# GoEX: Perspectives and Designs Towards a Runtime for Autonomous LLM Applications

**Shishir G. Patil**    **Tianjun Zhang**    **Vivian Fang**    **Noppapon C.**    **Roy Huang**    **Aaron Hao**

**Martin Casado**[1]    **Joseph E. Gonzalez**    **Raluca Ada Popa**    **Ion Stoica**

UC Berkeley    [1]Andreessen Horowitz

`shishirpatil@berkeley.edu`

## Abstract

Large Language Models (LLMs) are evolving beyond their classical role of providing information within dialogue systems to actively engaging with tools and performing actions on real-world applications and services. Today, humans verify the correctness and appropriateness of the LLM-generated outputs (e.g., code, functions, or actions) before putting them into real-world execution. This poses significant challenges as code comprehension is well known to be notoriously difficult. In this paper, we study how humans can efficiently collaborate with, delegate to, and supervise autonomous LLMs in the future. We argue that in many cases, "post-facto validation"—verifying the correctness of a proposed action after seeing the output—is much easier than the aforementioned "pre-facto validation" setting. The core concept behind enabling a post-facto validation system is the integration of an intuitive *undo* feature, and establishing a *damage confinement* for the LLM-generated actions as effective strategies to mitigate the associated risks. Using this, a human can now either revert the effect of an LLM-generated output or be confident that the potential risk is bounded. We believe this is critical to unlock the potential for LLM agents to interact with applications and services with limited (post-facto) human involvement. We describe the design and implementation of our open-source runtime for executing LLM actions, Gorilla Execution Engine (GoEx), and present open research questions towards realizing the goal of LLMs and applications interacting with each other with minimal human supervision. We release GoEx at `https://github.com/ShishirPatil/gorilla/`.

## 1 Introduction

Large Language Models (LLMs) are evolving from serving knowledge passively in chatbots to actively interacting with applications and services. Enabling agents and software systems to interact with one another has given rise to new innovative applications. In fact, it is no longer science fiction to imagine that many of the interactions on the internet are going to be between LLM-powered systems. Agentic systems (Park et al., 2023; Wang et al., 2023b; Wu et al., 2023), co-pilots (Roziere et al., 2023), plugins (OpenAI, 2023), function calling and tool use (Achiam et al., 2023; Parisi et al., 2022; Patil et al., 2023; Qin et al., 2023; Yao et al., 2022), are all steps towards this direction.

*The logical next-step in this evolution is towards autonomous LLM-powered microservices, services, and applications. This paper is a first step towards realizing this goal, and addresses some of the key trustworthiness challenges associated with it.*

As a running example, consider an LLM-powered personal assistant that has access to a user's email account. The user asks the assistant to send an important email to their boss, but instead, the LLM sends a sensitive email to the wrong recipient. In designing such a system, several critical challenges must be addressed:

**Hallucination, stochasticity, and unpredictability.** LLM-based applications place an unpredictable and hallucination-prone LLM at the helm of a system traditionally reliant on trust. Currently, services and APIs assume a human-in-the-loop, or clear specifications to govern how and which tools are used in an application. In our running example, the user clicks the "Send" button after confirming the recipient and body of the email. In contrast, an LLM-powered assistant may send an email that goes against the user's intentions, and may even perform actions unintended by the user. LLMs are not only capable of being stochastic, but crucially, are capable of unpredictable and unbounded behavior even when trained not to do so (Anil et al., 2024).
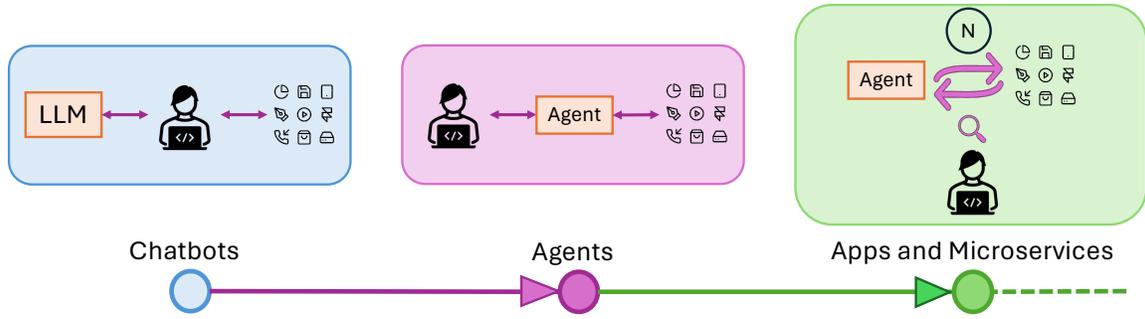
1

Figure 1: Evolution of LLMs powered applications and services from chatbots, to decision-making agents that can interact with applications and services with human-supervision, to autonomous LLM-agents interacting with LLM-powered apps and services with minimal and punctuated human supervision.

**Unreliability.** Given their unpredictability and impossibility to comprehensively test, it is difficult for a user to trust an LLM off the shelf. Consequently, LLM-powered applications are challenging for users to adopt as LLMs are untrusted components that would be running within a trusted execution context. Trivially, one can ensure safety by restricting the LLM to have no credentials, at the expense of losing utility. Given that prior work (Achiam et al., 2023; Patil et al., 2023; Schick et al., 2023) has surfaced the growing utility of LLM-based systems, mechanisms are needed to express the safety-utility tradeoff to developers and users.

**Delayed feedback and downstream visibility.** Lastly, from a system-design principle, unlike chatbots and agents of today, LLM-powered systems of the future will not have immediate human feedback. This means that the intermediate state of the system is not immediately visible to the user and often only downstream effects are visible. An LLM-powered assistant may interact with many other tools (e.g., querying a database, browsing the web, or filtering push notifications) before composing and sending an email. Such interactions before the email is sent are invisible to the user.

## 1.1 A Runtime that Enables Autonomous LLMs

The question this paper tries to address is: How do we enable an untrusted agent to take sensitive actions (e.g., code generation, API calls, and tool use) on a user's behalf and then verify that those actions aligned with the intent of that user's request?

