# TAAF: A Trace Abstraction and Analysis Framework Synergizing Knowledge Graphs and LLMs

Alireza Ezaz
Brock University
Computer Science
St. Catharines, Ontario, Canada
sezaz@brocku.ca

Ghazal Khodabandeh
Brock University
Computer Science
St. Catharines, Ontario, Canada
gkhodobandeh@brocku.ca

Majid Babaei
Mcgil University
Computer Science
Montreal, Quebec, Canada
majid.babaei@mcgill.ca

Naser Ezzati-Jivan
Brock University
Computer Science
St. Catharines, Ontario, Canada
nezzatijivan@brocku.ca

## Abstract

Execution traces are a critical source of information for understanding, debugging, and optimizing complex software systems. However, traces from OS kernels or large-scale applications like Chrome or MySQL are massive and difficult to analyze. Existing tools rely on predefined analyses, and custom insights often require writing domain-specific scripts, which is an error-prone and time-consuming task. This paper introduces **TAAF** (Trace Abstraction and Analysis Framework), a novel approach that combines time-indexing, knowledge graphs (KGs), and large language models (LLMs) to transform raw trace data into actionable insights. TAAF constructs a time-indexed KG from trace events to capture relationships among entities such as threads, CPUs, and system resources. An LLM then interprets query-specific subgraphs to answer natural-language questions, reducing the need for manual inspection and deep system expertise. To evaluate TAAF, we introduce **TraceQA-100**, a benchmark of 100 questions grounded in real kernel traces. Experiments across three LLMs and multiple temporal settings show that TAAF improves answer accuracy by up to 31.2%, particularly in multi-hop and causal reasoning tasks. We further analyze where graph-grounded reasoning helps and where limitations remain, offering a foundation for next-generation trace analysis tools.

## 1 Introduction

Execution traces are a foundational source of evidence for understanding, debugging, and optimizing software systems. They capture fine-grained temporal telemetry from across the stack, ranging from user-level function calls and inter-service messages to low-level kernel events such as system calls, context switches, and interrupt handling. Tools such as Zipkin [4], Jaeger [1], LTTng [2], and Trace Compass [3] provide access to such data, but interpreting it at scale remains difficult. Effective analysis requires deep expertise in OS internals, hardware scheduling, or file system behavior. Existing tools offer limited flexibility. Users must rely on rigid charts or write custom scripts using low-level APIs, a process that is error-prone, and inaccessible to non-experts.

Several challenges make trace analysis hard to scale. Traces are massive. Even short runs can produce hundreds of millions of events, which cannot fit in memory or in the context window of a large language model (LLM). They are also multidimensional. Events involve threads, processes, files, sockets, and CPUs interacting over time. Finally, reasoning often requires correlating events across long time intervals and interpreting raw numerical fields with little semantic guidance. Even simple questions like "Why was this file closed?" may demand deep manual effort. These limitations make trace analysis hard to democratize. To ground the problem we begin with a short operational scenario that reflects a common trace triage scene.

> **Motivating scenario.** A performance engineer investigates a recurring latency spike on CPU_1 during a load test. Warm up has passed and the spike repeats within a short window. The first triage step is to decide whether CPU_1 is dominated by one thread or by many. The engineer isolates the interval $[t_1, t_2]$ where the spike appears and asks:
> - Which thread executed most on CPU_1 between $t_1$ and $t_2$?
> - Within the same window, how many distinct threads ran on CPU_1 and did the top thread account for the majority of runtime?

These checks are representative of common and time-consuming tasks because they expose imbalance and migration patterns that decide the next diagnostic step. Here teams follow playbooks that first identify the top consumer on the affected CPU then measure diversity to decide whether to pursue one thread or investigate system wide pressure. Both checks are time consuming when done by hand since the analyst must sum per thread run time over the exact interval and account for migrations. Prior work reports 15–45 minutes per investigation [23]. This is one example of such investigations. Other incidents may favor analogous interval scoped checks such as the same pair on CPU_0 or the same procedure in a 10s window.

Building on this motivation, we propose **TAAF** (*Trace Abstraction and Analysis with Foundation models*), a three-layered framework for semantically meaningful and scalable trace analysis. Instead of

feeding raw data into an LLM, TAAF builds a compact structured abstraction. A temporal state system summarizes millions of events into structured, time-indexed transitions; a query-specific knowledge graph organizes entities and their relationships; and a large language model interprets the graph to answer natural-language questions. This modular design enables explainable reasoning, improves scalability, and reduces hallucinations in complex system scenarios.

Our research investigates how structured abstractions like state systems and knowledge graphs can bridge low-level event data and high-level questions. Specifically, we ask: (1) How can knowledge graphs model the temporal and multidimensional nature of traces? (2) Can indexed trace outputs improve LLM reasoning over raw logs? (3) How does combining KGs with LLMs affect the accuracy and explainability of trace analysis? These questions guide the design and evaluation of TAAF.

**We make the following contributions:**

- We present **TAAF**, a three-layered framework combining temporal abstraction, knowledge graphs, and LLM reasoning for semantic trace analysis.
- We propose a method to transform kernel traces into time-indexed graphs that encode entity interactions and temporal structure.
- We introduce **TraceQA-100**, a benchmark of 100 expert-authored trace-analysis questions grounded in real Linux traces, focusing on thread–CPU interactions and activity patterns related to scheduling behavior.
- We empirically show that graph-grounded reasoning improves answer accuracy (up to 31%) across GPT-4o, GPT-4.1 nano and o4-mini (reasoning), relative to raw or flattened inputs.

We evaluate our approach using real Linux traces by segmenting the trace into varying temporal locations and durations to assess scalability. Our experiments involve multiple LLMs and grounding strategies. TAAF consistently outperforms the baselines, demonstrating that structured grounding enhances both reasoning accuracy and robustness. We also anonymously release all code, data, and evaluation artifacts to support reproducibility and future research[1].

## 2 Related Work

We review prior work in four key areas relevant to TAAF: trace abstraction and stateful modeling, knowledge graph construction for system understanding, large language models for software reasoning, and integration strategies that bridge structured and neural techniques.

### 2.1 Trace Abstraction and Stateful Modeling

Execution traces are notoriously large and semantically dense, making raw analysis costly and inaccessible. Prior work has proposed several abstraction techniques to reduce trace size while preserving interpretability. Pirzadeh et al. [25] segment large traces into execution phases based on Gestalt principles and apply stratified sampling to enable scalable exploration. Feng et al. [11] introduce

_Sage_, a hierarchical abstraction framework that captures high-level behaviors from recurring low-level patterns. Hamou-Lhadj and Lethbridge propose summarization methods based on cue phrases and frequency analysis [12], as well as the _Utilityhood_ metric [13] to rank components by structural relevance. Cornelissen et al. [9] conduct a quantitative evaluation of multiple trace reduction techniques, such as sampling, stack depth limitation, and grouping. They examine the trade-offs of these methods in preserving behavioral fidelity.

Complementing these methods, the Trace Compass platform [14, 17, 29] provides a stateful abstraction via its _State System_, which maintains a time-indexed database of attribute-value intervals for efficient querying. This system constructs a persistent model of trace state using event replays and interval trees, enabling scalable extraction of system-level behaviors [31, 34]. While these techniques offer powerful indexing and summarization, they lack semantic structuring and are limited in supporting flexible, user-driven queries over multidimensional trace entities.

### 2.2 Knowledge Graphs for Execution Context Modeling

Knowledge graphs (KGs) offer a semantically rich abstraction to represent system entities and their relationships. In trace contexts, they enable modeling of complex interactions among processes, threads, files, and CPUs in a form that is interpretable and machine-readable. Liang et al. [21] categorize KGs into static, temporal, and multimodal types. Static KGs represent triples $(e_1, r, e_2)$ without time evolution; temporal KGs introduce evolving snapshots; and multimodal variants incorporate metrics or other data types such as images and videos. KG-based reasoning may be transductive, constrained to known entities, or inductive, supporting generalization. While KGs have been applied in domains like cybersecurity and software architecture, few systems construct them dynamically from execution traces. TAAF fills this gap by deriving time-indexed, per-query KGs from stateful trace outputs.

### 2.3 LLMs for Structured Software Reasoning

Large Language Models (LLMs) such as GPT-3 [8] and GPT-4 [6] have demonstrated strong capabilities in NLP tasks like translation [18], summarization [20], and question answering [28]. However, their application to trace analysis faces two major barriers. First, trace data volumes often exceed the token limits of LLMs, making end-to-end ingestion infeasible. Second, LLMs encode knowledge implicitly and opaquely, leading to difficulties in validating results or understanding reasoning steps [24]. This often manifests as _hallucinations_ which refers to plausible-sounding but incorrect answers [16]. Moreover, LLMs struggle with structured inputs such as graphs or tabular state histories, limiting their utility in domains requiring explicit causal reasoning.

