# Safe, Untrusted, "Proof-Carrying" AI Agents: toward the agentic lakehouse

Jacopo Tagliabue
*Bauplan Labs*
NYC, USA
jacopo.tagliabue@bauplanlabs.com

Ciro Greco
*Bauplan Labs*
NYC, USA
ciro.greco@bauplanlabs.com

*Abstract*—**Data lakehouses run sensitive workloads, where AI-driven automation raises concerns about trust, correctness, and governance. We argue that API-first, programmable lakehouses provide the right abstractions for *safe-by-design*, agentic workflows. Using Bauplan as a case study, we show how data branching and declarative environments extend naturally to agents, enabling reproducibility and observability while reducing the attack surface. We present a proof-of-concept in which agents repair data pipelines using correctness checks inspired by proof-carrying code. Our prototype demonstrates that untrusted AI agents can operate safely on production data and outlines a path toward a fully agentic lakehouse.**

*Index Terms*—**AI agents, lakehouse, data pipelines, versioning**

## I. Introduction

The data lakehouse is the *de facto* cloud architecture for analytics and Artificial Intelligence (AI) workloads [2], [3], thanks to storage-compute decoupling, multi-language support, and unified table semantics. As reasoning and tool usage in Large Language Models (LLMs) improve [4], autonomous decisions ("AI agents") are both supported by, and targeted at, cloud infrastructure: to what extent can agents manage the data lifecycle in a lakehouse?

*Prima facie*, the question appears both too hard and too broad. On one hand, lakehouses are distributed systems built for the collaboration of human teams on sensitive production data, not point-wise tasks immediately suitable for end-to-end automation. On the other, it is unclear how to prioritize agentic use cases across such heterogeneous platforms. *This* paper is a preliminary answer to these challenges: we detail lakehouse abstractions suitable for automation, and operationalize a prototype for an important use case: *repairing data pipelines*.

Pipelines are a compelling case study for three reasons: first, they cover a large portion of lakehouse workloads, measured both by developer time [5] and overall compute [6]. Second, data engineers spend a significant amount of their time fixing them [7], [8]. Finally, repairing pipelines is a canary test for agent penetration in high-stakes non-trivial scenarios, which are often hard even for expert humans [9], [10]. We summarize our contributions as follows:

1) we introduce abstractions to model the data life-cycle in a programmable lakehouse [11], i.e. building and executing cloud pipelines entirely through *code*. We argue that traditional systems resist automation mostly because of heterogeneous interfaces and complex access patterns, while code is the *lingua franca* suitable for agents, cloud systems, and human supervisors;

2) we review common objections to automating high-stakes workloads in light of the proposed abstractions: in particular, we argue that our model promotes trustworthiness and correctness both in data and code artifacts;

3) we release working code[1], showing a proof of concept for self-repairing pipelines using `Bauplan` as a lakehouse and an agentic loop. Starting from this prototype, we conclude by outlining practical next steps for a full agentic lakehouse.

The paper is organized as follows. After reviewing agent-friendly abstractions (Section II), we address key safety objections for high-stakes scenarios (Section III). Once safety is established, we describe a ReAct [12] loop built on these abstractions (Section IV). We put forward our working prototype as a feasibility demonstration of safe-by-design data agents, not as a full-fledged experimental benchmark.

We believe that sharing working code is of great value to the community, especially in times of quickly shifting mental models. However, it is important to remember that our fundamental insights – programmability and safety – can be replicated independently of the chosen APIs. For these reasons, we believe our paper to be valuable to a wide range of practitioners: on one hand, those looking for a new mental map of this uncharted territory; on the other, those looking to be inspired by tinkering with existing implementations and inspecting systems working at scale.

## II. A programmable lakehouse

In a *programmable* lakehouse, the *entire* data life-cycle – data, user and infrastructure management, pipeline and query execution, runtime observability – is exposed through code abstractions: server-side APIs, SDK methods, CLI shortcuts. In the rest of the paper, `Bauplan` snippets will be used as sample implementation, but the platform's composable nature makes it easy to replicate these functionalities in different

---

[1]Open source code is available at https://github.com/BauplanLabs/the-agentic-lakehouse.

architectures.[2] We break down the pipeline life-cycle into two major components: pipeline definition and pipeline execution.

