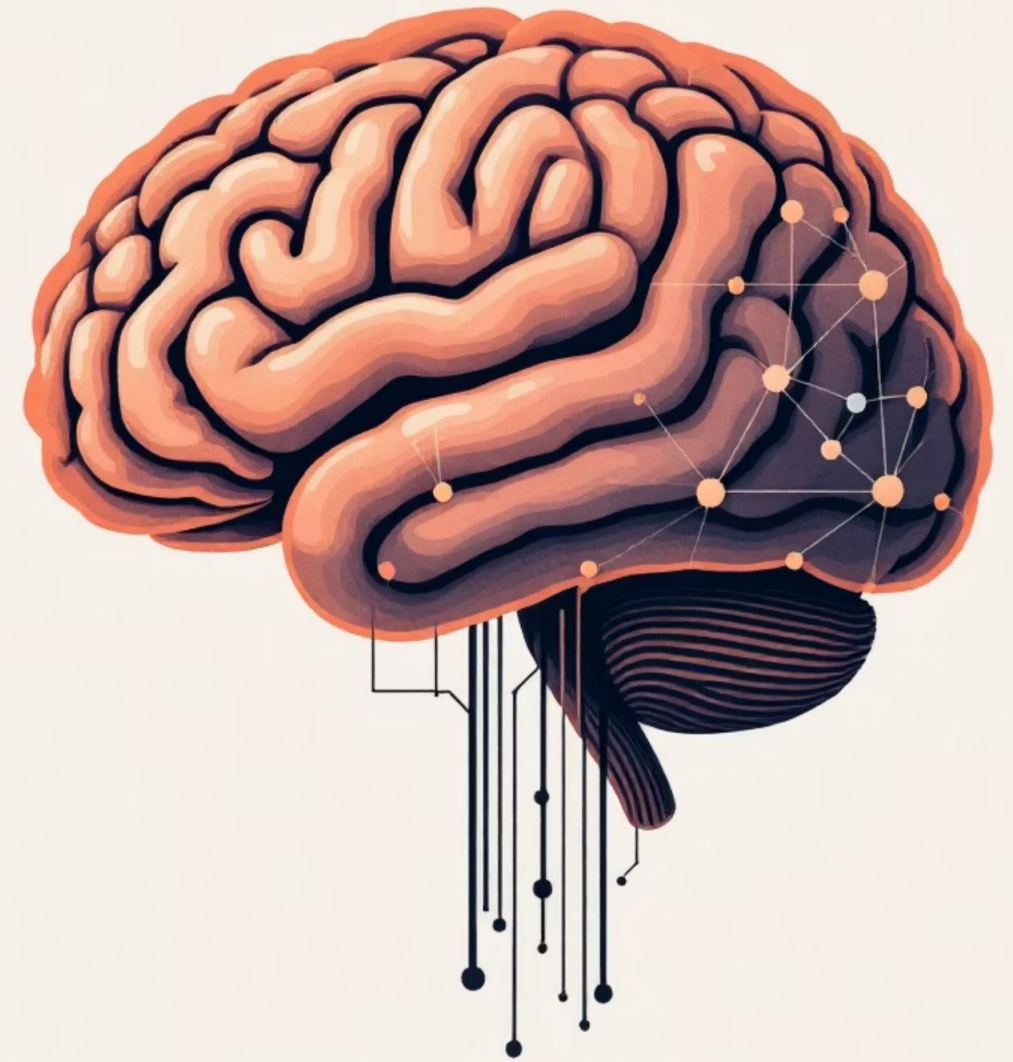# Day 8
# Ingestion & Chunking

# The Art of Chunking: Optimizing RAG Architecture

Understanding how to break down information is crucial for efficient Retrieval Augmented Generation (RAG) systems. This presentation explores the vital role of "chunking" in enhancing LLM performance and overall RAG architecture.

# Why is Chunking Essential for LLMs?

### LLMs' Processing Limits

Large Language Models cannot process entire documents at once. They require content to be broken into manageable, smaller pieces.

### Meaningful Segmentation

Documents must be split into "meaningful chunks" that retain coherence and context, not just arbitrary divisions.

### Enhanced Retrieval

Effective chunking directly leads to better information retrieval, which is the cornerstone of a superior RAG system.

