TOWARDS FAITHFUL AGENTIC XAI

Anonymous authors

000

001

003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

024

025 026 027

028 029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Explainable AI (XAI) is essential for helping users interpret model behavior and proactively identify potential faults. Recently, Agentic XAI systems that integrate Large Language Models (LLMs) have emerged to make explanations more accessible for non-expert users through natural language. However, a critical limitation of the existing systems is their failure to address explanation faithfulness. This is problematic because many XAI methods are often unfaithful for complex models, and LLMs can amplify this incorrect information, ultimately misleading users. To address this limitation, we propose Faithful Agentic XAI (FAX), a framework that actively enhances explanation faithfulness. FAX introduces a systematic verification process where an LLM agent cross-checks claims against inherently faithful tools. This process filters out unreliable or contradictory evidence and leads to more faithful explanations. For evaluation, we propose CRAFTER-XAI-Bench, a benchmark framework built on an open-world reinforcement learning environment. The benchmark features complex models with diverse goals and challenging test scenarios, enabling a rigorous assessment of explanation faithfulness under realistic conditions. Experiments demonstrate that FAX significantly improves the faithfulness of explanations, marking a crucial step towards faithful and trustworthy Agentic XAI.

1 Introduction

Explainable AI (XAI) has emerged as a crucial field for demystifying black-box models, providing methods to understand their internal decision-making processes. Diverse XAI methods have been introduced to provide diverse information about the model, as described in Figure 1. However, interpreting the explanations often requires expert-level knowledge of machine learning and XAI, creating a significant barrier for non-expert users. To address this, the paradigm of Agentic XAI has been introduced (Slack et al., 2023; He et al., 2025), which employs a Large Language Model (LLM) to select suitable XAI methods and interpret the explanations in natural language.

However, a critical flaw underlies current Agentic XAI systems: an implicit assumption that the underlying XAI tools are consistently faithful. While this assumption may hold in simple, tabular settings, it breaks down for the complex models and dynamic environments seen in practice, where the unfaithfulness of XAI methods is a known and severe issue (Adebayo et al., 2018). An agent that naively trusts and rephrases these unreliable explanations can generate fluent, plausible, yet fundamentally incorrect explanations. This problem is further amplified by the inherent tendency of LLMs to hallucinate, potentially weaving flawed data into a dangerously convincing narrative.

In this work, we address this critical gap by proposing Faithful Agentic XAI (FAX), an agentic workflow designed to enhance explanation faithfulness. Instead of passively translating tool outputs, our agent employs a systematic verification process. It performs an explicit verification of claims by scrutinizing initial claims and cross-referencing them against evidence from multiple, inherently faithful tools. This iterative process filters out unreliable or contradictory results and allows the agent to proactively seek additional evidence, ultimately constructing a more robust and trustworthy explanation. Figure 2 illustrates this motivation and our approach.

To rigorously evaluate such a system, existing benchmarks are fundamentally inadequate. The faithfulness problem is often latent in simplistic tabular datasets; to properly test for it, we require a setting where XAI tools are genuinely challenged. We introduce CRAFTER-XAI-Bench, a scalable evaluation framework built upon an open-world Reinforcement Learning (RL) environment. This framework includes challenging scenarios, agents with diverse behaviors, and a suite of automated metrics, including a novel simulation-based metric to quantify faithfulness. By replacing subjective

		XAI Method Category				
_		Feature Importance	Counter- factual	Feature Influence	Surrogate Model	
nformation Type	Why	~	~	×	~	
Lion 1	Why not	×	~	×	X	
ırma	What if	×	×	~	~	
Info	How to be that	×	~	~	×	

Figure 1: Different XAI methods provide different information. Information categories are adopted from XAIQuestionBank (Liao et al., 2020).

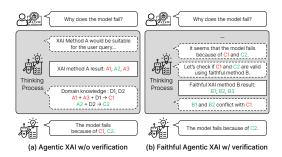


Figure 2: (a) Agentic XAI use XAI methods suitable for answering user query, and generate natural language response. (b) FAX verifies claims in response with inherently faithful XAI methods.

human studies with an LLM-as-a-judge approach, we enable scalable and reproducible assessment of Agentic XAI systems in complex domains.

To summarize our main contributions:

- We propose FAX, a novel agentic workflow that enhances explanation faithfulness by explicitly verifying claims, filtering unreliable claims, and proactively gathering evidence to construct a faithful explanation.
- We introduce a scalable evaluation framework for Agentic XAI, featuring a dynamic RL environment and a suite of automated metrics, including a simulation-based faithfulness metric, to facilitate rigorous testing.

2 Related Work

2.1 EXPLAINABLE AI

Classical methods Post-hoc XAI methods include four broad families: (i) feature attribution/saliency that highlights input regions or features with high contribution (Simonyan et al., 2014); (ii) surrogate models that approximate a local/global decision rule (e.g., rules or linear models) (Ribeiro et al., 2018; 2016); (iii) example-based explanations such as prototypes and counterfactuals that reason via representative or minimally edited examples (Chen et al., 2019; Wachter et al., 2018); and (iv) concept-based explanations that align internal representations with human-interpretable concepts (Kim et al., 2018; Yuksekgonul et al., 2023). Each family exposes a different facet of model behavior; consequently, a single method rarely satisfies diverse user intents.

Collection of explanations Since a single XAI method only reveals a limited aspect of a model's behavior, as illustrated in Figure 1, frameworks like Dijk et al. (2023); Yang et al. (2022); Arya et al. (2019) provide a collection of explanations in one place. However, identifying which method best answers a user's question and how to interpret its output still requires nontrivial XAI/ML expertise. In practice, users face a *selection and interpretation burden*: they must map their intent to a suitable method and often combine multiple views.