### 2.4 LLM–KG Integration Strategies

To address LLM limitations, recent work explores hybrid approaches that integrate KGs into the language model pipeline. One dominant strategy is _retrieval-augmented generation (RAG)_, which grounds LLMs in factually accurate subgraphs retrieved at inference time [7, 10]. Another line of work embeds KG structure directly into model

---

[1]https://anonymous.4open.science/r/TAAF-LLM-KG-State-System--60EC/README.md

architectures using entity-aware embeddings [32, 35] or graph neural modules such as QA-GNN [33]. A third family of techniques focuses on parameter-efficient domain adaptation using adapters [30], low-rank updates [15], or instruction-tuned prompts [26]. These integrations improve factual accuracy and interpretability, but introduce latency and alignment challenges especially when operating on large, temporal traces [19, 27]. Recent efforts such as compressed history trees [17] aim to reduce memory and runtime costs of retrieval, motivating approaches like TAAF that combine high-performance indexing with structuring and LLM reasoning.

## 3 Methodology

### 3.1 Problem Setting and Motivation

Modern kernel execution traces log millions of timestamped events across diverse entities such as threads, processes, CPUs, files, and network interfaces, offering rich runtime insights. However, answering practical trace-analysis questions from such voluminous and low-level data presents three core challenges. First, analysts must reason across varying temporal granularities, from instantaneous states to fine-grained intervals such as 1-second windows and full-trace trends. Second, many queries involve multidimensional interactions, for instance, reasoning across scheduling behavior, I/O flows, and resource contention, which requires interpretation of complex inter-entity relationships. Third, trace data lacks high-level semantic annotations, while analyst queries range from simple factual lookups to comparative or causal questions, demanding both flexible and precise semantic understanding. We also wanted to see how far we could go in trace analysis with LLMs and knowledge graphs.

To address these challenges, we introduce the Trace Abstraction and Analysis with Foundation Models (TAAF) framework, which integrates scalable data abstraction, structured knowledge representation, and natural-language reasoning to enable flexible, accurate, and explainable trace analysis. Figure 1 illustrates the framework architecture overview.

### 3.2 Architectural Overview and Trace-to-Answer Pipeline

TAAF is a three-layered abstraction framework for structured trace analysis. Rather than analyzing raw kernel trace events directly, TAAF systematically transforms them into structured, interpretable representations suitable for structured reasoning and LLM-based answering.

Overall, the system defines a trace-to-answer transformation pipeline $\mathcal{T} \to \mathcal{S} \to \mathcal{G} \to \mathcal{A}$, where $\mathcal{T}$ is the raw trace, $\mathcal{S}$ is the State System, $\mathcal{G}$ is the knowledge graph, and $\mathcal{A}$ is the final answer. Each layer serves to constrain scope, preserve structure, and expose trace semantics to the model in a usable form.

- **Layer 1: From Trace to State.** The raw event stream $\mathcal{T}$ is transformed into a structured, time-indexed State System $\mathcal{S} = (Q, H)$ using a deterministic encoding function:

$$\mathcal{T} \xrightarrow{\phi_1} \mathcal{S}$$

This transformation enables scalable temporal reasoning by organizing system attributes and their state transitions as synchronized trees.

- **Layer 2: From State to Graph.** Given a natural-language query $q$, TAAF extracts relevant temporal intervals from $\mathcal{S}$ and constructs a semantic knowledge graph $\mathcal{G}_q$ that represents system interactions within the scope of the query:

$$\mathcal{S}, q \xrightarrow{\phi_2} \mathcal{G}_q$$

This graph is query-specific, minimal, and enriched with temporal and semantic metadata.

- **Layer 3: From Graph to Answer.** The graph $\mathcal{G}_q$ and an associated schema prompt $C_q$ are serialized and passed to a large language model, which produces a grounded natural-language answer $\mathcal{A}_q$:

$$\mathcal{G}_q, C_q, q \xrightarrow{\phi_3} \mathcal{A}_q$$

This LLM-based reasoning phase supports quantitative analysis, comparisons, and explanations grounded in trace semantics.

These intermediate representations are structured and semantically annotated; they capture system behavior and support reasoning, composability, and extensibility across trace types.

We revisit this layered design in our evaluation (Section 4) to assess its impact on scalability, semantic grounding, and answer quality.

We describe each layer in more details below.

### 3.3 Layer 1: From Raw Trace to State System Abstraction

A kernel trace $\mathcal{T}$ is defined as a temporally ordered set of fine-grained system events:

$$\mathcal{T} = \{(e_i, t_i, c_i) \mid e_i \in \mathcal{E}, \ t_i \in \mathbb{T}, \ c_i \in C\}$$

Each tuple consists of an event type $e_i$, a timestamp $t_i$, and a set of involved system components $c_i$ such as threads, CPUs, locks, or I/O devices. These events reflect low-level operating system behavior, including context switches, file operations, and synchronization events. While highly expressive, such traces are both voluminous and semantically sparse. Even a few seconds of system activity may yield tens of millions of events, rendering direct reasoning infeasible due to the combined temporal complexity, multidimensionality of interacting entities, and lack of structured semantics.

To address this, we define a formal transformation $\phi_1$ that converts a raw trace $\mathcal{T}$ into a time-indexed structured abstraction known as the *State System*. Originally proposed by Montplaisir et al. [23], the State System enables scalable querying and semantic enrichment by organizing the trace into two synchronized tree structures: an *attribute tree* and a *history tree*.

Formally, the State System is defined as $\mathcal{S} = (Q, H)$, where:

- $Q = \{q_1, q_2, \ldots, q_n\}$ is the set of stable integer identifiers (*quarks*) corresponding to hierarchical attribute paths (e.g., `/CPU/2/Thread/5130`),
- $H$ is the history tree, recording value transitions as:

$$H = \{\langle t_{\text{start}}, t_{\text{end}}, q, v \rangle \mid q \in Q, v \in \mathcal{V}\}$$
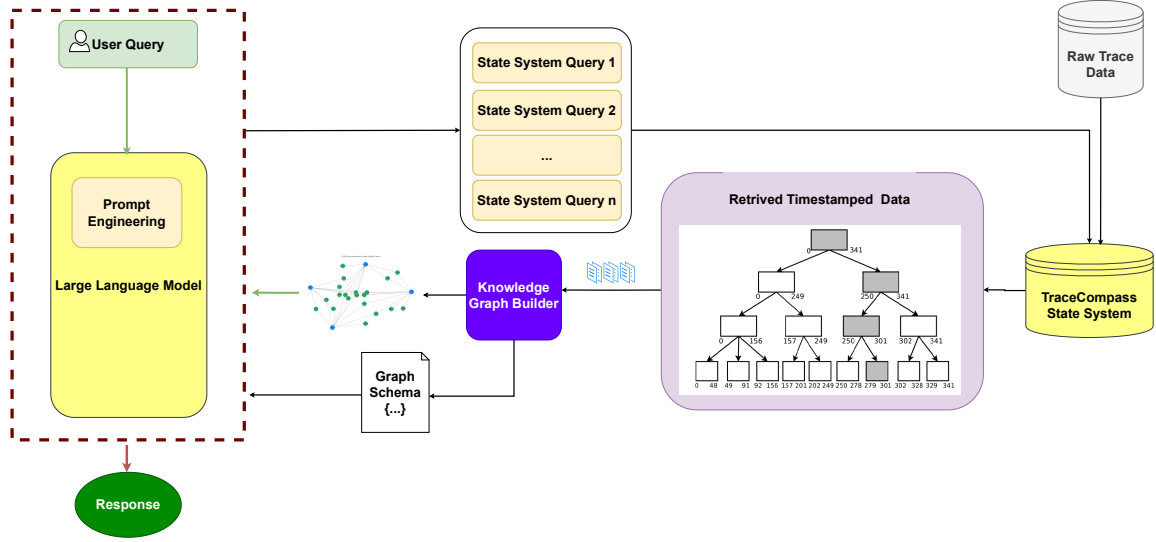
Figure 1: End-to-end architecture overview of the Trace Abstraction and Analysis Framework (TAAF).

Each entry indicates that attribute $q$ held value $v$ during the interval $[t_{start}, t_{end})$.

The transformation from raw events to this structure is governed by a predefined set of event-type-specific functions. For each event type $e \in \mathcal{E}$, we define a mapping:

$$\phi_1^{(e)} : \mathbb{T} \times C \to \mathcal{P}(H)$$

which encodes how an event of type $e$, given its timestamp and involved components, contributes to zero or more entries in $H$. The complete transformation $\phi_1$ is the union of these specialized rules:

$$\phi_1(\mathcal{T}) = \bigcup_{(e_i, t_i, c_i) \in \mathcal{T}} \phi_1^{(e_i)}(t_i, c_i)$$

This formulation captures the deterministic nature of event-to-state updates while preserving modularity across heterogeneous event types.

