

Model	Developer	Context Window Size	Description	Key Use Cases
GPT-4 (turbo)	OpenAI	Up to 128k tokens	Allows processing of long documents, chat history, and complex conversations.	Document summarization, code analysis
Claude 3	Anthropic	Up to 200k tokens	Optimized for large inputs like books, research papers, and comprehensive datasets.	Literature analysis, extended research
LLaMA 2	Meta	4k-32k tokens (varies by version)	Open-source LLM designed for research and performance scaling.	Custom applications, R&D
Mistral 7B	Mistral AI	4k tokens	Compact yet efficient, with high performance for general-purpose tasks.	Conversational AI, summarization
Falcon 180B	Technology Innovation Institute (TII)	4k tokens	High-performing model designed for industry-specific applications.	Business-specific AI solutions
Cohere Command R+	Cohere	16k tokens	Fine-tuned for retrieval-augmented tasks and enhanced for longer documents.	Search and knowledge retrieval
PaLM 2	Google	8k-32k tokens (varies by version)	Designed for a wide range of applications, including multilingual and multimodal tasks.	Multimodal AI, complex logic tasks
ChatGLM 2	Tsinghua University	8k tokens	Bilingual (Chinese-English) model optimized for extended dialogue and reasoning.	Multilingual dialogue, academic tasks

Key Notes:

1. What is the Context Window?

- The context window represents the number of tokens (words, punctuation, etc.) the model can process at one time.
- Larger context windows are critical for handling long documents, coding tasks, or maintaining conversational continuity.

2. Current Trends:

- Context windows are rapidly expanding to enable LLMs to handle more complex and lengthy inputs, such as books, multi-step reasoning, and extensive histories.

3. Choosing the Right Model:

- **Short Tasks:** Models like **Mistral 7B** or **Falcon 180B** work well.
- **Longer Inputs:** Use models like **Claude 3** or **GPT-4 Turbo**.
- **Specialized Domains:** Models like **Cohere Command R+** or **ChatGLM 2** shine for targeted applications.

Historical Chat Reference and context window

- **Current Context:** Most LLMs do not inherently "remember" past interactions across sessions unless explicitly programmed with memory mechanisms or external databases.
- **Within-Session Context:**
 - The chat history is passed back to the model as part of the prompt. If the context window is limited (e.g., 4,000 tokens), older parts of the chat history are truncated.
 - Truncating history reduces memory usage but at the cost of losing context.

Persistent Memory Solutions to use context window optimally

- To extend "memory" across sessions, systems integrate **external memory mechanisms**:
 - **Vector Databases:** Store embeddings of prior conversations for retrieval.
 - **Summarization:** Compress past chats into summaries to reduce token load in the context window.

Key Trade-offs

- **Larger Context Window:**
 - Pros: More immediate memory of conversation history.
 - Cons: Higher memory usage, longer inference time, and increased computational costs.
- **Smaller Context Window:**
 - Pros: Faster processing and lower resource consumption.
 - Cons: Limited historical context retention, requiring external solutions for long-term memory.

Memory Usage Factors

The memory required depends on the following factors:

a. Tokenization Overhead

- Text is tokenized into smaller units (tokens). For instance, "ChatGPT" might tokenize into two tokens, while "hello" is one.

- The number of tokens in the context window directly affects memory usage. More tokens mean higher memory consumption.

b. Attention Mechanism Complexity

- LLMs like GPT utilize a transformer architecture, where the attention mechanism computes relationships between all token pairs in the context.
- **Computational Complexity:** Attention has $O(n^2)$ complexity, where n is the number of tokens in the context window. Doubling the context window size roughly quadruples memory usage.

c. Model Parameters

- Each token in the context interacts with the model's parameters. The larger the model (more parameters), the higher the memory consumption for processing the same number of tokens.

Conclusion

The relationship is **quadratic**, meaning memory usage scales sharply with an increase in the context window. Optimizations like summarization, selective retrieval, or hybrid memory systems can help balance context size and memory constraints in practical applications.

For details about the Generative AI workshop, visit:

<https://www.nitinkapse.com/generative-ai>

Register here to secure your spot:

<https://forms.gle/PrzkmvYh5yvEWUKZ6>