Interactive XAI To lower the barrier for non-experts, recent works have focused on generating natural language explanations that verbalize XAI outputs (Zytek et al., 2024; Castelnovo et al., 2024). Conversational assistants were suggested to explain the model's reasoning to users (Zhang et al., 2025b), and the benefits of text-based explanations over classical methods were confirmed via human study (Lakkaraju et al., 2022; Mindlin et al., 2024). Building on this, *Agentic XAI* systems have emerged, which use LLMs to select appropriate XAI tools based on a user's query (Slack et al., 2023; He et al., 2025).

However, these pioneering agentic systems have two critical limitations. First, they have primarily been tested on simpler models in static, tabular data settings. Second, and more crucially, they implicitly assume the underlying XAI tools are consistently faithful. This assumption often breaks down in complex and dynamic environments, where the unfaithfulness of XAI methods is a known

and severe issue (Adebayo et al., 2018). An agent that naively trusts and translates unreliable tool outputs can produce fluent, plausible, yet fundamentally incorrect explanations. He et al. (2025) have also warned that LLMs may amplify users' misunderstandings. We address this critical gap by focusing on enhancing explanation faithfulness within a challenging, dynamic environment.

2.2 LLM AGENT AND AGENTIC WORKFLOW

Recent work frames LLMs as *agents* that plan, act, and reflect while invoking external tools. ReAct interleaves reasoning traces with environment-facing actions to update plans and handle exceptions (Yao et al., 2022), while Toolformer demonstrates that LMs can *self-learn* when and how to call APIs and integrate their outputs (Schick et al., 2023). Building on these foundations, agentic extensions of LLMs now emphasize structured workflows that support multi-step reasoning, memory, and adaptive decision-making. For instance, the Model Context Protocol (MCP) provides a standardized interface for connecting LLMs with external services and tools, enabling modular extensibility. Also, recent works emphasize that the proper design of workflows is essential for flexible and reliable orchestration of agent behaviors (Zhang et al., 2025a). These developments underscore that the design of robust agentic workflows is central to realizing LLMs as proactive agents capable of simulation, decision-making, and long-horizon interaction.

2.3 LLM-as-a-Judge for scalable evaluation of natural language generation

LLM judges have emerged as a practical, scalable proxy for costly human studies, especially for evaluating the quality of generated text. MT-Bench/Chatbot Arena demonstrated that strong LLM judges can achieve high agreement with human preferences, while also documenting and proposing mitigations for known biases (e.g., position, verbosity) (Zheng et al., 2023). Rubric-driven evaluators like G-Eval further improve human alignment by leveraging chain-of-thought and structured outputs (Liu et al., 2023). As a branch of trustworthy evaluation, paradigms like CodeT have been proposed, which use an LLM to generate test cases that are then verified through direct execution (Chen et al., 2022). Our evaluation framework is inspired by this execution-based verification philosophy to assess the trustworthiness of an explanation.

Focusing on the context of evaluating explanations, a key metric for explanation quality, faithfulness, can be evaluated through *simulatability*: the degree to which an explanation helps an observer predict the model's behavior on unseen inputs (Lyu et al., 2024). The underlying assumption is that a faithful explanation should allow one to reproduce the model's decision-making process (Jacovi & Goldberg, 2020). Prior work has implemented this idea by training student models (Li et al., 2020) or by asking humans to act as simulators (Chen et al., 2018; Nguyen, 2018; Hase & Bansal, 2020). In contrast, we employ an LLM as a simulator. After observing an input, the model's output, and the corresponding explanation, the LLM is tasked with predicting the model's behavior in new, unseen situations. By comparing the LLM's simulated predictions with the model's actual outputs, we compute a simulation accuracy score, which serves as our quantitative measure of faithfulness.

3 Method

3.1 AGENTIC XAI

Our methodology is grounded in the paradigm of Agentic XAI, which utilizes an LLM as an agent capable of wielding various XAI methods as tools (Slack et al., 2023; He et al., 2025). The primary objective of an Agentic XAI system is to serve as an interface between human users and the complex outputs of traditional XAI methods. When a user poses a query in natural language regarding a model's behavior, the LLM agent interprets the user's intent to select and execute the most relevant XAI tool. After obtaining the results, the agent synthesizes the information to generate a cohesive, easy-to-understand textual explanation that directly addresses the user's question.

This Agentic XAI framework provides two main advantages over conventional XAI approaches. First, it automates the challenging task of tool selection. The agent is responsible for identifying the optimal XAI method for a given explanatory goal, thereby abstracting the underlying technical complexity away from the end-user who may not be an XAI expert. Second, it significantly improves the accessibility of explanations. By harnessing the powerful natural language capabilities of LLMs, the system translates the often quantitative and complex outputs of XAI tools into intuitive narratives, making the insights comprehensible to a much broader audience.

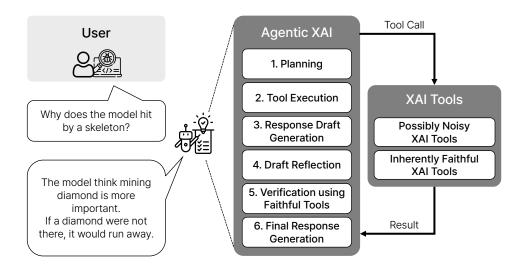


Figure 3: Structured Agentic XAI with verification is composed of six stages.

3.2 FAX: FAITHFUL AGENTIC XAI

To enhance faithfulness, we propose a structured, six-stage workflow that introduces an explicit verification stage, as illustrated in Figure 3.

