

## Solving Context Length Issue of Coding LLMs

Ashrey Ignise<sup>1\*</sup> and Yashika Vahi<sup>2</sup><sup>1</sup>Chief Executive Officer, ArtusAI Workspaces Pvt Ltd, Boston, USA<sup>2</sup>Research Scientist, ArtusAI Workspaces Pvt Ltd, Vancouver, British Columbia, Canada**\*Corresponding Author**

Ashrey Ignise, Chief Executive Officer, ArtusAI Workspaces Pvt Ltd, Boston, USA.

Submitted: 2024, Oct 17; Accepted: 2024, Nov 11; Published: 2024, Nov 19

**Citation:** Ignise, A., Vahi, Y. (2024). Solving Context Length Issue of Coding LLMs. *J Electr Comput Innov*, 1(1), 01-06.**Abstract**

Large Language Models (LLMs) face significant challenges related to context length when generating code, often resulting in incoherent or incomplete outputs. This paper aims to explore the context length issue, present technical solutions, and suggest future directions for improving context retention in LLMs used for code generation.

**Keywords:** Intelligent Agents, Artificial Intelligence, Distributed Artificial Intelligence, Large Language Models, LLM Limitations**1. Introduction**

When it comes to code generation, keeping things in context is essential to generating results that are accurate and logical. Even while LLMs are strong, they frequently have problems with context duration, which makes it hard to maintain long-term dependencies that are essential for complicated coding jobs. The significance of context in code generation is discussed in this paper, along with the unique difficulties caused by context length restrictions in LLMs. The paper's structure is also presented, covering technical solutions, hybrid approaches, evaluation techniques, and potential avenues for future research. Through comprehension and resolution of these problems, we can improve LLMs' ability to produce high-caliber code.

**2. Understanding Context Length Issues****2.1. Definition of Context Length****Explanation of Context Length in LLMs**

In the context of Large Language Models (LLMs), context length is the maximum amount of text or code that the model can efficiently process at once while producing an output. Usually, its length is expressed in tokens, which are the textual units that the model interprets. Tokens can be words or subwords in natural language processing, for example, and keywords, operators, and identifiers in code creation.

Due to their fixed context windows, LLMs like as GPT-3 are limited in how much of the previous text they can utilize to guide their responses. For instance, the maximum context window for GPT-3 is 2048 tokens. This indicates that the model only takes into account the last 2048 tokens and ignores anything that comes before this window when producing text or code.

**Importance of Maintaining Context for Coherent Code Generation** Maintaining context is critical in code generation for several reasons:

- Variable Scope:** Understanding the scope of variables and functions is essential for generating code that correctly references them. If the context is lost, the model might generate code that uses variables incorrectly or out of scope.
- Logical Flow:** For complex coding tasks, the logical flow and structure of the code are paramount. Losing context can lead to disjointed and logically inconsistent code snippets.
- Dependency Management:** Many programming tasks require understanding dependencies between different parts of the code. Losing context can cause the model to overlook these dependencies, leading to incomplete or incorrect code.

**Example:** Consider a function in Python that initializes variables, processes data, and then outputs results. If the context window is too short, the model might forget the initialization step by the time it needs to reference those variables, resulting in errors or incorrect outputs.

**2.2. Challenges Posed by Context Length****Difficulty in Retaining Long-Term Dependencies**

LLMs struggle with retaining long-term dependencies due to their fixed context windows. As the complexity and length of the input increase, maintaining a coherent understanding of the entire context becomes increasingly challenging. This difficulty is particularly pronounced in tasks that require understanding and generating code over multiple lines, functions, or files.

**Example:** In a large codebase, understanding how a function defined in one file interacts with a function in another file requires maintaining context across potentially thousands of lines of code. LLMs with limited context windows may lose track of these dependencies, leading to errors.