## A. Pipeline definition

A pipeline is a DAG of transformations. A DAG starts from source tables, which are progressively cleaned, augmented, aggregated through the transformation code (expressed in SQL or Python). A successful execution produces intermediate and final data assets, which are then consumed by downstream systems [15]. Fig. 2 shows a pipeline (hence, P) used as a recurring example throughout the paper. Taking two tables from the NYC taxi dataset [16], P defines two new tables (A and B), based on two transformations (1 and 2). The Bauplan implementation for P is as follows:[3]

Listing 1. P as the Python file *p.py*

```
@bauplan.model(materialization="REPLACE", name="A")
@bauplan.python("3.10", pip={"pandas": "2.0"})
def join_and_filter(
    trips=bauplan.Model("taxi_trips"),
    zones=bauplan.Model("taxi_zones")
):
    # some transformation here...
    return trips.join(zones).do_something()

@bauplan.model(materialization="REPLACE", name="B")
@bauplan.python("3.11", pip={"pandas": "1.5.3"})
def clean_and_transform(
    data=bauplan.Model("join_and_filter")
):
    return data.do_something()
```

Two important design choices are worth highlighting in connection to our safety discussion (Section III):

- **Function-as-a-service (FaaS) abstractions**: business logic is expressed in the body of plain vanilla functions with the signature $Table(s) \rightarrow Table$. DAGs are functions chained through naming convention. These abstractions naturally map to a serverless runtime, which can execute the requested computation efficiently [11];
- **declarative I/O and infrastructure**: functions are fully isolated (e.g. two functions, two versions of `pandas`) and their Python environment is specified declaratively [17]. Reading tables and writing artifacts back to the lake is also fully declarative: users specify the needed inputs and desired output, the platform performs the corresponding physical operations.

## B. Pipeline execution

A human (or an agent) with the proper access can execute `p.py` by simply installing the `bauplan` package, and running
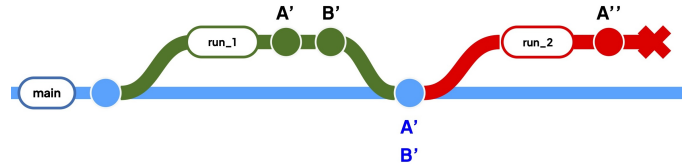


Fig. 1. **Transactional pipelines**: a successful execution of P, followed by one that failed after A materialization. Runs happen on data branches: input data is from sources in `main`, but writes are sandboxed so that materialized tables hit `main` *atomically* only on success (no half-written pipeline).

it from the terminal, without any additional steps – no Docker, no Terraform, no JDBC clients:[4]

Listing 2. Bauplan CLI

```
$ pip install bauplan
$ bauplan run --project_dir P_folder
```

While simplicity is a virtue, for our present concerns we focus on the *transactional* nature of the `run` API, which is obtained by re-purposing Git concepts to table evolution, and by managing data and compute as a logical whole even if physically decoupled. Consider now the two runs depicted in Fig. 1: `run 1` (successful) and `run 2` (failed). The branches and merges in the picture are a graphical representation of the underlying "Git-for-Data" abstractions [18]: if every change to the lake corresponds to a *commit*, a *branch* is the HEAD of a sequence of commits, and a *merge* atomically combines commits from two branches. If the `main` branch represents production, *merge* operations between on-branch writes and *main* mimic software best practices, and provides a natural hook for testing and reviews.

When `run 1` starts, the execution is automatically moved to a copy-on-write branch: source tables will have the same data as production, but the tables materialized by the run will be written to this branch first and then merged to `main` with an atomic operation, promoting A' and B' as the new tables for downstream consumers. The importance of coupling runs with a *branch-then-merge* pattern becomes evident with `run 2`, which failed before a new version of B was materialized: no merge happens for `run 2`, so no dirty read can occur in `main` – downstream systems will still read a consistent pair, A' and B', not A'' and B'. In other words, *branches* allow sandboxed writes starting from production reads, i.e. working with production data without the risk of destroying production.