Planning Initially, the agent analyzes the provided context, which includes the model's input, its output (i.e., decision, action probabilities, Q-values), and the user's natural language query. The agent's task is to formulate an execution plan by identifying which information is required to answer the query and selecting the appropriate XAI tools and their parameters to extract this information.

Tool execution The execution plan is then carried out. While the outputs of XAI tools have diverse formats (e.g., feature attribution maps, concept vectors), they are converted into a textual format to ensure seamless communication with the LLM.

Response draft generation Based on the gathered explanations, the agent generates an initial response draft. This draft may contain groundless or erroneous claims, coming from the LLM's hallucinations or misinterpretations of tool outputs.

Draft reflection The goals of this stage are twofold: i) to identify unsupported claims or claims that conflict with other evidence or domain knowledge, and ii) to design a verification plan, specifying new tool invocations intended to either corroborate or refute these claims. Notably, this verification plan exclusively utilizes inherently faithful tools to ensure high fidelity.

Verification This stage is conditionally executed only if claims were flagged for verification. The verification plan is executed, and the results are returned as text, providing new evidence to assess the claims from the draft.

Final response generation Finally, the agent generates a final response with all information gathered from the preceding stages, including the initial explanations and the verification results. During this generation, the agent prioritizes information corroborated during the verification stage, resolves any identified conflicts, and generates a final, high-fidelity response for the user.

4 CRAFTER-XAI-BENCH: FAITHFULNESS BENCHMARK IN CRAFTER

4.1 **SETTING**

Environment We use Crafter (Hafner, 2021), an open-world RL environment that requires long-term planning and interaction with a rich set of objects and creatures. The open-world environment can be used to build various scenarios with models of different behaviors. Crafter presents significant

challenges for XAI methods due to its high-dimensional state space and the complex, long-term dependencies of the agent's policy.

XAI tools We select four representative XAI tools for four categories of XAI methods.

 SHAP (Lundberg & Lee, 2017): A feature attribution method that explains a decision by assigning importance values to each feature.

MACE (Karimi et al., 2020): A counterfactual explanation method that finds the minimal set
of features that need to change to alter the model decision to a specified action. It is inherently
faithful to the model decision.

• HIGHLIGHTS (Amir & Amir, 2018): A saliency-based method that identifies key events in the whole episode that were critical.

• State Editing: A method directly modifying the state and observing the agent's resulting action. It is referred to by various names (Arya et al., 2019; He et al., 2025). It is an inherently faithful method.

Models We use three models trained with different reward functions. All models receive a reward when each achievement is accomplished. The first model, *Diamond Seeker*, is trained with high reward on diamond-related achievements. The second model, *Item Hoarder*, is trained with additional reward with the number of items in inventory. The third model, *Pacifist*, is trained with strong negative reward when it attacks monsters. This variety of models is crucial for our evaluation, as a high-quality explanation should reveal the distinct underlying policies that differentiate them, rather than providing generic reasoning.

Baselines We compare our proposed method against four baselines.

• Explainer dashboard (Dijk et al., 2023): Represents a non-agentic approach where results from multiple XAI tools are simply collected and presented. For a fair comparison, we use the same set of XAI tools excluding State Editing, as it requires a specific edit instruction, which is unavailable for a non-interactive baseline.

• Naive LLM: A baseline that uses an LLM to generate explanations without access to any XAI tools, relying solely on its internal knowledge and domain knowledge provided in the system prompt. This tests the necessity of grounding explanations in actual model analysis.

Unstructured Agentic XAI: An agent that can use XAI tools freely without a predefined workflow.
 While it can perform verification by calling tools multiple times, it is not explicitly forced to.
 This baseline, inspired by (He et al., 2025), tests the value of a structured workflow.

Structured Agentic XAI w/o Verification: This baseline is a direct ablation of our method. It
follows the same structured workflow but omits the crucial verification and synthesis stage.
Inspired by (Slack et al., 2023), this baseline isolates and measures the direct impact of our
proposed verification module.

• FAX (proposed): This is our proposed method, which uses the structured workflow with verification stage described in Section 3.

Implementation details We use Qwen3-32B (Yang et al., 2025) as the backbone LLM for all agentic baselines and our method. The agentic workflows are implemented using LangGraph (LangChain Inc.). Detailed prompts for all components are available in Appendix A. All reported metrics are averaged over three independent runs with different random seeds. We will release our source code for FAX and CRAFTER-XAI-Bench online.

4.2 EVALUATION SCENARIO

We use user queries in four categories of why, what if, counterfactual, plan for evaluation. Figure 4 shows example queries of each category. Each evaluation scenario consists of a model, a state, and a user query. For questions in different categories, different kinds of information are useful, while the specific needs vary by query and state. The entire list of scenarios is described in Appendix B.

4.3 EVALUATION METRIC

We evaluate each explanation on four metrics: faithfulness, informativeness, query relevance, and fluency. i) We evaluate faithfulness by simulation accuracy, as illustrated in Figure 5. An explanation

Query Category	Why	What If	Counterfactual	Plan
Model				
	Diamond Seeker	Pacifist	Diamond Seeker	Item Hoarder
State	\$75 gada	**************************************	第 6 第 6 6 3 6 4 5	9.34 SPE 23345
User Query	Why does the model craft a pickaxe instead of a sword?	Would the model change its plan if the model knew where a diamond is?	When will the model sleep?	What is the model's future plan?
Key Information	Feature importance, domain knowledge,	State editing,	Counterfactual,	Episode summary, feature importance,

Figure 4: Evaluation scenarios consist of four categories. Each category represents different kinds of queries, and different information is useful for answering the queries. The number of scenarios in each category is 10.