# Navigating the Context Window Limit

Every Large Language Model operates within a defined "context window" – a fixed number of tokens it can process simultaneously.
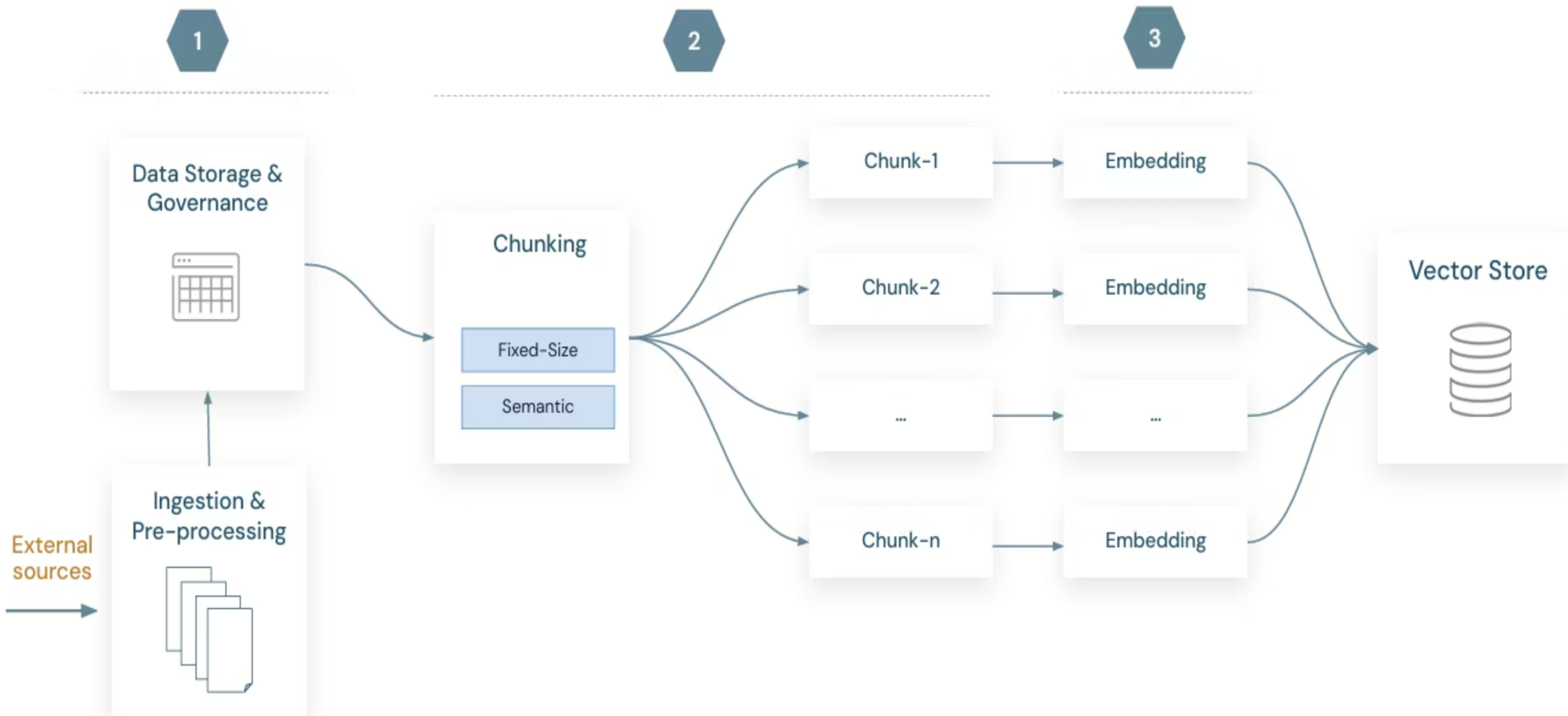
- Models can only "remember" and reason with information within this window.

- To provide more context to the LLM, large documents must be efficiently chunked and relevant pieces retrieved.

- Proper chunking ensures that the most pertinent information fits within the context window, preventing critical details from being overlooked.

When a query requires information beyond this limit, chunking and retrieval become indispensable.

# Data Prep Process Overview

A simple data prep process

# Defining a "Chunk"

## A Small Text Segment

A chunk is not just any arbitrary split. It's a carefully defined section of text.

## Containing One Clear Idea

The fundamental principle: each chunk should encapsulate a single, coherent thought or concept.

## Beyond Random Splitting

Avoid merely cutting text at fixed intervals; this often disrupts meaning and reduces the utility of the chunk.

# Key Chunking Strategies

01

## Fixed Size Chunking

Simple and straightforward, often used as a baseline.

02

## Overlap Chunking

Introduces redundancy to preserve context across boundaries.

03

## Recursive Text Splitter

Hierarchical splitting that respects document structure.