**Examples of Context Loss in Generated Code**

- Incomplete Function Implementations:** When generating

---

code for a complex function, an LLM might lose track of the initial parameters or variables, leading to incomplete or incorrect function implementations. For instance, an LLM might forget to return a value at the end of a function or incorrectly reference a variable that was defined earlier.

2. **Logical Inconsistencies:** In longer code sequences, the model might produce logically inconsistent code. For example, it might generate a loop that references variables not defined within the loop's scope, or it might include redundant or contradictory statements due to lost context.
3. **Misplaced Comments and Documentation:** When generating code along with comments, the model might place comments inappropriately or generate documentation that does not align with the code logic, leading to confusion and errors during code review.

### Impact on Code Quality and Accuracy

Context length constraints have a substantial effect on the precision and calibre of the code that is produced:

1. **Decreased Code Reliability:** When code is created without the necessary context, it frequently has faults and errors, which lowers its reliability and decreases its usefulness in practical applications.
2. **Increased Debugging Effort:** The productivity gains promised by automated code creation are diminished because developers have to spend more time troubleshooting and fixing the code produced by LLMs.
3. **Inconsistent Code Style:** When context is lost, disparate sections of the generated code may adhere to disparate standards or patterns, which makes the codebase more difficult to comprehend and maintain.

*Example:* In a collaborative coding environment, one developer might use an LLM to generate a new feature. If the LLM loses context and generates code with a different style or logic than the rest of the codebase, it can create integration challenges and increase the workload for the team.

### Conclusion

Understanding the context length issue is crucial for improving LLMs in code generation. Addressing these challenges requires both technical innovations and practical strategies to enhance the models' ability to maintain and utilize context effectively, ensuring higher quality and more reliable code outputs.

## 3. Technical Solutions to Context Length Issues

### 3.1. Memory-Augmented Networks

#### Overview of Memory-Augmented Neural Networks

A class of models known as memory-augmented neural networks (MANNs) was created to increase the capabilities of conventional neural networks by adding an external memory component. Compared to the typical context window of LLMs, this external memory enables the network to store and retrieve data across longer sequences. By doing this, MANNs hope to improve LLMs' capacity to handle long-term dependencies and preserve important context information, hence mitigating the drawbacks caused by fixed context lengths.

Generally, MANNs are made up of two primary parts:

1. **Neural Network Controller:** A component of the model that

communicates with the external memory. It has the ability to write to and read from memory, and it will eventually learn how to use these functions most effectively to increase performance on tasks requiring long-term context.

2. **External Memory Module:** Information that is not immediately related to the neural network controller can be stored in this differentiable memory bank. This module enables the network to retain pertinent data across longer sequences; it can be compared to an extension of its short-term memory.

#### Benefits of External Memory Components

1. **Extended Context Retention:** By utilizing an external memory, MANNs can maintain relevant information over extended sequences, reducing the likelihood of context loss and improving the coherence and accuracy of generated outputs.
2. **Enhanced Problem Solving:** The ability to store and retrieve information from an external memory allows MANNs to handle more complex tasks that require understanding and integrating information over longer periods, such as generating cohesive code that spans multiple functions or files.
3. **Improved Learning Efficiency:** With the external memory component, MANNs can learn to use their memory more efficiently, focusing on storing critical information and discarding irrelevant details, which can lead to better performance on tasks requiring long-term dependencies.

#### Examples and Case Studies

1. **Differentiable Neural Computers (DNCs):** Developed by DeepMind, DNCs are a kind of MANN intended to handle jobs requiring long-term knowledge retrieval and complicated thinking. It has been demonstrated that DNCs function well on a range of activities where preserving long-term context is essential, including as pathfinding in mazes and answering questions.
2. **Neural Turing Machines (NTMs):** They provide a further instance of MANNs. By integrating a memory matrix that the neural network controller can access and manipulate, NTMs merge the advantages of neural networks and Turing machines. NTMs have proven their ability to manage extensive context dependencies by succeeding in tasks like sorting, associative recall, and copying lengthy sequences.