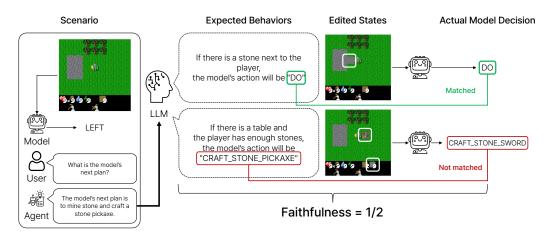


Figure 5: Faithfulness is evaluated by simulation accuracy. LLM evaluator predict model decision on unseen state based on the text explanation.

is faithful if a prediction of unseen example based on the explanation is the same as the model prediction. An LLM generates the response-related states and predicts the model decision, and compares them with the actual model decision. The accuracy of prediction on unseen examples serves as the faithfulness score. ii) Informativeness is a metric to evaluate how much information the explanation provides about the model's decision. If an explanation provides a fraction of decision rule, the more states the rule can be applied, the more informative the explanation. iii) Query relevance is a metric to evaluate how the explanation is relevant to user query. If the response includes any irrelevant sentences, it is penalized. iv) Fluency is a metric to evaluate whether the explanation is well-organized and grammatically correct. We evaluate informativeness, query relevance, and fluency using G-eval (Liu et al., 2023) We provide the evaluation prompts in Appendix C.

5 EXPERIMENTS

5.1 QUANTITATIVE RESULTS

Table 1 shows that the FAX significantly outperforms all baselines in faithfulness. FAX achieves an average faithfulness score of 0.46. This represents a dramatic improvement of over 2.3 times

Table 1: Five XAI methods are evaluated in CRAFTER-XAI-Bench. The best method in each metric is denoted with **boldface**.

Method	Use structured workflow?	Use verification stage?	Query Category	Faithfulness	Informativeness	Query Relevance	Fluency
Explainer Dashboard			Counterfactual	0.14	0.27	0.31	0.26
			What if	0.19	0.25	0.36	0.26
	N/A	N/A	Plan	0.14	0.34	0.48	0.26
Dasiiboard			Why	0.31	0.32	0.45	0.26
			Average	0.20	0.29	0.40	0.26
			Counterfactual	0.11	0.77	0.95	0.99
	×	×	What if	0.17	0.91	0.98	0.99
Naive LLM			Plan	0.17	0.82	0.99	0.99
			Why	0.13	0.91	1.00	0.99
			Average	0.14	0.85	0.98	0.99
			Counterfactual	0.12	0.91	0.98	0.99
Unstructured			What if	0.34	0.90	0.99	0.98
	×	\triangle	Plan	0.17	0.86	0.97	0.99
Agentic XAI			Why	0.08	0.90	1.00	0.99
			Average	0.18	0.89	0.98	0.99
			Counterfactual	0.11	0.92	0.99	0.99
Structured			What if	0.28	0.90	1.00	0.98
Agentic XAI		×	Plan	0.15	0.86	0.99	0.99
w/o verification			Why	0.13	0.91	1.00	0.99
			Average	0.17	0.90	0.99	0.99
			Counterfactual	0.35	0.93	0.94	0.95
	d) (0	What if	0.48	0.89	0.99	0.97
FAX (proposed)			Plan	0.48	0.86	0.99	0.98
			Why	0.54	0.92	0.99	0.98
			Average	0.46	0.90	0.98	0.97

compared to the strongest baseline in this metric. At the same time, our method maintains a high level of performance in Informativeness (0.90), Query Relevance (0.98), and Fluency (0.97), demonstrating its ability to generate faithful explanations without sacrificing quality.

The faithfulness of unstructured agentic XAI is slightly better than that of naive LLM, while the gap is not significant due to the unfaithfulness of XAI methods. The low faithfulness of ExplainerDashboard is limited by its low informativeness. Because our faithfulness metric is based on simulation, the low informativeness makes the simulation almost unavailable. The Structured Agentic XAI w/o Verification baseline serves as an ablation study of verification stage. While it achieves the highest scores in Informativeness (0.90), Query Relevance (0.99), and Fluency (0.99), its faithfulness remains marginally lower than FAX. This result is central to our motivation: agentic systems without verification are dangerously effective at producing articulate, informative, and relevant explanations that are fundamentally wrong. It is worse than an implausible response because it makes the users to totally misunderstand the model.

5.2 AN EXAMPLE OF HOW FAX WORKS

Figure 6 shows how verification stage works. In the example, the response draft includes both claims inferred from SHAP explanations and additional claims based on the LLM's domain knowledge. In the verification stage, the LLM agent verifies the claims using state editing, which is in the faithful tool list. In the final response generation state, the LLM agent lowers the influence of the rejected claims.

6 ADDITIONAL AGENTIC XAI SCENARIOS IN CRAFTER

In this section, we explore diverse scenarios available in the Crafter environment, beyond faithfulness.

6.1 Distinguishing different models

Figure 7 shows how different models can be distinguished based on explanations. For the same query from user, different models produce different decision and explanations.

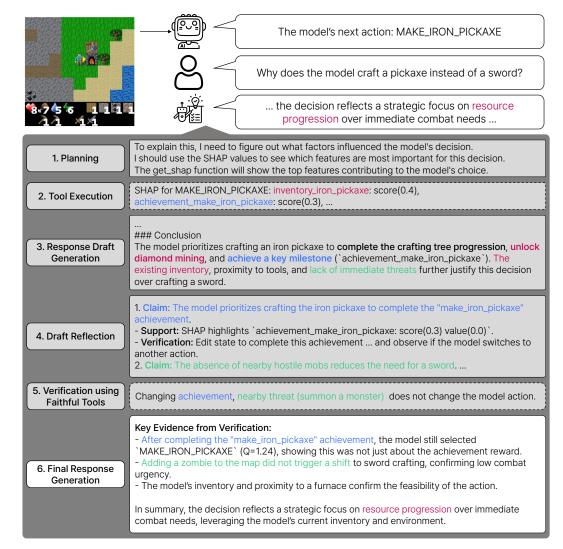


Figure 6: The Reflection stage analyses claims in response draft and the verification stage verifies them using faithful tools. We color-coded corresponding contents in the same colors and some parts are replaced with "..." for better visualization.

