

Do LLMs Dream of Discrete Algorithms?

Claudionor N. Coelho Jr, PhD/MBA^{1,2}, Yanen Li, PhD¹, and
Philip Tee^{1,3,4}

¹Zscaler Inc

²ECE Department, Santa Clara University

³Department of Informatics, University of Sussex

⁴The Beyond Center for Fundamental Science, Arizona State
University

July 1, 2025

Abstract

Large Language Models (LLMs) have rapidly transformed the landscape of artificial intelligence, enabling natural language interfaces and dynamic orchestration of software components. However, their reliance on probabilistic inference limits their effectiveness in domains requiring strict logical reasoning, discrete decision-making, and robust interpretability. This paper investigates these limitations and proposes a neurosymbolic approach that augments LLMs with logic-based reasoning modules, particularly leveraging Prolog predicates and composable toolsets. By integrating first-order logic and explicit rule systems, our framework enables LLMs to decompose complex queries into verifiable sub-tasks, orchestrate reliable solutions, and mitigate common failure modes such as hallucination and incorrect step decomposition. We demonstrate the practical benefits of this hybrid architecture through experiments on the DABStep benchmark, showing improved precision, coverage, and system documentation in multi-step reasoning tasks. Our results indicate that combining LLMs with modular logic reasoning restores engineering rigor, enhances system reliability, and offers a scalable path toward trustworthy, interpretable AI agents across complex domains.

1 Introduction

Large Language Models (LLMs) have rapidly advanced the field of artificial intelligence, enabling machines to generate fluent text, answer questions, and even write code at a level that often rivals human performance [28]. Yet, a fundamental question remains: do LLMs truly possess reasoning capabilities, or are they simply leveraging statistical correlations learned from massive datasets [22, 23, 38]?

This question is not just theoretical. As LLMs are starting to play an important role in decision making systems where wrong outcomes may involve potential financial or life loss, it is imperative to characterize their limitations, so that we can make these systems more robust. A particularly demanding challenge arises when LLMs are tasked with planning and orchestration—interpreting user intent, decomposing complex queries, and coordinating multiple software components to deliver robust solutions. These scenarios require a deeper capacity for understanding, abstraction, and structured problem-solving. While LLMs excel at tasks that rely on probabilistic inference—such as translation, summarization, or sentiment analysis—they often struggle with problems that require strict logical reasoning or adherence to formal rules [8, 36, 37, 42]. Tasks like mathematical proofs, logic puzzles, or precise legal interpretations expose the limitations of current models, which may generate plausible but ultimately incorrect or inconsistent answers.

Recent research has explored augmenting LLMs with external tools, code execution, or verification mechanisms to address these gaps [7, 12, 13, 27, 30]. However, unconstrained code generation introduces new risks—such as security vulnerabilities and unpredictable behaviors—and does not fundamentally solve the challenge of robust, interpretable reasoning [9, 24].

To bridge this gap, we propose a new approach: equipping LLMs with a broad set of specialized tools that can be composed to solve problems in smaller, manageable steps. Rather than aiming for unrestricted open-domain problem solving, our goal is to expand LLM capabilities across a wider range of topics by enabling them to orchestrate solutions through the intelligent composition of predefined components. This strategy, illustrated in Figure 1, represents a spectrum: on one end, systems are limited to fixed functionality, replicating conventional software engineering development processes; on the other, LLMs are used generate arbitrary code. Our approach seeks a balance, empowering LLMs to generate solutions by stitching together modular, well-understood components—thus increasing reliability, interpretability, and security.

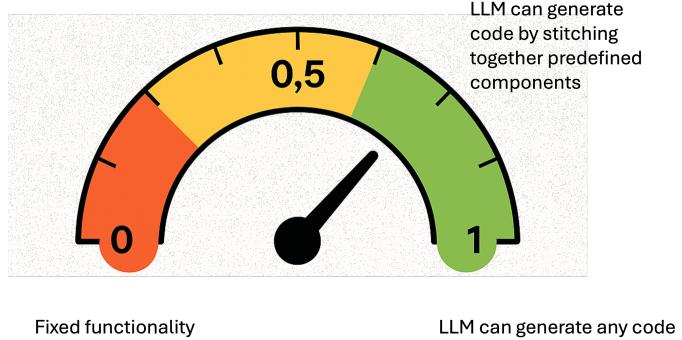


Figure 1: Trade-off between fixed functionality and complete freedom to generate any code.

This paper is organized as follows. Section 2 discusses the what we mean by logic, that provides a discrete way to reasoning about the world. Section 3 discusses the limitations of LLM reasoning, and what may happen if you mix discrete reasoning with probabilistic reasoning. Section 4 discusses how we can enhance reasoning skills by adding facts, rules and a first order proof system to assist LLMs. In Section 5, we show how AI Agents can benefit from the Neurosymbolic flow. In Section 6, we show that LLM-based systems with strong reasoning agents change the way we develop systems. In Section 8, we show some experiments, followed by conclusions.

2 L is for Logic

Logical reasoning [49] is foundational to structured problem-solving in mathematics, science, and computer science. At its core, logic provides a formal framework for analyzing situations, evaluating arguments, and drawing valid conclusions based on well-defined rules and relationships. This kind of reasoning is not only essential for human cognition but also underpins many aspects of artificial intelligence, including the design of systems that can interpret, manipulate, and reason about information.

There are two primary forms of logical reasoning: deductive and inductive. Deductive reasoning starts with general premises and derives specific, guaranteed conclusions—if the premises are true, the conclusion must be true. Inductive reasoning, in contrast, involves generalizing from specific observations, leading to conclusions that are probable but not certain. While both forms are valuable, deductive reasoning is particularly important for

tasks that require absolute correctness and reliability.

To formalize logical reasoning, two widely used systems are propositional logic and first-order logic (FOL). Propositional logic operates on statements that are either true or false, using logical connectives such as "and" (\wedge), "or" (\vee), and "not" (\neg). For example, consider the following scenario:

- $p \rightarrow r$ (if p then r)
- $r = \neg\text{day}$ (r is true if it is not day)
- day (it is day)

From these, we can deduce using traditional logic deduction systems, such as *modus tollens*, that $\neg p$ (not p) holds. While propositional logic is powerful for simple truth-functional relationships, it cannot express relationships between objects or quantify over domains. Moreover, determining the validity of propositional logic statements is computationally challenging (NP-complete), which can make large-scale reasoning intractable.

First-order logic extends propositional logic by introducing quantifiers like "forall" (\forall) and "exists" (\exists), enabling reasoning over individuals, their properties, and relationships. For example, a query such as "Give me the top-5 countries with anomalous purchases" can be represented as:

$$\forall x \in \text{Purchases} \mid \text{IsAnomalous}(x) \wedge \text{InState}(x, s) \wedge \text{Accumulate}(x, s, y) \wedge \text{Sort}(y, z, 5)$$

Here, x ranges over purchases, and predicates like IsAnomalous, InState, Accumulate and Sort capture domain-specific logic. FOL is much more expressive than propositional logic, making it suitable for modeling complex domains. However, this expressiveness comes at a cost: reasoning in FOL can be undecidable, meaning that no algorithm can guarantee a solution for every possible case.

Translating user queries into logical statements, and then decomposing them into actionable steps, presents several challenges for LLMs:

1. Breaking down complex queries into smaller, logically coherent steps;
2. Selecting or constructing functions to perform each step;
3. Composing these steps into an overall solution that is both correct and interpretable.

Current LLM architectures, while powerful in generating plausible text, often lack the explicit structure needed for rigorous logical reasoning and reliable problem decomposition. In the following sections, we examine why this is the case and how integrating formal logic systems with LLMs can help overcome these limitations. Our approach emphasizes the use of pre-defined, secure functions and limits arbitrary code generation to reduce security risks and improve reliability [2, 41].

