

# The Next Leap in AI: How Titans are Overcoming Transformer Limitations

AIML Team, Kanaka Software

20 Jan 2025

As Gen AI developers, we're all familiar with the power of **Transformers**, those incredible models that have revolutionised how machines understand and generate language. But like any technology, they have limitations, especially when dealing with complex, real-world tasks. Today, let's explore these issues and how a new architecture called **Titans** aims to tackle them, and what this means for the future of AI.

## 1 The Transformer Bottleneck

Current [transformer](#) models, while powerful, face some key challenges:

- **Limited Context Window:** Transformers use "attention" to understand relationships between words. However, this attention is limited to a short context window, making it difficult to remember and use information from earlier parts of a long text. This is because the attention mechanism has a quadratic cost.
- **Quadratic Complexity:** The computational cost of attention grows dramatically with the length of the input text, leading to high memory usage and slow processing. This makes them inefficient for very long sequences of text.

- **Memory Inefficiency:** Transformers store all past information within the context window without compression. This leads to excessive memory consumption and slower processing.
- **Reasoning and Generalisation Issues:** Transformers can struggle with tasks that require complex reasoning, length extrapolation, or the ability to generalise to new situations. They lack distinct components for short-term and long-term memory and the ability to actively learn from data.

## 2 Titans: A New Approach to Memory

**Titans** is a new architecture designed to overcome these limitations. It introduces a **neural long-term memory** module that learns to memorise historical context, working more like human memory with distinct but interconnected systems. Here's how Titans addresses the issues:

- **Neural Long-Term Memory:** Titans incorporate a special module that learns to store important information over time. This module doesn't just remember the exact words but also the key abstractions and ideas of what has been processed. It also learns to forget less important information to manage its memory capacity. This is a deep neural network which is trained using a loss function.
- **Short-Term Memory (Attention):** Titans still use attention for short-term focus on the current text, ensuring accurate modelling of dependencies within a limited context.
- **Three Branches:** Titans consist of three different branches each responsible for different aspects of the task: A '**Core**' branch that uses short-term memory to process the main flow of the data. A '**Long-term Memory**' branch that houses the neural long-term memory module. A '**Persistent Memory**' branch to store general knowledge about the task.

- **Learning to Memorise at Test Time:** Titans are unique because they continue to learn and adapt even while being used (at test time). They learn what to remember, how to store it, and how to retrieve it effectively, improving their performance over time. They measure surprise of an input to decide what to memorise.
- **Efficiency:** Titans are more efficient than Transformers, capable of handling longer sequences with less computational cost. This is due to the compression of information and parallelizable training.

### 3 What This Means for the Future of Gen AI

The implications of Titans are huge for the future of Gen AI:

- **Better Long-Form Content:** Titans' improved memory will enable AI to create much better long-form articles, stories, and even code.
- **More Coherent Conversations:** AI chatbots will be able to have more natural and context-aware conversations by remembering past exchanges.
- **Improved Reasoning:** Enhanced memory will lead to AI that is better at reasoning and problem-solving, enabling more complex tasks like medical diagnosis and financial analysis.
- **Handling Larger Contexts:** Titans are capable of scaling to larger context windows with higher accuracy than current transformer models, making them suitable for handling very large documents.
- **More Human-like AI:** By incorporating distinct short and long-term memory modules, Titans aim to mimic human learning and cognitive processes, potentially resulting in more human-like AI.

**In Simple Terms**

Transformers are like having a fast computer with a very short attention span. Titans are like a computer with the ability to learn, remember important information, and use it to handle more complex and long context tasks. This means future AI could be much more creative, conversational, and better equipped to tackle real-world problems.

**Titans represent a significant step forward, offering a way to overcome the memory and efficiency limitations of current Transformer models.** This advancement is crucial for the next generation of AI applications, promising a future where AI is more capable, more human-like, and better suited for complex, long-term interactions.