# GateLens: A Reasoning-Enhanced LLM Agent for Automotive Software Release Analytics

Arsham Gholamzadeh Khoee, Shuai Wang, Yinan Yu, Robert Feldt
Chalmers University of Technology
Gothenburg, Sweden
{arsham.khoee, shuaiwa, yinan, robert.feldt}@chalmers.se

Dhasarathy Parthasarathy
Volvo Group
Gothenburg, Sweden
dhasarathy.parthasarathy@volvo.com

*Abstract*—Ensuring reliable software release decisions is critical in safety-critical domains such as automotive manufacturing. Release validation relies on large tabular datasets, yet manual analysis is slow, costly, and error-prone. While Large Language Models (LLMs) offer promising automation potential, they face challenges in analytical reasoning, structured data handling, and ambiguity resolution. This paper introduces GateLens, an LLM-based system for analyzing tabular data in the automotive domain. GateLens translates natural language queries into Relational Algebra (RA) expressions and generates optimized Python code. Unlike traditional multi-agent or planning-based systems that can be slow, opaque, and costly to maintain, GateLens emphasizes speed, transparency, and reliability. Experimental results show that GateLens outperforms the existing Chain-of-Thought (CoT) + Self-Consistency (SC) based system on real-world datasets, particularly in handling complex and ambiguous queries. Ablation studies confirm the essential role of the RA layer. Industrial deployment shows over 80% reduction in analysis time while maintaining high accuracy across test result interpretation, impact assessment, and release candidate evaluation. GateLens operates effectively in zero-shot settings without requiring few-shot examples or agent orchestration. This work advances deployable LLM system design by identifying key architectural features—intermediate formal representations, execution efficiency, and low configuration overhead—crucial for safety-critical industrial applications.

*Index Terms*—Large Language Models, Tabular Question Answering, Software Release Analytics, Automotive Software Testing, Test Result Interpretation, Interpretable Reasoning

## I. INTRODUCTION

Validating software releases in the automotive industry is a multifaceted challenge, particularly for embedded software in safety-critical systems. Modern vehicles integrate numerous subsystems, making this process complex and resource-intensive. Each integration phase involves *gating* steps—critical checkpoints where tests verify compliance with predefined quality standards. Failures at these gates can ripple through the system, delaying dependent subsystems regardless of their individual quality. Release managers tasked with safeguarding quality must analyze vast quantities of test results and validation data. While essential for ensuring safety and reliability, this process is time-consuming and prone to human error in data interpretation.

The software industry's transition from manual to automated processes has entered a new era with the emergence of Large Language Models (LLMs) (Liu et al., 2024; Chang et al., 2024). Companies are increasingly integrating these AI agents into their workflows, seeking more cost-effective and optimized solutions for complex software engineering tasks (Leung and Murphy, 2023). However, direct application of LLMs to software release validation is hindered by limitations in interpretable reasoning and understanding of technical specifications (Marques, 2024; Austin et al., 2021).

To address these challenges, we introduce *GateLens*, a reasoning-enhanced LLM agent (Miehling et al., 2025) to support release validation in the automotive domain. GateLens integrates structured relational analysis with domain-specific expertise, leveraging a reasoning layer built on Relational Algebra (RA) to break down complex validation tasks into systematic analytical steps. GateLens simplifies three critical aspects of release validation:

**1. Test Result Analysis:** Analyzing test execution outcomes is foundational to release validation. This involves analyzing pass/fail patterns across comprehensive test suites, identifying recurring failures, and validating test coverage metrics. In automotive software, where a single release might involve a large number of test cases across multiple vehicle functions, this task becomes particularly demanding. Release engineers must not only identify failed tests but also understand their patterns, assess coverage adequacy, and evaluate test execution stability.

**2. Impact Assessment:** Impact Assessment is a systematic process for evaluating how software issues affect vehicle functionality and safety during release validation. It involves three phases: first, a critical failure analysis identifies the root cause and immediate effects of an issue, such as an ABS module causing a 200ms brake signal delay that exceeds the 100ms threshold. Second, a component-level impact evaluation traces how the issue propagates through interconnected systems, assessing both direct effects, like problematic emergency braking, and indirect effects, such as reduced stability control performance. Finally, an integration risk assessment quantifies the severity of these impacts against safety thresholds and functional requirements, categorizing issues like the ABS delay as system-wide risks with critical severity. This structured process enables engineers to understand system-wide effects, ensuring all safety and functionality requirements are met before release.

**3. Release Candidate Analysis:** The final quality gate involves evaluating Release Candidates (RCs) against predefined quality gates and criteria. This encompasses analyzing whether a particular RC meets all quality thresholds, identifying poten-

tial release blockers, and validating compliance with release requirements. In automotive software, where releases must meet stringent safety and quality standards, this analysis requires careful validation of each RC against established criteria, ensuring all prerequisites for a safe and reliable release are satisfied.

The traditional release validation process demands extensive manual effort. Release engineers meticulously analyze test results, assess impacts, verify RCs against quality gates, and report findings to stakeholders, such as release managers. As automotive software systems grow increasingly complex, these manual workflows become more challenging, time-consuming, and error-prone.

This work aims to streamline release validation by automating key analysis workflows, enabling engineers to focus on high-value analysis and discussion. By providing deeper analytical insights, the proposed approach reduces the time needed to deliver accurate validation results, empowering release managers to make informed decisions more efficiently. Our **contributions** can be summarized as follows:

- We design an *architecture optimized for time- and safety-critical environments*, minimizing LLM invocations while preserving reasoning depth. GateLens operates in a zero-shot setting, avoiding the need for few-shot examples or multi-agent coordination, which improves generalizability, execution speed, reduces maintenance overhead, and enhances transparency.
- We introduce a *scalable and maintainable framework for automotive software release validation*, developed in response to observed limitations in traditional planning-based multi-agent system. GateLens handles diverse user queries with higher robustness and clarity, supporting effective decision-making across a wider range of release engineering tasks.
- We conduct a *comprehensive empirical evaluation*, including comparisons with a CoT+SC system, ablation of the RA module, and performance across multiple LLMs (GPT-4o and Llama 3.1 70B). These experiments demonstrate GateLens's performance advantages in complex and ambiguous queries.
- We report on *real-world industrial deployment* in a partner automotive company, where GateLens reduces analysis time by over 80% and demonstrates strong generalization across user roles, highlighting its practical value and deployment-readiness in safety-critical release validation workflows.

## II. BACKGROUND AND MOTIVATION

Software release decisions in the automotive industry involve multiple stakeholders and extensive data analysis. Modern vehicles integrate hundreds of software components, each requiring rigorous testing and validation. The process advances through distinct phases: component-level testing, integration testing, system-level validation, and vehicle validation testing. Component-level testing verifies individual software modules,

integration testing ensures proper interaction between components, system-level validation examines the complete system behavior, and vehicle validation testing evaluates software performance under real vehicle conditions on closed tracks.

The development cycle grows in complexity with each integration phase. This increasing complexity presents challenges in managing large-scale test results, tracking interdependencies between components, correlating test failures across different subsystems, and maintaining historical context for recurring issues. The iterative nature of software testing and validation further expands this data ecosystem.

