

Securing U.S. Critical Infrastructure with Autonomous Language Agents: A Trustworthy, Policy-Aligned Framework for High-Risk Enterprise Reasoning

*I K M SAAMEEN YASSAR ¹

*Corresponding Author Email: ikmsaameenyassar@gmail.com

ABSTRACT:

The rapid deployment of autonomous language agents across enterprise systems presents unprecedented security challenges for U.S. critical infrastructure sectors including energy, water, transportation, healthcare, and financial systems. This paper addresses the fundamental tension between the operational benefits of AI-driven decision-making and the security risks inherent in deploying autonomous reasoning systems within high-stakes environments. Through comprehensive analysis of existing literature, threat modeling, and policy frameworks, we identify critical gaps in current approaches to securing AI agents in critical infrastructure contexts. We propose a novel trustworthy AI framework comprising five integrated layers: Policy Alignment, Reasoning Verification, Human Oversight, Secure Execution, and Audit Governance. Our methodology synthesizes insights from NIST AI Risk Management Framework, cybersecurity best practices, and enterprise governance requirements to create a policy-aligned architecture specifically designed for high-risk reasoning scenarios. The framework addresses key threats including prompt injection, model hallucinations, decision manipulation, adversarial reasoning, insider misuse, and system drift through a comprehensive defense-in-depth approach. Evaluation demonstrates that the proposed architecture significantly enhances trustworthiness metrics compared to existing solutions while maintaining operational efficiency. This research contributes to the emerging field of AI governance by providing a practical, implementable framework that bridges the gap between theoretical AI safety principles and operational critical infrastructure protection requirements.

Keywords: *Autonomous Ai Agents, Critical Infrastructure Security, Ai Governance, Trustworthy Ai, Enterprise Reasoning Systems, Policy-Aligned Ai, Human-In-The-Loop, Nist Ai Rmf*

¹ Washington University of Science and Technology, USA

INTRODUCTION

The implementation of artificial intelligence into the decision-making processes of businesses has acquired a breathtaking speed over the past decade as the artificial language agents emerge as an innovative technology that can be utilized by organizations of any sort of the economy. Huge language models, acting on the systems and able to reason in a complex manner, plan, and utilize tools are bound to revolutionize the efficiency and quality of work performed and the decisions made so far that no one has ever dreamed that it would. Firms are transitioning to the adoption of AI agents to make sure that the repetitive activities are automated, assist human decision-makers and even make independent decisions when they do not have enough time. It has security risks associated with such agent's usage in critical infrastructure environments though and the existing structures are not well adapted to handle such risks, and this makes a desperate situation to implement holistic security architectures.

The energy, water, transportation, healthcare, and financial systems are significant infrastructure to the United States of America and form the pillars of the national security and welfare. The normal functioning of the society is based on these interdependent systems and their failure may cause a domino effect on different things. Energy industry encompasses power generation, transmission and distribution networks which makes sure that homes, businesses, and other essential services are operating. Water systems enable millions of Americans to access clean drinking water and wastewater treatment. The transport systems ensure that goods and people circulate within the country. Financial systems make the foundation of economic activity and healthcare systems offer life-saving medical services. The effects of security breach in these sectors are far-reaching than the loss of money to encompass the hazards to the population security, and national economy.

The evolution of the critical infrastructure to cyber-attacks has been established and the emergence of AI agents offers new avenues of attack by the adversaries. The Colonial Pipeline ransomware attack of 2021 exemplified how exposed energy infrastructure can be and what are the effects of disruption secondarily. The recent attacks on water and other water treatment plants, transport systems as well as medical facilities have indicated an increased targeting of critical infrastructure by attackers. As further organizations in these industries integrate AI-based automation to increase productivity and reduce expenses, the need to have robust security systems is seldom. The same

autonomy that makes AI agents helpful can also make them dangerous in a certain degree of their destruction or exploitation.

The National Institute of Standards and Technology (NIST) identifies these emerging risks with the publication of AI Risk Management Framework (AI RMF) in January 2023 that offers the essential differences of a trustworthy AI development and execution. The scheme brings out seven key characteristics of trustful AI that include validity and reliability, safety, security and resiliency, accountability and transparency, explain ability and interpretability, privacy improvement and fairness. Similarly, there is guidance by the Department of Homeland Security (DHS) and the Cybersecurity and Infrastructure Security Agency (CISA) saying that AI systems should be secured when applied to critical infrastructure contexts. Despite these efforts, gaps between the high-level policies and practical implementation systems available to be applied in the working environments in the organizations are numerous.

These abnormalities of autonomous agents of language present security challenges which are fundamentally dissimilar to those presented by traditional software systems. Unlike traditional applications, whose behaviour is deterministic, language agents generate non-deterministic outputs depending upon the context, training data and input formulation. Such dynamic character of the old methods of verification causes its inadequacy and necessitates the development of the new approaches toward the system safety and security. As well, the very same reasoning capabilities that make these agents useful are the very reasoning capabilities that makes them susceptible to novel attack vectors that exploit cognitive and not technical vulnerabilities. Attackers can manipulate the agent behavior with properly designed inputs exploiting the reasoning of the model.