Traditional methods such as using containers—which guarantee isolation through virtualization—falter when adjustments to the environment are required within the user's context. For example, a user might want to modify the state of their operating system. Further, isolation alone cannot ensure the final state aligns with the `intended` actions, especially under ambiguity or execution errors and ambiguity of stated intent. Both are common flaws when intent is being specified in natural-language as opposed to a more precise domain-specific language.

In this paper, we introduce the notion of **"post-facto LLM validation"** as opposed to "pre-facto LLM validation". While in both scenarios humans are the ultimate arbitrators, in "post-facto validation" human's arbitrate the output produced by executing the actions produced by the LLM, as opposed to the process or the intermediate outputs. A natural risk arising from "post-facto validation" is that the actions we execute may have unintended consequences. While the benefits of evaluating the output might justify the risks involved with unintended execution, we recognize that this could be a non-starter in many applications. To remedy this, we introduce the intuitive abstractions of "undo" and "damage confinement" or "damage confinement". The "undo" abstraction allows LLMs to back-track an action that may be unintended, for example, delete a message that was sent in Slack. And for those actions which may not have an "undo", we present "damage confinement" semantic. "Damage confinement" can be considered as a quantification of the user's risk appetite. For example, a user may tolerate the risk of the LLM agent delivering pizza to the wrong address, but perhaps they might not want to allow their LLM to interact with their bank.

As a step towards realizing our vision, we developed GoEx, a runtime for executing actions generated by LLMs. By designing GoEx to exclusively utilize readily available off-the-shelf software components, we aim to assess the readiness of current resources and provide an ecosystem to empower developers.

In summary, we make the following contributions:

1. We first make the case for the future of deeply-embedded LLM systems with LLMs powering microservices, applications, etc. In this paradigm, LLMs are not merely used for information compression, but also as decision makers (§2). We identify the key challenges associated with integrating LLMs into existing systems, including the inherent unpredictability of LLMs, the lack of trust in their execution, and the difficulty in detecting and mitigating their failures in real-time (§3).

2. We introduce the concept of "post-facto LLM validation" as an approach to ensuring the safety and reliability of LLM-powered systems, focusing on validating the results of LLM-generated actions rather than the process itself (§4.1).

3. We introduce "undo" (§4.1.1) and "damage confinement" (§4.1.2) abstractions as mechanisms for mitigating the risk of unintended actions taken in LLM-powered systems.

4. We propose the *Gorilla Execution Engine (GoEx)*, a runtime designed to enable the autonomous interactions of LLM-powered software systems and agents by safely executing LLM-generated actions and striking a tradeoff between safety and utility (§5).

## 2  EVOLUTION OF LLM POWERED AGENTS

We first present the background on the evolution of LLM-powered-systems to mean applications, microservices, and other systems that integrate- or interface- with LLMs. Based on this trend, we then speculate on what such a future would look like.

### 2.1  LLM-HUMAN INTERACTION: CHATBOTS, AND SEARCH

LLMs have transformed the landscape of human-computer interaction. With early adoption as chatbots, these models were designed to mimic human conversation, allowing users to interact with computers in natural language. This era of LLMs focused on understanding and generating answers, serving as a bridge for humans to interact with vast amounts of web data in a more intuitive way. Early implementations were primarily used in the *read-only* model for information retrieval where the LLM did not make any stateful changes (Burns et al., 2022), customer service agents in enterprises, and as educational tools, laying the foundation for more sophisticated applications.

### 2.2  FROM CHATBOTS TO AGENTS: THE RISE OF ACTIONABLE LLMS

Increasing trustworthiness (Wang et al., 2023a), the availability of adapters, and novel-techniques (Patil et al., 2023; Schick et al., 2023), have expanded the role of LLMs from passive providers of information to active agents capable of executing simple tasks. These agents, powered by LLMs, can interact with applications, services, and APIs to perform actions on behalf of the user. This shift represents a significant leap in the capabilities of LLMs, enabling them to contribute actively to workflows and processes across various domains.

However, this evolution also brings challenges. The complexity of understanding context, intent, and the subtleties of human language make it difficult to ensure the accuracy and appropriateness of the actions taken by LLMs. As a result, *human oversight* remains crucial to manage and validate the set of actions proposed by the LLM.

### 2.3  TOWARDS UBIQUITOUS LLM INTEGRATION

Looking ahead, we expect the integration of LLMs into daily workflows and systems to deepen. This future envisions LLMs not just as tools or assistants but as pervasive agents embedded in a myriad of workflows, enhancing functionality and adaptability across the board. This vision is materialized through the development of advanced LLM-powered microservices, LLM-powered applications, and LLM-powered workflows all interacting with each other, constantly, with limited-to-no human interaction.

From an application developer's perspective, these can be categorized into personalized systems, hosted agents for collective use, and third-party integrations.

**Personalized LLM-powered workflows for individuals.**  In the personal domain, LLMs are anticipated to become deeply integrated with individual user experiences, offering tailored assistance that understands

and anticipates the unique preferences and needs of each user. Imagine a personalized LLM-powered version of Siri or Google Voice, not merely responding to queries but proactively managing schedules, filtering information based on user preferences, and even performing tasks across a range of applications and services. Such personalized systems would mark a significant departure from generic assistants to truly personalized digital companions.

**Hosted agents for enterprise and group applications.**   Within professional and enterprise environments, hosted LLM-powered agents will take on specific roles, such as managing database queries or automating routine administrative tasks, tailored to the unique needs of an organization. These specialized agents would be operating within the confines of an organization but would serve a wide set of users.

**Third-party agents: expanding the ecosystem of services.**   The expansion of LLM capabilities is also expected to include extensive collaboration with third-party service providers (e.g., Slack, Gmail, Dropbox, etc.), enabling seamless interactions between users and services through personalized LLM workflows. These third-party agents would allow users to communicate with and through them using customized LLM-powered workflows, and integrate a wide range of services.