*Case Study:* When it comes to code generation, a business that a complex software system used a DNC to maintain context across multiple modules and functions. By leveraging the external memory component, the DNC was able to generate code that correctly referenced variables and functions defined earlier in the sequence, significantly reducing the error rate and improving the coherence of the generated code. This approach not only enhanced the overall quality of the code but also reduced the time developers spent on debugging and integration.

### 3.2. Hierarchical Models

#### Explanation of Hierarchical Model Structures

Because they can process data at many levels of abstraction, hierarchical models are better able to handle lengthy sequences

---

by segmenting them into smaller, easier-to-manage pieces. This structure is especially helpful for operations involving long sequences, like code creation, when it's important to maintain context over longer periods of time.

In hierarchical models, the data is processed in a layered manner:

1. **Lower-Level Processing:** At this level, the model handles smaller, localized segments of data. For code generation, this might involve understanding and generating individual lines or small blocks of code.
2. **Higher-Level Processing:** This level processes larger chunks of data by aggregating information from the lower level. In code generation, this could involve understanding the structure of entire functions or classes.
3. **Top-Level Processing:** At the highest level, the model integrates and understands the overall context, such as the entire codebase or a complete project. This enables the model to maintain coherence and context over long sequences.

#### Advantages for Managing Long Sequences of Code

1. **Better Context Retention:** The model's ability to keep context across longer sequences is enhanced by its hierarchical data processing. Higher-level components make sure that the larger context is preserved, while lower level components deal with the current context.
2. **Scalability:** Because each level of the hierarchy only needs to manage a fraction of the data, hierarchical models may handle big datasets or lengthy sequences more effectively.
3. **Improved Understanding:** The multi-level strategy enables the model to comprehend both abstract and particular facets of the input, producing outputs that are more logical and correct in relation to the context.
4. **Modular Processing:** This method allows for modularity, which allows for the independent re-tuning of various levels of the hierarchy, resulting in more effective updates and enhancements.

#### Implementation Examples

1. **Transformer-Based Hierarchical Models:** Hierarchical transformers enhance the standard transformer design by including hierarchical processing. A transformer model, for instance, might process distinct functions or methods using a local transformer and combine them into a coherent whole using a global transformer. With the help of this structure, the model can produce cohesive code that maintains dependencies and context across several functions or files. *An illustration* of a hierarchical transformer model would be the local transformer, which would concentrate on producing a particular function inside a class and comprehending the variables and logic unique to that function. In contrast, the global transformer would make sure that the code created is in line with the project's overall structure by integrating this function into the class's and the codebase's bigger context.
2. **Hierarchies of Recurrent Neural Networks (RNNs):** Layers of RNNs are used in hierarchical RNNs, with each layer processing data at a distinct abstraction level. Higher layers collect these facts to preserve a wider context, while lower layers manage more specific details. *Example:* The

lower layer of a hierarchical RNN for code creation may handle individual lines of code, guaranteeing instantaneous context and proper syntax. After that, the top layer would combine these lines to create cohesive functions or methods, preserving dependencies and logical flow throughout the whole codebase.

3. **Hierarchical Attention Networks (HANs):** The model can concentrate on various tiers of the data hierarchy thanks to HANs, which expand attention methods to hierarchical processing. HANs are especially useful for operations like code creation when long-term dependencies need to be maintained. *An illustration* of this would be in a HAN for code creation, where the higher-level attention mechanism takes into account the function's overall structure and function inside the class or module, while the lower-level attention mechanism would concentrate on particular tokens or statements within a function. This method guarantees that the resulting code is cohesive at the macro level and contextually correct at the micro level.

#### Conclusion

Hierarchical models offer a robust solution to the context length issues in LLMs by breaking down long sequences into manageable chunks and processing them at multiple levels of abstraction. This approach enhances context retention, scalability, and understanding, leading to more coherent and accurate code generation. As research and development in hierarchical models continue, they hold significant potential for improving the capabilities of LLMs in various applications, particularly in handling long sequences of data.