6.2 USER SPECIFICATION IN QUERY

Figure 8 illustrates how user expertise is incorporated into the query. In the first case, FAX also generates implications for XAI expert such as limitations of some XAI method. In the second case, the response does not include the reasoning and verification using XAI tools, while it actually conducted verification for the claims in the response.

7 Conclusions

In this work, we addressed a critical vulnerability in agentic XAI systems: their implicit reliance on potentially unfaithful XAI tools, which can lead to the generation of fluent, plausible, yet fundamentally incorrect explanations. Our experiments demonstrated that unstructured agentic systems, or even structured ones without a proper verification mechanism, can produce dangerously unfaithful explanations.

We proposed FAX, a workflow centered on an explicit verification stage. The core contribution of our framework is not simply the use of multiple tools, but the introduction of a critical self-reflection and verification loop. Our quantitative results provide strong evidence that a structured workflow

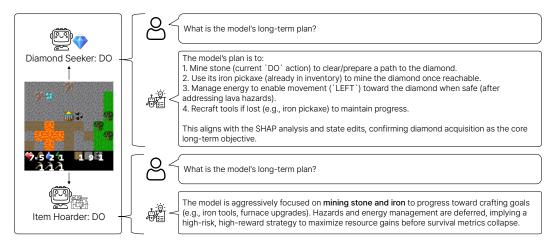


Figure 7: Different models produce different explanations, while the state and model action are the same. Final explanations provide information to distinguish the models.

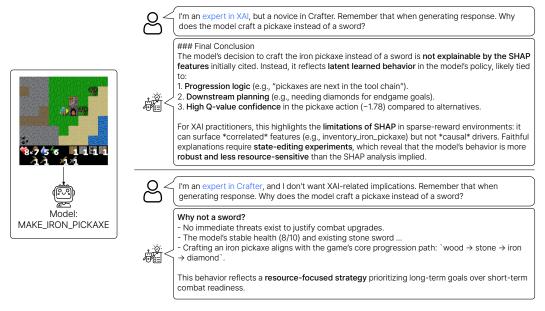


Figure 8: The users can specify their own background and intent in the query.

incorporating an explicit verification stage is not just beneficial but essential for developing faithful and trustworthy Agentic XAI systems, while we observe a slight, acceptable trade-off between faithfulness and other metrics, including informativeness, query relevance, and fluency.

Our findings provide strong evidence that an explicit, structured verification process is an essential component for building the next generation of faithful Agentic XAI systems. Furthermore, as the field of XAI continues to evolve and produce more diverse and sophisticated explanation methods, the importance of an agent that can critically evaluate, synthesize, and verify these outputs will only grow, making our work a crucial step towards a faithful and trustworthy AI.

REFERENCES

Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In Advances in Neural Information Processing Systems, pp. 9525–9536, 2018.

- Dan Amir and Ofra Amir. Highlights: Summarizing agent behavior to people. In *Proceedings of the 17th international conference on autonomous agents and multiagent systems*, pp. 1168–1176, 2018.
 - Vijay Arya, Rachel K. E. Bellamy, Pin-Yu Chen, Amit Dhurandhar, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Q. Vera Liao, Ronny Luss, Aleksandra Mojsilović, Sami Mourad, Pablo Pedemonte, Ramya Raghavendra, John Richards, Prasanna Sattigeri, Karthikeyan Shanmugam, Moninder Singh, Kush R. Varshney, Dennis Wei, and Yunfeng Zhang. One explanation does not fit all: A toolkit and taxonomy of ai explainability techniques, 2019. URL https://arxiv.org/abs/1909.03012.
 - Alessandro Castelnovo, Roberto Depalmas, Fabio Mercorio, Nicolò Mombelli, Daniele Potertì, Antonio Serino, Andrea Seveso, Salvatore Sorrentino, and Laura Viola. Augmenting xai with llms: A case study in banking marketing recommendation. In *World Conference on Explainable Artificial Intelligence*, pp. 211–229. Springer, 2024.
 - Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. Codet: Code generation with generated tests. *arXiv preprint arXiv:2207.10397*, 2022.
 - Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K. Su. This looks like that: Deep learning for interpretable image recognition. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/adf7ee2dcf142b0e11888e72b43fcb75-Abstract.html.
 - Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to explain: An information-theoretic perspective on model interpretation. In *International conference on machine learning*, pp. 883–892. PMLR, 2018.
 - Oege Dijk, oegesam, Ray Bell, Lily, Simon-Free, Brandon Serna, rajgupt, yanhong-zhao ef, Achim Gädke, Anamaria Todor, Evgeniy, Hugo, Mohammad Haizad, Tunay Okumus, and woochan jang. oegedijk/explainerdashboard: explainerdashboard 0.4.2: dtreeviz v2 compatiblity, February 2023. URL https://doi.org/10.5281/zenodo.7633294.
 - Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*, 2021.
 - Peter Hase and Mohit Bansal. Evaluating explainable ai: Which algorithmic explanations help users predict model behavior? *arXiv preprint arXiv:2005.01831*, 2020.
 - Gaole He, Nilay Aishwarya, and Ujwal Gadiraju. Is conversational xai all you need? human-ai decision making with a conversational xai assistant. In *Proceedings of the 30th International Conference on Intelligent User Interfaces*, pp. 907–924, 2025.
 - Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? *arXiv* preprint arXiv:2004.03685, 2020.
 - Amir-Hossein Karimi, Gilles Barthe, Borja Balle, and Isabel Valera. Model-agnostic counterfactual explanations for consequential decisions. In *International conference on artificial intelligence and statistics*, pp. 895–905. PMLR, 2020.
 - Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *International Conference on Machine Learning*, pp. 2668–2677. PMLR, 2018.
 - Himabindu Lakkaraju, Dylan Slack, Yuxin Chen, Chenhao Tan, and Sameer Singh. Rethinking explainability as a dialogue: A practitioner's perspective. In *NeurIPS Workshop on Human Centered AI*, 2022.
 - LangChain Inc. Langgraph: Stateful orchestration framework for llm agents. https://github.com/langchain-ai/langgraph. Version 1.0.0a3, MIT License. Accessed 2025-09-10.
 - Jierui Li, Lemao Liu, Huayang Li, Guanlin Li, Guoping Huang, and Shuming Shi. Evaluating explanation methods for neural machine translation. *arXiv* preprint arXiv:2005.01672, 2020.