This dual-tree design supports efficient structured queries such as:

$$\text{query}(q, t) \to v \quad \text{and} \quad \text{query}(q, [t_1, t_2]) \to \{v_i\}$$

enabling localized reasoning over fine-grained execution intervals without requiring full trace scans.

By compactly encoding temporal dynamics in this form, the State System allows higher-level semantic queries to be executed in logarithmic time relative to the trace size [17, 23, 34]. This layer thus transforms $\mathcal{T}$ into a semantically enriched structure $\mathcal{S}$ that forms the foundation for downstream reasoning tasks in TAAF.

## 3.4 Layer 2: From State System to Query-Specific Knowledge Graph

Given the structured state representation $\mathcal{S} = (Q, H)$ produced by $\phi_1$, the second layer of TAAF constructs a dynamic, query-specific knowledge graph that grounds trace semantics in a structured and interpretable form. This transformation is handled by a function
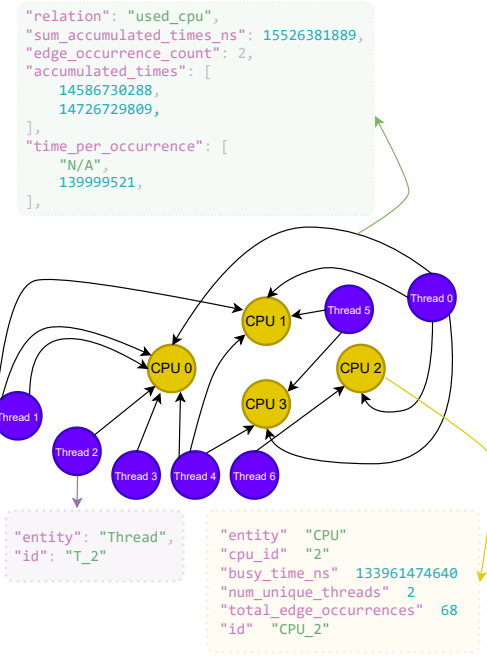


Figure 2: Example of a query-specific knowledge graph generated from a 1-second trace interval.

$\phi_2$, which takes the State System $\mathcal{S}$ and a natural language query $q$ as input, and produces a query-specific knowledge graph $\mathcal{G}_q$.

Formally, we define $\mathcal{G}_q$ a 4-tuple:

$$\mathcal{G}_q = (V_q, E_q, \mu_q, \tau_q)$$

where $V_q$ denotes the set of typed nodes representing system entities, and $E_q \subseteq V_q \times \mathcal{L} \times V_q$ denotes the set of labeled, directed edges that encode semantic relationships between those entities. The label

space $\mathcal{L}$ includes relation types such as `executes_on`, `reads_from`, and `holds_lock`. Each edge $e \in E_q$ is assigned a weight via the function $\mu_q : E_q \rightarrow \mathbb{R}$, typically representing quantitative values such as duration or frequency. Temporal scope is encoded through a mapping $\tau_q : V_q \cup E_q \rightarrow \mathbb{T}$, which associates each entity or interaction with a time window derived from $H$.

This knowledge graph is constructed on demand, scoped narrowly to the query $q$ and the relevant interval within the trace. For example, if $q$ asks, "How long did Thread 5130 run on CPU 2 between $t_1$ and $t_2$?", the system identifies the relevant quarks from $Q$, queries their associated intervals in $H$, and maps the result to a graph with two nodes (representing Thread 5130 and CPU 2) and a single edge labeled `executes_on`, weighted by total runtime (e.g., $\mu(e) = 5.25$ seconds) and scoped to the interval $\tau(e) = [t_1, t_2]$. Figure 2 shows an example KG generated from question querying a 1-second trace window.

By extracting only the entities and relations needed for $q$, $\mathcal{G}_q$ avoids the overgeneralization and interpretability issues of global or static graphs. It supports structured reasoning in a lightweight, interpretable, and semantically aligned form. Edges are directed to capture causal relationships and weighted with quantitative metrics from the State System. The graph is tailored to each query and expandable with additional metadata to enable focused, quantitative, and adaptable analysis grounded in domain knowledge.

This step completes the semantic grounding layer of TAAF, translating time-series trace data into a structured graph format suitable for large language model reasoning.

## 3.5 Layer 3: From Knowledge Graph to Answer via LLMs

The final layer of TAAF performs semantic reasoning over the structured representation $\mathcal{G}_q$ to generate a natural-language answer $\mathcal{A}_q$. This transformation is governed by the function $\phi_3$, which takes as input the knowledge graph $\mathcal{G}_q$, a natural-language query $q$, and a schema prompt $C_q$ that explicitly encodes the meaning of graph elements. The output is a grounded, human-readable answer:

$$\mathcal{A}_q = \phi_3(\mathcal{G}_q, C_q, q)$$

To maximize the factual accuracy and interpretability of LLM responses, TAAF follows a structured prompting strategy. The schema prompt $C_q$ formally defines the types of nodes and edges in $\mathcal{G}_q$, their attributes (e.g., execution time, file path), and any additional domain-specific hints. The graph $\mathcal{G}_q$ itself is serialized in a structured JSON format, preserving the graph topology and annotated metadata. Together, the components are passed as input to the large language model alongside the query $q$ in a templated format:

```
{ "schema": {...},
"graph": {...},
"user query": "Which thread executed most
on CPU 1 between t1 and t2?"
}
```

The model is thus explicitly grounded in a well-scoped semantic context, minimizing the need for open-ended inference. The result is the answer $\mathcal{A}_q$, which may include the direct outcome. Moreover, it can contain justifications or comparative insights, depending on the query type and model capabilities.

For example, given the above user query, the model may return:

> "Thread 9127 executed for 6.4 seconds on CPU 1 between $t_1$ and $t_2$, which is the highest among all threads during that interval."

This output demonstrates the core strength of the TAAF framework. By separating trace understanding into a structured abstraction layer and deferring reasoning to an LLM that operates over this structured input, we support robust, explainable trace analysis without relying on static rules or brittle heuristics.

Importantly, the LLM interface in TAAF is designed to be model-agnostic. In our implementation, we employ both general-purpose and reasoning-tuned foundation models (e.g., `o4-mini`), and observe that schema-conditioned inputs significantly improve factual precision. Furthermore, the prompt construction process is transparent and auditable, making it easier to debug reasoning failures or extend the system to new query types or domains.

By combining the structure with foundation model reasoning, this final layer completes the trace-to-answer pipeline and enables TAAF to flexibly support diverse analytic needs across performance diagnosis, behavioral explanation, and root cause investigation.

## 4 Empirical Evaluation

This section details the experimental design, benchmark construction, evaluation protocol, scoring metrics, research questions, and the results that collectively assess the accuracy, robustness, and explainability of TAAF.

### 4.1 Evaluation Framework: TRACEQA-100

Despite progress in using AI and LLM for trace analysis, the community lacks a public, ground-truth dataset for reasoning over kernel-level execution traces. Existing LLM benchmarks such as MMLU or Big-Bench do not address the unique challenges of trace analysis, including fine-grained time, multi-entity interactions, and numeric aggregation. To fill this gap, we introduce TRACEQA–100 a curated benchmark designed to exercise all three dimensions.

**Trace provenance and slicing.** We use LTTng traces of the SciMark 2.0 Java benchmark [5], which generates approximately 34 M kernel events. From each run, we extract 1s, 10s, and 100s slices around three canonical temporal locations, *start*, *mid*, and *end*, yielding a pool of trace segments that serve as the factual basis for question construction.

**Question authorship.** Questions were generated through a four-step, double-blind process. Two experts independently inspected trace slices and drafted questions requiring time-aware, multi-entity reasoning. Drafts were peer-reviewed and normalized into three answer formats: explanatory, multiple-choice, or true/false. A third expert, blind to the drafts, produced reference answers using hand-written Trace Compass scripts; disagreements were resolved through adjudication. Each question was then tagged as *single-hop* or *multi-hop* based on the reasoning scope.

The resulting set includes 100 questions, evenly split into 40 explanatory, 30 multiple-choice, and 30 true/false items, and a 50/50 split across hop scopes. Since each prompt is instantiated on nine distinct trace slices (3 temporal locations × 3 time-window lengths), the benchmark comprises $100 \times 9 = 900$ unique question–trace-segment pairs. Representative examples are shown in Table 1.

**Schema and reproducibility.** Each item in TRACEQA-100 includes the textual prompt, trace identifier, answer format, reference answer, and relevant metadata. All answer scripts are shipped and re-executed as part of our CI pipeline to verify trace-consistency. Detailed labeling guidelines are included to support extensions beyond SciMark or Linux.