3 Embeddings, Attention and Logic Reasoning

The development of embeddings has fundamentally changed how large language models (LLMs) represent and manipulate language [21]. Instead of treating words as isolated symbols, embeddings map words, phrases, and even complex concepts into high-dimensional vector spaces. This mapping allows LLMs to capture subtle semantic relationships—such as the association between "sunny" and "day"—by measuring distances or angles (for example, cosine similarity) between their respective vectors. As a result, LLMs can generate text that is contextually rich and coherent, reflecting the nuanced ways in which humans use language.

However, this statistical approach to language understanding comes with inherent limitations, especially when it comes to logical reasoning. Embeddings are exceptionally good at capturing patterns and correlations in data, but they do not encode the rigid, rule-based structures that underpin formal logic. For instance, while "sunny" and "day" may be closely related in many contexts, they are not logically equivalent. This distinction becomes critical when a task requires strict logical inference rather than probabilistic association.

The introduction of the Transformer architecture [51] further advanced LLM capabilities by enabling models to dynamically focus on different parts of an input sequence. Through self-attention mechanisms, models can weigh the relevance of each token in context, allowing for more sophisticated pattern recognition. In practice, this means that when an LLM encounters a paragraph that can be mapped into the following logic statements:

- $p \rightarrow r$ (if p then r)
- $r = \neg\text{day}$ (r is true if it is not day)
- sunny (it is sunny)

An LLM may use the similarity between "sunny" and "day" to infer $\neg p$, as shown in the previous deduction. Yet, this inference is based on statistical proximity in the embedding space, not on a formal logical equivalence. As illustrated in Figure 2, such assumptions can break down in real-world scenarios—like the phenomenon of the midnight sun in Sweden—where "sunny" and "day" do not always align. This brings a fundamental question when using LLMs for reasoning. We humans, understand clearly when to switch between probabilistic reasoning and logical reasoning. Can LLMs do the same?



Figure 2: Midnight sun in Sweden, during the summer, according to [31].

This example highlights a broader challenge: LLMs are fundamentally probabilistic systems. They excel at generalizing from data and generating plausible responses, but their reasoning remains approximate. When faced with problems that require discrete, binary decisions—such as those found in mathematics or formal logic—LLMs may produce answers that are statistically likely but logically unsound. This is particularly problematic for tasks where correctness is non-negotiable, such as legal reasoning, scientific proof, or safety-critical decision-making.

Recent research [1, 17, 34, 56] has sought to bridge this gap by combining the strengths of embeddings and neural architectures with the rigor of

symbolic logic. Hybrid approaches introduce mechanisms for translating natural language queries into formal logical representations, decomposing complex problems into smaller steps, and verifying solutions using symbolic reasoning tools like SAT solvers or rule-based engines. These neurosymbolic systems offer a path forward, enabling LLMs to harness both the flexibility of probabilistic reasoning and the precision of formal logic.

In summary, while embeddings and attention-based algorithms have unlocked remarkable advances in language understanding and generation, they are not a substitute for true logical reasoning. To build AI systems capable of robust, reliable problem-solving across a broad range of domains, it is essential to integrate these statistical methods with explicit logical frameworks. This integration not only enhances the reasoning capabilities of LLMs but also helps ensure that their solutions are both interpretable and trustworthy.

4 Facts, Rules and Full First Order Reasoning

If you ask a person from the Knowledge Graph [19] domain what you need to complement the LLMs’ reasoning ability, he/she will say graphs or meta-annotations, i.e. vertices specifying portions of the domain expertise or knowledge connecting to ideas and concepts, and edges connecting them. In the most basic form, you may want to represent ideas that may be difficult to the LLM, either because the LLM has not been trained on them, or because you need a crisp rather than a stochastic answer in the matter.

Dealing with discrete problems has been the center of Computer Science [32], which has existed for over 70 years, evolving from early theoretical foundations to a vast and dynamic field. Since its inception, much of its focus has been on discrete aspects of computation, dealing with structured, rule-based systems rather than continuous processes found in physical sciences. Topics such as algorithms, data structures, logic, artificial intelligence, and complexity theory form the backbone of computer science, emphasizing discrete mathematical principles. This discrete nature is evident in areas like programming languages, which operate on finite sets of instructions, and digital circuits, which rely on binary states to perform computations. While continuous paradigms like machine learning and signal processing have become increasingly influential, the fundamental role of discrete reasoning remains indispensable. Human cognition itself is often structured around discrete, categorical distinctions—logical inference, symbolic manipulation, and decision-making processes—all aligning closely with the discrete nature of computation.

We will be using throughout this article the programming language **Prolog** [6, 43], although there are several other programming language alternatives that deliver the similar capabilities, such as HOL [45], Lean [11] or Z3 [10], although with Z3, the representation would be at a much lower level.

The reasons why we chose Prolog as opposed to other languages, such as the ones mentioned above, are the following:

- We want to be able to represent programs as predicates without side effects, as side effects of programming language makes it harder to evaluate how information flows.
- We want an efficient way to represent not only the logical aspects modeling user's questions, but also a small to medium size database representing facts and rules.
- Even though LLMs have been trained extensively in any programming language, it has been our experience that they may struggle if the interfaces of the language or libraries change frequently (or example in different versions of C++). In this case, we may benefit from using a n more stable and older programming languages, such as Prolog and C.
- We wanted a language that is rich, but not too rich that will cause hallucination to become a problem, as we want to define programs as composable objects.

4.1 Facts

Facts express statements of a domain, such as:

```
acquirer_country(gringotts, gb).
```

This fact expresses the property that Gringotts is an acquirer located in GB. In Prolog, variables beginning with lowercase express atoms or predicates, i.e., `gringotts` and `gb` are atoms and `acquirer_country` is a predicate, expressing a relation between the atoms `gringotts` and `gb`.

In general, predicates express relations that are n-ary, such as:

```
merchat_data(yoga_masters, 2, gringotts, 7997, f)
```

You can now represent several a database containing several facts:

```

acquirer_country(gringotts, gb).
acquirer_country(the_savings_and_loan_bank, us).
acquirer_country(gringbank_of_springfieldottts, us).
acquirer_country(dagoberts_vault, 'nl').
acquirer_country(dagoberts_geldpakhuis, 'nl').
acquirer_country(lehman_brothers, us).
acquirer_country(medici, it).
acquirer_country(tellsons_bank, fr).

```

The reader may be asking the following questions.

- Why can't we represent these facts in a context to an LLM? In fact you can, but we are using a more compact representation. Furthermore, prolog enables us to establish relations that are not probabilistic, such as when you specify these facts in an LLM. When we say Gringrotts is an acquirer located in GB, we are not saying that Gringrotts is an acquirer located in GB with 90% chance, but always. Second, when we state that GB is located in Europe, but it is not part of the European Union, although GB was once part of the European Union, it may have this relation probabilistically, but we are stating a precise definition of what we are defining, and that does not include GB as part of the European Union.
- Why can't we just use facts from an SQL database? In fact, you can. Having large tables represented in database systems will be used later on, but being able to represent part of the data or even caching important data in memory as facts pose some advantages. In addition, as we will see it later, at a minimum, you must be able to represent domains of data in Prolog to enable negation.
- Facts allow us to represent knowledge graphs and meta information, such as $u \xrightarrow{\text{acquirer_country}} v$.

4.2 Rules

Once we define facts, we can define rules, which are compositional rules using facts. But to first define rules, we need to define the domain of the variables, as to perform the operation $\exists X$ or $\forall X$, we have first to define the domain of possible values of X .

```
acquirer(gringotts).
```

```

acquirer(the_savings_and_loan_bank).
acquirer(gringbank_of_springfieldott).
acquirer(dagoberts_vault).
acquirer(dagoberts_geldpakhuis).
acquirer(lehman_brothers).
acquirer(medici).
acquirer(tellsons_bank).
country(gb).
country('nl').
country(us).
country(it).
country(fr).

```

Once you define a domain, we can represent rules that applies to the domains.

```

acquirers_in_same_country(X,Y) :-
    acquirer(X), acquirer(Y), country(Z),
    X \= Y,
    acquirer_country(X,Z), acquirer_country(Y,Z).

