# LLM as OS, Agents as Apps: Envisioning AIOS, Agents and the AIOS-Agent Ecosystem

**Yingqiang Ge**
Rutgers University

**Yujie Ren**
Rutgers University

**Wenyue Hua**
Rutgers University

**Shuyuan Xu**
Rutgers University

**Juntao Tan**
Rutgers University

**Yongfeng Zhang**[*]
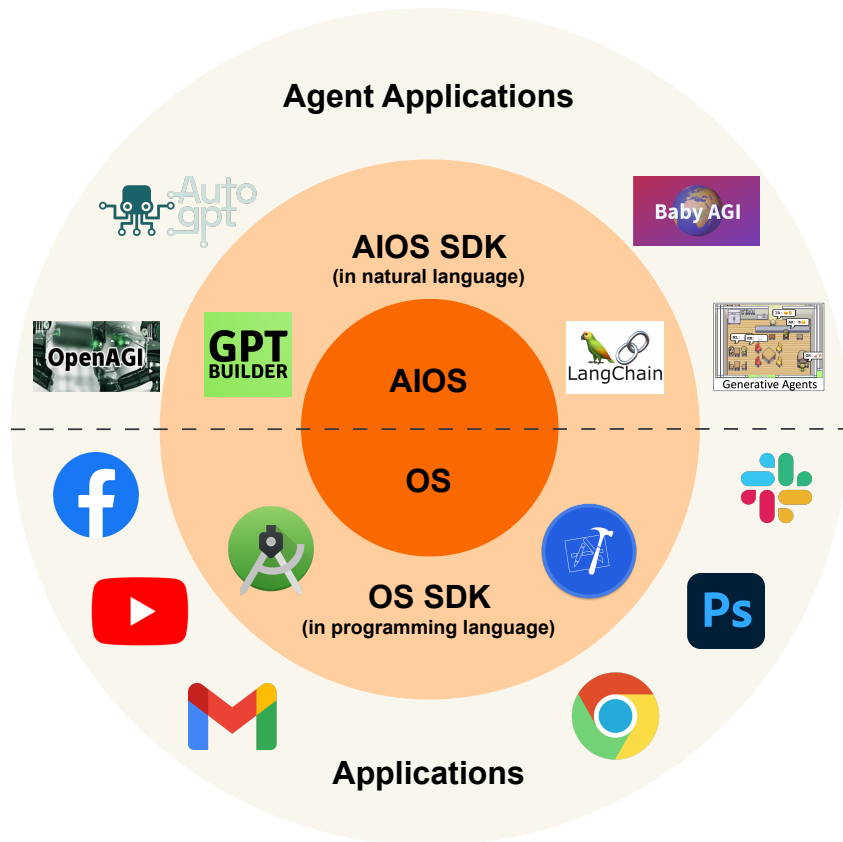Rutgers University

Figure 1: OS-APP ecosystem vs. AIOS-Agent ecosystem.

## Abstract

This paper envisions a revolutionary AIOS-Agent ecosystem, where Large Language Model (LLM) serves as the (Artificial) Intelligent Operating System (IOS, or AIOS)–an operating system "with soul". Upon this foundation, a diverse range of LLM-based AI Agent Applications (Agents, or AAPs) are developed, enriching the AIOS-Agent ecosystem and signaling a paradigm shift from the traditional OS-APP ecosystem. We envision that LLM's impact will not be limited to the AI application level, instead, it will in turn revolutionize the design and implementation

[*]**Author Affiliation**: Department of Computer Science, Rutgers University, New Brunswick, NJ, 08854, US;
**Author Emails**: {yingqiang.ge,yujie.ren,wenyue.hua,shuyuan.xu,juntao.tan,yongfeng.zhang}@rutgers.edu

of computer system, architecture, software, and programming language, featured by several main concepts: LLM as OS (system-level), Agents as Applications (application-level), Natural Language as Programming Interface (user-level), and Tools as Devices/Libraries (hardware/middleware-level).

In this paper, we begin by introducing the architecture and historical evolution of traditional Operating Systems (OS). Then we formalize a conceptual framework for AIOS through "LLM as OS (LLMOS),"[2] drawing analogies between AIOS components and traditional OS elements: LLM is likened to OS kernel, context window to memory, external storage to file system, hardware tools to peripheral devices, software tools to programming libraries, and user prompts to user commands. Subsequently, we introduce the new AIOS-Agent Ecosystem, where users and developers can easily program Agent Applications (AAPs) using natural language, democratizing the development of and the access to computer software, which is different from the traditional OS-APP ecosystem, where desktop or mobile applications (APPs) have to be programmed by well-trained software developers using professional programming languages. Following this, we explore the diverse scope of Agent Applications. These agents can autonomously perform diverse tasks, showcasing intelligent task-solving ability in various scenarios. We delve into both single agent systems and multi-agent systems, as well as human-agent interaction. Lastly, we posit that the AIOS-Agent ecosystem can gain invaluable insights from the development trajectory of the traditional OS-APP ecosystem. Drawing on these insights, we propose a strategic roadmap for the evolution of the AIOS-Agent ecosystem. This roadmap is designed to guide the future research and development, suggesting systematic progresses of AIOS and its Agent applications.

## Contents

---

[2]For convenience, LLMOS may be pronounced as "el-mos".

# 1 Introduction

In the evolving landscape of information technology, Operating Systems (OS) such as Windows[3], MacOS[4], iOS[5], and Android[6] have become cornerstones of our digital lives. On top of the operating systems, a diverse range of applications (APPs) are developed, helping with users' diverse tasks and enriching the OS-APP ecosystem. For example, Microsoft Word[7] and Google Docs[8] excel in drafting documents, while Microsoft Outlook[9] and Gmail[10] offer efficient email management.

Operating systems have advanced significantly, becoming more intuitive and user-friendly, yet their core remains rooted in static rules and predefined logic flows, without the intelligent, creative, and emergent task-solving abilities. The applications built on top of such OS, on the other hand, are also limited to their designed purposes, unable to transcend beyond their individual scopes. Whenever individual applications need to incorporate intelligent abilities, they have to implement their own AI methods or functionalities, sometimes based on third-party libraries. This isolated framework underscores a significant shortfall in the current OS-APP ecosystem and highlights the pressing need of infusing (artificial) intelligence into operating systems, so that intelligence can be natively distributed to the various applications built on top of it.

As a result, this paper envisions (Artificial) Intelligent Operating System (IOS, or AIOS), an operating system "with soul". Furthermore, a diverse scope of intelligent Agent applications are built on top of the AIOS, leading to the new AIOS-Agent ecosystem, in comparison to the traditional OS-APP ecosystem, as shown in Figure 1. Due to the versatile and remarkable capabilities they demonstrate, Large Language Models (LLMs) (Radford et al., 2019; Brown et al., 2020; Touvron et al., 2023; Taori et al., 2023) are regarded as potential sparks for Artificial General Intelligence (AGI) (Bubeck et al., 2023; Morris et al., 2023; Ge et al., 2023), offering hope as foundational elements for the development of AIOS. There are several reasons confirming LLMs' general capability and feasibility for building AIOS:

- First, LLMs have demonstrated exceptional language understanding abilities as well as reasoning/planning abilities to solve complex tasks, which can divide the tasks into several sub-tasks and conquer them one-by-one, sometimes with the assistance of external tools (Ge et al., 2023; Wei et al., 2022; Huang and Chang, 2022).

- Second, LLMs offer a highly flexible platform to process virtually any prompt, instruction, or query expressed in natural language, making it possible for a diverse range of Software Development Kits (SDKs) and/or applications to be built on top of them.

- Third, LLMs offer a more intuitive and user-friendly interface, since they can understand and respond to user prompts or instructions in natural language (Brown et al., 2020; Touvron et al., 2023). This sheds light on the future where natural language serves as the programming language, making technology more accessible, especially for those who may not be familiar with traditional computer interfaces and programming languages.

- Fourth, LLMs can be programmed to learn from interactions and customize their responses based on user preferences and past interactions, providing a more personalized experience (Safdari et al., 2023; Durmus et al., 2023).

As a result, infusing intelligence into the OS-level through LLM makes it possible for easily distributing intelligent abilities into the application-level, providing a promising way to democratize intelligence across various applications.

Inspired by the architecture of traditional OS (as shown in Figure. 2a), we present a general framework for **LLM as OS (LLMOS)** with several key components in Section 3 (as shown in Figure. 2b): LLM as Kernel, Context Window as Memory, External Storage as File, Tools as Devices/Libraries, User

---

[3]https://www.microsoft.com/en-us/windows/

[4]https://www.apple.com/macos/

[5]https://www.apple.com/ios/

[6]https://www.android.com/

[7]https://www.microsoft.com/en-us/microsoft-365/word/

[8]https://www.google.com/docs/about/

[9]https://www.microsoft.com/en-us/microsoft-365/outlook/
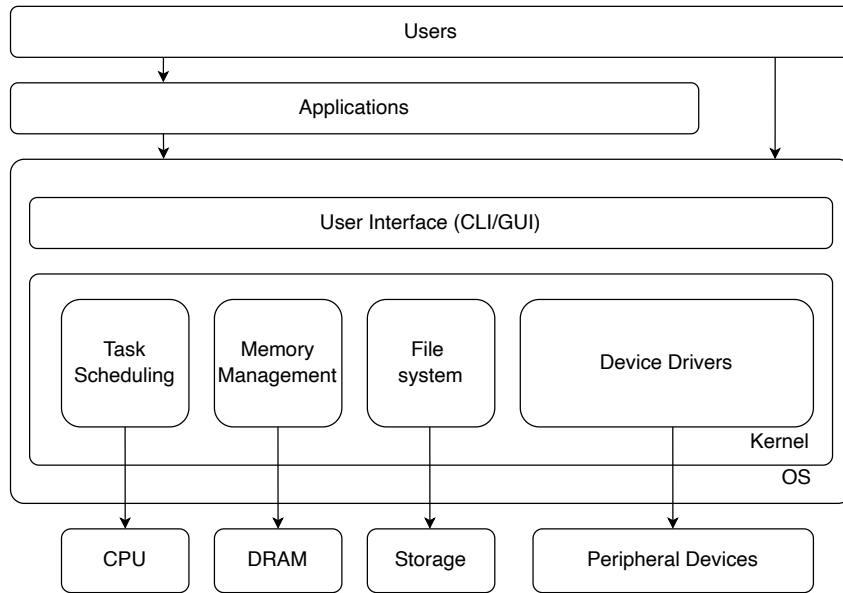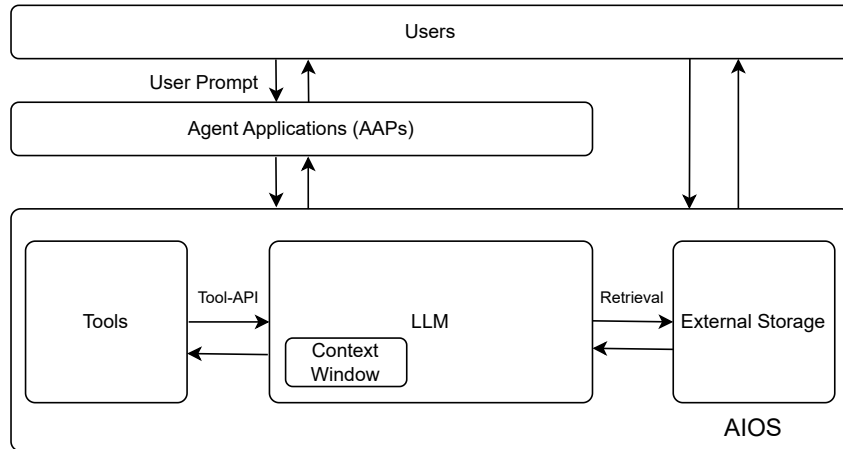
[10]https://www.google.com/gmail/about/

(a) Architecture of Operating System (OS).



(b) Architecture of Large Language Model as OS (LLMOS) for AIOS.

Figure 2: Illustrations of the architectures of OS and AIOS (LLMOS).

Prompt/Instruction as User Interface (UI), and Agents as Applications. The conceptual framework presented draws an analogy between an LLM as OS (LLMOS) and a traditional Operating System (OS), mapping various components of the LLMOS to elements of an OS. The LLM itself is likened to the kernel, the central part managing the system's core functions. The LLM's context window is compared to the memory of an OS, handling immediate context selection and data processing (Shi et al., 2023). External storage for the LLM is analogous to the files of an OS, allowing for long-term data storage. Meanwhile, there are corresponding data retrieval methods to enable retrieval-augmented LLMs (Guu et al., 2020), which serve as the file system of OS to manage and find relevant files. Besides, LLM can make use of various tools for task solving (Ge et al., 2023), including both hardware tools and software tools. Hardware tools in the LLMOS framework are equated to peripheral devices in a traditional OS, each offering specific functionalities to help interact with the physical world, while software tools in the LLMOS framework serve as the programming libraries in traditional OS, enabling Agent applications to interact with the virtual/digital world. User prompts or instructions for the LLM are akin to the user interface (UI) in a traditional OS, facilitating interaction between the user and the system. The user prompts or instructions can be direct natural language instructions provided by the user, and they may also be (sometimes semi-structured) natural language instructions converted from users' non-natural language instructions such as clicking on icons. This mapping provides a systematic way to understand the operational similarities between LLMOS-based AIOS and traditional OS.