The wide range of stakeholders in the release process creates additional challenges in data interpretation and presentation. Project managers need high-level progress indicators, verification engineers require detailed technical insights, quality engineers focus on trend analysis and improvement metrics, and release engineers need specific release-readiness indicators. This variety of perspectives necessitates different views of the same underlying data, making the analysis process more complex.

Release managers function as gatekeepers in the software deployment pipeline. They handle test result analysis, cross-system impact assessment, decision-making, stakeholder coordination, and safety compliance verification. The manual workflow introduces vulnerabilities: time-intensive processing, potential errors in interpretation, decision delays, and communication barriers between technical and business teams.

Within this process, statisticians provide an overall view of the data to project managers and quality engineers for future business decisions. The existing manual approach faces several limitations, particularly regarding time and resource constraints. These include labor-intensive data analysis, delayed response to critical issues, limited capacity for comprehensive analysis, and bottlenecks in the release pipeline. Communication challenges further complicate the process, with misalignment between technical analysts and statisticians, varying interpretations of project requirements, inconsistent reporting formats, and knowledge transfer gaps.

Internal Testing on Closed Track represents a crucial validation phase. Release managers must analyze extensive datasets to evaluate progression readiness. The manual query process for report generation can impact release timelines, business objectives, and subsystem integration schedules.

The deployment of intelligent assistants presents opportunities to address these challenges through automated data processing and analysis, standardized reporting frameworks, real-time insights generation, and stakeholder-specific view generation. However, current automation solutions, including basic LLM implementations, face limitations in understanding complex technical specifications, maintaining structured analysis steps, handling domain-specific requirements, and processing automotive validation data systematically.

These limitations highlight the need for enhanced solutions combining domain expertise with advanced analytical reasoning capabilities. Such solutions must facilitate efficient decision-making while maintaining high safety and quality
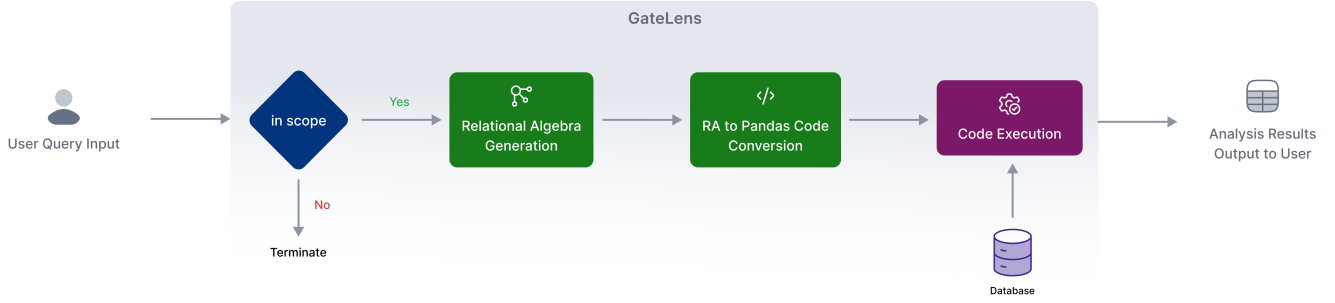
Fig. 1: GateLens top-level architecture: The system processes high-level queries from the end user, generates the necessary data manipulation code using the enhanced reasoning layer with the help of RA, executes it, and outputs the result table as a decision-support resource.

standards in automotive software development (Wang et al., 2024a). Effective intelligent assistants can transform the release decision process, enabling release managers to prioritize result interpretation and strategic decision-making over routine data analysis.

## III. APPROACH AND METHODOLOGY

The complexity of automotive software release validation demands a system that can bridge the gap between human-centric inquiries and precise technical analysis. GateLens addresses this challenge by utilizing LLM agents to transform natural language queries into actionable insights through systematic analysis. At its core, GateLens must fulfill three fundamental requirements:

**1. Query Understanding:** The system must accurately interpret diverse user queries, ranging from high-level management questions to detailed technical inquiries about specific test cases.

**2. Query Transformation:** The system needs to transform these interpretations into structured formal expressions, ensuring consistent and verifiable reasoning structures.

**3. Analysis Execution:** The system must generate and execute precise analysis code that processes validation data according to these formal expressions.

The architecture of GateLens is driven by these fundamental requirements, establishing a systematic pipeline for transforming user queries into analytical results. The system leverages RA to enhance LLMs' reasoning capabilities and transform user queries into formal relational operations. To support this transformation, GateLens employs domain-specific data schemas that guide its relational modeling and enhance the generation of RA expressions. This approach ensures both accessibility for non-technical stakeholders and the precision required for automotive software validation.

### A. System Overview

The primary objective of GateLens is to generate executable code that performs precise test data analysis based on user queries. The system's workflow consists of two main phases: query interpretation and code generation. As illustrated in Figure 1, GateLens first processes user queries through an LLM agent that translates natural language inputs into formal RA expressions. This translation incorporates a detailed

relational model of the test data, ensuring precise specification of automotive domain concepts. The resulting RA expressions serve as a pivotal intermediate representation that is more transparent to both LLM agents and humans. In the second phase, these formal expressions are passed to the coder agent, which generates executable desirable code, such as SQL or Python code, to perform the required analysis on the test data and produce results.

### B. Core Components

The system architecture consists of two primary components that work in tandem to transform natural language queries into executable code: the query interpreter agent and the coder agent. The query interpreter agent first translates user queries into RA expressions, providing a structured framework for analytical reasoning. The coder agent then converts these formal expressions into executable code, completing the transformation pipeline. This two-stage approach ensures both analytical precision and efficient implementation, where the prompt engineering flow and prompt structure are presented in Figure 2.

*1) Query Interpreter:* The query interpreter agent is responsible for converting user queries into formal RA expressions, providing a precise framework for analytical reasoning. Before initiating this translation, the agent consults the knowledge base, comprising the data schema and domain-specific context. The data schema provides a detailed understanding of the dataset, including its relational modeling, detailed field descriptions, and data types. Using this information, the agent verifies whether the query is relevant and within the scope of the dataset (Manik et al., 2021). This validation step ensures that only supported and meaningful queries are processed, improving both accuracy and efficiency. Once the query is confirmed to be in scope, the agent leverages the data schemas to interpret and decompose the query into formal RA expressions.

The agent's primary function is to map natural language queries into formal RA expressions, enhancing LLM reasoning through structured decomposition (Khot et al., 2022). This approach extends traditional Chain-of-Thoughts (CoT) (Wei et al., 2022) reasoning by constraining the model to think within a formal system framework (Zhang et al., 2023a). Instead of generating free-form solutions, the agent must

express analytics through a limited set of standard operations: selection, projection, union, set difference, cartesian product, and rename as basic operations, as well as derived operations such as join, intersection, and division and complemented by aggregation functions like average, minimum, maximum, sum, and count.

By limiting operations to this standard set, the agent effectively handles ambiguous queries through formal translation, ensures technical precision, and prevents deviation from analytical requirements. The formal nature of RA enables query optimization, which the agent incorporates by prioritizing data reduction operations early in the expression chain. This optimization strategy involves applying filters first, then performing expensive operations on the reduced dataset, thereby minimizing processing time and resource utilization.