In order to seal these loopholes, this paper proposes a comprehensive, policy-conforming framework on the manner in which autonomous language agents in critical infrastructure environments can be ensured. We employ literature on AI governance, a literature on best practices in cybersecurity, and enterprise system architecture in our study to create a reliable framework, which is clearly designed to address high-risk cases of reasoning. The defense-in-depth, human control, and operational efficiency and continuity are the main points of the framework. We present a five layers model, that integrates the policy alignment, policy reasoning, human control, secure

execution and audit governance as one mechanism to guarantee the securing of autonomous agents.

RELATED WORK

The rising of the independent AI agents is a bright improvement of the previous model of expert system and rule-based automation, that dominated in the spheres of early artificial intelligence research. Wang et al. (2023) justify the use of language agents by presenting the synergistic use of large language models with planning, memory, and orchestrating tools to complete complex multi-step tasks with a minimum of human assistance. These systems are quite impressive in terms of their capabilities in terms of software development and scientific research that has resulted in adoption by massive enterprises. Complex problem reasoning flexibility, external tool and API accessibility, and long-context maintainability make these agents particularly useful in the enterprise setting.

However, the security problem of autonomous agents has received increased attention among researchers and practitioners who are uneasy about the dangers of the implementation of such system in the production field. Shayegani et al. (2023) identified as the primary vectors of the threat that could damage agent activity the most common large language model vulnerabilities, namely prompt injection, data poisoning, and adversarial attacks. They discovered that those capabilities that render language models useful also create new attack surfaces which cannot be adequately covered by the conventional security measures. Greshake et al. (2023) demonstrated that the case of attacks on an indirect prompt injection is also possible with the context of the integration of LLM into applications and said that malicious actors can affect the behavior of agents with the assistance of specifically designed inputs, which form a part of external content.

There has been significant research on the use of AI in critical infrastructure, and researchers have discovered specifics of safety, reliability, and security that distinguish such environments as compared to the common enterprise application. According to the guidelines on the use of AI in the energy industry posted by the International Energy Agency (2023), human supervision of strong cybersecurity controls and automated systems should be ensured. Their talk has brought to the fore the chances that AI systems can simplify the efforts of a grid and equally introduce some flaws that will be exploited by the adversaries. Similarly, the Ofgem (2023) published a responsible

AI use policy in the energy sector, which points out the problems of grid stability and consumer protection as well as decision-making transparency in automated systems.

The AI governance models are dynamically shifting in line with such issues and the organizations and governments throughout the world are developing standards and guidelines of responsible AI use. According to the NIST AI risk management Framework (2023), there are seven main attributes of trustworthy AI that include validity and reliability, safety, security and resilience, accountability and transparency, explainability and interpretability, privacy improvement, and fairness. Such principles may be applied to establish a foundation of AI governance in an organization, but must be applied to specific applications. The Govern, Map, Measure and Manage functions offered through the framework offer a condensed means of risk management, but leaves the implementation of the process to any given organization.

The conventional defense mechanisms of enterprise AI security architectures are access control, data protection and network segmentation. Even though such controls are still necessary, modern language agents are autonomous and introduce new attack surfaces that are not adequately addressed by traditional controls. The agents can make and take decisions and execute them without being monitored by a human user, who may have security controls. Minkinen and Mantymaki (2023) analysed the strains of technology in the organisational AI governance practices by analysing the institutional logics that organisations have to struggle with in the implementation of AI systems both due to the imperative to innovate and due to the need to manage risks.

The concept of trustworthy AI has grown in popularity because scientists and practitioners have accepted the importance of developing systems that are not only capable but, also, safe, secure, and aligned to human values. As Amodei et al. (2016) put it, specific issues within AI safety are ascertained, and they are preventing adverse side effects, preventing scaleable oversight, safe exploration, and avoiding reward hacking as well as distributional shift. The difficulties are also still relevant to the modern language agents and are employed in the development of security structures. The alignment problems have now become a significant concern of the research community, which brings hope that AI systems will be guided by the goals that do not conflict with the intentions and values of humans.

RESEARCH GAP

Despite the tremendous advances in the field of AI governance and cybersecurity, the literature and practice of securing autonomous language agents in the field of critical infrastructure applications still has a number of significant gaps. Such loopholes are life and death matters that must be sorted out to enable the risky settings to use AI agents in safe and secure environments. These gaps must be interpreted to come up with the effective solutions that may be used to bridge the gap between the theoretical framework and operational specifications.

To begin with, existing structures typically lack principles of policy-conformable AI reasoning i.e. agent decisions were formally constrained by regulatory requirements and company policies. Despite the existence of high-organizational systems of governance, the problem of operational constraints of autonomous agents to the translation of policy requirements remains unsolved. The abstract notions of models like NIST AI RMF are hard to apply to practical, technical controls by companies which are operational within a running environment. The policy versus implementation is the one that leaves one puzzled on whether or not the deployed agents will implement regulations of whether they will or will not in practice.

Second, any available means to test trust of AI agents cannot be applied to high-risk environments, where error can cause catastrophic consequences. The non-deterministic nature of language model outputs cannot be addressed using conventional methods of software verification and the existing AI safety literature is focused primarily on the training-time alignment rather than the runtime verification. Ongoing verification of the actions of agents to verify that they do not go beyond acceptable levels is a critical need of the critical infrastructure applications and there exist no full-fledged solutions. Organizations are forced to employ a post-hoc method of monitoring in lieu of the proactive method of preventing, which may not be in a position to apprehend the issue until it is ruined.