```

This rule states that X and Y are acquirers in the same country if there is a country Z (in first order logic, $\exists Z \in \text{Domain}$) such that X is an acquirer Z and Y is an acquirer in Z . Using `acquirer(X)` and `acquirer(Y)` ensures that in the event the user did not specify X or Y , e.g., when they are used as output variables, we will be doing $\forall X \in \text{Acquirer_Domain}$ and $\forall Y \in \text{Acquirer_Domain}$, respectively. In the definition of `acquirers_in_same_country(X,Y)`, $\backslash=$ means not equal.

Prolog, having first order representation of symbolic programs, allows you to specify more complex relations, such as negations.

```

not_in_same_country(X, Y) :-
    acquirer(X), acquirer(Y), country(Z),
    X \= Y,
    acquirer_country(X, Z),
    \+ acquirer_country(Y, Z).

```

In this case, `\+ acquirer_country(X, Z)` means `acquirer_country(X, Z)` fails, which may be considered a form of negation. Being able to express negation is what differentiates Prolog from traversals in knowledge graphs, and negation makes the language a first-order logic proof system.

In this example, if you query the database for `not_in_same_country(lehman_brothers, Y)`, we get as answer the following, using SWIPL Prolog interpreter [46].

```
?- not_in_same_country(lehman_brothers, Y).
Y = gringotts ;
Y = dagoberts_vault ;
Y = dagoberts_geldpakhuis ;
Y = medici ;
Y = tellsons_bank.
```

You can see now the advantages of representing facts and rules as opposed to representing just knowledge graphs. Whereas facts are equivalent to knowledge graphs, rules represent paths, but negation adds another level of complexity that is very hard to represent in pure knowledge graphs.

One important fact here is that expressing domains for very large databases may be infeasible. We have to decide on the feasibility of expressing negation and quantifiers for variables to be considered inputs or outputs, or just only accepting instantiated variables as inputs (in case we cannot represent the full domain of variables).

The reader should note that we are not saying that probabilistic search is bad. If you have a graph, or a database table with transactions, and you want to know which pair of connected nodes or which translation may be faulty with high probability, you are referring to a probabilistic search. On the other hand, if you want to know if a node can reach a destination node of a graph, you do not care if the probability is 1 or 99%, you want to know if there is a path to that destination. That's a binary answer (yes/no), and that's where a language representing facts, rules and proof system such as Prolog sits in.

Prolog is a rich language, and we will mention a few operators that are worth mentioning (non-exhaustive):

- $P1(X), !, P2(X, Y)$: finds an assignment for variables in $P1$ (e.g. X in our example), then for all solutions $P2$ finds a solution (e.g. Y in our example), but does not search for an alternative X in $P1$. It can be thought of as converting a $\forall X | P1(X)$ into $\exists X | P1(X)$, i.e., it finds the first match to $P1(X)$ and it never re-evaluates it.

In our example, it becomes, $\exists X | P1(X) \wedge \forall Y | P2(X, Y)$ with the caveat that once a solution for $P1(X)$ is satisfied, we will only look for values in the domain of Y that satisfies $P2(X, Y)$.

For example, the operator `\+ P` can be defined as:

```
negate(P) :- call(P), !, fail.  
negate(_).
```

In other words, if calling `P` succeeds, `negate(P)` fails, otherwise it succeeds.

- `findall(Variable, Goal, Result)`: equivalent to `Result = [V1, V2, ..., Vn] :- Goal(V1), ..., Goal(Vn)`.
- `maplist(Function, List, Result)`: equivalent to `Result = [F(X1), F(X2), ..., F(Xn)] :- List = [X1, X2, ..., Xn], Function(X1, R1), ..., Function(Xn, Rn), Result = [R1, R2, ..., Rn]`.
- `Term =.. [Func|Args]`, `call(Term)`: allows you to do partial matching or dynamic matching of functions.
- `current_op(Precedence, Type, Name)`: given an operator, what's its precedence and type (the three variables can be inputs or outputs).
- `current_predicate(Name/Arity)`: this predicate searches for name of the function and its arity.

These last functions are useful to create dynamic infrastructure for facts, rules, and proofs, and searching for their information. The reader should know that Prolog is a much richer function with native support for lists, with `[1,2,3]` representing a list of three elements, and `[X|Xs]` representing a list where `X` is the first element of the list, and `Xs` representing the remaining elements of the list.

5 Intelligent Agents

We will be discussing in this section how we can transition from a use-case or MVC software development to intelligent agents using LLMs as interfaces. In the following example, we want to examine the complexity of adding the user story of ordering a Margherita gourmet pizza in 20 minutes to a food app, as an optimization to the the flow presented in Figure 3.

We have to assume that to implement this use case, we need access to the following data sources and algorithms:

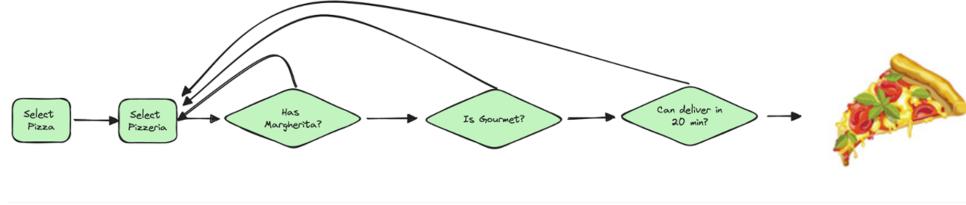


Figure 3: User Story to Order Pizza for a Food Delivery App

1. Restaurant database that can be searched by location and by type of food.
2. Menu database, where user can search for types of food served by the restaurant.
3. Algorithm that computes the delivery time from the restaurant to your location.

The reader should notice that this use case implements a **single** type of user interaction, and if we decide to modify the interaction, we will need to change the user story, or create another implementation that accommodates a different user story.

In general, given as sample of user’s queries about a system, we want to be able to generate complete systems that generalize and are able to answer a very large number of questions. In [20], we addressed the problem industry is facing right now that with LLMs, and more accurately, with AI Agents, specifications for software systems are appearing as a set of questions the system needs to answer. We addressed this problem by using LLMs to enlarge the few questions people provide to hundreds of questions of similar nature in an environment, and later generating maps to **algorithms, data sources** and **UI/UX interactions or visualizations**.

5.1 AI Agents

An AI Agent [3, 50, 54] encompasses a system that employs an LLM to process and reason about a specific domain. To generate specific answers (often related to the domain), the AI Agent leverages auxiliary systems in conjunction with the LLM. These auxiliary systems support the agent in comprehending the domain and facilitating the creation of accurate responses.

AI Agents consist of four major components. The *agent core* forms the central component and is responsible for orchestrating the agent's overall functionality. The *memory* module enables the agent to store and retrieve relevant information, enhancing its ability to retain context and make informed decisions. The *planner* component guides the agent's actions by formulating a strategic course of actions based on the given problem or task. Finally, the set of *tools* encompasses various external components and resources that assist the agent in performing specific tasks or functions within the defined domain, such as algorithms and data queries. These components collaboratively enable AI Agents to effectively process information, reason, adapt to query changes and generate responses in a manner aligned with their designated purpose.

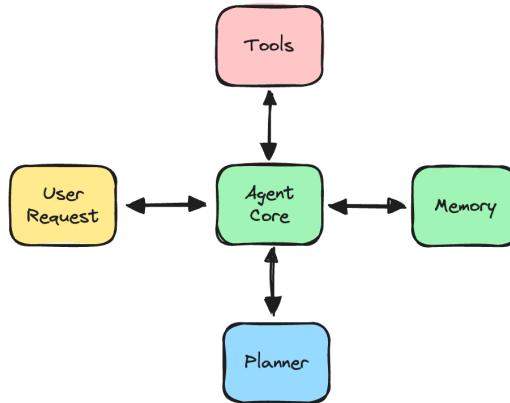


Figure 4: AI Agent from [50]

5.1.1 Agent Core

The agent core is a crucial component within an AI Agent that plays a central role in orchestrating the agent's overall functionality. It receives a query from the user. Consequently, it manages the decision-making processes, communication, and coordination of various modules and subsystems within the agent. Finally, it aggregates the information and generates a response.