Upon establishing a robust conceptual framework for LLMOS-based AIOS, we introduce the AIOS-Agent Ecosystem, akin to the traditional OS-APP Ecosystem, in Section 4. We begin by introducing the concept of Agent Applications (AAPs) within LLMOS, analogous to traditional applications (APPs) based on an operating system. These AAPs represent a diverse scope of specialized tasks for users to execute based on LLMOS. As illustrated in Figure 4, by integrating the LLMOS layer with the OS layer, Hardware layer, and the Agent Application layer, we can construct an autonomous AI Agent system. Such AI Agent system responds to user prompts or instructions in natural language and is capable of performing a multitude of tasks based on its interaction with the physical or digital environment. Since a diverse scope of Agents can be developed on top of the shared AIOS foundation, this eventually leads to an AIOS-Agent ecosystem. Moreover, in the new AIOS-Agent Ecosystem, the users and developers can easily program Agent Applications (AAPs) using natural language, democratizing the development of and the access to computer software, which is different from the traditional OS-APP ecosystem, where desktop or mobile applications (APPs) have to be programmed by well-trained software developers using professional programming languages.

Following this, we further delve into the practical LLMOS-based Agent Applications in Section 5. This section explores the potential of enhancing LLMOS's functionality by developing various agents in the real world, and further investigates the dynamic interactions between multiple agents and humans. LLMOS-based agents are characterized by their creative autonomy, enabling them to generate novel ideas, narratives, or solutions not pre-programmed into them (Chase, 2022; Gravitas, 2023; Ge et al., 2023; Li et al., 2023; Yao et al., 2022b,a), which is indicative of an advanced level of creative intelligence (Yuan et al., 2022; Franceschelli and Musolesi, 2023). Specifically, Section 5.1 discusses applications of single agents, Section 5.2 explores multi-agent systems, and Section 5.3 focuses on human-agent interactions.

Finally, in Section 6, we explore several crucial future research directions for LLMOS and AIOS in general, drawing parallels and learning from the evolution of traditional operating systems. These directions span a wide array of areas, aiming to enhance the capabilities and applications of LLMOS: 1) **Resource Management.** For example, OS employs virtual and shared memories to address the issue of limited physical memory. LLMOS can inspire from these ideas to mitigate its own problem of limited context window challenges; 2) **Communication.** Different OSes and applications communicate using standardized protocols (such as Domain-Specific Languages); LLMOSes and Agents can build and use similar standard protocols for exchanging data and instructions with various systems, ensuring compatibility and smooth interaction across diverse platforms; 3) **Security.** Security vulnerabilities in OS are important issues. State-of-the-art approaches aim to detect and capture malware and viruses at various levels. Similarly, LLMOS can implement detection and intervention mechanisms to regulate and monitor the implementation of third-party tools and Agent Applications.

# 2 Aligning LLM and OS

## 2.1 OS and Connections with LLM

The von Neaumann architecture, laying the foundation of modern computer hardware system, manipulate electrons and gates in the binary world, whereas human beings communicate with natural languages. Such gigantic semantic gaps between human users and computer hardware motivates an intermediary software layer interacting with users with a protected and abstracted view of underlying hardware resources such as CPU (Central Processing Unit), GPU (Graphics Processing Unit), RAM (Random Access Memory), storage and various other devices, which is called the Operating System (OS). The modern operating systems, over the past few decades, have evolved in a multi-layered architecture with modular components in each layer. This design not only enhances the efficiency and functionality of the systems but also facilitates easier management, scalability, and integration of diverse hardware and software elements.

### 2.1.1 Kernel

The kernel, as its name suggests, encapsulates a set of core functionalities of managing hardware resources (e.g., CPU, GPU, DRAM, storage, and devices) as a nutshell.[11]

---

[11]We limit our scope to traditional monolithic kernel design.

- **CPU Management** (Process/Thread, scheduling). To manage the CPU resources, modern operating systems abstract the execution of user programs or applications on a physical CPU as processes or threads. When the user launches an application, the operating system kernel loads the executable binary files of this application or program into DRAM, create the necessary data structures to book-keep any running states for this program, and allocate necessary resources.

  However, the power wall (Agarwal et al., 2006) limits the number of physical CPUs to be integrated into a single chip. To provide users with the illusion that they possess physical CPUs without sharing with others, modern operating systems multiplex the running processes or threads with limited number of CPUs with time-sharing (Ritchie and Thompson, 1974) and more dedicated policies (Molnár, 2007; Stoica and Abdel-Wahab, 1995), much like multiple users share the same LLM backbone when processing their prompts or instructions. In such cases, if the inputs from multiple users are beyond the capacity of the resources, such as the tools, then the LLM may perform scheduling of the user prompts.

- **Memory Management.** The physical memory in a computer, also known as DRAM (Dynamic Random Access Memory), is the key component to store the instructions and data of both the OS and applications. The OS is responsible for managing and allocating free space in the physical memory from application's requests.

  As described above, the physical memory, depicted as "dynamic," cannot store data persistently when power is off, as it needs to refresh periodically to prevent the loss of data. In addition to its volatile nature, the evolution of memory falls behind the CPU for a long time, which is known for "The memory wall". This fact lays in two orthogonal aspects. First, the data transfer rate between DRAM and CPU can no longer catch up with the speed that CPU processes data. Second, the memory capacity on a single node stops scaling. Although the emergence of Compute eXpress Link (CXL, 2023) alleviates the memory capacity wall, it still cannot keep the rapid pace of data growth at the artificial intelligence era.

  This is much like the context window of LLMs, which is usually limited by the maximum number of tokens that the LLM can handle. Besides, LLM usually needs to select the relevant information from the input context, since not all contexts are relevant for the current task and LLM may be easily distracted by irrelevant contexts (Shi et al., 2023). Context selection can be either implicitly realized through attention mechanisms, or be explicitly implemented by selecting relevant segments from the input context, much like the memory management process by traditional OS.

- **Storage Management.** Storage devices, which store data persistently, provide much more density than memory with fewer cost, but much slower. The operating systems abstract the storage as file, and organize files in a systematic way with a component called the "file system". The file system contains the meta-data to index the actual data stored on the storage media.

  In addition to the file abstract, modern operating systems reserve a small portion of storage as an extension of memory, which is called the "swap area". It is the operating system that tracks the hotness and coldness of the data from user applications, and swaps code data from physical memory to the swap area on the storage devices.

  Similarly, LLM often has access to external data storage for retrieval-augmented language modeling (Guu et al., 2020). These external data can be free text data, structured tabular data, semi-structured graph data, or others. Furthermore, the external data are usually properly indexed for efficient and accurate retrieval, much like the storage management process of traditional operating systems.

- **Device Management.** Peripheral devices, often excluded from the core computing system of CPU and DRAM, form an important set of functionalities for user input and output. Those devices range from mouse, keyboard, to GPU and network interface cards (NIC). The operating systems take the responsibility to manage those devices to orchestrate with other core components in the OS kernel.

  There are thousands of different devices for different purposes and from different vendors all around the world. Hence, it is not possible to implement the driver program for all of those devices. Instead, modern operating systems expose a generic interface as device driver APIs for device vendors, and thus shifts the responsibility of developing device driver programs from OS maintainers to vendors of those devices. To provide a 'plug-and-play' feature, modern operating systems such as Linux include some universal and necessary device drivers for some common devices; for example, the drivers for GPU and USB devices (Corbet et al., 2005).

  Similarly, large language models are not only text-in-text-out models, instead, they have the ability of leveraging various tools for solving complex tasks (Schick et al., 2023; Ge et al., 2023). There can be two types of tools, hardware tools and software tools, which help the LLM to interact with

the physical world and the digital world, respectively. In this context, the hardware tools for LLM are similar to the devices for traditional OS. Furthermore, just like driver programs connect devices with OS, Tool-Drivers can be developed to connect LLM and hardware tools, so that LLM can easily leverage the tools for task-solving. We will discuss software tools in the following.

- **SDK and Programming Libraries.** The SDK (Software Development Kit) and programming libraries of an operating system are crucial tools that enable developers to easily create applications. They serve as the backbone of application development for operating systems, not only enabling and simplifying the creation of applications, but also ensuring that these applications are secure, efficient and compatible with the OS, which significantly contribute to the vitality and growth of the OS-APP ecosystem.

  Similarly, LLM and the various AI Agents built on top of it can make use of various software tools such as searching on the web, checking for weather conditions, and booking for flight tickets (Ge et al., 2023). These software tools serve as reusable functionalities that can be leveraged by LLM and Agents for complex task-solving, and they can be provided as SDK or programming libraries so that users or developers can easily use them.

### 2.1.2 User Interface

As detailed previously, the kernel manages the hardware resources with proper abstraction for user to utilize. In order to enhance the accessibility of virtualized hardware resources, it is crucial to establish an interface between user and OS.

- **System Call.** The system call, as the channel between OS kernel and users, defines a set of core functions to allocate and use the virtualized hardware resources. For instance, in a POSIX-compliant operating system (IEEE and Group, 2018), the `mmap` system call allocates and manipulates memory resources. The `fork` and `exec` system call family deals with process and thread creation. The `read` and `write` system call are used to interact with storage devices. In terms of LLMs, the system calls can be formulated as natural language prompts into instruct the LLM for task execution.

- **Command-line Interface (CLI).** The command-line interface defines a set of utility programs built on top of system calls to facilitate users to operate on the computer hardware in an interactive manner. Users interact with the OS in those utility programs as commands. For instance, the `cd` and `ls` commands implement the action of entering a directory and list all the files and folders in a directory. In the context of LLM, natural language prompts naturally serve as the interface for users to interact with LLMs. Furthermore, the LLM may also pre-define some foundational and commonly used functionalities as standard prompt templates for users to use, similar to the standard commands as in traditional OS.

- **Graphic User Interface (GUI).** The graphic user interface is a visual way for users to interact with computers and electronic devices using graphical elements such as icons, buttons, windows, and menus, as opposed to a text-based interface such as a command-line interface. GUIs make it easier for users to interact with complex systems by representing actions through visual elements and providing a more intuitive user experience, especially with the increasing demand and use of mobile devices such as smart phones discussed in section 2.1.4. In terms of LLMs, graphic user interfaces can also be developed for LLM and Agents so that users can more conveniently interact with them without the need to writing long prompts. Instead, these GUIs will convert user's non-language instructions (such as clicking on icons) into (sometimes semi-structured) natural language prompts based on pre-defined prompt templates, and these converted natural language prompts will be sent to LLM for executing the user instruction.

### 2.1.3 Operating System Ecosystem

The operating system ecosystem functions as an extension of the operating system, providing a comprehensive set of developer tool-kits (OS SDK) and runtime libraries, shown in Figure 1. These tools empower application developers to efficiently design, implement, and run their applications within the operating system environment. For instance, the well-known iOS ecosystem includes a dedicated application development toolkit known as Xcode[12], alongside an application publishing platform called the AppStore[13], complementing the core iOS ecosystem. In this ecosystem, the OS

---

[12]https://developer.apple.com/xcode/
[13]https://www.apple.com/app-store/

provides a bunch of resources to support APP development and also services as the platform for deploying and hosting these APPs, which eventually leads to a prospering OS-APP ecosystem.

Similarly, we envision an AIOS-Agent ecosystem, where LLM serves as the operating system and hosts a diverse range of AI Agent applications, as shown in Figure 1. The LLM as OS (LLMOS) environment shall also provide a comprehensive set of AIOS SDKs and/or libraries, predominately supporting programming in natural language, to help developers or even average users without any knowledge on professional programming languages, to easily develop and deploy Agent applications in the LLMOS-based AIOS-Agent ecosystem.

### 2.1.4 Evolution History of Operating Systems

- **Batch Processing System.** Early batch processing systems were a fundamental aspect of the early days of computing, dating back to 1950s (UW:CSE451, 2023). These systems were characterized by a sequential execution of tasks, where jobs were submitted in batches for processing. Early batch processing systems laid the groundwork for subsequent developments in operating systems. While lacking the interactivity and responsiveness of modern systems, they played a crucial role in advancing computing capabilities and setting the stage for more interactive and user-friendly computing environments in the years to come.

- **Time Sharing.** Time-sharing systems (UW:CSE451, 2023), as proposed in *Multics* (Corbató et al., 1971) represent a significant advancement in the history of operating systems, providing a departure from traditional batch processing systems and introducing the concept of shared, interactive computing. Many concepts introduced in time-sharing systems, such as multitasking, interactive interfaces, and dynamic resource allocation, have become integral parts of modern operating systems, which laid the groundwork for user-friendly computing environments, enabling efficient resource utilization and interactive computing.

- **Multitasking.** As the hardware evolved to multiple cores, planning user tasks on available multi-core CPU is critical to maximize the CPU utiliztion. Multitasking involves scheduling processes to run on the CPU in a way that gives the appearance of concurrent execution. Process scheduling algorithms determine the order in which processes are executed. The UNIX operating system (Ritchie and Thompson, 1974), developed in the late 1960s and early 1970s at Bell Labs, introduced the concept of processes, each with its own address space, and implemented a simple and efficient multitasking model.

- **Visualization (GUI).** As described previously in section 2.1.2, the command-line interface used to be the narrow bridge between users to interact with OS. Since command lines are highly professional, it prevented broader groups of users from easily and efficiently operating the computers. From Xerox Alto introduced by Palo Alto Research Center (PARC) in 1973[14], to Apple Macintosh introduced in 1984[15], to Microsoft Windows[16] introduced in the early 1990s, and to a wide range of open-source GUIs for Linux such as GNOME[17], KDE[18] and UNITY[19], the development of GUIs has significantly influenced the accessibility and usability of computers, making computing more intuitive and user-friendly.