## 3   NAVIGATING THE NEW FRONTIER: CHALLENGES IN UBIQUITOUS LLM DEPLOYMENTS

With advancements in LLM capabilities and their applications come new challenges: How do we contend with the inaccuracies inherent to even the best models currently available? How do we ensure security for client information? How do we handle system reliability and quantify risk tolerance to a user? In the following sections, we elaborate on these new challenges.

### 3.1   DELAYED SIGNALS

**The challenge of timely feedback.** In traditional software development, immediate feedback through error messages or direct outputs enables quick recovery mechanisms to be triggered. However, when embedding LLMs, especially in complex systems or applications interfacing with real-world data and actions, feedback can be significantly delayed. This delay in obtaining relevant signals to assess the performance or correctness of LLM actions (especially with text-in, text-out modality) introduces challenges in rapidly iterating and refining model outputs.

**Impact on system development.** The lag between action and feedback complicates the identification of errors and the assessment of system performance, potentially leading to additional state being built on top of the system. This necessitates designing systems that can accommodate these delays and implement strategies for asynchronous feedback collection.

### 3.2   AGGREGATE SIGNALS

In the realm of LLMs, particularly when applied to large-scale systems or microservices, individual LLM actions may not provide clear insights into overall system performance, nor assist in diagnosing the cause of the error. Instead, the true measure of success (or failure) often emerges from aggregated outcomes, necessitating a shift in how developers and stakeholders evaluate and account for, in LLM-driven applications.

### 3.3   THE DEATH OF UNIT-TESTING AND INTEGRATION-TESTING

The integration of LLMs into software systems challenges traditional paradigms of unit testing and integration testing. While closed-source and continuous-pre-trained LLMs from third-parties pose a new challenge of the model changing constantly, in-house LLMs are not a panacea either. Given the dynamic, often unpredictable nature of LLM outputs, establishing a fixed suite of tests that accurately predict and verify all potential behaviors becomes increasingly difficult, if not impossible.

## 3.4 VARIABLE LATENCY

LLMs' auto-regressive text-generation by the very nature means inference time may vary as the LLM can either output a long or short response. This is an important consideration for hard-deadline real-time systems (RTS) (Quigley et al., 2009).

## 3.5 PROTECTING SENSITIVE DATA

In order for LLMs to interact with a user's accounts across multiple applications, the LLM must be able to reason about having access to credentials granting access to the user's accounts. When an LLM is hosted by an (untrusted) external service, it is desirable to not directly pass any credentials or sensitive data to the LLM while still preserving functionality of the LLM-powered system. LLMs can also generate (untrusted) code and are susceptible to hallucinations (Rawte et al., 2023; Zhang et al., 2023), which can result in running potentially malicious code or inadvertently performing actions that were not intended by the user.

## 4 DESIGNING A RUNTIME

To address the challenges introduced by the new paradigm of LLM-powered applications, we propose an *LLM runtime*—designed specifically to execute actions generated by LLMs—as a compelling solution. This section discusses the necessity of such a runtime, its envisioned properties, and strategies to mitigate the risks associated with LLM-generated actions.

**LLMs started the problem, can LLMs solve the problem?**   Despite the rapid advancements in LLMs, expecting LLMs to self-correct and eliminate all potential errors or unintended actions through existing training techniques——such as pre-training, instruction tuning, DPO (Rafailov et al., 2024), or RLHF (Ziegler et al., 2019)—is promising (§5) but challenging. The challenges are manifold, primarily due to ill-defined metrics and the inherent complexities of accurately predicting the real-world impacts of actions suggested by LLMs. Thus, while LLMs are at the heart of these challenges, their current evolutionary trajectory suggests they cannot entirely solve the problem without external frameworks to guide their execution.

## 4.1 POST-FACTO LLM VALIDATION

In the realm of LLM-powered-systems, we introduce "post-facto LLM validation," which contrasts with traditional "pre-facto" methods. In "post-facto validation," humans evaluate the outcomes of actions executed by the LLM, rather than overseeing the intermediate processes. This approach assumes that validating results over processes, acknowledging that while verifying outcomes is crucial, understanding and correcting processes based on those outcomes is equally important.

Forgoing "pre-facto validation" means execution of actions without prior validation, which introduces risks and potentially leads to undesirable outcomes. We propose two abstractions to mitigate the risk associated with post-facto validation: undoing an action (§4.1.1), and damage confinement (§4.1.2).

## 4.1.1 REVERSIBILITY

When possible, actions executed by an LLM should give users the right to *undo* an action. This approach may require maintaining multiple versions of the system state, leading to high costs in terms of memory and computational resources. Furthermore, the feasibility of implementing undoing an action is often dependent on the level of access granted to the system. For instance, in file systems or databases, where root access is available, undoing actions is possible. However, in scenarios where such privileged access is not granted, such as in email clients like Gmail, the ability to undo an action may be limited or require alternative approaches. One potential solution is for the runtime to make a local copy of the email before deleting, which introduces additional state to the runtime but enables *undo* for email deletion.

To account for the resource costs associated with maintaining multiple system states, we adopt the notion of a *commit*, also called a "watermark" in streaming data-flow systems (Akidau et al., 2015; Carbone et al., 2015). By grouping together sets of actions based on their associativity, commutativity, and distributive properties, it may be possible to define checkpoints at which the system state can be saved or rolled back. This approach would enable selective undoing of actions within a defined scope, rather than maintaining the ability to undo every individual action.

**Atomiticy.** In agent-driven systems, the option for users to demand atomicity of operations can be crucial. Atomicity ensures that either all of the operations within a task are successfully completed, or, in the event of a failure in any step, the system is reverted to its initial state before any operation was applied. This binary outcome—success or a clean-slate reset—provides a clear, predictable framework for managing tasks, increasing the system's reliability and user trust in the LLM agent executing complex sequences of actions.

### 4.1.2 DAMAGE CONFINEMENT

Not all applications or tools provide the ability to undo an action. For example, emails currently cannot be unsent after some time has elapsed. In scenarios like these, we fall back to "damage confinement" or "blast-radius confinement", as it is necessary to provide users with mechanisms to quantify and assess the associated risks of the actions their LLM-powered application may take.