The agent core is also responsible for managing the agent's internal state. It maintains a representation of the agent's assets and internal state, allowing it to reason, plan, and adapt its behavior accordingly. The core oversees the update and retrieval of information from the agent's memory, enabling it

to access relevant knowledge and contextual information during decision-making processes.

5.1.2 Memory

The memory module within an AI Agent encompasses two important aspects: historical memory [53] and contextual memory, usually represented as a knowledge-graph [58].

Historical memory serves as a repository for past interactions and experiences of the AI Agent. It stores a record of previous inputs, outputs, and the outcomes of actions taken by the agent. This historical data is valuable as it enables the agent to learn from past interactions and avoid repeating mistakes. Through the historical memory the agent gains insights about effective strategies, successful outcomes/patterns enabling an informed decision making process. Historical memory is usually split as short term memory and long term memory. Short term memory is usually retrieved from temporal most recent interactions, whereas long term memory is usually stored and retrieved through closeness to the topic being discussed.

Contextual memory maintains an understanding of the current and coherent context, and the world as we perceive, in the context of the objective on the AI Agent. It stores relevant context that provides the necessary background for the agent to interpret and respond appropriately to the present state. This can include information about the environment, the user's preferences or intentions, and any other contextual factors that influence the agent's behavior. Contextual memory allows the agent to adapt its action and responses to specific circumstances, thereby enhancing its ability to interact intelligently with changing environments.

Together, historical and contextual memories allow the AI agent to combine past experiences and current context for an efficient decision making process.

5.1.3 Planner

The planner component within an AI Agent plays a crucial role in guiding the agent's actions and formulating a strategic course of action based on the given problem or task. It is responsible for generating a sequence of steps or actions that lead the agent towards achieving its objectives. The planner analyzes the current state of the environment, along with any available information or constraints, to determine the most effective sequence of actions (by calling functions or making queries to databases) to achieve

the desired outcome. It also takes into account other factors such as goals, resources, rules, and dependencies to generate a plan that optimizes the agent’s decision-making process.

When the user asks the question *I want to order a Margherita pizza in 20 min*, the planner needs to consolidate information about user’s temporal requests, context (e.g. the location where the user is located), tools and data sources available, and finally develop a detailed plan of action. In later sections, this is where we will focus our attention, i.e. on how to assist AI Agents with logical reasoning in order to be able to answer complex questions, or deal with complex scenarios.

5.1.4 Tools

In an AI Agent, the set of tools encompasses various resources and functionalities that assist in performing specific tasks or functions within the defined domain. Here is a non-exhaustive list of possible tools that can be utilized in an AI Agent:

- *RAG (Retrieval-Augmented Generation [26])* — Combines retrieval-based methods with generative language models. It enables the agent to retrieve relevant information from a knowledge base and use it to generate coherent and contextually appropriate responses. Common data sources for RAG include Question-Answer databases, documentation and web pages.
- *Structured data access* — Connect to databases and allow the AI Agent to dynamically access and retrieve information from structured external data sources.
- *Unstructured data access* — Connect to a plethora of multi-modal sources of information (text, pictures, video streams, PDF, power-point presentations, spreadsheets, to name a few), some of which can be indexed by RAG systems, some of which may be accessed and retrieved by some other means, such as search.
- *Machine learning frameworks* — Provide tools and algorithms for training and deploying machine learning models. These frameworks enable the agent to use various machine learning techniques, including supervised learning, unsupervised learning, or reinforcement learning, to enhance its capabilities.

- *Visualization tools* — Help represent and interpret data or model outputs in a visual format. These tools can help the agent understand complex patterns, relationships, or trends in the data, helping in decision-making and analysis.
- *Simulation environments* — Provide a controlled virtual environment where the AI Agent can interact and learn without impacting the real world. These tools allow the agent to practice and refine its skills, test different strategies, and evaluate the potential outcomes of its actions.
- *Monitoring and logging frameworks* — Monitoring and logging frameworks facilitate the tracking and recording of agent activities, performance metrics, or system events. These tools assist in evaluating the agent’s behavior, identifying potential issues or anomalies, and supporting debugging and analysis.
- *Data preprocessing tools* — Help in cleaning, transforming, and preparing raw data before feeding it into the AI Agent. These tools may include techniques for data cleaning, normalization, feature selection, or dimensionality reduction, ensuring the quality and relevance of data used by the agent.

These tools enhance LLM capabilities by providing it with specialized functionalities for specific domains that would be outside the domain of LLMs. Any system using LLMs as interfaces to user interactions, such as the Copilot in [40] will incorporate these components into the implementation.

In this paper, we will use the classification of [20] to label tools as data sources, visualization artifacts and algorithms.

6 LLM is the New UI/UX

With advent of LLMs and intelligent agents, it has become natural for people to desire LLMs to being able to answer questions like the following question.

I want to order a gourmet Margherita pizza in 20 minutes.

When used as the front end to human interaction, LLM-based user interfaces enable systems to dynamically change queries and algorithm executions, depending on user’s intent, and to route the correct visualization to the user. As mentioned in Gartner’s keynote [39], *the most obvious value of LLMs is to put the human in the center of the interactions of systems*,

and this provides a radical simplification of the use of such systems. This simplification is achieved by being able analyze information, and the adaptation comes from being able to change the flow of control and access to data to accommodate what the user requests. In other words, LLMs can write their own code and access data sources in patterns we may not have thought before.

The problem with that this previous statement is that it carries a great deal of uncertainty that software engineers are required to answer before they can implement the system. In user story development, for example, as follow-up questions one would need to document in the development plan include:

- Which data sources should we connect to?
- Which algorithms do we need to invoke to solve this request?
- Which interfaces are required to implement this user story?
- How do we ensure levels of privacy and security, given the user's query may trigger a dynamic program and data access patterns ?
- Which other questions do we want to be able to solve?

Together with that simple question, and in order to make systems more generalizable [16], we need to understand the scope of such systems, and possibly being able to answer a broader set of questions, such as:

1. Can this restaurant deliver food in 20 min?
2. Give me the list of all restaurants that deliver gourmet pizza in 20 min.
3. Give me the 20 top evaluated restaurants that can deliver gourmet pizza in 20 minutes.

The reader can easily see that the first question requires just a simple yes/no answer. The second question requires a summarization or visualization agent to provide the answer. The third query will require getting data from possibly an additional table from the backend. Without fully understanding the scope of the system, and understanding what are the data sources, algorithms and visualizations that are needed, generalizing user's intent can easily become an almost impossible task.

Because of the complexity of such systems, a better representation of an AI Agent interaction with different types of tools and data sources is given in 5

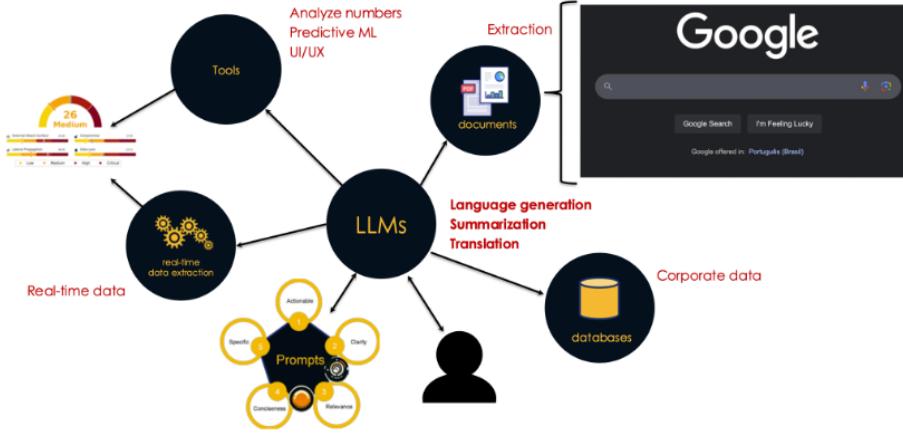


Figure 5: AI Agent Representation From Interaction View

7 The Role of Prolog to Assist Planner Activities

In order to discuss the role of logic reasoning in assisting planning activities, we have first to understand what may go wrong when a user asks a question.