# How to Chunk Data?

## How should we organise it?



**Title**

**Section**

**Diagram**

**Semantic Chunking**:

- Chunk by sentence/paragraph/section
- Leverage special punctuation (i.e. '.', '\n')
- Include/Inject metadata/tags/title(s)

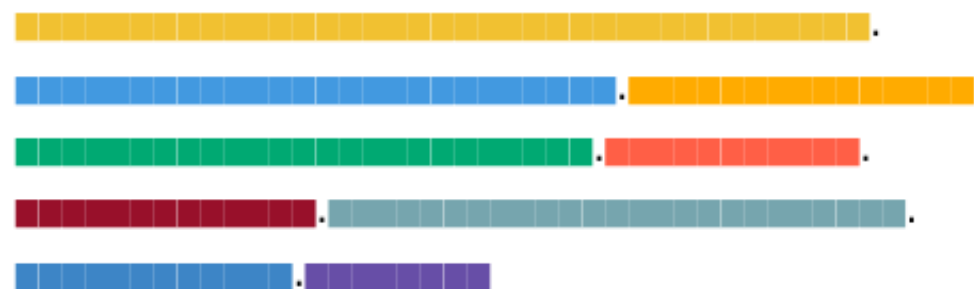### &/OR

**Fixed-size Chunking**:

- Divide by a specific number of tokens
- Simple and computationally cheap method
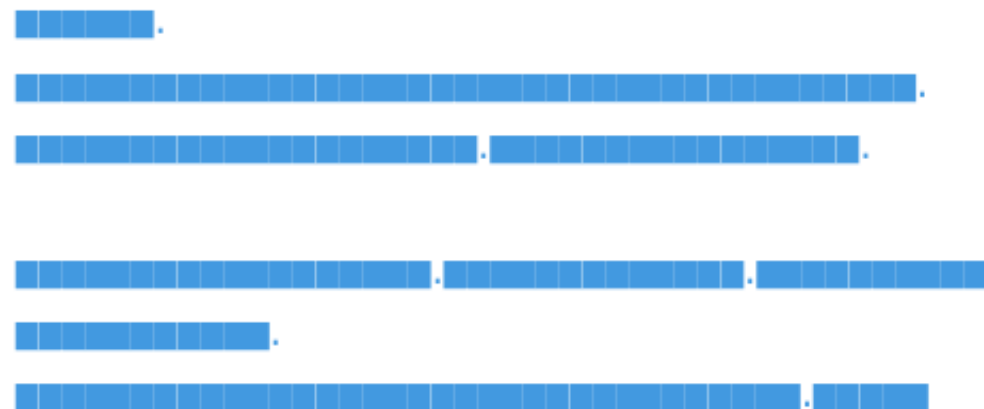
# Chunking Strategy is Use-Case Specific

Another iterative step! Experiment with different chunk sizes and approaches

- How long are our documents?
  - 1 sentence?
  - N sentences?

- If 1 chunk = **1 sentence**, embeddings focus on specific meaning

- If 1 chunk = **multiple paragraphs**, embeddings capture broader theme
  - How about splitting by headers?
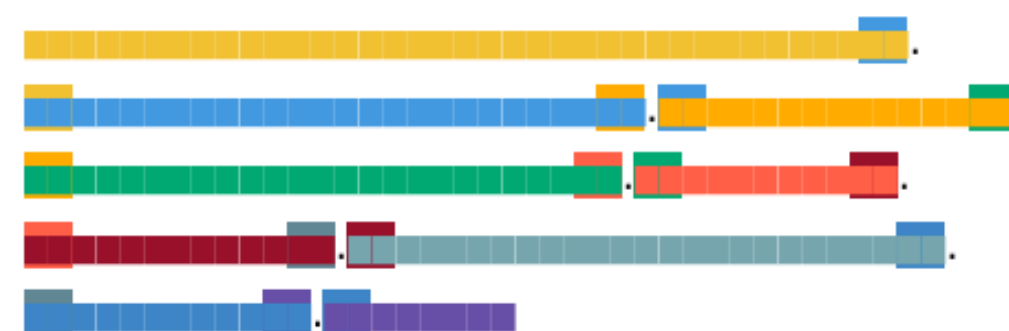
**Chunking by sentence:**
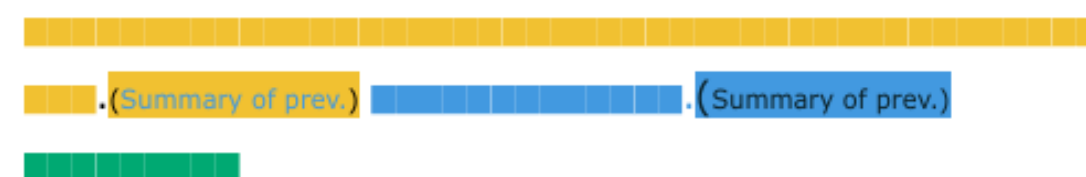
**Chunking by Paragraph:**