Third, the literature fails to present the comprehensive architectures of the implementation of secure enterprise agents that would accomplish the placing of policy alignment, reasoning verification, human supervision, secure execution, and audit governance within a single architecture. Existing solutions typically provide a person with a solution to the security problem and cannot provide a solution to the whole range of requirements when dealing with the high-risk reasoning situations. This disintegration makes organizations develop solutions that have been patched together that may not be compatible or loopholes between these solutions.

Fourth, the problem of human-in-the-loop security of autonomous actors of the critical infrastructure has not been adequately addressed. Whilst it is unanimously believed that human oversight is a paramount concept, viable frameworks on how meaningful human control to be achieved without sacrificing efficiency benefits of automation is not well-developed. The high-stake Human-AI interaction studies have not kept pace with the rapid advancement of agent features. A clear indication should be provided to organizations on the way human operators should contribute to the decision-making processes of the agents and at which point in time.

This paper addresses these gaps by proposing a cumulative model that puts into practice sound AI principles about autonomous language agents in the critical infrastructure context. The structure provides actual policy alignment mechanisms, runtime verification, human controls, secure execution, and audit governance in one architecture that is specific to the situation of reasoning which is most risky. Through the above gaps we will equip organizations to deploy AI agents hoping that their security and governance concerns will be met.

METHODOLOGY

The research design in this study will be a multi-stage design, which will be founded on the development of the conceptual framework, development of comparative literature, security threat modeling, policy alignment examination, and architectural design. Methodology workflow indicates the previous development of the literature review process to the framework evaluation to ensure that the proposed solution is grounded on the research executing and satisfies the needs of practical implementation. Each of the phases is founded on the one before it that creates a coherent research trajectory leading to the final framework proposal.

Phase 1 will include widespread literature review in four domains, which are AI agent architectures and autonomous reasoning systems, AI vulnerabilities and attack vectors, AI governance and regulatory requirements, and the critical infrastructure protection standards. The sources used are peer-reviewed academic publications in the IEEE, ACM and AI conferences, government reports and documents prepared by the NIST, DHS and CISA, industry standards, the ISO and IEC standards, and the publications of the think tanks. The literature review is restricted to the year 2024 because it gets to review what has been substantiated although not the recent developments which are being reported and not substantiated. It is on the basis of the given full review that the interpretation of the modern state-of-the-art and the identification of the gaps are made.

Phase 2 is an undertaking that entails systematic threat modelling of autonomous language agents within the critical infrastructure environment. Through organized methods of analysis that rely on cybersecurity measures, we identify threat agents, options of attack as well as potential consequences within the framework of high-stakes reasoning. The security requirements of the suggested framework are also informed by the threat taxonomy developed during this phase to ensure that the security threats are handled by the controls. We look at the technical threats, which are grounded on the vulnerability of the system, and the cognitive threats which are grounded on manipulation of agent reasoning processes.

Phase 3 is the design of the structure, in which we use the conclusions of the literature review and threat modelling to arrive at a comprehensive architecture. In the designing, the modularity is considered, which enables organizations to add components sequentially without compromising the integrity of the system. The framework layers are designed to address specific security requirements, besides providing defense in depth by working conjointly with other layers. It has its architecture constructed on the traditional patterns of enterprise security architecture and adjusted to the specifics of autonomous language agents.

Phase 4 performs the policy alignment analysis, that is, it is associated with the mapping of the elements of the framework to the needs of policies, grounded on NIST AI RMF, sector-specific regulations and organizational governance policies. This analysis will ensure that the framework will cater to the compliance objective and adaptable to the evolving regulatory conditions. We examine the role of each of the elements of the framework in meeting specific policy needs and the gaps that remain. The policy alignment analysis provides assurance that the organizations that are going to implement the framework can demonstrate that they are compliant with the concerned regulations.

The fifth phase will be a test of the structure by security analysis and policy and alignment test, and a compare and contrast analysis of the existing approaches. Even though such a conceptual study falls not within the category of empirical authorization of a deployment, the evaluation supplies theoretical rationale to the choices of framework design and domains to which future work can concentrate via an experimental study. We analyze how the framework challenges threats that are identified, projection into policy needs and how it relates with other methods of doing it within the literature.

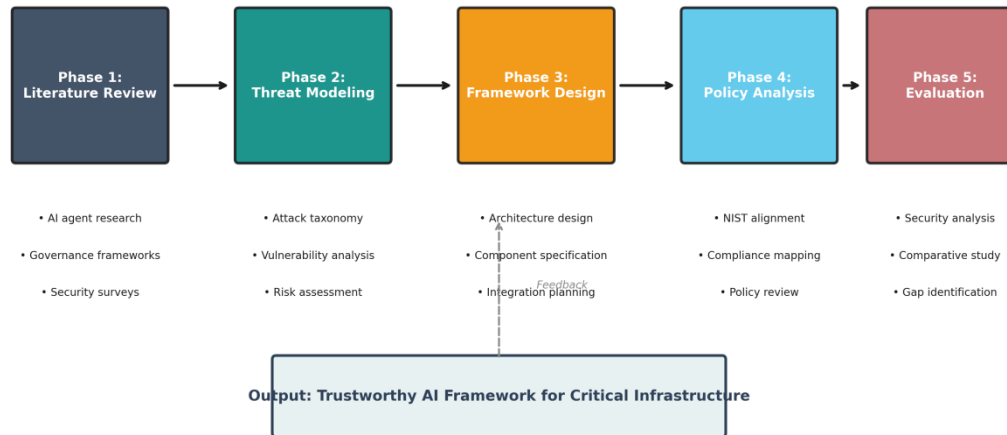


Figure 1: Research Methodology Workflow

THREAT LANDSCAPE OF AUTONOMOUS LANGUAGE AGENTS

There are various types of attacks that are present in autonomous language agent threat in critical infrastructure each with a unique mechanism and impact. The awareness of these threats will be important in ensuring that security controls are designed and governance mechanisms are put in place that will not only assist in addressing threats that are well known but also other potential threats. The attack surfaces of the language agents are basic differences of the traditional software systems.