The translation to RA offers two significant advantages. First, it makes the analytics more transparent in technical terms, allowing for clear interpretation and validation of the reasoning process. Second, it ensures that every solution generated is precisely defined and feasible for implementation, preventing the agent from proposing impractical or undefined analytical approaches.

*2) Coder:* The coder agent is responsible for generating executable code from given RA expressions. Upon receiving an RA expression, the agent follows precise instructions to produce code that delivers the final analytical results. This capability allows the agent to generate complete, self-contained code at once, eliminating the need for step-by-step generation and execution phases.

To ensure the generated code meets quality standards, we designed prompts that specifically instruct the LLM to include essential validation mechanisms. The prompts explicitly require data type validation instructions for numerical and categorical field handling, null value checking procedures to maintain data integrity, validation of numerical operations, and validation of join conditions to ensure proper matching of key columns between tables.

By generating the entire code in a single pass rather than through iterative refinement, the agent significantly reduces processing overhead, system response time, and resource consumption while minimizing potential errors that could arise from multiple execution steps. This streamlined yet precise approach ensures both efficiency and reliability in the analysis pipeline, fully aligned with the formal rigor of RA expressions and improving overall system responsiveness to user queries.

### C. Data Handling

A key architectural decision in GateLens is its indirect interaction with test data. Rather than exposing sensitive test data directly to LLM agents, which could raise privacy concerns (Boudewijn et al., 2023) and exceed input context limitations, the system operates on data schemas and relational models. This approach serves multiple critical purposes:

- **Privacy Protection:** Sensitive automotive test data remains secure within the organization's infrastructure



Fig. 2: Overview of the prompt engineering flow and prompt structure within the GateLens system architecture.

- **Scalability:** The system can handle large-scale test datasets that would exceed LLM context windows
- **Knowledge Integration:** Data schemas and relational models serve as an essential knowledge base, providing necessary structural understanding without raw data exposure

The final execution of the generated code runs on the test data in the target environment, maintaining data privacy while delivering precise analytical results.

### IV. EXPERIMENTAL EVALUATION

In this section, we aim to answer the following research questions:

RQ1: How effectively does GateLens address user queries and deliver accurate results across various query categories?

RQ2: How robust is GateLens in handling out of scope queries and imprecise queries?

RQ3: How does the RA reasoning procedure contribute to the overall performance of GateLens?

RQ4: How does the number of few-shot examples impact GateLens's performance?

### A. Experimental Setup

To address the research questions introduced in Section IV, we designed and conducted extensive experiments to evaluate the performance of GateLens.

The experimental data comprises two distinct benchmarks. The first benchmark consists of 50 queries designed with the assistance of release engineers, quality engineers, and verification engineers. To assess GateLens's performance across a spectrum of query complexities, these queries are categorized into four difficulty levels. The four levels of query difficulty are defined as follows:

<u>**Level 1**</u> Simple queries involving a single operation such as filtering or sorting.

<u>**Level 2**</u> Queries combining two or three basic operations, such as multiple filtering followed by sorting.

<u>**Level 3**</u> Queries involving more than three operations, potentially including grouping and aggregating.

<u>**Level 4**</u> Complex queries requiring multiple advanced operations beyond basic filtering and sorting, such as grouping and aggregating for statistical calculations.

The second benchmark is derived from real-world user queries collected from production logs at our partner company. These queries were sourced from the historical logs of an agentic system that employed a well-established tabular data reasoning approach combining CoT (Wei et al., 2022) prompting with Self-Consistency (SC) (Wang et al., 2022). This system was used in production to support software release analytics, and the collected queries reflect a wide range of user roles, query types, and domain-specific requirements. While this system performs effectively in many scenarios, its limitations become apparent as the range of roles and users expands, leading to a significant diversification of queries. This broader query diversity exposes the system's reliance on few-shot examples, making it less capable of handling highly complex, ambiguous, or ill-defined queries that require greater flexibility and adaptability. Nevertheless, this preliminary system played a critical role in data collection for GateLens by providing query logs used to develop and validate our approach. From these logs, we filtered out near-duplicates and selected 244 frequently repeated unique queries, which we then organized into eight functional categories based on their purposes.

In order to assess GateLens's performance, we run experiments with two large language models (LLMs): GPT-4o, a leading commercial model, and Llama 3.1 70b, a recently released open-source model. We also benchmark GateLens against the CoT+SC agentic system currently used to support the company's release decisions. Our comparative analysis is designed to quantify the improvements introduced by GateLens's novel architecture.

To address the challenge of handling out-of-scope user queries during real-time interactions, GateLens incorporates an in-scope filtering mechanism as explained in Section III-B. This mechanism ensures that the system only attempts to process queries that fall within its scope, thereby improving reliability and reducing errors. Performance evaluation focused on two key aspects:

1) **Quality of responses**: Measured using precision, recall, and F1 Score, which reflect the system's ability to address relevant queries correctly.

2) **Coverage of relevant queries**: Ensuring the system does not reject a significant proportion of valid queries, thus maintaining broad applicability.

Additionally, an ablation study is conducted to examine the contribution of the RA reasoning mechanism of GateLens.

In our experiments, the evaluation of system performance is based on the following definitions: A **True Positive (TP)** occurs when the system produces a result that matches the manually generated ground-truth result, specifically when the final output of the executed code matches the expected ground-truth output. A **False Positive (FP)** occurs when the system provides an incorrect result, meaning the executed code produces output that differs from the ground-truth. A **False Negative (FN)** occurs when the system fails to provide any result for a query. **True Negatives (TNs)** are not applicable since we focus on valid queries producing meaningful output.

Based on these definitions, we calculate Precision, Recall, and F1 scores to assess the performance of the system. Precision ensures that incorrect results are minimized, recall ensures relevant queries are addressed, and the F1 score balances the two to provide an overall assessment of system performance.

### B. Performance in Addressing User Queries (RQ1)

We conducted experiments to compare the performance of GateLens across the two introduced benchmarks. The first benchmark, consisting of 50 queries categorized by difficulty levels, was used to evaluate and compare the performance of GateLens and CoT+SC. Both systems were tested using GPT-4o and Llama 3.1 70B as their underlying LLMs. The results are summarized in Table I.

The results demonstrate that GateLens with GPT-4o significantly outperforms GateLens with Llama 3.1 70B, indicating GPT-4o's superior capability for interpreting and generating RA. Similarly, CoT+SC with GPT-4o outperforms its Llama 3.1 70B variant, with the performance gap growing as query complexity increases. CoT+SC performance declines with query complexity. This underscores the importance of the RA reasoning mechanism in GateLens, which enables effective handling of complex, unstructured queries by decomposing them into logical, structured expressions. Most notably, GateLens with GPT-4o achieved optimal performance on this benchmark, maintaining 100% accuracy across all difficulty levels. This stems from integrating RA reasoning into our framework. By translating queries into RA expressions, GateLens explicitly captures the logical structure of operations, enhancing both the clarity and precision of the generated code. The intermediate RA conversion allows the system to focus on the relevant table operations while filtering out irrelevant elements in the query, greatly enhancing the problem-solving capabilities of the LLM agent.

