Shawn Chumbar
CMPE 258: Deep Learning
09 April 2024

# Short Story on AIOS: LLM Agent Operating System

# INTRODUCTION

In this research paper, the authors introduce a new architecture called AIOS that integrates large language models (LLMs) into operating systems. The primary objective of AIOS is to create an intelligent operating system that efficiently manages and facilitates the execution of LLM-powered agents. By embedding LLMs at the operating system level, AIOS aims to address the challenges associated with the complexity and quantity of intelligent agents. The authors use an LLM-specific kernel to segregate the OS duties, related resources, and development kits. Through this segregation, the LLM kernel tries to improve the management and coordination of LLM-related activities.

# ARCHITECTURE

The proposed architecture is designed with three distinct layers:

1. Application Layer
2. Kernel Layer
3. Hardware Layer

### *Application Layer*

The application layer focuses on agent development using the AIOS SDK. This layer allows developers to dedicate their focus to the essential logic and functionalities of their agents by allowing development of agent applications that abstract away the complexities of lower-level system functions.

### *Hardware Layer*

The hardware layer consists of physical components, like CPU, GPU, memory, disk, etc. It is important to note that the hardware layer allows the LLM kernel to leverage hardware resources without requiring direct hardware management, which allows the system to maintain integrity and efficiency.

The kernel layer is split into two components:

1. OS kernel for non-LLM tasks
2. LLM kernel for LLM-specific operations.

      The OS kernel serves the requirements of non-LLM specific operations, while the LLM kernel serves the requirements of LLM-specific operations. This allows the LLM kernel to focus on LLM specific tasks (i.e. context management and agent scheduling). The LLM kernel is composed of multiple modules such as the agent scheduler, context manager, memory manager, storage manager, tool manager, and access manager. These modules allow the system to optimize resource allocation, enable context switching, support concurrent agent execution, provide tool services, and enforce access control policies.

# LLM KERNEL MODULES

The LLM Kernel Modules include the following modules:
1. **Agent Scheduler**: Prioritizes and schedules agent requests to optimize LLM utilization.
2. **Context Manager**: supports snapshot and restore the intermediate generation status in LLM and content window management of LLM.
3. **Memory Manager**: Provides short-term memory for each agent's interaction logs.
4. **Storage Manager**: Persists agent interaction logs to long-term storage for future retrieval.
5. **Tool Manager**: Manages agent's calling of external API tools (e.v., search, scientific computing, etc.)
6. **Access Manager**: Enforces privacy and access control policies between agents.

## *Agent Scheduler*

      The Agent Scheduler manages the agent requests in an efficient way. The agent scheduler uses strategies such as First-In-First-Out (FIFO), Round Robin (RR), and other scheduling algorithms to optimize the agent requests. The scheduler handles balancing waiting time and turnaround time of each agent as tasks from different agents are executed in parallel.

*Context Manager*

The context manager handles managing the context that is provided to LLM and the generation process given certain context. The two main critical functions for this manager are the context snapshot & restoration, and context window management. Context snapshot and restoration is the mechanism employed by the context manager module to handle situations where the execution of an LLM-based agent is interrupted or suspended before completing its response generation. The context restoration would work in reverse, allowing us to resume a suspended agent's execution.

An example of where the context snapshot would come into play is given below. When an agent's execution is about to be suspended, the context manager takes a snapshot of the current state of the LLM's generation process. This snapshot would capture all the necessary information such as the current token, the generated tokens, the probabilities of different possible next tokens, and any other relevant information regarding the internal states of the LLM. The snapshot would then be saved and associated with a specific agent.

The Context Restoration would then resume the execution of the suspended agent. The context manager retrieves the previously saved snapshot associated with that agent, and the LLM's state is restored based on the snapshot, essentially setting it back to the exact point where the generation process was interrupted. The response generation is then continued from where it left off.

The second critical function of the context manager, the context window management, supports basic text summarization and incorporates other expansion techniques to manage the context window.

*Memory Manager*

The next LLM module, the memory manager, is used to manage short term memory within an agent's lifecycle. This module ensures that the data is stored and accessible only when the agent is active. This module also supports the storing of each agent's memory independently.

*Storage Manager*

The storage manager allows for the long-term preservation of data. The medium of this data includes local files, databases, or cloud-based storage. The storage manager supports retrieval augmentation and enriching the agent knowledge update and enhancing long-term user experience.