One approach to address this challenge is through the implementation of coarse-grained access control mechanisms. A user could permit their LLM to only read emails instead of sending emails, thus confining the blast radius to an tolerable level. Such permissioning has already been explored preliminarily by Wu et al. (2024), in the context of a user authorizing independent LLM applications to interact with one another.

### 4.2 SYMBOLIC CREDENTIALS AND SANDBOXING

As mentioned in §3.5, the LLM could be (1) hosted by an (untrusted) external provider and (2) susceptible to hallucinations (Rawte et al., 2023; Zhang et al., 2023), resulting in code that may be unsafe to execute.

In order to protect sensitive user information from and untrusted LLM, the sensitive information in the input prompt can be substituted with a symbolic credential (e.g., a dummy API key), similar to the anonymization approach taken by Presidio (Microsoft), and sending this sanitized prompt to the LLM. Then, the LLM will never see the user's sensitive information in its input.

We can mitigate this risk of running potentially unsafe code by executing the generated code in a sandboxed environment, whether it be a container or a bare-metal VM. With this approach, we only grant the code access to the required dependencies and necessities, such as a specific API key used for service access, and nothing more. If utilizing a container, we only mount the necessary files and impose appropriate network restrictions.

### 4.3 STORING KEYS AND ACCESS CONTROL

As LLMs are inherently untrusted, a user would feel uneasy permitting the LLM based system to store their credentials. There are two challenges an LLM runtime should address: how to store and manage the user's credentials; and determining and projecting the minimal set of permissions that an LLM needs in order to accomplish its task. More formally, mapping an action to the least privileges that need to be granted to perform the task. While solutions such as asking a ML model have the benefits of generalizability and scalability, pre-computing this permission set manually provides strong security guarantees. Finding a common ground between these two techniques, remains an interesting open area for research. Further, in an enterprise use case in which an LLM-powered application is managing many user credentials, recording an audit trail of credential access is critical.

## 5 GoEx: LLM RUNTIME

GoEx represents a first step towards building a runtime for executing LLM-generated actions within a secure and flexible runtime environment. Central to GoEx are its abstractions for "undo" (§4.1.1) and "damage confinement" or "blast-radius confinement" (§4.1.2), which provide developers of apps and services the flexibility to tailor policies to their specific needs, recognizing the impracticality of a one-size-fits-all policy given the varied contexts in which LLMs are deployed. GoEx supports a range of "actions" including RESTful API requests (§5.1), databases operation (§5.2), and filesystem actions (§5.3). Each action type, while initiated from a unified GoEx interface, are handled uniquely as described below.

### 5.1 RESTFUL API CALLS

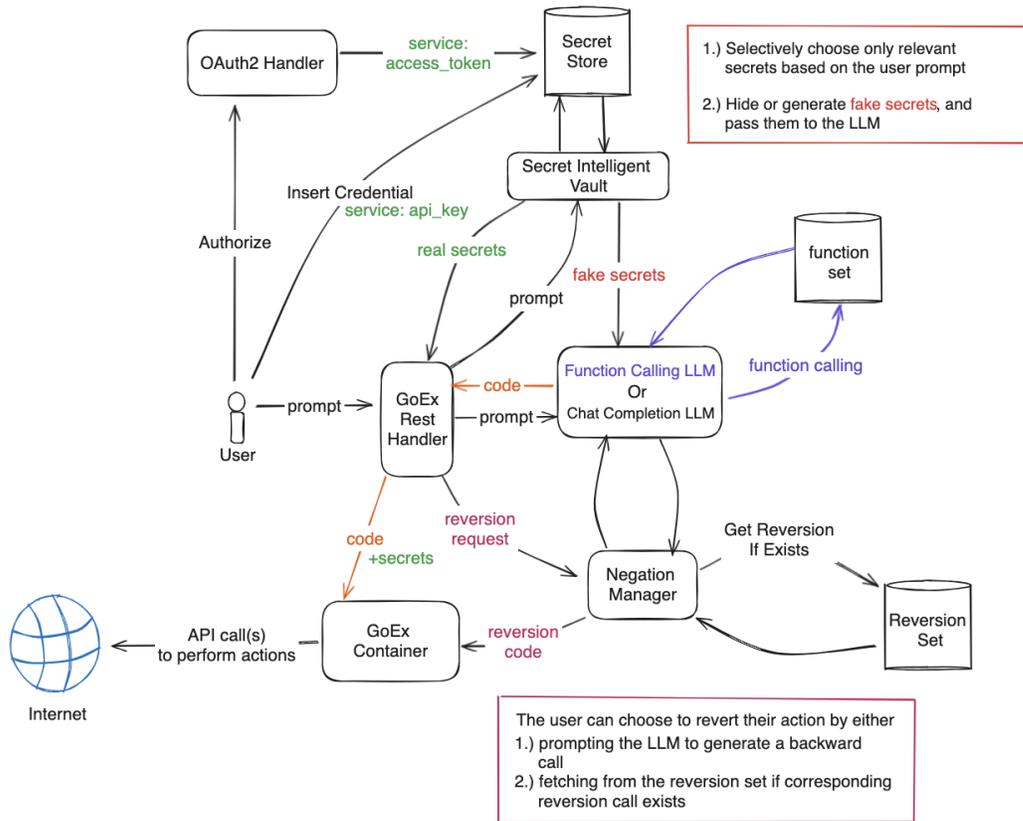We first describe how GoEx handles RESTful API calls (illustrated in Fig. 2).

Figure 2: GoEx's runtime for executing RESTful API calls. Upon receiving the user's prompt, GoEx presents two alternatives. First, an LLM can be prompted to come up with the (`Action, Undo-Action`) pair. Second, the application developer can provide tuples of actions and their corresponding undo-actions (function calls) from which the LLM can pick amongst.

**Authentication.** GoEx provides a secure way to handle user secrets, whether using OAuth2 for token-based authentication or API keys for direct service access. GoEx acts as the secure intermediary to facilitate authenticated actions across various services. For OAuth2, GoEx sits between the user and services, facilitating the necessary relay to retrieve access tokens. These tokens allow users to delegate the GoEx system to perform actions on their behalf. For other services that authenticate accounts through API keys, GoEx provides an interface that allows users to insert and retrieve them.