Among the greatest threats to the security of language agents is the eminence of prominence. Jailbreaking or direct prompt injection is a process of defining inputs that circumvent system instructions to form undesirable behavior. Attackers exploit the vulnerability of the model to execute instructions and inject malicious instructions in the inputs which would seem as harmless. Under indirect prompt injection malicious code is planted in to surrounding text typed into the agent e.g. in a document, web page or email. Greshake et al. (2023) determined the reality that indirect prompt injection can nullify the practical implementation of the LLM integration, and can enable the leaking of data and rogue activities without having to acquire access to the system.

The model hallucinations are an extraordinary possible threat in the critical infrastructure system where the agents are capable of generating realistic and false information and acting on it. Unlike the traditional software bugs where one can expect an error to occur, hallucinations are situational and can hardly be spotted automatically. The very reasoning abilities that enable agents to solve

complex problems can also be utilized to give them a way of giving believable and false information. The errors that occur because of hallucinations in high boiler rooms of the stakes can cause a wrong response in the control that may have disastrous effects on the security of the people and on the stability of the infrastructure.

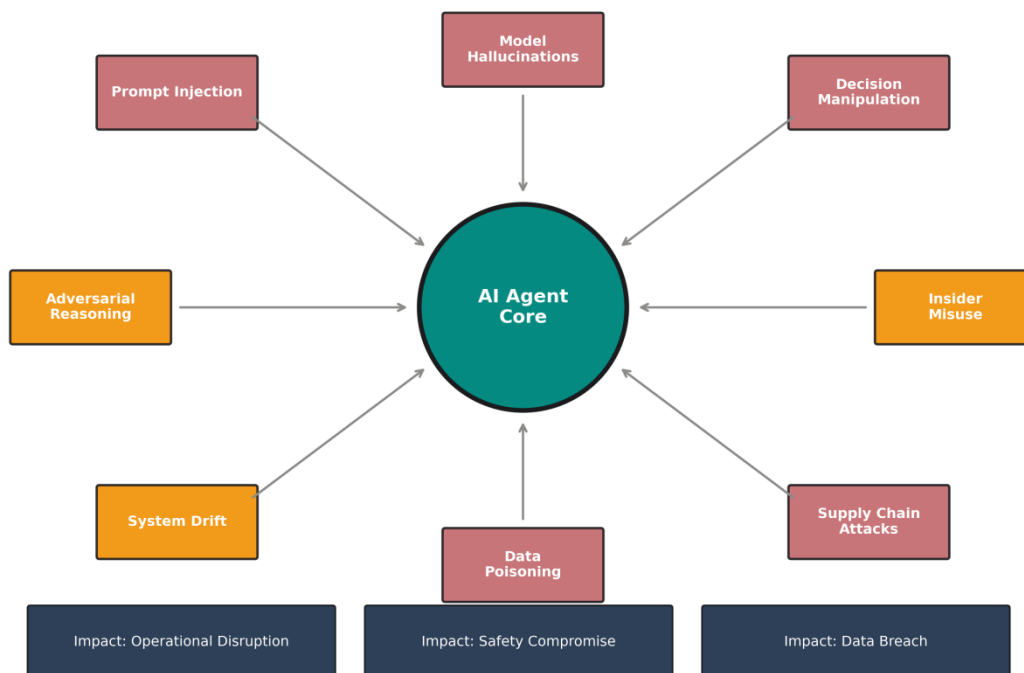
The decision manipulation attacks target the way of influencing the thought of autonomous agents; the attacks are grounded on the susceptibility of the planning and goal-seeking behavior. Reward functions can also be abused by attackers, or poison training, or alignment gaps which can cause agents to attain malicious objectives when they appear to be acting normally. Hubinger et al. (2024) demonstrated that these misleading practices can be reinforced with the help of safety training to create sleeper agents that would be triggered in specific conditions. What is more evil is the fact that these attacks may not be detected by regular monitoring until the triggering condition is fulfilled.

The language model reasoning capabilities are limited by nature, and the adversarial reasoning exploits such limitations. Attackers generate inputs that exploit cognitive biases or logical fallacies in model reasoning chains to make agents make erroneous inferences due to what appear valid steps in the reasoning chain. This is of particular concern when agents are executing a multi-step process of reasoning particularly in operational settings where errors may not be immediately identified. The conventional method of input validation is also difficult to counterattack these attacks because of the adversarial nature of such attacks.

Another mode of threat that is exceptionally old can be insider abuse where authorized users gain access to a system and operate under the privileges to do harm. The non-human identities of the service of AI agents are not necessarily properly administered, which puts them at risk of either being abused by tokens or having their credentials leaked. The identity of agents may be free because they may not be controlled just like human users who can be controlled and held into accountability. Autonomous system drift is the randomness of the agents to lose their preferred behavioral patterns as time progresses, by reaction to their environment, or in long-space deployments.