To our knowledge, TRACEQA-100 is the first public benchmark that pairs large-scale kernel traces with expert-verified QA tasks. By targeting temporal, structural, and numeric reasoning, it offers a reusable reference point for future structured or neural approaches to trace understanding.

## 4.2 Evaluation Setup and Metrics

Each item in TRACEQA-100 is a standalone QA task, combining a natural-language prompt, a temporal slice of the SciMark 2.0 trace, and a single ground-truth answer. We evaluate each task under three trace–representation settings (Table 2): the *events-only* variant streams raw kernel events and serves as a qualitative baseline; the *Baseline* uses the numeric output of the temporal index; and **TAAF** augments this with a query-specific knowledge graph built by $\phi_2$.

To account for LLM stochasticity, we follow best practices [22] and sample each (question, representation, model) combination three times. With 100 benchmark questions, this yields 300 scored responses per configuration. Across three models, three time windows (1 s, 10 s, 100 s), and two structured inputs (Baseline, TAAF), the core grid totals 5,400 labeled outputs. It should be noted that events-only runs are included only for qualitative comparison, as such runs typically exceed the LLM context window and thus are not part of the scalable evaluation.

All runs execute on a single GPU node. End-to-end processing of the largest 100-second trace, including indexing, KG construction, and inference, completes in under 40 seconds, confirming feasibility for interactive use.

Each response is labeled on a 3-point scale: a score of **0** indicates the answer is incorrect or irrelevant; **0.5** denotes a partially correct response with minor errors or incomplete reasoning; and **1** signifies a fully correct answer that is consistent with the ground truth.

We report two metrics. **Accuracy** is the weighted average over scores:

$$\text{Accuracy} = \frac{(N_0 \cdot 0) + (N_{0.5} \cdot 0.5) + (N_1 \cdot 1)}{N_{\text{total}}} \times 100 \quad (1)$$

where $N_0, N_{0.5}, N_1$ denote the count of each label (out of $N_{\text{total}} = 300$).

**Consistency** captures the model's stability over repeated samples. We compute entropy:

$$E = - \sum_{i \in \{0, 0.5, 1\}} P_i \log_2 P_i \quad (2)$$

$$\text{Consistency} = \left(1 - \frac{E}{\log_2 3}\right) \times 100 \quad (3)$$

where $P_i$ is the empirical frequency of score $i$. Higher consistency indicates lower variance across trials.

*Model Configuration.* We evaluate three models via API: **GPT-4.1 nano**, a compact high-throughput model; **GPT-4o**, OpenAI's current flagship; and **o4-mini (Reasoning)**, optimized for high

reasoning accuracy. All models use the same prompt template and few-shot examples. Decoding temperature is fixed at 0.5, except in RQ7. For o4-mini, we activate `reasoning=high`.

*Cross-family probe.* To address generalizability beyond a single vendor, we also evaluate Google and Anthropic models. From Google, we use **Gemini 2.5 Pro** as the flagship reasoning model and **Gemini 2.5 Flash** as the fast, cost-efficient workhorse. From Anthropic, we use **Claude Opus 4.1** as the flagship reasoning model and **Claude Haiku 4.5** as the fast, cost-efficient model. We selected these pairs because each represents its vendor's top reasoning model alongside a smaller, speed- and value-oriented model, and both vendors publicly position them as such. All four models use the same prompt template and scoring protocol as the GPT models.

## 4.3 Research Questions

To guide our evaluation, we formulate seven research questions grouped into two phases. Phase 1 explores core performance dimensions of TAAF across query types, time windows, and LLM variants. Phase 2 isolates specific factors such as schema inclusion, temporal placement, and sampling behavior.

**Phase 1 – Full-grid evaluation**

- **RQ1:** How does TAAF's accuracy vary across different user query types, such as multiple-choice, true/false, and explanatory, and across graph question types, such as single-hop versus multi-hop?
- **RQ2:** What is the effect of time-window length on the quality of model responses?
- **RQ3:** To what extent does incorporating a knowledge graph improve accuracy compared to raw State System output?
- **RQ4:** How does performance vary across LLM backends including a focused cross-family probe?

**Phase 1** spans a 3×3 grid of evaluation settings formed by three LLMs and three time-window lengths (1s, 10s, 100s), all sampled from the midpoint of the trace. For each cell in the grid, we run both the TAAF pipeline with its query-specific knowledge graph and the baseline pipeline without graph grounding. Each configuration covers 100 benchmark questions sampled three times, yielding $9 \times 2 \times 300 = 5400$ labeled responses that inform RQ1 through RQ4.

**Phase 2 – Focused factors**

- **RQ5:** Does providing the LLM with the graph schema, such as node types and features, enhance the accuracy of its responses?
- **RQ6:** To what extent does the temporal location of the window (early, mid, late) affect accuracy?
- **RQ7:** How does the temperature parameter(sampling randomness) influence both accuracy and response consistency?

**Phase 2** holds the model (GPT-4o) and query grounding fixed, then varies key factors. RQ5 toggles schema support in the mid-trace 1s slice, contributing 300 responses. RQ6 compares temporal window placement, early versus late in the trace, at the 10s granularity, adding 600 responses. RQ7 examines the impact of sampling temperature by re-running the mid-trace 10s slice at four values (0.1, 0.3, 0.7, 0.9) for an additional 1,200 responses. In total, Phase 2 adds 2,100 outputs, bringing the complete dataset to 7,500 labeled responses used in the following analysis.

| Question type | hop scope | Example question | Temporal loc | Time window | Ref answer |
|---|---|---|---|---|---|
| Explanation | Single | What is the total accumulated CPU time for thread 5130 on CPU 0? | mid | 1 s | 5,247,200 ns |
| Explanation | Multi | Which CPU served the most number of distinct threads? | mid | 100 s | CPU_0 |
| Multiple choice | Single | Thread 2559 primarily uses: (A) CPU_0 (B) CPU_1 (C) CPU_2 (D) None | start | 10 s | B |
| Multiple choice | Multi | Thread 5236 primarily uses: (A) CPU_0 (B) CPU_1 (C) CPU_2 (D) CPU_3 | end | 10 s | C |
| True or False | Single | True or False? CPU 2 total busy time > 1e11 ns | mid | 10 s | False |
| True or False | Multi | True or False? The busiest CPU also runs the most unique threads | start | 10 s | True |

Table 1: Sample items from TraceQA-100 covering question formats, hop scopes, temporal locations, and windows.

| Variant | Input passed to LLM | Practical outcome |
|---|---|---|
| Events-only | Raw kernel events ($\geq 10^7$) | Exceeds context limits; produces incoherent answers. Retained only to illustrate the need for abstraction. |
| Baseline | State-System numeric values | Fits budget but remains low-level. |
| TAAF | Query-specific temporal knowledge graph | Adds explicit entities and relations; grounds the model while staying compact. |

Table 2: Trace representations supplied to the language model. Only the latter two fit within contemporary context windows.

The empirical results for these research questions are presented in the following.

## 4.4 Phase 1: Full Grid Experiments

**RQ1: How does TAAF's accuracy vary across different user query types, and across graph question types, such as single-hop versus multi-hop?**

We begin the analysis by probing how user and graph question styles influence TAAF's success. Our TraceQA-100 benchmark mixes three natural-language formats: (1) True/False, (2) Multiple Choice, and (3) Explanatory. Each format is asked in two graph-question settings: (1) single-hop questions that focus on one central entity, such as a CPU, and its local neighbourhood; (2) multi-hop questions that require reasoning about multiple distinct centres, for example two CPUs that share or migrate threads.

To better understand how user's query style and graph query structure interact, we grouped all model responses by question category and visualized the average accuracy in each bucket. Specifically, we aggregate scores across all three LLM models, then compute mean accuracy for every combination of *user question type* and *graph question type*. Figure 3 summarizes the result of this experiment. Adding KG in TAAF improves accuracy in every cell. The largest gain appears in True/False multi-hop queries, which rise by +30.12 percentage to reach 79.63%. Explanatory questions remain the most challenging across the board, but TAAF still boosts their accuracy by 23.89% in baseline setting. Looking at the absolute values, the highest observed accuracy is 90.99% for True/False single-hop queries under the TAAF setting. On the other hand, explanatory multi-hop queries are the lowest at 37.22% under the Baseline setting, reflecting the compound reasoning required. A full breakdown by model and configuration is provided in Table 3.

Table 3: Accuracy results for LLM models across different intervals and query types, with and without a KG.