**Storing secrets.** User secrets and keys are stored locally on the user's device in a Secret Intelligent Vault (SIV). SIV maps `service_name` to `key` and `format`. When user wishes to interact with specific service(s), the corresponding keys are requested from the SIV. The `format` specifies how the keys are store, that is, in a file, or as a string, etc. The role of the SIV is to selectively retrieve just the required keys for a given execution. For example, if a user wants to send an email invite to their friend for lunch, the agent only needs their OAuth2 token for their email provider, and not, for example, their bank account's API keys. The policy used for SIV is user-defined and highly flexible; it could be as simple as parsing through the user prompt to detect which service's keywords are present, or as complex as a fine-tuned prompt-to-service retrieval model.

**Generating actions.** The GoEx framework pupports two techniques to generate the APIs. In the Chat Completion case, assuming the user prompt is, "send a Slack message to `gorilla@yahoo.com`," the user must initially authorize GoEx to use their access token through the Slack browser. After receiving the user prompt, GoEx requests the SIV for the necessary secrets from the Secret Store. Slack secrets (OAuth2) are inherently hidden because they are stored as a file, so GoEx passs the file path along with the prompt directly to the LLM. GoEx mounts the Slack secret file and passes the LLM-generated code to be executed in the GoEx container. If the user wishes to revert the execution, the reversion call will be retrieved from the reversion set if it exists; otherwise, the handler prompts the LLM to generate it. If the user chooses Function Calling, instead of asking the LLM to come up with a command to satisfy the user's prompt, GoEx asks it to select a function from a user-defined function set and populate the arguments. Secrets will be chosen from the SIV similarly, and execution occurs in the GoEx container. If the user wishes to revert, another function from the function set will be chosen by the LLM.

**Generating undo actions.** Identifying the 'undo' action for RESTful APIs, includes the following steps. First, we check if the reverse call for the action API is in the database `Reversion Set` as shown in figure 2. GoEx presents the systems abstractions, while developers are free to define the policies for mapping. For some APIs it might be critical to check for exact match for all parameters of the API, on the other hand for some other APIs, perhaps just the API name might be sufficient to uniquely identify what the reverse API would be. For example it is *not* sufficient to say the reverse of `send_slack_message` is `delete_slack_message`, since number of messages to be deleted could be one of the arguments.

To populate such a mapping, first, we instruct the LLM to generate a reverse API call whenever the user attempts to perform an action. We recognize that this gives no guarantees, but the philosophy is that we allow the LLM to be wrong at most once. Post each new API, the table is then if the reversion worked or not making this information available for future invocations. For applications that need guarantee, developers can pre-populate this table and combined with function-calling mode of operation, the system can be forced to only use those API's that are 'guaranteed' by the developers to be reversible.

**Damage confinement.** Often reversibility cannot be guaranteed. For examples sending an email isn't really reversible. For such scenarios, GoEx presents abstraction to bound the worst case. Currently, the way blast-radius-containment is implemented is through coarse-grained access control, and exact string match. First, GoEx looks at the user's prompt to determine the end service that they are then authorized to use. For example, a prompt of *I would like to send a slack message* would only need credentials for slack, and not, say, their bank. GoEx currently does this, through a simple sub-string check of the prompt, while giving developers the flexibility to adopt any mapping they might choose.

**Execution.** Once the API, and the set of credentials required are determined, the APIs are then executed in a Docker container for isolation.

## 5.2 DATABASE OPERATIONS

GoEx leverages the mature transaction semantics offered by databases. This section describes the abstractions available, and the two default policies.

### 5.2.1 ABSTRACTIONS

GoEx relies on the LLM to generate database operations, but there are two prerequisites needed to execute database operations: (1) knowledge of the current database state, and (2) knowledge on how to access the database. To provide these, `DBManager` class is used. This allows the database to readily minimally query for the database state (e.g. only the schema) to provide additional info to the LLM during prompting without leaking sensitive data. It also tracks the connection configuration to the database so that connections can be established without leaking credentials to the LLM as an untrusted third-party by asking the user to store the credentials locally, and after the LLM generates the operation, GoEx then executes the operation.

`DBManager` also assists the user store with storing a previous state. Here, the *commit* and *undo* actions are introduced where a *commit* means the user permanently saves the executed changes, and an *undo* reverses the aforementioned changes. Most modern databases also provide ACID guarantees (Haerder & Reuter, 1983), including NoSQL databases like DynamoDB and MongoDB, which we leverage to implement committing and undoing actions.
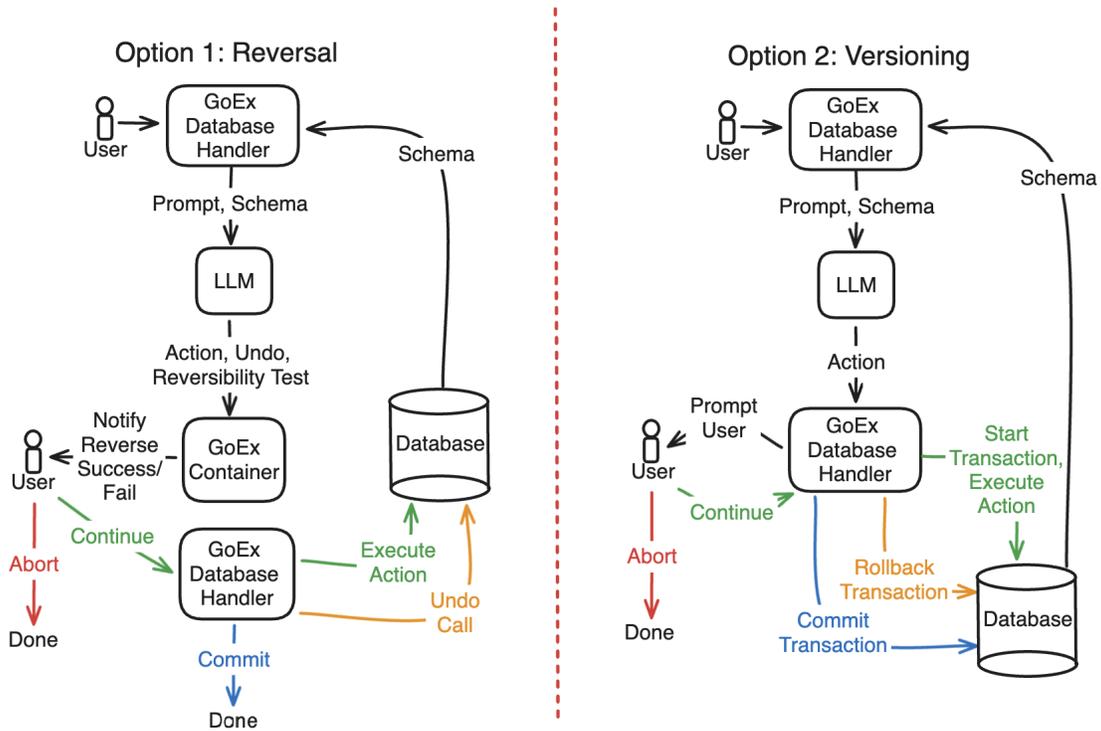
Figure 3: Runtime for executing actions on a database. We present two techniques to determine if a proposed action can be undone. On the left, for non-transactional databases like MongoDB, and for flexibility, we prompt the LLM to generate (`Action`, `Undo-Action`, `test-bed`) tuples, which we then evaluate in a isolated container to catch any false (`Action`, `Undo-Action`) pairs. On the right, we can provide a deterministic undo with guarantees by employing the transaction semantics of databases.