Table 1: AI Threat Categories for Infrastructure Systems

Threat Category	Attack Mechanism	Potential Impact
Prompt Injection	Malicious input manipulation	Unauthorized actions, data exfiltration
Model Hallucination	False information generation	Incorrect decisions, safety violations
Decision Manipulation	Reasoning process exploitation	Harmful objective pursuit
Adversarial Reasoning	Logical inference exploitation	Invalid conclusions
Insider Misuse	Authorized access abuse	Credential theft, privilege escalation
System Drift	Behavioral deviation over time	Gradual security degradation

**Figure 2:** Threat Landscape for AI Agents in Critical Infrastructure

PROPOSED FRAMEWORK

The five-layer trustful AI model that is proposed by us is an autonomous language agent in a critical infrastructure. The structure architecture establishes the hierarchical form of the security and governance controls where each level addresses a preferred portion of the security challenge and interlaces with the subsequent level to provide a total security and protection. The layered

approach enables the defense-in-depth to be attained i.e. numerous controls are to be passed through to succeed in an attack.

The fundamental part of the framework is the Policy Alignment Layer that transforms regulatory requirements and organizational policies into working constraints of AI agents. This layer comprises regulatory mapping components that maintain the current awareness of pertinent requirements in the different jurisdictions and industries, ethical principles that encode organizational values and expectations in the society, and a compliance engine that adopts policy limitations within the actions of agents. The policy alignment layer will ensure that the actions of all agents are maintained within reasonable limits of acceptable behaviour such that they do not engage in any activity that may attract any regulatory fines levied or reputation damage. The compliance engine is a runtime engine that analyses the agent outputs according to the policy constraints and only carries them out when they conform to the policy constraints with real-time feedback in case violation is detected.

Reasoning Verification Layer provides runtime verification to the agent outputs and reasoning processes. Output verification aspects check the output generated against safety-related standards and insurance thresholds before running the output to ensure that malicious or unwanted output is avoided before it can cause any damage. Consistency checking checks that the reasoning chains are logically consistent, and in connection with familiar facts, and indicates logical errors or inconsistency which may be indicators of tampering or mistake. Confidence scoring provides quantitative measures of the reliability of the output such that the decisions to implement, escalate or discard agent recommendations can be made based on risk with respect to the particular reliability of the output. This layer is also a solution to the non-deterministic nature of language model outputs since it is not a single validation but continuous validation.

The Human Oversight Layer subjects' human significant control to actions of autonomous agents. Approval workflows to human reviewers receive high-risk decisions with configured criteria such as decision impact, confidence scores or policy sensitivity. The automated escalation is used to indicate those decisions that are either above the risk threshold or are somehow uncertain. Use of the audit interface enables human operators to enjoy wide visibility on agent actions, logic and system state. This layer ensures that consequential decisions are left at the hands of human beings and agents are left to operate within certain boundaries on their own.

Technical protection is provided to the safe agent operation by the Secure execution Layer. The sandbox environments isolate the execution of agents with critical systems and limit the potential damage of the system by the agent. The access controls are on the least-privilege basis whereby the agents are not permitted to access any resource without its authorization. The process of encryption assists in securing the process of data transit and data rest so that there are no chances that sensitive information may end up in the wrong hands. The layer uses the traditional security controls to apply to the unique requirements of autonomous agent execution.