- First, and foremost, the LLM may hallucinate in generating the answer, as such systems are based on probabilistic generation of text [18].
- Second, if the AI Agent is searching for information using RAG, the search mechanism may not be enough to retrieve the answer to the user [14]. In its simplest forms, RAG utilizes a proximity metric that may not align well with the user’s intent.
- Third, if the planner decided on the wrong set of steps to execute answer a user’s question, we may end up executing the wrong query from the user. This is the problem we want to attack by adding logic reasoning [22, 23].

It is important to note that considering symbolic and neural techniques has been proposed in the past, such as in [15, 33, 47], and the use of Prolog in conjunction with LLMs have been proposed in several recent works [4, 5, 29, 35, 44, 48, 55].

One of the major roles here is played by the planner stage, as it is responsible for outlining the steps to solve the query from the user. As we mentioned earlier, we need a solution that is able to:

- Define the steps as a sequence of positive and negative predicates (remember, Prolog has negation of predicates).
- Utilize a large number of predicates that are able to cover technical questions from the users.
- Create compositions of steps, such as and-ing or or-ing predicates.
- Either train an LLM to understand the underlying language to specify these steps, or use a language that has already been trained in the LLM.
- Have a very low cost of establishing ground truth, and even fix small mistakes from the LLM (if we use a programming language such as Prolog, ground truth can be established by compiling the code, and a parse tree can understand how good or bad the solution is).

Our approach here is to exercise what LLMs know best, i.e. how to process grammars, as LLMs have excelled in being able to generate text following grammatical rules.

We will follow the steps defined in [20] to create the tools we need to solve a given set of problems, i.e., starting from a set of seed questions from the user, we enlarge these questions and attempt to generate a list representative questions that correspond to the problem to be solved, generate a set of MVC components (model-view-controller) [52], as the MVC framework nicely represents the different types of components that exists between interactions of users to systems, as shown in Figure 6. Finally, from the MVC components, we identify the gaps of our implementation, and outline which functions need to be implemented.

As we want to cast this problem as a grammar problem with strong typing system, we need first to define the basic types we will use in this system, as seen below.

- **StrOrList:** Union[str,List[StrOrList]]: meta type to be used by **Filter**.
- **TableName:** str: table according to the list of tables the system can access.
- **Filter:** StrOrList: filter expressions, e.g. [and, [issuing_country, 'GB'], [ip_country, 'FR']].

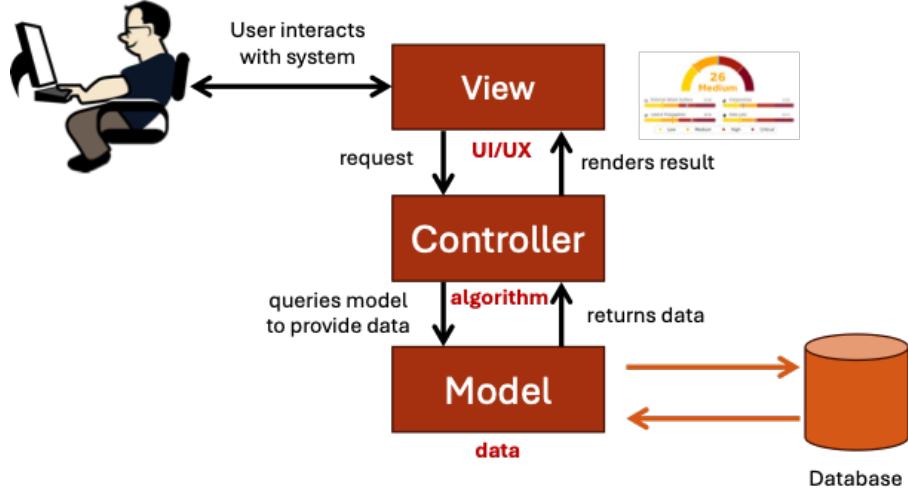


Figure 6: Model-View-Controller system

- **Projection:** `List[str]`: list of names from the table `TableName` that we will return to the user, according to the field names in the database.
- **Header:** `List[str]`: list of field names from `TableName` or `Projection` if the projected list is not empty.
- **Data:** `List[List[str]]`: list of data fields, where the size of the inner list is the same as the `Header`, and each element of the inner list of `Data` corresponds a value of each field of `Header`.

Once we define data types used by the system, our objective is to define predicates that manipulate data. We provide at minimum functions that can be called to retrieve data and filter data in the system. Some examples of functions that we can use are presented in Table 1. This first set of functions correspond to functions available in SQL language (for now, we will keep these functions separate, although the user may see similarity between them and SQL query artifacts), as bundling them together depends on how we pay for the access to the data - e.g. do we pay per query or per byte transmitted? In addition, we may have certain domains that can be easily expressed as facts in Prolog, if the data sizes are not too large. The second set of functions correspond to algorithms, such as traditional predictive ML algorithms, complex arithmetic functions, etc. The third set of functions correspond to UI/UX artifacts, such as visualization of maps.

Table 1: List of functions available for the logic reasoner.

Function
query_data(-TableName, -Filter, -Projection, +[Header Data])
filter(-[Header Data], -Filter, +[Header1 Data])
project(-[Header Data], -Projection, +[Projection Data])
count(-[Header Data], +[Header Data])
anomaly(-[Header Data], +[[is_anomaly Header] Data])
visualize_map(-[Header Data])

It is clear that we may still need to comment the functions defined above, explaining which fields are expected or are generated. For example, an anomaly detection algorithm may require additional information, such that notifying the Planner that an amount paid by the customer and the state where the purchase were made are required input fields for the anomaly detection algorithm to work. Sometimes, even small examples of usage of these algorithms may be required.

The user may be asking the question, 1) how many functions is enough, 2) whether it is possible to dynamically instantiate functions in this infrastructure.

To answer the first question, we will refer the user to [20]. In this paper, we have shown that software engineering has changed from a structured design process to giving a sample of questions by product managers when using LLMs as interfaces. We argued that we need in this case to extract specification by following the steps below.

1. Expand the list of questions from the users using the original sample questions as seed.
2. Extract the model-view-components from the questions, always reusing the previous components whenever possible

The full algorithm proposed in [20] is described in Algorithm 1. It is worth noting that as LLMs can hallucinate, and since extraction of the specification is executed in the beginning of the project, we recommended manual review of all enlarged questions and MVC components.

7.1 Planning by Stitching Together Logical Components

Now we can list the general structure of the prompt for the planner to generate a planner strategy in Prolog.

Algorithm 1 Algorithm to extract MVC components from enlarged number of questions

Require: List of questions Q

```
AllTasks ← ∅  
for  $q \in Q$  do  
    Generate  $N$  related questions  $Q_q$  from  $q$ .  
     $T_q \leftarrow Planner($   
         $\{q\} \cup Q_q, current\_tools = AllTasks,$   
         $minimize = True)$   
    for  $t \in T_q$  do  
         $t['task'] \leftarrow model|view|controller$   
    end for  
    Manually validate the set  $T_q$ .  
     $AllTasks \leftarrow AllTasks \cup T_q$   
end for  
Manually validate final set of  $AllTasks$ 
```

You will interpret the user's query and generate a prolog code that best matches the user's intent to the prolog code. In doing so, you will evaluate the prolog code according to the 'evaluation' metric, and you will try to maximize the 'evaluation' metric.

```
<query>  
user query  
</query>
```

Previous interactions with the user are as follows:

```
<history>  
short term, long term interactions.  
</history>
```

You base use all functions available in the standard prolog language, with the addition of the following list of prolog assert database and foreign functions.

```
<assert>  
list of prolog database
```

```
</assert>
```

```
<foreign-functions>
list of foreign-functions
</foreign-functions>
```

Data accesses must follow the following schema.

```
<database-schema-description>
Definition of database schema in English
</database-schema-description>
```

You must use the following examples to guide your reasoning:

```
<examples>
list Chain-of-Thought examples
</examples>
```

Your response must be in the following JSON format:

```
'''json
{
    "explanation": str,
    "gaps": List[str],
    "findings": List[str],
    "plan": List[str],
    "action": List[str],
    "result": str,
    "evaluation": float,
}
'''
```

Explanation of each field, asking the LLM to align the action list to the CoT examples, and on how to increase and decrease the evaluation by specific quantifiable amounts.