### 5.2.2 POLICY

`DBManager` implements reversibility in two ways. The user chooses which one to use when they execute a prompt in GoEx.

1. **Option 1 (Reversal).** Makes use of a reverse database operation to perform the *undo*. It is done by prompting the LLM with the original operation (action call) along with the schema to generate the reversal operation (undo call). Committing would require no action, and undoing would just be performing the undo call after the action call is done. This option scales better as additional users can continue perform database actions without needing to wait for the previous user to finish their transaction at the cost of relying on the LLM to come up with an undo call, which may or may not have unexpected behaviors.

2. **Option 2 (Versioning).** Makes use of the traditional ACID transaction guarantees of the database and holds off on completing a transaction until the user specifies to do so, or rolls back to the previous state. Committing would involve committing the transaction, and undoing is synonymous to a rollback transaction. This branch is able to provide reversal guarantees that branch 1 cannot, at the expense of higher performance overhead.

**Reversibility testing.** Within Option 1, GoEx also performs a reversibility test to verify that the generated reversal operation indeed reverses the original operation. This requires a containerized environment to be separate from the original database to maintain the original database state. Since copying over the database into the container is very expensive, the approach is to ask the LLM to generate a bare-bones version of the database for reversibility testing, given the action, undo calls, and the database schema. The outcome of the test is sent back to the user for final confirmation before committing or undoing the operation. This method allows for efficient testing by decoupling the testing runtime from being scaled by the number of entries in the database.
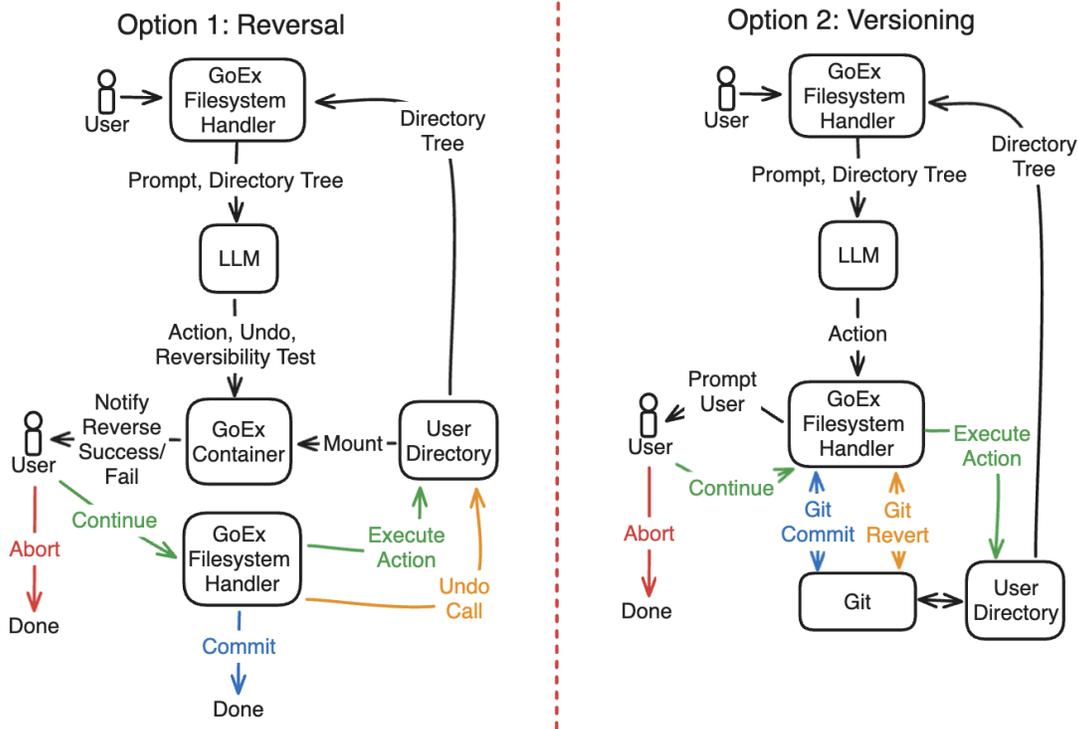
Figure 4: Runtime for executing actions on a filesystem. GoEx presents two abstractions. On the left, the LLM is prompted to come up with an (`Action`, `Undo-Action`, `test-bed`) which GoEx evaluates in a isolated container to catch any false (`Action`, `Undo-Action`) pairs. On the right presents deterministic guarantees by using versioning control system like Git or Git LFS.