- **Cloud Computing.** Performing data computing and storage on a single computing node was no longer sufficient as the data scale grew significantly in the early 2010s. The development of client-server architecture in the early 1990s, where multiple clients connect to centralized servers, laid the foundation for distributed computing with the support of networked environments and the communication between clients and servers that emerged as OS core functionalities. More technologies such as virtualization and resource disaggregation empower the modern cloud computing community and market. Though cloud computing is not directly tied to the history of operating systems, the evolution of cloud computing has had a profound impact on how operating systems are designed, deployed, and utilized.

- **Mobile and Embedded Systems.** Aligning with the Moore's law, the size and the computing power of a general-purpose CPU redefines the capability of modern embedded devices. Tailored

---

[14]https://spectrum.ieee.org/xerox-alto

[15]http://apple-history.com/128k

[16]https://winworldpc.com/product/windows-3/31

[17]https://www.gnome.org/

[18]https://www.kde.org/

[19]https://unityd.org/

9

for operating on a relatively constrained computing resource, modern mobile OS are specifically focused on power efficiency, network efficiency, optimized graphics for touch screen, and the prospering application ecosystem (detailed in section 2.1.3). Similarly, embedded devices ranging from IoT to robotics pose the aspect of real-time responses and more strict resource management in OS, which has been reflected in many successful embedded OSes such as VxWorks[20] and QNX[21].

- **AI-powered Features.** In the past decade, artificial intelligence (AI) technologies have experienced explosive growth. Specifically, breakthroughs in several areas of AI, such as natural language processing, computer vision, and speech recognition, are extending the interactive interface between users and OS to a new level. As examples of early implementations, Siri[22] from Apple and Cortana[23] from Microsoft, the AI-powered virtual assistants provide a rich set of capabilities to understand the requests from the users, such as message sending, call making, question answering, web search, recommendation, and controlling smart home devices.

## 2.2 AIOS, LLMOS and AI Agents

Recent advances in foundation models, such as Large Language Models (LLMs), have significantly changed how AI applications are designed, trained, and used, including but not limited to NLP, CV, Search, RecSys, Multimedia, and Game applications. We envision that the influence of LLMs will not be limited to the application level, instead, it will in turn revolutionize the design and implementation of computer system, architecture, software, and the system-application ecosystem. This is realized through the following key concepts.

### 2.2.1 LLM as OS (system-level)

LLM serves as the foundation AIOS platform that provides intelligent computing, APIs, and services that support various applications and manage various computing resources; Different from traditional OS, which aims at precision and efficiency, AIOS is characterized by its "intelligence," which enables its intelligent and creative interaction with various applications and computing resources for task solving, leading to an operating system "with soul".

### 2.2.2 Agents as Applications (application-level)

Based on the LLM-driven AIOS, various AI Agents can be developed as AIOS-native applications, such as trip planning agent, medical consulting agent, financial management agent, etc.; these agents make use of the intelligent computing power of the LLM and the various tools provided by the AIOS SDK to solve various tasks; except for single-agent applications, agents may also communicate and interact with each other, enabling multi-agent applications; agent can even be created when needed and released after use, enabling on-the-fly creation of agents. There are several reasons that many specialized agents will be developed instead of integrating all of the functionality into one LLM: 1) the tasks that may be initiated by users are diverse, unbounded and unable to predetermine, 2) tasks may require specialized tools, reasoning and planning abilities, accessing the external physical or digital world, or domain-specific knowledge to complete, which the LLM-based AIOS platform may not support and need to be developed at the agent-level, 3) though language is a very powerful medium for describing tasks, objects, information or ideas, it may not be able to describe everything, and the completion of certain tasks may need creative forms of interaction beyond the "language input, language output" paradigm, such as the processing of vision, sound, touch, smell, and the diverse types of signals from various sensors. The processing of these unbounded and unpredictable types of information needs to be handled at the agent-level instead of the AIOS-level.

### 2.2.3 Natural Language as Programming Interface (user-level)

In the AIOS-Agent ecosystem, average users can easily program Agent Applications (AAPs) using natural language, democratizing the development of and the access to computer software, which is very different from traditional OS-APP ecosystem, where desktop or mobile applications (APPs) have to be programmed by well-trained software developers using professional programming languages.

---

[20]https://www.windriver.com/products/vxworks
[21]https://blackberry.qnx.com/en
[22]https://www.apple.com/siri/
[23]https://www.microsoft.com/en-us/cortana

This trend is consistent with the evolving history of programming languages, which are becoming more and more user-friendly, beginning from binary codes to assembly language to various high-level languages such as C, C++, Java and Python; Natural Language as Programming Language is a natural extension of this evolving history, making it possible for average users to program agent applications and interact with computers without special training on professional programming languages.

### 2.2.4 Tools as Devices/Libraries (hardware/middleware-level)

Just like traditional OS can take advantage of various input and output devices such as keyboards and printers to support various APPs for task fulfillment, AIOS can also provide various tools as services to support the agents' task fulfillment. Such tools include both hardware tools such as devices and software tools such as plugins or libraries; the tools also include both input tools such as collecting signals from a sensor and output tools such as sending messages to other agents; the AIOS shall include the basic and foundational tools that are frequently used by many agents as part of the platform, and provide Tool-Drivers (for hardware tools) and tool-APIs (for software tools) to enable easy calling of such tools by agents; Besides, each agent may also include its own specialized tools used for its own tasks.
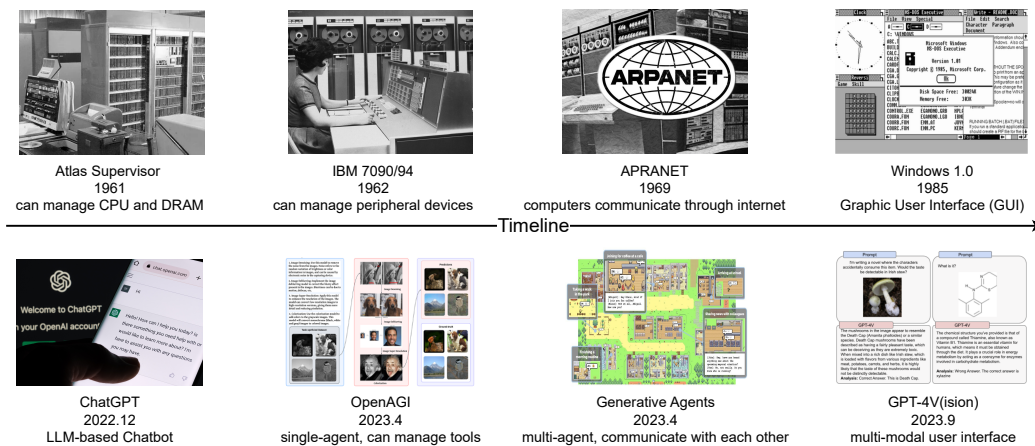
## 2.3 Development of OS and AIOS Aligned



Figure 3: Parallel Development of OS and AIOS/LLMOS.[24]

Figure 3 draws a parallel between the evolution of Operating Systems (OS) and LLM as OS (LLMOS). This comparison highlights their parallel progression in terms of enhanced functionality and user interaction. Initially, operating systems such as Atlas Supervisor were designed to manage basic computer resources such as CPU and DRAM. This evolved into more sophisticated version, exemplified by UNIX which can manage various peripheral devices. In a similar vein, LLMs have progressed from text-in-text-out chatbots to intricate LLM-based agents capable of managing various tools for complex task solving, as seen in OpenAGI (Ge et al., 2023). The figure also emphasizes the pivotal role of ARPANET in the development of TCP/IP protocols, which laid the foundation for today's Internet and connects multiple computers for communicating with each other. This is paralleled by the advancement of LLM-based multi-agent systems where agents can communicate with each other, signifying a trend towards more interconnected and collaborative LLM-based computing environments (Park et al., 2023). Additionally, the evolution of operating systems is marked by the development

---

[24]Image Sources:
https://history-computer.com/atlas-computer/,
https://en.wikipedia.org/wiki/History_of_Unix,
https://www.magzter.com/stories/technology/Gadgets-Philippines/ARPANET,
https://en.wikipedia.org/wiki/Windows_1.0x,
https://www.brookings.edu/articles/chatgpt-educational-friend-or-foe/,
https://github.com/agiresearch/OpenAGI,
https://github.com/joonspk-research/generative_agents,
https://openai.com/research/gpt-4v-system-card

of sophisticated graphical user interfaces (GUIs), such as those in Windows and Apple Macintosh. Similarly, LLM is evolving to include multi-modal interfaces, as demonstrated by GPT-4V(ision) (OpenAI, 2023). This comparison underscores the role of both OS and LLM in revolutionizing user interaction with technology, transitioning from managing fundamental components to facilitating more intuitive, user-centric experiences.

# 3 Architecture of AIOS

In this section, we establish the conceptual framework of "LLM as OS (LLMOS)," creating parallels between LLMOS components and traditional OS elements, as is shown in Table 1. In this analogy, the LLM is compared to the OS kernel, the context window to memory, and external storage to the file system. Tools and user prompts are equated with devices/libraries and the user interface, respectively. We will delve into the specifics of this correlation in the following discussion.

| OS-APP Ecosystem | AIOS-Agent Ecosystem | Explanation |
| --- | --- | --- |
| Kernel | LLM | The core of AIOS, providing supportive services for Agent Applications (AAPs). |
| Memory | Context Window | Short-term memory of the current session, such as the prompt-response history of this session. |
| Memory Management | Context Selection and Management | Select relevant context for the session, and manage the context such as adding, deleting and changing context information. |
| File | External Storage | Long-term storage of the AIOS's history sessions, user profiles, factual knowledge, etc. |
| File System | Retrieval-augmentation | Retrieve relevant information from long-term storage. |
| Devices/Libraries | Hardware/Software Tools | Help systems interact with the external world (both physical and virtual/digital world). |
| Driver/API | Tool-Driver/Tool-API | Serves as the interface for LLM/Agents to access and use the tools, usually in the form of prompts. |
| OS SDK | AIOS SDK | Assist users in developing Agent Applications. |
| User Interface (UI) | User Prompt/Instruction | Instruction(s) given by user to the system in (sometimes semi-structured) natural language (NL) to perform specific tasks. The NL instruction may be converted from users' non-NL instructions such as clicks. |
| Application (APP) | Agent Application (AAP) | A diverse world of AI Agents. |

Table 1: Comparison of OS-APP Ecosystem and AIOS-Agent Ecosystem

## 3.1 LLM (as AIOS Kernel)

The kernel is a set of computer programs at the core of a computer's operating system and generally has complete control over everything in the system. When a user command is issued, the OS parses the command and translates it into one or more system calls to enter the kernel, which are precise requests for the kernel to perform tasks such as process creation, memory allocation, and file manipulation. It then manages these tasks by scheduling the processes, allocating the necessary resources, interacting with device drivers for hardware access, and enforcing security measures. Throughout this process, the kernel also handles error checking and provides appropriate feedback. Upon completion, the kernel ensures that outputs are directed back to the user and that all resources are cleaned up to maintain system stability, effectively abstracting the intricate hardware details from the user.

Similarly, LLM acts as the operational core, akin to an OS kernel, responsible for planning and decomposing user requests, selectively calling tools, performing retrieval, and integrating all the

information from previous steps to generate the final response. An LLM handles a complex prompt or instruction by first interpreting the input to understand its context and intent, and then uses its internal parameters to process the information, decomposing the instruction into manageable sub-tasks. The LLM generates responses for these sub-tasks, ensuring coherence and relevance to the original instruction. Finally, it synthesizes these into a cohesive output delivered to the user.

### 3.1.1   Reasoning and Planning

In LLMOS, the role of LLM as the kernel subtly echoes the dynamics illustrated in the Dining Philosophers problem from computer science, a scenario that highlights the challenges of resource allocation and synchronization. In this problem, philosophers seated around a table must navigate shared resource usage–forks–in a way that avoids deadlock, a situation where no progress is made due to mutual blocking. Reflecting these challenges, the kernel in an operating system is adept at managing and scheduling resources to avert system deadlocks, ensuring smooth and conflict-free operations. This essential function of the kernel in maintaining system stability and efficiency finds a parallel in the role of the LLM within the AIOS framework. Here, the LLM is responsible for navigating complex informational environments and managing external resources. Crucial to the LLM's role is its planning ability, which involves generating a series of actions to achieve a specific goal (Ge et al., 2023). Central to the planning ability is the capability of reasoning (Bratman et al., 1988; Russell and Norvig, 1995; Fainstein and DeFilippis, 2015; Huang and Chang, 2022). Through reasoning, LLMOS deconstructs complex tasks from user instructions into more manageable sub-tasks, devising appropriate plans for each. Next, we will explore several representative planning strategies that illustrate this capability.

- **Single-path Planning.** In this strategy, the final task is decomposed into several intermediate steps, which are connected in a cascading manner, with each step leading to only one subsequent step. For example, CoT (Chain of Thoughts) prompting (Wei et al., 2022), as one of the celebrated capabilities of recent LLMs, is a pivotal breakthrough for performing complex multi-step reasoning when provided with limited examples. For example, by providing the models with "chain of thoughts", i.e., reasoning exemplars, or a simple prompt "Let's think step by step", these models are able to answer questions with explicit reasoning steps (Kojima et al., 2022). ReAct (Reasoning + Acting) (Yao et al., 2022b) is a prompt-based paradigm to synergize reasoning and acting in language models. It enriches the action space with self-thinking (or thoughts), which compose useful information by reasoning over the current context to support future reasoning or acting. RAFA (Reason for Future, Act for Now) (Liu et al., 2023a) studies how to complete a given task provably within a minimum number of interactions with the external environment.