#### 3.3. Sliding Window Mechanisms

##### Concept of Sliding Window Techniques

The sliding window technique is a method used to manage large sequences of data by processing them in overlapping pieces, or "windows," of a given size. To preserve some context from earlier portions, each new window that glides across the series incorporates a piece of the preceding window. With this method, context is preserved over lengthy sequences without overburdening the model with all of the input at once.

A sliding window can be used in the context of LLMs for code generation to divide lengthy code sequences into more manageable, smaller pieces that the model can handle one after the other. A predetermined amount of tokens or lines of code are included in each window, and when the window glides, it overlaps with the preceding window to guarantee continuity of context.

##### Balancing Context Length and Computational Efficiency

Sliding window technique's main benefit is its ability to strike a balance between computing efficiency and context length. This is how it operates:

1. **Context Preservation:** The model can save some of the prior context by overlapping windows, which aids in preserving coherence and continuity in the code that is created. For jobs requiring comprehension of relationships and dependencies between various code segments, this is essential.

- 2. Manageable Processing:** By limiting the amount of data the model processes at once, the fixed window size helps to lower memory and computational strain. Long sequences can now be handled without taxing the model's capabilities.
- 3. Flexibility:** The size of sliding windows can be changed to suit the particular needs of the job. Greater window size provide more context at the cost of higher computational demand, while smaller windows are more efficient but might capture less context.

#### Practical Applications:

- 1. Code Review and Refactoring:** Sliding window techniques can be used in automated code review tools to analyze large codebases incrementally. By processing the code in overlapping windows, the tool can maintain context and continuity, ensuring that changes are coherent and dependencies are respected. *Example:* An automated code review tool might use a sliding window of 50 lines to scan through a large code file. As it reviews each window, it can flag issues, suggest improvements, and ensure that changes do not introduce errors or inconsistencies.
- 2. Documentation Generation:** For generating documentation for extensive codebases, sliding windows can help the model understand and describe sections of code incrementally. This ensures that the generated documentation is coherent and contextually accurate. *Example:* A documentation tool might process 100-line windows of code to generate descriptions and summaries. By overlapping windows, it can maintain context and produce detailed, accurate documentation that covers the entire codebase.
- 3. Bug Detection and Fixing:** Sliding windows can be applied in tools designed to detect and fix bugs in large codebases. By examining code in overlapping segments, the tool can identify issues that span multiple sections of the code and suggest comprehensive fixes. *Example:* A bug detection tool might use sliding windows of 20 lines to scan for errors. If a bug affects multiple functions, the overlapping windows ensure that the tool captures the full context and provides an effective solution.

#### Limitations:

- 1. Context Loss:** Although sliding windows contribute to the preservation of context, context may still be lost at the window's borders. The model may ignore important information that falls outside of the current window, which could result in mistakes or inconsistent results. *Example:* If a variable is defined at the beginning of a long function and used at the conclusion, a sliding window may not fully capture its usage, which could result in misinterpretations or improper code generation.
- 2. Computational Overhead:** Although sliding windows reduce the overall computational load compared to processing the entire sequence, overlapping windows can still introduce some redundancy and overhead. Each window requires separate processing, which can increase the total computation time. *Example:* Overlapping windows of 50 lines each with a 10-line overlap might lead to redundant processing of certain lines, slightly increasing the overall computation time and resource usage.

- 3. Optimal Window Size:** Determining the optimal window size is a challenge. Too small a window might lead to insufficient context, while too large a window could negate the computational efficiency benefits. *Example:* Finding the balance between capturing enough context and maintaining computational efficiency requires experimentation and tuning based on the specific use case and model capabilities.

#### Conclusion

Sliding window mechanisms offer a practical solution to the context length issues in LLMs for code generation by processing long sequences in manageable, overlapping segments. While this approach helps maintain context and reduces computational load, it also comes with challenges such as potential context loss and computational overhead. By carefully balancing window size and overlap, and applying this technique to appropriate tasks, developers can leverage sliding windows to enhance the effectiveness of LLMs in various applications.