## 5.3 FILE SYSTEMS

GoEx tries to present expressive abstractions to let LLM-powered systems to interact with file-systems using Git version control. To track the directory tree, on every GoEx filesystem-type execution, GoEx does an exhaustive, recursive walk of the directory and its subdirectories and stores the directory structure as a formatted string.

### 5.3.1 ABSTRACTIONS

Filesystems operation support in GoEx uses abstractions similar to what is used to support database operations. `FSManager`, is a filesystems manager that tracks (1) the directory tree structure with all filenames, and (2) the directory path that the user wishes to execute the filesystem's operations in. The tree structure, which is updated with executions, enables the LLM to generate operations that reflect the actual state of the user's filesystem.

Utilizing the relevant abstractions presented by journaling and log-structured filesystem for undo-semantics is left as future work, as the current GoEx system aims for compatibility.

### 5.3.2 POLICY

The options are similar to the database case, where Option 1 is for reversals and Option 2 is for versioning. The largest differences are how `FSManager` carries out reversibility testing and that versioning is accomplished using Git.

**Git.** GoEx uses Git to perform versioning. Since Git is already a version-control system for files, it is a straightforward solution to use, but has several limitations. Git does not have the ability to version track outside of the directory that it was initialized in. GoEx limits the user execution scope to the specified path in `FSManager`—which is always inside of a Git repository—and its subdirectories in accordance to our blast-radius confinement abstraction to prevent the LLM from performing arbitrary actions in undesired parts of the user's system. With larger directories, Git versioning can be expensive space-wise. GoEx leverages Git LFS for larger directories as an optimization. A threshold is defined for directory size that GoEx would then check whether or not to initialize Git LFS (200 MB by default).

**Reversibility testing.** Similar to supporting databases operations, the LLM generates the testing code using the action and undo calls, along with the directory tree. Inside the container, the specified path is mounted in read-only mode to again do blast radius containment. GoEx begins by duplicating the directory contents in the container, then run the action and undo calls on the copied directory, and finally compare contents. Depending on the original operation, the content comparison can just be a check of filenames or an exhaustive file content comparison of all the files. We rely on the LLM to come up with the test-case. Unsurprisingly, here GoEx allows you to trade off guarantees for performance.

## 6  DISCUSSION

### 6.1  IS POST-FACTO LLM VALIDATION ALWAYS PREFERABLE?

It is not lost on us that while post-facto LLM validation has many benefits, our advocacy of it is also somewhat philosophical. For example, if one were to bake a cake it's probably better to taste the cake than check the recipe. But on the the contrary, if one were to produce an audit report, it might be preferable to check the process. We acknowledge both options—post-facto LLM validation, and pre-facto LLM validation—as two techniques to evaluate the LLM's actions, however, this paper focuses on verifying results as this is, perhaps, more appropriate in most scenarios we consider, which are complex microservice settings.

### 6.2  DESIGNING LLM-FRIENDLY APIS

The conversation around LLM-powered systems design is predominantly centered around designing systems to conform with the API semantics of existing applications and services. However, an equally interesting question is what API design in an LLM-centric world would looks like.

LLMs introduce a paradigm where applications and services can anticipate and adapt to the intricacies of LLM interactions. A notable feature that embodies this adaptability is the implementation of "dry-run" semantics, akin to the functionality commonly visible in infrastructure products such as AWS, Kubernetes, and where API calls can be tested to predict their success without executing any real changes. This concept can be extended beyond mere prediction, serving as a bridge between LLMs' proposed actions and user consent. By repurposing "dry-run" operations, service providers can offer a preview of the uncommitted state resulting from an LLM's actions, allowing users to evaluate and approve these actions before they are finalized. This process adds an essential layer of user oversight, ensuring that actions align with user expectations and intentions.

**Chaining-aware.** Applications and Services should by-default expect their APIs to be chained with each other when used by agents. To support such a scenario, there needs to be a way to express which APIs can be commutative, associative or distributive with a given set of APIs.

### 6.3  TRACKING LLM AGENTS

The introduction of a nonce mechanism (i.e., a session identifier) would enable LLMs to present their identity and facilitate smoother interactions with API providers. This could serve various purposes, such as identifying a session initiated by an LLM or providing a context for transactions. A transaction ID, for example, can enable a system to identify and potentially rollback actions based on this ID. This not only improves the traceability of interactions but also contributes to the overall robustness and reliability of the system by providing an auditable framework for LLM-powered systems.

# 7 RELATED WORK

**Isolation.** Prior work on enabling automated LLM-powered systems (Wu et al., 2024) draws on concepts from existing computer systems (Cohen & Jefferson, 1975; Linden, 1976; Reis et al., 2019; Wilkes & Needham, 1979) and emphasizes isolation between LLM-powered applications in order to secure the overall system. Isolation is one facet of safely executing LLM-powered systems, as it alone cannot ensure that the final execution outcome aligns with the user's intended action.

**Trusthworthiness in LLMs.** There is a rich body of work benchmarking LLMs on their robustness (Chang et al., 2023). Recently, trustworthiness (Wang et al., 2023a) has been introduced as a multifaceted benchmark encapsulating robustness, stereotype bias, toxicity, privacy, ethics, and fairness. Wang et al. (2023a) found that more advanced LLMs (e.g., GPT-4) exhibit higher, albeit still imperfect, trustworthiness.

**Attacks on LLMs and defenses.** LLMs are also susceptible to prompt hacking attacks, of which include prompt injection (Greshake et al., 2023; Liu et al., 2023; Perez & Ribeiro, 2022; Schulhoff et al., 2023; Yu et al., 2023) and jailbreaking (Anil et al., 2024; Kang et al., 2023). Such attacks can lead to unpredictable and malicious decisions made by the LLM. There is also an active line of work on defending against such attacks on LLMs (Chen et al., 2024; Piet et al., 2023; Suo, 2024; Toyer et al., 2023; Yi et al., 2023). These attacks and defenses will continue to evolve, and consequently the potential of LLMs being susceptible to having their trustworthiness undermined necessitates a runtime that can provide execution of LLM-decided actions while limiting risk.