*Tool Manager*

The next LLM module, the Tool Manager, manages a diverse array of API tools that improves the functionality of the LLMs. This can include giving LLMs access to different APIs, such as system calls, context system calls, memory syscalls, and storage syscalls. An example of possible API calls would be web search, scientific computing, database retrieval, and image processing.

*Access Manager*

The Access manager handles access control operations among different agents by administering privilege groups for each agent. This allows only certain agents access to certain resources. Additionally, the access manager also maintains audit logs that capture information about access requests, agent activities, and modifications made to the access control parameters. This allows the access manager to protect against potential privilege attacks.

# LLM SYSTEM CALL INTERFACE

The LLM system call interface within the LLM kernel offers basic LLM call operation functions that act as a bridge between complex agent requests and the execution of different kernel's modules. LLM system calls offer a suite of basic functions that span across the multiple kernel modules (i.e. agent manager, content handler, memory and storage managers, access manager).

# AIOS SDK

The AIOS SDK is a versatile toolkit for creating agent applications within the AIOS. It offers various functionalities, from initialization of agents and management of agent lifecycles to resource management and generating plans for agent tasks. The AIOS SDK aims to provide developers with the ability to create tools to harness the full potential of their agent applications within the AIOS framework.

# FUTURE WORK

In regards to future work, there are many areas that can offer significant improvements. For starters, future research should focus on advanced scheduling algorithms that perform dependency analysis among agent requests to optimize the allocation of resources. Another area of improvement is context management efficiency. More efficient mechanisms can be used to assist context management (i.e. time-efficient content management techniques that make the context snapshot and restoration faster). Additionally, there could be various compression

algorithms that could be used to create space-efficient solutions for content snapshotting. The last main area that could use an improvement is the safety and privacy portion of AIOS. Due to the various LLM agents being used, it is possible for malicious attacks to occur, such as jailbreaking or accessing restricted resources. The exploration of advanced encryption techniques will help safeguard data within AIOS, and can maintain the confidentiality of agent communications.

# PERSONAL ANALYSIS

AIOS represents a significant step forward in the realm of intelligent agent management and execution. By integrating LLMs at the operating system level, AIOS tackles the challenges associated with the growing complexity and scale of LLM-based agents. The modular architecture and well-defined layers provide a robust foundation for efficient agent development, deployment, and execution.

The introduction of the LLM kernel is particularly noteworthy, as it segregates LLM-specific tasks from general OS operations. This separation allows for focused optimization and management of LLM-related activities, ensuring that the unique requirements of LLM-based agents are met.

The agent scheduler and context manager are two standout features that demonstrate the thoughtfulness and innovation behind AIOS. The scheduler's ability to prioritize and manage agent requests ensures optimal LLM utilization, while the context manager's snapshotting and restoration capabilities enable seamless agent execution even in the face of interruptions.

The integration of tool management and access control further enhances the practicality and security of AIOS. By providing agents with access to external APIs and enforcing strict access control policies, AIOS creates a secure and extensible environment for agent operation.

Overall, AIOS represents a significant milestone in the evolution of operating systems and the integration of AI agents. As the complexity and prevalence of LLM-based agents continue to grow, AIOS provides a comprehensive framework for efficient management and execution. The modular architecture, intelligent scheduling, context management, and security features make AIOS a promising platform for the future of intelligent agent computing.

While the research presents a solid foundation, there is still room for further exploration and refinement. Future work could focus on advanced scheduling algorithms, more efficient context management techniques, and enhanced security measures. As AIOS continues to evolve and mature, it has the potential to revolutionize the way we develop, deploy, and interact with intelligent agents, paving the way for a new era of AI-powered computing. It will be very

interesting to see a world in which malicious actors will be attempting to jailbreak an AI Operating System. This may even lead to various agents attempting to perform attacks on AI operating systems, and the AI operating systems fighting back.

## Works Cited

Mei, K., Li, Z., Xu, S., Ye, R., Ge, Y., & Zhang, Y. (2024). AIOS: LLM Agent Operating System. Retrieved from arXiv preprint server (ID: 2403.16971).

WorldofAI. "AIOS: LLM Agent Operating System! Create Software, Automate Tasks, and More!" YouTube, 10 Apr. 2024, https://www.youtube.com/watch?v=7WVP0yWVrLI&ab_channel=WorldofAI.