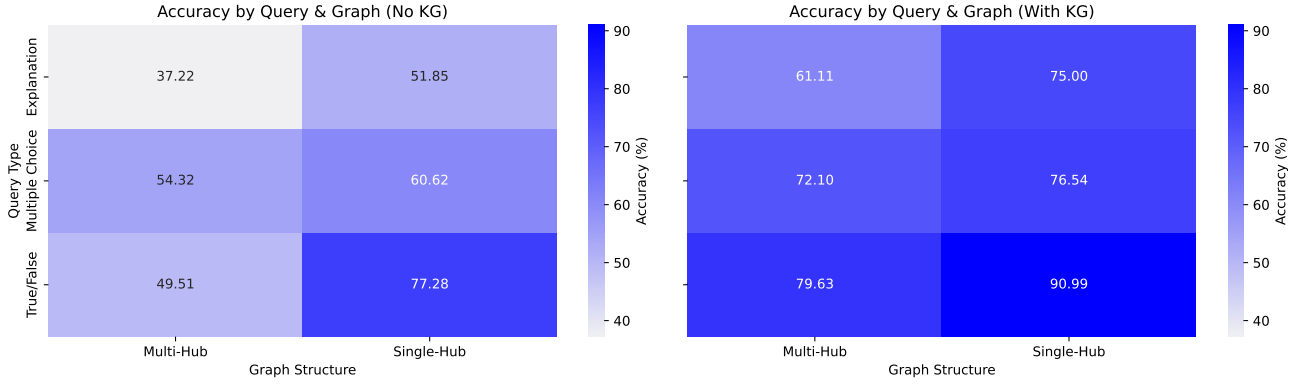| Model | Interval | UQT | GQT | No KG (Baseline) | | | | With KG (TAAF) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 0% | 0.5% | 1% | Acc% | 0% | 0.5% | 1% | Acc% |
| GPT 4.1 nano | 1 s | Expl. | S-hop | 45.00 | 13.33 | 41.67 | 48.33 | 33.33 | 8.33 | 58.33 | 62.50 |
| GPT 4.1 nano | 1 s | Expl. | M-hop | 61.67 | 20.00 | 18.33 | 28.33 | 46.67 | 20.00 | 33.33 | 43.33 |
| GPT 4.1 nano | 1 s | MC | S-hop | 35.56 | 17.78 | 46.67 | 55.56 | 11.11 | 5.56 | 83.33 | **86.11** |
| GPT 4.1 nano | 1 s | MC | M-hop | 40.00 | 11.11 | 48.89 | 54.44 | 26.67 | 4.44 | 68.89 | 71.11 |
| GPT 4.1 nano | 1 s | T/F | S-hop | 20.00 | 8.89 | 71.11 | 75.56 | 10.00 | 2.22 | 87.78 | **88.89** |
| GPT 4.1 nano | 1 s | T/F | M-hop | 41.11 | 21.11 | 37.78 | 48.33 | 16.67 | 12.22 | 71.11 | 77.22 |
| GPT 4.1 nano | 10 s | Expl. | S-hop | 46.67 | 16.67 | 36.67 | 45.00 | 43.33 | 5.00 | 51.67 | 54.17 |
| GPT 4.1 nano | 10 s | Expl. | M-hop | 56.67 | 21.67 | 21.67 | 32.50 | 46.67 | 13.33 | 40.00 | 46.67 |
| GPT 4.1 nano | 10 s | MC | S-hop | 50.00 | 13.33 | 36.67 | 43.33 | 29.17 | 5.83 | 65.00 | 67.92 |
| GPT 4.1 nano | 10 s | MC | M-hop | 55.67 | 13.00 | 31.33 | 37.83 | 37.14 | 4.29 | 58.57 | 60.71 |
| GPT 4.1 nano | 10 s | T/F | S-hop | 27.78 | 9.44 | 62.78 | 67.50 | 11.11 | 4.44 | 84.44 | 86.67 |
| GPT 4.1 nano | 10 s | T/F | M-hop | 55.56 | 16.67 | 27.78 | 36.11 | 29.17 | 9.17 | 61.67 | 66.25 |
| GPT 4.1 nano | 100 s | Expl. | S-hop | 49.03 | 15.97 | 34.99 | 42.98 | 45.10 | 7.84 | 47.06 | 50.00 |
| GPT 4.1 nano | 100 s | Expl. | M-hop | 59.15 | 21.13 | 19.72 | 30.29 | 49.30 | 19.72 | 30.99 | 40.85 |
| GPT 4.1 nano | 100 s | MC | S-hop | 43.40 | 19.81 | 36.79 | 46.70 | 26.61 | 6.42 | 66.97 | 70.18 |
| GPT 4.1 nano | 100 s | MC | M-hop | 48.50 | 13.50 | 38.00 | 44.75 | 33.58 | 6.72 | 59.70 | 62.06 |
| GPT 4.1 nano | 100 s | T/F | S-hop | 35.63 | 14.38 | 50.00 | 57.19 | 14.13 | 6.38 | 79.50 | 82.69 |
| GPT 4.1 nano | 100 s | T/F | M-hop | 55.10 | 17.65 | 27.25 | 36.07 | 26.13 | 10.92 | 62.95 | 67.41 |
| GPT 4o | 1 s | Expl. | S-hop | 30.00 | 8.00 | 62.00 | 66.00 | 18.00 | 2.67 | 79.33 | **80.67** |
| GPT 4o | 1 s | Expl. | M-hop | 46.00 | 12.00 | 42.00 | 48.00 | 29.33 | 9.33 | 61.33 | 66.00 |
| GPT 4o | 1 s | MC | S-hop | 26.00 | 6.00 | 68.00 | 71.00 | 7.00 | 1.33 | 91.67 | **92.33** |
| GPT 4o | 1 s | MC | M-hop | 33.33 | 8.33 | 58.33 | 62.50 | 13.33 | 3.33 | 83.33 | 85.00 |
| GPT 4o | 1 s | T/F | S-hop | 13.33 | 5.00 | 81.67 | **84.17** | 4.00 | 0.67 | 95.33 | **95.67** |
| GPT 4o | 1 s | T/F | M-hop | 32.00 | 14.67 | 53.33 | 60.67 | 12.67 | 3.33 | 84.00 | 85.67 |
| GPT 4o | 10 s | Expl. | S-hop | 34.48 | 11.49 | 54.02 | 59.76 | 33.33 | 3.45 | 63.22 | 64.94 |
| GPT 4o | 10 s | Expl. | M-hop | 50.00 | 12.50 | 37.50 | 43.75 | 38.46 | 11.54 | 50.00 | 55.77 |
| GPT 4o | 10 s | MC | S-hop | 34.55 | 12.73 | 52.73 | 58.09 | 13.64 | 2.73 | 83.64 | 85.00 |
| GPT 4o | 10 s | MC | M-hop | 46.43 | 11.90 | 41.67 | 47.62 | 22.32 | 3.57 | 74.11 | 75.89 |
| GPT 4o | 10 s | T/F | S-hop | 18.27 | 6.73 | 75.00 | 78.37 | 6.25 | 2.50 | 91.25 | 92.50 |
| GPT 4o | 10 s | T/F | M-hop | 41.67 | 16.67 | 41.67 | 50.00 | 17.44 | 10.47 | 72.09 | 76.32 |
| GPT 4o | 100 s | Expl. | S-hop | 39.39 | 12.12 | 48.48 | 54.55 | 29.54 | 8.41 | 62.05 | 66.26 |
| GPT 4o | 100 s | Expl. | M-hop | 55.17 | 15.52 | 29.31 | 36.07 | 35.79 | 14.74 | 49.47 | 56.84 |
| GPT 4o | 100 s | MC | S-hop | 39.22 | 14.71 | 46.08 | 53.44 | 11.76 | 3.92 | 84.31 | 86.27 |
| GPT 4o | 100 s | MC | M-hop | 50.88 | 14.04 | 35.09 | 42.11 | 23.81 | 5.95 | 70.24 | 72.22 |
| GPT 4o | 100 s | T/F | S-hop | 26.92 | 9.62 | 63.46 | 68.27 | 7.21 | 2.88 | 89.90 | **91.35** |
| GPT 4o | 100 s | T/F | M-hop | 47.83 | 17.39 | 34.78 | 43.48 | 15.85 | 6.71 | 77.44 | 80.80 |
| GPT o4-mini | 1 s | Expl. | S-hop | 17.78 | 15.56 | 66.67 | 74.44 | 4.44 | 5.00 | 90.56 | **93.06** |
| GPT o4-mini | 1 s | Expl. | M-hop | 35.24 | 22.86 | 41.90 | 53.33 | 8.57 | 14.29 | 77.14 | 84.29 |
| GPT o4-mini | 1 s | MC | S-hop | 17.02 | 8.51 | 74.47 | 78.72 | 3.19 | 1.06 | 95.74 | 96.27 |
| GPT o4-mini | 1 s | MC | M-hop | 21.21 | 14.14 | 64.65 | 71.72 | 5.41 | 2.70 | 91.89 | 93.24 |
| GPT o4-mini | 1 s | T/F | S-hop | 8.33 | 7.22 | 84.44 | **88.06** | 2.22 | 1.11 | 96.67 | **97.22** |
| GPT o4-mini | 1 s | T/F | M-hop | 21.75 | 13.75 | 64.50 | 71.88 | 4.37 | 2.46 | 93.17 | 94.35 |
| GPT o4-mini | 10 s | Expl. | S-hop | 9.52 | 13.10 | 77.38 | **83.93** | 5.95 | 5.95 | 88.10 | **91.07** |
| GPT o4-mini | 10 s | Expl. | M-hop | 27.40 | 19.86 | 52.74 | 62.67 | 9.59 | 13.70 | 76.71 | 83.56 |
| GPT o4-mini | 10 s | MC | S-hop | 13.68 | 10.26 | 76.06 | 81.19 | 3.08 | 3.08 | 93.85 | 95.38 |
| GPT o4-mini | 10 s | MC | M-hop | 18.10 | 11.43 | 70.48 | 76.19 | 5.24 | 2.86 | 91.90 | 93.33 |
| GPT o4-mini | 10 s | T/F | S-hop | 3.33 | 5.00 | 91.67 | **94.17** | 0.56 | 0.56 | 98.89 | **99.17** |
| GPT o4-mini | 10 s | T/F | M-hop | 14.29 | 12.14 | 73.57 | 79.64 | 1.78 | 3.56 | 94.67 | 96.44 |
| GPT o4-mini | 100 s | Expl. | S-hop | 12.50 | 12.50 | 75.00 | 81.25 | 6.15 | 6.15 | 87.69 | **90.76** |
| GPT o4-mini | 100 s | Expl. | M-hop | 25.64 | 17.95 | 56.41 | 65.38 | 7.69 | 12.82 | 79.49 | 85.90 |
| GPT o4-mini | 100 s | MC | S-hop | 17.39 | 11.59 | 71.01 | 76.81 | 4.96 | 2.48 | 92.56 | 94.05 |
| GPT o4-mini | 100 s | MC | M-hop | 22.22 | 11.11 | 66.67 | 72.22 | 5.88 | 2.94 | 91.18 | 92.65 |
| GPT o4-mini | 100 s | T/F | S-hop | 6.58 | 7.89 | 85.53 | **89.47** | 1.54 | 1.92 | 96.54 | **97.50** |
| GPT o4-mini | 100 s | T/F | M-hop | 18.25 | 13.14 | 68.61 | 75.18 | 2.55 | 3.63 | 93.83 | 95.64 |