For the second benchmark, results in Table II show that GateLens (GPT-4o) and CoT+SC (GPT-4o) significantly outperformed GateLens with Llama 3.1 70B. This performance disparity is primarily due to the strict code generation requirements of the task, including table filtering, merging

TABLE I: Performance comparison across different models and difficulty levels on the first benchmark, which consists of 50 designed queries with annotated difficulty levels.

| Level | # Queries | GateLens with GPT-4o | | | GateLens with Llama 3.1 70B | | | CoT+SC with GPT-4o | | | CoT+SC with Llama 3.1 70B | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Precision | Recall | F1 Score | Precision | Recall | F1 Score | Precision | Recall | F1 Score | Precision | Recall | F1 Score |
| 1 | 16 | 100% | 100% | **100%** | 100% | 43.75% | 60.87% | 93.33% | 87.5% | 90.32% | 100% | 93.75% | 96.77% |
| 2 | 16 | 100% | 100% | **100%** | 100% | 62.5% | 76.92% | 100% | 81.25% | 89.66% | 92.31% | 75% | 82.76% |
| 3 | 12 | 100% | 100% | **100%** | 100% | 50% | 66.67% | 91.67% | 91.67% | 91.67% | 90.91% | 83.33% | 86.96% |
| 4 | 6 | 100% | 100% | **100%** | 100% | 33% | 49.62% | 66.67% | 66.67% | 66.67% | 60% | 50% | 54.55% |
| **Total** | **50** | **100%** | **100%** | **100%** | **100%** | **47.31%** | **63.52%** | **87.91%** | **81.77%** | **84.57%** | **85.81%** | **75.52%** | **80.26%** |

TABLE II: Performance comparison of GateLens and the CoT+SC system across different categories on the second benchmark, which consists of 244 real-world queries.

| Category | # Queries | GateLens with GPT-4o | | | GateLens with Llama 3.1 70B | | | CoT+SC with GPT-4o | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Precision | Recall | F1 Score | Precision | Recall | F1 Score | Precision | Recall | F1 Score |
| Column Operations | 17 | 64.7% | 64.7% | 64.7% | 50% | 11.76% | 19.04% | 76.47% | 76.47% | **76.47%** |
| Complex Multi-Condition Queries | 77 | 86.3% | 81.82% | **84%** | 100% | 24.68% | 39.59% | 84.75% | 64.94% | 73.53% |
| Conditional Calculations | 8 | 100% | 87.5% | **93.3%** | 100% | 37.5% | 54.55% | 87.5% | 87.5% | 87.5% |
| Data Filtering | 32 | 89.66% | 81.25% | **85.25%** | 90.91% | 31.25% | 46.51% | 86.21% | 78.13% | 81.97% |
| Duplicate Removal | 78 | 87.67% | 82.05% | **84.77%** | 100% | 23.08% | 37.5% | 75.93% | 52.56% | 62.12% |
| Grouping and Aggregation | 10 | 80.0% | 80.0% | **80.0%** | 100% | 30% | 46.15% | 83.33% | 50% | 62.5% |
| Metadata Queries | 13 | 91.67% | 84.61% | **88.0%** | 100% | 15.38% | 26.66% | 46.15% | 46.15% | 46.15% |
| Table Generation | 9 | 88.89% | 88.89% | **88.89%** | 100% | 44.44% | 61.53% | 88.89% | 88.89% | **88.89%** |
| **Total** | **244** | **86.02%** | **81.14%** | **83.51%** | **92.61%** | **27.26%** | **41.44%** | **83.15%** | **63.52%** | **70.61%** |

strategies, and key-value mapping operations, where GPT-4o demonstrated markedly superior capabilities.

GateLens with GPT-4o outperformed CoT+SC (GPT-4o) across most categories, particularly evident in Metadata Queries (those seeking basic table information). For example, when processing the query "Give me the list of release candidates," the CoT+SC system often fails to identify the correct field. A common failure mode in CoT+SC occurred when user queries included typographical errors or incorrect casing in field names, with the system directly using the erroneous fields without correction. GateLens addresses this limitation through its query-to-RA transformation process, which incorporates the database's relational model, adjusts query fields to match table formats, and can handle fuzzy matching to detect and correct field names, enabling the system to resolve typographical errors and ambiguous queries effectively. This approach improves accuracy and resilience, particularly in real-world scenarios where user queries may not always adhere to strict formatting standards.

Across both datasets, GateLens with GPT-4o consistently delivered the best performance, achieving perfect accuracy across difficulty levels in the first benchmark and superior handling of diverse query categories in the second benchmark. This superior performance stems from the system's ability to extract logical conditions from natural language queries and transform them into RA expressions, clarifying query logic, ensuring accurate, complete condition handling and minimizing errors and omissions in the generated code. These results demonstrate the framework's enhancement of LLM performance on complex queries and improved real-world reliability

### C. Robustness: Handling Out of Scope and Imprecise Queries (RQ2)

To assess the robustness of our approach in handling diverse user queries under real-world conditions, we conducted further experiments focusing on filtering out-of-scope queries as well as processing imprecise queries. For this purpose, the data analysis team at our industrial partner company manually selected 37 out-of-scope queries and 50 imprecise queries from the historical logs of the first-generation system, which are used to perform targeted evaluations.

Out-of-scope queries are those that cannot be meaningfully answered using the available data. For example, a query like "What is the most beautiful truck?" requires subjective judgment and cannot be resolved through database operations; it should be identified and filtered as out of scope. On the other hand, imprecise queries are those that can be answered using the database but contain ambiguous or inexact terms. For instance, a query such as "Find some trucks for cases that are NOK" is considered imprecise because while it seeks truck names where test results are "NOK" (failed), it uses ambiguous terminology - referring to "trucks" instead of the actual database field "name", and mentions "NOK" without specifying the "test_result" field. Such imprecise queries require mapping informal language to precise database fields and conditions for proper execution.

*1) Handling Out of Scope Queries:* We compared GateLens with other models; the results can be found in Table III. The results demonstrate that GateLens with GPT-4o achieved the best performance, particularly in terms of precision, which is approximately 40% higher than other models, indicating GateLens' ability to avoid generating incorrect results.

The superior precision of GateLens with GPT-4o can be attributed to two key aspects of its design. First, its robust filtering mechanism ensures that out-of-scope queries are identified and excluded early in the processing pipeline, preventing irrelevant results. Second, the conversion of raw natural language queries into structured RA expressions enables the model to isolate and capture task-relevant components of a query. This structured approach considerably decreases erroneous outcomes and enhances the model's ability to handle complex and diverse query formulations in real-world scenarios.

CoT+SC showed significantly lower precision due to the

TABLE III: Model comparison for out-of-scope queries.

| Model | Precision | Recall | F1 Score |
|---|---|---|---|
| **GateLens with GPT-4o** | 92.5% | 100% | **96.10%** |
| **GateLens with Llama 3.1 70B** | 52.94% | 97.30% | 68.57% |
| **CoT+SC with GPT-4o** | 51.10% | 89.19% | 64.97% |

TABLE IV: Model comparison for imprecise queries.