#### 4. Enhancing Context Retention with Hybrid Approaches

##### 4.1. Combining Symbolic AI and Neural Networks

###### Overview of Hybrid Approaches

Neural networks (machine learning models) and symbolic AI (rule-based systems) are combined in hybrid techniques. Neural networks are good at managing unstructured tasks and learning from data, while symbolic AI is best at handling organized, logical tasks with explicit rules and representations. Through the combination of these two methods, the advantages of each can be utilized to boost context retention and optimize code creation as a whole.

###### Benefits of Integrating Symbolic AI for Context Retention

- 1. Structured Reasoning:** Symbolic AI can help maintain long-term dependencies and logical structures within the code. By representing rules and constraints explicitly, it ensures that the generated code adheres to specific guidelines and standards.
- 2. Enhanced Context Management:** Neural networks can struggle with retaining context over long sequences. Symbolic AI can bridge this gap by providing a structured way to manage and recall important information across different segments of code.
- 3. Error Reduction:** Symbolic AI can validate and correct the output of neural networks, reducing errors and improving the accuracy of the generated code. This hybrid approach ensures that the generated code is both contextually relevant and logically correct.

###### Case Studies Demonstrating Effectiveness

- Case Study 1: Automated Code Refactoring**

In an automated code refactoring tool, symbolic AI can be used to identify patterns and rules for refactoring, while neural networks handle the generation of new code segments. For instance, the symbolic AI can ensure that variable names are consistent and follow naming conventions, while the neural network generates the actual refactored code. This combination improves the quality and coherence of the refactored code.

- Case Study 2: Intelligent Code Completion**

Neural networks predict and recommend code snippets depending on the present context, whereas symbolic AI can

---

give guidelines for syntactic and semantic accuracy in an integrated development environment (IDE) with intelligent code completion. By using a hybrid method, it is ensured that the code recommendations follow standards and best practices for programming in addition to being contextually relevant.

## 4.2. Multi-Stage Generation Processes

### Explanation of Multi-Stage Generation

Multi-stage generation involves breaking down the code generation process into multiple stages, each handling a specific aspect of the task. This approach helps manage long sequences by focusing on smaller, manageable segments and progressively building up the final output. Each stage adds context and detail, ensuring that the overall coherence and quality of the code are maintained.

### Breaking Down Code Generation into Manageable Segments

**Stage 1: High-Level Planning:** The first stage involves generating a high-level plan or outline of the code. This includes defining the main functions, classes, and modules, and their interactions. This stage sets the foundation for the subsequent stages.

**Stage 2: Detailed Code Generation:** In the second stage, the high-level plan is expanded into detailed code. This involves generating the actual code for each function, class, and module, ensuring that the logic and structure are correctly implemented.

**Stage 3: Context Integration:** The third stage focuses on integrating context and ensuring coherence across the generated code. This involves linking different segments, resolving dependencies, and maintaining consistency.

**Stage 4: Streamlining and Enhancement:** The resulting code must be improved and optimized in the last step. This entails making the code more readable, performing better, and adhering to standards and best practices.

### Examples and Impact on Context Retention

#### • Example 1: Web Application Development

The multi-stage generation technique can be used to create the general architecture (Stage 1) and the detailed code for each component (Stage 2) of a web application. Refinement increases performance and maintainability, while context integration guarantees that the front-end and back-end components are correctly integrated (Stage 4). This method guarantees maintainability, efficiency, and coherence in the code generated.

#### • Example 2: Data Processing Pipeline

For a data processing pipeline, the high-level plan defines the main stages of data ingestion, processing, and output (Stage 1). Detailed code is generated for each stage, including data transformation and analysis (Stage 2). Context integration ensures that data flows correctly between stages (Stage 3), and refinement optimizes performance and resource usage (Stage 4). This multi-stage approach ensures that the pipeline is robust, efficient, and scalable.