- **Multi-path Planning.** In this strategy, the reasoning steps for generating the final plans are organized into a tree-like or graph-like structure. This mirrors the nature of human decision-making, where individuals often encounter a range of choices at each step of their reasoning process. For example, CoT-SC (Self-consistent CoT) (Wang et al., 2022) uses the CoT approach to produce multiple reasoning paths and answers for a complex problem, selecting the most frequently occurring answer as the final output. Tree of Thoughts (ToT) (Yao et al., 2023) assumes that humans tend to think in a tree-like structure when making decisions on complex problems for planning purposes, where each tree node is a thinking state. It uses LLM to generate evaluations or votes of thoughts, which can be searched using breadth first search (BFS) or depth first search (DFS). These methods improve the performance of LLMs on complex reasoning tasks. DFSDT (Depth-first search-based decision tree) (Qin et al., 2023b) employs a tree structure with each node representing a state that includes instruction context, prior API calls, and observations. The model not only reasons and moves to child nodes based on API calls but can also backtrack to explore alternative paths, ensuring a more diversified search and preventing cumulative errors. Graph of Thoughts (GoT) (Besta et al., 2023) further extends the reasoning paths from tree-structure to graph-structure, representing data produced by an LLM in the form of a flexible graph with individual units of information as nodes.

Valmeekam et al. (2022) concluded that "LLMs are still far from achieving acceptable performance on common planning/reasoning tasks which pose no issues for humans to do." Similarly, Wei et al. (2022) also acknowledged this limitation, noting that "we qualify that although chain of thought emulates the thought processes of human reasoners, this does not answer whether the neural network is actually reasoning." As a result, extensive efforts are still needed from the community to enhance the reasoning and planning ability of large language models.

### 3.1.2 Self-Improving

Just as kernel updates are driven by human feedback, focusing on bug reports and performance issues, leading to improvements in functionality, security, and stability, LLMOS also requires continuous enhancement to elevate its performance. This process involves the LLMs learning from interactions, refining their algorithms based on user experiences and feedback. By doing so, LLMs can develop new capabilities and skills, adapting to changing requirements and expectations. This iterative process of improvement ensures that LLMs remain effective, relevant, and efficient in handling diverse tasks and queries, mirroring the evolving nature of technology and user needs. After the LLM has been pre-trained, the learning strategies of LLM can be broadly categorized into two main types, as summarized and exemplified in the following.

- **Learning from Feedback.** LLMs can learn from feedback about the consequences of its actions to the environment. For example, Reflexion, as proposed by Shinn et al. (2023), is a framework to reinforce language agents through linguistic task feedback rather than update weights, enabling them to learn from prior failings. Similarly, Recursively Criticize and Improve (RCI) (Kim et al., 2023) is a prompting approach, which involves prompting the LLMs to identify issues in their output and subsequently enhance it based on the identified problems. Furthermore, these learning processes can be framed within the paradigm of Reinforcement Learning (RL). In this context, the LLM is trained to select actions that maximize a reward signal, aligning with the principles of RL. For example, Ouyang et al. (2022) presents Reinforcement Learning from Human Feedback (RLHF) to align LLMs with the feedback from human users; OpenAGI (Ge et al., 2023) presents Reinforcement Learning from Task Feedback (RLTF), which learns from the performance of executing the LLM-generated task solution to fine-tune the planning strategy of the LLM.

- **Learning from Examples.** Recently, there has been a surge of interest in fine-tuning LLMs to perform tasks in a supervised way. For example, ToolLLaMA (Qin et al., 2023b) is created by collecting various real-world APIs and generating instructions for their practical use, with solutions annotated using a language model such as ChatGPT and an efficient Depth-First Search Decision Tree. This process results in a dataset of instruction-solution pairs, on which a large language model such as LLaMA is fine-tuned to execute APIs according to instructions. Moreover, Gorilla (Patil et al., 2023) is trained by extracting key fields from API documentation and using GPT-4 to create instruction-API pairs, which are then used in a conversational format to fine-tune a model such as LLaMA, incorporating both retriever-aware and non-retriever training methods.

### 3.2 Context Window (as Memory)

In LLMOS, memory functions similarly to the context window in an LLM, defining the scope of information that the LLM can utilize and learn from while producing outputs. Increasing the amount of memory is desirable but always at a high cost.

Within the framework of an LLM, an expansion of the context window precipitates increased computational demands. This escalation is attributed to the quadratic computational complexity that is a characteristic of the attention mechanism employed in these models (Vaswani et al., 2017). During the training phase, one widely adopted strategy to avoid the exorbitant costs associated with lengthy context is to conduct an initial pre-training phase using a relatively limited context window, and subsequently followed by a fine-tuning stage employing an expanded context window (Chen et al., 2023b). Consequently, two primary challenges emerge: 1) the reduction of computational costs associated with processing long contexts, and 2) the development of a flexible position encoding mechanism suitable for extended contexts.

- **Efficient attention**: Various methods have been proposed to reduce the complexity of the multi-head attention mechanism, which can be categorized into three primary types: 1) Sparse attention mechanism (Zaheer et al., 2020; Gray et al., 2017; Kitaev et al., 2020; Beltagy et al., 2020; Ainslie et al., 2020; Wang et al., 2020) redefines the traditional computation of attention weights. In this approach, each query tensor is limited to engaging with only a subset of key tensors, as opposed to the entire set. This method effectively zeroes out other attention weights, thereby diluting the computation and reducing the overall computational burden; 2) Linear attention mechanism (Katharopoulos et al., 2020), which modifies the tensor multiplication process to be linear with respect to the sequence length without reducing the interaction between query tensors and key tensors; 3) Traditional multi-head attention uses multiple heads to split the query, key, and value

tensor. Another method of reducing complexity aims at sharing heads for keys and values to optimize memory usage. By sharing heads between the keys and values in the multi-head attention setup (Shazeer, 2019; Ainslie et al., 2023), it reduces the storage requirements for the key-value (KV) cache and improves efficiency.

- **Position encoding**: Position encoding can be classified into absolute position encoding and relative position encoding. The original Transformers paper (Vaswani et al., 2017) utilizes absolute position encoding, in which information is encoded in a combination of sin/cos functions whose wavelength increases from low- to higher-order dimensions of the embedding. Various methods, such as RoPE (Su et al., 2021), are proposed later to emphasize the relative position information between tokens. To extend the context window size, Alibi (Sun et al., 2022), LeX (Sun et al., 2022), and XPos (Press et al., 2021) are proposed to enable length extrapolation so that a model can conduct inference on long context while being trained only on short context. To increase context length, methods such as position interpolation (Chen et al., 2023b) and NTK-aware (Rozière et al., 2023) position interpolation based on RoPE can be used.

Even though technically, the long context can be processed, there is no guarantee that the information in the long context can be correctly used and learned by LLM. Recent research has found that LLM does not robustly make use of information in long input contexts (Liu et al., 2023b; Shi et al., 2023). In particular, performance is often highest when relevant information occurs at the beginning or end of the input context and significantly degrades when models must access relevant information in the middle of long contexts. Similar findings are also found in chat-based LLM (Yang and Ettinger, 2023), where model's ability to track and enumerate environment states is unsatisfying. Given long documents, various prompt compression methods, such as LLMLingua (Jiang et al., 2023), have been proposed to reduce context length by removing unimportant tokens from the text.

## 3.3 External Storage (as Files)

In addition to generating content based on knowledge acquired during pre-training, LLMs also utilize externally stored information for several purposes. These include enhancing predictions in domain-specific areas, generating more current information based on recent updates, and improving long-term memory retention. The external storage functions similarly to a file system within an operating system and supports various data formats. Certain formats enable the LLMs to swiftly query data as required, while others act as auxiliary knowledge bases. These knowledge are accessible for instant access to provide information in response to queries requiring high precision or expert knowledge. This section will commence with an explanation of how data is stored, followed by a discussion on the methodologies employed to retrieve pertinent information from the data storage.

### 3.3.1 Data Formats

A file system in OS stores, organizes, and retrieves information in different modalities, including but not limited to plain text, images, audios, or videos. Similarly, LLM stores and retrieves data in different formats, such as natural language, embedding vectors, or graphs. It should be emphasized that these formats are not mutually exclusive, and many models incorporate multiple formats to concurrently harness their respective benefits. In the following, we introduce several representative data formats, each with its distinct features and applications:

- **Embedding Vectors.** Embedding vectors represent words, phrases, or even entire documents as dense vectors in high-dimensional spaces. This format is crucial for processing information in machine-readable forms, such as natural language and images, enabling LLMs to efficiently and accurately retrieve necessary information for content generation. We note that dense vector retrieval, as exemplified by Karpukhin et al. (2020), generally serves as a necessary intermediate step for various data formats. Many different data formats are represented as dense vectors during this retrieval process. Specifically, when introducing embedding vector-based data formats, we refer to methods that explicitly store information in vector format. This approach is commonly used in conversational agents to maintain long-term memories. For example, Zhong et al. (2023) identified the lack of long-term memory as a significant limitation in current LLM-based applications such as ChatGPT. This issue is particularly noticeable in scenarios that require sustained interactions, such as personal companionship, psychological counseling, and secretarial tasks. To equip LLMs with the ability to effectively access long-term memory, Zhong et al. (2023) proposed "MemoryBank," a

novel vector-retrieval mechanism designed specifically for LLM-based applications. MemoryBank stores user-system past interactions, such as dialogues, as time-aware dense vectors known as memory pieces, with a dual-tower dense retriever and a memory updater inspired by Ebbinghaus' Forgetting Curve theory. A chatbot powered by MemoryBank demonstrates the effectiveness of this mechanism in long-term conversation. Similarly, Zhao et al. (2023) stores long-term conversational data as vectors for open-domain dialogue applications, where each piece of dialogue or summarized memory is transformed into a high-dimensional vector that captures its semantic meaning. Contrastive representation learning is used to increase the accuracy of memory retrieval. Lee et al. (2023) proposed a memory-enhanced chatbot to achieve long-term consistency and flexibility in open-domain conversations. It uses a summarizer model to summarize dialogues and store them as dense vectors in a memory pool after a certain number of rounds. Then, the relevant memory is retrieved to help generate responses. Later, Lu et al. (2023) introduced "Memochat," which proposed an iterative "memorization-retrieval-response" cycle to maintain consistency in long conversations covering multiple topics. In this cycle, LLMs are trained to create structured memos, which help to keep track of the conversation context and topics.

- **Plain-Text Documents.** Language models that generate responses by retrieving external plain-text documents can greatly improve the quality of responses. This is especially true for tasks that demand domain-specific knowledge or up-to-date information, such as question answering (Guu et al., 2020; Borgeaud et al., 2022) and fact verification (Lewis et al., 2020; Izacard et al., 2022). Plain-text documents, which consist of unstructured text data, are often vast in size and contain a wealth of information. Examples of such retrieval include Wikipedia articles (Guu et al., 2020), textbooks (Wang et al., 2023d), and programming code (Zhou et al., 2022). Retrieving information from these documents poses a challenge due to their dense and extensive nature. Therefore, advanced retrieval methods such as Dense Passage Retrieval (DPR) (Karpukhin et al., 2020) are frequently utilized in the retrieval process. Retrieval-augmented LLM based on external document collections is an important extension of LLMs, where the system is not just reliant on pre-trained knowledge but also actively retrieves external knowledge for better response generation. The augmentation allows for more accurate, up-to-date, and context-relevant responses, particularly crucial for tasks involving real-time data or specialized knowledge. Research in this field is often referred to as Retrieval-Augmented Generation (RAG), focusing on critical issues such as pre-training language models in retrieval contexts and fine-tuning them for downstream tasks.

- **Structured Data.** Structured data is a commonly used format for storing external, useful knowledge for LLMs. Since much information naturally exists in structural formats, this enables models to access a wide range of information within complex knowledge structures. The sources of structured data are diverse, with knowledge graphs being one of the most pertinent formats for augmenting LLMs' generative abilities: Pan et al. (2023) explicitly discuss the use of knowledge graphs to address the hallucination issues in LLMs and enhance their interpretability. Furthermore, several popular LLM tools, such as LlamaIndex (Liu, 2022), offer options to leverage existing knowledge graphs to improve knowledge generation. Another potentially valuable data format is tabular data. As introduced in Sundar and Heck (2023), a novel concept of Conversational Tables (cTBLS) is designed to enhance the capabilities of conversational AI by integrating tabular data. This work utilizes a dense table retrieval method to rank relevant table cells. Subsequently, the responses of LLMs are grounded in the tabular information and the conversational context.

### 3.3.2 Data Retrieval Methods

The ability to access relevant and accurate information from external storage through sophisticated data retrieval methods is crucial for LLMOS's execution of targeted actions. A primary challenge in this process is selecting the most appropriate data file from an extensive repository of information. These data retrieval methods operate automatically, leveraging advanced algorithms and machine learning techniques.

For instance, Zhu et al. (2023) show that LLM can store past accomplished sub-goals of video games using a dictionary, where the sub-goals are keys, and the corresponding action sequences are values. When encountering familiar objectives, the first action is easily retrieved using the name of the goal. Conversely, Park et al. (2023) propose a more sophisticated "memory stream" to record agents' past experiences in a list, labeled with text descriptions and timestamps of creation or last interaction. This data storage strategy enables agents to effectively retrieve useful experiences based on their current situation, using scores of Recency, Importance, and Relevance. Similarly, Zhong et al. (2023) discuss

storing detailed records of daily conversations, summaries of past events, and assessments of user personalities as vector representations indexed by FAISS (Johnson et al., 2019), a library used for efficient similarity search in stored vectors. Furthermore, Hu et al. (2023) highlight the limitations of conventional neural memory mechanisms, which are not symbolic and rely on vector similarity calculations, being prone to errors. It suggests the use of databases as an external symbolic memory for LLMs. Complex problems are decomposed into a series of SQL-based memory operation steps, greatly simplifying the retrieval process.