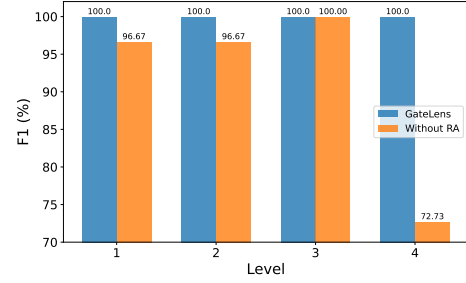| Model | Precision | Recall | F1 Score |
|---|---|---|---|
| **GateLens with GPT-4o** | 92.86% | 78% | **84.78%** |
| **GateLens with Llama 3.1 70B** | 92.86% | 26% | 40.63% |
| **CoT+SC with GPT-4o** | 90% | 36% | 51.43% |

variability of real-world queries and the inconsistency of user narratives, which often contain a mix of relevant and irrelevant content. This variability increases uncertainty and poses challenges for models that struggle to identify task-relevant information. Although all models demonstrated high recall, this did not translate into accurate processing.

*2) Handling Imprecise Queries:* To further assess the robustness of our approach, we evaluated its performance in handling imprecise queries, which posed two primary challenges. First, some queries are informal and conversational in style, appearing unrelated to data analysis but actually carrying relevant intent. Second, many queries referred to fields using terms differing from the column headers.
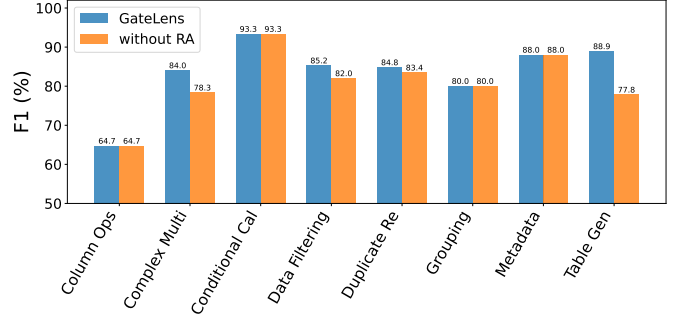
The results of these experiments are presented in Table IV. As shown, GateLens with GPT-4o demonstrates the best overall performance. In terms of precision, all methods performed relatively well, indicating that when results are generated, they are likely to be correct. However, our method significantly outperformed the others in recall, highlighting its ability to handle a larger portion of the imprecise queries. As a result, GateLens with GPT-4 achieved a substantially higher F1 score compared to other methods, demonstrating that it not only processes most queries but also produces accurate results for them.

The observed performance gap between GateLens with GPT-4o and the other models can be attributed to their inherent limitations. Specifically, the Llama 3.1 70B model struggled to interpret user queries that deviated from the exact column header descriptions in the database schema. In such cases, Llama 3.1 70B often converted only the clearly defined parts of the query into RA, leading to incomplete execution and reduced accuracy. On the other hand, CoT+SC exhibits low recall, as it is highly susceptible to confusion by ambiguous query elements. This causes CoT+SC to frequently generate incorrect code that fails execution, significantly lowering its recall rate.

Overall, our approach demonstrated a strong ability to handle imprecise queries by maintaining high precision and significantly improving recall. This robustness ensures that the system is adaptable to variations and ambiguities in query formulations, allowing users to freely express their queries without being constrained by strict adherence to column header names. Such flexibility enhances the system's practicality and applicability in real-world scenarios, making it a reliable and user-friendly tool for diverse interactions.



(a) The first benchmark with annotated difficulty levels.



(b) The second benchmark with real-world user queries

Fig. 3: Comparison of the original method and the method without the RA module across different datasets.

### D. Effectiveness of the RA module (RQ3)

To evaluate the impact of the RA module that converts user queries into RA expressions, we conducted experiments by removing the RA module from the framework and comparing the results to the original system. The outcomes, shown in Figure 3, demonstrate significant performance degradation across both benchmarks when operating without the RA module.

In the first benchmark, performance declined most notably for Level 4 queries, showing a drop exceeding 27%. These queries, which involve advanced operations like grouping, aggregating, and statistical calculations, demonstrated that RA translation is particularly crucial for handling queries with multiple, intricate operations. Similarly, the second benchmark shows decreased performance in complex tasks such as multi-condition filtering, duplicate data removal, and table generation, further emphasizing RA's effectiveness in managing complex database operations.

The RA module maintained consistent performance for simpler queries, demonstrating its versatility across varying complexity levels. By transforming natural language into precise, logical representations, the RA module serves as a critical bridge between user intent and code execution. This translation process enables the code generator to produce accurate, efficient executable code for data analysis tasks.

These results establish RA as a fundamental component in our system's architecture, significantly enhancing its reasoning capabilities and ensuring reliable code generation across diverse query types and complexity levels. This enables the system to deliver accurate and efficient solutions for tabular data analysis tasks.

### E. Impact of Few-Shot Examples (RQ4)

To investigate the effect of including few-shot examples in prompts, we conducted experiments by varying the number of examples provided to both GateLens and CoT+SC. This experiment is performed on the first benchmark containing 50 designed queries, with results illustrated in Figure 4.

The results demonstrate that GateLens relies heavily on its RA translation process, achieving optimal performance even in a 0-shot setting without any examples. Interestingly, when a small number of examples are added, GateLens becomes slightly biased toward them, leading to a minor degradation in performance. However, it regains optimal performance with additional examples as the system learns to generalize better. In contrast, CoT+SC's performance heavily depends on in-context examples, showing suboptimal results without sufficient few-shot examples. While CoT+SC's performance improves with more examples, this approach is fundamentally inefficient and impractical. Increasing example count expands input context size, leading to higher latency due to the quadratic complexity of Transformer-based models, particularly problematic for real-time applications. This also escalates operational costs through increased token usage (higher cloud API fees) and computational demands. Moreover, the larger context risks exceeding the maximum length, which can lead to truncation or lost information, compromising the model's response quality. In long contexts, critical content may be ignored, resulting in a phenomenon known as "lost in the middle," which adversely affects overall performance (Liu et al., 2023).

The results further highlight the critical role of RA in enabling effective query handling. GateLens, through its RA-based translation process, achieves robust performance without requiring extensive in-context learning. GateLens's independence from examples not only ensures consistent performance across diverse queries but also offers practical benefits by mitigating issues related to context length limitations and accordingly reducing computational and financial resource consumption, enhancing system scalability, and minimizing maintenance needs for adaptation to new tasks. These advantages position GateLens as a robust and efficient solution for automating the analysis of tabular data in dynamic environments.

## V. Industrial Deployment: Lessons Learned

The deployment of GateLens at a partner automotive company has provided valuable insights into integrating AI-assisted analytics into complex industrial workflows, specifically for streamlining decision-making in automotive software release validation.