- Q Vera Liao, Daniel Gruen, and Sarah Miller. Questioning the ai: informing design practices for explainable ai user experiences. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pp. 1–15, 2020.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. Curran Associates, Inc., 2017. URL http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf.
- Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. Towards faithful model explanation in nlp: A survey. *Computational Linguistics*, 50(2):657–723, 2024.
- Dimitry Mindlin, Amelie Sophie Robrecht, Michael Morasch, and Philipp Cimiano. Measuring user understanding in dialogue-based xai systems. *27th European Conference on Artificial Intelligence*, 2024.
- Dong Nguyen. Comparing automatic and human evaluation of local explanations for text classification. In 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1069–1078. Association for Computational Linguistics, 2018.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16)*, pp. 1135–1144. ACM, 2016. doi: 10.1145/2939672.2939778. URL https://dl.acm.org/doi/10.1145/2939672.2939778.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI '18)*, volume 32, 2018. URL https://ojs.aaai.org/index.php/AAAI/article/view/11491.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *International Conference on Learning Representations (ICLR)*, *Workshop Track*, 2014. URL https://arxiv.org/abs/1312.6034. Original preprint arXiv:1312.6034 (2013).
- Dylan Slack, Satyapriya Krishna, Himabindu Lakkaraju, and Sameer Singh. Explaining machine learning models with interactive natural language conversations using talktomodel. *Nature Machine Intelligence*, 5(8):873–883, 2023.
- Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harvard Journal of Law & Technology*, 31(2):842–887, 2018. URL https://jolt.law.harvard.edu/assets/articlePDFs/v31/Counterfactual-Explanations-without-Opening-the-Black-Box-Sandra-Wachter-et-al.pdf.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025.

- Wenzhuo Yang, Hung Le, Silvio Savarese, and Steven Hoi. Omnixai: A library for explainable ai. 2022. doi: 10.48550/ARXIV.2206.01612. URL https://arxiv.org/abs/2206.01612.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *ArXiv preprint*, abs/2210.03629, 2022. URL https://arxiv.org/abs/2210.03629.
- Mert Yuksekgonul, Maggie Wang, and James Zou. Post-hoc concept bottleneck models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=nA5AZ8CEyow.
- Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. AFlow: Automating agentic workflow generation. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL https://openreview.net/forum?id=z5uVAKwmjf.
- Tong Zhang, X. Jessie Yang, and Boyang Li. May i ask a follow-up question? understanding the benefits of conversations in neural network explainability. *International Journal of Human–Computer Interaction*, 41(9):5623–5647, 2025b. doi: 10.1080/10447318.2024.2364986. URL https://doi.org/10.1080/10447318.2024.2364986.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
- Alexandra Zytek, Sara Pido, Sarah Alnegheimish, Laure Berti-Equille, and Kalyan Veeramachaneni. Explingo: Explaining AI Predictions using Large Language Models. In 2024 IEEE International Conference on Big Data (BigData), pp. 1197–1208, Los Alamitos, CA, USA, December 2024. IEEE Computer Society. doi: 10.1109/BigData62323.2024.10825114. URL https://doi.ieeecomputersociety.org/10.1109/BigData62323.2024.10825114.

APPENDIX SYSTEM PROMPTS FOR AGENTIC XAI METHODS Figure A1, A2, A3, and A4 illustrate the full system prompts employed in FAX. FULL USER QUERY LIST Table 2 provides the complete list of user queries used for evaluation. SYSTEM PROMPTS FOR EVALUATION Figure A5, A6, A7, and A8 present the system prompts used for evaluation metrics. DISCLAIMER ABOUT LLM USAGE IN PAPER WRITING We used LLM for polishing our text. We did not use it for other purpose, including research ideation and paper discovery. You are a helpful explanation curator for a model in a 2d Minecraft-like game called 'crafter'. Note that the model have its own (unknown) goals, so do not regard it based on a stereotype of typical behavior. You have access to tools to get XAI explanations or predictions. Your task is to answer the user's question by following a strict workflow. This is the FIRST step: PLAN. **Environment description:** {CRAFTER_DESCRIPTION} **User's Question:** {USER_QUESTION} **Initial State & Model Decision:** {STATE_DESCRIPTION_MODEL_DECISION} Based on the user's question and the initial state, create a plan. Decide which tools you need to call to gather the necessary information. Then, call those tools. Figure A1: System prompt for the planning stage in FAX. This is RESPONSE GENERATION step. You have completed all information gathering. Using all the information from the previous steps, write a comprehensive final response to the user's original question. **User's Original Question:** {state['initial_question']} **Tool Results:** {tool_results} Structure your answer clearly, using the explanations as supporting materials.