On the other hand, retrieving information from extremely large external documents heavily relies on existing methods in Information Retrieval (IR) research (Croft et al., 2010). Modern Information Retrieval involves two key stages: retrieval and ranking. The retrieval stage focuses on fetching relevant documents based on user queries using algorithms such as vector space models (Salton, 1975; Robertson et al., 2009), or pre-trained models such as BERT (Devlin et al., 2018; Karpukhin et al., 2020). LLM extensively applies these methods for applications that depend on generating comprehensive domain knowledge or accurate information. To ensure better accuracy in retrieving information from a vast amount of documents, effective index representations are usually learned during pre-training or fine-tuning (Guu et al., 2020; Borgeaud et al., 2022; Lewis et al., 2020; Izacard et al., 2022; Hua et al., 2023b).

### 3.4 Tools (as Devices/Libraries)

In a traditional OS, peripheral devices, such as keyboards, mice, and printers, extend the system's capabilities, allowing for diverse forms of input and output that enhance the overall functionality of the computer. There are also various programming libraries, which include a diverse set of reusable functionalities that can be leveraged by applications through API calls. Similarly, in the context of AIOS, hardware tools can be seen as analogous to these peripheral devices, and software tools can be seen as analogous to these libraries and APIs. These tools can range from data analysis modules to interactive interfaces, each adding a unique dimension to the LLM's processing and response abilities. They allow the LLM to interact with different data types, environments, and user requirements, significantly enriching its functionality. As a result, these hardware and software tools help the LLM to interact with the physical and digital worlds, expanding LLM's capabilities, and can be leveraged by agents for complex task solving. In the upcoming sections, we will delve into the specifics of these tools, illustrating how they enrich the LLM's capacity within the AIOS.

#### 3.4.1 Tool Categories

Though LLMs are adept at handling numerous tasks, they encounter limitations in complex tasks that require deep domain knowledge or interaction with the external world. External tools enable LLMs to harness various resources and specialized knowledge, bolstering their capabilities. In the following, we present several representative tools for LLMs, as discussed in the existing literature.

- **Software Tools** are domain expert models or APIs that help LLMs to finish a specialized sub-task, such as searching a query on the Web through a search API or checking the weather through a third party weather service API. Recent trends show a growing integration of APIs with LLMs, serving as interfaces for external programs to interact with LLMs, and acting as a bridge between the LLM and other software applications, thus extending LLMs' capabilities across various applications and services. For instance, OpenAGI (Ge et al., 2023) trains LLMs to use various domain expert models as tools for reasoning, planing, and complex task solving based on reinforcement learning from task feedback (RLTF). TPTU (Ruan et al., 2023a) interfaces with both Python interpreters and LaTeX compilers for mathematical computations. Gorilla (Patil et al., 2023), a fine-tuned LLM, is engineered to generate precise API call arguments and prevent hallucinations. ToolLLM (Qin et al., 2023b) presents a general framework for tool use, including data construction, model training, and evaluation. It also provides an instruction-tuning dataset for tools, collected from over 16,000 real-world APIs. TaskMatrix.AI (Liang et al., 2023b) connects foundational models with millions of APIs for diverse task completion, facilitating user interaction in an interactive manner. ChemCrow (Bran et al., 2023) integrates several expert-designed models to augment LLMs in chemistry-related tasks such as organic synthesis, drug discovery, and materials design. MM-REACT (Yang et al., 2023a) combines LLMs with various vision expert models for advanced multi-modal reasoning and actions. Using expert models as tools, LLM agents can tackle advanced tasks necessitating expert knowledge.

- **Hardware Tools.** While the aforementioned tools enhance LLMs in the digital world, physical tools such as robotics and embodied AI serve as pivotal means to connect LLMs with the physical world. These tools enable LLMs to actively observe, understand, and interact with their physical surroundings. Observations allow LLMs to gather various inputs from the physical world, converting them into multi-modal signals to augment actions such as navigation and manipulation. For example, Soundspace (Chen et al., 2020) explores observing physical space geometry using reverberating audio sensory inputs. Physical tools enable LLMs to execute user commands, transcending the role of merely providing natural language instructions. SayCan (Ahn et al., 2022), for instance, incorporates physical tools into LLMs for real-world tasks such as cleaning tabletops or retrieving items. In essence, physical tools act as the "hands, ears, eyes" etc. of LLMs to interact with the physical world, fostering real-world grounding.

- **Self-made Tools.** Current tools are primarily designed by humans. Recently, there are increasing interest in the use of LLMs for automated tool making. This involves generating executable programs or enhancing existing tools to create more powerful solutions, guided by appropriate instructions and demonstrations (Qin et al., 2023a; Qian et al., 2023b; Chen et al., 2021). For example, a large code model, as evaluated in Chen et al. (2021), is capable of generating executable programs based on user-provided instructions. These programs can then serve as specialized tools to address particular tasks. Furthermore, these LLMs can also acquire the ability to self-debug, which is an essential skill for maintaining and improving the tool functionality, as detailed in Chen et al. (2023a).

### 3.4.2 Tool-Driver and Tool-API

In traditional OS, Drivers or APIs play a pivotal role in enabling the system to interact with specific Devices or Libraries. The drivers provide interfaces to connect the OS with hardware devices, while the APIs provide interfaces to connect the OS or application with software libraries. In the context AIOS, where hardware tools are viewed as devices and software tools are viewed as libraries, Tool-Drivers and Tool-APIs are required, which serve as the interface for AIOS or agents to use these hardware and sofware tools, respectively.

Existing literature usually defines the Tool-Drivers and Tool-APIs in the form of prompts. Specifically, these prompts are composed of two essential elements: application scenarios and invocation methods. Much like how Drivers or APIs in a traditional OS control access to specific Devices or Libraries, commonly used prompts in AIOS should equip LLMs with an in-depth understanding of application scenarios. This enables LLMs to judiciously determine the appropriateness of a particular tool for a given task. Moreover, in parallel to how Drivers or APIs facilitate the communication between the OS and devices, tool instruction prompts should clearly outline the invocation methods. This is crucial for LLMs to comprehend the inputs and outputs of the tools, ensuring their effective execution and integration into the system.

Utilizing the inherent zero-shot and few-shot learning capabilities of LLMs (Radford et al., 2019; Brown et al., 2020), agents can gain insights about tools through zero-shot prompts that elucidate tool functionalities and parameters, or few-shot prompts offering demonstrations of particular tool usage scenarios and methodologies (Schick et al., 2023; Parisi et al., 2022). Usually, a single tool is inadequate for complex tasks. Hence, agents must adeptly decompose these tasks into manageable subtasks, where their understanding of the tools is pivotal. After understanding individual tools, LLMs should determine their application in addressing complex tasks. One approach is to generate actions by extracting pertinent information from memory relevant to the current task. For instance, Generative Agents (Park et al., 2023) maintain a memory stream, consulting it for recent, pertinent, and crucial information before each action to guide their decisions. In GITM (Ghost in the Minecraft) (Zhu et al., 2023), to achieve specific sub-goals, the agent probes its memory to identify if any similar tasks have been successfully executed before. If so, it replicates those successful actions for the current task. In collaborative frameworks such as ChatDev (Qian et al., 2023a) and MetaGPT (Hong et al., 2023), agents engage in mutual communication, retaining the conversation history in their memories. Another strategy involves executing a pre-formulated plan. For example, in DEPS (Describe, Explain, Plan and Select) (Wang et al., 2023a), given a task, if there are no indicators of the plan's failure, the agent proceeds with actions based on that plan.

# 4 AIOS-Agent Ecosystem

## 4.1 Agents as Applications

On top of the LLMOS-based AIOS, a diverse scope of AI Agent applications can be developed, akin to traditional OS-based applications such as browsers and photo-editing softwares. These Agent Applications (AAPs) generally comprise three components: the agent profile, an accessible external database, and task-specific tools. On one hand, the agent profiles are written into the prompt and used to indicate the agent functionality, influencing the LLM behaviors to exhibit certain roles such as a coder agent, a teacher agent, or a travel planning agent, as described in sources such as Qian et al. (2023a); Li et al. (2023); Ge et al. (2023). Typical agent profiles may encompass basic information such as age, gender, and career (Park et al., 2023), or psychology information reflecting the personalities of the agents (Serapio-García et al., 2023), as well as social information detailing the relationships between agents (Li et al., 2023). The nature of the agent's profile is tailored to the application's context; for example, applications focused on human cognitive processes will emphasize psychological information. Furthermore, the inclusion of an external database and toolset within the prompt enhances the LLM's interaction with the external world, specifying elements such as file locations and tool descriptions.

As illustrated in Figure 4, by merging the LLMOS layer with the OS layer, Hardware layer, and the Agent Application layer, we can establish an autonomous LLMOS-based Agent system. This system will respond to natural language instructions from users, capable of executing a variety of tasks through its interactions with the environment and its inherent knowledge.
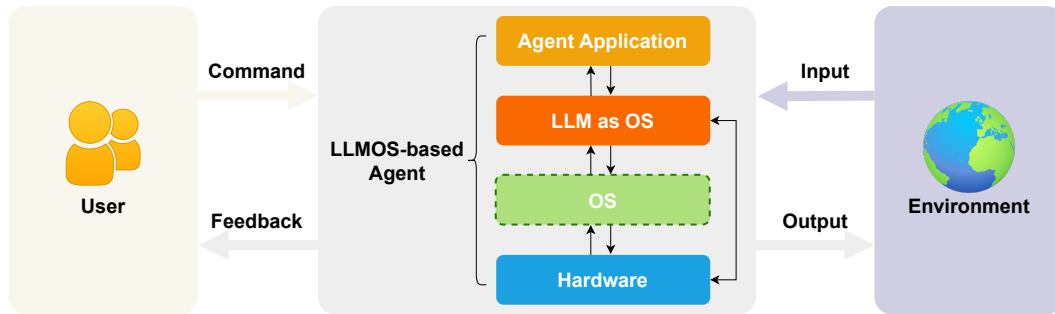


Figure 4: An illustration of LLMOS-based AI Agent.

## 4.2 Natural Language Programming for Agents

Natural Language Programming (NLProg) acts as a crucial intermediary between LLMOS and Agents within the AIOS-Agent ecosystem. This innovative approach allows average users to easily program Agent Applications (AAPs) using natural language. This development democratizes the creation and accessibility of computer software, representing a significant shift from the traditional OS-APP ecosystem. Traditionally, desktop or mobile applications (APPs) required programming by skilled software developers using professional programming languages. In contrast, NLProg empowers even those without formal training in these languages to develop applications.

This trend aligns with the historical evolution of programming languages, which have progressively become more user-friendly. The journey from binary codes to assembly language and then to various high-level languages such as C, C++, Java, and Python reflects this trend. The introduction of Natural Language as a Programming Interface is a natural progression in this evolution. It simplifies programming for the average user, allowing them to program agent applications and interact with computers without the need for specialized training in conventional programming languages.

Moreover, this shift has broader implications for the field of computer science and technology. It opens up new possibilities for innovation and creativity, as a wider range of individuals can contribute to software development. This inclusivity can lead to the development of more diverse and tailored applications, as people from different backgrounds and with varying expertise can bring their unique perspectives to software design. Additionally, NLProg in the AIOS ecosystem fosters a more intuitive interaction between humans and computers, enhancing user experience and potentially leading to

more efficient and effective use of technology in various sectors. As this approach gains traction, it could significantly alter the landscape of technology development, making it more accessible and aligned with the natural human way of communication and understanding.

### 4.3 The Ecosystem

In the AIOS-Agent ecosystem, the potential for collaboration and networking among agents heralds a new era of digital assistance and decision-making. These agents, with their diverse skill sets and role profiles, can work in tandem to address complex challenges, offering a multifaceted approach to problem-solving. For instance, an agent with a background in data analysis can collaborate with another agent specialized in creative design, leading to innovative solutions that a single agent might not conceive.

Moreover, the AIOS ecosystem is designed to be scalable, accommodating an increasing number of agents as user needs grow and evolve. This scalability ensures that the system remains efficient and effective, even as the complexity of tasks increases. The AIOS also emphasizes the importance of learning and adaptation. Agents in this ecosystem can learn from their interactions, both with users and other agents, continuously evolving and enhancing their capabilities. This feature is crucial for keeping up with the rapidly changing landscape of technology and user expectations.

In the broader context, the AIOS-Agent ecosystem can significantly impact various industries, from healthcare, where agents can assist in patient care and medical research, to education, where agents can offer personalized learning experiences. The versatility and adaptability of this ecosystem make it a valuable asset in any field where decision making, data analysis, and creative problem solving are crucial. As this technology matures, it holds the promise of transforming our interaction with the digital world, making it more intuitive, efficient, and responsive to our needs.

## 5 LLMOS in Practice: AI Agents

In this section, we present a comprehensive overview of the current applications of LLMOS-based agents, aiming to provide a wide-ranging perspective on their practical deployment scenarios. Specifically, we mainly introduce three scenarios: single agent applications, multi-agent applications, and human-agent applications.

### 5.1 Single Agent Applications

Single agent applications mainly focus on utilizing a single agent to address various user tasks. We categorize single agent applications into two distinct types based on the external environments they interacted with, i.e., physical environment and virtual/digital environment.