Automotive software integration at the company typically occurs across three hierarchical stages: subsystem (control unit), system (multiple control units), and full vehicle levels. Each stage involves extensive testing, with results stored in a central database. Critical Go/No-Go decisions are made at these stages to determine whether a release meets quality thresholds. However, stakeholders from diverse
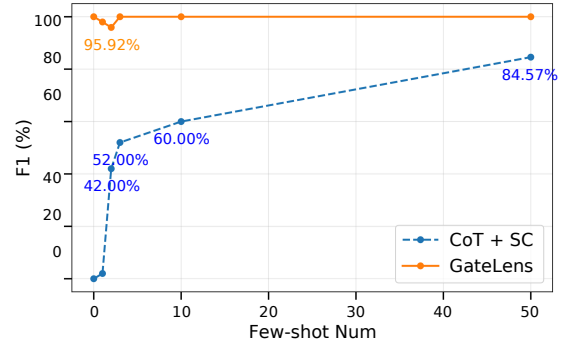


Fig. 4: Comparison of GateLens against CoT+SC across different numbers of few-shot examples.

backgrounds—including project managers, mechanical engineers, and software engineers—must query the raw data to evaluate product quality. Many lack expertise in data analytics, creating bottlenecks and delays in the decision-making process.

Previously, these analytics were managed by a small team of 2–3 full-time analysts, who were often overwhelmed by the volume and diversity of requests. Scaling the team to meet the current demand would have required tripling its size. GateLens addresses this challenge by automating much of the workload, enabling more efficient decision-making. Currently, GateLens is in an extended pilot phase, supporting a pool of 60-80 users. The analytics team has transitioned to a support role, helping stakeholders articulate their needs into clear, actionable prompts for the system. User adoption of GateLens has progressed in phases:

- **Small-Scale Pilot**: The initial deployment within the analytics team established benchmarks.
- **Expanded Pilot**: Five additional users from varied backgrounds contributed to refining the benchmarks.
- **Wider Rollout**: The current phase involves a larger group of 60–80 users. Feedback has been highly positive, with stakeholders recognizing GateLens's ability to simplify and accelerate complex analyses.

Since the launch, the number of both new and recurring users has grown, encompassing diverse roles and types of queries, thereby demonstrating the tool's increasing utility and trust. GateLens significantly reduces the time and effort required for complex analyses, but the shift towards automation also requires users to take on more responsibility in defining and clarifying their needs. The transition from a primarily supportive tool to a more fully automated system is ongoing, demanding a gradual approach with careful calibration to ensure the tool continues to meet evolving needs.

The diversity in stakeholders' needs and backgrounds is an inevitable factor, leading to a wide range of requirements and queries when interacting with analytics systems. Our preliminary agentic system based on CoT and SC has performed well when serving a specific group of stakeholders. However, as we opened the system to a broader audience, the variety of queries increased substantially. Systems that rely heavily on few-shot learning face significant challenges in these situations,

TABLE V: GateLens (zero-shot) vs CoT+SC (few-shot) performance across different roles on the second benchmark. For each role tested, CoT+SC was trained using examples from the other two roles only (leave-one-role-out approach).

| Roles | # Queries (244 in total) | GateLens F1 Score | CoT+SC (Few-shot with All Roles) F1 Score | CoT+SC (Few-shot with Leave-One-Role-Out) | | |
|---|---|---|---|---|---|---|
| | | | | Without Mechanic | Without Project | Without Software |
| Mechanically-oriented | 36 | 94.25% | 80.60% | 73.85% ↓ (-6.75%) | 86.57% | 80.60% |
| Project-oriented | 193 | 80.21% | 78.93% | 78.93% | 76.22% ↓ (-2.71%) | 78.40% |
| Software-oriented | 15 | 100% | 100% | 100% | 100% | 82.76% ↓ (-17.24%) |

as they are limited by the relatively small set of few-shot examples, making it difficult to handle a wider range of potentially unexpected inputs. Similarly, fine-tuned models often struggle to adapt to dynamic and diverse environments. Improving generalization is essential to ensure the scalability and reliability of analytics systems in these complex, domain-specific contexts. In our design and development work, we prioritized the system's ability to generalize effectively and meet the needs of a broader and more diverse group of users.

To explore the system's generalizability, we categorized the roles within the company into three groups: mechanically-oriented, project-oriented, and software-oriented roles. Mechanically-oriented roles typically focus on truck-specific data filtering. Project-oriented roles often combine meta-queries with conditional filters for release management and statistical analysis. Software roles emphasize truck software applications and user functions. We can see from Table V, both GateLens (zero-shot) and CoT+SC (few-shot) exhibit differences in system performance across these groups, which is likely stemming from the complexity and variety of their typical queries. Nevertheless, the results demonstrate that GateLens is capable of supporting all groups to a high degree. To further evaluate CoT+SC's dependency on few-shot examples, we conducted a **leave-one-role-out** experiment. In this approach, examples from a specific role are excluded in each iteration. For instance, 'without software' indicates that all examples from the software-oriented role have been removed, while the total number of examples is maintained by substituting them with examples from other roles. This highlights the potential challenges with the robustness and generalizability of techniques that rely on few-shot examples. This is a crucial factor to consider in industries where diverse teams collaborate and a wide range of queries may arise.

The impact of automated systems, such as GateLens, on the release process has been substantial. Compared to the previous manual process, GateLens has reduced the time required for Go/No-Go analytics by more than 80%, significantly improving operational efficiency. Stakeholders can now focus on high-level decision-making, freed from the burden of data preparation and analysis.

A key advantage of GateLens lies in its domain-specific design. Unlike general-purpose tools like TaskWeaver (Qiao et al., 2023) or AutoGen (Wu et al., 2023), GateLens is tailored to automotive workflows, making it easier to understand, debug, and adapt to automotive common procedures. This focus on domain relevance ensures that the system aligns more closely with stakeholders' needs while providing reliable and nuanced support.

In summary, the deployment of GateLens demonstrates how domain-specific AI solutions can transform critical workflows in the automotive sector. By automating labor-intensive processes and enhancing decision-making, GateLens has delivered measurable improvements in efficiency and user satisfaction. However, its success depends on ongoing refinement and careful management of the transition to full(er) automation. Balancing automation with user empowerment remains crucial, particularly in a complex industry like automotive, where diverse stakeholder needs must be met.

GateLens represents a promising step forward, showcasing the potential of AI-driven systems to improve not only the automotive domain but also other industries requiring robust, scalable solutions for intricate processes.

## VI. RELATED WORK

General-purpose LLMs are primarily designed for and trained on natural languages. Working with tabular data requires specialized adaptations to effectively handle its structured and heterogeneous nature (Fang et al., 2024; Wang et al., 2024b; van Breugel and van der Schaar, 2024). First, the structured tabular data is typically transformed into serialized text. The performance of the LLM may depend on this transformation (Min et al., 2024). Subsequently, the serialized text data is used as input to the LLM for various tasks, such as question-answering, summarization, or logical reasoning. Common approaches to improve LLM performance include prompt engineering, pre-training, fine-tuning, and Retrieval-Augmented Generation (RAG).

Pre-training and fine-tuning (Zhang et al., 2023b; Parthasarathy et al., 2024; VM et al., 2024; Dong et al., 2022; Hegselmann et al.) often face scalability concerns. Although resource-efficient training techniques have been proposed to mitigate the substantial computational demands of LLMs (Han et al., 2024; Lin et al., 2024), in safety-critical applications with evolving data and requirements, training LLMs presents significant challenges due to the constant need for rigorous validation and verification. This ongoing necessity substantially increases resource demands for development and maintenance, potentially exceeding the capacities of many companies. Techniques such as RAG have been employed to dynamically integrate external knowledge bases during inference, reducing the need for frequent model updates (Zhao et al., 2024; Gao et al., 2023). However, such methods can pose challenges in safety-critical industrial settings as well, since both retrieval modules and model components must undergo synchronized updates to maintain relevance, reliability, and compliance with validation and verification requirements. Costs would also be especially high with fine-tuning since re-tuning would be needed when new and improved base LLMs are released and should be incorporated.