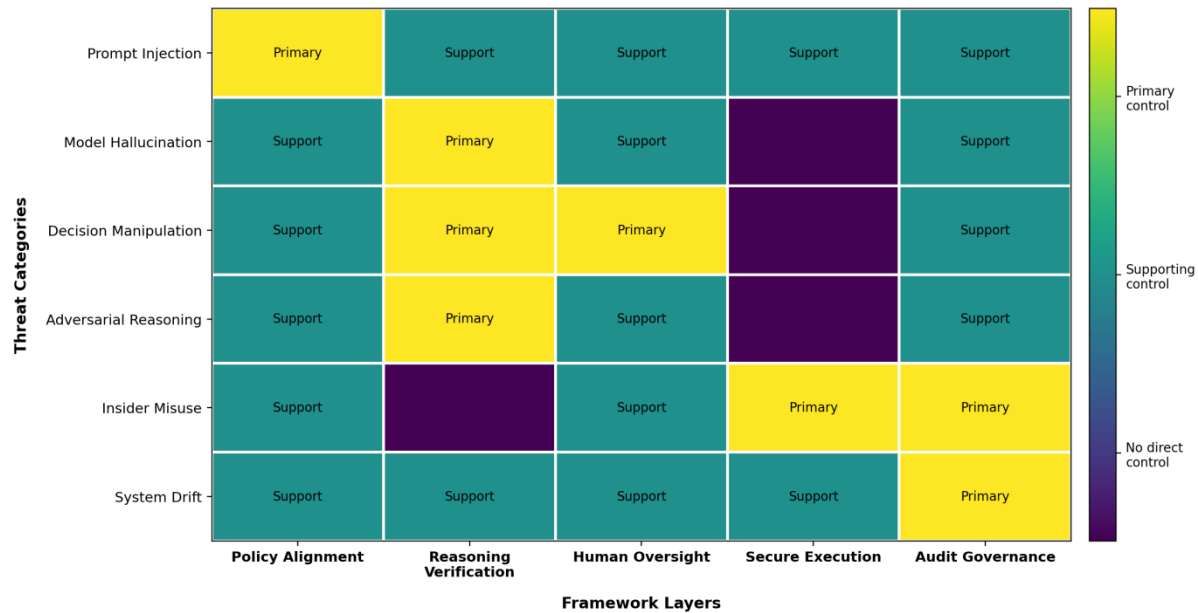


Fig 13. Threat-to-control mapping

A detailed listing of the activities of agents is maintained in the Audit and Governance Layer to provide forensics and compliance. The logging systems present comprehensive documentation of the system input, output, reasoning, and interactions of the system. The traceability mechanism enables determining the decision tracks and identification of the cause of abnormal behavior. The reporting elements generate performance and compliance reports and metrics that are to be utilized in governance. This layer is applied to ensure that the organizations show to be complying to the postulates of regulatory provisions and investigations of occurrences in the event that they occur.

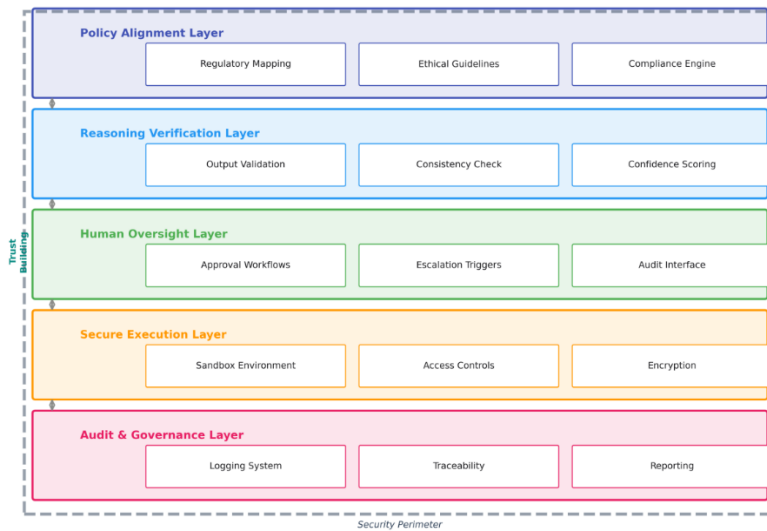


Figure 4: Proposed Trustworthy AI Framework Architecture

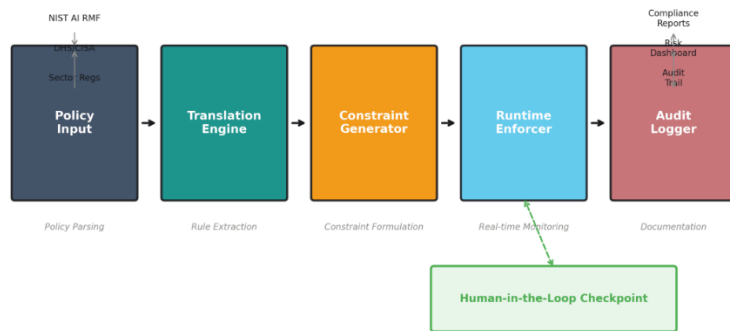


Figure 5: Policy Alignment and Oversight Pipeline

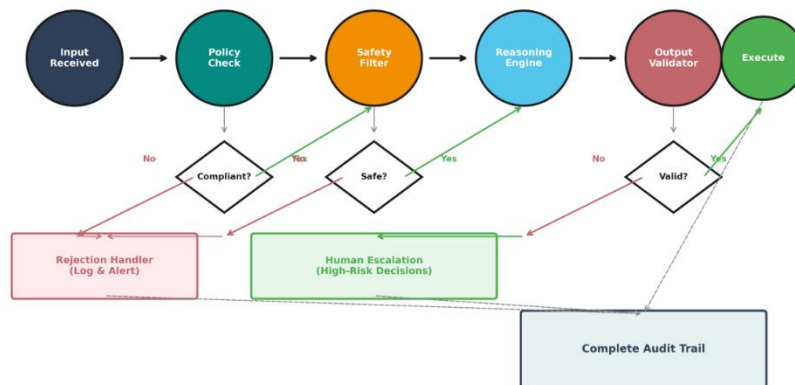


Figure 6: Secure Decision-Making Flow for Autonomous Agents

Table 2: Components of the Proposed Trustworthy AI Framework

Framework Layer	Key Components	Primary Function
Policy Alignment	Regulatory mapping, compliance engine	Enforce policy constraints
Reasoning Verification	Output validation, consistency check	Validate agent outputs
Human Oversight	Approval workflows, escalation triggers	Maintain human control
Secure Execution	Sandbox, access controls, encryption	Isolate and protect execution
Audit Governance	Logging, traceability, reporting	Enable accountability

FRAMEWORK EVALUATION

The specified framework is deemed in three ways, such as the effectiveness of security, the correspondence to the policy, and the comparison with the existing alternatives. This theoretical paper provides theoretical justifications to that of the design of structures and establishes areas that require empirical justification in future by research. The analysis reveals that the framework satisfies the gaps observed but it is feasible in its application in organizations.