#### 5.1.1 Physical Environment

Unlike virtual or simulated environments, the physical environment is made up of tangible elements that an AI Agent must navigate or manipulate. This concept is particularly relevant in the field of robotics and embodied AI, where the agents are not just software algorithms but have a physical presence or are integrated with physical systems. Following this, we present a range of exemplary scenarios, as detailed in existing literature.

- **Scientific Research.** Agents have the capacity to function autonomously, undertaking experiments independently, while also serving as invaluable resources for scientists engaged in research projects (Boiko et al., 2023; Bran et al., 2023). For instance, Boiko et al. (2023) propose an end-to-end platform that automates scientific experimentation, integrating AI modules, reasoning capabilities, software tools, and laboratory hardware. It autonomously performs tasks ranging from online information gathering and experiment design to running protocols and controlling robotic equipment, while adapting to errors and refusing unethical requests. ChemCrow (Bran et al., 2023) is a suite of 17 specialized tools designed to aid chemical research, offering methodological suggestions and highlighting safety hazards based on the input objectives. Huang et al. (2023) also propose MLAgentBench, a suite of ML tasks for benchmarking AI research agents.
- **Robotics.** Recent studies have employed LLM-based agents in the fields of Robotics. For example, ChatGPT for Robotics (Vemprala et al., 2023) employs ChatGPT for a wide array of robotics tasks

through strategic prompt engineering, showing its ability to comprehend and respond to natural language instructions in the context of robotics applications. SayCan (Ahn et al., 2022) comprises two integral components: the "Say" part of the LLM, responsible for task-grounding by identifying pertinent actions for a high-level objective, and the "Can" part, which encompasses the learned affordance functions that offer a world-grounding, thereby determining the feasible actions for the plan's execution. It ensures not only the relevance of the actions chosen for the specified task but also their feasibility in the real-world scenario. VOYAGER (Wang et al., 2023e) presents lifelong learning with prompting mechanisms, skill library, and self-verification, which are based on LLM. These three modules aim to enhance the development of more complex behaviors of agents.

- **Autonomous Driving.** Recent studies have harnessed LLMs to enhance self-driving car technologies. For example, Fu et al. (2023a) propose an agent-based autonomous driving system, which widely adopts a four-module framework: Environment, Agent, Memory, and Expert. Within this framework, the Environment sets the interaction stage, the Agent perceives the environment and makes decision, the Memory accumulates experience for action, and the Expert provides training advice and inconsistency feedback.

### 5.1.2 Virtual/Digital Environment

Agents in virtual or digital environment mainly includes the manipulation of APIs (Schick et al., 2023; Yao and Narasimhan, 2023; Ge et al., 2023; Parisi et al., 2022; Tang et al., 2023), accessing the Internet and websites (Nakano et al., 2022), executing codes (Zhang et al., 2023), and simulation in historical settings (Hua et al., 2023a). Such digital grounding is cheaper and faster than physical or human interaction. It is thus a convenient test bed for language agents and has been studied with increasing intensity in recent years. In the following, we present several representative scenarios as studied in the existing literature.

- **Coding.** This category focuses on leveraging the capabilities of agents to generate programs. For example, ToolCoder (Zhang et al., 2023) is a system that merges API search tools with existing models to facilitate code generation and API selection, using a two-step approach. Initially, an automated data annotation technique involving ChatGPT embeds tool usage data into the source code, followed by fine-tuning the code generation model; during inference, the API search tool is integrated to autonomously suggest API choices, optimizing code generation and improving API selection decision-making. Moreover, Lemur-series models (Xu et al., 2023) are meticulously pre-trained and instruction fine-tuned to demonstrate balanced language and coding capabilities.

- **Web Service.** This category primarily revolves around utilizing agents to address web-based tasks through diverse APIs. For example, Auto-GPT (Gravitas, 2023) is an automated agent designed to set multiple objectives, break them down into relevant tasks, and iterate on these tasks until the objectives are achieved. OpenAGI (Ge et al., 2023) is an LLM-based agent designed for reasoning, planing, and executing tools to achieve complex tasks, accompanied with a benchmark to evalute the agent's task-solving performance. BMTools (Qin et al., 2023a) is an open-source repository that extends LLMs with tools and provides a platform for community-driven tool building and sharing. It supports various types of tools, enables simultaneous task execution using multiple tools, and offers a simple interface for loading plugins via URLs, fostering easy development and contribution to the BMTools ecosystem. Mind2Web (Deng et al., 2023) provides a benchmark for developing and evaluating generalist agents for the Web, which are agents that can follow language instructions to complete complex tasks on websites. MusicAgent (Yu et al., 2023) integrates numerous music-related tools and an autonomous workflow to address user requirements. Auto-UI (Zhan and Zhang, 2023) is a multi-modal solution that directly interacts with the interface, bypassing the need for environment parsing or the dependence on application-dependent APIs.

- **Games.** This includes agents interacting in the game environments (Hausknecht et al., 2020; Côté et al., 2019; Shridhar et al., 2020). For example, MineClip, introduced by Fan et al. (2022), is a novel agent learning algorithm that leverages large pre-trained video-language models as a learned reward function. Based on it, the authors further proposed MineDojo, a framework built on the popular Minecraft game that features a simulation suite with thousands of diverse open-ended tasks and an internet-scale knowledge base with Minecraft videos, tutorials, wiki pages, and forum discussions.

- **Recommendation.** Recent literature demonstrates the efficacy of employing LLM and Agents in recommender systems (Geng et al., 2022; Wang et al., 2023c; Feng et al., 2023; Wang et al.,

2023f). For instance, RecMind (Wang et al., 2023c) develop an LLM-based recommender system agent, which provides personalized recommendations based on planning, use of tools, obtaining external knowledge, and leveraging individual user's personalized data. LLMCRS Feng et al. (2023) is a conversational recommendation agent that utilizes Large Language Models (LLMs) for efficient sub-task management during the recommendation process. It combines LLMs with expert models for specific sub-tasks and employs LLMs as a language interface for generating improved user responses, thereby enhancing overall performance and response quality in conversational recommendation systems.

## 5.2 Multi-Agent Applications

Multi-Agent Systems (MAS) (Wooldridge and Jennings, 1995) emphasize the coordination and collaboration among a group of agents to effectively solve problems. The existing LLM-based MAS landscape is broadly categorized into two types: Collaborative Interaction and Adversarial Interaction.

### 5.2.1 Collaborative Interaction

As the scope and complexity of tasks amenable to Large Language Models (LLMs) increase, a logical strategy to augment the effectiveness of these agents is to employ cooperative multi-agent systems. Such systems, prevalently utilized in practical applications, operate on the principle where each agent evaluates and understands the requirements and capabilities of its peers, thereby fostering a collaborative environment conducive to shared actions and information exchange (Li et al., 2023). In the specific area of Non-Player Characters (NPCs), the concept of Generative Agents (Park et al., 2023) emerges as a compelling simulation of human behavior within interactive applications. This approach is exemplified by the deployment of twenty-five agents in a sandbox environment akin to The Sims, allowing users to engage with and influence the agents as they execute daily routines, interact socially, establish relationships, and organize group activities. Furthermore, the Humanoid Agents system (Wang et al., 2023b) enhances the realism of Generative Agents by incorporating three fundamental aspects of "System 1" processing: the fulfillment of basic needs (such as hunger, health, and energy), emotional responses, and the dynamics of interpersonal relationships. In the field of Software Development, the MetaGPT system (Hong et al., 2023) represents a specialized LLM application that leverages a multi-agent conversational framework. This innovative framework facilitates automatic software development by assigning distinct roles to various GPT models, enabling them to collaborate effectively in the creation of software applications. Additionally, BOLAA (Liu et al., 2023c) introduces a controller module that orchestrates the coordination and communication among multiple collaborative agents, thereby streamlining the selection and interaction processes between different labor agents. CHATDEV (Qian et al., 2023a) proposes an advanced software development framework that utilizes agents to foster enhanced collaboration among the diverse roles integral to the software development cycle. In the domain of Conversational AI, research exemplified by Fu et al. (2023b) delves into the potential of LLMs to autonomously refine their negotiation skills. This is achieved through engaging the models in bargaining games against one another, complemented by the integration of natural language feedback from an AI critic. This study underscores the evolving capabilities of LLMs in complex, interactive settings.

### 5.2.2 Adversarial Interaction

Traditionally, collaborative methods have been extensively explored in multi-agent systems. However, researchers are increasingly recognizing that introducing concepts from game theory into systems can lead to more robust and efficient behaviors. For example, Du et al. (2023) introduce the concept of debate, endowing agents with responses from fellow peers. When these responses diverge from an agent's own judgments, a "mental" argumentation occurs, leading to refined solutions. ChatEval (Chan et al., 2023) establishes a role-playing-based multi-agent referee team. Through self-initiated debates, agents evaluate the quality of text generated by LLMs, reaching a level of excellence comparable to human evaluators. Corex (Sun et al., 2023) is constituted by diverse collaboration paradigms including Debate, Review, and Retrieve modes, which collectively work towards enhancing the factuality, faithfulness, and reliability of the reasoning process. These paradigms foster task-agnostic approaches that enable LLMs to "think outside the box," thereby overcoming hallucinations and providing better solutions. MAD (Multi-Agent Debate) (Liang et al., 2023a) is a framework wherein several agents engage in a "tit-for-tat" exchange of arguments under the oversight of a judge

who steers the discussion towards a conclusive solution. Furthermore, WarAgent (Hua et al., 2023a) considers each country as an LLM-based agent and simulates the international conflicts among the countries using World War I, World War II, and the Warring States Period in Ancient China as examples, which showcases possible approaches towards LLM multi-agent based policy simulation and answering the "what if" questions for historical analysis.

## 5.3 Human-Agent Applications

Most existing agent frameworks often limit themselves to defining and controlling agent behavior through system prompts, allowing the agent to independently plan and act. A notable shortcoming of this approach is the restricted, and sometimes non-existent, capacity for meaningful interaction between human users and agents, including multi-agent setups. Addressing this gap, AutoGen (Wu et al., 2023) offers an open-source solution enabling developers to construct LLM applications through multiple agents. Specifically, AutoGen distinguishes itself with agents that are not only customizable and conversable, but also versatile in their operational modes, which incorporate a blend of LLMs, human inputs, and tools, enhancing the interaction capabilities and efficiency of the agents. Furthermore, AGENTS (Zhou et al., 2023) introduces a novel approach to creating controllable agents. This method involves the use of symbolic plans or standard operating procedures (SOPs), which can be generated by an LLM and subsequently modified by the user. This feature allows for greater customization and fine-tuning of agents, providing a more user-centric and adaptable agent framework. Collectively, these developments represent a significant shift in the landscape of agent frameworks, moving towards systems that not only automate tasks but also facilitate a more interactive and collaborative environment between humans and agents.

# 6 OS-inspired Future Directions

The evolution history of operating system over the past half century has witnessed the continuous development of computer hardware and the explosive growth of data, which empowers the fast iteration of Artificial Intelligence in the past decade. In this section, we enumerate the lessons learned from the history of operating systems and provide envisions of the future directions for AIOS.

## 6.1 Resource Management

### 6.1.1 Memory Management

Physical memory (DRAM) has always been an insufficient resource from the beginning age of computer systems until now. To reduce the tension between users' requirements and the fact of DRAM resource shortage, several approaches have been proposed in modern operating systems.

- **Swapping to external storage.** A dedicated partition in external storage is reserved for swapping unused memory regions to external storage in a user-transparent way, as detailed in section 2.1.1. This approach enlarges the available DRAM resources in the system.

- **Memory sharing.** To support data sharing across applications, modern operating systems provide memory sharing between applications, which provides both sharing and also reduces the extra copies of shared data.

- **Memory disaggregation.** Entering the terabyte scale data, fitting the memory requirements for an application in a single operating system on a sole machine is becoming challenging. To remedy this, disaggregated memory techniques were proposed to allow an application to use memory on another machine via network. Moreover, recent Compute eXpress Link (CXL) technique (CXL, 2023) significantly reduces the software overheads and hardware latency of remote memory access.

LLMs have revolutionized abilities, but are constrained by limited context windows, hindering their utility in tasks such as extended conversations and document analysis. To mitigate the limited context window problem in LLMs, approaches inspired by the above operating system memory management methods can be developed. A swapping mechanism, similar to OS swapping to external storage, can be implemented in LLMs to temporarily store inactive parts of their context, effectively enlarging the context window. This would require efficient retrieval systems to minimize latency. For example, MemGPT (Packer et al., 2023) intelligently manages different memory tiers in order to effectively

provide extended context within the LLM's limited context window, and utilizes interrupts to manage control flow between itself and the user. Additionally, mirroring OS memory sharing, LLMOS can share context or learned patterns across different LLMOS-based Agent instances, reducing redundancy and allowing access to a larger shared knowledge pool for agents. Finally, drawing from memory disaggregation in modern OS, agents can leverage networked memory resources, enabling them to access and process data stored across multiple LLMOSes, thus expanding their abilities significantly. Each of these strategies, while offering potential solutions, also presents unique challenges such as managing latency, ensuring consistency in shared contexts, and handling the complexities of distributed memory systems.