Prompt engineering techniques are among the most resource-efficient methods for improving LLM output (Sahoo et al., 2024; Jin and Lu, 2023). From a user standpoint, when the input is natural language, prompting techniques can be broadly categorized based on the type of language generated by the LLM. These include outputs in natural language, structured languages (Li et al., 2023), or symbolic languages. When the generated language is natural language, LLMs often fail to consistently follow instructions, particularly when the instructions are complex or require precise, step-by-step execution (Pham et al., 2024). This inconsistency arises because natural language, while flexible and expressive, can be ambiguous and prone to misinterpretation by LLMs. Structured languages include general-purpose languages (e.g. Python) (Ye et al., 2024), query languages (e.g. SQL) (Li et al.; Dong et al., 2023; Mouravieff et al., 2024), configuration formats (e.g. YAML or JSON), or other Domain-Specific Languages (DSLs) (Glenn et al., 2024; Dai et al., 2024). These languages are subsequently interpreted and/or executed by either external tools, the same LLM, or another LLM agent. This approach offers significant advantages, as it enables precise execution of tasks. Another popular type of output is symbolic languages. Literature shows that symbolic representations provide a more rigorous framework for articulating premises and intent, which can enhance reasoning capabilities (Pan et al., 2023).

In this paper, we introduce a novel prompt-only (training-free) approach that bridges natural language and executable code through RA, a symbolic formalism designed for relational modeling and ideally suited for analyzing tabular data. Unlike prior work that often relies on complex multi-agent planning, our approach leverages RA as a lightweight intermediate representation to enable precise query normalization, disambiguation of natural language input, and efficient code generation. RA acts as an abstraction layer that can target multiple execution backends (e.g., Python, SQL), providing adaptability across systems. In GateLens, we generate Python code to support practical industrial deployment and high-performance execution. GateLens is *training-free*, *feed-forward* (single-pass, without looping or multi-agent orchestration), and thus easier to verify, trace, maintain, and trust – qualities critical for safety-critical industrial applications.

## VII. DISCUSSION AND CONCLUSIONS

This study introduced GateLens, a reasoning-enhanced LLM agent designed for automotive software release validation. By introducing RA as an intermediate representation before code generation, Gatelens addresses the "Unfaithful Chain-of-Thought (CoT) reasoning" problem in code generation, where reasoning steps in CoT explanations do not accurately reflect the model's actual thought process (Turpin et al., 2023). Specifically, GateLens divides the analysis into two steps: (1) natural language queries are first translated into RA expressions, and (2) these RA expressions are then converted into executable code. We use Python as our target language in step (2) due to its widespread use in our partner company and the model's strong performance in Python, which benefits

from more extensive training data. However, this step is also compatible with RA-to-SQL generation, enabling flexibility across backends. The inclusion of the RA reasoning module is a critical factor in improving robustness and scalability, as evidenced by superior F1 scores in both benchmarking and industrial evaluations.

In real-world deployment serving 60-80 users, GateLens demonstrates significant practical value through its user-friendly interface and robust query processing capabilities, with users particularly appreciating the flexibility to input, debug, and refine queries easily. This marks a substantial advancement in industrial data interaction, successfully handling complex and ambiguous queries while providing practical support for faster decision-making in safety-critical software release processes. GateLens demonstrates significant practical advancements by reducing analysis time by over 80% while maintaining high accuracy in test result interpretation, impact assessment, and release candidate evaluation.

Our implementation insights highlight the advantages of focusing on training-free and single-pass agent systems by foregrounding the perception phase in the code generation pipeline. This modular architecture opens opportunities for incorporating emerging LLM capabilities while preserving the system's practical utility in safety-critical industrial applications. A phased deployment program is ongoing and shows that multiple stakeholder groups can be supported, though the evolving roles and analytical needs require continued refinement.

Future work will address current limitations by expanding domain applicability beyond automotive software and testing alternative LLM configurations to enhance reliability across other safety-critical industries. This approach demonstrates the potential for LLM-based solutions to transform industrial data analysis workflows while maintaining the reliability standards required for critical applications.

## VIII. THREATS TO VALIDITY

The validity of our findings is subject to several potential threats. First, this framework depends on well-defined, static schemas for accurate query decomposition, which fundamentally limits its effectiveness in environments with incomplete or frequently changing schemas. Second, the generalizability of GateLens beyond the automotive software release domain is limited, as the system relies on domain-specific relational modeling and datasets, necessitating further validation in other safety-critical industries. Third, the benchmarks and query scenarios used for evaluation, derived from historical data and real-world queries, may not fully capture the diversity and complexity of potential use cases, which could impact the robustness of the system in broader deployments. Finally, the system's performance depends on the specific LLM configurations used, such as GPT-4o and Llama 3.1 70B, and their ability to interpret and generate RA expressions. Future work will address these limitations through broader domain testing, expanded evaluation scenarios, and alternative LLM configurations.

## REFERENCES

M. Liu, J. Wang, T. Lin, Q. Ma, Z. Fang, and Y. Wu, "An empirical study of the code generation of safety-critical software using llms," *Applied Sciences*, vol. 14, p. 1046, 2024.

Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang *et al.*, "A survey on evaluation of large language models," *ACM Transactions on Intelligent Systems and Technology*, vol. 15, no. 3, pp. 1–45, 2024.

M. Leung and G. Murphy, "On automated assistants for software development: the role of llms," *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pp. 1737–1741, 2023.

N. Marques, "Using chatgpt in software requirements engineering: a comprehensive review," *Future Internet*, vol. 16, p. 180, 2024.

J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le *et al.*, "Program synthesis with large language models," *arXiv preprint arXiv:2108.07732*, 2021.

E. Miehling, K. N. Ramamurthy, K. R. Varshney, M. Riemer, D. Bouneffouf, J. T. Richards, A. Dhurandhar, E. M. Daly, M. Hind, P. Sattigeri *et al.*, "Agentic ai needs a systems theory," *arXiv preprint arXiv:2503.00237*, 2025.

L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin *et al.*, "A survey on large language model based autonomous agents," *Frontiers of Computer Science*, vol. 18, no. 6, p. 186345, 2024.

L. P. Manik, Z. Akbar, H. F. Mustika, A. Indrawati, D. S. Rini, A. D. Fefirenta, and T. Djarwaningsih, "Out-of-scope intent detection on a knowledge-based chatbot." *International Journal of Intelligent Engineering & Systems*, vol. 14, no. 5, 2021.

T. Khot, H. Trivedi, M. Finlayson, Y. Fu, K. Richardson, P. Clark, and A. Sabharwal, "Decomposed prompting: A modular approach for solving complex tasks," *arXiv preprint arXiv:2210.02406*, 2022.

J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.

Z. Zhang, Y. Yao, A. Zhang, X. Tang, X. Ma, Z. He, Y. Wang, M. Gerstein, R. Wang, G. Liu *et al.*, "Igniting language intelligence: The hitchhiker's guide from chain-of-thought reasoning to language agents," *arXiv preprint arXiv:2311.11797*, 2023.