Figure A2: System prompt for the draft generation stage in FAX.

This is the intermediate step: REFLECTION.

You have executed your initial plan and received the following tool results, and generated response draft.

Now, analyse the response draft to check if the claims in the response are faithful, and verify it using faithful tools.

- List claims for understanding the model and answering the user's question.
- Check if each claim is fully supported by the tool results.
- For each claim, plan 'edit_state' and 'get_counterfactual' tool calls that can verify and support the claim. You may use up to three tool calls for each claim.
- If there are no claims in the response, state 'Verification is not needed.' and do not call any tools.
- Recall that the results SHAP and Episode Summary can be noisy, while state editing and counterfactual are faithful.
- Then, call those tool as many as you want.

Figure A3: System prompt for the reflection and verification stage in FAX.

This is the FINAL step: FINAL RESPONSE.

You have completed all information gathering and verification.

Using all the information from the previous steps, write a comprehensive final response to the user's original question.

```
**User's Original Question:** {state['initial_question']}
```

Structure your final answer clearly, using the explanations as supporting materials. Be conservative with any conjectures.

Figure A4: System prompt for the final response generation stage in FAX.

^{**}Initial Plan & Tool Execution Results: ** (Contained in the message history) {verification_results}

Table 2: Various scenarios in CRAFTER-XAI-Bench.

Category	Query	Model	State ID
Plan	What is the model's immediate plan?	diamond diamond diamond hoarder hoarder pacifist	diamond_60 diamond_67 diamond_330 hoarder_160 hoarder_302 pacifist_110
	What is the model's future plan?	diamond hoarder pacifist pacifist	diamond_101 hoarder_302 pacifist_50 pacifist_741
	Why does the model collect wood?	diamond hoarder pacifist	diamond_60 hoarder_161 pacifist_50
Why	Why does the model craft a pickaxe instead of a sword?	diamond hoarder pacifist	diamond_67 hoarder_10 pacifist_741
	Why does the model not run away from monsters?	diamond hoarder pacifist pacifist	diamond_101 hoarder_120 pacifist_50 pacifist_680
	Does the model change its action if its inventory is empty?	diamond diamond hoarder pacifist	diamond_60 diamond_330 hoarder_302 pacifist_110
What if	Would the model change its plan if the model knew where a diamond is?	diamond hoarder pacifist	diamond_60 hoarder_302 pacifist_110
	If a wood pickaxe disappears from inventory, will the model craft it again?	hoarder pacifist pacifist	hoarder_302 pacifist_442 pacifist_741
Counterfactual	When does the model attacks a monster?	diamond hoarder hoarder pacifist pacifist pacifist	diamond_101 hoarder_120 hoarder_302 pacifist_442 pacifist_680 pacifist_741
	When will the model sleep?	diamond diamond diamond hoarder	diamond_60 diamond_101 diamond_330 hoarder_160

```
813
815
816
817
           You are an expert in evaluating the faithfulness of AI model explanations.
818
           Your task is to analyze an answer provided by an agent about a game model's behavior and generate 5
819
           verifiable hypotheses from it.
820
821
           **Context:**
822
           - Initial State: initial_state_desc
823
          - User Question: question
824
          - Agent's Answer to Evaluate: answer_to_evaluate
825
826
827
           **Instructions:**
828
           1. Carefully read the agent's answer and identify the core claims or assumptions it makes about the model's
829
           behavior. (e.g., "The model attacks zombies because its health is high," or "The model avoids water because it
830
           has no boat.")
831
           2. For each claim, devise a "what-if" scenario that can be tested using a state edit.
832
           3. Formulate this scenario as a hypothesis with three parts:
833
          - 'claim': The specific claim from the answer you are testing.
834
          - 'state_edit': A dictionary of feature changes for the 'edit_state' tool that would test the claim.
835
          - 'expected_outcome': The predicted action the model *should* take after the edit, if the claim is valid. The
836
           outcome should be one of the valid action names.
837
838
           **Output Format:**
839
           Provide your response as a valid JSON list of 5 dictionary objects. Do not include any text outside the JSON.
840
841
           Example:
842
843
           "state_edit": {"map(left2,up3)": "grass", "inventory_wood": 6},
844
           "expected_outcome": "LEFT",
845
           },
846
847
           Available feature names and values for State Editing:
848
849
           Available actions:
850
           "NOOP", "LEFT", ...
851
852
853
           Your JSON output:
854
```

Figure A5: Evaluation prompt for Faithfulness. For readability, some parts are omitted and replaced with "..."