**Figure 3: Accuracy by query type and hop count. Left: baseline; right: TAAF. White text shows accuracy gain (+p.p.).**

**RQ1 Key observations.**
- True/False single-hop questions achieve the highest accuracy overall (90.99%).
- Multi-hop queries are harder in general but benefit more from graph grounding.
- Explanatory prompts remain the most difficult, especially on multi-hop questions, despite strong gains.

### RQ2: What is the effect of time-interval length on the performance of TAAF's answers?

Next, we examine how the time interval length affects performance. To isolate the effect of time-interval length, we contrasted the three time windows 1s, 10s, and 100s across all three LLM settings with and without the knowledge graph (TAAF vs. Baseline). Figure 4 plots the resulting accuracy curves. Solid markers denote the TAAF variant and dashed markers trace the baseline. The results (TAAF) reveal a clear, model-agnostic trend: accuracy falls as the window widens, but the rate of decline depends on model strength.



**Figure 4: Accuracy at for each model and interval**

**Numerical drop-off (TAAF Setting).**
- **o4-mini** 95.5 % → 92.0 % → 90.17 % (−5.33 pp overall)
- **GPT-4o** 79.83 % → 79.33 % → 76.83 % (−3.0 pp)

- **GPT-4.1 nano** 60.5 % → 50.5 % → 51.33 % (−9.17 pp)

**Numerical drop-off (baseline setting).**

- **o4-mini** 73.83 % → 65.0 % → 59.0 % (−14.83 pp overall)
- **GPT-4o** 58.17 % → 53.5 % → 47.67 % (−10.5 pp)
- **GPT-4.1 nano** 49.33 % → 37.5 % → 42.67 % (−6.66 pp)

Two forces compete here. Longer windows supply richer relational evidence, yet they also enlarge the State System output or the graph. Stronger models like o4-mini and GPT-4o exploit the extra context with only modest degradation, whereas the lightweight GPT-4.1 nano loses grounding signal more rapidly. The baseline decline is steeper because it must process ever-larger volumes of raw State System data, while TAAF mostly needs to update node and edge features as the graph grows. We also observe two small upticks in the blue GPT-4.1 nano curves; one with KG on (50.5 → 51.33%) and one in the baseline (37.5 → 42.67%) when moving from the 10s to the 100s window. While the exact cause remains speculative and needs further investigation, one plausible explanation is that the wider temporal window enables the model to access more cumulative evidence, which benefits some of queries that depend on aggregated behavior. Additionally, GPT-4.1 nano may struggle with the intermediate 10s input which is too complex to reason over fully but not large enough to highlight dominant patterns.

**RQ2 Key observations.**
- *Shorter windows are the safest.* All models peak on the 1-second slice; noise accumulation outweighs extra context beyond 10s.
- *o4-mini remains robust.* Even at 100s it retains 90%+ accuracy, suggesting that careful reasoning objectives can handle graph expansion.
- *Baseline deteriorates faster.* In baseline settings, accuracy drops roughly two-to-three times as quickly as in TAAF, underscoring how raw State-System output becomes unwieldy as the interval grows.

### RQ3: To what extent does the incorporation of a Knowledge Graph improve the accuracy of the TAAF when answering trace-related queries?

Third, we quantify how much explicit structure within KG helps over purely numeric State System data. Figure 5 quantifies the absolute accuracy lift obtained by adding the knowledge-graph context in TAAF compared with the baseline. Across all nine model–interval combinations the KG consistently helps, ranging from a modest + 8.67% for GPT-4.1 nano on the 100-second slice up to a sizeable + 31.17% for o4-mini on the 100-second slice. The mean gain over the entire grid is 21.5%, confirming that graph grounding is the most influential factor in our pipeline.

Figure 6 explains how the accuracy rises in more details. For every model–interval pair the left stacked bar is the baseline and the right stacked bar is the TAAF variant. The KG chiefly converts partial or wrong answers into fully correct ones. The fraction of 0.5-scores also falls, indicating the model is not merely guessing and it reaches a decisive, correct conclusion more often.
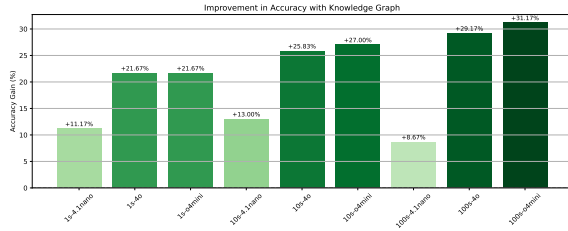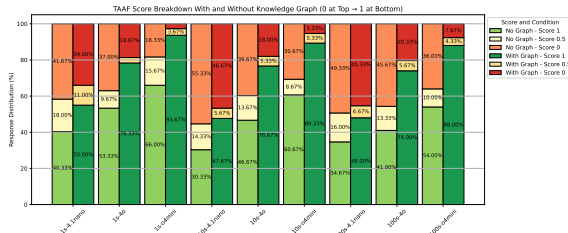
**RQ4: How does the accuracy and quality of TAAF responses vary across different LLM backends (including a focused cross-family probe)?**

Finally within Phase 1, we ask whether the observed gains depend on the underlying LLM. Because the same experiment grid spans all three LLMs, we can read the impact of interval length directly from Figure 4.

Overall, o4-mini with TAAF settings tops the chart at every interval, peaking at 95.5% on the 1-second slice and still clearing 90.17% on the 100-second slice. GPT-4o trails by roughly 15%, while the compact GPT-4.1 nano stays in the 50–60% band even with graph support. With baseline settings the spread narrows, confirming that stronger reasoning models capitalize on explicit structure more fully than weaker ones.
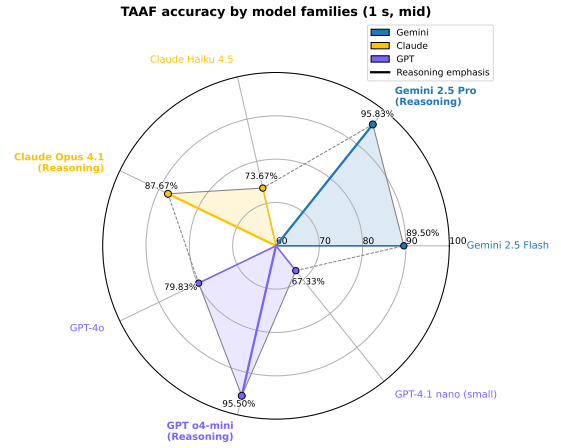


Figure 5: Accuracy gain from KG in TAAF .