# 8 CONCLUSION

The evolution of LLMs from chatbots to deeply embedding them in applications and services for autonomous operation among themselves and other agents presents an exciting future. In this paper, we introduce the concept of "post-facto LLM validation," as opposed to pre-facto LLM validation, to enable users to verify and roll back the effects caused by executing LLM generated actions (e.g., code, API invocations, and tool use). We propose GoEx, a runtime for LLMs with an intuitive undo and damage confinement abstractions, enabling the safer deployment of LLM agents in practice. We hope our attempt to formalize our vision and present open research questions towards realizing the goal of autonomous LLM-powered systems in the future, is a step towards a world where LLM-powered systems can independently, with minimal human verification, interact with other tools and services.

## REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Tyler Akidau, Robert Bradshaw, Craig Chambers, Slava Chernyak, Rafael J Fernández-Moctezuma, Reuven Lax, Sam McVeety, Daniel Mills, Frances Perry, Eric Schmidt, et al. The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing. *VLDB*, 2015.

Cem Anil, Esin Durmus, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Nina Rimsky, Meg Tong, Jesse Mu, Daniel Ford, et al. Many-shot jailbreaking. 2024.

AWS. Testing your AWS KMS API calls. `https://docs.aws.amazon.com/kms/latest/developerguide/programming-dryrun.html`.

Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. *arXiv preprint arXiv:2212.03827*, 2022.

Paris Carbone, Asterios Katsifodimos, Stephan Ewen, Volker Markl, Seif Haridi, and Kostas Tzoumas. Apache Flink: Stream and batch processing in a single engine. *The Bulletin of the Technical Committee on Data Engineering*, 2015.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM TIST*, 2023.

Sizhe Chen, Julien Piet, Chawin Sitawarin, and David Wagner. StruQ: Defending against prompt injection with structured queries. *arXiv preprint arXiv:2402.06363*, 2024.

Ellis Cohen and David Jefferson. Protection in the Hydra operating system. *SOSP*, 1975.

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. Not what you've signed up for: Compromising real-world LLM-integrated applications with indirect prompt injection. In *ACM AISec*, 2023.

Theo Haerder and Andreas Reuter. Principles of transaction-oriented database recovery. *ACM CSUR*, 1983.

Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. Exploiting programmatic behavior of LLMs: Dual-use through standard security attacks. *arXiv preprint arXiv:2302.05733*, 2023.

Kubernetes. kubectl usage conventions. `https://kubernetes.io/docs/reference/kubectl/conventions/`.

Theodore A Linden. Operating system structures to support security and reliable software. *CSUR*, 1976.

Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, and Yang Liu. Prompt injection attack against LLM-integrated applications. *arXiv preprint arXiv:2306.05499*, 2023.

Microsoft. Presidio - data protection and de-identification SDK. `https://github.com/microsoft/presidio`.

OpenAI. ChatGPT plugins, 2023. `https://openai.com/blog/chatgpt-plugins`.

Aaron Parisi, Yao Zhao, and Noah Fiedel. TALM: Tool augmented language models, 2022.

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *UIST*, 2023.

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

Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*, 2022.

Julien Piet, Maha Alrashed, Chawin Sitawarin, Sizhe Chen, Zeming Wei, Elizabeth Sun, Basel Alomair, and David Wagner. Jatmo: Prompt injection defense by task-specific finetuning. *arXiv preprint arXiv:2312.17673*, 2023.

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master 16000+ real-world APIs, 2023.

Morgan Quigley, Ken Conley, Brian Gerkey, Josh Faust, Tully Foote, Jeremy Leibs, Rob Wheeler, Andrew Y Ng, et al. Ros: an open-source robot operating system. In *ICRA workshop on open source software*, volume 3, pp. 5. Kobe, Japan, 2009.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *NeurIPS*, 2024.

Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.

Charles Reis, Alexander Moshchuk, and Nasko Oskov. Site isolation: Process separation for web sites within the browser. In *USENIX Security*, 2019.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code Llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.

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

Sander Schulhoff, Jeremy Pinto, Anaum Khan, Louis-François Bouchard, Chenglei Si, Svetlina Anati, Valen Tagliabue, Anson Liu Kost, Christopher Carnahan, and Jordan Boyd-Graber. Ignore this title and hackaprompt: Exposing systemic vulnerabilities of LLMs through a global scale prompt hacking competition. *arXiv preprint arXiv:2311.16119*, 2023.

Xuchen Suo. Signed-Prompt: A new approach to prevent prompt injection attacks against LLM-integrated applications. *arXiv preprint arXiv:2401.07612*, 2024.

Sam Toyer, Olivia Watkins, Ethan Adrian Mendes, Justin Svegliato, Luke Bailey, Tiffany Wang, Isaac Ong, Karim Elmaaroufi, Pieter Abbeel, Trevor Darrell, et al. Tensor trust: Interpretable prompt injection attacks from an online game. *arXiv preprint arXiv:2311.01011*, 2023.

Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. DecodingTrust: A comprehensive assessment of trustworthiness in GPT models. *arXiv preprint arXiv:2306.11698*, 2023a.

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

Maurice Vincent Wilkes and Roger Michael Needham. The Cambridge CAP computer and its operating system. 1979.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen LLM applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.

Yuhao Wu, Franziska Roesner, Tadayoshi Kohno, Ning Zhang, and Umar Iqbal. SecGPT: An execution isolation architecture for llm-based systems. *arXiv preprint arXiv:2403.04960*, 2024.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

Jingwei Yi, Yueqi Xie, Bin Zhu, Keegan Hines, Emre Kiciman, Guangzhong Sun, Xing Xie, and Fangzhao Wu. Benchmarking and defending against indirect prompt injection attacks on large language models. *arXiv preprint arXiv:2312.14197*, 2023.

Jiahao Yu, Yuhang Wu, Dong Shu, Mingyu Jin, and Xinyu Xing. Assessing prompt injection risks in 200+ custom gpts. *arXiv preprint arXiv:2311.11538*, 2023.

Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. Siren's song in the AI ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219*, 2023.

Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.