In the evaluation, use the following guidelines:

- start with an evaluation score of 1.0 after generating the code.
- if you utilize an input to a foreign function before

```
instantiating its value, set evaluation score to 0.0.  
- if you use an LLM written assert predicate, reduce the  
evaluation score by 0.2.  
- if you call native Prolog functions such as sort, findall  
or aggregate with objects containing both the header and the  
data fields, reduce the evaluation score by 0.4.  
...
```

Remember, ...

Now generate the JSON file, maximizing the evaluation score
as instructed.

In this prompt, the list of `assert` functions corresponds to predicates written in Prolog, and `foreign-functions` corresponds to predicates written in other programming languages.

The idea of the evaluation metric is that we ask the LLM to self reflect on its reasoning steps by giving rewards and penalties. It is worth noting that in [57], a similar approach has been proposed to self-generate evaluation metrics that are fed into GRPO fine-tuning loop.

To answer the second question, we have to consider the three scenarios:

1. In the past [52], people have resorted to the specification of the user stories determine how the whole interaction with the user is driven, such as in the case of the Pizza ordering App, as presented in Figure 3.
2. In [13], we are able to dynamically generate, and keep trying until code passes. As advantage, it connects directly to Python infrastructure. However, because the Python language is extremely complex (e.g. users have to specify which libraries they intend to use).
3. We provide a large number of small components, and expect the LLM to "compose" a solution by stitching together small portions to comprise a full solution. As advantage, this allows one to have better control over hallucination and security, and even use smaller or simpler models to generate the code. On the negative side, our ability to answer very complex questions reside on the initial steps to explore the solution space.

We consider our approach to be a meet in the middle from a completely fixed functionality specification, such as the ones presented in traditional

software engineering strategies, and the completely loose application to let the LLM freely try different implementations until one passes, which may lead to answering open-domain questions.

8 Dataset Overview

8.1 DABStep Benchmark Overview

To develop our approach, we leveraged the **DABStep Benchmark**¹, a dataset specifically designed to assess the multi-step reasoning capabilities of AI agents. DABStep (Data Agent Benchmark for Stepwise Evaluation of Planning) comprises over 450 real-world, multi-hop reasoning tasks focused on domains such as financial transactions, risk analytics, merchant behavior, and operational metrics.

Each task includes structured data (in CSV or JSON format), optional unstructured context (such as documentation or manuals), and a natural language question that must be answered with high factual precision. The benchmark is engineered to challenge agents on planning, data selection, transformation, and reasoning, and it uses an exact-match evaluation metric that does not tolerate partial correctness or hallucination.

Table 2 shows a snapshot of the released data files that comprise the DABStep benchmark.

Table 2: Core data files included in the DABStep benchmark

File Name	Description
payments.csv	138k anonymized transactions with fraud/risk signals
payments-readme.md	Documentation for interpreting the payments dataset
acquirer.countries.csv	List of acquiring banks and their corresponding countries
fees.json	Dataset of 1000 structured fee schemes from financial networks
merchant.category.codes.csv	Lookup table of Merchant Category Codes (MCCs)
merchant.data.json	JSON file describing merchant identities and properties
manual.md	Condensed business handbook used for domain rules and constraints

The dataset tables are summarized as follows.

- acquirer_countries.csv
 - **acquirer**: string. acquirer name
 - **country_code**: The location (country) of the acquiring bank (Categorical) - SE, NL, LU, IT, BE, FR, GR, ES.
- fees.json

¹<https://huggingface.co/blog/dabstep>

- **ID**: identifier of the fee rule within the rule fee dataset
- **card_scheme**: string type. name of the card scheme or network that the fee applies to
- **account_type**: list type. list of account types according to the categorization ‘Account Type’ in this manual
- **capture_delay**: string type. rule that specifies the number of days in which the capture from authorization to settlement needs to happen. Possible values are ’3-5’ (between 3 and 5 days), ’>5’ (more than 5 days is possible), ’<3’ (before 3 days), ’immediate’, or ’manual’. The faster the capture to settlement happens, the more expensive it is.
- **monthly_fraud_level**: string type. rule that specifies the fraud levels measured as ratio between monthly total volume and monthly volume notified as fraud. For example ’7.7
- **monthly_volume**: string type. rule that specifies the monthly total volume of the merchant. ’100k-1m’ is between 100.000 (100k) and 1.000.000 (1m). All volumes are specified in euros. Normally merchants with higher volume are able to get cheaper fees from payments processors.
- **merchant_category_code**: list type. integer that specifies the possible merchant category codes, according to the categorization found in this manual in the section ‘Merchant Category Code’. eg, ’062, 8011, 8021’.
- **is_credit**: bool. True if the rule applies for credit transactions. Typically credit transactions are more expensive (higher fee).
- **aci**: list type. string that specifies an array of possible Authorization Characteristics Indicator (ACI) according to the categorization specified in this manual in the section ‘Authorization Characteristics Indicator’.
- **fixed_amount**: float. Fixed amount of the fee in euros per transaction, for the given rule.
- **rate**: integer. Variable rate to be specified to be multiplied by the transaction value and divided by 10000.
- **intracountry**: bool. True if the transaction is domestic, defined by the fact that the issuer country and the acquiring country are the same. False are for international transactions where the issuer country and acquirer country are different and typically are more expensive.

- merchant_category_codes.csv
 - **mcc**: integer that specifies the possible merchant category code.
 - **description**: textual description of merchant category code.
- merchant_data.json
 - **merchant**: Merchant name (Categorical), eg Starbucks or Netflix.
 - **capture_delay**: string type. rule that specifies the number of days in which the capture from authorization to settlement needs to happen. Possible values are '3-5' (between 3 and 5 days), '≥5' (more than 5 days is possible), ' <3 ' (before 3 days), 'immediate', or 'manual'. The faster the capture to settlement happens, the more expensive it is.
 - **acquirer**: list of string. acquirer names allowed for merchant. eg, "dagoberts_geldpakhuis", "bank_of_springfield".
 - **merchant_category_code**: list type. integer that specifies the possible merchant category codes, according to the categorization found in this manual in the section 'Merchant Category Code'. eg, '8062, 8011, 8021'.
 - **account_type**: string, one of R,D,H,F,S,O, where R = Enterprise - Retail, D = Enterprise - Digital, H = Enterprise - Hospitality, F = Platform - Franchise, S = Platform - SaaS, and O = Other.
- payments.csv
 - **psp_reference**: Unique payment identifier (ID).
 - **merchant**: Merchant name (Categorical), eg Starbucks or Netflix.
 - **card_scheme**: Card Scheme used (Categorical) - MasterCard, Visa, Amex, Other.
 - **year**: Payment initiation year (Numeric).
 - **hour_of_day**: Payment initiation hour (Numeric).
 - **minute_of_hour**: Payment initiation minute (Numeric).
 - **day_of_year**: Day of the year of payment initiation (Numeric).
 - **is_credit**: Credit or Debit card indicator (Categorical).
 - **eur_amount**: Payment amount in euro (Numeric).

- **ip_country**: The country the shopper was in at time of transaction (determined by IP address) (Categorical) - SE, NL, LU, IT, BE, FR, GR, ES.
- **issuing_country**: Card-issuing country (Categorical) - SE, NL, LU, IT, BE, FR, GR, ES.
- **device_type**: Device type used (Categorical) - Windows, Linux, MacOS, iOS, Android, Other.
- **ip_address**: Hashed shopper’s IP (ID).
- **email_address**: Hashed shopper’s email (ID).
- **card_number**: Hashed card number (ID).
- **shopper_interaction**: Payment method (Categorical) - E-commerce, POS. POS means an in-person or in-store transaction.
- **card_bin**: Bank Identification Number (ID).
- **has_fraudulent_dispute**: Indicator of fraudulent dispute from issuing bank (Boolean).
- **is_refused_by_adyen**: Adyen refusal indicator (Boolean).
- **aci**: Authorization Characteristics Indicator (Categorical).
- **acquirer_country**: The location (country) of the acquiring bank (Categorical) - SE, NL, LU, IT, BE, FR, GR, ES.