Git-for-Data abstractions at the execution level therefore complement the functional abstractions at the definition level, solving three important use cases through simple APIs:

- **reproducibility**: runs are immutably, deterministically identified through a pointer to the starting *commit* and a copy of the code;

---

[2]We refer the interested reader to existing papers, in particular [11], [13], [14].

[3]Snippets are simplified in the interest of space: full code is available in the open source repository.

[4]To get a first-person perspective on the developer experience, the reader is invited to pause and watch a recorded run before continuing: https://www.loom.com/share/99ac0d5b5f944fc9aeef132bfaea0881

- **transactionality**: the *branch-then-merge* pattern allows transactional pipelines, i.e. doing MVCC in a "deconstructed database" [19], [20];
- **reversibility**: writes are immutable but never final, as we can always revert to a previous state through the relevant commit.

### C. Code as the universal interface

Agents are a combination of reasoning (provided by LLMs) and tool usage. In a lakehouse, tools should allow *observability* of present and past runs and *reproducibility* of existing pipelines: if a lakehouse offers these capabilities in typed, documented methods, exposing them for agentic usage is as trivial as wrapping those methods in MCP routes. In other words, from the point of view of the relevant abstractions – functions, decorators, data branches –, a programmable lakehouse *is already* an agentic lakehouse.

In this sense, while serving overlapping use cases, the design distance between a programmable lakehouse and traditional systems could not be greater. Practitioners today have to context-switch between one-off Spark clusters, development notebooks, SQL editors, cloud warehouses and more [5]; *code*, on the other hand, is a universally understood interface – easy for agents to master, immediate for the cloud to run *and* safe for humans to properly supervise.

### III. THE SAFETY CHECKLIST

Making agents able to "code the lakehouse" through purpose-built abstractions is a necessary, but not sufficient condition to repair pipelines in production. We also need to make sure that the proposed abstractions address *trust and correctness* concerns over malicious or (simply wrong) usage (Table I).

**Trust in data**: can agents access data they are not supposed to? No. I/O is always mediated by the platform: agents have no access to the physical data layer (S3), so reads and writes happen always in platform space, not user space. More generally, since *all* operations are API operations, RBAC over API keys provides fine-grained permissions with a minimal attack surface.

**Trust in code**: can agents run untrusted packages or access malicious web resources? No. Functions are run as independent processes isolated from their host and other functions, with *no internet access* (since I/O is in platform space, users do not even access S3!). The declarative syntax makes (dis)allowing packages as trivial as checking a decorator against a whitelist. Once again, a concise surface comes in handy, as there is only one entry point for dependency management.

**Correctness in data**: can agents damage production data? No. First, incomplete pipelines will not affect downstream systems (Section II-B); second, removing merge-to-main permissions will make a human review necessary to reach production; third, humans can *always* revert tables using past commits.

**Correctness in code**: can agents push to production silent bugs? No (as in, no more than humans can). The key insight

| Concern | Mode | Abstraction |
|---|---|---|
| Trust | Data | Declarative I/O |
| Trust | Code | FaaS runtime |
| Correctness | Data | Transactional runs |
| Correctness | Code | Verify-then-merge |

is that even untrusted code can be useful provided a hard-to-fake correctness test. The aptly titled [1] defines an interesting protocol: a "consumer" specifies safety conditions upfront, a "producer" provides evidence that its work satisfies them. If the consumer is satisfied, the work is now "trusted". We can imagine a similar protocol, where the "consumer" is the lakehouse owner (i.e. a data engineer), and the "producer" is the agent repairing the pipeline on her behalf. The evidence will revolve around the pipeline output (A and B in P) meeting certain criteria. In the prototype below, the owner came up with a verifier, i.e. a function $Branch \rightarrow bool$ to allow an agent branch to be merged. Unlike the original protocol, data verifiers are less about formal properties and more about the business context for the generated tables: the same principles apply though, here strengthened by the "pull-request" flow that Git-for-Data enables. Importantly, verifiers use the same APIs as the agents, leveraging once again the unified access to data and compute, with no semantic or infrastructure drift between lakehouse clients.

### IV. A PROOF OF CONCEPT

`run 2` failed (Fig. 1): can an agent repair it? Now that we know it is safe to do so, we combine what we learned so far into an agentic loop and draw some preliminary conclusions from the prototype.