### 6.1.2   Tool Management

Serving as external resources outside LLM, the hardware tools (e.g., robotics) and software tools (e.g., Web Search API) are the counterpart of devices and libraries in modern operating systems as shown in Table 1, respectively. Taking the example of the ecosystem in modern Linux operating system–the most widely-used open source operating system community by now, here are the successful experiences in its evolving history. Specifically, a rich set of built-in and third-party libraries from thousands of experts and open-source developers leads to the success of the Linux ecosystem. Managing the install/uninstall and dependencies of those libraries, along with the versioning for tracking software development, is critical. The Linux ecosystem, over the past few decades, provides library management tools such as *dpkg*[25] in Debian-based Linux distributions, and *yum*[26] in RHEL-based Linux distributions. *Git*[27], a distributed version control system that plays a crucial role in the development of the Linux ecosystem, allows collaborative and parallel development, boost the speed of development cycles in building the Linux ecosystem.

The process of code comparison in Git, particularly during merging and rebasing, typically analyzes changes at the line level rather than understanding the semantic meaning of the text. However, as discussed in Section 4.2, the development of agents using natural language is increasingly achievable. Adapting existing version control software, which is based on the code or texts without knowing the semantics and contexts, for use with natural languages poses distinct challenges. First, natural languages, while structured by grammars, exhibit a loosely coupled relationship with the varied expressions of different users. Second, this loose coupling can result in natural language statements that are semantically equivalent but differ in their wording. For example, in the context of an AI agent system, the instructions "Analyze the latest sales data and generate a report" and "Generate a report based on the analysis of the most recent sales data" convey the same instruction to the agent but are phrased differently. Incorporating the ability to recognize and reconcile these semantic nuances in natural language into version control software is essential. This is especially critical in collaborative development of complex AI agent systems, where such a feature can greatly streamline development cycles by effectively managing and merging diverse natural language inputs.

### 6.2   Communication

Domain-Specific Languages (DSLs) are widely used in both native operating systems and cloud environments, which address specific requirements or tasks within a particular domain. Operating systems are often equipped with scripting languages on top of the command-line interfaces that allow users to perform various tasks collaboratively. Unix-like operating systems use shell scripting languages such as *Bash*[28], which is designed for automating tasks in a command-line environment. In the cloud computing, DSLs are commonly used to define infrastructure as code (Sledziewski et al., 2010). It allows users to describe and provision cloud infrastructure resources, including virtual machines, networks, and storage; It can also be used for scheduling tasks on infrastructure resources. DSLs strike a balance between the underlying OSes and users that they are readable for both.

Leveraging the wisdom gleaned from OSes, multi-AIOS communication can be enhanced by adopting structured communication protocols that are analogous to the function of DSLs in simplifying complex tasks. Just as DSLs provide a medium for users to interact effectively with the operating system and its resources, a similar specialized protocol can be established for LLMs to communicate with

---

[25]https://man7.org/linux/man-pages/man1/dpkg.1.html
[26]http://yum.baseurl.org/
[27]https://git-scm.com/
[28]https://www.gnu.org/software/bash/

one another. This protocol would standardize interactions, allowing for the clear transmission of context, tasks, and goals between LLMs, thus facilitating a more coordinated and coherent multi-agent operation. It would ensure that despite the varied functionalities and knowledge bases of individual LLMs, there is a common language or method through which they can collaborate, share insights, and synchronize their learning processes. This approach mirrors the way operating systems manage resources and processes, ensuring harmonious and efficient functionality across various components of a system.

In current AI Agents, the task-solving plan is still represented by natural language in most cases. However, in the future, we can even develop DSLs or semi-structured natural language grammars for representing Agent's task-solving plans. This involves the following breakthroughs:

- Defining Basic Operations: This would standardize the basic operations, tools, or commands that an AI Agent understands and can execute for task-solving.

- Sequence of Operations: Task-solving plans would then be sequences of these basic operations, making them more structured and potentially more efficient.

Using LLM to interpret users' natural language instructions into a DSL-composed plan brings several benefits, including 1) Improved Consistency: Standardized operations would lead to more predictable and consistent outcomes; 2) Easier to Interpret: A well-defined DSL makes it easier for AI Agents to interpret and execute plans; and 3) Scalability: With a standard DSL, it is easier to scale solutions across different platforms and applications.

However, achieving this goal also meets some important challenges that need future research attention: 1) Complexity of Natural Language: Natural language is inherently ambiguous and context-dependent, making it challenging to convert instructions into structured plans consistently; 2) Flexibility vs. Standardization: Striking a balance between the flexibility of natural language and the rigidity of a standardized DSL; 3) Interoperability: Ensuring the DSL works well with a wide range of tools and platforms; and 4) Adaptation and Learning: The system needs to continuously learn and adapt to new instructions, tools, and tasks.

Overall, the use of a large language model as an interpreter (LLM as Interpreter), which can translate natural language described plans to DSL described plans, can be an important direction to create task-solving plans, and to bridge LLM as OS (LLMOS) and Agent Applications (AAP). The development of a DSL for such plans could further streamline and standardize the process, though it would come with its own set of challenges. This approach has the potential to make complex task execution more accessible and efficient, paving the way for more advanced AI systems in the future.

### 6.3 Security

Security has been an important issue as the wide-spread of operating systems from labs to our daily life in the 1980s. The consequences of the operating system vulnerabilities can be roughly categorized as following.

- **Breaking down the system.** In its early stages, viruses like the *Morris Worm* (Orman, 2003) were primarily created to compromise operating systems, serving as a demonstration of individual programmers' prowess in hacking. While these attacks did not result in direct financial losses, users faced potential harm through the compromise of personal data or critical workplace documents.

- **Racketeering.** In later stages, vulnerabilities within operating systems became targets for illicit activities, including the exploitation of users for racketeering purposes. An example of this is the *WannaCry Ransomware*[29], which encrypts users' files and demands a ransom for their release. Additionally, banking trojans like *Zeus*[30] exploit vulnerabilities by intercepting communication channels between users and banking systems, resulting in direct financial losses for the affected users.

- **Stealing Resource.** Malicious software, such as *Coinhive*[31], is crafted to harness the computing power of other machines for cryptocurrency mining, including Bitcoin. While it may not inflict

---

[29]https://nvd.nist.gov/vuln/detail/cve-2017-0143
[30]https://nvd.nist.gov/vuln/detail/cve-2010-0188
[31]https://krebsonsecurity.com/2018/03/who-and-what-is-coinhive/

direct harm on operating systems, the substantial utilization of CPU and memory resources can significantly impair the overall system performance. In cloud environments, this slowdown has the potential to translate into financial losses, making it imperative to address such threats proactively.

Security vulnerabilities in operating systems are hard to prevent. The state-of-the-art approaches detect and capture those malware and virus at different levels.

- **Static Analysis.** This approach conducts code-level or binary-level analyses by examining the code or binary image of an application or part of the OS. It is often performed when a third-party application is published to the cloud, or when a file is downloaded to the file system in a user's operating system.
- **Fuzzing.** Fuzzing (Fuzz Testing) involves the automated generation of a large number of random or semi-random inputs to a program to discover vulnerabilities, crashes, or unexpected behaviors.

Similarly, adversarial attacks have been the subject of extensive study in LLM research (Wei et al., 2023; Zou et al., 2023; Qi et al., 2023a; Yang et al., 2023c), representing a significant threat to the security of AIOS. Furthermore, research (Yang et al., 2023b; Qi et al., 2023b) has revealed that an aligned LLM can be broken using a very small dataset comprising only a few hundred data points. This vulnerability is not only a significant concern in terms of system integrity but also raises alarming implications for the safety of agents interacting with these systems. How to be robust and defend against such attacks has yet to be studied in the context of AIOS, LLMOS, and Agents.

Looking into the future, as agent applications in AIOS evolve to be programmed by natural language, the intricacy of scanning the code of such applications will increase due to the expansive and less structured nature of natural language compared to the constrained syntax of programming languages. This necessitates a robust and effective scanning tool for AIOS agents. Moreover, fuzzing can be seen as an initial step towards creating a Red-teaming dataset for LLMOS-based Agents. For instance, Ruan et al. (2023b) have developed a multi-agent system that generates red-teaming scenarios for LLM-based agents, using GPT-4 to simulate adversarial environments within textual scenarios. The study provides a rich array of scenarios that serve as a valuable resource for future research into the alignment of LLM-based Agents.

## 7 Conclusions

This paper presents a novel vision for the future of computing within the AIOS-Agent ecosystem, where the LLM functions as the core of AIOS. This innovative approach marks a significant departure from the conventional OS-APP ecosystem, heralding a new era in technology where AI and traditional computing systems merge seamlessly. The AIOS-Agent ecosystem envisaged here is not just an incremental change but a fundamental shift in how we interact with technology. By positioning LLM at the system level, Agents as applications, Tools as devices/libraries, and Natural Language as the Programming Interface, we redefine the interaction between users, developers, and the digital world. This paradigm shift promises to democratize software development and access, allowing users and developers to program Agent Applications (AAPs) using natural language. This accessibility contrasts sharply with the traditional ecosystem, where software development is confined to those with specialized programming skills. Moreover, the discussion of single and multi-agent systems, as well as human-agent interactions, illustrates the potential of AIOS in enhancing productivity, creativity, and decision-making processes across various domains. Looking ahead, the proposed strategic roadmap, informed by the developmental trajectory of the traditional OS-APP ecosystem, offers a pragmatic and systematic approach to the evolution of AIOS and its Agent Applications. This roadmap not only guides future development and research in this field but also anticipates the challenges and opportunities that lie ahead.

26

# References

Amit Agarwal, Saibal Mukhopadhyay, Arijit Raychowdhury, Kaushik Roy, and Chris H Kim. 2006. Leakage power analysis and reduction for nanoscale circuits. *IEeE Micro* 26, 2 (2006), 68–80.

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. 2022. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691* (2022).

Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints. *arXiv preprint arXiv:2305.13245* (2023).

Joshua Ainslie, Santiago Ontanon, Chris Alberti, Vaclav Cvicek, Zachary Fisher, Philip Pham, Anirudh Ravula, Sumit Sanghai, Qifan Wang, and Li Yang. 2020. ETC: Encoding long and structured inputs in transformers. *arXiv preprint arXiv:2004.08483* (2020).

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150* (2020).

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687* (2023).

Daniil A Boiko, Robert MacKnight, and Gabe Gomes. 2023. Emergent autonomous scientific research capabilities of large language models. *arXiv preprint arXiv:2304.05332* (2023).

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*. PMLR, 2206–2240.

Andres M Bran, Sam Cox, Andrew D White, and Philippe Schwaller. 2023. ChemCrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376* (2023).

Michael E Bratman, David J Israel, and Martha E Pollack. 1988. Plans and resource-bounded practical reasoning. *Computational intelligence* 4, 3 (1988), 349–355.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of Artificial General Intelligence: Early experiments with GPT-4. arXiv:2303.12712 [cs.CL]

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. ChatEval: Towards Better LLM-based Evaluators through Multi-Agent Debate. arXiv:2308.07201 [cs.CL]

Harrison Chase. 2022. *LangChain*. https://github.com/hwchase17/langchain

Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman. 2020. Soundspaces: Audio-visual navigation in 3d environments. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*. Springer, 17–36.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023b. Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595* (2023).

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023a. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128* (2023).

F. J. Corbató, J. H. Saltzer, and C. T. Clingen. 1971. Multics: The First Seven Years. In *Proceedings of the May 16-18, 1972, Spring Joint Computer Conference* (Atlantic City, New Jersey) *(AFIPS '72 (Spring))*. Association for Computing Machinery, New York, NY, USA, 571–583. `https://doi.org/10.1145/1478873.1478950`

Jonathan Corbet, Alessandro Rubini, and Greg Kroah-Hartman. 2005. *Linux Device Drivers, 3rd Edition*. O'Reilly Media, Inc.

Marc-Alexandre Côté, Akos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, et al. 2019. Textworld: A learning environment for text-based games. In *Computer Games: 7th Workshop, CGW 2018, Held in Conjunction with the 27th International Conference on Artificial Intelligence, IJCAI 2018, Stockholm, Sweden, July 13, 2018, Revised Selected Papers 7*. Springer, 41–75.

W Bruce Croft, Donald Metzler, and Trevor Strohman. 2010. *Search engines: Information retrieval in practice*. Vol. 520. Addison-Wesley Reading.

CXL. 2023. Compute Express Link Specification. `https://www.computeexpresslink.org/`.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2Web: Towards a Generalist Agent for the Web. *arXiv preprint arXiv:2306.06070* (2023).

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. Improving Factuality and Reasoning in Language Models through Multiagent Debate. arXiv:2305.14325 [cs.CL]

Esin Durmus, Karina Nyugen, Thomas I Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, et al. 2023. Towards measuring the representation of subjective global opinions in language models. *arXiv preprint arXiv:2306.16388* (2023).

Susan S Fainstein and James DeFilippis. 2015. *Readings in planning theory*. John Wiley & Sons.

Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems* 35 (2022), 18343–18362.

Yue Feng, Shuchang Liu, Zhenghai Xue, Qingpeng Cai, Lantao Hu, Peng Jiang, Kun Gai, and Fei Sun. 2023. A Large Language Model Enhanced Conversational Recommender System. *arXiv preprint arXiv:2308.06212* (2023).

Giorgio Franceschelli and Mirco Musolesi. 2023. On the creativity of large language models. *arXiv preprint arXiv:2304.00008* (2023).

Daocheng Fu, Xin Li, Licheng Wen, Min Dou, Pinlong Cai, Botian Shi, and Yu Qiao. 2023a. Drive Like a Human: Rethinking Autonomous Driving with Large Language Models. arXiv:2307.07162 [cs.RO]

Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023b. Improving language model negotiation with self-play and in-context learning from ai feedback. *arXiv preprint arXiv:2305.10142* (2023).