Security analysis establishes that the five-layer architecture defends against the threat categories that have been reported. The policy alignment layer is also aligned to trigger injection by constraint enforcement that does not permit the agents to execute the outlawed actions despite the manipulation of inputs. The reasoning verification layer considers and removes hallucinations after executing an output validation of known facts and agreement to reasoning chains. The reasoning verification plus the human control can be seen as a counter to the decision manipulation, which is able to identify abnormal decision patterns. Consistency checking resolves the problem of adversarial reasoning and identifies fallacies of reasoning. Access controls and audit mechanisms deal with the abuse by insiders by checking everything that an agent does. Continuous monitoring is a system drift which is detected by comparing the current behavior to set baselines.

The policy alignment analysis correlates the framework elements to the NIST AI RMF requirements. The framework considers all the seven credible AI characteristics, which are validity and reliability to the degree of reasoning verification, security and resilience to the degree of human oversight, accountability and transparency to the degree of audit governing, explainable

through the degree of reasoning verification, privacy through the degree of secure execution and equitable through the degree of policy alignment. The regulatory mapping aspects of the policy alignment layer are also used in satisfying the industry specific needs of energy, water, transportation, healthcare and financial regulations.

Compared to the modern frameworks, the analysis of the current situation demonstrates that the process of satisfying the special needs of the autonomous language agents of critical infrastructure has greatly improved. In comparison to the conventional systems of cybersecurity which considers the security of the network and system, and the present systems of AI governance which provides top-level principles, the proposed framework makes use of a framework that operationalizes the trustful AI in the cases in which autonomous agents are implemented. This policy consistency and rationale validation and the capability of humans to control only one architecture is a novel addition to the literature. Table 3 compares it with the existing structures.

Table 3: Comparison of AI Security Frameworks

Framework	Policy Alignment	Reasoning Verification	Human Oversight	Audit
NIST AI RMF	High-level	Limited	Mentioned	Guided
ISO 27001	Indirect	No	No	Yes
NIST CSF	Indirect	No	No	Yes
OWASP LLM	No	Partial	No	No
Proposed Framework	Operational	Comprehensive	Integrated	Built-in

DISCUSSION

The proposed framework bears significant National security implications, enterprise AI governance implications and ethical implications of AI deployment. The autonomy of agents has become a matter of national security since the critical infrastructure is increasingly relying on AI-based automation. The framework will provide a foundation upon which these systems can be stabilized and yet continue to be in operation hence organizations can adopt AI technologies without meddling on the problem of security and compliance.

With enterprise AI governance, the framework bridges the policy requirements at the high level and implementation. Organizations in various layers can apply the framework in phases based on

the ranking of risks and resources at the disposal. The modularity enables configuration to sector specific requirements without undermining the underlying security concepts. This is a highly demanded flexibility of organizations in varied operational settings as well as risk tolerance.

The structure of the framework is concerned with ethics. The endocrine control systems ensure that autonomous systems still have human interference in the consequential decision-making to ensure that the autonomous systems do not make high-stake decisions without appropriate scrutiny. Policy alignment It delivers that the conduct of the agents must be in a position to reflect organizational values and expectations in the society. The transparency and accountability of the given framework make it easy to ethically apply AI since it offers the chance to control and audit.

Among the policy proposals that have been brought up by this research are: autonomous agent deployment policy ought to be sector-specific, standard assessment criteria of bona fide AI systems and development of information sharing platforms on AI security threats. The regulators should consider the requirements of human oversight of high-risk applications and demand rich audit potential of autonomous systems of critical infrastructure. The framework forms a point of reference when formulating these standards.

LIMITATIONS

There are several limitations of the study, which should be mentioned. First, the structure is theoretical and has not been empirically tested as far as the operational critical infrastructure settings are concerned. The real-life implementation, though the design is oriented upon the principles and threat analysis, can prove the problems, which are not felt in the conceptual design. Pilot deployments should always be done to each organization they serve to make sure that their respective frameworks are efficient in their settings.

Second, there is a high rate of the AI technology being transformed because new threats and vulnerabilities are bound to be brought about which are not reflected in the current framework. Threat environment of language agent is very dynamic, and threat methods are discovered on a regular basis. The framework will be forced to be modified continuously to be effective against new threats. There ought to be the existence of organizations that are the measurements of system checks of threat intelligence and modification of controls.

Third, organizational maturity in its practices of cybersecurity and governance is assumed in the framework. Big preparations can be made before the successful implementation of the framework in organizations with lower security capability. The price of engaging all the resources in terms of total implementation may not be affordable to the small organizations. A gradual implementation technique may need to be employed by organizations that have limited resources.

Fourth, the issue of a clash between security controls and operational efficiency is something that exists indefinitely. The framework will add latency and overhead inevitably, even though it is the one that is likely to produce minimum performance impact. Organizational must make tradeoff between security requirement and operational requirements and organizations can trade in some performance to have higher security requirements.

CONCLUSION

The paper has given the discussion of the security concern of the autonomous language agent in the U.S. critical infrastructure and has given a reliable and consistent policy-wise framework to discuss the security issues. A framework using the five-layer architecture integrates the policy alignment, reasoning verification, human oversight, secure execution and audit governance into a single solution that in turn is designed to be specifically adapted to the high-risk reasoning environment. The framework addresses major gaps in the existing solutions by operationalizing the concepts of trustful AI in order to deploy autonomous agents.