### A. The self-repairing setup

Our setup is as follows:

- Bauplan as the programmable lakehouse;
- the Bauplan MCP[5], exposing as tools the lakehouse APIs: tools can be used for *observability* (e.g. get failed jobs and their logs), *data exploration* (query tables, check types), *execution* (create branches, start a run)
- `smolagents`[6] is the ReAct framework, handling the loop, tool calls and logs: `smolagents` performs reasoning in Python, making it effective at pipeline reasoning when compared to JSON-based tool calling [21], [22]
- LLM inference provided by OpenAI, Anthropic and TogetherAI through a configurable LiteLLM[7] interface;
- a verifier function $Branch \rightarrow bool$, which is the "proof-checking" step before merging into main.

---

[5] https://github.com/BauplanLabs/bauplan-mcp-server
[6] https://huggingface.co/docs/smolagents/index
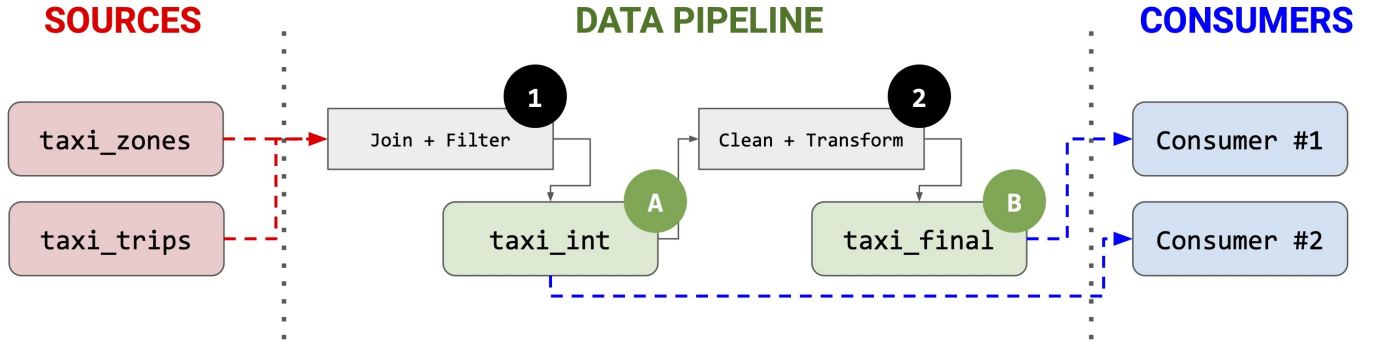[7] https://docs.litellm.ai/

Fig. 2. **Sample pipeline**: upstream from the pipeline, two source tables containing taxi trips and location data, downstream, multiple consumers. The pipeline itself is a two-node DAG, with compute steps in *gray* – 1 and 2 –, and tables in *green* – A and B.
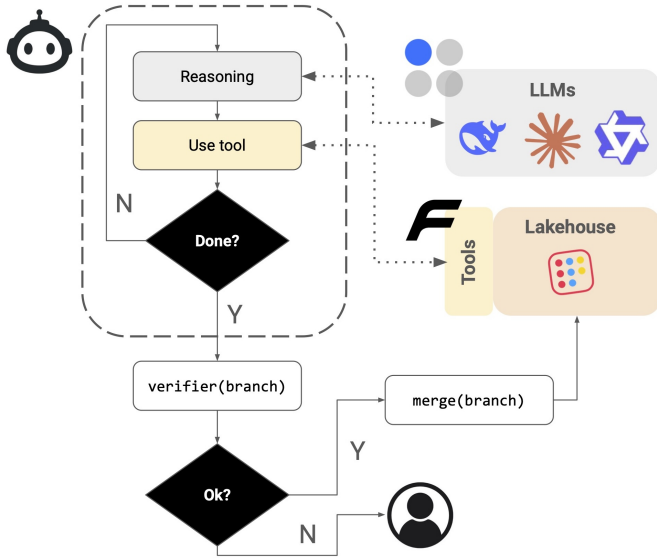


Fig. 3. **Safe, untrusted lakehouse agents**: the agent leverages LLMs and tools to repair a data pipeline. When an answer is produced, a deterministic check in the "outer loop" verifies that it is safe to merge to production.