# Chunking Strategy is Use-Case Specific

Another iterative step! Experiment with different chunk sizes and approaches

- Chunk **overlap** defines the amount of overlap between consecutive chunks, ensuring that no contextual information is lost between them.

- **Windowed summarization** is a 'context-enriching' chunking method where each chunk includes a 'windowed summary' of previous few chunks.

**Chunk overlap:**



**Windowed summarization:**



(Summary of prev.) (Summary of prev.)

- Prior knowledge of user's query patterns can be helpful (*i.e. query length?*)

  - While long queries may have better aligned embeddings to returned chunks, shorter queries could be more precise

# Advanced Chunking Strategies

## Summarization

# Advanced Chunking Strategies

## Summarization with metadata

# Data Extraction and Chunking Challenges

Working with complex documents



**Other challenges:**

- Text mixed with image
- Irregular placement of text
- Color encoded focus (*Important for context*)

# Data Extraction and Chunking Challenges

Working with complex documents



Flow Chart with order info

Multi-column text

Image with related info

## Other challenges:

- Chart with hierarchical information. Keeping the order of the information is critical.

- Multi-column text and the order of columns if crucial.

- Keeping images with related information is crucial.

# Fixed Size Chunking: Pros and Cons

## Pros:

- **Ease of Implementation:** It's the simplest chunking method to set up and execute.
- **Predictable Output:** Generates chunks of a consistent length, which can be useful for certain models.

## Cons:

- **Loss of Meaning:** Can frequently cut sentences or ideas mid-flow, leading to incoherent chunks.
- **Contextual Gaps:** Important connections between cut text segments can be lost, hampering retrieval accuracy.

# The Advantage of Overlapping Chunks

Overlapping chunks address a critical flaw in fixed-size methods by maintaining contextual continuity.

- **Preserves Sentence Integrity:** By allowing a small portion of text to repeat in subsequent chunks, it ensures that sentences and short ideas are not cut off abruptly.

- **Reduces Contextual Gaps:** The overlap acts as a bridge, linking related information and improving the chances of retrieving a complete thought.

- **Minor Redundancy for Major Gain:** While it introduces a small amount of redundant information, the benefit of improved coherence far outweighs this drawback.

# Recursive Text Splitter: Structure-Aware Chunking

### Prioritizes Document Structure

This advanced method first attempts to split text based on hierarchical structures like headings.

### Breaks Down to Paragraphs

If the text is still too large, it then splits into individual paragraphs.

### Final Split by Sentences

As a last resort, it breaks down paragraphs into sentences, ensuring the smallest meaningful units.

### Most Accurate for RAG

By respecting the inherent organization of the document, it provides the most contextually relevant chunks for RAG systems.

# Optimal Chunk Configuration for Q&A RAG

For most Question-Answering (Q&A) based RAG applications, a carefully tuned chunk size and overlap can significantly improve performance.

### 300

**Chunk Size (Tokens)**

This size generally allows for sufficient context without overwhelming the LLM's context window. It captures enough detail for most queries.

### 50

**Overlap (Tokens)**

A 50-token overlap effectively bridges potential gaps, ensuring that key information at chunk boundaries remains connected and retrievable.

This configuration balances detail, context preservation, and processing efficiency, making it ideal for robust Q&A interactions.

# Good vs. Bad Chunks: A Visual Comparison

## Good Chunk

A "good" chunk provides clear context, focusing on a single, complete topic. It answers a potential question without ambiguity and retains all necessary surrounding information.

- **Clear context:** All relevant information for one idea.
- **Single topic:** Focused and coherent content.
- **Well-bounded:** Starts and ends logically.

## Bad Chunk

A "bad" chunk might be cut mid-sentence, contain fragmented ideas, or blend unrelated lines. This leads to confusion and hinders effective retrieval, making it difficult for the LLM to understand or answer queries accurately.

- **Mid-sentence cut:** Incomplete thoughts.
- **Unrelated lines:** Jumbled, incoherent information.
- **Ambiguous context:** Hard for LLM to interpret.

# Hands on tasks !