In addition to its dataset structure, DABStep categorizes tasks by difficulty and content complexity. Table 3 summarizes key properties of the benchmark.

Table 3: Summary statistics of the DABStep benchmark dataset

Property	Value
Number of tasks	456
Average number of tables per task	3.2
Median rows per table	5,000
Percent of tasks with joins	64%
Percent requiring external docs	28%
Easy/Hard task split	55% / 45%
Eval metric	Exact match
Domain types	Payments, Logistics, Inventory

8.2 Sample Tasks

To illustrate the benchmark's complexity and structured reasoning requirements, we present two representative examples:

- **Question:** Which card scheme had the highest average fraud rate in 2023?

Guidance: Answer must be the name of the scheme.

Answer: SwiftCharge

- **Question:** For the year 2023, focusing on the merchant *Crossfit Hanna*, if we aimed to reduce fraudulent transactions by encouraging users to switch to a different Authorization Characteristics Indicator through incentives, which option would be the most cost-effective based on the lowest possible fees?

Guidance: Answer must be the selected ACI to incentivize and the associated cost rounded to 2 decimals in this format: {card_scheme:fee}

Answer: E:346,49

These questions require the agent to reason over multiple dimensions: temporal filtering, column selection, aggregation, and domain knowledge interpretation. Although this dataset is targeted at multi-step reasoning tasks, we find it useful as it provides good documentation for the tables and the tasks, and that provides good insights.

Given the schema of the database, the instructions presented in dataset DABstep, samples of the data (extracted from the data itself), and the dev-set of the dataset, we can obtain an enlarged set of questions that result in the MVC diagram presented in Figure 7.

In this paper, we extended the work of [20] by instructing the generation of the MVC also to generate the interfaces to controllers and to the views, as shown in Figure 7.

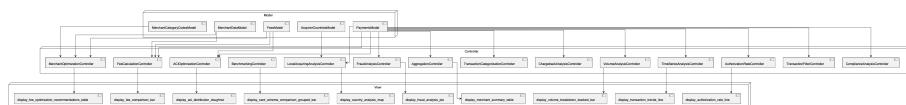


Figure 7: Model-View-Controller system for DABset expanded questions

Table 4 presents the list of tools that were inferred from the system. Please note that we have not listed here the full interface of the functions due to lack of space. We have also instructed the LLM to utilize the information from [25] when recommending UI artifacts in the views section.

That can be seen as the suffix for the `display` functions, such as `doughnut` in `display_aci_distribution doughnut`.

9 Conclusions

This paper examined the evolving role of Large Language Models (LLMs) in software systems, focusing on their integration as reasoning engines and user interfaces, and the challenges this presents for software engineering, system reliability, and interpretability. We identified that while LLMs excel at probabilistic language tasks, they exhibit critical limitations when confronted with problems requiring strict logical reasoning, discrete decision-making, or adherence to formal rules. These shortcomings are particularly evident in domains that include safety-critical applications, where correctness and transparency are paramount.

To address these gaps, we proposed a neurosymbolic approach that augments LLMs with modular, composable tools—especially logic-based components such as Prolog predicates. This strategy allows LLMs to decompose complex queries into smaller, verifiable steps, orchestrating solutions by stitching together specialized, well-understood functions. By leveraging first-order logic and explicit rule systems, our methodology enables precise representation of facts, rules, and negation—capabilities that are difficult to achieve with purely probabilistic models or knowledge graphs alone.

We further demonstrated how this approach supports the development of intelligent agents, where LLMs serve as the core orchestrators, supported by structured memory, planning modules, and domain-specific tools. This architecture facilitates dynamic, context-aware reasoning and allows agents to adapt to a wide range of user queries, while maintaining reliability and interpretability. The integration of logic reasoning at the planning stage helps mitigate common failure modes of LLMs, such as hallucination and incorrect step decomposition, by providing a framework for self-reflection and evaluation of reasoning chains.

Our experiments, conducted using the DABStep benchmark, validated the effectiveness of this approach in real-world, multi-step reasoning tasks. By expanding initial user questions into comprehensive sets of sub-tasks and mapping them to Model-View-Controller (MVC) components, we were able to regain the precision and coverage previously afforded by traditional user stories and use cases. This process not only aids in system documentation and effort estimation but also supports the systematic identification of functional gaps and the design of robust, generalizable AI-driven systems.

Table 4: List of functions by Algorithm 1, using 40 questions/query and extraction of MVC functions from MVC diagram

Function
get_payments_data()
get_fee_rules()
get_merchant_config()
get_acquirer_countries()
get_mcc_descriptions()
calculate_transaction_fee()
calculate Merchant_fraud_rate()
calculate_authorization_rate()
calculate_chargeback_rate()
categorize_transactions()
generate_optimization_recommendations()
analyze_volume_trends()
identify_local_acquiring_opportunities()
analyze_aci_usage()
check_compliance_metrics()
filter_transactions_complex()
aggregate_transaction_data()
analyze_time_series()
benchmark_merchant_performance()
display_fee_comparison_bar()
display_transaction_trends_line()
display_fraud_analysis_pie()
display_merchant_summary_table()
display_aci_distribution_doughnut()
display_volume_breakdown_stacked_bar()
display_fee_optimization_recommendations_table()
display_authorization_rate_line()
display_country_analysis_map()
display_card_scheme_comparison_grouped_bar()

Our work showed that combining LLMs with logic-based reasoning modules and modular tool orchestration provided a scalable path toward building AI agents capable of reliable, interpretable, and secure problem-solving across complex domains. This hybrid approach restored engineering rigor to LLM-driven systems and offers a practical methodology for both system designers and end-users seeking trustworthy AI solutions.

10 Acknowledgments

Portions of this document used Perplexity-AI to improve readability. We acknowledge the valuable feedback from Luis Lamb and Ajay Kumar in reviewing earlier versions of this document.

Reference

- [1] Nikhil Abhyankar, Vivek Gupta, Dan Roth, and Chandan K. Reddy. H-star: Llm-driven hybrid sql-text adaptive reasoning on tables, 2025.
- [2] Sakana AI. The ai scientist: Towards fully automated open-ended scientific discovery. <https://github.com/SakanaAI/AI-Scientist>, 2024. Accessed: 2025-06-20.
- [3] Praveen Akkiraju, Sophie Beshar, and Hunter Korn. Ai agents are disrupting automation: Current approaches, market solutions and recommendations, 2024.
- [4] Nasim Borazjanizadeh and Steven T. Piantadosi. Reliable reasoning beyond natural language, 2024.
- [5] Chris Clark. Enhancing llm reasoning through prolog: A breakthrough in symbolic logic processing, March 2025.
- [6] William F. Clocksin. *Programming in Prolog: Using The Iso Standard*. Springer, 5th edition, 2003.
- [7] CrewAI Contributors. Crewai: Orchestrate autonomous ai agents as a crew. <https://github.com/joaomdmoura/crewai>, 2024. Accessed: 2025-06-20.
- [8] Polyaná Costa, Pedro Santos, José Boaro, Daniel Moraes, Júlio Duarte, and Sergio Colcher. Experimenting with planning and reasoning in ad

hoc teamwork environments with large language models. In *Proceedings of the 16th International Conference on Agents and Artificial Intelligence - Volume 1: ICAART*, pages 438–445. INSTICC, SciTePress, 2024.