### B. Running an experiment

Fig. 3 illustrates the high-level flow of an experiment: the agentic loop (reasoning and tool usage), the underlying programmable lakehouse (entirely accessible through APIs, wrapped by the MCP server), the proof-checking step at the end. In the example scenario, the script first launches a faulty pipeline to create the conditions for agentic repair: following prior reports [10] and industry experience, we simulate a package mismatch around the release of `NumPy 2.0` that caused crashes in containers using `pandas 2.0`.

A sample run[8] using a model such as *Sonnet 4.5* hits exactly all the abstractions discussed above: retrieving logs, querying the state of the lake, using declarative code to specify infrastructure changes, creating debug branches from production data, running code safely. Because this paper focuses on map-

[8]https://www.loom.com/share/fc79b52601074a5ba06e2f0272be3c62

ping abstractions for the agentic lakehouse and demonstrating feasibility, a thorough exploration of this experimental space is beyond scope. However, a few considerations are in order:

- as a testament to the complexity of the task, frontier model performance vary greatly in success rate, token usage and number of tool calls; as system builders downstream from LLMs, the crucial takeaway is that even when models failed (e.g., *GPT-5-mini*), the lakehouse exhibited no disruption or unsafe behavior;
- industry-leading traditional stacks – such as Snowflake with dbt – do not support agentic repair, even if they both have MCP servers and serve overlapping use cases. MCPs are a necessary but not sufficient condition for automation;
- because switching models is a single configuration change, we can easily imagine a budget-constrained scenario in which model selection is step-dependent, or a time-constrained situation in which data branches support concurrency control at scale for models in parallel.

## V. CONCLUSION AND FUTURE WORK

In a moment when most infrastructure agents are geared toward specific tasks [23], programmable lakehouses have the ambition to support agentic reasoning across the full life-cycle of data. To the best of our knowledge, we addressed for the first time the open-ended challenge of repairing cloud pipelines. We argued that a programmable lakehouse is already agentic, and that declarative DAGs and Git-like data management are ideally suited to support safe-by-design agentic usage.

To move beyond the hype, it is necessary to build novel systems and share working implementations. However, we also recognize that the infrastructure underlying agentic abstractions may need to evolve accordingly. Although massive parallelism is outside the scope of our prototype, it is arguably the primary challenge for OLAP systems in the age of agentic data exploration [24]. We look forward to continuing to share with the community our research journey toward a fully agentic lakehouse.

## REFERENCES

[1] G. C. Necula and P. Lee, "Safe, untrusted agents using proof-carrying code," in *Mobile Agents and Security*. Berlin, Heidelberg: Springer-Verlag, 1998, p. 61–91.

[2] D. Mazumdar, J. Hughes, and J. Onofre, "The data lakehouse: Data warehousing and more," 2023. [Online]. Available: https://arxiv.org/abs/2310.08697

[3] M. A. Zaharia, A. Ghodsi, R. Xin, and M. Armbrust, "Lakehouse: A new generation of open platforms that unify data warehousing and advanced analytics," in *Conference on Innovative Data Systems Research*, 2021.

[4] Z. Shen, "Llm with tools: A survey," *arXiv preprint arXiv:2409.18807*, 2024.

[5] C. G. Tapan Srivastava, Jacopo Tagliabue, "Eudoxia: a faas scheduling simulator for the composable lakehouse," *Proceedings of Workshops at the 51th International Conference on Very Large Data Bases*, 2025.

[6] A. van Renen, D. Horn, P. Pfeil, K. E. Vaidya, W. Dong, M. Narayanaswamy, Z. Liu, G. Saxena, A. Kipf, and T. Kraska, "Why tpc is not enough: An analysis of the amazon redshift fleet," in *VLDB 2024*, 2024.

[7] Data World, "Burned-out data engineers are calling for dataops," 2021. [Online]. Available: https://data.world/reports-and-tools/data-engineering-survey-2021-burned-out-data-engineers-are-calling-for-dataops/

[8] Z. Wang, T.-H. P. Chen, H. Zhang, and S. Wang, "An empirical study on the challenges that developers encounter when developing apache spark applications," *Journal of Systems and Software*, vol. 194, p. 111488, 2022.