Figure 6: Score breakdown (0/0.5/1) with vs. without KG.



Figure 7: TAAF accuracy across model families.

> **RQ3 Key observations.**
> - *Longer intervals benefit more.* The lift grows steadily from 1s to 100s because multi-hop relations become denser and harder to reconstruct from raw numerics alone.
> - *Stronger models still improve.* Even the larger models (GPT-4o) gain roughly + 22% at 1s and + 29% at 100s, showing that parametric knowledge does not replace explicit structure.
> - *Reasoning-tuned models exploit the KG best.* o4-mini shows the highest incremental benefit, suggesting that lightweight reasoning objectives align well with graph-grounded setting.

*Cross-family probe.* To check that TAAF's accuracy and findings do not depend on a single vendor family, we also ran a focused probe on two additional model families under the same prompt and scoring protocol. We selected one flagship and one reasoning-tuned model per family to mirror the GPT tiers. The probe uses the **TAAF** setting on a **1 s mid-trace** slice to bound runtime and isolate the model family factor. Accuracies are Gemini 2.5 Pro (Reasoning) **95.83%**, Gemini 2.5 Flash **89.50%**, Claude Opus 4.1 (Reasoning) **87.67%**, Claude Haiku 4.5 **73.67%**. The pattern mirrors the GPT tiering. Reasoning-oriented flagships top their family, while compact fast models trade a modest drop for speed. Fig. 7 better presents the numbers and indicates that the accuracy comes from the TAAF graph grounding rather than vendor specifics. This provides evidence that TAAF carries across model families and supports general use.

**RQ4 Key observations.**
- *o4-mini with TAAF settings leads.* Its accuracy exceeds 90% at every interval. This shows the value of explicit reasoning objectives when graph context is available.
- *Gap widens with structure.* The performance delta between GPT-4o and GPT-4.1 nano grows once the KG is supplied, suggesting larger models extract richer relational cues.
- *Vendor-agnostic pattern.* The cross-family probe mirrors the GPT tiering. Gaps stay within a single-digit band when tiers are matched, which points to TAAF's graph grounding as the main source of the lift.

## 4.5  Phase 2: Focused Factors

**RQ5: Does providing the LLM with the graph schema (node types and features) enhance the accuracy of its responses?**

Having established the KG's value, we test whether giving the model a lightweight JSON schema improves grounding further. Before running the full suite of experiments, we first asked whether the LLM actually *uses* the structural hints contained in the schema (node types and edge features). We therefore repeated the GPT-4o mid-trace, 1-second experiment with the schema strings removed from the prompt but kept the graph triples intact. Table 4 shows the outcome.
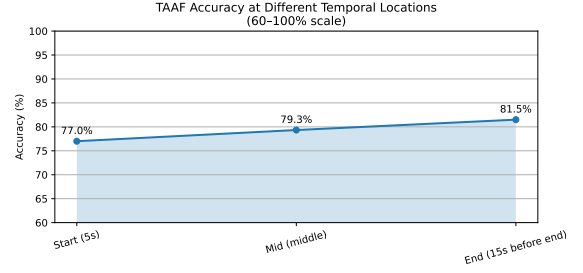
**Table 4: Impact of passing the graph schema.**

| Prompt Variant | 0 | 0.5 | 1 | TAAF Accuracy |
|---|---|---|---|---|
| No schema | 20.0 % | 16.7 % | 63.3 % | 71.7 % |
| With schema | 18.7 % | 3.0 % | 78.3 % | **79.8 %** ↑ |

**RQ5 Key observation.** Passing the lightweight schema yields an extra + 8.1% absolute accuracy; therefore all remaining experiments keep the schema enabled.

**RQ6: How does the temporal location, that is the position of the time window within the trace, affect TAAF accuracy?**

Timing matters: we now check if accuracy shifts when the time window is taken early, mid-way, or late in the trace. For RQ6 we fixed GPT-4o in the TAAF setting and applied a 10-second analysis window at three temporal positions. We tested three placements of that window in each trace. The **start** location covers from 5 seconds into the trace until 15 seconds in. The **mid** location spans from the midpoint of each trace to 10 seconds beyond that midpoint. The **end** location covers from 15 seconds before the trace ends until 5 seconds before the end. Figure 8 shows the details of this experiment.

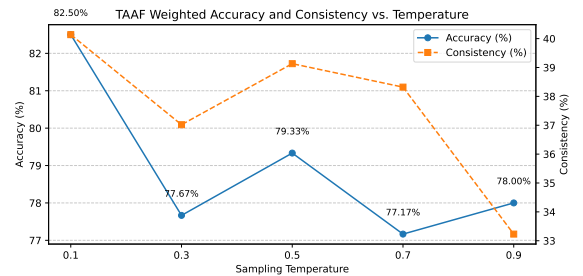**Figure 8: TAAF accuracy at different temporal locations**

As shown in Figure 8 accuracy climbs from 77.0% at the start, to 79.3% in the middle, and reaches 81.5% near the end. Later windows often contain longer, steadier bursts of activity, making scheduling patterns easier to spot. Early windows cover the warm-up phase of the workload and include more transient behaviour. The overall spread is small, about 4%, so users can place the window almost anywhere in the trace and still receive reliable answers.

**RQ6 Key observations.**
- Accuracy slightly improves as the window moves from the start toward the end, plus 4.5% overall.
- Late segments feature steadier execution phases, giving clearer patterns for the model to follow.
- The small gap shows that TAAF remains robust to window placement, so users need not fine-tune the timestamp.

**RQ7: How does the sampling temperature influence both accuracy and consistency?** Finally, we explore the sensitivity of both accuracy and stability to the LLM's sampling temperature.

Like in RQ6, we kept GPT-4o in TAAF setting and the 10 second mid-trace window, then varied the sampling temperature that controls randomness in the LLM output. Five values were tested: 0.1, 0.3, 0.5, 0.7, and 0.9. For each value we computed accuracy and the consistency metric defined in Section 4.2. Figure 9 plots the results.

**Figure 9: Accuracy and consistency as a function of sampling temperature.**

Based on the results, accuracy is highest at the lowest temperature, 82.5% at 0.1, then drops to 77.67% at 0.3, recovers slightly at 0.5, and declines again at 0.7. Consistency follows a similar but

smoother curve: best at 0.1, stable near 0.5, and lowest at 0.9. Low temperatures reduce hallucination but risk deterministic errors, while very high temperatures introduce too much noise.

> **RQ7 Key observations.**
> - Best trade-off appears at 0.1 to 0.3, high accuracy with the top consistency scores.
> - Temperatures above 0.7 lower accuracy and consistency, showing the model drifts when sampling is too random.
> - A moderate setting near 0.5 balances diversity and reliability, which is why earlier experiments used this as the default.

## 4.6 Discussion

The two-phase evaluation shows that TAAF raises accuracy in every setting while remaining stable across models, window sizes, and query styles. Graph grounding supplies explicit structure, short windows limit noise, and a lightweight schema boosts grounding at little extra cost. These results justify the design choices made in Section 3. The State System gives fast, time-indexed access, the query-specific knowledge graph keeps the context small and relevant, and the LLM converts that context into actionable insight.

We have also highlighted where performance still falls, such as explanatory multi-hop questions and very long windows. These gaps suggest clear avenues for future optimization and help frame the limitations that follow.

With the empirical evidence in place, we now examine possible threats to the validity of our study.

## 5 Threats to Validity

We outline potential limitations that may affect the interpretation of our results and suggest directions for future improvement.

### 5.1 Construct Validity

**Question diversity.** Our benchmark includes 100 hand-crafted questions from SciMark 2.0 traces. While diverse in format and scope, it may miss edge cases like rapid context switching or deeply nested interrupt chains. Future work can extend coverage via automated generation.

**Scoring resolution.** Our 3-level rubric {0, 0.5, 1} captures broad correctness but cannot distinguish minor errors (e.g., off-by-one times). More fine-grained metrics (e.g., BLEU, numeric tolerance) could improve sensitivity.

### 5.2 Internal Validity

**Annotation bias.** Author-labeling may bias scores, especially on borderline outputs. Using multiple independent raters and reporting inter-annotator agreement would increase objectivity.