- [9] Cybersecurity and Infrastructure Security Agency (CISA). Secure by design alert: Eliminating os command injection vulnerabilities. <https://www.cisa.gov/resources-tools/resources/secure-design-alert-eliminating-os-command-injection-vulnerabilities>, 2024. Accessed: 2025-06-20.
- [10] Leonardo de Moura and Nikolaj Bjørner. Z3: an efficient smt solver, 2008. Accessed: June 5, 2025.
- [11] Leonardo de Moura, Sebastian Ullrich, et al. *The Lean Language Reference Manual*. Lean Community, 2024. Accessed: 2025-06-20.
- [12] Hao Ding, Ziwei Fan, Ingo Guehring, Gaurav Gupta, Wooseok Ha, Jun Huan, Linbo Liu, Behrooz Omidvar-Tehrani, Shiqi Wang, and Hao Zhou. Reasoning and planning with large language models in code development. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*, pages 6480–6490. ACM, 2024.
- [13] Hugging Face. Smolagents: Agents that think in code. <https://github.com/huggingface/smolagents>, 2024. Accessed: 2025-06-20.
- [14] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models, 2024.
- [15] Artur S. D’Avila Garcez, Luis C. Lamb, and Dov M. Gabbay. *Neural-Symbolic Cognitive Reasoning (Cognitive Technologies)*. Springer, 2009.
- [16] Sarah Gibbons, Tarun Mugunthan, and Jakob Nielsen. The 4 degrees of anthropomorphism of generative ai. *Nielsen Norman Group*, October 2023.
- [17] David Herron, Ernesto Jiménez-Ruiz, and Tillman Weyde. On the potential of logic and reasoning in neurosymbolic systems using owl-based knowledge graphs. *Neurosymbolic Artificial Intelligence*, 1:29498732251320043, 2025.

- [18] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, January 2025.
- [19] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and Philip S. Yu. A survey on knowledge graphs: Representation, acquisition and applications. arXiv:2002.00388 [cs.CL], 2020.
- [20] Claudionor N. Coelho Jr, Hanchen Xiong, Tushar Karayil, Sree Koratala, Rex Shang, Jacog Bollinger, Mohamed Shabar, and Syam Nair. Effort and size estimation in software projects with large language model-based intelligent interfaces. arXiv:2402.07158 [cs.SE], 2024.
- [21] Dan Jurafsky and James H. Martin. *Speech and Language Processing*. <https://web.stanford.edu/~jurafsky/slp3>, 3rd draft edition, January 2025.
- [22] Subbarao Kambhampati, Kaya Stechly, and Karthik Valmeekam. (how) do reasoning models reason? Annals of the New York Academy of Sciences, April 2025.
- [23] Subbarao Kambhampati, Kaya Stechly, Karthik Valmeekam, Lucas Saldyt, Siddhant Bhambri, Vardhan Palod, Atharva Gundawar, Soumya Rani Samineni, Durgesh Kalwar, and Upasana Biswas. Stop anthropomorphizing intermediate tookens as reasoning/thinking traces arXiv:2504.09762 [cs.AI], 2025.
- [24] Andrej Karpathy. Software is changing (again), 2025. Keynote at AI Startup School, published by Y Combinator. Accessed: 2025-06-20.
- [25] Randy Krum. Data visualization reference guides, 2025. Accessed: 2025-06-20.
- [26] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021.
- [27] Fangru Lin, Emanuele La Malfa, Valentin Hofmann, Elle Michelle Yang, Anthony Cohn, and Janet B. Pierrehumbert. Graph-enhanced large language models in asynchronous plan reasoning, 2024.

- [28] Kyle Lo, Akshita Bhagia, and Nathan Lambert. Large language modeling. Tutorial at Neural Information Processing Systems (NeurIPS), December 2024.
- [29] Lu Mao. Llm and prolog: the logical alternative to chain-of-thought reasoning, April 2025.
- [30] Dang Nguyen, Viet Dac Lai, Seunghyun Yoon, Ryan A. Rossi, Handong Zhao, Ruiyi Zhang, Puneet Mathur, Nedim Lipka, Yu Wang, Trung Bui, Franck Dernoncourt, and Tianyi Zhou. Dynasaur: Large language agents beyond predefined actions, 2024.
- [31] Sweden’s official website for tourism and travel information. Visit sweden, 2025.
- [32] Gerard O’Regan. *A Brief History of Computing*. Springer, 3rd edition, 2021.
- [33] Liangming Pan, Alon Albalak, Xinyi Wang, and William Yang Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning, 2023.
- [34] Hyun Ryu, Gyeongman Kim, Hyemin S. Lee, and Eunho Yang. Divide and translate: Compositional first-order logic translation and verification for complex logical reasoning, 2025.
- [35] Eugene Sahara. Prolog’s role in the llm era – part 1, August 2024.
- [36] Jintian Shao and Yiming Cheng. Cot is not true reasoning, it is just a tight constraint to imitate: A theory perspective, 2025.
- [37] Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce C. Ho, Carl Yang, and May Dongmei Wang. EHRA-agent: Code empowers large language models for few-shot complex tabular reasoning on electronic health records. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 22315–22339, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [38] Parshin Shojaee*†, Iman Mirzadeh*, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity, 2025.

- [39] Jen Singleton and Jim Hare. The ripple effect: How innovation disrupts, 2024.
- [40] Vikas Srivastava. Step into the future of zdx with 3 exciting new features: Zdx copilot, data explorer, and hosted monitoring. *Zscaler Blog*, April 2024.
- [41] BleepingComputer News Staff. Asana warns mcp ai feature exposed customer data to other orgs, 2025. Accessed: 2025-06-20.
- [42] Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. On the self-verification limitations of large language models on reasoning and planning tasks, 2024.
- [43] Leon Sterling and Ehud Shapiro. *The Art of Prolog: Advanced Programming Techniques (Mit Press Series in Logic Programming)*. MIT Press, 2nd edition, 1986.
- [44] Xiaoyu Tan, Yongxin Deng, Xihe Qiu, Weidi Xu, Chao Qu, Wei Chu, Yinghui Xu, and Yuan Qi. Thought-like-pro: Enhancing reasoning of large language models through self-driven prolog-based chain-of-thought, 2024.
- [45] HOL team. Hol - interactive theorem prover, 2025.
- [46] SWI-Prolog team. Robust, mature, free. prolog for the real world.
- [47] Charanraj Thimmisetty, Praveen Tiwari, Didac Gil de la Iglesia, Nandini Ramanan, Marjorie Sayer, Viswesh Ananthakrishnan, and Claudionor Nunes Coelho, Jr. Log2NS: Enhancing deep learning based analysis of logs with formal to prevent survivorship bias. <https://arxiv.org/abs/2105.14149>, 2021.
- [48] Priyesh Vakharia, Abigail Kufeldt, Max Meyers, Ian Lane, and Leilani H. Gilpin. *ProSLM: A Prolog Synergized Language Model for explainable Domain Specific Knowledge Based Question Answering*, page 291–304. Springer Nature Switzerland, 2024.
- [49] Dirk van Dalen. *Logic and Structure*. Springer-Verlag, 5th edition, 2013.
- [50] Tanay Varshney. Introduction to llm agents. <https://developer.nvidia.com/blog/introduction-to-l1m-agents/>, November 2023. Accessed on December 19, 2023.

- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <https://arxiv.org/abs/1706.03762>, 2017.
- [52] D.P. Voorhees. *Guide to Efficient Software Design: An MVC Approach to Concepts, Structures, and Models*. Texts in Computer Science. Springer International Publishing, 2021.
- [53] Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang, Yongyue Zhang, Huifeng Guo, Ruiming Tang, and Yong Liu. From human memory to ai memory: A survey on memory mechanisms in the era of llms, 2025.
- [54] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey, 2023.
- [55] Xiaocheng Yang, Bingsen Chen, and Yik-Cheung Tam. Arithmetic reasoning with llm: Prolog generation & permutation, 2024.
- [56] Jinghan Zhang, Xiting Wang, Weijieying Ren, Lu Jiang, Dongjie Wang, and Kunpeng Liu. Ratt: A thought structure for coherent and correct llm reasoning, 2024.
- [57] Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Song. Learning to reason without external rewards, 2025.
- [58] Kaiwen Zuo, Yirui Jiang, Fan Mo, and Pietro Lio. Kg4diagnosis: A hierarchical multi-agent llm framework with knowledge graph enhancement for medical diagnosis, 2025.