[9] J. Yasmin, J. Wang, Y. Tian, and B. Adams, "An empirical study of developers' challenges in implementing workflows as code: A case study on apache airflow," *ArXiv*, vol. abs/2406.00180, 2024. [Online]. Available: https://api.semanticscholar.org/CorpusID:270213226

[10] H. Foidl, V. Golendukhina, R. Ramler, and M. Felderer, "Data pipeline quality: Influencing factors, root causes of data-related issues, and processing problem areas for developers," *Journal of Systems and Software*, vol. 207, p. 111855, 2024. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0164121223002509

[11] J. Tagliabue, T. Caraza-Harter, and C. Greco, "Bauplan: Zero-copy, scale-up faas for data pipelines," in *Proceedings of the 10th International Workshop on Serverless Computing*, ser. WoSC10 '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 31–36. [Online]. Available: https://doi.org/10.1145/3702634.3702955

[12] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao, "React: Synergizing reasoning and acting in language models," 2023. [Online]. Available: https://arxiv.org/abs/2210.03629

[13] J. Tagliabue, R. Curtin, and C. Greco, " FaaS and Furious: abstractions and differential caching for efficient data pre-processing ," in *2024 IEEE International Conference on Big Data (BigData)*. Los Alamitos, CA, USA: IEEE Computer Society, Dec. 2024, pp. 3562–3567. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/BigData62323.2024.10825377

[14] J. Tagliabue and C. Greco, "Reproducible data science over data lakes: replayable data pipelines with bauplan and nessie," in *Proceedings of the Eighth Workshop on Data Management for End-to-End Machine Learning*, ser. DEEM '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 67–71. [Online]. Available: https://doi.org/10.1145/3650203.3663335

[15] S. Salami, "Hub star modeling 2.0 for medallion architecture," 2025. [Online]. Available: https://arxiv.org/abs/2504.08788

[16] NYC.Gov, "Tlc trip record data," 2025. [Online]. Available: https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page

[17] J. Tagliabue, H. Bowne-Anderson, V. Tuulos, S. Goyal, R. Cledat, and D. Berg, "Reasonable scale machine learning with open-source metaflow," 2023. [Online]. Available: https://arxiv.org/abs/2303.11761

[18] J. Tagliabue, C. Greco, and L. Bigon, "Building a serverless data lakehouse from spare parts," *ArXiv*, vol. abs/2308.05368, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:260775634

[19] A. Cerone, G. Bernardi, and A. Gotsman, "A Framework for Transactional Consistency Models with Atomic Visibility," in *26th International Conference on Concurrency Theory (CONCUR 2015)*, ser. Leibniz International Proceedings in Informatics (LIPIcs), L. Aceto and D. de Frutos Escrig, Eds., vol. 42. Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2015, pp. 58–71.

[20] Jacopo Tagliabue et al., "Under Review," 2025.

[21] K. Yang, J. Liu, J. Wu, C. Yang, Y. R. Fung, S. Li, Z. Huang, X. Cao, X. Wang, Y. Wang, H. Ji, and C. Zhai, "If llm is the wizard, then code is the wand: A survey on how code empowers large language models to serve as intelligent agents," 2024. [Online]. Available: https://arxiv.org/abs/2401.00812

[22] X. Wang, Y. Chen, L. Yuan, Y. Zhang, Y. Li, H. Peng, and H. Ji, "Executable code actions elicit better llm agents," 2024. [Online]. Available: https://arxiv.org/abs/2402.01030

[23] Y. Gu, Y. Xiong, J. Mace, Y. Jiang, Y. Hu, B. Kasikci, and P. Cheng, "Argos: Agentic time-series anomaly detection with autonomous rule generation via large language models," 2025. [Online]. Available: https://arxiv.org/abs/2501.14170

[24] S. Liu, S. Ponnapalli, S. Shankar, S. Zeighami, A. Zhu, S. Agarwal, R. Chen, S. Suwito, S. Yuan, I. Stoica, M. Zaharia, A. Cheung, N. Crooks, J. E. Gonzalez, and A. G. Parameswaran, "Supporting our ai overlords: Redesigning data systems to be agent-first," 2025. [Online]. Available: https://arxiv.org/abs/2509.00997