Yingqiang Ge, Wenyue Hua, Kai Mei, jianchao ji, Juntao Tan, Shuyuan Xu, Zelong Li, and Yongfeng Zhang. 2023. OpenAGI: When LLM Meets Domain Experts. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as Language Processing (RLP): A Unified Pretrain, Personalized Prompt & Predict Paradigm (P5). In *Proceedings of the 16th ACM Conference on Recommender Systems*. 299–315.

Significant Gravitas. 2023. *AutoGPT*. https://news.agpt.co/

Scott Gray, Alec Radford, and Diederik P Kingma. 2017. Gpu kernels for block-sparse weights. *arXiv preprint arXiv:1711.09224* 3, 2 (2017), 2.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*. PMLR, 3929–3938.

Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. 2020. Interactive fiction games: A colossal adventure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 7903–7910.

Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352* (2023).

Chenxu Hu, Jie Fu, Chenzhuang Du, Simian Luo, Junbo Zhao, and Hang Zhao. 2023. ChatDB: Augmenting LLMs with Databases as Their Symbolic Memory. *arXiv preprint arXiv:2306.03901* (2023).

Wenyue Hua, Lizhou Fan, Lingyao Li, Kai Mei, Jianchao Ji, Yingqiang Ge, Libby Hemphill, and Yongfeng Zhang. 2023a. War and Peace (WarAgent): Large Language Model-based Multi-Agent Simulation of World Wars. *arXiv preprint arXiv:2311.17227* (2023).

Wenyue Hua, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. 2023b. How to Index Item IDs for Recommendation Foundation Models. In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region*. 195–204.

Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403* (2022).

Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2023. Benchmarking Large Language Models As AI Research Agents. *arXiv preprint arXiv:2310.03302* (2023).

IEEE and The Open Group. 2018. The POSIX standard. https://pubs.opengroup.org/onlinepubs/9699919799.2018edition/.

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299* (2022).

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. Llmlingua: Compressing prompts for accelerated inference of large language models. *arXiv preprint arXiv:2310.05736* (2023).

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.

Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906* (2020).

Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*. PMLR, 5156–5165.

Geunwoo Kim, Pierre Baldi, and Stephen McAleer. 2023. Language Models can Solve Computer Tasks. *arXiv preprint arXiv:2303.17491* (2023).

Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451* (2020).

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems* 35 (2022), 22199–22213.

Gibbeum Lee, Volker Hartmann, Jongho Park, Dimitris Papailiopoulos, and Kangwook Lee. 2023. Prompted LLMs as Chatbot Modules for Long Open-domain Conversation. *arXiv preprint arXiv:2305.04533* (2023).

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. CAMEL: Communicative Agents for "Mind" Exploration of Large Scale Language Model Society. arXiv:2303.17760 [cs.AI]

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. 2023a. Encouraging Divergent Thinking in Large Language Models through Multi-Agent Debate. arXiv:2305.19118 [cs.CL]

Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, et al. 2023b. Taskmatrix. ai: Completing tasks by connecting foundation models with millions of apis. *arXiv preprint arXiv:2303.16434* (2023).

Jerry Liu. 2022. *LlamaIndex*. https://doi.org/10.5281/zenodo.1234

Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023b. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172* (2023).

Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang. 2023a. Reason for Future, Act for Now: A Principled Framework for Autonomous LLM Agents with Provable Sample Efficiency. *arXiv preprint arXiv:2309.17382* (2023).

Zhiwei Liu, Weiran Yao, Jianguo Zhang, Le Xue, Shelby Heinecke, Rithesh Murthy, Yihao Feng, Zeyuan Chen, Juan Carlos Niebles, Devansh Arpit, et al. 2023c. BOLAA: Benchmarking and Orchestrating LLM-augmented Autonomous Agents. *arXiv preprint arXiv:2308.05960* (2023).

Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. 2023. MemoChat: Tuning LLMs to Use Memos for Consistent Long-Range Open-Domain Conversation. arXiv:2308.08239 [cs.CL]

Ingo Molnár. 2007. Linux CFS Scheduler. https://docs.kernel.org/scheduler/sched-design-CFS.html.

Meredith Ringel Morris, Jascha Sohl-dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksandra Faust, Clement Farabet, and Shane Legg. 2023. Levels of AGI: Operationalizing Progress on the Path to AGI. *arXiv preprint arXiv:2311.02462* (2023).

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2022. WebGPT: Browser-assisted question-answering with human feedback. arXiv:2112.09332 [cs.CL]

OpenAI. 2023. GPT-4V(ision) System Card. (2023).

H. Orman. 2003. The Morris worm: a fifteen-year perspective. *IEEE Security & Privacy* 1, 5 (2003), 35–43. https://doi.org/10.1109/MSECP.2003.1236233

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.

Charles Packer, Vivian Fang, Shishir G Patil, Kevin Lin, Sarah Wooders, and Joseph E Gonzalez. 2023. MemGPT: Towards LLMs as Operating Systems. *arXiv preprint arXiv:2310.08560* (2023).

Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2023. Unifying Large Language Models and Knowledge Graphs: A Roadmap. *arXiv preprint arXiv:2306.08302* (2023).

Aaron Parisi, Yao Zhao, and Noah Fiedel. 2022. Talm: Tool augmented language models. *arXiv preprint arXiv:2205.12255* (2022).

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–22.

Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334* (2023).

Ofir Press, Noah A Smith, and Mike Lewis. 2021. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409* (2021).

Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. 2023a. Visual Adversarial Examples Jailbreak Aligned Large Language Models. arXiv:2306.13213 [cs.CR]

Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023b. Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! *arXiv preprint arXiv:2310.03693* (2023).

Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023a. Communicative agents for software development. *arXiv preprint arXiv:2307.07924* (2023).

Cheng Qian, Chi Han, Yi R Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. 2023b. CREATOR: Disentangling Abstract and Concrete Reasonings of Large Language Models through Tool Creation. *arXiv preprint arXiv:2305.14318* (2023).

Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. 2023a. Tool learning with foundation models. *arXiv preprint arXiv:2304.08354* (2023).

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023b. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789* (2023).

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.

Dennis M. Ritchie and Ken Thompson. 1974. The UNIX Time-Sharing System. *Commun. ACM* 17, 7 (jul 1974), 365–375. https://doi.org/10.1145/361011.361061

Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.

Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code Llama: Open Foundation Models for Code. arXiv:2308.12950 [cs.CL]

Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Xingyu Zeng, and Rui Zhao. 2023a. Tptu: Task planning and tool usage of large language model-based ai agents. *arXiv preprint arXiv:2308.03427* (2023).

Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J. Maddison, and Tatsunori Hashimoto. 2023b. Identifying the Risks of LM Agents with an LM-Emulated Sandbox. arXiv:2309.15817 [cs.AI]

Stuart Russell and Peter Norvig. 1995. *Prentice Hall series in artificial intelligence*. Prentice Hall Englewood Cliffs, NJ:.

Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184* (2023).

Gerard Salton. 1975. A vector space model for information retrieval. *Journal of the ASIS* (1975), 613–620.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761* (2023).

Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality Traits in Large Language Models. arXiv:2307.00184 [cs.CL]

Noam Shazeer. 2019. Fast transformer decoding: One write-head is all you need. *arXiv preprint arXiv:1911.02150* (2019).

Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*. PMLR, 31210–31227.

Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366* (2023).

Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2020. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768* (2020).

Krzysztof Sledziewski, Behzad Bordbar, and Rachid Anane. 2010. A DSL-Based Approach to Software Development and Deployment on Cloud. In *2010 24th IEEE International Conference on Advanced Information Networking and Applications*. 414–421. `https://doi.org/10.1109/AINA.2010.81`

I. Stoica and H. Abdel-Wahab. 1995. *Earliest Eligible Virtual Deadline First: A Flexible and Accurate Mechanism for Proportional Share Resource Allocation*. Technical Report. USA.

Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864* (2021).

Qiushi Sun, Zhangyue Yin, Xiang Li, Zhiyong Wu, Xipeng Qiu, and Lingpeng Kong. 2023. Corex: Pushing the Boundaries of Complex Reasoning through Multi-Model Collaboration. arXiv:2310.00280 [cs.AI]

Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. 2022. A length-extrapolatable transformer. *arXiv preprint arXiv:2212.10554* (2022).

Anirudh S Sundar and Larry Heck. 2023. cTBL: Augmenting Large Language Models for Conversational Tables. *arXiv preprint arXiv:2303.12024* (2023).

Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, and Le Sun. 2023. ToolAlpaca: Generalized Tool Learning for Language Models with 3000 Simulated Cases. *arXiv preprint arXiv:2306.05301* (2023).

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA model. `https://github.com/tatsu-lab/stanford_alpaca`.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).

UW:CSE451. 2023. History of Operating Systems. `https://courses.cs.washington.edu/courses/cse451/16wi/readings/lecture_readings/LCM_OperatingSystemsTimeline_Color_acd_newsize.pdf`.

Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2022. Large Language Models Still Can't Plan (A Benchmark for LLMs on Planning and Reasoning about Change). In *NeurIPS 2022 Foundation Models for Decision Making Workshop*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. 2023. Chatgpt for robotics: Design principles and model abilities. *Microsoft Auton. Syst. Robot. Res* 2 (2023), 20.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023e. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291* (2023).

Lei Wang, Jingsen Zhang, Xu Chen, Yankai Lin, Ruihua Song, Wayne Xin Zhao, and Ji-Rong Wen. 2023f. RecAgent: A Novel Simulation Paradigm for Recommender Systems. *arXiv preprint arXiv:2306.02552* (2023).

Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. 2020. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768* (2020).

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *The Eleventh International Conference on Learning Representations*.

Yancheng Wang, Ziyan Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Xiaojiang Huang, Yanbin Lu, and Yingzhen Yang. 2023c. RecMind: Large Language Model Powered Agent For Recommendation. *arXiv preprint arXiv:2308.14296* (2023).

Yubo Wang, Xueguang Ma, and Wenhu Chen. 2023d. Augmenting Black-box LLMs with Medical Textbooks for Clinical Question Answering. *arXiv preprint arXiv:2309.02233* (2023).

Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023a. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560* (2023).

Zhilin Wang, Yu Ying Chiu, and Yu Cheung Chiu. 2023b. Humanoid Agents: Platform for Simulating Human-like Generative Agents. arXiv:2310.05418 [cs.CL]

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How Does LLM Safety Training Fail? arXiv:2307.02483 [cs.LG]

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.

Michael Wooldridge and Nicholas R Jennings. 1995. Intelligent agents: Theory and practice. *The knowledge engineering review* 10, 2 (1995), 115–152.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2023. AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation. arXiv:2308.08155 [cs.AI]

Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, et al. 2023. Lemur: Harmonizing Natural Language and Code for Language Agents. *arXiv preprint arXiv:2310.06830* (2023).

Chenghao Yang and Allyson Ettinger. 2023. Can You Follow Me? Testing Situational Understanding in ChatGPT. *arXiv preprint arXiv:2310.16135* (2023).

Haomiao Yang, Kunlan Xiang, Hongwei Li, and Rongxing Lu. 2023c. A Comprehensive Overview of Backdoor Attacks in Large Language Models within Communication Networks. *arXiv preprint arXiv:2308.14367* (2023).

Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023b. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949* (2023).

Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023a. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381* (2023).

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems* 35 (2022), 20744–20757.

Shunyu Yao and Karthik Narasimhan. 2023. Language Agents in the Digital World: Opportunities and Risks. *princeton-nlp.github.io* (Jul 2023). `https://princeton-nlp.github.io/language-agent-impact/`

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601* (2023).

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022b. ReAct: Synergizing Reasoning and Acting in Language Models. In *The Eleventh International Conference on Learning Representations*.

Dingyao Yu, Kaitao Song, Peiling Lu, Tianyu He, Xu Tan, Wei Ye, Shikun Zhang, and Jiang Bian. 2023. MusicAgent: An AI Agent for Music Understanding and Generation with Large Language Models. arXiv:2310.11954 [cs.CL]

Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. Wordcraft: story writing with large language models. In *27th International Conference on Intelligent User Interfaces*. 841–852.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems* 33 (2020), 17283–17297.

Zhuosheng Zhan and Aston Zhang. 2023. You Only Look at Screens: Multimodal Chain-of-Action Agents. *arXiv preprint arXiv:2309.11436* (2023).

Kechi Zhang, Ge Li, Jia Li, Zhuo Li, and Zhi Jin. 2023. ToolCoder: Teach Code Generation Models to use APIs with search tools. *arXiv preprint arXiv:2305.04032* (2023).

Kang Zhao, Wei Liu, Jian Luan, Minglei Gao, Li Qian, Hanlin Teng, and Bin Wang. 2023. UniMC: A Unified Framework for Long-Term Memory Conversation via Relevance Representation Learning. *arXiv preprint arXiv:2306.10543* (2023).

Wanjun Zhong, Lianghong Guo, Qiqi Gao, and Yanlin Wang. 2023. MemoryBank: Enhancing Large Language Models with Long-Term Memory. *arXiv preprint arXiv:2305.10250* (2023).

Shuyan Zhou, Uri Alon, Frank F Xu, Zhengbao Jiang, and Graham Neubig. 2022. Docprompting: Generating code by retrieving the docs. In *The Eleventh International Conference on Learning Representations*.

Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. 2023. Agents: An Open-source Framework for Autonomous Language Agents. arXiv:2309.07870 [cs.CL]

Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. 2023. Ghost in the Minecraft: Generally Capable Agents for Open-World Enviroments via Large Language Models with Text-based Knowledge and Memory. *arXiv preprint arXiv:2305.17144* (2023).

Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and Transferable Adversarial Attacks on Aligned Language Models. arXiv:2307.15043 [cs.CL]