A. T. P. Boudewijn, A. F. Ferraris, D. Panfilo, V. Cocca, S. Zinutti, K. De Schepper, and C. R. Chauvenet, "Privacy measurements in tabular synthetic data: State of the art and future research directions," in *NeurIPS 2023 Workshop on Synthetic Data Generation with Generative AI*, 2023.

X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," *arXiv preprint arXiv:2203.11171*, 2022.

N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang, "Lost in the middle: How language models use long contexts," *arXiv preprint arXiv:2307.03172*, 2023.

B. Qiao, L. Li, X. Zhang, S. He, Y. Kang, C. Zhang, F. Yang, H. Dong, J. Zhang, L. Wang *et al.*, "Taskweaver: A code-first agent framework," *arXiv preprint arXiv:2311.17541*, 2023.

Q. Wu, G. Bansal, J. Zhang, Y. Wu, S. Zhang, E. Zhu, B. Li, L. Jiang, X. Zhang, and C. Wang, "Autogen: Enabling next-gen llm applications via multi-agent conversation framework," *arXiv preprint arXiv:2308.08155*, 2023.

X. Fang, W. Xu, F. A. Tan, J. Zhang, Z. Hu, Y. J. Qi, S. Nickleach, D. Socolinsky, S. Sengamedu, C. Faloutsos *et al.*, "Large language models (llms) on tabular data: Prediction, generation, and understanding-a survey," 2024.

Z. Wang, H. Zhang, C.-L. Li, J. M. Eisenschlos, V. Perot, Z. Wang, L. Miculicich, Y. Fujii, J. Shang, C.-Y. Lee *et al.*, "Chain-of-table: Evolving tables in the reasoning chain for table understanding," *arXiv preprint arXiv:2401.04398*, 2024.

B. van Breugel and M. van der Schaar, "Why tabular foundation models should be a research priority," *arXiv preprint arXiv:2405.01147*, 2024.

D. Min, N. Hu, R. Jin, N. Lin, J. Chen, Y. Chen, Y. Li, G. Qi, Y. Li, N. Li *et al.*, "Exploring the impact of table-to-text methods on augmenting llm-based question answering with domain hybrid data," *arXiv preprint arXiv:2402.12869*, 2024.

T. Zhang, X. Yue, Y. Li, and H. Sun, "Tablellama: Towards open large generalist models for tables," *arXiv preprint arXiv:2311.09206*, 2023.

V. B. Parthasarathy, A. Zafar, A. Khan, and A. Shahid, "The ultimate guide to fine-tuning llms from basics to breakthroughs: An exhaustive review of technologies, research, best practices, applied research challenges and opportunities," *arXiv preprint arXiv:2408.13296*, 2024.

K. VM, H. Warrier, Y. Gupta *et al.*, "Fine tuning llm for enterprise: Practical guidelines and recommendations," *arXiv preprint arXiv:2404.10779*, 2024.

H. Dong, Z. Cheng, X. He, M. Zhou, A. Zhou, F. Zhou, A. Liu, S. Han, and D. Zhang, "Table pre-training: A survey on model architectures, pre-training objectives, and downstream tasks," *arXiv preprint arXiv:2201.09745*, 2022.

S. Hegselmann, A. Buendia, H. Lang, M. Agrawal, X. Jiang, and D. Sontag, "TabLLM: Few-shot Classification of Tabular Data with Large Language Models."

Z. Han, C. Gao, J. Liu, J. Zhang, and S. Q. Zhang, "Parameter-efficient fine-tuning for large models: A comprehensive survey," *arXiv preprint arXiv:2403.14608*, 2024.

X. Lin, W. Wang, Y. Li, S. Yang, F. Feng, Y. Wei, and T.-S. Chua, "Data-efficient fine-tuning for llm-based recommendation," in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 365–374. [Online]. Available: https://doi.org/10.1145/3626772.3657807

P. Zhao, H. Zhang, Q. Yu, Z. Wang, Y. Geng, F. Fu, L. Yang, W. Zhang, and B. Cui, "Retrieval-augmented generation for ai-generated content: A survey," *arXiv preprint arXiv:2402.19473*, 2024.

Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, and H. Wang, "Retrieval-augmented generation for large language models: A survey," *arXiv preprint arXiv:2312.10997*, 2023.

P. Sahoo, A. K. Singh, S. Saha, V. Jain, S. Mondal, and A. Chadha, "A systematic survey of prompt engineering in large language models: Techniques and applications," *arXiv preprint arXiv:2402.07927*, 2024.

Z. Jin and W. Lu, "Tab-cot: Zero-shot tabular chain of thought," *arXiv preprint arXiv:2305.17812*, 2023.

C. Li, J. Liang, A. Zeng, X. Chen, K. Hausman, D. Sadigh, S. Levine, L. Fei-Fei, F. Xia, and B. Ichter, "Chain of code: Reasoning with a language model-augmented code emulator," *arXiv preprint arXiv:2312.04474*, 2023.

C. M. Pham, S. Sun, and M. Iyyer, "Suri: Multi-constraint instruction following for long-form text generation," *arXiv preprint arXiv:2406.19371*, 2024.

J. Ye, M. Du, and G. Wang, "DataFrame QA: A Universal LLM Framework on DataFrame Question Answering Without Data Exposure," 2024, version Number: 1. [Online]. Available: https://arxiv.org/abs/2401.15463

J. Li, B. Hui, G. Qu, J. Yang, B. Li, B. Li, B. Wang, B. Qin, R. Geng, N. Huo, X. Zhou, C. Ma, G. Li, K. C. C. Chang, F. Huang, R. Cheng, and Y. Li, "Can LLM Already Serve as A Database Interface? A BIg Bench for Large-Scale Database Grounded Text-to-SQLs."

X. Dong, C. Zhang, Y. Ge, Y. Mao, Y. Gao, J. Lin, D. Lou *et al.*, "C3: Zero-shot text-to-sql with chatgpt," *arXiv preprint arXiv:2307.07306*, 2023.

R. Mouravieff, B. Piwowarski, and S. Lamprier, "Learning relational decomposition of queries for question answering from tables," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 10 471–10 485. [Online]. Available: https://aclanthology.org/2024.acl-long.564/

P. Glenn, P. P. Dakle, L. Wang, and P. Raghavan, "Blendsql: A scalable dialect for unifying hybrid question answering in relational algebra," *arXiv preprint arXiv:2402.17882*, 2024.

H. Dai, B. Wang, X. Wan, B. Dai, S. Yang, A. Nova, P. Yin, M. Phothilimthana, C. Sutton, and D. Schuurmans, "UQE: A query engine for unstructured databases," *Advances in Neural Information Processing Systems*, vol. 37, pp. 29 807–29 838, 2024.

L. Pan, A. Albalak, X. Wang, and W. Y. Wang, "Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning," *arXiv preprint arXiv:2305.12295*, 2023.

M. Turpin, J. Michael, E. Perez, and S. Bowman, "Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting," *Advances in Neural Information Processing Systems*, vol. 36, pp. 74 952–74 965, 2023.