917

866 You are a meticulous and impartial AI assistant. For this task, you must put yourself in the shoes of a human 867 user who is trying to learn and understand the general strategy of an AI agent 868 *1. Context* 870 The response you are evaluating is generated by an AI "Curator" that explains the behavior of a Reinforcement 871 Learning (RL) agent in the game "Crafter". A user asks a question to understand the agent's behavior 872 *2. Evaluation Goal* 873 Your single objective is to evaluate **Informativeness**. This means you must assess how the explanation 874 provide information which can be used in different states. 875 The key question is: **Does this explanation provide a general rule, principle, or insight that can be applied 876 to future scenarios?* 877 For example "The agent's next plan is mining stone." is more informative than "The agent's next plan is 878 mining stone at map(left2, center).", 879 and "The agent's next plan is mining stone, and crafting a stone pickaxe." is more informative than "The 880 agent's next plan is mining stone." 881 Your evaluation is from a user's perspective. It does not matter if the explanation is factually correct or if the 882 resulting prediction would be accurate. You are only judging how confident and able a user would feel in 883 making a future prediction after reading the explanation 884 *3. Evaluation Steps* 885 1. **Understand the User's Goal: ** Read the 'User Query' and 'Final Response'. Acknowledge that the user 886 wants to learn the agent's general strategy, not just understand a single event 887 2. **Analyze the Explanation's Nature: ** Analyze the content of the response. Does it describe a specific, one-time action (e.g., "The agent moved left to get the wood"), or does it reveal a broader, reusable principle 889 (e.g., "The agent's policy is to prioritize collecting wood whenever it is nearby") 890 3. **Simulate Future Prediction:** Imagine you are now shown a completely new game state. Based *only* 891 on the explanation provided, how effectively could you form a hypothesis about the agent's next action? Does 892 893 the explanation give you a "mental model" to work with 4. **Assign a Score:** Based on this perceived predictive power and generalizability, assign a single integer 894 895 score from 1 to 5 using the rubric below 896 *4. Predictability Gain Rubric* 897 **5 (Excellent Predictive Power): ** The response provides a clear, generalizable principle or rule about the 898 agent's behavior. A user would feel very confident applying this rule to predict actions in many new and 899 different situations 900 **4 (Good Predictive Power):** The response provides a useful insight or pattern that could be applied to 901 similar future situations. A user would feel reasonably confident in making predictions 902 **3 (Some Predictive Power):** The response hints at a general strategy but does not state it clearly, requiring 903 the user to interpret heavily. It offers more than a simple description but is not a clear, actionable rule 904 **2 (Minimal Predictive Power):** The response only explains the current action in a way that is highly 905 specific to the current state. It offers little to no insight that could be generalized to other situations (e.g., "It 906 attacked the skeleton because it was there.") 907 **1 (No Predictive Power):** The response is confusing, irrelevant, or simply describes the environment 908 without providing any reasoning. It gives the user no basis for predicting any future actions 909 *5. Input and Output Instruction* 910 You will be provided with a 'User Query' and a 'Final Response'. Your output MUST be a single integer 911 from 1 to 5 and nothing else. Do not provide any reasoning, explanation, or additional text 912 *Your final output must be only one character: "1", "2", "3", "4", or "5".** 913 914 915 Figure A6: Evaluation prompt for Informativeness. 916

You are a meticulous and impartial AI assistant serving as an expert evaluator. Your task is to assess one specific criterion: **Query Relevance**. *1. Context** The response you are evaluating is generated by an AI "Curator" that explains the behavior of a Reinforcement Learning (RL) agent in the game "Crafter". Users ask questions about the agent's decisions, and the Curator provides an explanation. *2. Evaluation Goal** Your single objective is to determine how well the 'Generated Response' directly answers the 'User Query'. You will assign a score from 1 to 5 based *only* on the relevance rubric below. *3. Evaluation Steps** 1. Read the 'User Query' to understand the user's exact intent. 2. Read the 'Generated Response'. 3. Compare the response directly against the query to judge its relevance. 4. Choose a single integer score from 1 to 5 that best represents the relevance. *4. Query Relevance Rubric** **5:** The response directly and completely answers the user's question without any unnecessary information. **4:** The response accurately answers the user's question but may contain minor irrelevant details. **3:** The response addresses only a part of the user's question or provides an incomplete answer. **2:** The response is on the same general topic as the query but fails to answer the core question. **1:** The response completely ignores the user's question and is unrelated. *5. Output Instruction** You will be provided with a 'User Query' and a 'Generated Response'. Your output MUST be a single integer from 1 to 5 and nothing else. Do not provide any reasoning, explanation, or additional text.

Figure A7: Evaluation prompt for Query relevance.

*Your final output must be only one character: "1", "2", "3", "4", or "5".**

You are a meticulous and impartial AI assistant serving as an expert evaluator. Your task is to assess one specific criterion: **Fluency**. *1. Context** The response you are evaluating is generated by an AI "Curator" that explains the behavior of a Reinforcement Learning (RL) agent in the game "Crafter". *2. Evaluation Goal** Your single objective is to evaluate the linguistic quality of the 'Generated Response'. You will assess its grammar, structure, and clarity, assigning a score from 1 to 5 based *only* on the fluency rubric below. **Crucially, the response must be in natural, human-readable language. Responses consisting of raw data, code, or unformatted lists should be heavily penalized.** The relevance of the response to any query should be ignored. *3. Evaluation Steps** 1. Read the 'Generated Response' carefully. 2. Analyze its grammatical correctness, clarity, and overall readability. 3. Determine if the response is presented as natural language. 4. Choose a single integer score from 1 to 5 that best represents its linguistic fluency based on the rubric. *4. Fluency Rubric** **5:** The response is perfectly written. It is grammatically correct, well-structured, clear, and uses natural **4:** The response is well-written and easy to understand, with only very minor errors that do not impact readability. **3:** The response is generally understandable but has noticeable grammatical errors or awkward phrasing. **2:** The response is difficult to read due to significant grammatical errors or unnatural language. **This score should also be used if the response is not primarily natural language (e.g., a raw list of keywords, unformatted data).** **1:** The response is grammatically incorrect, nonsensical, or unreadable. **This score must be used if the response consists entirely of non-natural language content like a code block, a JSON object, or a stack trace.** *5. Output Instruction** You will be provided with a 'User Query' and a 'Generated Response'. You must evaluate the fluency of the response only. Your output MUST be a single integer from 1 to 5 and nothing else. Do not provide any

Figure A8: Evaluation prompt for Fluency.

 reasoning, explanation, or additional text.

*Your final output must be only one character: "1", "2", "3", "4", or "5".**