### Conclusion

Hybrid approaches, including the integration of symbolic AI with neural networks and multi-stage generation processes, offer effective solutions to the context length issues in LLMs for code generation. By leveraging structured reasoning and breaking down code generation into manageable segments, these

approaches enhance context retention, improve code quality, and ensure coherence. As these techniques continue to evolve, they hold the potential to significantly advance the capabilities of LLMs in generating high-quality, contextually accurate code.

## 5. Future Directions and Research Opportunities

### 5.1. Emerging Technologies and Innovations

The field of artificial intelligence and machine learning is always changing, and new technologies are coming out that claim to solve the context length problems in LLMs. These include developments in innovative architectures such as transformers with improved context-handling capabilities, memory-augmented neural networks, and hierarchical models. Extended sequences of generated code are being investigated to preserve long-term dependencies and enhance coherence through the use of techniques like recurrent memory networks and sophisticated attention mechanisms.

These new technologies have the potential to greatly improve LLMs' code generating capabilities. They can result in the creation of more precise, logical, and contextually relevant code by enhancing context retention. As a result, LLMs will be able to execute increasingly difficult coding jobs with less need for human involvement, increasing overall efficiency and productivity of software development processes.

### 5.2. Collaborative Approaches

A coordinated strategy combining domain-specific knowledge, computer science, linguistics, and cognitive science expertise is needed to address the context length concerns in LLMs. Multidisciplinary research can help us comprehend problems more thoroughly and come up with better solutions. For the purpose of promoting innovation and raising the bar for LLMs for code production, cooperation between academic institutions, businesses, and research centres is essential.

Numerous cooperative efforts and projects are currently in progress with the goal of improving LLM capabilities and addressing their shortcomings. For instance, the OpenAI Codex effort and partnerships between top tech firms and universities are aimed at enhancing LLMs' context-handling capabilities. These programs are crucial for combining resources, exchanging information, and quickening the development of advanced AI technologies.

## 6. Conclusion

This paper has explored the challenges posed by context length issues in LLMs for code generation, highlighting the impact on accuracy, coherence, and overall performance. We have discussed various technical solutions, including memory-augmented networks, hierarchical models, sliding window mechanisms, and hybrid approaches, all aimed at enhancing context retention.

Resolving context length concerns is essential to maximizing LLMs' potential in code creation. Enhancing context preservation will guarantee more correct, consistent, and dependable code that is generated, which will ultimately increase software development processes' productivity and efficiency.

---

To overcome the shortcomings of the current LLMs and advance the sector, more research and development are needed. To overcome context length difficulties and improve the capabilities of LLMs in code creation, we encourage academics, developers, and companies to make investments in innovation, teamwork, and interdisciplinary research.

### References

1. Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., & Liang, P. (2024). Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12, 157-173.
2. An, C., Gong, S., Zhong, M., Zhao, X., Li, M., Zhang, J., ... & Qiu, X. (2023). L-eval: Instituting standardized evaluation for long context language models. *arXiv preprint arXiv:2307.11088*.
3. Bai, Y., Lv, X., Zhang, J., Lyu, H., Tang, J., Huang, Z., ... & Li, J. (2023). Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
4. Chen, S., Wong, S., Chen, L., & Tian, Y. (2023). Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*.
5. Ding, Y., Zhang, L. L., Zhang, C., Xu, Y., Shang, N., Xu, J., ... & Yang, M. (2024). Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753*.
6. Li, D., Shao, R., Xie, A., Sheng, Y., Zheng, L., Gonzalez, J., ... & Zhang, H. (2023). How Long Can Context Length of Open-Source LLMs truly Promise?. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.
7. Peng, B., Quesnelle, J., Fan, H., & Shippole, E. (2023). Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*.

**Copyright:** ©2024 Ashrey Ignise, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.