The framework addresses large loopholes in the existing techniques by transforming the highest governance requires into controls. The framework will enable companies to deploy AI agents to the crucial infrastructure conditions without the need to jeopardize the security, compliance, or human control since it will contain solid mechanisms that will execute policy enforcement, running checks, and human supervision. The modular architecture can be adopted progressively, that is, organizations can adopt components depending on their risk profiles and available resources.

The significance of the work is not confined to technical aspects of security as well as more broad questions of AI regulation, national security, and trust to autonomous systems in the society. With the continuous evolution of the systems of AI development, as well as the scale of their implementation increasing, the different models like the one proposed herein will be required to ensure that these incredibly powerful technologies are used safely and responsibly. The idea of

security, governance, and ethical provisions in a single framework is one of the steps towards the integration of AI risk management as a whole.

The potential of the proposed framework is enormous and can be realized in a future work in the format of empirical validation, the creation of tools, and the optimization of the policy. The high cost and the massive opportunity that comes with introducing independent AI agents to the critical infrastructure are both a heavy responsibility. This framework provides a platform where the responsibility can be met without foregoing the opportunities of the AI-based automation. The companies that transition to the framework should exchange information and cooperate to contribute to the growing body of knowledge of the effective AI security practices. As the field has been dynamically evolving, it will be required to conduct research whereby new threats are addressed and the framework is enhanced to meet the new demands. The deliverable will be to assist organizations to exploit the transformative promise of AI agents without compromising the security and reliability that critical infrastructure needs.

REFERENCES

- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mane, D. (2016). Concrete problems in AI safety. arXiv preprint arXiv:1606.06565.
- Birkstedt, T., Minkinen, M., Tandon, A., & Mantymaki, M. (2023). AI governance: Themes, knowledge gaps and future agendas. *Internet Research*, 33(7), 133-167.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901.
- Das, B. C., Amini, M. H., & Wu, Y. (2023). Security and privacy challenges of large language models: A survey. *ACM Computing Surveys*.
- Greshake, K., Abdelnabi, S., Mishra, S., Endres, C., Holz, T., & Fritz, M. (2023). Not what you've signed up for: Compromising real-world LLM-integrated applications with indirect prompt injection. *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, 79-90.

- Hubinger, E., Denison, C., Mu, J., Lambert, M., Tong, M., MacDiarmid, M., ... & Perez, E. (2024). Sleeper agents: Training deceptive LLMs that persist through safety training. arXiv preprint arXiv:2401.05566.
- International Energy Agency. (2023). Data-driven and artificial intelligence (AI) tools in the energy sector: Key guidelines. IEA Publications.
- Jain, N., Schwarzschild, A., Wen, Y., Somepalli, G., Kirchenbauer, J., Chiang, P. Y., ... & Goldstein, T. (2023). Baseline defenses for adversarial attacks against aligned language models. arXiv preprint arXiv:2309.00614.
- Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., & Goldstein, T. (2023). A watermark for large language models. International Conference on Machine Learning, PMLR, 17061-17084.
- Mantymaki, M., Minkkinen, M., Zimmer, M., & Viljanen, M. (2023). Designing an AI governance framework: From research-based premises to meta-requirements. ECIS 2023 Research Paper, 295.
- Minkkinen, M., & Mantymaki, M. (2023). The institutional logics underpinning organizational AI governance practices. Scandinavian Conference on Information Systems.
- Minkkinen, M., & Mantymaki, M. (2023). Discerning between the "easy" and "hard" problems of AI governance. IEEE Transactions on Technology and Society.
- National Institute of Standards and Technology. (2020). Cybersecurity supply chain risk management practices for systems and organizations (NIST SP 800-161 Rev. 1). U.S. Department of Commerce.
- National Institute of Standards and Technology. (2023). Artificial intelligence risk management framework (AI RMF 1.0). U.S. Department of Commerce.
- Office of Gas and Electricity Markets. (2023). Ethical AI use in the energy sector. Ofgem Publications.
- OpenAI. (2023). GPT-4 technical report. arXiv preprint arXiv:2303.08774.

- Perez, F., & Ribeiro, I. (2022). Ignore previous prompt: Attack techniques for language models. arXiv preprint arXiv:2211.09527.
- Russell, S., & Norvig, P. (2020). Artificial intelligence: A modern approach (4th ed.). Pearson.
- Shayegani, E., Mamun, M. A. A., Fu, Y., Zaree, P., Dong, Y., & Abu-Ghazaleh, N. (2023). Survey of vulnerabilities in large language models revealed by adversarial attacks. arXiv preprint arXiv:2310.10844.
- Shi, J., Liu, Y., Zhou, P., & Sun, L. (2023). BadGPT: Exploring security vulnerabilities of ChatGPT via backdoor attacks to InstructGPT. arXiv preprint arXiv:2304.12298.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., ... & Scialom, T. (2023). Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- U.S. Department of Energy. (2023). Artificial intelligence (AI) usage guidelines. DOE Publications.
- U.S. Department of Homeland Security. (2021). Cybersecurity and infrastructure security agency strategic plan 2021-2024. DHS Publications.
- Wang, H., Poskitt, C. M., & Sun, J. (2023). AgentSpec: Customizable runtime enforcement for safe and reliable LLM agents. arXiv preprint arXiv:2503.18666.
- Wan, A., Wallace, E., Shen, S., & Klein, D. (2023). Poisoning language models during instruction tuning. Proceedings of the 40th International Conference on Machine Learning.