**Prompt length limits.** Some large traces or graphs exceed LLM context limits. We truncate edge attributes as needed, which may reduce answer quality. Future solutions include hierarchical prompting or sliding windows.

### 5.3 External Validity

**Trace generalizability.** All questions are based on SciMark 2.0 under Linux. Other kernels or workloads could alter behavior. Broader evaluation using diverse traces (e.g., eBPF, HPC, mobile) is needed.

**Model stability.** LLMs accessed via APIs may evolve silently, affecting reproducibility. Caching model versions and outputs can help mitigate this threat.

**Scope limits.** TAAF is most effective when queries are scoped to a window. Global, unbounded queries can produce KGs too large for contemporary contexts or collapse temporal order if summarized aggressively. For such cases, global pre-aggregation or approximate sampling is more appropriate.

**LLM-as-reasoner limits.** We observed even with a correct KG context, models sometimes fail on multi-step aggregation. For example, the query "which CPU has the lowest difference between its highest and lowest accumulated times among threads" requires multiple aggregation steps. We observed that models may skip steps or make small arithmetic slips.

### 5.4 Conclusion Validity

**Sampling variability.** Each query was sampled three times. While standard practice, higher replication or paired tests would improve statistical confidence.

**Baseline scope.** We compare TAAF only to the raw State System. Intermediate designs, such as summaries or static graphs, could further contextualize TAAF's gains.

## 6 Conclusion and Future Work

This paper introduced **TAAF**, a framework that transforms low-level kernel traces into query-relevant answers by integrating temporal indexing, knowledge graphs, and large language models. TAAF addresses key challenges in trace analysis, including volume, temporal complexity, semantic ambiguity, and the need for domain expertise, by combining abstractions with structured context for reasoning.

Across 7,800 responses spanning three GPT models, three time windows, and six query types, graph-based grounding improved weighted accuracy by an average of 21.5%, reaching up to 95.5% on short intervals. These gains confirm that injecting explicit structure enhances both factuality and stability, even for high-performing models.

Evaluation also revealed that longer windows and multi-hop queries remain challenging, though reasoning-oriented models like o4-mini remain robust. TAAF's performance is stable across early, mid, and late trace slices, and sampling temperatures between 0.1 and 0.3 offer a strong trade-off between accuracy and diversity.

Looking forward, promising directions include integrating temporal knowledge graphs for continuous reasoning and embedding TAAF into autonomous agent workflows. Additional work will apply the framework to production-scale systems to surface real bugs and scale beyond current context limits through retrieval-augmented or hierarchical prompting. A complementary direction is a systematic latency review that profiles each pipeline stage from end to end while exploring optimizations such as graph partitioning and sampling, streaming or incremental KG construction, batching, and early-exit heuristics.

# References

[1] [n. d.]. Jaeger: Open Source, End-to-End Distributed Tracing. https://www.jaegertracing.io/. Accessed: 2025-05-22.

[2] [n. d.]. LTTng: Linux Trace Toolkit Next Generation. https://lttng.org/. Accessed: 2025-05-22.

[3] [n. d.]. TraceCompass: A Trace Visualization Tool. https://www.eclipse.org/tracecompass/. Accessed: 2025-05-22.

[4] [n. d.]. Zipkin. https://zipkin.io/. Accessed: 2025-05-22.

[5] 2020. SciMark 2.0 Traces. https://zenodo.org/record/437170. Accessed: 2025-05-01.

[6] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774* (2023). https://arxiv.org/abs/2303.08774

[7] Jinheon Baek, Chingyu Chung, Sung Ju Kim, and Seung-won Hwang. 2023. KAPING: Knowledge-Aware Prompting Improves Factuality in Abstractive Summarization. *arXiv preprint arXiv:2305.13952* (2023). https://arxiv.org/abs/2305.13952

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

[9] Bas Cornelissen and Leon Moonen. 2008. On large execution traces and trace abstraction techniques. *Software Engineering Research Group, Delft* (2008).

[10] Yihao Feng, Zeyu Chen, Yue Zhang, et al. 2023. Knowledge Solver: Teaching LLMs to Search for Domain Knowledge from Knowledge Graphs. *arXiv preprint arXiv:2305.16581* (2023). https://arxiv.org/abs/2305.16581

[11] Yang Feng, Kaj Dreef, James A Jones, and Arie van Deursen. 2018. Hierarchical abstraction of execution traces for program comprehension. In *Proceedings of the 26th Conference on Program Comprehension*. 86–96.

[12] Abdelwahab Hamou-Lhadj. 2005. The concept of trace summarization. *Proceedings of PCODA* (2005).

[13] Abdelwahab Hamou-Lhadj and Timothy Lethbridge. 2006. Summarizing the content of large traces to facilitate the understanding of the behaviour of a software system. In *14th IEEE International Conference on Program Comprehension (ICPC'06)*. IEEE, 181–190.

[14] Weijiang Hong, Zhenbang Chen, Yide Du, and Ji Wang. 2021. Trace abstraction-based verification for uninterpreted programs. In *Formal Methods: 24th International Symposium, FM 2021, Virtual Event, November 20–26, 2021, Proceedings 24*. Springer, 545–562.

[15] Edward Hu, Yelong Shen, Phil Wallis, et al. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *Proceedings of the International Conference on Learning Representations*. https://openreview.net/forum?id=nZeVKeeFYf9

[16] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *Comput. Surveys* 55, 12 (2023), 1–38.

[17] Herve Mbikayi Kabamba, Matthew Khouzam, and Michel Dagenais. 2023. Node Compass: Multilevel Tracing and Debugging of Request Executions in JavaScript-Based Web-Servers. *arXiv preprint arXiv:2401.08595* (2023). https://arxiv.org/abs/2401.08595

[18] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7871–7880.

[19] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kanal, Sebastian Riedel, Edward Grefenstette, Mike Lewis, et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, Vol. 33.

[20] 9459–9474.

[20] Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2024. Pre-trained language models for text generation: A survey. *Comput. Surveys* 56, 9 (2024), 1–39.

[21] Ke Liang, Lingyuan Meng, Meng Liu, Yue Liu, Wenxuan Tu, Siwei Wang, Sihang Zhou, Xinwang Liu, Fuchun Sun, and Kunlun He. 2024. A survey of knowledge graph reasoning on graph types: Static, dynamic, and multi-modal. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2024).

[22] Gili Lior, Eliya Habba, Shahar Levy, Avi Caciularu, and Gabriel Stanovsky. 2025. ReliableEval: A Recipe for Stochastic LLM Evaluation via Method of Moments. *arXiv preprint arXiv:2505.22169* (2025). https://arxiv.org/abs/2505.22169

[23] Alexandre Montplaisir, Naser Ezzati-Jivan, Florian Wininger, and Michel Dagenais. 2013. Efficient model to query and visualize the system states extracted from trace data. In *International Conference on Runtime Verification*. Springer, 219–234.

[24] Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering* (2024).

[25] Heidar Pirzadeh Tabari. 2012. *Trace abstraction framework and techniques*. Ph. D. Dissertation. Concordia University.

[26] Parth Sen, Anirudh Ramesh, Pritish Mukerji, et al. 2023. Knowledge-augmented Graph Machine Learning for Drug Discovery: A Survey from Precision to Interpretability. *arXiv preprint arXiv:2307.00559* (2023). https://arxiv.org/abs/2307.00559

[27] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. REPLUG: Retrieval-Augmented Black-Box Language Models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 8371–8384. https://aclanthology.org/2024.naacl-long.463

[28] Dan Su, Yan Xu, Genta Indra Winata, Peng Xu, Hyeondey Kim, Zihan Liu, and Pascale Fung. 2019. Generalizing question answering system with pre-trained language model fine-tuning. In *Proceedings of the 2nd workshop on machine reading for question answering*. 203–211.

[29] The Eclipse Foundation. 2025. *Generic State System*. The Eclipse Foundation. https://archive.eclipse.org/tracecompass/doc/stable/org.eclipse.tracecompass.doc.dev/Generic-State-System.html

[30] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. 2021. K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. 1405–1418. https://aclanthology.org/2021.findings-acl.121

[31] Florian Wininger, Naser Ezzati-Jivan, and Michel R Dagenais. 2017. A declarative framework for stateful analysis of execution traces. *Software Quality Journal* 25 (2017), 201–229.

[32] Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. LUKE: Deep Contextualized Entity Representations with Entity-aware Self-attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*. 6442–6454.

[33] Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, et al. 2022. QA-GNN: Reasoning with Language Models and Knowledge Graphs for Question Answering. In *NAACL*. 582–609. https://aclanthology.org/2022.naacl-main.43

[34] Siwei Zhang. 2023. *Trace Visualization of Distributed and Centralized Applications*. Master's thesis. KTH Royal Institute of Technology, Stockholm, Sweden.

[35] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 1441–1451.