

# Context Windows: The Foundation of LLM Memory

Understanding How Much LLMs Can Remember

## What is a Context Window?

- **Definition:** The maximum amount of text (in tokens) an LLM can process at once.
- **Includes:** Your prompt + conversation history + generated response.
- **Limitation:** Fixed size determined by the model's architecture.
- **Behavior:** When exceeded, oldest tokens are "forgotten" (sliding window).

## Visual Metaphor

Think of it like a whiteboard:

Small Context

(4K)

✓

Fits a short essay or brief conversation

✓

Good for quick tasks and simple queries

Medium Context

(32K)

✓

Fits a research paper or long document

✓

Ideal for code projects & doc analysis

Large Context

(200K)

✓

Fits a small book or extensive dialogue

✓

Perfect for complex, long-form analysis

## Why Context Windows Matter



Task Capability



Cost Implications



Use Case Selection

💡 **Critical Insight:** Context window size is often the primary factor in choosing an LLM. It's not just about intelligence, but about memory capacity and its suitability for a specific task.